# Correlation and Beyond: Positive Definite Dependence Measures for Inference, Scenarios, and Stress Testing for Financial Portfolios

## JD Opdyke, Chief Analytics Officer, DataMineit, LLC

Monograph: 1<sup>st</sup> Draft, November, 2021; Current draft, July 2025

<u>Disclaimer</u>: The views presented in this monograph are those of the sole author, JD Opdyke, and do not necessarily reflect those of particular institutions.

<u>Acknowledgements</u>: This work is dedicated to my family – to my one and only daughter, Nicole, my one and only son, Ryan, and my one and only wife, Toyo, for whom Euler and Gauss will always be close to her heart: I extend my pride, love, and deepest gratitude for your unwavering support.

Professional Biography: JD Opdyke is a senior data scientist of over 30 years in the investment and risk analytics space. Currently Chief Analytics Officer at DataMineit, LLC, JD has strong and extensive experience across major financial verticals (capital markets, banking, and insurance) as well as decades of risk modeling and data science consulting expertise in related industries. JD has built and led several senior quant teams, published 14 peer reviewed journal papers and book chapters, several of which were voted 'Paper of the Year' by panels of experts, and is a frequently invited speaker and presenter at top quant and risk conferences globally.

JD earned his Bachelor's degree, with honors, from Yale University, his Master's degree from Harvard University where he was awarded multiple paid, competitive Fellowships, and he completed a postgraduate fellowship in MIT's graduate mathematics department as an Advanced Study Program Fellow. He serves as review editor of several journals, including Artificial Intelligence in Finance.

#### TABLE OF CONTENTS

- 1. Summary and Organization
- 2. Introduction and Background
  - a. NAbC: Summary of Methodology
  - b. Types of Dependence Measures
    - i. Monotonic Measures
    - ii. Tail Dependence Measures
    - iii. Distance-based and Other New Measures
    - iv. Asymmetric, Directional Measures
- 3. Estimation
  - a. Covariance and Pearson's Correlation
  - b. Other Dependence Measures

#### 4. NAbC: (Robust) Statistical Inference

- a. Brief Literature review of Pearson's Matrix: Distributional Results and Sampling Algorithms
  - i. Distributional Results
  - ii. Sampling Algorithms
  - iii. Distributional Results, More General Conditions
  - iv. Sampling Algorithms, More General Conditions
- b. NAbC: Pearson's Correlation, the Gaussian Identity Matrix
  - i. Correlations to Angles, Angles to Correlations
  - ii. Fully Analytic Angles Density, and Efficient Sample Generation
  - iii. Matrix-level p-values and Confidence Intervals
- c. NAbC: Pearson's Correlation, Real-World Financial Data, Any Matrix Values
  - i. Nonparametric Kernel Estimation
- d. NAbC: Any (Positive Definite) Dependence Measure, Any Data, Any Matrix Values
  - i. Spectral and Angles Distributions, Examples from Other Dependence Measures
- e. NAbC: Fully General Conditions, Statistical Comparison of Two Matrices
- f. NAbC Remains "Estimator Agnostic"

#### 5. NAbC: Granular, Fully Flexible Scenarios, Reverse Scenarios, and Stress Testing

- a. Review of Existing Methods
- b. A New Method for Fully Flexible Scenarios

## 6. NAbC Example: Kendall's Tau p-values & Confidence Intervals, Unrestricted & Scenariorestricted, One- and Two-Sample Tests

- 7. NAbC: Beyond 'Distance' to Generalized Entropy
- 8. NAbC: Future Research and Additional Applications
- 9. Conclusions
- 10. References

#### 1. Summary and Organization

We live in a multivariate world, and effective modeling of financial portfolios, including their construction, allocation, forecasting, and risk analysis, simply is not possible without explicitly modeling the dependence structure of their assets. Dependence structure can drive portfolio results more than many other parameters in investment and risk models - sometimes even more than their combined effects. But the literature provides relatively little to define the finite-sample distributions of dependence measures in useable and useful ways under challenging, real-world financial data conditions. Yet this is exactly what is needed to make valid inferences about their estimates, and to use these inferences for essential purposes such as hypothesis testing, dynamic monitoring, realistic and granular scenario and reverse scenario analyses, and mitigating the effects of correlation breakdowns during market upheavals. This work develops a new and straightforward method, Nonparametric Angles-based Correlation (NAbC), for defining the finite-sample distributions of any dependence measure whose matrix of pairwise associations is positive definite (e.g. Pearson's, Kendall's, Spearman's, the Tail Dependence Matrix, and others). The solution remains valid under marginal asset distributions characterized by notably different and varying degrees of serial correlation, non-stationarity, heavy-tailedness, and asymmetry. Importantly, NAbC provides p-values and confidence intervals at both the matrix level and the pairwise cell level, for both one and two-sample tests, with analytical consistency across levels. Finally, NAbC maintains validity even when selected cells in the matrix are frozen, thus enabling flexible, granular, and realistic scenarios and stress tests. NAbC stands alone in providing all of these capabilities simultaneously, and should prove to be a very useful means by which we can better understand and manage financial portfolios in our multivariate world.

This monograph is organized as follows: Section 2 below is an Introduction and Background that discusses various challenges of the problem to be addressed, along with an overview of the dependence measures that define the range of application of NAbC as its solution. The next section treats estimation of these dependence measures, making some suggestions for possible improvements in this area while also making clear that estimation is beyond the scope of NAbC's core methodology: NAbC provides the sampling distributions of these estimates to enable inferences about them, not the estimates themselves. Section 4 develops NAbC, first with a brief and relevant literature review, followed by NAbC's fully analytic derivation for a narrow but foundational special case. This solution serves to make seamless the transition to the general case at the end of Section 4. Section 5 demonstrates how the general solution can be applied in fully flexible scenarios within the framework of the all-pairwise matrix, something that has not been done in the current literature. This is followed in Section 6 by a full empirical example covering all of NAbC's inferential capabilities. Section 7 explains how NAbC can be used to define a new, generalized entropy with probabilistic meaning and motivation, giving it many advantages over its distance-based competitors. Finally, I conclude in Section 9 with directions for future research, as well as direct applications of NAbC in related areas, such as causal modeling.

#### 2. Introduction and Background

Dependence structure is widely acknowledged as a central driver of portfolio results in investment and risk models.

"Correlation is one of the most important, if not the most important, risk factor in finance, driving everything" (Packham & Woebbeking, 2023, p.1).

"Having a few good uncorrelated return streams is better than having just one, and knowing how to combine return streams is even more effective than being able to choose good ones (though of course you have to do both)." (Dalio, 2017).

"...choosing an asset pool consisting of (as many as possible) assets with pairwisely uncorrelated or even negative-correlated returns...becomes a primary objective..." (Yu et al., 2025).

Despite this, the literature provides relatively little to define the finite-sample distributions of commonly applied dependence measures, like (Pearson's) correlation, in useable and useful ways <u>under challenging, real-world financial data conditions</u>.<sup>1</sup> Yet this is exactly what is needed to make valid inferences about their estimates, and to use these inferences for a myriad of essential purposes, such as hypothesis testing, dynamic monitoring, realistic and granular scenario and reverse scenario analyses, as well as mitigating the effects of correlation breakdowns during, and preferably before, market upheavals (which is when we need valid inferences the most).

The goal of this monograph is to fill this gap, both in the literature and in practice, by developing a new and straightforward method – Nonparametric Angles-based Correlation ("NAbC") – defined by eight critically important characteristics listed below. When satisfied simultaneously, these characteristics not only elevate the analytical rigor applied to dependence structure, placing it on par with that applied to the other parameters in investment and risk models, but also allow for practical, non-textbook application to real-world portfolios under conditions where other methods simply cannot be applied (due to their unrealistic assumptions and/or overly restrictive requirements). Yet NAbC's foundations rest squarely on very well established results in the relevant literatures, making its methodology transparent and intuitive, and its application straightforward.

1. NAbC remains valid under challenging, real-world data conditions, with marginal asset distributions characterized by notably different and varying degrees of serial correlation, non-stationarity, heavy-tailedness, and asymmetry.

<sup>&</sup>lt;sup>1</sup> I take 'real-world' financial returns data to be multivariate with marginal distributions that can vary notably from each other, and change in time, in their degrees of heavy-tailedness, serial correlation, asymmetry, and non-stationarity. These obviously are not the only defining characteristics of such data, but from a distributional and inferential perspective, they remain some of the most challenging, especially when occurring concurrently as they do in non-textbook settings.

- 2. NAbC can be applied to ANY dependence measure with a matrix of all-pairwise associations that is positive definite,<sup>2</sup> including long established measures for which positive definiteness is analytically proven (e.g. the foundational Pearson's product moment correlation matrix (Pearson, 1895), rank-based measures like Kendall's Tau (Kendall, 1938) and Spearman's Rho (Spearman, 1904), and the tail dependence matrix (see Embrechts, Hofert, and Wang, 2016, and Shyamalkumar and Tao, 2020). This also includes newer measures for which positive definiteness must be empirically validated, such as Chatterjee's correlation (Chatterjee, 2021) and its variants (Pascual-Marqui et al., 2024), the improved Chatterjee's correlation (Xia et al., 2024), Lancaster's correlation(s) (Holzmann and Klar, 2024), and Szekely's distance correlation (Szekely, Rizzo, and Bakirov, 2007) and its variants (such as Sejdinovic et al., 2013, and Gao and Li, 2024).
- 3. NAbC remains "estimator agnostic," that is, valid regardless of the sample-based estimator used to estimate any of the above-mentioned dependence measures. So to be clear, NAbC is not an estimator of the correlation matrix or other dependence measures: rather, it is a method for obtaining the finite-sample distribution of the estimates generated by various estimators, so that inferences can be made about their estimated values.
- 4. NAbC provides valid confidence intervals and p-values at both the matrix level and the pairwise cell level, with analytic consistency between these two levels (i.e. the confidence intervals for all the individual cells define that of the entire matrix, and the same is true for the p-values; this effectively facilitates, and in many cases makes possible, granular and targeted attribution analyses).
- 5. NAbC provides valid confidence intervals and p-values not only for one-sample tests against matrices of fixed, assumed 'true' values, but also for two-sample tests comparing two matrices, so that we can assess inferentially whether dependence structures truly are different, for example, across different sectors or segments of our businesses.
- 6. NAbC provides a one-to-one quantile function, translating a matrix of all the cells' cumulative distribution function (cdf) values to a (unique) correlation/dependence measure matrix, and back again, enabling precision in reverse scenarios and stress testing, as well as informed and targeted 'what if' analyses.
- 7. All the above results remain valid even when selected cells in the matrix are 'frozen' for a given scenario or stress test that is, unaffected by the scenario thus enabling flexible, granular, and realistic scenarios.

<sup>&</sup>lt;sup>2</sup> Note that "positive definite" throughout this monograph refers to the dependence measure calculated on the matrix of all pairwise associations in the portfolio, that is, calculated on a bivariate basis. Some of the dependence measures addressed in this monograph (e.g. Szekely's correlation, variants of Chatterjee's, and others) can be applied on a multivariate basis (sometimes even in arbitrary dimensions), for example, to test the hypothesis of multivariate independence. But "positive definite" herein is not applied in this sense (see for example Cardin, 2009), and I explain below some of the reasons for using the dependence framework of all pairwise associations, which is highly flexible, and allows for more precise attribution and intervention analyses.

8. NAbC remains valid not just asymptotically, i.e. for sample sizes presumed to be infinitely large, but rather, for the specific sample sizes we have in reality (for full-rank matrices with n>p)<sup>3</sup>, enabling inferentially reliable application in actual, real-world, non-textbook settings.

The alternative to a method satisfying the eight objectives above, simultaneously, is to use a piecemeal, incomplete patchwork of disparate derivations of distributions, some asymptotic, some not, valid under typically restrictive, differing, and unrealistic data conditions for only a few of the widely used dependence measures (as distributions for many have not yet been derived). This patchwork approach not only materially limits the scope of possible comparative analyses, but also the degree to which it can be truly ceteris paribus. Since differing assumptions are confounded with the capabilities of the methods themselves, it is impossible to know where the effects of the assumptions end and those of the different methods begin. This is exacerbated by the unwieldy, opaque, and difficult-to-implement nature of many of these solutions.

NAbC circumvents all of these problems with a single, unified, and straightforward method for dependence structure inference that, compared to its more limited and narrowly defined competitors, simultaneously and dramatically increases i. robustness, ii. scenario flexibility, iii. accuracy in attribution analyses, and iv. targeted precision in 'what if' intervention analyses, all while enabling v. ceteris paribus analyses across a very broad range of dependence measures (including those listed in 2. above).

Before explaining how NAbC's methodology accomplishes this, however, it is important to ask why there is a dearth of "real-world effective" methodology in this setting, as it will inform and clarify the explanations throughout this monograph. If we were to define a problem statement here, it would be: define the finite sample distributions of all positive definite measures of dependence structure, robustly under real-world financial data conditions, that remain valid regardless of the estimators used, and even if the co-movement of selected pairs of variables is 'frozen', i.e. scenario-restricted. While this objective admittedly remains broad, asking why this hasn't been done previously remains a fair and important question.

Financial markets certainly have seen more extreme downturns in recent decades than many would have predicted ex ante (e.g. Black Monday (1987), Tech Bubble (2000), Housing Bubble (2008), Covid (2020)), during which correlation breakdowns have been well documented, their very material effects measured and assessed (see for example Feng & Zeng, 2022, and Packham & Woebbeking, 2023), and the importance of mitigation efforts widely discussed, considered, and acted upon (see Greenspan, 1999; BIS, 2011a; and EBA-CRR, 2013). What's more, practitioners, academics, and regulators have a long history of bringing analytic and probabilistic rigor to bear when analyzing and estimating *the other* parameters of our portfolio risk and investment models. There is no shortage of empirical research defining, for example, various estimators of the tail indices of a portfolio's marginal distributions, and deriving their associated p-values, confidence intervals, and statistical power and level. When rigorously and properly estimated, these tools are highly actionable, providing invaluable guidance in decision-

<sup>&</sup>lt;sup>3</sup> Recall that this condition is required for the all-pairwise matrix to be positive definite.

making and mitigation efforts. But what do we find when we look for those same tools, say, for an entire matrix of pairwise Kendall's tau values, to make decisions based on answers to questions such as, "Has this Kendall's matrix shifted in the past week? What is the probability of observing the movement we observed, given that our baseline estimate is true? Does this meet our probabilistically defined threshold for a 'breakdown'? What are the two (upper and lower) Kendall's matrices that capture 95% of the conditional sample variation in this setting? How far beyond these bounds, <u>probabilistically</u>, do the Kendall's matrices for each of our scenarios lie? Given our distributions of losses/returns, does a tail dependence matrix better capture what we are trying to measure here, and can we conduct a ceteris paribus analysis, using the exact same distribution-defining methodology, to compare the statistical power of these two dependence measures under the various relevant data conditions?" If we require the p-values and confidence intervals and rigorous, <u>probabilistic</u> answers to these questions to be valid under challenging, real-world financial data conditions, the current literature provides relatively little. Given the need, as well as the rigor applied to other areas of portfolio analytics, this arguably is surprising.

On the other hand, the possible explanations for this dearth of useable and useful methodology are not entirely unreasonable. First, the multivariate nature of this problem arguably makes it more challenging than those related to modeling some of the other (univariate) parameters of investment and risk portfolios. Even though each cell value of the dependence matrix is a bivariate association, we are measuring all the pairwise associations in the portfolio simultaneously, and the values of the cells are, in non-trivial cases, all interrelated, making this a complex, multivariate problem. Immutable mathematical requirements for this setting, such as positive definiteness, arise, and make deriving and simulating the distribution of the all-pairwise matrix a non-trivial task. This is especially true if we require, as we should, that the p-values and confidence intervals of each and every cell of the matrix are consistent with those of the entire matrix when taken as a whole. Additionally, requiring that the finite sample distribution of the matrix (which makes possible the calculation of the p-values and confidence intervals) remain valid <u>under challenging, real-world financial data conditions</u> adds significantly to the nontrivial nature of the problem. Distributions of dependence measures are more readily derived when we can assume that returns are, say, multivariate normal, or at least independent and identically distributed (iid). It is another matter entirely when the portfolio's marginal distributions vary notably from each other, while also changing over time in their degrees of heavy-tailedness, serial correlation, asymmetry, and non-stationarity. Yet this is exactly the empirical challenge of actual financial portfolios. In fairness, the literature does provide many solutions under mathematically convenient conditions, which more narrowly define and restrict both in the distributional characteristics of the underlying returns data as well as the assumptions made regarding the values of the all-pairwise matrix. But these largely unrealistic assumptions limit practical, real-world application, which is exactly the motivation for this monograph.

Another complicating factor is the requirement that the method defining the finite sample distribution of the dependence measures is the same across all those in practical usage: in this case, all those dependence measures for which the all-pairwise matrix is positive definite. This arguably covers all that

JD Opdyke, Chief Analytics Officer

Page **7** of **92** 

**Correlation and Beyond** 

could conceivably be used and useful in the financial setting. This universality is certainly desirable, but it also increases the challenge of deriving the methodology. Still, this remains a crucial requirement as it provides the ability to conduct all-else-equal analyses comparing the performance of different dependence measures under controlled conditions: we can be certain that material and statistically significant differences in results are due to the dependence measures themselves, rather than the very distinct methodologies we typically would have to use (based on the current literature) to define their distributions. But nothing in the extant literature provides this broad ceteris paribus capability.

Finally, one of the major uses of dependence measures and their all-pairwise matrices is in defining scenarios and reverse scenarios.<sup>4</sup> These remain central for and critical to all manner of risk analyses, and fully flexible scenarios require the ability to 'freeze' groups of selected cells of the all-pairwise matrix while allowing others to vary. For example, many of the pairwise cells that will change dramatically, in both direction and magnitude, under a Covid-like scenario will be completely unaffected under a housing bubble (see Feng & Zeng, 2022, and Pramanik, 2024), and scenario analysis must be able to validly define the finite sample distribution of the all-pairwise matrix under both types of scenario-restricted conditions. However, no existing method allows for this without inadvertently affecting the other 'peripheral' cells of the all-pairwise matrix, rendering the associated inferences for the scenario(s) invalid, and decisions based on them potentially harmful. Granular flexibility in scenario definition, at the level of the pairwise cells, and *the valid distribution of the associated, scenario-restricted matrix*, is a necessity if we are to accurately capture the fundamentally different nature of disparate correlation breakdowns, and accurately assess, forecast, and mitigate their impacts.

So given the breadth of the problem statement, perhaps it is not so surprising that we have comparatively little in the way of real-world solutions to this problem. This is true not only of the extant literature, but also in financial practice, which often relies on ad hoc, largely qualitative, and 'judgmental' approaches to specifying and utilizing dependence structure. When quantitative approaches are used, those selected typically are the most conveniently implemented methods that remain valid only within narrowly defined boundaries and/or requiring unrealistic but mathematically convenient assumptions, such as i. the distributions derived are only asymptotically valid (i.e. assume infinitely large sample sizes); ii. they require very restrictive and/or unrealistic assumptions about the marginal returns distributions of the portfolio (e.g. that they are multivariate Gaussian, or elliptical; or even that they are independent and identically distributed ("iid"), or all symmetric, or all stationary, or not serially correlated, etc.); iii. they require very restrictive and/or unrealistic assumptions about the values of the dependence measures themselves (e.g. the cells are all zeros, or all have the same value, or follow very discrete and limited block structures); iv. they estimate the all-pairwise matrix in ways that do not guarantee its positive definiteness, or violate other fundamental mathematical requirements (e.g. unit diagonals); or v. most

<sup>&</sup>lt;sup>4</sup> Scenarios typically are designed to answer questions of the type, "What loss is associated with, say, the 99.5% tile of the loss distribution?" while reverse scenarios answer questions of the type, "What percentile of the loss distribution produces a loss of \$X?" The dollar amounts referenced in the latter typically are associated with specific extreme or catastrophic events, such as insolvency or the failure of a major business line or geography.

typically, they require multiple of these restrictive and/or unrealistic assumptions combined. Many of these more narrow solutions are mathematically elegant, but our goal herein is to obtain an actual solution that works and remains inferentially valid under the challenging data conditions of actual financial portfolios, which are 'real-world' messy and decidedly less elegant.

Importantly, note that the original statement in this section, "effective modeling of financial portfolios, including their construction, allocation, forecasting, and risk analysis, simply is not possible without explicitly modeling the dependence structure of their assets" applies to all frameworks for portfolio analysis, even those that may not always make explicit their estimation of, or their reliance on, dependence structure. For example, some path dependent approaches generate distributions of portfolio results based in large part, or even primarily, on (usually subjectively defined) probabilities associated with various scenarios, without explicitly defining dependence structure. But such approaches still make many *implicit* assumptions regarding dependence structure, such as that it does not change from one period to the next, or that it does not change under one scenario versus another, or that, even if (Pearson's) correlations may be controlled via 'views' specified in the model, other measures of dependence, such as tail dependence, are not unwittingly changed from one period to another (even if this is unlikely when 'views' on volatilities are changed). Whether implicit, explicit, or indirectly explicit via 'views' on other parameters, all such assumptions about dependence structure will affect simulated results, which consequently always should be made fully explicit in any model (see Meucci, 2010b, and Vorobets, 2025, for examples), even if only for expost testing using NAbC to ensure that the effects of (possibly changing) dependence structure are not (unknowingly) confounding results.

But beyond and in addition to simply avoiding confounding, NAbC provides such models with statistical control and inferential validity when specifying dependence matrix values based on scenario 'views'. As opposed to ad hoc or judgement-based matrix values, the 'view' of an extreme correlation/dependence matrix should be defined probabilistically, based directly on its finite sample distribution, which NAbC provides. For example, the 'view' of a correlation matrix corresponding to an extreme scenario, when used as an input to a path dependent simulation, should be a percentile (say, 99%tile) of the distribution of the all-pairwise matrix, as provided by NAbC's quantile function. All that NAbC needs to define the values of this matrix are the data generating mechanism and the null hypothesis (i.e. the baseline values of the matrix). Conversely, NAbC also can provide the cdf value (percentile) of an all-pairwise matrix whose values are specified for a specific scenario, thus ensuring that it is sufficiently 'extreme,' or in the alternative, not too 'extreme,' (e.g. is it the 95%tile? Or the 99.999%tile?) for the scenario being tested. NAbC provides both: the matrix. Only in this way does a 'view' on dependence structure retain objective meaning regarding its relative size, not to mention its inferential validity, as opposed to being informed by qualitative, subjective judgements or ad hoc procedures.

#### 2.a. NAbC: Summary of Methodology

Both conceptually and in its implementation, NAbC remains a straightforward method based directly on very well-established results in the relevant literatures. Its innovation and originality come less from new derivations and more from the careful assembly of these well-established results, which contributes to its very broad range of application.

NAbC is not an estimator. It does not, for example, provide estimates of the values of a Pearson's correlation matrix. Rather, given a well-estimated matrix, and its known or well-estimated data generating mechanism, NAbC provides the sampling distribution of the matrix. It does this for any positive definite dependence measure, and under very general conditions, based directly on challenging, real world financial returns data without restrictive assumptions or distorting data transformations. This finite sample distribution allows statistical inferences to be made about its values. NAbC provides confidence intervals and p-values at both the level of the entire matrix, and the level of the individual cells, simultaneously, and these results are consistent across these two levels, as the former are based directly on the latter.

NAbC accomplishes the above by obtaining the distributions of the ANGLES between each of the pairwise data vectors in a portfolio, rather than focusing on the values of the dependence measures themselves.<sup>5</sup> Every positive definite (dependence measure) matrix can be translated, cell-to-cell, to a matrix of such angles, which is simply the matrix analogue to the well-known cosine similarity formula (see Pourahmadi and Wang, 2016; Ghosh et al., 2021; Rapisarda et al., 2007; Tsay and Pourahmadi, 2017; and Zhang et al., 2015).

These angles have very useful properties in this setting. First, unlike the dependence measure values themselves, they are random variables whose multivariate relationship is one of independence.<sup>6</sup> This allows for the very straightforward construction of the multivariate distribution of the entire matrix of angles, and consequently, that of the entire dependence measure matrix. In addition, the method for obtaining the angles relies on calculation of the Cholesky factor, which automatically enforces positive definiteness as the sample space remains on the unit hyper-hemisphere (in contrast, ex post enforcement of positive definiteness typically distorts the sample space and consequently, invalidates inference). Finally, we lose no information when using angles here as they contain all the information that

<sup>&</sup>lt;sup>5</sup> Note that the geometric relationship between data vectors (or more generally, 'flats') and the angles between them is long established, going back at least to Jordan (1875) but more recently associated with canonical correlation analysis (see Afriat, 1957; Björck and Golub, 1973; and Knyazev and Argentati, 2002).

<sup>&</sup>lt;sup>6</sup> This independence is well established in the literature. Zhang et al. (2015) (supplementary material) and Rapisarda et al. (2007) use a geometric interpretation of the correlation matrix, based on (orthogonal) Givens rotations, to explain in detail the relationship between correlations and angles as well as why the angles distributions are multivariate independent.

is in the original data, sans scale, and scale does not, and should not, matter for measuring dependence.<sup>7</sup>

In Section 4.b below, I derive NAbC's fully analytic solution for a special but foundational case – that of the Gaussian identity matrix. This provides, together for the first time, the probability density function (pdf), cumulative distribution function (cdf), and quantile function (inverse cdf) of the angles distributions, and consequently, those of the correlation matrix (shown in (29) below). These provide confidence intervals and p-values for both the individual cells and the entire matrix, simultaneously (shown in (30)-(33) below). This all is presented in an interactive spreadsheet (url link provided below), with fully transparent formulae, in which the user can input A. sample size (as long as n>p, i.e., the matrix remains full rank, which is required for positive definiteness); B. the correlation matrix values (to obtain a correlation matrix); and D. alpha critical values (to obtain confidence intervals on the individual cells and the entire matrix, simultaneously). As long as the matrix is positive definite, NAbC provides p-values and confidence intervals, at both levels, as well as a measure of generalized entropy described in Section 7 below.

Beyond this specific case, the fully general solution provided by NAbC conceptually is the same, for any real-world financial data conditions, for any values of the matrix, and for any positive definite dependence measure: the only difference is that the angles distributions now are defined nonparametrically, via non-parametric kernels (shown in (35)-(37) below). These are obtained via a set of simulations using the estimated dependence measure matrix, and its data generating function, as described in the steps below (these are described in more detail in Section 4.c):

## 5 Steps for Obtaining Angles Distributions

- Simulate N samples (N=10,000 typically is sufficient) based on the dependence matrix and the data generating mechanism (each can either be specified/known, or well estimated). For example, N samples from a correlated multivariate gaussian distribution (with p variables representing the p assets in the portfolio), with correlation matrix R, and n observations in each sample (note than n>p always for our purposes as non-full rank matrices will not be positive definite).
- 2. Calculate the corresponding N all-pairwise dependence matrices, and their Cholesky factorizations, and transform each of these factorizations into a lower triangle matrix of angles (this is a straightforward and well-established calculation shown in Section 4.b in (21), (22), and Table A).
- 3. Fit a kernel density to each cell of the matrix of angles based on the N values obtained from the N samples in 2 (there will be p(p-1)/2 cells, where p is the dimension of the matrix).<sup>8</sup>
- 4. Generate N samples, with n observations each, based on the kernel densities in 3.

<sup>&</sup>lt;sup>7</sup> Scale invariance is proved and widely cited for Pearson's rho, Kendall's tau, and Spearman's rho (see Xu et al., 2013, and Schreyer et al., 2017 for examples).

<sup>&</sup>lt;sup>8</sup> Algorithms for sample generation based on commonly used kernels (e.g. the Gaussian and Epanechnikov) are widely known. An example of the latter is simply the median of three uniform random variates (see Qin and Wei-Min, 2024).

5. Convert each of the N samples from 4. back to a re-parameterized Cholesky factorization, and then multiply it by its transpose to obtain a set of N validly sampled dependence matrices (shown in Section 4.b in (21), (22), and Table A). Positive definiteness is enforced automatically as the Cholesky factor places us on the unit hyper-hemisphere. All sample generation hereafter uses just 4. and 5.

The samples of correlation/dependence measure matrices from 5. will follow the same distribution as those generated in 2., but after the kernel densities are fit once in 3., generating samples based on 4. and 5. is orders of magnitude faster than relying on direct simulations in steps 1. and 2. More importantly, using 4. and 5., rather than 1. and 2., allows for valid probabilistic inference, both at the cell level and at the matrix level, due to the independence of the angles distributions. NAbC always translates dependence measure matrices to matrices of angles, then obtains angles distributions, and then the sampled angle matrices are converted back to correlation/dependence measure matrices. The latter are never directly perturbed because they provide no inferential capability, because the distributions of their cells are not multivariate independent. Neither direct data simulation (step 1.) nor a cavalier 'bootstrap' of samples generated from step 1. can change this: it is the independence of the corresponding angles distributions that allows for probabilistic inference here.

This reliance on angles, and their subsequent transformation to correlation/dependence measure values, allows us to isolate the distribution of the correlation/dependence measure matrix, for probabilistic inference, without touching any other distributional aspect of the data, which is the point of the methodology. Several papers in the literature also use spherical angels for similar purposes in this setting (see Lan et al., 2020, and Ghosh et al., 2021), but as described in detail in Section 4.a below, they have notable limitations relative to NAbC. The most important of these is NAbC's unmatched ability to implement fully flexible scenarios that allows us to 'freeze' the values of any combination of cells within the framework of the all-pairwise matrix, and still obtain the inferentially valid finite sample distribution of the (rest of the) matrix. I am not aware of any other method in the extant literature that can provide this capability, which is arguably necessary for granular and realistic (reverse) scenarios and stress testing. This is described in detail in Section 5, and is one of the reasons NAbC utilizes the framework of the all-pairwise matrix, and stress multivariate dependence structures.

The more detailed description of NAbC's implementation in Sections 4 and 5 below demonstrates how it simultaneously satisfies all eight of the critical objectives listed above. An important ancillary benefit of its broad range of application is that a single methodology now allows for ceteris paribus analyses, both for comparing different estimators of the same dependence measure, and for comparing different dependence measures, in many cases where such analyses previously were not possible (or extremely unwieldy). I review many of the most relevant, useful, and commonly used dependence measures in this setting below.

#### 2.b. Types of Dependence Measures

Measures of association, otherwise known as dependence measures, are as old as modern statistics itself (see Pearson, 1895). They provide a quantitative assessment of how variables move together or in opposite directions over time. I address their relation to causal mechanisms in later sections, and merely note here that they remain distinct from what are often called 'metrics' or 'distance metrics,' even though the two are sometimes confused.<sup>9</sup>

For reasons discussed more thoroughly below in the following sections, the dependence measures covered in this monograph, and for which NAbC's application remains valid, include those for which the all-pairwise matrix is positive definite.<sup>10</sup> One could argue that all dependence measures in common usage are, in fact, positive definite, at least under relevant, real-world data conditions, making this requirement de facto non-restrictive. It could even be viewed as a test of appropriateness of usage in applied finance, since many situations for which it demonstrably does not hold (e.g. cases of perfect linear dependence) are degenerate cases in other ways as well. But I avoid such debates herein, and merely state that the validity of NAbC's application does require positive definiteness, and that this includes both those measures for which positive definiteness has been proven analytically, and those for which such proofs do not (yet) exist and thus, which require testing and verification of positive definiteness empirically. The former group includes long established measures such as Pearson's product moment correlation matrix (Pearson, 1895), rank-based measures like Kendall's Tau (Kendall, 1938) and Spearman's Rho (Spearman, 1904), as well as measures designed to capture highly non-linear dependence such as the tail dependence matrix (see Sabato et al., 2007, for proofs of the first three, and Embrechts et al., 2016, for a proof of the latter). The second group includes newer measures designed to capture cyclical and other types of non-linear dependence such as Chatterjee's correlation (Chatterjee, 2021), Lancaster's correlation (Holzmann and Klar, 2024), and Szekely's distance correlation (Szekely, Rizzo, and Bakirov, 2007) and their many variants (such as Sejdinovic et al., 2013, and Gao and Li, 2024). In the end, however, as long as the values of the matrix being evaluated and used by NAbC render it positive definite, NAbC will 'work.' In the extensive empirical analyses performed herein on the second

<sup>&</sup>lt;sup>9</sup> Even though 'metrics' or 'distance metrics' often are built directly on dependence measures, they typically do not share many of their characteristics (e.g. their spaces typically are not positive definite (see Alpay & Mayats-Alpay, 2023; and Meckes, 2013)). In finance they often are used non-inferentially and mechanistically in hierarchical portfolio construction models (see Tumminello et al., 2005; and Dom et al., 2024) where they have received decidedly mixed reviews (see Trucíos Maza, 2025; Aznar, 2023; Cota, 2019; and Ciciretti & Pallotta, 2023), especially under correlation breakdowns (see Marti et al., 2021). As they stand, they are not designed to answer the inferential questions posed herein. In fact, I show in later sections how NAbC provides a generalized entropy that has many useful advantages over an entire class of metrics most commonly used in this setting, called 'norms.'

<sup>&</sup>lt;sup>10</sup> It is worth reiterating here that "positive definite" throughout this monograph refers to the dependence measure calculated on the matrix of all pairwise associations in the portfolio, that is, calculated on a bivariate basis. While some of the dependence measures addressed in this monograph (e.g. Szekely's correlation, as well as some variants of Chatterjee's (see Pascual-Marqui et al., 2024)), can be applied on a multivariate basis, sometimes in arbitrary dimensions, the term "positive definite" in this monograph is not applied in this sense (see for example Cardin, 2009). For surveys of related multivariate methods, see Chatterjee (2024) and Han (2021), in addition to Grothe et al. (2014), Latif and Morettin (2014), Reddi et al. (2015), Li and Joe (2024), Yu et al. (2021) and Puccetti (2022) for some approaches not covered herein.

group of measures, not a single matrix, of the many millions simulated, was ever found to be non-positive definite, making the distinction between these two groups arguably moot, at least for the empirical testing performed. However, until positive definiteness is proven analytically for a dependence measure, responsible analysis requires that this always is verified empirically.<sup>11</sup>

#### 2.a.i. Monotonic Measures

The oldest and most widely used and known dependence measure is Pearson's product moment correlation (see Pearson, 1895), which is what is usually referenced when the word "correlation" alone is mentioned. Taking two variables, say, the financial returns of two assets X and Y, Pearson's measures how often and to what degree they deviate from their respective sample means in the same or in opposite directions, as shown in (1) below.<sup>12</sup>

(1) 
$$r_{X,Y} = \frac{\sum_{i=1}^{n} \left( X_{i} - \frac{1}{n} \sum_{j=1}^{n} X_{j} \right) \left( Y_{i} - \frac{1}{n} \sum_{j=1}^{n} Y_{j} \right) / (n-1)}{\sqrt{\sum_{i=1}^{n} \left( X_{i} - \frac{1}{n} \sum_{j=1}^{n} X_{j} \right)^{2} / (n-1)} \sqrt{\frac{\sum_{i=1}^{n} \left( Y_{i} - \frac{1}{n} \sum_{j=1}^{n} Y_{j} \right)^{2} / (n-1)}} = \frac{C \hat{o} v(X,Y)}{s_{X} s_{Y}}$$

The numerator is the (sample) covariance of X and Y, and the denominator – the product of the (sample) standard deviations of X and Y – has the effect of scaling the (sample) covariance to a (maximum) range of -1 to 1.<sup>13</sup> So Pearson's is just the scaled covariance between X and Y.

Another of the most commonly used dependence measures is Spearman's Rho (see Spearman, 1904), which is exactly the same formula as Pearson's but instead of using the values of the returns of X and Y, their ranks are used instead:

(2a) 
$$sr_{X,Y} = \frac{\sum_{i=1}^{n} \left( R_{X_i} - \frac{1}{n} \sum_{j=1}^{n} R_{X_j} \right) \left( R_{Y_i} - \frac{1}{n} \sum_{j=1}^{n} R_{Y_j} \right) / (n-1)}{\sqrt{\sum_{i=1}^{n} \left( R_{X_i} - \frac{1}{n} \sum_{j=1}^{n} R_{X_j} \right)^2 / (n-1)} \sqrt{\frac{\sum_{i=1}^{n} \left( R_{Y_i} - \frac{1}{n} \sum_{j=1}^{n} R_{Y_j} \right)^2 / (n-1)}}$$

<sup>13</sup> Note that this range can be tighter under specific circumstances, such as for equicorrelation matrices where

$$\left[-1/(p-1)\right] \le r \le 1, \ p = \dim(r).$$

JD Opdyke, Chief Analytics Officer

<sup>&</sup>lt;sup>11</sup> This empirical verification remains advisable even for the first group of dependence measures, for which positive definiteness has been proven analytically, as numerical issues always can arise in cases of matrices that approach singularity (i.e. those with values that, if changed just slightly, would not be positive definite).

<sup>&</sup>lt;sup>12</sup> Importantly, all formulae of estimators herein, unless otherwise noted, refer to those based on sample data, where "*n*" indicates the number of observations in the sample, as opposed to those based on an entire population of data.

If there are no ties in the data, (2a) can be shortened to

(2b) 
$$sr_{X,Y} = 1 - \frac{6\sum_{i=1}^{n} \left(R_{X_i} - R_{Y_i}\right)^2}{n^3 - n}$$
 (see Zar, 1999)

Using ranks can make Spearman's less sensitive than Pearson's to extreme data values under some data conditions, just like another rank-based dependence measure, Kendall's Tau.

Also called a measure of concordance, Kendall's Tau (see Kendall, 1938) is the sum of all pairwise comparisons of every data point of X and Y, divided by the total number of pairs.<sup>14</sup> The pairwise comparisons are given values of 1, 0, or -1, respectively, if both from one period to another are in increasing/decreasing order, if the values from both periods are tied for either of the assets, or if the assets are NOT both in increasing/decreasing order; it thus gives the number of pairs in concordance minus the number in discordance relative to the total number of pairs, as shown below.

(3a) 
$$\tau(X,Y) = \frac{\#\text{concordant pairs} - \#\text{discordant pairs}}{\text{total } \#\text{ pairs}} = \frac{2}{n(n-1)} \sum_{i=1}^{n-1} \sum_{j=i+1}^{n} \text{sgn}(x_i - x_j) \text{sgn}(y_i - y_j)$$

where  $\operatorname{sgn}(z) = 1$  if z > 0,  $\operatorname{sgn}(z) = -1$  if z < 0,  $\operatorname{sgn}(z) = 0$  if z = 0, for both N and n

However, ties in the values of either of the pairs,  $(x_i \text{ and } x_j)$  or  $(y_i \text{ and } y_j)$ , will restrict the range from achieving -1 or +1, even under otherwise perfect discordance or concordance, respectively, so a commonly used variant of Kendall's Tau that avoids this drawback when ties exist is:

(3b) 
$$\tau_b(X,Y) = \frac{1}{\sqrt{(n_0 - n_1)(n_0 - n_2)}} \left[ \sum_{i=1}^{n-1} \sum_{j=i+1}^n \operatorname{sgn}(x_i - x_j) \operatorname{sgn}(y_i - y_j) \right]$$

where 
$$n_0 = n(n-1)/2$$
;  $n_1 = \sum_{i=1}^{i \text{ grps ties}} t_i(t_i-1)/2$ ;  $n_2 = \sum_{j=1}^{j \text{ grps ties}} u_j(u_j-1)/2$ ;

 $t_i = \#$  ties in i-th group of x;  $u_i = \#$  ties in j-th group of y

The "big 3" dependence measures – Pearson's, Kendall's, and Spearman's – are by far the most widely used in practice.<sup>15</sup> Although widely held myths persist regarding Pearson's as a measure strictly of *linear* monotonic relationships (see van den Heuvel & Zhan, 2022), all three measure monotonic association (i.e. the direction of the association, positive or negative, does not change within the covered time period) that is symmetric, or non-directional in the variable order (i.e. the measured dependence of X on Y is

Page **15** of **92** 

<sup>&</sup>lt;sup>14</sup> Note that the pairs are based on time-ordered data, but most of the periods in each pair are not contiguous.

<sup>&</sup>lt;sup>15</sup> Other long-established measures include Hoeffding's D (see Hoeffding, 1948), Blomqvist's coefficient (see Blomqvist, 1950), and Gini's gamma (see Gini 1914; and Genest et al., 2010).

assumed to be the same as that of Y on X). It is important to recall here that as measures of monotonic dependence, values of zero generally do not necessarily imply independence between X and Y,<sup>16</sup> but independence between X and Y does imply values of zero for the big 3.<sup>17</sup> Some of the dependence measures treated below avoid this limitation under many conditions.

The properties of the big 3 have been studied extensively in the literature, but real gaps remain. Our interest in this monograph lies not just in a single bivariate relationship between X and Y, but rather, in all pairwise relationships of all assets in a portfolio, simultaneously: X may be strongly, positively associated with Y, which also may be positively associated with Z, which also may be negatively associated with A, B, and C, while B and C may be modestly but negatively associated with X again! So we have a matrix of dependence measure values with rows and columns identifying the pairwise relationships between all the asset pairs, as shown in (4) for assets 1, 2, 3, and 4. As above, I refer to this matrix herein as the all-pairwise matrix.

(4) 
$$R = \begin{bmatrix} 1 & r_{1,2} & r_{1,3} & r_{1,4} \\ r_{2,1} & 1 & r_{2,3} & r_{2,4} \\ r_{3,1} & r_{3,2} & 1 & r_{3,4} \\ r_{4,1} & r_{4,2} & r_{4,3} & 1 \end{bmatrix}$$

Some of the characteristics of this matrix, for all of the big 3 and many of the other measures presented below, include:

- i. Symmetry:  $r_{i,j} = r_{j,i}$
- ii. Unit diagonal entries:  $r_{i=j} = 1$

iii. Bounded non-diagonal entries, with maximum range of:  $-1 \le r_{i,j} \le 1$ 

iv. The matrix is positive definite, i.e. all eigenvalues  $\lambda_i > 0$ 

For completeness, and for reference throughout this monograph, I define eigenvalues for  $R \in \mathbb{R}^{p \times p}$  here:<sup>18</sup>

If there exists a nonzero vector v such that  $Rv = \lambda v$  then  $\lambda$  is an eigenvalue of R and v is its corresponding eigenvector.  $\lambda$  and v can be obtained by solving

 $det(\lambda I - R) = 0$ , then  $det(\lambda I - R)v = 0$ , where *I* is the identity matrix and det is the determinant. The

<sup>&</sup>lt;sup>16</sup> However, an exception occurs when data is distributed as bivariate normal, in which a Pearson's value of zero *does* indicate independence.

<sup>&</sup>lt;sup>17</sup> This is easy to visualize with a non-monotonic relationship like  $y=x^2+\mathcal{E}$ , which on average will yield big 3 values close to zero. But the relationship is non-linear and u-shaped, which most certainly is not one of independence.

<sup>&</sup>lt;sup>18</sup> Financial returns are real numbers, and so this definition holds for all relevant dependence measures in this setting.

eigenvalue can be thought of as the magnitude of the (portfolio) variability in the direction of the eigenvector. With actual, real-world financial data (i.e. values that are not imaginary or complex), this variability can never be negative,<sup>19</sup> so computational numeric issues aside,<sup>20</sup> proper measures of dependence should be positive definite,<sup>21</sup> either via analytical proof, or in the absence of such, then via empirical verification of all relevant cases (although empirical verification is advisable even in the first instance).

The main point here is that we need to understand the characteristics of the estimators we use to estimate the values of the all-pairwise matrix, based on our sample of financial returns data. This is the only way we can define the finite-sample distribution of the estimator, which is the only way we will be able to make *inferences* about the true population values of these estimates.

## 2.a.ii. Tail Dependence Measures

Another important and time-tested dependence measure, especially for risk analyses, is the tail dependence matrix (TDM). Conceptually, TDM measures the probability of a variable value residing in the tail of one variable's distribution given that the value of the other variable (asset return) resides in the tail of its distribution. More precisely, TDM provides the probability of a variable exceeding a quantile of its distribution conditional on the other variable in the pair exceeding the same quantile of its distribution. Hence, the tail dependence matrix consists of conditional probabilities of quantile exceedance, so each value can range from zero to one, rather than -1 to 1 like the "big 3." But otherwise the matrix conditions listed in (4) above all hold (its positive definiteness was proven by Embrechts et al., 2016). The upper tail dependence matrix only is equal to the lower tail dependence matrix if data distributions are perfectly symmetric: otherwise, the two metrics have distinct values, as shown below in (5) and (6):

(5) 
$$TDMU_{X,Y} = \lim_{q \to 1^{-}} P(Y > F_Y^{-1}(q) | X > F_X^{-1}(q))$$

(6) 
$$TDML_{X,Y} = \lim_{q \to 0^+} P(Y \le F_Y^{-1}(q) | X \le F_X^{-1}(q))$$

where quantile function = inverse cdf =  $F^{-1}(q) = \inf \{x \in \mathbb{R} : F(x) \ge q\}$ 

<sup>&</sup>lt;sup>19</sup> This can be seen most easily when the covariance (or equivalently, Pearson's correlation) is the dependence measure used: the covariance is the (expected value of a) sum of squared <u>real</u> numbers (as no imaginary or complex values are observed in financial returns). Because a squared, real number (other than zero) is always greater than zero, the sum of such numbers can never be negative.

<sup>&</sup>lt;sup>20</sup> Numerical calculations based on positive definite matrices can sometimes render slightly negative estimates of specific eigenvalues, but as shown in later sections herein, NAbC is designed specifically to be more robust to such numerical errors than the more common approaches related to eigen decompositions in the extant literature.

<sup>&</sup>lt;sup>21</sup> If any  $\lambda$  = 0, and none are negative, the matrix is said to be positive semi-definite, although herein this is treated as a textbook border case as returns would have to exhibit perfect linear dependence for an eigenvalue to be exactly zero.

Straightforward empirical estimators for (5) and (6) are presented in Garcin and Nicolas (2023) as below:

$$(5a) \hat{\lambda}_{U_{X,Y}}(i/n) = \frac{1-2(i/n) + \hat{C}_n(i/n, i/n)}{1-(i/n)}$$

$$(6a) \hat{\lambda}_{L_{X,Y}}(i/n) = \frac{\hat{C}_n(i/n, i/n)}{(i/n)}$$

$$where (i/n) = q \text{ for } q < 0.5 \text{ and } (i/n) = (1-q) \text{ for } q \ge 0.5,$$

$$\hat{C}_n(u,v) = \frac{1}{n} \sum_{j=1}^n \mathbf{1} \{ \hat{F}_{X,n}(X_j) \le u \} \mathbf{1} \{ \hat{F}_{Y,n}(Y_j) \le v \} \text{ is the empirical copula, and}$$

$$\hat{F}_{X,n}(x) = \frac{1}{n} \sum_{j=1}^n \mathbf{1} \{ X_j \le x \} \text{ and } \hat{F}_{Y,n}(y) = \frac{1}{n} \sum_{j=1}^n \mathbf{1} \{ Y_j \le y \} \text{ are the empirical cumulative distribution functions.}$$

These are shown in Schmidt and Stadtmüller (2006) to have good statistical properties (i.e. strong consistency and asymptotic normality). Many other measures of tail dependence exist (see AghaKouchak et al., 2013, Babić et al., 2023, Manistre, 2008, Li and Joe, 2024, Krupskii and Joe, 2014, Lauria et al., 2021, and intriguingly, Siburg et al., 2024), but (5) and (6) are the oldest, most widely used, and best understood. Tail dependence is especially important in the risk analytics of financial portfolios because "tail events" often represent the most material financial impacts, are typically associated with non-linear effects and associations, and are closely tied to correlation breakdowns: as is commonly and rightly stated, "when things go bad they go bad together." The phenomenon of "correlation breakdowns" is treated in more detail later in this monograph, but note that the tail dependence matrix has been one of the principal tools used in both the literature and by practitioners to quantitatively estimate it and mitigate its effects.

## 2.a.iii. Distance-Based and Other New Measures

The design of Szekely's distance correlation (Szekely et al., 2007) seeks to better handle dependence that is both non-linear and non-monotonic. It uses two matrices: the matrix of pairwise distances between all X values in the sample, and the same matrix calculated from all Y values. To the extent that these matrices vary together, Szekely's distance correlation will approach a value of 1, and to the extent they do not, it will approach a value of zero. So its range is zero to one and a value of zero, unlike the "big 3," does indicate independence between X and Y. Also unlike the "big 3," its value does not indicate with a positive or negative sign whether dependence between X and Y is positive or negative. Notably, the distance correlation can be calculated in arbitrary – and different – dimensions, so the sample from X can be drawn, for example, from a three dimensional distribution, and the sample from Y can be drawn from a six dimensional distribution.<sup>22</sup>

 <sup>&</sup>lt;sup>22</sup> Note that the primary focus of the development of many multivariate dependence measures is on testing the null hypothesis of multivariate independence, and thus, on the level and power of this specific test for these measures. While this objective is
 JD Opdyke, Chief Analytics Officer Page 18 of 92 Correlation and Beyond

(7) first, create n x n distance matrices a and b by letting

$$a_{i,j} = ||x_i - x_j||$$
 and  $b_{i,j} = ||y_i - y_j||$ ,  $i, j = 1, 2, 3, ..., n$  where  $||vector z_n|| = \sqrt{z_1^2 + z_2^2 + \dots + z_n^2}$ 

Next, subtract from *a* and *b* their row and column means, and add their respective matrix means, as shown below:

$$A_{i,j} = a_{i,j} - a_{*,j} - a_{i,*} + a_{*,*} \text{ and } B_{i,j} = b_{i,j} - b_{*,j} - b_{i,*} + b_{*,*}$$
  
Then Szekly's distance correlation =  $dcorr = \sqrt{\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n A_{i,j} B_{i,j}} / \sqrt{\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n A_{i,j}^2 \cdot \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n B_{i,j}^2}$ 

Another recent dependence measure – Lancaster's correlation (see Holzmann and Klar, 2024) – shares several characteristics with Szekely's: its values range from zero to one, a value of zero indicates independence, and it does not indicate with a positive or negative sign whether the dependence between X and Y is positive or negative. Lancaster's correlation was designed not only to handle non-linear and non-monotonic dependence, but also to improve upon, via increased robustness and generalizability and ease of computation, another dependence measure, the maximal correlation (see Hirschfeld (1935) and Gebelein (1941)).

(8) 
$$lan = \max\left(\left|r\left(\tilde{X}, \tilde{Y}\right)\right|, \left|r\left(\tilde{X}^{2}, \tilde{Y}^{2}\right)\right|\right)$$
 where  $\tilde{X} = \Phi^{-1}\left(F_{X}\left(X\right)\right)$  and  $\tilde{Y} = \Phi^{-1}\left(F_{Y}\left(Y\right)\right)$ , where *r* is Pearson's correlation,  $\left| \right|$  is the absolute value function,  $\Phi^{-1}$  is the quantile (inverse cdf) function of the standard normal distribution, and *F* is the (empirical) cdf of each variable.

A second version is called linear Lancaster's correlation:

$$lanL = \max\left(\left|r(X,Y)\right|, \left|r\left(\bar{X}^2, \bar{Y}^2\right)\right|\right) \text{ where } \bar{X} = \left(X - \bar{X}\right) / \sqrt{\sum_{i=1}^n \left(X_i - \bar{X}\right)^2 / (n-1)} \text{ and } \bar{Y} = \left(Y - \bar{Y}\right) / \sqrt{\sum_{i=1}^n \left(Y_i - \bar{Y}\right)^2 / (n-1)} \text{ and } \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \text{ and } \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

Holzmann and Klar (2024) conduct empirical analyses comparing Szekely's distance correlation and both Lancaster's correlations under a wide range of data conditions. They also compare these to another new, but directional dependence measure, called Chatterjee's correlation coefficient.

JD Opdyke, Chief Analytics Officer

foundational, that of this monograph is on dependence as measured using bivariate associations. Consequently, I use the framework of the all-pairwise matrix to focus on dependence measures in the literature with strong results related to their statistical power, level control, ease of implementation, low computational complexity, and attainment of the full range of values they are meant to attain under the relevant sample spaces (these are the measures I select to review in this section, and to which I have applied NAbC). More importantly, relying on the bivariate, as opposed to multivariate, relationships measured in the all-pairwise matrix is critical to the scenario flexibility provided by NAbC, as explained in later sections.

#### 2.a.iv. Asymmetric, Directional Measures

The concept of asymmetric, directional dependence is not new. Recent research on such measures goes back over a dozen years (see Zheng et al., 2012), but has its direct origins in work done at the end of the nineteenth century (see Yule, 1897, and Allena and McAleerb, 2018). In later sections I will go over examples of how these measures are being used effectively in causal frameworks, but only present them in this section. A recent example, Chatterjee's correlation coefficient, has garnered much attention upon its publication in 2021. This is largely due to its simplicity and ease of implementation as a measure of non-linear, non-monotonic, regression-based, and cyclical dependence. If X and Y pairs are ranked according to X values, with no ties on the X values, so that  $((X_{(1)}, Y_{(1)}), (X_{(2)}, Y_{(2)}), \dots, (X_{(n)}, Y_{(n)}))$  then:

(10) 
$$chcorr = \xi_n(X,Y) := 1 - \frac{3\sum_{i=1}^{n-1} |r_{i+1} - r_i|}{n^2 - 1}$$
 where  $r_i = \text{rank of } Y_i$ 

Under ties for some of the X values, break ties uniformly at random, and

(11) 
$$chcorr = \xi_n(X,Y) := 1 - \frac{n \sum_{i=1}^{n-1} |r_{i+1} - r_i|}{2 \sum_{i=1}^{n} l_i (n - l_i)}$$
 where  $l_i = \# j$  such that  $Y_{(j)} \ge Y_{(i)}$ 

Unlike the big 3, Chatterjee's new correlation coefficient ranges from zero to one asymptotically (it can exceed these bounds slightly under finite samples), and a value of zero does indicate independence. Also, no positive or negative dependence is indicated by a positive or negative sign on the measure value. Most notably, this is an asymmetric dependence measure, that is, the order of X and Y matters:

 $\xi_n(X,Y)$  does not necessarily equal  $\xi_n(Y,X)$ . In other words, the dependence of Y on X is not assumed to be identical to the dependence of X on Y, respectively: dependence is *directional*.<sup>23</sup> However, note that Chatterjee's can be made to be symmetric by simply taking the maximum of two measures, one in each direction as in (12):

(12) chcorr\_sym = max  $\left[ \xi_n(X,Y), \xi_n(Y,X) \right]$ 

Chatterjee's breakthrough has spawned many variants (see Lin & Han, 2023, Pascual-Marqui et al., 2024, and especially Gao and Li, 2024). One of these is the "improved Chatterjee's correlation" derived by Xia

<sup>&</sup>lt;sup>23</sup> It is important to note that herein, when using dependence measures that are asymmetric/directional, the corresponding allpairwise matrix remains symmetric. So when using, say, Chatterjee (2021), on the returns of two particular assets in the portfolio, say, X3 and X4, the value in cell row 3, column 4 of the matrix is  $\xi_n(X3, X4)$ , and the value in cell row 4, column 3 of the matrix is identical, that is,  $\xi_n(X3, X4)$ ; it is NOT  $\xi_n(X4, X3)$ .

et al. (2024), the motivation of which is to increase power by using inverse distance weightings of all neighboring data values as opposed to just one.

(13) 
$$ichcorr = \xi_n^{IM}(X,Y) = 1 - \frac{\sum_{i\neq j}^n |r_i - r_j|/|i - j|}{\frac{n+1}{3}\sum_{i\neq j}^n |i - j|}$$

Xia et a. (2024) test the power and level of "improved Chatterjee" against both Chatterjee and Szekely's correlations in an empirical study under widely varying data conditions. Both Chatterjee's and improved Chatterjee's coefficients exhibit power under non-monotonic, non-linear, and cyclical dependence, with the latter usually winning.<sup>24</sup>

Interestingly, Zhang (2024a) has proposed combining Chatterjee's and Spearman's in an effort to obtain the best of both worlds: a dependence measure that has reasonable power under cases of nonmonotonic, non-linear, and/or cyclical dependence (where Spearman's has little to no power, especially compared to Chatterjee's) as well as reasonable power under monotonic dependence (where Chatterjee's has less power than Spearman's).

(14) 
$$zcorrsp = I_{n_{sp}}(X,Y) = \max\left\{ |sr_{X,Y}|, \sqrt{5/2}\xi_n(X,Y) \right\}$$

Zhang's (2024a) combined correlation ranges from 0 to 1, where zero indicates independence. This dependence measure also is asymmetric due to its inclusion of Chatterjee's coefficient. Zhang (2024b) later derived the symmetric version of this test as (14a):

(14a) 
$$zcorrsp\_sym = \max\left\{\left|sr_{X,Y}\right|, \sqrt{5/2}\xi_n(X,Y), \sqrt{5/2}\xi_n(Y,X)\right\}$$

Another measure similar to Chatterjee's is the Differential Distance Correlation (DDC) of Liu and Shang (2025). DDC's values range from 0 to 1, with 0 indicating independence. Notably, like Szekely's distance

<sup>&</sup>lt;sup>24</sup> Note that both Chatterjee (2021) and Xia et a. (2024) test against two dependence measures not explored further herein: the HSIC measure of Gretton et al. (2007), and the HHG measure of Heller et al. (2013). Both appear to have excellent power under circular and heteroskedastic data, and the former maintains reasonably large power under other conditions where both Chatterjee statistics outperform it. While both HSIC and HHG are much more computationally intensive than either Chatterjee dependence measure (as is Zheng et al., 2012, as well), Sejdinovic et al. (2013) intriguingly prove that "reproducing kernel Hilbert space (RKHS)-based dependence measures are precisely the formal extensions of the [Szekely et al. (2007)] distance covariance." So HSIC is a generalized version of Szekely et al. (2007) that circumvents "the problem of nonintegrability of weight functions by using translation-invariant kernels called distance-induced kernels." RKHS-based dependence measures and intriguing area of continuing research (see Ke, 2019; Mitchell et al., 2022; Wahba, 2017; and Zhang & Songshan, 2023), especially because they all are positive definite by design and definition (see Tripathi et al., 2022, for an example), and thus, are suitable for NAbC. A related and similarly intriguing approach is that of Pascual-Marqui et al. (2024) who combine Szkely's and Chatterjee's measures in directional regressions within a causal modeling framework. This is discussed further below, alongside NAbC's potential use within causal frameworks.

correlation, DDC can be multidimensional, but when X is univariate so that DDC can be used in an allpairwise matrix, it is defined as (14b) below:

(14b) 
$$DDC_n(X | Y) = 1 - \frac{1}{(n-1)} \sum_{i=1}^{n-1} ||X_{(i)} - X_{(i+1)}|| / \left[ \binom{n}{2}^{-1} \sum_{i=1}^n (2i - n - 1) X^{(i)} \right]$$

where  $\left\{ \left( X_{(i)}, Y_{(i)} \right) \right\}_{i=1}^{n}$  are ordered to satisfy  $Y_{(i)} \leq \cdots \leq Y_{(n)}$ , and

 $X^{(i)}$  are ordered to satisfy  $X_{(i)} \leq \cdots \leq X_{(n)}$ . Liu and Shang (2025) show in an empirical study that DDC has power similar to that of Chatterjee's measure, but with slightly more power under damped oscillator data. And like Chatterjee's measure, DDC is directional, so in the general case,

 $DDC_n(X|Y) \neq DDC_n(Y|X)$ , although also like Chatterjee's measure, a symmetric, non-directional version may be obtained via (14c)<sup>25</sup>:

(14c) 
$$DDC_{n-sym}(X,Y) = \max\left[DDC_n(X|Y), DDC_n(Y|X)\right]$$

Finally, asymmetric, directional dependence measures also can be applied only to the tails of X and Y, and it is important to note that correlation breakdowns often are associated specifically with (asymmetric) tail dependence: "Extensive evidence has been gathered showcasing the prevalence of heavy-tailed distributions and asymmetric tail interdependence within equity and foreign exchange markets, particularly during times of crisis. ...This phenomenon causes markets that typically exhibit minimal or no correlation to behave similarly, often in opposition to fundamental principles." (Pramanik, 2024).<sup>26</sup> One straightforward example of an asymmetric tail dependence measure is that of Deidda et al, (2023) which is essentially Kendall's Tau applied conditionally, only when the percentile, q, of X (or Y) is exceeded:

(15) 
$$\hat{\tau}_{X,Y}\left(q\right) = \frac{1}{\binom{k}{2}} \sum_{1 \le i \le j \le n} \operatorname{sgn}\left(X_i - X_j\right) \operatorname{sgn}\left(Y_i - Y_j\right) I\left(X_i, X_j > X_{(n-k)}\right)$$

where q = 1 - k/n, and  $k \le n$  is the number of exceedences used in the tail, and I() is the indicator function (one when true, zero otherwise) ensuring that only the k largest observations of X are used. Note again that generally,  $\hat{\tau}_{X,Y}(q) \ne \hat{\tau}_{Y,X}(q)$ , that is, this tail dependence measure is directional, and the effect of X's tail on Y's tail is not assumed to be the same as that of Y's tail on X's tail.

<sup>&</sup>lt;sup>25</sup> Based on email correspondence with author Yixiao Liu, July 9, 2025.

<sup>&</sup>lt;sup>26</sup> See also Ito and Yoshiba (2025): "We provide new evidence that lower tail dependence coefficients increased compared to upper ones for all pairs in the COVID-19 crash..."

Other directional, asymmetric dependence measures include the dynamic asymmetric tail dependence measure of Ito and Yoshiba (2025), the QAD measure of Junker et al. (2021), the generalized correlation of Zheng et al. (2012) and the measures of Vinod (2022), and others described in Jondeau (2016).<sup>27</sup>

It remains notable that NAbC's broad scope allows for its application to these asymmetric, directional dependence measures as readily as it is applied to the big 3. As seen in a later section, this gives NAbC great utility in some surprising settings. Even as it is designed fundamentally as a method for robust statistical inference, and is thus associational, when used on these directional dependence measures NAbC can be applied to increase the power of causal models to accurately recover directed acyclic graphs (DAGs). It also can be used effectively to robustify other causal frameworks. These are areas of continuing research, but they serve as examples of how NAbC's breadth of application can be useful even beyond ceteris paribus comparisons of the inferential power of competing dependence measures. Yet this remains invaluable as it stands, as such comparisons often would not be possible without NAbC. All dependence measures have strengths and weaknesses, not only under different data conditions, but also depending on the specific questions applied researchers and practitioners need to answer. So we need to be able to test them using the same unifying method under controlled conditions if we are to determine which is most appropriate for a given situation.

With this brief but important review of dependence measures aside, I address their estimation in the next section before turning to the derivation of NAbC in subsequent sections.

## 3. Estimation

## 3.a. Covariance and Pearson's Correlation

Regarding estimation of the all-pairwise matrix, the lion's share of the literature focuses on estimators for the covariance matrix and Pearson's correlation matrix. This is not terribly surprising given the relatively long history and widespread usage of Markowitz's portfolio framework (see Markowitz, 1952) and related models.

"Accurate covariance matrix prediction is crucial for portfolio optimization and risk management because it captures the relationships and co-movements between asset returns." (Lee et al., 2024)

But fortunately, some see the bigger picture, that these analyses can and should be broadened to ALL positive definite dependence measures:

<sup>&</sup>lt;sup>27</sup> Note that under certain conditions, such as when categorical and ordinal data are being analyzed and the number of categories between the two variables differs dramatically, even Pearson's correlation can be unambiguously directional (see Metsämuuronen, 2022, for details).

"Modeling covariance matrices – or more broadly, positive definite (PD) matrices – is one of the most fundamental problems in statistics" (Lan et al., 2020).

So I first focus below on estimation of Pearson's matrix, and then on the other dependence measures discussed in the previous section.

The first of the two major challenges of estimating the all-pairwise matrix of any dependence measure is sample size, because we are not just estimating a single parameter, say, a volatility or a beta of a single

asset, but rather,  $\binom{p}{2} = \frac{p(p-1)}{2}$  pairwise associations. To do this accurately and with reasonable precision, we need more data than is needed for a single estimate. Regarding accuracy, the sample covariance matrix, and thus, the sample Pearson's matrix are consistent estimators, that is, they are asymptotically unbiased. But regarding precision, their estimates will be way too variable to be useable or useful, not to mention biased in the finite-sample case, in the absence of large(r) data samples. A widely recognized rule of thumb is that the sample size needs to be at least ten times the dimension of the matrix (N $\geq$ 10p; see Bongiorno et al., 2023), but this arguably depends on the method used to estimate it. For example, Bun et al. (2016) devise a rotationally invariant estimator that "cleans" or "denoises" the estimate of Pearson's matrix<sup>28</sup> using functions of its eigen values, a method for which they argue that N $\geq$ 2p is sufficient. Note that the estimators of all the dependence measures presented herein are at least asymptotically unbiased, and that some researchers believe the sample size issue has been addressed as well as it can be, especially if the best methods are being used (see Bouchaud, 2021: "Now the data problem is solved as best as possible..." referring to Bun et al., 2023, among others). With this in mind, currently it would appear that Bun et al. (2016) is the state-of-the-art estimator for the unconditional estimate of Pearson's matrix (see du Plessis & van Rensburg (2020) for a comparison study). However, in most cases in this setting, obtaining accurate and realistic estimates of dependence structure requires conditioning on time, because in reality, time series of financial returns do change over time, and so does their dependence structure. And this is the second major challenge when estimating any all-pairwise dependence matrix: non-stationarity.

As Bouchard (2021) rightly points out, portfolio frameworks like that of Markowitz (1952), and really any in applied usage, require knowledge of the dependence measure to be representative of the future realized correlations, because financial data is non-stationary (i.e. its distribution changes over time; the term 'conditional distribution' refers to the distribution *conditional* on a specific time period). Therefore, we need a forecast, into the near-term future, of the *conditional* dependence matrix. And a very compelling one for Pearson's and the covariance matrix, dubbed "Average Oracle" (AO), is exactly what is provided by Bongiorno et al. (2023) (see also Bongiorno & Challet, 2023a, for an extensive empirical study against

<sup>&</sup>lt;sup>28</sup> See also Palomar (2025), 3.5.3 (pp. 52-55), for several robust estimators of the covariance (Pearson's correlation) matrix that can be very useful when the marginal returns distributions are decidedly non-Gaussian, which is the rule rather than the exception for most financial markets. For additional recent, innovative robust estimators of the covariance (Pearson's correlation) matrix, see Centofanti et al. (2025), Besson (2025), Carrara et al. (2025), Casa and Cappozzo (2025), and Sun and Huang (2025).

competitors). Conceptually AO is very straightforward: based on the eigen decomposition of Pearson's matrix, it uses the (oracle) covariance of the next-period 'future' with the eigenvectors of the adjacent past period to obtain eigenvalues that, when averaged over many samples, embed the desired, dynamic time effects for a robust forecast of Pearson's matrix. Somewhat surprisingly, this intuitive method outperforms all flavors of advanced "shrinkage," both non-linear (see Ledoit & Wolf, 2017) and quadratic (see Ledoit & Wolf, 2022a, 2022b) as well as DCC and NLS combinations (see Engle et al., 2019). It is fast, straightforward to understand and implement, and importantly, fully nonparametric. AO's outperformance of the widely used NLS approach perhaps should not be so surprising given that Bongiorno & Challet (2023b) recently proved that NLS is not optimal for portfolio optimization, as was widely believed, because it does not optimize under non-stationarity.<sup>29</sup> So I recommend AO as the current state-of-the-art conditional estimator of Pearson's matrix. However, this literature is vast, comprising easily many hundreds of papers when both covariance and Pearson's correlation estimation are included, and given the current rapid pace of research in this area, it is certainly possible that new, worthy competitors exist, especially under specific data conditions (see for example Zhang et al. (2022), Zhang et al. (2023), Vanni et al. (2024), and Zhangshuang et al. (2025)).

#### 3.b. Other Dependence Measures

For estimation, conditional or unconditional, of the all-pairwise matrices based on the other dependence measures listed in the previous section, the literature has little to offer beyond the fact that all of the sample estimators presented in the previous section are at least asymptotically unbiased.<sup>30</sup> So as long as sample sizes are sufficiently large these estimators will retain good statistical properties. However, I offer two additional suggestions below for possible improvements under specific conditions. The first is simply the inverse of a common robustification technique using a well established relationship between Pearson's and the rank-based measures, Kendall's and Spearman's. Estimates of (bivariate) Kendall's Tau or Spearman's Rho often are used to robustify those of Pearson's using the widely known relationships of  $r = \sin(\tau \pi/2)$  and  $r = 2\sin((sr)\pi/6)$ , respectively, which are valid under iid elliptical data distributions (see Sheppard, 1899; Greiner, R., 1909; Lindskog et al., 2003; Heinen & Valdesogo, 2022; McNiel et al., 2005; and Hansen & Luo, 2024; and for advanced methods on this, see Barber & Kolar, 2018, and Niu et al., 2020). Yet under specific, known data conditions that are non-elliptical, it may be demonstrable that these transformations remain reasonably accurate (see Hamed (2011) and Hansen & Luo (2024) for examples). In this case, given a strong estimator of Pearson's from an improved estimation method like those described above (e.g. Average Oracle of Bongiorno et al. (2023)), the inverses of these functions could be used to obtain estimates for the all-pairwise matrices of Kendall's and Spearman's that likely would share some of the benefits of an improved estimate for Pearson's matrix, especially

<sup>30</sup> However, see Zhao et al. (2014) for an intriguing exception.JD Opdyke, Chief Analytics Officer Page 25 of 92

<sup>&</sup>lt;sup>29</sup> For those that view shrinkage favorably in general, an improved shrinkage competitor with arguably better properties than NLS is that of Kelly et al. (2024).

when it is conditional. Note, however, that these transformations would require verifying, and sometimes enforcing, positive definiteness ex post (see McNeil et al. (2015) and Higham (2002)). Also, they apply only to Kendall's tau and Spearman's rho.

An approach with a much broader range of application would be the implementation of Average Oracle on ANY of the above-mentioned dependence measures directly, simply replacing Pearson's matrix with the measure of choice, but keeping the methodology otherwise identical. Because their all-pairwise matrices will be verified to be positive definite, any of these dependence measures will have valid eigen decompositions wherever the covariance matrix does (and even in some cases where the covariance matrix is singular): the 'training' eigenvectors of the adjacent past can be used in exactly the same manner with next-period 'future' all-pairwise dependence matrices to obtain (averaged) eigenvalues that embed the measured, average, empirical time dependency. This approach would presumably share the benefits of AO's conditional correlation estimates (i.e. its robustness and fully nonparametric nature) to other measures of dependence in this setting.

The above discussion regarding estimation should clarify and reemphasize the fact that NAbC is not an estimator of any of these dependence measures: rather, it provides the finite-sample distribution of the estimator, as long as its estimate is positive definite, for any dependence measure under any real-world data conditions. This allows us to make actionable inferences about dependence structure in a unified way, allowing for comparative, ceteris paribus analyses. The literature to date provides such distributions in a highly piecemeal fashion for some of the dependence measures under some (often very limited and/or unrealistic) data conditions for some (often very limited and/or unrealistic) ranges of values. These derivations also can be extremely complex and unwieldy, making them unusable for all intents and purposes for many practitioners. NAbC sidesteps all of these problems with a single, unified, and straightforward method that allows for ceteris paribus comparisons of different measures, or different estimators of a particular measure. Estimation of the all-pairwise matrix is the only thing out-of-scope for NAbC, but in a sense this is a strength of the approach since it permits NAbC to remain "estimator agnostic," allowing for its application on any reasonable and relevant estimator of the all-pairwise matrix. It thus provides the flexibility to use those that are most robust and/or most precise and/or most accurate - or any combination thereof - under different conditions. So we do not need to reinvent the already wellestablished wheel of estimation here:<sup>31</sup> NAbC can provide the finite sample distribution of any estimator, and thus, make statistically valid and actionable inferences about a portfolio's dependence structure with those estimator(s) that are 'best' under specific conditions. Derivation and application of NAbC follows below in the next section.

<sup>&</sup>lt;sup>31</sup> The possible exception here, mentioned above as the topic of current research, is applying Average Oracle to dependence measures beyond Pearson's, which could be a notable improvement to their estimation over other methods for forecasting their conditional values.

#### 4.a. Brief Literature review of Pearson's Matrix: Distributional Results and Sampling Algorithms

I begin with Pearson's product moment correlation matrix, the oldest and arguably most widely used measure of dependence (see Pearson, 1895), in part because I derive below a fully analytic solution for NAbC for a special case of Pearson's matrix. Although its limitations often are mischaracterized or misunderstood, especially as they relate to widely held views classifying it strictly as a measure of linear association (see van den Heuvel & Zhan, 2022), in many settings Pearson's remains either optimal or centrally relevant for wide-ranging purposes.<sup>32</sup> These include robust asset allocation (Welsch and Zhou, 2007), Black-Litterman variants (Meucci, 2010a, Qian and Gorman, 2001), entropy pooling with fully flexible views (Meucci, 2010b), portfolio optimizations combined with random matrix theory (Pafka and Kondor, 2004), stress testing (Bank for International Settlements, Basel Committee on Banking Supervision, 2011a), and even non-linear, tail-risk-aware trading algorithms (Li et al., 2022, and Thakkar et al., 2021) to name a few. Consequently, Pearson's is the foundational dependence measure we start with (see also Rodgers & Nicewander (1988) for a broad, useful, and applied introduction to Pearson's).

When it comes to statistical inference and simulation-based decision-making, the extant literature on Pearson's matrix can be placed roughly into two categories: 1. distributional derivations that preserve inferential capabilities, and 2. sample-generating algorithms. The former typically are limited by unrealistic distributional assumptions, while the latter attempt to generate stylized, real-world distributions, but fail to preserve inferential (probabilistic) validity. Of course, we want both worlds: robust, fast, straightforward algorithms to generate samples when needed (i.e. in the absence of fully analytic solutions), that also preserve inferential capabilities, so that we can base decisions on rigorously defined probabilities. Unfortunately, none of the methods reviewed below provide both, which was a strong motivator for NAbC's development.

I begin with a brief and admittedly non-comprehensive, but well-targeted literature review of both 1. And 2. under more restricted cases. I follow with a review of 1. snd 2. under more general conditions. Subsequently I develop NAbC under both a narrowly defined but foundational case, and then under fully general conditions that satisfy the eight original objectives listed in the Introduction and Background. Defining NAbC under a narrow case provides a fully analytic version that very transparently shows how NAbC accomplishes both objectives listed above – useful sample generation and valid statistical inference, simultaneously – while also serving as a helpful, transparent referential baseline for NAbC's generalization to all positive definite dependence measures, under all data conditions.

<sup>&</sup>lt;sup>32</sup> In addition to the linearity 'myth' effectively addressed in van den Heuvel & Zhan (2022), note also that while Pearson's, under dependence, does not retain invariance under marginal transforms generally, the set of cases where it *does* retain invariance is broader than previously thought (see Koike et al., 2024).

#### 4.a.i. Distributional Results

Derivations of the distribution of Pearson's matrix go all the way back to one of the fathers of modern statistics, Sir Ronald A. Fisher (see Fisher, 1915, 1928). Intriguingly, Fisher (1928) recognizes the relationship between the Pearson's correlation formulae and the cosine between the angles of the two data vectors in the bi-variate case, i.e. what is now widely referred to as "cosine similarity" (described in more detail later).<sup>33</sup> He builds on this in his derivations (as does NAbC below), and although without closed forms some of the mathematical results prove unwieldy, they are foundational for those (re)derived below. Joarder and Ali (1992) replicate some of Fisher's earlier results (see Fisher, 1915), and more generally derive the distribution of Pearson's for any dimension when the underlying data is elliptically distributed (which includes the case of Gaussian data). Their density, however, requires iterated integration on the order of the dimension of the matrix, so like many of Fisher's results, while mathematically correct, it remains unscalable and less readily implemented.

For more recent results, below I start with narrowly defined cases and then expand. Restrictions on the narrow cases include i. on the underlying data (e.g. only Gaussian); ii. on the dimension of the matrix (e.g. only the bivariate case of p=2); iii. on the values of the matrix (e.g. only the identity matrix, where all correlations equal zero); and iv. with a priori known, rather than estimated, parameter values (e.g. known variances).

#### Gaussian data, any matrix, p=2

For Gaussian data with matrix dimension p=2, i.e. the bi-variate case, Taraldsen (2021) derives the exact confidence distribution of Pearson's correlation:

(17) 
$$\pi(\rho | r) = \frac{(1-r^2)^{(\nu-1)/2} \cdot (1-\rho^2)^{(\nu-2)/2} \cdot (1-r\rho)^{(1-2\nu)/2}}{\sqrt{2}B(\nu+1/2,1/2)} \cdot {}_2F_1\left(\frac{3}{2}, -\frac{1}{2}; \nu+\frac{1}{2}; \frac{1+r\rho}{2}\right) \text{ where }$$

 $B(X,Y) = \left[\Gamma(X)\Gamma(Y)\right]/\Gamma(X+Y)$  the Beta function, v = n - 1 > 1, and F is the Gaussian hypergeometric function where

$${}_{2}F_{1}[a,b;c;z] = \sum_{n}^{\infty} \frac{(a)_{n}(b)_{n}}{(c)_{n}} \cdot \frac{z^{n}}{n!} \text{ where } (h)_{n} = h(h+1)(h+2)\cdots(h+n-1), n \ge 1, (h)_{0} = 1^{34}$$

$$\cos\left(\hat{\theta}\right) = \frac{\text{inner product}}{\text{product of norms}} = \frac{\langle \mathbf{X}, \mathbf{Y} \rangle}{\|\mathbf{X}\| \|\mathbf{Y}\|} = \frac{\sum_{i=1}^{n} \left(X_{i} - \frac{1}{n} \sum_{j=1}^{n} X_{j}\right) \left(Y_{i} - \frac{1}{n} \sum_{j=1}^{n} Y_{j}\right)}{\sqrt{\sum_{i=1}^{n} \left(X_{i} - \frac{1}{n} \sum_{j=1}^{n} X_{j}\right)^{2}} \sqrt{\sqrt{\sum_{i=1}^{n} \left(X_{i} - \frac{1}{n} \sum_{j=1}^{n} Y_{j}\right)^{2}}} = \frac{C\hat{o}v_{X,Y}}{s_{X}s_{Y}} = r_{X,Y}, \ 0 \le \hat{\theta} \le \pi$$

<sup>34</sup> Interestingly, the Gaussian hypergeometric function makes many appearances in this and related settings: i. in derivations of the distribution of individual (bivariate) correlations (besides Taraldsen, 2021, see also Muirhead, 1982); ii. in moments of the spectral distribution under some conditions (see Adams et al. 2018, and

https://reference.wolfram.com/language/ref/MarchenkoPasturDistribution.html); iii. in the cumulative distribution function of JD Opdyke, Chief Analytics Officer Page 28 of 92 Correlation and Beyond

<sup>&</sup>lt;sup>33</sup> Briefly, this widely used mathematical relationship recognizes that the cosine of the angle between two mean-centered data vectors equals Pearson's (bivariate) correlation coefficient, as shown below:

It is important to note here, as discussed later in this monograph, that Taraldsen (2021) shows that the Fisher's Z-transformation (Fisher, 1921), a widely used approximation of this distribution, loses accuracy as correlation values approach one or negative one, especially for smaller samples.

## Gaussian data, identity matrix only, p≥2

For the Gaussian identity matrix (all correlations of zero) with  $p\ge 2$ , Gupta & Nagar (2000) derive the density

(18) 
$$f(R) = \frac{\left[\Gamma\left(\left[n-1\right]/2\right)\right]^p |R|^{(n-p-2)/2}}{\Gamma_p\left(\left[n-1\right]/2\right)}, \quad -1 \le r_{i,j} = r_{j,i} \le 1, \ r_{i,i} = 1, \ 1 \le i, j \le p \text{ and } \Gamma_p\left(\left[n-1\right]/2\right) = \pi^{\left[p(p-1)/4\right]} \prod_{i=1}^p \Gamma\left(\left[n-i\right]/2\right)$$

Although Pham-Gia & Choulakian (2014) claim this is a new result, it actually is identical to that of Joarder and Ali (1992) (see (4.2)) under these conditions, and after some manipulation, that of Fisher (1915) for the bivariate case (see (4.2), (4.3), and (3.1) in Joarder & Ali, 1992).

## Gaussian data, any matrix, with p≥2

Under Gaussian data, with p≥2, Pham-Gia & Choulakian (2014) provide the distribution of the sample Pearson's matrix under any values, not just the identity matrix:

(19)  
$$f(R) = \frac{\left[\Gamma\left(\left[n-1\right]/2\right)\right]^2 \exp\left\{-\sum_{i< j} \frac{\lambda_{i,j} S_{i,j}}{\sqrt{\sigma_{i,i}\sigma_{j,j}}}\right\}}{\Gamma_p\left(\left[n-1\right]/2\right)\left[|\Lambda||\prod_{i=1}^p \lambda_{i,i}\right]^{\left[\left[n-1\right]/p\right)}} |R|^{(n-p-2)/2}$$

with sample covariance  $\{s_{i,j}\}$ ,  $1 \le i < j \le p$ ,  $\Gamma_p([n-1]/2) = \pi^{\lfloor p(p-1)/4 \rfloor} \prod_{i=1}^p \Gamma([n-i]/2)$ ,

known variance  $\{\sigma_{i,i}\}$ , | | is the determinant function,  $\Lambda$  is the true correlation matrix, and  $\lambda_{i,i}$  the diagonals of  $\Lambda^{-1}$ 

The limitations of Pham-Gia & Choulakian (2014) include the requirement of a priori knowledge of true (not estimated) variances, and of course the fact that it is valid only for normally distributed data. It also, arguably, is quite cumbersome to implement.

## 4.a.ii. Sampling Algorithms

Moving now to sample generation under various 'non-generalized' conditions, i.e. conditions that are not generalized to those common in financial portfolios, the literature provides a number of methods, many of which are quite involved. Note that I have focused on more recent ones, as these usually explicitly

Pearson's under the Gaussian identity matrix of any dimension (see Opdyke, 2022, 2023, and 2024a); and iv. in the definition of positive definite functions (see Franca & Menegatto, 2022).

subsume previously published algorithms, and many of the below are even compared against each other in their own empirical simulation studies. Note that none of these are generalized to include the stylized empirical characteristics observed in financial portfolios, and hence are labelled here as 'nongeneralized'.

- i. The onion and c-vine methods, the former of which can generate random correlation matrices with the joint density of the correlations being proportional to a power of the determinant of the correlation matrix, and the latter of which is based on partial correlations specified in a vine copula (specifically, a c-vine copula). (Lewandowski et al., 2009)
- the chordal sparsity method of Kurowicka (2014), which generalizes Lewandowski et al. (2009), although "it is not clear whether it is possible to extend them to other patterns of unspecified correlations" beyond those with chordal sparsity patterns.
- iii. The restricted Wishart distribution approach of Wang et al. (2018), which is equivalent to Lewandowski et al. (2009) but somewhat more efficient.
- iv. The hyperspherical coordinate approach of Pourahmadi et al. (2015)
- v. The Cholesky-Metropolis method of Cordoba et al. (2018), which claims to be faster than the previously listed methods.
- vi. The direct formulation method of Madar (2015)
- vii. The flexible bijection method of Veleva (2017)
- viii. The rejection algorithm of Makalic and Schmidt (2018), which is based on the polar (hyperspherical) angles representation of Pearson's matrix<sup>35</sup>

Makalic and Schmidt (2018) is treated in more detail below. Implementation of all but vi., vii., and viii. above arguably remain quite involved, but one of the self-described focuses of most of these is computational efficiency (which is not surprising as they are sampling algorithms). From a close read of the runtime results of the successively published and compared algorithms above, it appears that Makalic and Schmidt (2018) is the fastest among them (excepting those of Madar (2015) and Veleva (2017), which have not been compared to the others). However, as discussed in more detail below, Rubsamen (2023) shows that for the case of the Gaussian identity matrix, when NAbC is used in its sample generation capacity, it is over 30% faster than Makalic and Schmidt (2018); when NAbC is used for this specific case analytically, its results are, for all intents and purposes, instantaneous, as can be seen in the excel workbook at the following link (see http://www.datamineit.com/JD%20Opdyke--The%20Correlation%20Matrix-

Analytically%20Derived%20Inference%20Under%20the%20Gaussian%20Identity%20Matrix--02-18-24.xlsx).

Runtimes of NAbC under the fully general case, i.e. that of fully realistic financial data conditions (which

JD Opdyke, Chief Analytics Officer

<sup>&</sup>lt;sup>35</sup> See Joarder & Ali (1992), Pinheiro & Bates (1996), Rebonato & Jaeckel (2000), Rapisarda et al. (2007), and Pourahmadi & Wang (2015). The use of spherical angles for analysis of Pearson's matrix goes back at least to Fisher (1915), but Joarder & Ali (1992) and Rapisarda et al. (2007) provide geometrically motivated, thorough, and clear descriptions of its derivations, and Rebonato & Jaeckel (2000) appears to have been the first to propose its application in financial settings.

none of the algorithms above claim to cover), as well as any valid matrix values, are discussed in later sections below.

But beyond speed, the more important issue regarding these sampling algorithms is that none preserve inferential validity, on a sample-by-sample basis, by providing a readily calculated cumulative distribution function (cdf) value (probability density function (pdf) values will be more cumbersome and less useable here). In other words, to make these simulation results truly useful for valid hypothesis testing and other inferential purposes, we need to know where in the distribution of correlation matrix samples a particular sample sits: given a known/true correlation matrix (our null hypothesis, if hypothesis testing), what probability is associated with observing a particular sample correlation matrix, or one more extreme, relative to the known matrix? Is it larger than 5% of all randomly generated samples? Or 99%? Calculating this cdf value is tortured, if not impossible for these methods, although this arguably should be the primary focus of such algorithms.

The counter argument is that the group of samples that these algorithms generate, taken as a whole, is a valid representation of the data generating mechanism behind the specified correlation matrix. This group of samples can then be used as inputs to comprehensive and arguably real-world portfolio simulations. While this is certainly true, at best the group of sample correlation matrices, then, only provide indirect inferential value, with what is arguably a notable lack of control. For example, the group of samples cannot be used to specify, for controlled portfolio simulations, the two matrices representing the 95% confidence interval under a given null hypothesis for the dependence structure, which is the kind of targeted, controlled capability needed for precise, powerful testing and consequent, targeted decision-making. Some may argue that the group of sample matrices can be used with ad hoc measures of 'distance' from a hypothesized matrix of 'true' values (e.g. a Euclidean 'norm' distance from, say, the identity matrix), yet such multivariate distances can be measured in many different and equally valid ways under various conditions (this is addressed in more detail in the Section 7 below). The same 'distances' also can have different interpretations under different conditions, and even widely used ones can be 'wrong' when applied to very commonly used dependence measures in this setting (see Zhang et al., 2024, for a compelling example). So making inferences based on them remains arguably as ad hoc, at best, as the arbitrary choice of how to measure distance between a sample matrix and its null hypothesis. Neither can such distances be used to rank order the sampled matrices to obtain, say, an empirical cdf, because different distances will yield different rank orderings. The only 'distance' that avoids these issues is probability itself, most conveniently and rightly represented as a cdf value.

In the end, for inferential capability and subsequent probability-based decision-making ability, what is necessary here is an analytically rigorous connection between a specific sampled correlation matrix and its associated, properly defined cdf value, and none of these sampling methods provide this. Fortunately, NAbC does, as is discussed further below. But first, I finish reviewing results from the extant literature that cover 1. distributional derivations of Pearson's matrix, but under the more general case of realistic, financial returns data, as well as 2. sample generation algorithms of Pearson's matrix under these same conditions (which should correspond to their stylized, empirical characteristics). For 1., I cover three recent and intriguing methods below.

JD Opdyke, Chief Analytics Officer

Page **31** of **92** 

**Correlation and Beyond** 

## 4.a.iii. Distributional Results, More General Conditions

Archakov & Hansen (2021) introduce an original parameterization of Pearson's matrix that maps uniquely, one-to-one, to the positive definite space, thus providing a density for inference. It is valid under general conditions, based on the Fisher Z transformation, remains invariant to reorderings of the variables/assets, i.e. the rows and columns of the matrix, and is accompanied by an algorithm that provides the inverse mapping from the parameterization to the correlation matrix (i.e. a matrix level quantile function).

This approach is original, but the method still has limitations. As the authors state, "This makes the transformation potentially useful for ... inference. These attributes tend to deteriorate as C approaches singularity. This is not unexpected, because it is also true for the Fisher transformation when the correlation is close to ±1." As previously noted, Taraldsen (2021) similarly shows that the approximate density of the pairwise correlation using Fisher's Z-transformation loses accuracy as correlation values approach ±1, especially for smaller samples. This is consistent with the authors' comments here, but they state this may only be material under extreme conditions. However, extreme conditions are exactly when correlation breakdowns occur, and when we most need robust, accurate inferences. All else equal, it would be preferrable and important in practice to have a method that avoided this non-robustness issue altogether.

In addition, the method only provides the distribution of the entire correlation matrix: it does not appear to be able to modify correlation matrices, cell-by-cell, probabilistically, based on their individual celllevel distributions, for things like scenarios and 'what if' analyses. While this may not be a stated objective of the method, all else equal it would be a very useful feature for stress testing and scenario analysis, as well as attribution analyses. The same holds true for larger submatrices of the matrix, i.e. submatrices larger than one cell. Note that NAbC shares none of these limitations, but shares all of the method's advantages listed above, in addition to many others.

Lan et al. (2020) take a fully Bayesian approach to this problem for both covariance and Pearson's correlation matrices. Similar to NAbC, they use the Cholesky factorization to automatically enforce positive definiteness, and by defining distributions on spheres as NAbC does, utilize a large class of flexible prior distributions. This method includes *estimation* of the correlation/covariance matrix, which as described above, NAbC does not. But it also lacks very important advantages that NAbC provides. As the authors state, "The priors for correlation matrix specified through the sphere-product representation are in general dependent among component variables. For example, the method we use to induce uncorrelated prior between  $y_i$  and  $y_j$  (i < j) by setting  $l_{jk} \approx 0$  for  $k \leq i$  has a direct consequence that  $\operatorname{Cor}(y'_i, y_j) \approx 0$  for  $i' \leq i$ . In another word [sic], more informative priors (part of the components are correlated) may require careful ordering in  $\{y_i\}$ . To avoid this issue, one might consider the inverse of covariance (precision) matrices instead. This leads to modeling the conditional dependence, or *Markov network* ... Our proposed methodology applies directly to (dynamic) precision matrices/processes, which will be our future direction."

JD Opdyke, Chief Analytics Officer Page 32 of 92

Fortunately, NAbC does not suffer from this order-dependence problem. Like Archakov & Hansen (2021) its results are invariant to the ordering of the rows and columns of the matrix, but unlike Archakov & Hansen (2021) or Lan et al. (2020) it can 'freeze' any submatrix of the correlation matrix, even if it is non-contiguous, as dictated by any particular scenario or stress test, and still obtain a valid, finite-sample distribution for the (rest of the) matrix. In all cases within the framework of the all-pairwise matrix, there are no unintended 'dependencies' between cells that confound these results. As discussed further below, the 'unintended dependencies' problem is one that other researchers have struggled with, but which NAbC avoids.

Like Lan et al. (2020), Ghosh et al. (2021) also take a fully Bayesian approach to this problem, and just like NAbC, they reparameterize Cholesky factors in terms of hyperspherical coordinates where the angles vary freely in the range  $[0, \pi)$ . Their focus is on estimation, although as a Bayesian approach it is comprehensive and provides credible regions. Among its arguable limitations, however, is that its use is restricted to parametric priors, which given the non-small dimensions of most financial portfolios (e.g. Bongiorno & Challet (2023a) call p=100, which has p(p-1)/2=4,950 pairwise cells, 'mid-sized') it is hard to see how this would not limit its implementation under complex, real world financial data conditions (e.g. 4,950 distributions with different and varying degrees of serial correlation, asymmetry, non-stationarity, and heavy-tailedness). In other words, it is hard to imagine a fully parametric multivariate distribution of dimension p=100 or greater that was analytically tractable but simultaneously able to adequately incorporate all of the distributional characteristics listed above, for 4,950 distributions (or more). In contrast, NAbC makes use of flexible nonparametric kernels that fit ANY angles distribution resulting from ANY data generating mechanism (with finite mean and variance required for Pearson's matrix). Also, like both Archakov & Hansen (2021) and Lan et al. (2020), the approach of Ghosh et al. (2021) does not appear to have the capability of modeling submatrices while leaving select cells of the correlation matrix 'untouched.' This is absolutely essential for flexible and realistic scenario modeling and (reverse) stress testing, and one of the many advantages NAbC provides. Now I treat some of the more recent, generalcase sample generation algorithms.

## 4.a.iv. Sampling Algorithms, More General Conditions

Marti (2019) proposes using generative adversarial networks (GANs) to incorporate the stylized empirical characteristics of financial portfolios' correlation matrices into an algorithm (CorrGAN) that directly generates samples of the all-pairwise matrix, as opposed to generating samples of returns from which correlation matrices are then estimated. This appears to be the first method to attempt this approach. The stylized characteristics include positive-shifted correlations, Marchenko-Pastur distributed correlations excepting a few large eigenvalues, Perron-Frobenius property (positive entries of the first eigenvector), hierarchical correlation structure, and scale-free property of the corresponding minimum spanning tree. Marti (2019) does not address how computationally intensive is the method, but apparently it is not prohibitively so as he implements it on 100x100 matrices in a follow-on blog post (see Marti, 2020).

JD Opdyke, Chief Analytics Officer

Page **33** of **92** 

**Correlation and Beyond** 

There are three main limitations to this approach. First, as the author notes, while it appears to capture most of the identified distributional stylized facts, it does not capture the tails well. This arguably is the most important part of the distribution, as it is critically related to portfolio risk analytics, and many if not most scenarios, especially those that incorporate events like correlation breakdowns. In addition to use in trading algorithms, these are the stated purposes of the method, so difficulty estimating the tails of the distribution is not a minor limitation. Secondly, the method generates samples that "are not exactly correlation matrices" with non-unit diagonals and negative eigenvalues. Marti (2019) states that positive definiteness is enforced ex post using Higham (2002). I have used Higham (2002) extensively in my research in this setting, closely examining its effects on both the spectral distribution and the distribution of the correlation matrix itself, and have found that both can be dramatically distorted when Higham (2002) is used.<sup>36</sup> However, simply discarding non-positive definite samples is not the answer as this, too, will distort its true sample distribution. The only way to simulate the matrix and generate a valid, nondistorted, representative pool of sample matrices is to use a method that automatically enforces positive definiteness, ex ante. Finally, the last limitation is that the samples generated by this method do not retain inferential validity generally, that is, they are not associated with a probability of occurrence or a cdf value, as discussed above.

Papenbrock et al. (2021) develop a novel and intriguing approach to simulating correlation matrices for financial markets using evolutionary algorithms. These allow for the flexible yet robust incorporation of many observed features of real-world financial correlation matrices (the list is similar to that of Marti (2019), with some enrichments). The algorithm scales well and can be used for backtesting, pricing, and hedging correlation-dependent investment strategies and financial products. However, it has several limitations: the first relates to how upper and lower barriers are established for the sampled correlation matrices, which the authors describe as "This neighborhood could be defined in a static way or by expert knowledge." Regarding the latter option, making this criterion (strictly) subjective arguably defeats the purpose of objective, quantitative analysis in this setting. Regarding the former option, Papenbrock et al. (2021) suggest using the most extreme values of the matrices, although none of the implementations listed (e.g. random matrix denoising, shrinkage, or exponential weighting) are inferentially valid in themselves, that is, they do not allow for probabilistic inference. Arguably, if the range of the matrices sampled needs to be restricted at all, it should be restricted based on rigorously defined probabilistic bounds, say, 99% confidence intervals. Fortunately, this is a capability that NAbC provides.

The second limitation of Papenbrock et al. (2021) is shared with Marti (2019) in that the algorithm does not enforce positive definiteness ex ante. The authors do acknowledge the importance of positive definiteness in this setting, but do not explain how their algorithm handles non-positive definite samples. Again, both ignoring/eliminating them from consideration, and/or 'fixing' them with algorithms like that of Higham (2002), distort the distribution of the correlation matrix in non-trivial ways, and thus invalidate

<sup>&</sup>lt;sup>36</sup> To be clear, this is not a critique of Higham (2002), which is seminal and extremely useful in wide-ranging, applied settings. Rather, it is only to say that 'fixing' non-positive definite matrices that are generated by non-trivially complex algorithms typically, if not always in practice, strongly distorts the **distribution** of the sample matrices, as well as the associated spectral distribution. This is readily verified empirically.

inferences based on it. Finally, the sample correlation matrices generated by the evolutionary algorithms are not inferentially valid in themselves, i.e. each is not associated with a cdf probability. Again, none of these limitations – subjective or ad hoc restriction of the sample space, ex post enforcement of positive definiteness, or lack of inferential validity – apply to NAbC.

A sophisticated and more recent attempt at directly generating sample correlation matrices with stylized characteristics of real-world financial data is that of Kubiak et al. (2024). They develop denoising diffusion probabilistic models (DDPMs) that, across multiple asset classes and market regimes, compare favorably against a number of alternate modern approaches, including CorrGAN of Marti (2019) and other GANs approaches, <sup>37</sup> variational autoencoders, and more traditionally, block bootstraps. Limitations of the approach are similar to those of the other 'direct sampling' algorithms: i. the matrices generated are not true correlation matrices, lacking unit diagonals and true asymmetry; ii. the matrices are not guaranteed to be positive definite, and when they are not, positive definiteness is enforced ex post using Higham (1988); and iii. the sampled matrices do not retain inferential validity, that is, they are not associated with a probability of occurrence (a cdf value), as described above.

This is an interesting area of research and the work of these approaches, real limitations notwithstanding, is encouraging. This is especially true because NAbC, which shares none of these limitations, can be applied to the realistic groups of samples that they generate to give them inferential validity (as long as the sample distribution is not distorted). If researchers can find a way for these algorithms to generate true correlation matrices and automatically enforce positive definiteness ex ante, or convincingly prove that deviations from either are truly numerically de minimis along any dimension of analysis (which certainly has not been done to date), NAbC can be applied to their samples to provide inferential validity, that is, to associate a cdf value with each and every sample matrix generated so that the distribution can be used inferentially. Again, the starting point for NAbC's application is the known or well-estimated dependence matrix, and the known or well-estimated data generating mechanism (in these cases, the *matrix* generating mechanism), and these methods provide both (again, as long as the samples are representative of the true distribution). But until these two 'fixes' can be applied (i.e. use only non-'fixed' matrices for which positive definiteness always holds ex ante), with proof that this truly has been achieved numerically, if not analytically, real inferential challenges will remain as obstacles for this path of 'direct matrix simulation' research.

It also is notable that none of this work has demonstrated that generating samples of correlation matrices by first generating synthetic *returns* data that have all of the empirical, stylized characteristics of actual *returns* data, is not sufficient, if carefully done, to generate the desired correlation sample matrices (even if not based on historically realized data, but rather, for plausible future scenarios). This connection needs to be established, mathematically and explicitly, because in reality the sample matrices are only and exclusively based on the sample returns, and without mathematically defining the path from the returns data to the stylized sample matrices, something is missing, if not wrongly specified, on one end or the other. It is not that we cannot or should not jump right to directly sampling the

<sup>&</sup>lt;sup>37</sup> Kubiak et al. (2024) believe this is due to "the standard instability issues commonly associated with GAN training. (p.5)" JD Opdyke, Chief Analytics Officer Page 35 of 92 Correlation and Beyond

correlation matrix per se; only that the connection between all of the stylized characteristics of the correlation matrix, and those of the returns data on which it is based, needs to be rigorously established if we are to have real insight into the mechanics, provably appropriate simulations, and distribution-based inferences (if not causal drivers) of the former. This connection, in fact, likely holds the key to solving at least two of the three problems these methods face: i. the matrices generated are not true correlation matrices, lacking unit diagonals and true asymmetry; and ii. the matrices are not guaranteed to be positive definite.

In the absence of an explanation as to why it is not preferable to start with the returns data itself (aside from computational considerations), I hypothesize that part of the motivation of taking the 'direct simulation' route, even if not explicitly stated, is the right-minded desire to separate and isolate the distribution of the correlation/dependence matrix from other characteristics of the distribution of the returns data. And this is exactly what NAbC does, but it does so while managing to preserve inferential validity (not to mention valid matrices). In some cases, it does this fully analytically, as I show in the next section. I start with a narrow but foundational special case, which provides the fully analytical result, and then expand NAbC's application to very general conditions, because the first provides a useful referential baseline for understanding the mechanics of the latter.

## 4.b. NAbC: Pearson's Correlation, the Gaussian Identity Matrix

## 4.b.i. Correlations to Angles, Angles to Correlations

I continue with Pearson's for the first derivation and implementation of NAbC, and the data and correlation structure I initially presume is Gaussian data under no correlation: that is, Pearson correlation values of zero off the diagonal of the matrix as below.

$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$		1	0	0	0	
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$		0	1	0	0	
-0001 for $p = 4$ assorts		0	0	1	0	
	=	0	0	0	1	for p = 4 assets

identity matrix

The key to the NAbC approach rests in its use of the spherical ANGLE,  $\theta$ , between the two meancentered data vectors of X and Y, as opposed to directly and only using of the values of the correlations themselves. As mentioned above, using angles to understand the distribution of Pearson's goes back at least to Fisher (1915), and it turns out to be a much more general framework applicable to any dependence measure whose all-pairwise matrix is positive definite, not just Pearson's. But to start with Pearson's, for a single pair of variables, providing a single bivariate correlation value, the relationship between angle value and correlation value is most readily seen in the widely known "cosine similarity," where the cosine of the angle equals the inner product divided by the product of the two vectors' (Euclidean) norms as in (16), which I show again below for the reader's convenience.
(16) 
$$\cos\left(\hat{\theta}\right) = \frac{\text{inner product}}{\text{product of norms}} = \frac{\langle \mathbf{X}, \mathbf{Y} \rangle}{\|\mathbf{X}\| \|\mathbf{Y}\|} = \frac{\sum_{i=1}^{n} \left(X_{i} - \frac{1}{n} \sum_{j=1}^{n} X_{j}\right) \left(Y_{i} - \frac{1}{n} \sum_{j=1}^{n} Y_{j}\right)}{\sqrt{\sum_{i=1}^{n} \left(X_{i} - \frac{1}{n} \sum_{j=1}^{n} X_{j}\right)^{2}} \sqrt{\sum_{i=1}^{N} \left(X_{i} - \frac{1}{n} \sum_{j=1}^{n} Y_{j}\right)^{2}}} = \frac{C \hat{o} v_{X,Y}}{s_{X} s_{Y}} = r_{X,Y}, \ 0 \le \hat{\theta} \le \pi$$

If a portfolio has p assets, the number of its pairwise relationships is npr=p(p-1)/2. For all these npr relationships, the matrix analogue to (16), as long as the matrix is symmetric positive definite,<sup>38</sup> is well established in the literature (Joarder and Ali, 1992, Pinheiro and Bates, 1996, Rebonato and Jackel, 2000, Rapisarda et al., 2007, Pouramadi and Wang, 2015, and Cordoba et al., 2018) and shown below, formulaically in (20)-(22) and in computer code (SAS/IML) in Table A. The steps for translating between correlations and angles, in both directions, are straightforward and shown in A.-C. below.

A. estimate the correlation matrix from sample data

B. obtain the Cholesky factorization of the correlation matrix

C. use inverse trigonometric and trigonometric functions on B. to obtain corresponding spherical angles

and in reverse:

- C. start with a matrix of spherical angles
- B. apply trigonometric functions to obtain the Cholesky factorization
- A. multiply B. by its transpose to obtain the corresponding correlation matrix

(see Rebonato & Jaeckel, 2000, Rapisarda et al., 2007, and Pourahmadi & Wang, 2015, but note a typo in the formula in Pourahmadi & Wang, 2015, for the first 3 steps)

Central to this correlation-angle translation mechanism is obtaining the Cholesky factor of the correlation/dependence matrix, which is usually a built-in function in most statistical and mathematical software. The relevant formulae are included below for completeness.

(20) A correlation matrix *R* will be real, symmetric positive-definite,<sup>39</sup> so the unique matrix B that satisfies

 $R = BB^T$  where B is a lower triangular matrix (with real and positive diagonal entries), and  $B^T$  is its transpose, is the Cholesky factorization of R. Formulaically, B's entries are as follows:

$$B_{j,j} = (\pm) \sqrt{R_{j,j} - \sum_{k=1}^{j-1} B_{j,k}^2}, \quad B_{i,j} = \frac{1}{B_{j,j}} \left( R_{i,j} - \sum_{k=1}^{j-1} B_{i,k} B_{j,k} \right) \text{ for } i > j$$

<sup>&</sup>lt;sup>38</sup> Note again that this is true not only for Pearson's, but also for all relevant (i.e. positive definite) dependence measures in this setting, as will be discussed below.

<sup>&</sup>lt;sup>39</sup> Semi-positive definiteness includes the case of eigenvalues exactly equal to zero, which I largely ignore herein as a border case relevant mainly for textbook examples since returns would have to exhibit perfect linear dependence for an eigenvalue to be exactly zero.

The Cholesky factor can be viewed as a matrix analog to the square root of a scalar, because like a square root the product of it and its transpose yields the original matrix. Importantly, the Cholesky factor places us on the UNIT hyper-(hemi)sphere (where scale does not matter) because the sum of the squares of its rows always equals one. Next, we recursively apply inverse trigonometric and trigonometric functions to each cell of the Cholesky factor to obtain each cell's angle value per (21); or in reverse, we obtain a correlation/dependence value from trigonometric functions applied to each cell's angle value per (22) (see both Joarder & Ali, 1992, and Rapisarda et al., 2007, for meticulous derivations of these formulas). Note that this relationship is one-to-one, with a unique correlation/dependence matrix yielding a unique angles matrix, and vice versa.

$$R = \begin{bmatrix} 1 & r_{1,2} & r_{1,3} & \cdots & r_{1,p} \\ r_{2,1} & 1 & r_{2,3} & \cdots & r_{2,p} \\ r_{3,1} & r_{3,2} & 1 & \cdots & r_{3,p} \\ r_{4,1} & r_{4,2} & r_{4,3} & \cdots & r_{4,p} \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ r_{p,1} & r_{p,2} & r_{p,3} & \cdots & 1 \end{bmatrix}$$

(21) For *R*, a p x p correlation matrix,

 $R = BB^t$  where B is the Cholesky factor of R and

$$B = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0\\ \cos(\theta_{2,1}) & \sin(\theta_{2,1}) & 0 & \cdots & 0\\ \cos(\theta_{3,1}) & \cos(\theta_{3,2})\sin(\theta_{3,1}) & \sin(\theta_{3,2})\sin(\theta_{3,1}) & \cdots & 0\\ \cos(\theta_{4,1}) & \cos(\theta_{4,2})\sin(\theta_{4,1}) & \cos(\theta_{4,3})\sin(\theta_{4,2})\sin(\theta_{4,1}) & \cdots & 0\\ \vdots & \vdots & \vdots & \cdots & \vdots\\ \cos(\theta_{p,1}) & \cos(\theta_{p,2})\sin(\theta_{p,1}) & \cos(\theta_{p,3})\sin(\theta_{p,2})\sin(\theta_{p,1}) & \cdots & \prod_{k=1}^{n-1}\sin(\theta_{p,k}) \end{bmatrix}$$

for i > j angles  $\theta_{i,j} \in (0,\pi)$ .

To obtain an individual angle  $\,\theta_{\scriptscriptstyle i,j}\,$  , we have:40

For 
$$i > 1$$
:  $\theta_{i,1} = \arccos(b_{i,1})$  for  $j=1$ ; and  $\theta_{i,j} = \arccos\left(\frac{b_{i,j}}{\prod_{k=1}^{j-1}\sin(\theta_{i,k})}\right)$  for  $j > 1$ 

(22) To obtain an individual correlation,  $r_{i,j}$ , we have, simply from  $R = BB^T$ :

$$r_{i,j} = \cos(\theta_{i,1})\cos(\theta_{j,1}) + \prod_{k=2}^{i-1}\cos(\theta_{i,k})\cos(\theta_{j,k})\prod_{l=1}^{k-1}\sin(\theta_{i,l})\sin(\theta_{j,l}) + \cos(\theta_{j,i})\prod_{l=1}^{i-1}\sin(\theta_{i,l})\sin(\theta_{j,l}) \quad \text{for } 1 \le i < j \le n$$

<sup>&</sup>lt;sup>40</sup> Note that as shown in Madar (2015), a similar recursive relationship exists between partial correlations, and Madar (2015) generalizes this result beyond Pearson's to any positive definite matrix.

SAS/IML code translating correlations to angles and angles to correlations is shown in Table A below.

The above all is well-established and straightforward,<sup>41</sup> and demonstrates, as we know intuitively, that **scale does not (and should not) matter when it comes to dependence measures;**<sup>42</sup> again, in this setting, this is because geometrically, the Cholesky factor places us on the UNIT hyper-(hemi)sphere. Importantly, the Cholesky factor also ensures that sampling based directly on the resulting angles will yield only positive definite matrices, as the Cholesky factor remains undefined otherwise. This **automatic enforcement of positive definiteness makes this approach much more efficient than others** that require ex post verification of positive definiteness, and subsequent resampling or enforcement when this requirement is violated (examples of this, discussed above, include Makalic & Schmidt, 2018; Cordoba et al. 2018; Marti, 2019; Papenbrock et al., 2021; and Kubiak et al., 2024). This

#### TABLE A:

Correlations to Angles	Angles to Correlations
* INPUT rand_R is a valid correlation matrix;	* INPUT rand_angles is a valid matrix of correlation angles;
<pre>cholfact = T(root(rand_R, "NoError")); rand_corr_angles = J(nrows,nrows,0); do j=1 to nrows; do j=j to nrows; if j=j then rand_corr_angles[i,j]=.; else do; cumprod_sin = 1; if j=1 then rand_corr_angles[i,j]=arcos(cholfact[i,j]); else do; do kk=1 to (j-1); cumprod_sin = cumprod_sin*sin(rand_corr_angles[i,kk]); end; rand_corr_angles[i,j]=arcos(cholfact[i,j]/cumprod_sin); end; end; end; end; </pre>	<pre>Bs=J(nrows, nrows, 0); do j=1 to nrows; do i=j to nrows; if j&gt;1 then do; if j&gt;j then do; sinprod=1; do gg=1 to (j-1); sinprod = sinprod*sin(rand_angles[i,gg]); end; Bs[i,j]=cos(rand_angles[i,j])*sinprod; end; else do; sinprod=1; do gg=1 to (i-1); sinprod = sinprod*sin(rand_angles[i,gg]); end; Bs[i,j]=sinprod; end; do gg=1 to (i-1); sinprod = sinprod*sin(rand_angles[i,gg]); end; end; end; end; end; end; end; end</pre>
* OUTPUT rand corr angles is the corresponding matrix of angles;	else do; if i>1 then Bs[i,j]=cos(rand_angles[i,j]); else Bs[i,j]=1; end; end; end; rand_R = Bs*T(Bs);
SAS/IML code (v9.4)	* OUTPUT rand_R is the corresponding correlation matrix;

<sup>&</sup>lt;sup>41</sup> Reliance on spherical angles and hypersphere parameterizations is increasingly common in quantitative finance (see for some examples Li, Q., 2018; Helton, 2020; Golts & Jones, 2009; Zhang, 2022; Saxena et al., 2023; and Zhang & Yang, 2023), in large part due to its scale invariance: it has even been used to define entire financial markets (see Kim and Lee, 2016).

JD Opdyke, Chief Analytics Officer

Page **39** of **92** 

<sup>&</sup>lt;sup>42</sup> Scale invariance is proved and widely cited for Pearson's, Kendall's, and Spearman's (see Xu et al., 2013, and Schreyer et al., 2017 for examples).

inefficiency grows very rapidly with the size of the matrix/portfolio, as shown in the ratio below in (23) (see Bohn and Hornik, 2024, and Pourahmadi & Wang, 2015).

(23) 
$$\Pr(rand "R" \sim PosDef) = X = \frac{\prod_{j=1}^{p-1} \left[ \sqrt{\pi} \Gamma\left(\frac{j+1}{2}\right) \right]^j}{2^{p(p-1)/2}} < \prod_{j=1}^{p-1} \left[ \frac{\sqrt{\pi}}{2} \right]^j = \left[ \frac{\sqrt{\pi}}{2} \right]^{p(p-1)/2}; \lim_{p \to \infty} \left[ X \right] = 0$$

Even for relatively small matrices of dimension p=25, the odds of successfully randomly generating a single valid positive definite correlation matrix, by uniformly sampling the off-diagonal correlation values themselves across values ranging from -1.0 to 1.0, are less then 2 in 10 quadrillion, leading to prohibitively inefficient sampling. Consequently, even when sampling-rejection algorithms achieve some efficiency gains, realistically the sampling approach in this setting should possess automatic enforcement of positive definiteness, ex ante. Conceptually, an imperfect but apt analogy is to a rubik's cube: the colored stickers on the cube cannot simply be peeled off and repasted, even some of the times, to solve the cube. The valid solution must be obtained by (always) following the rules governing shifts in the cube, and every move of each of the small cubes (correlation cells) affects the positions of many of the other cubes (correlation cells), not just the one we need to reposition. Similarly with sampling the correlation/dependence matrix: converting to the Cholesky factor (en)forces positive definiteness by forcing the matrix onto the UNIT hyper-(hemi)sphere, where we can subsequently use the distributions of the angles to perturb it and obtain, after re-translation, the distribution of the original correlation/dependence matrix, without violating positive definiteness. This is done simply by following steps A., B., and C., and C., B., and A., above. Importantly, aside from efficiency issues, this avoids distortion of the *distribution* of these samples via ex post enforcements of positive definiteness using algorithms like Higham (2002), and thus preserves inferential validity.

Another crucial characteristic of these angles is that **they are random variables whose multivariate relationship is one of independence** (see Pourahmadi and Wang, 2016; Ghosh et al., 2021; Rapisarda et al., 2007; Tsay and Pourahmadi, 2017; and Zhang et al., 2015).<sup>43</sup> This is critically important for practical usage as it enables the straightforward construction of the multivariate distribution of a matrix of angles, which is the more important objective here (vs merely sampling) and essential for the application of NAbC below.

Finally and most critically, the above demonstrates that **the angles between pairwise data vectors contain ALL the information that exists regarding dependence between the two variables** because the only information we lose by translating to the unit hyper(hemi-)sphere is scale (see Fernandez-Duren & Gregorio-Dominguez, 2023, and Zhang & Songshan, 2023, as well as Opdyke, 2022). This will be covered more extensively below.

<sup>&</sup>lt;sup>43</sup> This independence is well established in the literature. Zhang et al. (2015) (supplementary material) and Rapisarda et al. (2007) use a geometric interpretation of the correlation matrix, based on (orthogonal) Givens rotations, to explain in detail the relationship between correlations and angles as well as why the angles distributions are multivariate independent.

So with all this in mind we proceed with the use of the angles as described and defined above.<sup>44</sup> The goal is to use the angles as the basis for 1. sample generation of the correlation/dependence matrix; and more importantly, 2. definition of the multivariate distribution of the correlation/dependence matrix.

# 4.b.ii. Fully Analytic Angles Density, and Efficient Sample Generation

Once we have the matrix of angles (per (21) and Table A above), one angle for each value in the allpairwise correlation/dependence measure matrix, we use the well-established finding that, to sample uniformly from the space of positive definite matrices, the probability density function (pdf) must be proportional to the determinant of the Jacobian of the Cholesky factor as in (24) (see Cordoba, 2018, Pourahmadi & Wang, 2015, and Lewandowski et al., 2009).

(24) det
$$[J(U)] = 2^{p} \prod_{i=1}^{p-1} u_{ii}^{i}$$
 where U is the Cholesky factorization of correlation matrix  $R = UU^{t}$ 

We see directly from (24) that  $\sin^{k}(x)$ , suitably normalized in (25), satisfies this requirement (see Pourahmadi & Wang, 2015, and Makalic & Schmidt, 2018).

(25) 
$$f_X(x) = c_k \cdot \sin^k(x), x \in (0,\pi), k = 1,2,3,..., (\# \text{columns}-1), \text{ and } c_k = \frac{\Gamma(k/2+1)}{\sqrt{\pi}\Gamma(k/2+1/2)}$$

Although not explained in Makalic & Schmidt (2018), importantly note that k = #columns – column# (so for the first column of a p=10x10 matrix, k=9; for the second column, k=8, etc.).

Beyond (25), however, we need both the cumulative distribution function (cdf) and its inverse, the quantile function, to make use of this density for sampling and other purposes. The most widely used and straightforward method of sampling is inverse transform, whereby the values of a uniform random variate are passed to the quantile function to generate sampled values. Yet regarding the cdf corresponding to (25) above, Makalic & Schmidt (2018) state, "Generating random numbers from this distribution is not straightforward as the corresponding cumulative density [sic] function, although available in closed form, is defined recursively and requires O(k) operations to evaluate. The nature of the cumulative density [sic] function makes any procedure based on inverse transform sampling computationally inefficient, especially for large k."

Fortunately, that turns out not to be the case, as Opdyke (2022, 2023, and 2024a) derived an analytic, non-recursive expression of the cdf below in (26):

<sup>&</sup>lt;sup>44</sup> Note that even for estimation, the spherical (angles) parameterization of the covariance (correlation) matrix is a preferred choice. Per Pinheiro and Bates (1996): "Of the five parameterizations considered here, the spherical parameterization presents the best combination of performance and statistical interpretability of individual parameters."

(26)

$$F_{X}(x;k) \sim \frac{1}{2} - c_{k} \cdot \cos(x) \cdot {}_{2}F_{1}\left[\frac{1}{2}, \frac{1-k}{2}; \frac{3}{2}; \cos^{2}(x)\right] \text{ for } x < \frac{\pi}{2},$$
  
$$\sim \frac{1}{2} + c_{k} \cdot \cos(x) \cdot {}_{2}F_{1}\left[\frac{1}{2}, \frac{1-k}{2}; \frac{3}{2}; \cos^{2}(x)\right] \text{ for } x \ge \frac{\pi}{2}$$
  
where the Gaussian hypergeometric function  ${}_{2}F_{1}[a, b; c; r] = \sum_{n}^{\infty} \frac{(a)_{n}(b)_{n}}{(c)_{n}} \cdot \frac{r^{n}}{n!}$   
where  $(h)_{n} = h(h+1)(h+2)\cdots(h+n-1), n \ge 1, (h)_{0} = 1, \text{ and } |r| < 1, c \ne 0, -1, -2, ...$ 

As mentioned in a footnote above, the Gaussian hypergeometric function makes many interesting appearances in this setting, but it is admittedly cumbersome mathematically. Yet Opdyke (2022, 2023, and 2024a) has shown that (26) can be simplified further, based on some arguably obscure hypergeometric identities shown in (27) below:

(27) For c = a + 1, 0 < r < 1 simultaneously, which holds in this setting, then  ${}_{2}F_{1}[a,b;c;r] = B(r;a,1-b)(a/r^{a})$ 

where 
$$B(r;a,b) = \int_{0}^{r} u^{a-1} (1-u)^{b-1} du$$
 = the incomplete beta function (see DLMF, 2024)

In addition we have

$$F_{Beta}(r;a,b) = B(r;a,b)/B(a,b)$$
 where  $B(a,b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$  = the complete beta function, so

$$B(r;a,b) = F_{Beta}(r;a,b) \cdot B(a,b)$$
 (see Weisstein, E., 2024a and 2024b)

Combining terms we have

$$F_{X}(x;k) \sim \frac{1}{2} - c_{k} \cdot \cos(x) \cdot F_{Beta}\left[\cos^{2}(x);\frac{1}{2},\frac{1+k}{2}\right] \cdot \frac{\Gamma\left(\frac{1}{2}\right)\Gamma\left(\frac{1+k}{2}\right)}{\Gamma\left(\frac{2+k}{2}\right)} \cdot \left(\left[1/2\right]/\sqrt{\cos^{2}(x)}\right) \text{ for } x < \frac{\pi}{2},$$

$$F_{X}(x;k) \sim \frac{1}{2} + c_{k} \cdot \cos(x) \cdot F_{Beta}\left[\cos^{2}(x);\frac{1}{2},\frac{1+k}{2}\right] \cdot \frac{\Gamma\left(\frac{1}{2}\right)\Gamma\left(\frac{1+k}{2}\right)}{\Gamma\left(\frac{2+k}{2}\right)} \cdot \left(\left[1/2\right]/\sqrt{\cos^{2}(x)}\right) \text{ for } x \ge \frac{\pi}{2}$$

Recognizing that the complete Beta function is the inverse of the normalization factor of c(k) for these values, their product equals 1 and cancels, as do the two cosine terms, and we obtain the following signed beta cdf:

$$F_{X}(x;k) \sim \frac{1}{2} - \left(\frac{1}{2}\right) \cdot F_{Beta}\left[\cos^{2}\left(x\right); \frac{1}{2}, \frac{1+k}{2}\right] \text{ for } x < \frac{\pi}{2},$$

JD Opdyke, Chief Analytics Officer

**Correlation and Beyond** 

$$\sim \frac{1}{2} + \left(\frac{1}{2}\right) \cdot F_{Beta}\left[\cos^2\left(x\right); \frac{1}{2}, \frac{1+k}{2}\right] \text{ for } x \ge \frac{\pi}{2}$$

And now, with this straightforward, fully analytic, non-recursive cdf, we can obtain a straightforward, fully analytic quantile function of the angle distribution in (28):

(28) Let  $p = \Pr(x \ge X)$ . Then for  $x < \frac{\pi}{2}$ ,  $p = \frac{1}{2} - \left(\frac{1}{2}\right) \cdot F_{Beta} \left[\cos^{2}(x); \frac{1}{2}, \frac{1+k}{2}\right]$   $-2p = -1 + F_{Beta} \left[\cos^{2}(x); \frac{1}{2}, \frac{1+k}{2}\right]$   $1 - 2p = F_{Beta} \left[\cos^{2}(x); \frac{1}{2}, \frac{1+k}{2}\right]$   $F_{Beta}^{-1} \left(1 - 2p; \frac{1}{2}, \frac{1+k}{2}\right) = \cos^{2}(x)$   $\sqrt{F_{Beta}^{-1} \left(1 - 2p; \frac{1}{2}, \frac{1+k}{2}\right)} = \cos(x)$  $\operatorname{arcos} \left(\sqrt{F_{Beta}^{-1} \left(1 - 2p; \frac{1}{2}, \frac{1+k}{2}\right)}\right) = x$  (Not

(Note that arcos is arc-cosine, the inverse of the cosine function.)

We must reflect the symmetric angle density for p≥0.5, so we have

$$x = \arccos\left(\sqrt{F_{Beta}^{-1}\left(1 - 2p; \frac{1}{2}, \frac{1 + k}{2}\right)}\right) \text{ for } p < 0.5,$$
$$= \pi - \arccos\left(\sqrt{F_{Beta}^{-1}\left(1 - 2[1 - p]; \frac{1}{2}, \frac{1 + k}{2}\right)}\right) \text{ for } p \ge 0.5$$

Importantly, although often ignored in the sample-generation literature (see, for example, Makalic & Schmidt, 2018), note that properly adjusting for sample size, n, and degrees of freedom gives  $k \leftarrow k + n - \#cols - 2$ , so consequently, k = n - column # - 2.45

<sup>&</sup>lt;sup>45</sup> So notably, the angles distributions vary systematically based on their (column) position in the matrix, even though the distributions of the correlations themselves do not, as is discussed in later sections.

So now from (28) above we have for the angles distribution, under the Gaussian identity matrix, for the first time together, the pdf, cdf, and quantile function in (29):

(29)

$$f_X(x) = c_k \cdot \sin^k(x), x \in (0,\pi), k = 1,2,3... \# \text{columns} - 1, \text{ and } c_k = \frac{\Gamma(k/2+1)}{\sqrt{\pi}\Gamma(k/2+1/2)}$$

$$F_{X}(x;k) \sim \frac{1}{2} - \left(\frac{1}{2}\right) \cdot F_{Beta}\left[\cos^{2}(x);\frac{1}{2},\frac{1+k}{2}\right] \text{ for } x < \frac{\pi}{2},$$
  
$$\sim \frac{1}{2} + \left(\frac{1}{2}\right) \cdot F_{Beta}\left[\cos^{2}(x);\frac{1}{2},\frac{1+k}{2}\right] \text{ for } x \ge \frac{\pi}{2}$$
  
$$F_{X}^{-1}(p;k) = \arccos\left(\sqrt{F_{Beta}^{-1}\left(1-2p;\frac{1}{2},\frac{1+k}{2}\right)}\right) \text{ for } p < 0.5;$$
  
$$= \pi - \arccos\left(\sqrt{F_{Beta}^{-1}\left(1-2[1-p];\frac{1}{2},\frac{1+k}{2}\right)}\right) \text{ for } p \ge 0.5$$

Apparently the first (and only other) presentation of this quantile function result comes from an anonymous blog post in March, 2018, although it was obtained via a different derivation, which serves to further validate the result.<sup>46</sup>

The above (29) now provides a fully analytic solution,<sup>47</sup> and in fact is so straightforward as to be readily implemented in a spreadsheet, and one is provided for download via the link below.

http://www.datamineit.com/JD%20Opdyke--The%20Correlation%20Matrix-Analytically%20Derived%20Inference%20Under%20the%20Gaussian%20Identity%20Matrix--02-18-24.xlsx

So contrary to the assertions of Makalic & Schmidt (2018), the straightforward approach of inverse transform sampling *can* be used in this setting, for this narrow case, to very efficiently generate samples from the correlation matrix. And in fact, this is the most efficient way to sample it. Rubsamen (2023) has

<sup>&</sup>lt;sup>46</sup> See Xi'an, March, 2018 (<u>https://stats.stackexchange.com/questions/331253/draw-n-dimensional-uniform-sample-from-a-unit-n-1-sphere-defined-by-n-1-dime/331850#331850</u>

and <u>https://xianblog.wordpress.com/2018/03/08/uniform-on-the-sphere-or-not/</u>). In the interest of proper attribution, a reference on the website to the book "The Bayesian Choice" hints that the Xi'an pseudonym is Christian Robert, a professor of Statistics at Université Paris Dauphine (PSL), Paris, France, since 2000 (<u>https://stats.stackexchange.com/users/7224/xian</u>).

<sup>&</sup>lt;sup>47</sup> Note that I use the term 'analytic' as opposed to 'closed-form' because I am unaware of a closed-form algorithm for the inverse cdf of the beta distribution (see Sharma and Chakrabarty, 2017, and Askitis, 2017). However, for all practical purposes this is essentially a semantic distinction since this quantile function is hard-coded into all major statistical / econometric / mathematical programming languages.

compared Makalic and Schmidt (2018) to the above method and obtained over 30% decrease in runtime when using inverse transform sampling via (29). But of course, these results are instantaneous when used analytically (for example via the linked spreadsheet) as opposed to using inverse transform sampling, if, for example, only p-values and/or confidence intervals are needed, and generating many samples is not.

So sampling arguably is the less important of our two goals, because with a fully analytic finite-sample distribution, we can define, exactly for a given sample size, the p-value of a given cell, and the confidence interval of a given cell. The one-sided p-value simply is the cdf value for the lower tail, or [1 – (cdf value)] for the upper tail (30), and due to this pdf's symmetry, the two-sided p-value is simply two times either one-sided value. Correspondingly, the confidence interval for the critical value alpha is based on the quantile function as in (31).

(30) one-sided p-value =  $F_x(x;k)$  or  $1 - F_x(x;k)$  where k = n - column# - 2;

two-sided p-value = 2 x one-sided p-value

(31)  $F_X^{-1}(\alpha/2;k)$  and  $F_X^{-1}(1-\alpha/2;k)$  where, for a 95% confidence interval for example,  $\alpha = 0.05$ 

Notably, because the angles distributions are independent, the density of the entire matrix is simply the product of the densities of all the cells. This means we can readily define the p-value and confidence intervals of the entire matrix such that they are analytically consistent with those of the cells, because they are determined based directly on the cell level p-values and confidence intervals, respectively, as shown below.

## 4.b.iii. Matrix-level p-values and Confidence Intervals

As mentioned above, a key characteristic of the angles is that they are independent random variables, which makes defining their multivariate distribution straightforward: it is simply the product of all the angles' pdf's. More practically, the multivariate cdf is the product of every cell's cdf defined as (31a).

(31a) matrix cdf = 
$$F_{matrix} \begin{pmatrix} X_{i,j} \\ i > j \end{pmatrix}$$
 =  $\prod_{i>j} F_{X_{i,j}} (x;k)$  where  $X_{i,j}$  is the dependence measure of a

particular cell, i is row number, j is column number, and x is the value specified for that cell (and k = n – column# – 2, although this is determined by j)

But what does this mean for the p-value and confidence intervals for the entire matrix? Recall that a p-value is simply the probability of observing, based on a given data sample, a statistic value at least as extreme as what is observed, assuming the null hypothesis is true. The p-values defined in (31) above for each correlation/dependence cell are the probabilities of observing, for a given sample, angle values as large as what is observed assuming the null hypothesis is true. The fact that the angles are independent random variables, i.e. each is independent vis-à-vis all the other angles, makes obtaining the p-value for

## **Correlation and Beyond**

the entire matrix very straightforward. Note that the probability that <u>none</u> of the correlation/dependence cells are as extreme as what was observed is simply the product of one minus every p-value, because they are independent. So the probability that <u>one or more</u> of the correlation/dependence cells are as extreme or more extreme than what was observed is simply one minus this value, shown in (32) below, and this is the p-value for the entire matrix.

(32) matrix (2-sided) 
$$pvalue = \left[1 - \prod_{i=1}^{p(p-1)/2} (1 - p - value_i)\right]$$
 where  $p$ -value<sub>i</sub> is the 2-sided p-value.

Another way of conceptualizing this is that if the null hypothesis of just one of the matrix cells is not true, then the null hypothesis for the entire matrix is not true, and this is what the matrix-level p-value measures: the probability that <u>at least one</u> of the cell-level null hypotheses is not true. If instead of p-values we were using critical values in p(p-1)/2 hypothesis tests, this would be exactly consistent with controlling the familywise error rate (FWER) of the joint hypothesis including all the cells of the all-pairwise matrix.<sup>48</sup> And just as no other approach to estimating FWER would increase statistical power in this case due to the independence of the angles distributions,<sup>49</sup> no other definition of the matrix-level p-value will have greater power for the same reason: independence means there is no correlation structure to exploit to increase power. Of course, even though (32) provides a maximal power *p-value*,<sup>50</sup> it is arguably conservative when used as a test of the joint null hypothesis that there is no change in any of the cells of the matrix. To provide more power for this test, a false discovery rate (FDR) method could be used. FDR is the expected proportion of false discoveries among all the discoveries (i.e. among all the rejections of the null hypotheses of all the cells in the matrix). The goal is to control this proportion for a given α, and under our condition of multivariate independence here, the seminal Benjamini and Hochberg (1995) procedure provides a powerful and general-purpose FDR benchmark, as shown in (32a).

(32a) 
$$\alpha_{FDR} =: p\text{-}value_{(k)} \leq k\alpha/m$$

where  $\alpha$  = selected critical value; p-value<sub>(i)</sub> is from sorted p-values according to

p-value<sub>(1)</sub>  $\leq p$ -value<sub>(2)</sub>  $\leq \cdots \leq p$ -value<sub>(m)</sub>; m = p(1-p)/2 = the number of cells/tests; and k = the rank of the largest ordered p-value for which (32a) is true.

This controls the FDR at the level  $\alpha$ . For any of the m p-values, if p-value  $\leq \alpha_{FDR}$ , that hypothesis is rejected. So if any cells of the matrix have p-values equal to or less than  $\alpha_{FDR}$ , the joint hypothesis (for

<sup>&</sup>lt;sup>48</sup> Note that this approach has been used in the literature for addressing very closely related problems (see Fang et al., 2024).

<sup>&</sup>lt;sup>49</sup> Other approaches for calculating the FWER that rely on, for example, resampling methods (see Westfall and Young, 1993, and Romano and Wolf, 2016) exploit dependence structure to increase power; here, under independence, they would provide no power gain over the analogue to (32) because there is no dependence structure for them to exploit.

<sup>&</sup>lt;sup>50</sup> A maximal power p-value is the lowest p-value attainable for a given power level, significance level (alpha), and effect size.

the entire matrix) is rejected. In cases where k = m, all hypotheses are rejected and  $\alpha_{FDR} = \alpha$ , and when k = 0, i.e. (32a) is not true for any k ≤ m, no hypotheses are rejected and  $\alpha_{FDR} = 0$ . Both (32) and (32a) are readily calculated for matrix-level results, and I present both in the example results Section 6.

The calculation of confidence intervals for the entire matrix is essentially the same as the calculation of the p-value in (32), but of course the critical value is divided in half to account for each tail, and the root of one minus this critical value is calculated, rather than the product of one minus all the p-values. Otherwise, the calculations are identical, and (33) provides the critical alphas for these 'simultaneous confidence intervals.'

(33) 
$$\alpha_{crit-simult-LOW} = \left(1 - \left[1 - \alpha/2\right]^{\left(1/\left[p(p-1)/2\right]\right)}\right)$$
 and  $\alpha_{crit-simult-HIGH} = 1 - \alpha_{crit-simult-LOW}$ 

These critical alphas, when inserted in the quantile function (28) and applied to every cell, provide the two correlation matrices that define and capture, say,  $(1-\alpha) = (1-0.05) = 95\%$  of randomly sampled matrices under the null hypothesis, which up to this point has been the identity matrix. Again, it is the independence of the angles that makes these simultaneous confidence intervals very straightforward to calculate.

Importantly, again note that because we derived the quantile (inverse cdf) function in (28) and (29) above, we can go in either direction regarding these results: we can specify a correlation matrix and, under the null hypothesis (of the identity matrix), obtain its p-values, both for the individual cells and the entire matrix, simultaneously. We also can specify a matrix of cdf values and obtain its corresponding correlation matrix, which is extremely useful and straightforward when constructing both stress and reverse stress scenarios. Finally, we can use simultaneous confidence intervals to obtain the two correlation matrices that form the matrix-level confidence interval.

Note that all these calculations are included in the downloadable spreadsheet (link provided above), with visible formulae corresponding to each step of these calculations for full transparency. In the next section below I expand these results for Pearson's to apply to all data conditions, and all values of the null hypothesis (i.e. any values for the matrix, not just the identity matrix).

## 4.c. NAbC: Pearson's Correlation, Real-World Financial Data, Any Matrix Values

Currently, the extant literature does not provide analytic forms for the angles distributions under general conditions. Deriving these appears to be a non-trivial problem. Spectral (eigenvalue) distributions, which many researchers turn to in this setting, have been shown to vary dramatically when data is characterized by different degrees of heavy-tailedness (see Burda et al., 2004, Burda et al., 2006, Akemann et al., 2009; Abul-Magd et al., 2009, Bouchaud & Potters, 2015, Martin & Mahoney, 2018; Heiny and Yao, 2022, and Opdyke, 2022), as well as by different degrees of serial correlation (see Burda et al., 2004, 2011, Hisakado and Kaneko, 2023, and Opdyke, 2022), and the literature provides no general

analytic form under general, real-world financial data conditions – certainly nothing that is analogous to convergence to the Marchenko-Pastur distribution under iid conditions (see Marchenko and Pastur, 1967).<sup>51</sup> If angles distributions are of similar complexity, then deriving their analytic form under general conditions, if possible, currently appears to be a non-trivial, unsolved problem.

However, this should not be (and does not need to be) a showstopper for our purposes, in part because angles distributions have many characteristics that make them useful here generally, and more useful specifically than spectral distributions, by multiple criteria: empirically, distributionally, and structurally.

<u>Empirically</u>: If an angle distribution approaches degeneracy, most of its values typically will approach 0 or **π**. But the relevant trigonometric functions (sin and cos) of these values are stable, and will simply approach -1, 0, or 1. This makes  $R = BB^T$  a relatively stable calculation empirically, even if it produces an *R* matrix that is approaching non-positive definiteness (NPD). In contrast, eigenvalue/vector estimations are more numerically involved compared to the application of simple trigonometric functions, and this, combined with the fact that empirically, their upper bound is not well-bounded (in the general case),<sup>52</sup> makes their computation comparatively less numerically stable as matrices approach NPD.

<u>Distributionally</u>: As shown graphically below under challenging, real-world financial data conditions, the angles distributions are relatively "well behaved," both in the general sense and relative to spectral

distributions. They are relatively smooth and typically unimodal, and clearly bounded on  $\theta \in (0, \pi)$ . Spectral distributions, based on the same data, very often are spikey,<sup>53</sup> multimodal, and for practical purposes, empirically unbounded (at least in higher dimensions), all of which translates into larger variances and less tail accuracy. Simply put, they typically are much more complex and challenging to estimate precisely and accurately compared to individual angles distributions for a given correlation/dependence matrix *R* under real-world financial data. At least part of the reason for this is the much larger number of cells that need to be estimated compared to the relatively few eigenvalues that need to be estimated, which leads to more structural stability of the former when all are combined into a unified, estimated matrix.

<u>Structurally</u>: Aggregation level becomes relevant and important here. For a given correlation/dependence matrix R there are many more angles than eigenvalues (i.e. p(p-1)/2 cells vs p eigenvalues, for a factor of

<sup>&</sup>lt;sup>51</sup> Note that some exceptions to convergence to this celebrated distribution do exist (see Li and Yao (2018), Hisakado and Kaneko (2023), Heiny and Yao (2022), and Maltsev and Malysheva (2024) for examples).

<sup>&</sup>lt;sup>52</sup> Even though the largest eigenvalue is known to follow the Tracy-Widom distribution under certain sets of conditions, under others it can diverge, with unbounded support (see for example Li, 2025). Even when the latter cases do not hold mathematically, in practice, empirically, the largest eigenvalue can become so large that it is essentially unbounded.

<sup>&</sup>lt;sup>53</sup> In fact, one of the most commonly encountered covariance (correlation) matrices under real world financial data conditions is the spiked matrix (see Johnstone, 2001), where one or few eigenvalues dominate and the majority of eigenvalues are close to zero, i.e. not reliably estimated. This further demonstrates that spectral approaches are far too limited and limiting to effectively solve this problem under real-world conditions.

(p-1)/2 more). By capturing the pairwise dependence structures of p(p-1)/2 cells using only p eigenvalues, we unavoidably increase the complexity of each eigenvalue distribution, at least compared to those of the angles which can more accurately reflect each pairwise association. Consequently, as a matrix *approaches* singularity (NPD), which arguably is the rule rather than the exception for non-small investment portfolios, a much smaller *proportion* of angles distributions will approach degeneracy (i.e. minimum/maximum values of zero and  $\pi$ ) than is true for eigenvalue distributions (where more values will wrongly fall below zero). The overall construction, then, of the correlation matrix based on angle estimates via  $R = BB^T$  generally will remain much more stable than one based on eigenvalue estimates using a decomposition of  $R = VAV^{-1}$  where V is a matrix with column eigenvectors and A is a diagonal matrix of the corresponding eigenvalues.

All of this adds up to a more robust and granular basis for inference and analysis when relying on angles distributions as opposed to spectral distributions. The underlying reason for this is the fact that spectral distributions simply are at the wrong level of aggregation for these purposes: they remain at the (higher) level of the p assets of a portfolio – NOT at the granular level of the p(p-1)/2 pairwise associations of that portfolio, which is where both the angles distributions, and those of the correlations/dependence measure values themselves, lie. Consequently, while potentially very useful for things like portfolio factor analysis, spectral analysis simply is too blunt a tool for our purposes here. We need to be able to make inferences and monitor processes and conduct (reverse) scenario analyses and customized stress tests on ALL aspects of the dependence structure measured by the all-pairwise correlation/dependence matrix, at the granular level at which it is defined. The specific need for this in scenario and reverse scenario analyses is covered in more detail below.

So given the useful characteristics of the angles distributions, not to mention the fact that they remain at the right level of aggregation for granular analysis of the correlation/dependence matrix, we are able to proceed WITHOUT their analytic derivation. Rather, we can use a time-tested nonparametric approach, such as kernel estimation, to reliably define them. All this requires is a single simulation (say, N=10,000) based on the known or well-estimated values of the correlation/dependence matrix, and its known or well-estimated data generating mechanism. These are the two stated requirements for the application of NAbC under general conditions. Then, after translating all N simulated correlation matrices to N matrices of angles, we fit a kernel to each empirical angle distribution, i.e. the empirical distribution of each angle for each cell of the matrix. We now have not only the distributions of all the individual angles, but also the multivariate distribution of the matrix, which is just the product of all the individual distributions due to their independence. Note that this goes in both directions: we can perform 'lookups' on the empirically defined distribution to obtain the cdf value(s) corresponding to particular angle value(s), or use cdf value(s) to 'look up' corresponding angle (quantile) value(s). The subsequent kernel fitting smooths this empirical density to all (continuous) values, and sampling readily can be performed using these kernel densities. This process is described step by step as in Section 2.c (but with more brevity here as the steps are explained in more detail above).

## 5 Steps for Obtaining Angles Distributions

- 1. Simulate N samples (N=10,000 typically is sufficient) based on the dependence matrix and the data generating mechanism (each can either be specified/known, or well estimated).
- 2. Calculate the corresponding N all-pairwise dependence matrices, and their Cholesky factorizations, and transform each of these factorizations into a lower triangle matrix of angles.
- 3. Fit a kernel density to each cell of the matrix of angles based on the N values obtained from the N samples in 2.
- 4. Generate N samples based on the kernel densities in  $3.^{54}$
- 5. Convert each of the N samples from 4. back to a re-parameterized Cholesky factorization, and then multiply it by its transpose to obtain a set of N validly sampled dependence matrices. Positive definiteness is enforced automatically as the Cholesky factor places us on the <u>unit</u> hyperhemisphere. All sample generation hereafter uses just 4. and 5.

The samples of correlation/dependence measure matrices from 5. will follow the same distribution as those generated in 2., but after the kernel densities are fit once in 3., generating samples based on 4. and 5. is orders of magnitude faster than relying on direct simulations in steps 1. and 2. So one simulation gives us the distribution of each and every angle, corresponding to each and every correlation/dependence cell. And now going forward using 4.-5., rather than 1. and 2., allows for correct probabilistic inference, both at the cell level and at the matrix level, due to the independence of the angles distributions (remember the correlations themselves are NOT independent, so 1. and 2. provide no direct inferential capability). This reliance on the angles, and their subsequent transformation to correlations, allows us to isolate specifically the distribution of the entire correlation/dependence matrix, for probabilistic inference, without touching any other distributional aspect of the data, which is the point of the methodology. Of course, either a direct data simulation (step 1. above) or a cavalier 'bootstrap' of the matrices calculated based on step 1. fails at this objective, because the non-independence of the correlation cells undermines the validity of any empirically-based inference based on simple metrics (e.g. distances) across the group of sample matrices. In other words, a group of sample correlation matrices based on simulated data does not provide any inferential capabilities about the correlations, but a group of matrices based on simulated angles does.

So this framework is essentially identical to that for the specific case of the Gaussian identity matrix, with the only difference being it is based on nonparametrically defined, as opposed to parametrically defined, angles distributions. Before covering implementation details below, I show some examples of graphs of the angles distributions and the corresponding spectral distribution under challenging, simulated financial returns data (these are all generated based on the 5 Steps above). The multivariate returns distribution of the portfolio is generated based on the t-copula of Church (2012), with p=5 assets, varying degrees of heavy-tailedness (df=3, 4, 5, 6, 7), skewness (asymmetry parameter=1, 0.6, 0, -0.6, -1), non-stationarity (standard deviation= $3\sigma$ ,  $\sigma/3$ , and  $\sigma$ , each with n/3 observations), and serial correlation (AR1=-

<sup>&</sup>lt;sup>54</sup> Algorithms for sample generation based on commonly used kernels (e.g. the Gaussian and Epanechnikov) are widely known. An example of the latter is simply the median of three uniform random variates (see Qin and Wei-Min, 2024).

0.25, 0, 0.25, 0.50, 0.75), with a block correlation structure shown in (34) below and n=126 observations for a half year of daily returns.<sup>55</sup> The spectral distribution is compared against Marchenko-Pastur as a referential baseline that assumes independence (and identically distributed asset returns).

	1	-0.3	-0.3	0.2	0.2
	-0.3	1	-0.3	0.2	0.2
	-0.3	-0.3	1	0.2	0.2
	0.2	0.2	0.2	1	0.7
(34)	0.2	0.2	0.2	0.7	1

Several points are worth noting and reemphasizing based on these graphs. First, the graphs of the angles distributions, all of which are based on the 5 Steps above, contain three densities: A. one based on angles perturbation (i.e. sampling from the fitted kernel) as described above in steps 3.-4., B. one based on direct data simulations (steps 1.-2.), and C. the analytical density under the (Gaussian) identity matrix as a comparative baseline. Note that the only reason I include B. is to demonstrate the validity of A., and as expected, the angles distributions from A. and B. are empirically identical (with A. being orders of magnitude faster and more computationally efficient, not to mention providing a basis for valid inference). The spectral distributions based on the samples generated in both A. and B. also are identical, as are a wide range of additional aggregated metrics not presented herein (e.g. various norms, VaR-based economic capital, and 'generalized entropy' as described in a later section below). This

## Graph 1a:

# Spectral Distribution Based on i. NAbC Angles Kernel (Step 4) vs. ii. Data Simulations (samples from Step 2) vs. Marchenko Pastur Distribution



<sup>&</sup>lt;sup>55</sup> Note that this is only approximately Church's (2012) copula, which incorporates varying degrees of freedom (heavytailedness) and asymmetry, because I also impose ex post serial correlation and non-stationarity on the data (and subsequently rescale the marginal densities).

**Correlation and Beyond** 

#### Graphs 1-10:

Angles Distributions Based on i. NAbC Angles Kernel (Step 4) vs. ii. Data Simulations (samples from Step 2) vs. iii. Gaussian Identity Matrix



JD Opdyke, Chief Analytics Officer

**Correlation and Beyond** 

Page 52 of 92



empirically validates that the angles/kernel perturbation approach (steps 3.-5.) is an effective and correct method for isolating and generating the distribution of the correlation/dependence matrix, and unlike steps 1. and 2., one that preserves inferential capabilities. In other words, these results empirically validate that the angles contain all extant information regarding dependence structure (see Fernandez-Duren & Gregorio-Dominguez, 2023, and Zhang & Yang, 2023, as well as Opdyke, 2022).

Second, note again that a nonparametric approach works in practice here at least in part because the angles distributions are 'well behaved.' Since they are relatively smooth, typically if not always unimodal, and well bounded, N=10,000 simulations typically suffices to provide a precise and accurate measure of their distributions. Poorly behaved distributions that are very spikey, multi-modal, and essentially unbounded for all empirical, practical purposes could require numbers of simulations orders of magnitude larger. If N=10,000,000 or even 1,000,000 for example, this approach could be computationally prohibitive in many cases for real-world-sized portfolios, which often exceed p=100 with p(p-1)/2=4,950 pairwise associations/cells.

Finally, as described above, note the multi-modal and long tailed / high-upper-bounded nature of the spectral distribution for this portfolio compared to the angles distributions, where the biggest thing approaching an estimation challenge is a modest asymmetry. But this speaks only to estimation issues. More notable is the fact that on a cell-by-cell basis, the angles distributions deviate materially i. not only from central values of  $\pi/2$ , and less dramatically from perfect symmetry when compared to their (analytic) distributions under the (Gaussian) identity matrix, but also ii. from each other! Each angle's distribution can vary quite notably compared to the other angles' distributions, especially under smaller sample sizes. There simply is no way that one spectral distribution for a matrix, or even p distributions for each eigenvalue individually and even if perfectly estimated, will be able to capture and reflect all the richness of dependence structure reflected here at the granular level of all the p(p-1)/2 pairwise cells. This remains true regardless of their use in this setting, whether for cell-level attribution analyses, granular scenario and reverse scenario analyses, cell-level intervention 'what if' analyses, or customized

stress testing, let alone precise and correct inference at either the cell level OR the matrix level. I now leave comparisons to spectral distributions behind<sup>56</sup> to cover implementation issues below.

## 4.c.i. Nonparametric Kernel Estimation

Due to the bounded nature of the angles distributions on  $\theta \in (0, \pi)$ , the nonparametric kernel must be appropriately reflected at the boundary (see Silverman, 1986) via:

if  $\theta < 0$  then  $\theta \leftarrow -\theta$ ; if  $\theta > \pi$  then  $\theta \leftarrow (2\pi - \theta)$ , which is asymptotically valid.

As per the standard implementation, the kernel itself is defined as

(35) 
$$f_{h}(\theta) = \frac{1}{N} \sum_{i=1}^{N} K_{h}(\theta - \theta_{i}) = \frac{1}{hN} \sum_{i=1}^{N} K_{h}([\theta - \theta_{i}]/h) \text{ with}$$
  
Gaussian:  $K(\theta) = (1/\sqrt{2\pi}) \cdot e^{-\theta^{2}/2}$  Epanechnikov:  $K(\theta) = (3/4) \cdot (1 - \theta^{2}), |\theta| \le 1$ 

I have tested both the Gaussian and the Epanechnikov kernels extensively in this setting,<sup>57</sup> along with three different bandwidth estimators, *h*, from Silverman (1986) and one from Hansen (2014),

respectively:  $h = 1.06 \cdot \hat{\sigma} \cdot N^{-1/5}$ ,  $h = 0.79 \cdot \text{IQR} \cdot N^{-1/5}$ ,  $h = 0.9 \cdot \min(\text{IQR}/1.34, \hat{\sigma}) \cdot N^{-1/5}$ , and lastly,  $h = 2.34 \cdot \hat{\sigma} \cdot N^{-1/5}$  for Epanechnikov only, where  $\hat{\sigma}$  = sample standard deviation and IQR = sample interquartile range.

As with virtually all kernel implementations, the choice of kernel matters less than the choice of bandwidth, although in this setting, across a broad range of data conditions and correlation/dependence values, the Epanechnikov kernel appears to perform slightly 'better' (i.e. with slightly less variance, thus providing slightly more statistical power) than the Gaussian, perhaps because its sharp bounds require reflection at the boundary less often than the Gaussian kernel (although reflection at the boundary is quite uncommon, even for 'extreme' dependence matrices). The bandwidth that appears to perform best across wide-ranging conditions is  $h = 0.9 \cdot \min(IQR/1.34, \hat{\sigma}) \cdot N^{-1/5}$ . Additionally, for larger matrices (e.g. p=100), bandwidths need to be tightened by multiplying *h* by a factor of 0.15. When there are many cells (e.g. for p=100, #cells=p(p-1)/2=4,950) this tightening avoids a slight drift in metrics that are aggregated across all the cells (e.g. correlation matrix norms, spectral distributions, and LNP (a type of 'generalized entropy' defined in a later section below)). Multiplying by this factor for smaller matrices does not adversely affect the density estimation in any way, so this factor always is used. For matrices

<sup>&</sup>lt;sup>56</sup> Continued reliance on spectral approaches for this specific problem brings to mind a quotation attributed to John M. Keynes: "the difficulty lies not so much in developing new ideas as in escaping from old ones."

<sup>&</sup>lt;sup>57</sup> Note that the Epanechnikov kernel is used in very closely related problems in this setting (see Burda and Jarosz, 2022). JD Opdyke, Chief Analytics Officer Page 54 of 92 Correlation and Beyond

much larger than p=100, a further tightening of this factor may be required, and this is readily determined by empirically comparing the distributions of these aggregated metrics under direct data simulation (steps 1. and 2.) vs. NAbC's kernel-based sampling (steps 3., 4. and 5.).

Once the kernels have been estimated and the angles distributions generated by perturbing/sampling based on those kernels, the p-values and confidence intervals for both the individual correlation/dependence cells and the entire correlation/dependence matrix are the same as those derived for the Gaussian identity matrix. The only difference, aside from their now-nonparametric basis, is that the angles distributions are no longer symmetric by definition, as is true under the (Gaussian) identity matrix. This can be seen in the Graphs 1-10 of the angles distributions provided above. The p-value calculation, however, remains very straightforward, and it requires just a bit of care to properly account for asymmetry. The one-sided p-value remains simply as in (30), shown in (36) below:

(36) one-sided p-value =  $F_X(x;k)$  or  $1 - F_X(x;k)$  for lower and upper tails, respectively, where k = n - column# - 2

However, due to possible (probable) asymmetry, the two-sided p-value differs, and is not just two times the one-sided p-value as in (30), requiring first the calculation of the empirical mean correlation matrix from the simulations in step 2. of the five sampling steps above. This mean correlation matrix is then translated into a matrix of angles, and we obtain the empirical cdf values corresponding to these "mean angles" with a "look-up" on the entire angles distributions generated in step 4. These cdf's will be close to 0.5 when the angles distributions are close to symmetry, and they will deviate from 0.5 under asymmetry, and will serve as the baseline off of which the two-sided p-values are calculated. Specifically, the difference between the cdf values of each of the angles of the specified correlation matrix being 'tested,' and those of the corresponding "mean angles," defines the two-sided p-values, which are simply the sum of the probability in the tails BEYOND this difference. Formulaically this is shown in (37):

(37) two-sided p-value = max[0, Mcdf – d] + max[0, 1 – (Mcdf + d)], where d = abs(Mcdf – cdf), Mcdf = mean angle cdf, cdf = cdf of specified angle

This usually results in summing both tails, but under notable asymmetry, sometimes only one tail is used. Below is a graphical example of both cases, where the cdf of the "mean angle" is 0.6 and the cdf of the relevant angle in the specified correlation matrix (i.e. the correlation matrix for which we are obtaining pvalues) is cdf=0.1 in the single-tail case (Graph 11) and cdf=0.85 in the double-tail case (Graph 12). In the statistical sense, however, both cases remain two-sided p-values.

Graph 11: p-value for a specified angle with a more extreme cdf



#### Graph 12: p-value for a specified angle with a less extreme cdf



Note that while cdf=0.1 is hardly more 'extreme' than cdf=0.85 in absolute terms, relative to the mean angle cdf=0.6, it is twice as 'extreme,' i.e. twice as far from the mean cdf=0.6 with a distance of 0.5 for Graph 11, but a distance of only 0.25 for Graph 12. Moreover, the tail probability of Graph 11 (0.1) is only 1/5 that of Graph 12 (0.5) (compare the red shaded areas). This example demonstrates why asymmetry must be properly taken into account in this setting, but the two-sided p-value still remains a very straightforward calculation, and the "mean angles" matrix is used for additional, important purposes below, as discussed in Section 5.

Cell-level confidence intervals still are simply calculated as in (31), which automatically takes asymmetry into account as we are using the empirically fitted kernel. This is identical to the same calculation under the (Gaussian) identity matrix (sans the known symmetry of the cdf). And the matrixlevel p-value, again, is simply one minus the probability of observing the sample matrix that was observed, or one 'less extreme,' exactly as in (32). This, again, is analogous to the definition of controlling the family-wise error rate (FWER) where all the cells of the matrix comprise the joint null hypothesis. Finally, just as under the (Gaussian) identity matrix, calculation of the confidence interval for the entire matrix remains (33) as previously.

Importantly, again note that we can go in either direction regarding these results: we can specify a correlation/dependence matrix and, under the null hypothesis of the specified correlation matrix, obtain the p-values of an observed matrix, both for the individual cells and the entire matrix, simultaneously. We also have the matrix-level quantile function: we can specify a matrix of cdf values and obtain its corresponding, unique correlation/dependence matrix, which can be extremely useful and straightforward when constructing reverse (stress) scenarios. Finally, we can use simultaneous confidence intervals to obtain the two correlation matrices that form the matrix level confidence interval. An example with all these results is shown in Section 6 below, but first I extend NAbC's range of application beyond Pearson's to all dependence measures with positive definite all-pairwise matrices.

#### 4.d. NAbC: Any (Positive Definite) Dependence Measure, Any Data, Any Matrix Values

First, I reemphasize here that that the relationship between spherical angles and Pearson's matrix shown in (21) and (22) above holds for any positive definite matrix (see Joarder and Ali, 1992, Pinheiro and Bates, JD Opdyke, Chief Analytics Officer Page 56 of 92 Correlation and Beyond

1996; Rebonato and Jackel, 2000; Rapisarda et al., 2007; Pouramadi and Wang, 2015; Cordoba et al., 2018; and Lan et al., 2020), and so applies to all other positive definite dependence measures discussed herein. The Cholesky factor, defined in terms of these angles in (21) and (22), remains undefined under non-positive definiteness, and because NAbC relies on this formulation of the Cholesky factor, its application requires the positive definiteness of the matrix to which it is being applied. As mentioned above, for long-established dependence measures like Pearson's, Kendall's, Spearman's, and the tail dependence matrix, positive definiteness has been proven analytically (see Sabato, 2007, for the first three, and Embrechts et al., 2016, for the last). But when analyzing their empirical results in any given simulation or estimation, positive definiteness always should be verified, since for matrices that approach non-positive definiteness (which is not uncommon for non-small financial portfolios), their corresponding sample-based empirical estimates can sometimes be non-positive definite due strictly to numerical issues. Consequently, verifying the positive definiteness of estimates of those measures for which analytical proofs have not (yet) been derived does not require any special treatment: best practices dictate that positive definiteness always is verified empirically for all sample-based estimates of these matrices, regardless of the dependence measure being used. Notably, in the many millions of simulations conducted when testing and developing NAbC, none of the 'newer' measures<sup>58</sup> generated empirically non-positive definite estimates. This admittedly is driven by the range of data conditions being examined, and so scientific caution remains warranted, and positive definiteness should never be presumed in the absence of analytical proof. Yet the bottom line remains: as long as a measure's matrix is positive definite, NAbC is applicable and will 'work' to provide inferentially valid simulated samples, and/or inferential statistics related to its sample-based estimate. And all of the 'newer' measures tested herein always remained positive definite under a very broad range of challenging, real-world data conditions (see Opdyke, 2022, 2023, and 2024a).

But to flip this script on this requirement of positive definiteness, one reasonably could argue that if a dependence measure was analytically shown to be non-positive definite, at least under relevant conditions, and/or its empirical estimates were non-positive definite more often than could be attributable solely to numeric considerations, then researchers and practitioners might want to question the wisdom of using it. Non-positive definiteness also could be a function of unknown or mis-specified data conditions, such as perfect linear dependence unwittingly built into a simulation (although with actual market data, it could be a very useful flag for extreme multicollinearity). Or non-positive definiteness perhaps could be due to a combination of the dependence measure used and the specific data conditions being examined. Either way, the non-positive definite results could be serving as a

<sup>&</sup>lt;sup>58</sup> The 'newer' measures on which NAbC has been tested include most of those listed in the Introduction and Background section, such as Chatterjee's correlation (both the asymmetric and symmetric versions), the 'improved' Chatterjee's correlation of Xia et al. (2024) (both the asymmetric and symmetric versions), the combination of Chatterjee's+Speaman's due to Zhang (2024a) (both the asymmetric and symmetric versions), Szekely's (2007) distance correlation, the asymmetric tail dependence measure due to Diedda et al. (2023), both the symmetric and asymmetric versions of Liu and Shang's (2025) DDC method, and both Lancaster's correlation and Lancaster's linear correlation (see Holzmann and Klar, 2024). Including the 'big 3' and the tail dependence matrix (upper and lower tails counted separately), this makes a total of 19 dependence measures to which NAbC has been successfully applied.

correct and useful warning to avoid the dependence measure (and/or those simulated data conditions) altogether. In such cases, the requirement of positive definiteness is less a limitation of a method like NAbC and more a proper boundary on the right measures and conditions under which such analyses should be conducted in the first place.

On a separate but related note, the ranges of many of the 'newer' dependence measures (e.g. Szekely's, Lancaster's, and Chatterjee's) are (0, 1) instead of (–1, 1) like the big 3, but operationally, implementing NAbC on these measures does not change, even as it relates to how we reflect at the boundary when fitting the nonparametric kernel. This is because specific cells of the Cholesky factor can validly be negative, making the assignation in the last line of the "Correlations to Angles" code in Table A above sometimes assign an angle value slightly above  $\pi/2$ , even though  $\pi/2$  corresponds to a measure value of zero.<sup>59</sup> So this is a soft upper boundary in this case, even though the measure's range of (0,1) typically is not.<sup>60</sup> So when NAbC generates angle  $\theta$ , we continue to reflect based on:

if  $\theta < 0$  then  $\theta \leftarrow -\theta$ ; if  $\theta > \pi$  then  $\theta \leftarrow (2\pi - \theta)$ , since for measures with a (0,1) range, the upper bound of  $\pi$  will never be reached, and the lower bound of zero remains valid and hard. So NAbC applies in exactly the same way, for all of these dependence measures, whether their range of values is (-1, 1) or (0, 1).

Finally, again note that the condition of symmetric positive definiteness holds not only for all relevant dependence measures, as shown above, but also under all relevant real-world data conditions: that is, multivariate financial returns data whose marginal distributions typically are characterized by varying and different degrees of asymmetry, heavy-tailedness, non-stationarity, and serial correlation. So this is a very weak and general condition, allowing for the extremely wide-ranging application of NAbC.

## 4.d.i. Spectral and Angles Distributions, Examples from Other Dependence Measures

I present below graphs of the spectral and angles distributions for some of the dependence measures discussed above, beyond Pearson's, under simulated data reflecting challenging, real-world data conditions (see Opdyke, 2022, for the application of NAbC to a wide range of data conditions). As in the above example, the multivariate returns distribution of the simulated portfolio is generated based on the t-copula of Church (2012), with p=5 assets, varying degrees of heavy-tailedness (df=3, 4, 5, 6, 7), skewness (asymmetry parameter=1, 0.6, 0, -0.6, -1), non-stationarity (standard deviation= $3\sigma$ ,  $\sigma/3$ , and  $\sigma$ , each with n/3 observations), and serial correlation (AR1=-0.25, 0, 0.25, 0.50, 0.75), with a block

<sup>&</sup>lt;sup>59</sup> Note that angle values (which range from zero to  $\pi$  on the hyper-hemisphere) decrease while dependence measure values increase, so a measure value of -1 corresponds to an angle value of  $\pi$ , a measure value of zero corresponds to an angle value of  $\pi$ /2, and a measure value of 1 corresponds to an angle value of zero (see Zhang et al., 2015 and Lu et al., 2019).

<sup>&</sup>lt;sup>60</sup> On a related issue, note that Chatterjee's correlation, for example, is bounded by (0,1) only asymptotically, and finite sample results can exceed these bounds. However, when applying NAbC to this and other measures in millions of data simulations under widely varying conditions, as an empirical matter such finite sample exceedences never caused NAbC's angles distributions to deviate from those of direct data simulations, nor did they ever make empirical matrices non-positive definite.

correlation structure shown in (34) below and n=126 observations, for half a year of daily returns.<sup>61</sup>

	1	-0.3	-0.3	0.2	0.2
	-0.3	1	-0.3	0.2	0.2
	-0.3	-0.3	1	0.2	0.2
	0.2	0.2	0.2	1	0.7
(34)	0.2	0.2	0.2	0.7	1

For verification purposes only, I compare those angles distributions based on the data simulation directly against those based on NAbC's kernels, and in all cases the results are empirically indistinguishable. The same is true for the spectral distributions, which I also present below against the Marchenko-Pastur distribution as a referential baseline (that presumes iid data; see Marchenko and Pastur, 1967). The empirical results yield some expected, and some additional interesting findings.

First, note that the spread, and the spread and shifts, of both the spectral and angles distributions, respectively, are larger for Pearson's than for Kendall's, which is consistent with the former's relative sensitivity to more extreme values under many conditions. The shifts and spread of both measures are much larger than those of Chatterjee,<sup>62</sup> although this is largely due to the fact that while Chatterjee is generally more powerful under dependence that is cyclical or non-monotonic in some way, it is less powerful under associations that are more monotonic, and the data conditions of this example fall more (but not entirely) into the latter category. The story changes a bit when we use the dependence measure suggested by Zhang (2024a), which is essentially a maximum between Spearman's rho and Chatterjee's correlation, its objective being to obtain large, if not maximum power under both types of dependence structures (i.e. strong monotonic dependence as well as cyclical or otherwise non-monotonic dependence). This shows how readily NAbC can be applied to any (positive definite) dependence measure, and its utility for making cross-measure comparisons, all else equal, using the same, universally applicable method.

Before providing a complete example of NAbC's application below, i.e. one that provides both matrix and cell level p-values and confidence intervals, and checks all (but one) of the original objectives listed in the Introduction and Background, I discuss three of its additional and important capabilities: the first is its use as a two-sample test comparing two correlation/dependence matrices; the second is that fact that it remains "estimator agnostic"; and the third is its capability to perform fully flexible scenario analytics, providing granular, realistic scenario analytics far beyond what its competitors can provide.

<sup>&</sup>lt;sup>61</sup> Note again that this is only approximately Church's (2012) copula, which incorporates varying degrees of freedom (heavy-tailedness) and asymmetry, because I also impose serial correlation and non-stationarity ex post on the data (and subsequently rescale the marginal densities).

<sup>&</sup>lt;sup>62</sup> The symmetric version of Chatterjee's correlation coefficient is used here (see Chatterjee, 2021), with the finite sample bias correction proposed by Dalitz et al. (2024).





Graphs 13b: Angles Distributions, by Dependence Measure by Cell, Based on i. NAbC Angles Kernels vs. ii. Direct Data Simulations vs. iii. Identity Matrix

Page 61 of 92

**Correlation and Beyond** 

JD Opdyke, Chief Analytics Officer

## Graph 14: Spectral Distribution, by Dependence Measure, Based on i. NAbC Angles Kernel vs. ii. Direct Data Simulations vs. iii. Marchenko Pastur



Kendall's Tau



#### 4.e. NAbC: Fully General Conditions, Statistical Comparison of Two Matrices

The above development of NAbC has so far covered only hypothesis tests against a matrix of fixed values, i.e. a one-sample test. But the NAbC approach allows us to perform two-sample tests of one sampled matrix against another sampled matrix, say, from two different sectors or two different business lines, where the null hypothesis is no difference between the dependence structures of the two sectors. The implementation is very similar to the one sample case, except that N=10,000 samples based on the estimates of each of the two angle matrices are generated separately. Then the differences between the two groups of samples of angles (sample i from a particular cell of the first matrix minus sample i from the same cell of the second matrix) are calculated, and the N differences are tested against the values of the identity matrix, i.e. values of zero representing zero difference between the two matrices, similar to testing a single sample against the null hypothesis of the identity matrix. The only difference is that we must use the "mean angles" cdfs to account for asymmetry in the angles distributions slightly differently: instead of averaging the correlation matrices and then converting this average matrix into the matrix of "mean angles" cdfs (as described above in step 2. of the kernel-based angles simulations), we first take

JD Opdyke, Chief Analytics Officer

Page **62** of **92** 

#### **Correlation and Beyond**

the difference between angles and then calculate the mean of these differences before obtaining the corresponding cdf values for each cell. This avoids incorporating into our 'accounting for asymmetry' what are possibly true differences between the two matrices, i.e. the hypothesis we are testing. Otherwise, the approach is exactly the same as the one-sample test. The only obvious constraints on this approach are that the two matrices being compared should be the same type of dependence measure (e.g. Szekely's vs. Szekely's, not Szekely's vs. Chatterjee's) and have the same dimension.<sup>63</sup> I include in Section 6 below an empirical example of such a two-sample test, under both unrestricted and scenario-restricted conditions. Extending this approach to the multi-sample case is addressed in more detail in future research. Below, I briefly describe how NAbC remains "estimator agnostic" before moving on to NAbC's application to fully flexible scenarios.

## 4.f. NAbC Remains "Estimator Agnostic"

As mentioned briefly in the Introduction and Background, another important and useful characteristic of NAbC is that it remains "estimator agnostic," that is, valid for use with any reasonable estimator of any of the dependence measures being utilized. Different estimators will have different characteristics under different data conditions. For example, some will provide minimum variance / maximum power, while others may provide unbiasedness or less bias, while others may provide more robustness, and/or different combinations of these characteristics under different data conditions. Ideally, we would like to be able to use estimators that provide the best trade-offs for our purposes under the conditions most relevant to our given portfolio. Fortunately, NAbC "works" for any estimator, as the relationship between correlations/dependence measure values and angles requires only symmetric positive definiteness. NAbC's finite sample distribution and its resulting inferences obviously will inherit the advantages and disadvantages of the estimator being used, but this is generally an advantage as it provides flexibility to use the 'best' estimator under the widest possible range of conditions. And the ability to apply NAbC as a single, unifying method across very wide-ranging data conditions is what allows for very effective ceteris paribus analyses that otherwise may not possible, as when the inferential/distributional characteristics of an estimator are only known or derived under restrictive distributional assumptions. NAbC eschews such restrictions, thus permitting accurate, all-else-equal comparisons.

Note that all empirical results presented herein use the sample estimators specified in the Introduction and Background section, and sample sizes in every example all exceed 10p (10 times the dimension of the matrix), which is a widely used threshold for whether a more sophisticated, bias-correcting estimator

JD Opdyke, Chief Analytics Officer

<sup>&</sup>lt;sup>63</sup> Note that, as an empirical method, the ability to implement NAbC relies on the degree to which the empirical (kernel-based) simulations of the angles approximate continuous distributions over their entire sample spaces. This is largely controlled by the number of simulations run, and fortunately the "good behavior" of the angles distributions renders N=10,000 simulations, which is computationally feasible even for non-small matrices, more than sufficiently large in most cases. However, empirically challenging cases can arise. For example, if we are comparing two sample matrices where some cell values are very different, the distribution of the difference between the two angles distributions (corresponding to the same cell in each matrix) may not contain the value zero, in which case the empirically-based p-value would be exactly zero. Like any empirical method (e.g. bootstraps, permutation tests, etc.) care must be taken to ensure that the consequences of such results are noted, understood, and properly accounted for.

is needed, at least for Pearson's matrix (see Bongiorno et al., 2023). As mentioned in Section 3 above, I recommend for conditional (forecast) estimation the Average Oracle (AO) of Bongiorno et al. (2023) (see also Bongiorno & Challet, 2023a, for an extensive empirical study against competitors). Further testing may show that AO can be applied to all positive definite dependence measures as well, not just Pearson's, although this currently is the topic of continuing research. In the next section, I show how all of the previously derived characteristics of NAbC remain valid for the scenario-restricted case, that is, when selected cells of the all-pairwise matrix are 'frozen' as dictated by specific scenarios, while the rest are allowed to vary.

### 5. NAbC: Granular, Fully Flexible Scenarios, Reverse Scenarios, and Stress Testing

### 5.a. Review of Existing Methods

"... however, a unified and generally accepted correlation risk management framework does not yet exist" (Packham & Woebbeking, 2023, p.1).

The size, breadth, and surprise of the effects of correlation breakdowns are well documented in the literature (see Kim & Finger, 1998; Loretan & English, 2000; Li et al., 2024; BIS, 2011a, 2011b; Nawroth et al., 2014; Ng et al., 2014; Yu et al., 2014; Chmeilowski, 2014; Epozdemir, 2021; Feng & Zeng, 2022; and Parlatore and Philippon, 2024), if underappreciated during periods of relative market calm: "Furthermore, joint distributions estimated over periods without panics will misestimate the degree of correlation between asset returns during panics." (FRB Chairman, Alan Greenspan, 1999). And yet despite its importance, the ability to model, predict, and mitigate correlation breakdowns effectively across very different scenarios, in a fully flexible way, has remained elusive.

To start with, although many approaches do otherwise, it is not enough to stress only the inputs to a correlation/dependence matrix - the matrix itself must be stressed and evaluated under stressed conditions of a particular (extreme) scenario: "... in order to calculate stressed VaR accurately it is also necessary to stress the correlation matrix ... most correlations tend to increase during market crises, asymptotically approaching 1.0 during periods of complete meltdown, such as occurred in 1987, 1998 and 2008. ...Certain methods that could be meaningful [include e]mploying fat-tailed distributions for the risk factors and replacing the standard correlation matrix with a stressed one....." (BIS, 2011a). Secondly, if a method perturbs eigen decompositions and/or polar angles to obtain correlation/dependence measure distributions, this cannot be done on an ad hoc basis, using mathematically convenient distributions, like the Gaussian, to perturb eigen values, or arbitrary bounded functions, like inverse tangent, to perturb angles (see Galeeva, 2007). Spectral and spherical distributions follow specific and often known distributions under various conditions, and such approaches need methodological support, whether theoretical or empirical or both, to justify their use when taking what is otherwise a smart approach to generating scenario-specific correlation/dependence matrices. Additionally, such methods must remain cognizant of all the characteristics of the conditions they are attempting to generate. Hardin JD Opdyke, Chief Analytics Officer Page **64** of **92 Correlation and Beyond** 

et al. (2013), for example, utilize a normalized vector of independent gaussian random variables to perturb the observed correlation matrix, but correctly note that "The amount of noise that can be added to the original matrix is determined by its smallest eigenvalue. ... We provide the user with ... a general algorithm to apply to any correlation matrix for which the smallest eigenvalue can be reasonably estimated." (emphasis added). Unfortunately, as mentioned above, this eliminates what are arguably the most widely observed correlation matrices in finance – those based on a 'spiked' covariance matrix (see Johnstone, 2001) where one or a few eigenvalues dominate and the majority of eigenvalues are close to zero, which often indicates they cannot be reliably estimated. Robustness as dependence matrices approach singularity/NPD is an important quality of any method, but it remains especially critical in the analysis of financial portfolios.

Several other approaches avoid these limitations (see Packham and Woebbeking, 2021, Chmielowski, 2014, and Parlatore and Phillippon, 2024) but they have other arguable limitations. For example, Packham and Woebbeking, 2021, enforce positive definiteness ex post, which as described above distorts the desired distributions of dependence measures. And none of these provide granular, celllevel control to restrict perturbation on any combination of cells, while still obtaining a valid distribution of the remaining cells of the correlation/dependence matrix being used. Yet this is exactly what is needed for realistic scenario analytics and stress testing, let alone precise attribution analyses and 'what if' analysis capabilities. Correlation/dependence matrices under a tech market bubble (2000) vs those under a housing bubble (2008) vs those under Covid (2020) will change very different individual cells of the all-pairwise matrix, and very different combinations of cells, in very different ways, often in terms of both direction and magnitude, while leaving many cells strongly affected under one upheaval completely unaffected under another (see Feng & Zeng, 2022). But some combinations of cells might change similarly in all of these scenarios, and distributional analyses must be able to accommodate every possible combination of changes, in terms of both magnitude and direction. In other words, while correlation 'breakdowns' will occur under all of these extreme conditions, the granular nature of allpairwise matrices ensures that the fundamentally different (and sometimes similar) nature of these breakdowns will be captured and reflected empirically in all related analyses. Although some approaches settle for stretching across a few different covariance/correlation matrices, with fixed values, representing several different scenarios (see Parlatore and Phillippon, 2024), this arguably still is too rigid and discrete and limited for realistic analyses of the dynamic distributions of these matrices, and cannot remain truly robust across qualitatively different, and often as yet unobserved (future) breakdowns. Neither does this approach allow for flexible, all-else-equal, targeted 'what if' analyses, or granular attribution analyses. If we are to achieve the same level of flexibility in quantitatively modeling dependence matrices that has been attained for the other parameters in the risk and investment models of our portfolios, practitioners and applied researchers must be able to flexibly and realistically model dependence matrices at the most granular level - that of the individual correlation cells - without restriction.

Despite the research on correlation breakdowns listed above, I am aware of only two other limited attempts at this granular level of modeling specific groups of cells in the dependence matrix (see Saxena

Page **65** of **92** 

et al., 2023, and Veleva, 2017), and both are restricted in notable ways.<sup>64</sup> Fortunately, NAbC allows for specifying ANY combination of cells, within the framework of the all-pairwise matrix, to be 'frozen' at their current values while allowing all the rest to vary, providing full flexibility within this framework.

## 5.b. A New Method for Fully Flexible Scenarios

Several results allow for this full flexibility. First, 1. independence of the angles distributions allows us to vary individual cells without changing the structure of the multivariate distribution of the entire matrix. Second, 2. the distributions of individual correlation cells, as well as the distribution of the entire correlation matrix, both remain invariant to the ordering of the rows and columns of the matrix (see Pourahmadi & Wang, 2015, and Lewandowski et al., 2009). Third, 3. based on 1. and 2., we can exploit the simple mechanics of matrix multiplication so that only selected cells of the matrix are affected, and the rest frozen, as required for a given scenario.

To explain 3., I focus only on the lower triangle of the correlation matrices below in Graphs 15-17, since the upper triangle is just its reflection due to symmetry. Note again that using NAbC, we only perturb angles. We never perturb the correlation values directly. We must always convert correlations to angles, perturb the angle values using the fitted kernels, and then translate back to correlation values. In doing so, when multiplying the Cholesky factor by its transpose,  $R = BB^T$ , changing a given angle cell in matrix *B* will affect other cells, but only those cells to the right of it in the same row, and those below the diagonal of the corresponding column, as shown graphically for several examples in Graph 15 below.<sup>65</sup>

1						1						1						1						1						1					
	1						1						1						1						1						1				
		1						1						1						1						1						1			
			1						1						1						1						1						1		
				1						1						1						1						1						1	
					1						1						1						1						1						1

#### GRAPH 15: Correlation Cells Affected by Changing a Specific Angle Cell

This means that we can simply reorder the matrix so that the targeted cells we want to vary all end up in the rightmost triangle of the lower triangle of the matrix, according to the fill order in Graph 16 below.

<sup>&</sup>lt;sup>64</sup> Saxena et al. (2023) explores the possibility of restricting individual covariance/correlation terms to zero, although they are not always able to enforce this restriction while maintaining positive definiteness. Velena (2017) restricts the values of the correlation matrix being simulated to specified ranges, but only for all off-diagonal cells; in some cases, one algorithm allows for cell-level values to be specified, but without guarantees of positive definiteness.

<sup>&</sup>lt;sup>65</sup> Note that not all of these (orange) correlation cells will necessarily change if values of zero are involved, but none OTHER than these (orange) correlation cells CAN change when only the red angle cell changes.

#### GRAPH 16: 'Fill Order' for Resorting the Correlation Matrix to Change Only Specific Cells

#### **Rightmost Triangle Fill Order**

11					
12	7				
13	8	4			
14	9	5	2		
15	10	6	3	1	

If we only change in matrix *B* the angle values of cells 1, 2, and 3 above, no other cells in the correlation matrix *R* will be affected, simply by virtue of the mechanics of matrix multiplication from  $R = BB^T$ . Below I show another example. Reorder the correlation matrix so that rows 1-6 are now 6-1 and columns 1-6 are now 6-1, so that the original (green) cells 1,2 and 1,3 and 2,3 and 4,3 are now in the rightmost triangle of

# GRAPH 17: Example of Matrix Resorting Using the 'Fill Order' to Change Only Specific Correlation Cells

#### **Determine Targeted Change Cells**

1,2				
1,3	2,3			
		4,3		

Reorder Rows/Cols to Fill Rightmost Triangle with Targets According to Fill Order



Changes in These Angles Cells ONLY change the Same Cells in the Correlation Matrix

11					
12	7				
13	8	4,3			
14	9	5	2,3		
15	10	6	1,3	1,2	

the lower triangular matrix, in the fill order shown above.

Changes to the corresponding cells in the angles matrix *B* (the orange cells) will only change these same cells, after  $R = BB^{T}$ , in the resulting correlation matrix, *R*, leaving the rest unaffected. Note that the green cells to be targeted for change do not even have to be contiguous, nor do they have to completely 'fill' the rightmost (orange) triangle (note that cells 5 and 6 in the right matrix are not targeted): they only must fill the rightmost triangle according to the order of the center matrix above. Note also that the "rightmost triangle" rule is nested/hierarchical: if I wanted to perform 'what if' analyses on only one of those cells (e.g. cell "1,2") without changing the other three, I order the original correlation matrix to place that cell as the 'first' in the lower triangle of the *B* matrix, as shown. Then, subsequent changes to it will not affect the other (orange) cells, let alone any other non-orange cells. In contrast, changes to cell "4,3" will affect the values of the other orange cells (but not the non-orange cells). Readers are encouraged to test this in the interactive spreadsheet (url link provided above).

So we can exploit these four simultaneous conditions – 1. independence of the angles distributions; 2. distribution invariance of the correlation/dependence matrix, and its individual cells, to row and column order; 3. the mechanics of matrix multiplication; and 4. the granular, cell-level geometry of NAbC – to JD Opdyke, Chief Analytics Officer Page 67 of 92 Correlation and Beyond

obtain great flexibility in defining scenarios wherein some cells vary and some do not. To my knowledge, no other approach allows this degree of flexibility, which is what is required for defining correlation/dependence matrices for use in realistic, plausible, and sometimes extreme stress market scenarios. This also greatly simplifies attribution analyses, isolating and making transparent the identification of effects due to specific pairwise associations, which is something spectral and more aggregated analyses cannot do in this setting.

To be clear, the above allows for the specification of ANY scenario within the structure of the pairwise matrix. Note, however, that some scenarios can include combinations of cells which are forced to include (in the lower right triangle) one or a few cells not affected by the scenario. This is unavoidable due to the structure of the pairwise matrix: for example, in the matrix above, there are only p! (i.e.5!=120) ways to sort the rows and columns, but there are [p(p1-)/2]! (i.e. 15!= 1,307,674,368,000) ways to sort the 15 cells freely. The matrix obviously cannot accommodate freely sorting the individual cells in this way because it breaks the structure of the matrix. Some scenarios, therefore, could conceivably be required to include for perturbation some few additional cells in the lower rightmost triangle that are not relevant to the scenario and otherwise should be held constant. Fortunately, in practice, especially with large matrices, this appears to be a relatively rare occurrence, and when it happens, the effects are identifiable so that materiality can be assessed via 'what if' analyses on the specific cells. But dealing with these potential cases appears to be well worth the price of the unmatched flexibility that using NAbC and the all-pairwise matrix provides,<sup>66</sup> not to mention the other advantages it maintains over more complex, strictly multivariate dependence measures. For usage with actual market data, the latter typically are more difficult to estimate with the same levels of precision and statistical power, let alone to manipulate for purposes of intervention or mitigation. In contrast, pairwise associations are directly identifiable, typically more easily and accurately estimated,<sup>67</sup> and intervention 'what if' tests are more targeted and transparent.

To conclude this section, I deal with one final implementation issue. When the matrix is scenariorestricted, and we only perturb a subset of the matrix while keeping the remaining cells fixed, what values do we use for the angles of those 'frozen' cells? This is where the mean angles matrix, used to account for asymmetry when calculating the two-sided p-values in the previous section, comes into play. When the matrix angles are sampled using the fitted kernel densities, a sample is drawn from the entire matrix,

<sup>&</sup>lt;sup>66</sup> Most of the related scenario literature perturbs scenario-based cells and simply ignores their (often notable) effects on the rest of the matrix (which should remain 'frozen,' but isn't), not to mention the effects of the rest of the matrix on the scenario-related cells. These papers euphemistically refer to the former as 'peripheral' correlations (see Ng et al. (2014) and Yu et al. (2014)). NAbC is the only method that fully controls the values and thus, the indirect effects of these so-called 'peripheral' correlations.

<sup>&</sup>lt;sup>67</sup> They also can be *estimated* rigorously, and with targeted precision and flexibility, with well-established methods such as vine copulas (see Czado and Nagler, 2022)). Ironically, however, when used for *inference or sampling* for this problem specifically, vine copulas and similar methods become extremely unwieldy and much more complex and less transparent than NAbC, not to mention ungeneralizable beyond Pearson's (see the vine and extended onion algorithms of Lewandowski et al., 2009, and the similar chordal sparsity method of Kurowicka, 2014).

and if it is scenario restricted, the sampled values for those cells that are 'frozen' are simply overwritten with their means. So after N=10,000 samples are drawn, all 10,000 values of the 'frozen' cells unaffected

by the scenario will have the same mean value for that specific cell, and when translated via  $R = BB^{T}$  back into correlation matrices, all the correlation values for those cells will be the mean correlation values of the respective cells. In other words, their values will not change, and will remain 'frozen,' based on a reasonably robust estimator of their true value (note that these 'frozen' values are not based on a single estimated matrix, but rather, they each are the means of N=10,000 matrices, which is similar to the approach Bongiorno et al. (2023) use for their average oracle estimator, as well as the estimator used in Sun and Huang (2025)). The order of magnitude of empirical accuracy of these values is inversely related to the number of samples drawn, N. In the example in Section 6 below, we observe accuracy to the fourth decimal place for these frozen cells when using N=25,000 simulations, as expected. Alternately, the values could be treated as truly known constants from the beginning, but it is more conservative (and realistic) to use estimates based on the mean of all the samples.<sup>68</sup>

I end this section by reemphasizing that this matrix sorting method for providing fully flexible scenarios, within the framework of the all-pairwise matrix, applies not only to Pearson's, but also to all positive definite dependence measures, under the fully general data conditions for financial portfolios described above. One complete, empirical example of all of NAbC's inferential capabilities, covering all (but one) of its original objectives described in the Introduction and Background, is shown below.

# 6. NAbC Example: Kendall's Tau p-values & Confidence Intervals, Unrestricted & Scenariorestricted, One- and Two-sample Tests

Now, with NAbC's characteristics established and its broad range of application described above, I can present a complete empirical example of its implementation here.<sup>69</sup> This example will check seven of the eight original objectives boxes listed in the Introduction and Background above (objective 1. is ignored in this example, and the data generating mechanism used is simply multivariate standard normal, solely for the purpose of facilitating for the reader the replication of these results). The dependence measure chosen is Kendall's Tau, under two cases: unrestricted, and scenario-restricted. NAbC provides both p-values and confidence intervals, at both the cell level and matrix level, with N=25,000 simulations and the number of observations n = 160, representing about eight months of daily market returns. The values

<sup>&</sup>lt;sup>68</sup> Note that when NAbC is being used as a two sample test in the scenario-restricted setting, we are only testing, by design, the scenario-relevant cells. So the mean values of each matrix are inserted into the 'frozen' cells of each matrix just as in the one-sample test, but then those cells are ignored thereafter, i.e. when calculating cell-level p-values, as well as the (restricted) matrix-level p-values, just as in the one-sample test. This ensures that the two-sample test is conducted only for the scenario-relevant cells, even as we properly perturb each entire matrix (and insert means ex post) when generating the samples.

<sup>&</sup>lt;sup>69</sup> See Opdyke (2022, 2023, and 2024a) for extensive additional examples under wide ranging data conditions.

of the matrix [A] are based on a Pearson's matrix from A' below, translated to A via  $\tau = (2/\pi) \arcsin(r)$ , which is valid under elliptical data (see McNeil et al., 2005),<sup>70</sup> and approximately valid under some broader classes of distributions (see Hansen & Luo (2024) and Hamed (2011) for examples).

1				
0.2	1			
-0.1	0.3	1		
0.3	-0.3	-0.1	1	
0.6	0.4	0.0	0.1	1

[A'] = [

<u>UNRESTRICTED CASE</u>: Given a specified or well-estimated dependence matrix [A], and its specified or well-estimated data generating mechanism:

		[A]			
1					
0.13	1				
-0.06	0.19	1			
0.19	-0.19	-0.06	1		
0.41	0.26	0.00	0.06	1	

		[B]		
0.8				
0.7	0.8			
0.8	0.7	0.7		
0.7	0.8	0.8	0.7	

		[C]		
1				
0.3	1			
0.1	0.1	1		
0.05	-0.1	0.1	1	
0.5	0.25	0.2	0.15	1

- Q1. **Confidence Intervals**: What are the two dependence matrices that correspond to the lower– and upper–bounds of the 95% confidence interval for [A]? What are, simultaneously, the individual 95% confidence intervals for each cell of [A]?
- Q2. **Quantile Function**: What is the unique dependence matrix associated with [B], a matrix of cumulative distribution function values associated with the distribution of [A]?
- Q3. **p-values**: Under the null hypothesis that observed dependence matrix [C] was sampled from the data generating mechanism of [A], what is the p-value associated with [C]? And simultaneously, what are the individual p-values associated with each cell of [C]?
- Q4. **Two-sample p-values**: Under the null hypothesis that observed dependence matrix [A] and observed dependence matrix [C] each were sampled from the same population, and therefore have the same values, what is the matrix-level p-value? And simultaneously, what are the individual p-values associated with each cell of the matrix?

<u>SCENARIO-RESTRICTED CASE</u>: Under a specific scenario only selected pairwise dependence cells of [A] will vary (green), while the rest (red) are held constant, unaffected by the scenario (e.g. COVID). This is

<sup>&</sup>lt;sup>70</sup> See Koike et al. (2024) for a sophisticated paper defining the broader-than-expected conditions under which Pearson's retains the invariance property under marginal transformations.

#### matrix [D].



- Q5. <u>Confidence Intervals</u>: What are the two dependence matrices that correspond to the lower– and upper–bounds of the 95% confidence interval for [D] (holding constant the non-selected red cells)? What are, simultaneously, the individual 95% confidence intervals for the green cells of [D]?
- Q6. **Quantile Function**: What is the unique dependence matrix associated with [E], a matrix of cumulative distribution function values associated with the distribution of [D]'s green cells?
- Q7. **p-values**: Under the null hypothesis that observed dependence matrix [F] was sampled from the (scenario-restricted) data generating mechanism of [D], what is the p-value associated with [F] (with red cells held constant)? And simultaneously, what are the individual p-values associated with each (green) cell of [F]?
- Q8. <u>Two-sample p-values</u>: Under the null hypothesis that observed dependence matrix [D] and observed dependence matrix [F] each were sampled from the same population, and therefore have the same values (for their unrestricted green cells), what is the matrix-level p-value (for the unrestricted portion of the matrix)? And simultaneously, what are the individual p-values associated with each (green) cell of the matrix?

Answers to these questions require inference at both the cell and matrix levels, simultaneously and with cross-level consistency, as well as requiring the matrix-level quantile function, all under both the unrestricted and scenario-restricted cases. Only NAbC can simultaneously answer Q1.-Q8. above, as shown in Tables B1 and B2 below.

For Q1 and Q5, the two top matrices correspond to the first (matrix-level) question, and the bottom two matrices correspond to the second (cell-level) question. Note the wider intervals on a cell-by-cell basis for the matrix-level confidence intervals compared to the cell-level confidence intervals, as expected. Also note, for Q3 and Q7, the smaller p-values for the individual cells compared to the respective matrix-level p-values, which are larger, as expected, as they are analogous to the family-wise error rate (FWER) of a joint hypothesis covering all cells of the matrix. Note also that the green cells of Q6 differ from the corresponding cells in Q2: even though the (green) angles distributions themselves remain unaffected by scenario restrictions, the ultimate correlation values of those cells ARE affected due to the matrix multiplication of the Cholesky factor,  $R = BB^T$ . Comparing the two-sample test of Q4 to the one-sample test of Q3, we find, as expected, the increased variability from two samples increases all the cell-level p-values as well as the matrix-level p-value. Only when we double the sample size (as well as the number of simulations to more accurately account for smaller p-values) do we obtain similarly small p-values of

JD Opdyke, Chief Analytics Officer

Page **71** of **92** 

**Correlation and Beyond** 

				Γ					1										1
	5370	0043		-				1	0789		1285	÷ 2		0008				1	.0131
34	0=6	= 0.(	0=	_			1	.1148	.0043 0.	_	)=0.`	), N=50k — 0_1		= 0.			1	.0263	0 8000.
	alue	~~0.05)	(10.0-2) (10.0-2)			1	0723	1350 0	0 2260		'alu€	(n=320	α=0.05)	α=0.01)		н Н	.0072	.0319 0	.0110 0
	^-   0	$\alpha_{_{_{FDP }}}$	$\alpha_{_{FDP }}$		-	.0213	.0342 0.	.0593 0.	.1194 0.		v-d	2	C FDR	$\alpha_{_{FDR}}$	1	.0042	.0078 0.	.0056 0.	.0282 0.
						0	0	•		」 】						0	0	0	0
	503	5	900						8										
	0,11	= 0.0	= 0.0(					6	0.00										
Ø	lue=	0.05	(0.0	(10.0-			80	5 0.026	7 0.008										
	ev-c	ja va	$\alpha = \frac{1}{\alpha} \sum_{\alpha = 1}^{\alpha} \frac{1}{\alpha}$	mlun		10	0.021	0.031	0.015										
		α,	ά			0.000	0600.0	0.0222	0.0170										
		1																	
					H							1						Ч	
				1	0.1040 1						1	0.2789 1					1	0.2103 1	
Q2			1	-0.0392 1	0.0614 0.1040 1					1	0.1182 1	0.2335 0.2789 1				1	0.0570 1	0.1611 0.2103 1	
Q2		1	0.2369 1	-0.1510 -0.0392 1	0.3159 0.0614 0.1040 1				1	0.3475 1	0.0127 0.1182 1	0.4370 0.2335 0.2789 1			1	0.3013 1	-0.0525 0.0570 1	0.3849 0.1611 0.2103 1	
Q2		0.1729 1	-0.0355 0.2369 1	0.2374 -0.1510 -0.0392 1	0.4335 0.3159 0.0614 0.1040 1			1	0.2735 1	0.0910 0.3475 1	0.3323 0.0127 0.1182 1	0.5250 0.4370 0.2335 0.2789 1		1	0.2300 1	0.0424 0.3013 1	0.2920 -0.0525 0.0570 1	0.4929 0.3849 0.1611 0.2103 1	
Q2	1	0.1729 1	-0.0355 0.2369 1	0.2374 -0.1510 -0.0392 1	0.4335 0.3159 0.0614 0.1040 1			1	0.2735 1	0.0910 0.3475 1	0.3323 0.0127 0.1182 1	1 0.5250 0.4370 0.2335 0.2789 1		1	0.2300 1	0.0424 0.3013 1	0.2920 -0.0525 0.0570 1	1 0.4929 0.3849 0.1611 0.2103 1	
Q2		0.1729 1	-0.0355 0.2369 1	0.2374 -0.1510 -0.0392 1	0.4335 0.3159 0.0614 0.1040 1			1	0.2735 1	0.0910 0.3475 1	1 0.3323 0.0127 0.1182 1	0830 1 0.5250 0.4370 0.2335 0.2789 1		1	0.2300 1	0.0424 0.3013 1	1 0.2920 -0.0525 0.0570 1	0.0478 1 0.4929 0.3849 0.1611 0.2103 1	
21 Q2		0.1729 1	-0.0355 0.2369 1	0.2374 -0.1510 -0.0392 1	0.4335 0.3159 0.0614 0.1040 1			1	0.2735 1	1 0.0910 0.3475 1	.1602 1 0.3323 0.0127 0.1182 1	.1873 -0.0830 1 0.5250 0.4370 0.2335 0.2789 1		1	0.2300 1	1 0.0424 0.3013 1	0.1427 1 0.2920 -0.0525 0.0570 1	1.1396         -0.0478         1         0.4929         0.3849         0.1611         0.2103         1	
Q1 Q2		0.1729 1	-0.0355 0.2369 1	0.2374 -0.1510 -0.0392 1	0.4335 0.3159 0.0614 0.1040 1			1	1 0.2735 1	0626 1 0.0910 0.3475 1	3567 -0.1602 1 0.3323 0.0127 0.1182 1	0926 -0.1873 -0.0830 1 0.5250 0.4370 0.2335 0.2789 1		1	1 0.2300 1	.0986 1 0.0424 0.3013 1	.3131         -0.1427         1         0.2920         -0.0525         0.0570         1	.1410 -0.1396 -0.0478 1 0.4929 0.3849 0.1611 0.2103 1	
Q1 Q2	1	0.1729 1	-0.0355 0.2369 1	0.2374 -0.1510 -0.0392 1	0.4335 0.3159 0.0614 0.1040 1			1	0172 1 0.2735 1	2100 0.0626 1 0.0910 0.3475 1	0472         -0.3567         -0.1602         1         0.3323         0.0127         0.1182         1	2794         0.0926         -0.1873         -0.0830         1         0.5250         0.4370         0.2335         0.2789         1		1	0250 1 0.2300 1	.1661         0.0986         1         0.0424         0.3013         1	0926 -0.3131 -0.1427 1 0.2920 -0.0525 0.0570 1	3210 0.1410 -0.1396 -0.0478 1 0.4929 0.3849 0.1611 0.2103 1	

TABLE B1: NAbC's Complete Inferential Results, Example of Kendall's Tau – Cell and Matrix Level p-values and Confidence Intervals

**Correlation and Beyond** 

Page 72 of 92

JD Opdyke, Chief Analytics Officer
Q8	p-value=0.4436	$lpha_{_{FDR(lpha=0,05)}}=0$	$\alpha_{FDR(\alpha=0.01)}=0$			0.0328		0.3992 0.0290 0.0140		p-value=0.2251	(n=320, N=50k) $lpha_{_{FD,P(lpha=005)}}=0.0017$	$\alpha_{FDR(\alpha=0.01)} = 0.0017$			0.0017		0.2231 0.0008 0.0001
Q7	p-value=0.0436	$\alpha_{_{FD,R(\alpha=0,05)}} = 0.05$	$lpha_{FDR(lpha=0.01)}=0$			0.0047		0.0077 0.0148 0.0171									
					1						1					1	
					867						-					92	
				F	0.0					1	0.184				1	0.148	
Q6			1	-0.0639 1	0.0362 0.0				1	-0.0639 1	0.1150 0.184			1	-0.0639 <b>1</b>	0.0809 0.148	
Q6		1	0.2398 1	-0.1940 -0.0639 1	0.2895 0.0362 0.00			1	0.3492 1	-0.1940 -0.0639 1	0.3425 0.1150 0.184		1	0.3006 1	-0.1940 -0.0639 1	0.3165 0.0809 0.148	
Q6	1	0.1282 1	-0.0636 0.2398 1	0.1942 -0.1940 -0.0639 1	0.4098 0.2895 0.0362 0.0		1	0.1282 1	-0.0636 0.3492 1	0.1942 -0.1940 -0.0639 1	0.4098 0.3425 0.1150 0.184	1	0.1282 1	-0.0636 0.3006 1	0.1942 -0.1940 -0.0639 1	0.4098 0.3165 0.0809 0.148	
Q6	1	0.1282 1	-0.0636 0.2398 1	0.1942 -0.1940 -0.0639 1	0.4098 0.2895 0.0362 0.0		1	0.1282 1	-0.0636 0.3492 <b>1</b>	0.1942 -0.1940 -0.0639 1	1      0.4098      0.3425      0.1150      0.184	1	0.1282 1	-0.0636 0.3006 1	0.1942 -0.1940 -0.0639 1	1      0.4098      0.3165      0.0809      0.148	
Q6	1	0.1282 1	-0.0636 0.2398 1	0.1942 -0.1940 -0.0639 1	0.4098 0.2895 0.0362 0.0		1	0.1282 1	-0.0636 0.3492 <b>1</b>	1      0.1942      -0.1940      -0.0639      1	-0.0541 <b>1</b> 0.4098 0.3425 0.1150 0.184	1	0.1282 1	-0.0636 0.3006 1	<b>1</b> 0.1942 -0.1940 -0.0639 <b>1</b>	-0.0184      1      0.4098      0.3165      0.0809      0.148	
Q5 Q6	1	0.1282 1	-0.0636 0.2398 1	0.1942 -0.1940 -0.0639 1	0.4098 0.2895 0.0362 0.0		1	0.1282 1	1 -0.0636 0.3492 1	-0.0639 <b>1</b> 0.1942 -0.1940 -0.0639 <b>1</b>	-0.1144 -0.0541 <b>1</b> 0.4098 0.3425 0.1150 0.184	1	0.1282 1	1 -0.0636 0.3006 1	-0.0639 <b>1</b> 0.1942 -0.1940 -0.0639 <b>1</b>	-0.0789 -0.0184 <b>1</b> 0.4098 0.3165 0.0809 0.146	
Q5 Q6	1	0.1282 1	-0.0636 0.2398 1	0.1942 -0.1940 -0.0639 1	0.4098 0.2895 0.0362 0.0		1	1 0.1282 1	0.0478 1 -0.0636 0.3492 1	-0.1940 -0.0639 <b>1</b> 0.1942 -0.1940 -0.0639 <b>1</b>	0.1757 -0.1144 -0.0541 1 0.4098 0.3425 0.1150 0.18	1	1 0.1282 1	0.0904 1 -0.0636 0.3006 1	-0.1940 -0.0639 <b>1</b> 0.1942 -0.1940 -0.0639 <b>1</b>	0.2028 -0.0789 -0.0184 1 0.4098 0.3165 0.0809 0.146	

TABLE B2: NAbC's Complete Inferential Results, Example of Kendall's Tau – Cell and Matrix Level p-values and Confidence Intervals

**Correlation and Beyond** 

Page 73 of 92

JD Opdyke, Chief Analytics Officer

the same order of magnitude. Similar patterns hold for the one-sample vs two-sample test results under the scenario-restricted cases (Q7 and Q8, respectively). Finally, note that the empirical values of the red cells in Q5-Q6 differ slightly from those in [D] and [F]. This is due to NAbC's conservative use of the mean of the estimated angles (correlation) matrices, rather than presuming we know the absolute 'true' values of these cells (although this is justified in some specific cases).

In terms of actual runtimes, note that NAbC is somewhat computationally intensive, but not prohibitively so. Implementing NAbC on synthetic data representing real-world data conditions (e.g. margins with different and varying degrees of asymmetry, non-stationarity, serial correlation, and heavy-tailedness) for non-small portfolios of dimension 100x100, on a commodity laptop purchased in 2019 with 32GB of RAM but no multi-threading, NAbC generates a full set of results, based on N samples = 10,000, in about 2.4 hours. However, in a multithreaded environment, let alone one with more memory, NAbC could be applied on similarly non-small matrices in minutes. For the specific case of the gaussian identity matrix, applying inverse probability transform sampling as described above, on a 100x100 matrix with N = 10,000 samples, NAbC takes less than 25 minutes to run on the same laptop. Notably, Rubsamen (2023) benchmarked NAbC's sampling under the gaussian identity matrix against that of Makalic & Schmidt (2018) and obtained up to a 30% reduction in runtime under NAbC. But of course, NAbC's analytic (non-sampling) solution under these conditions is instantaneous (see url for excel workbook above). So while NAbC is not a result that currently can be used "real-time" for, say, high frequency trading (except for when the fully analytic solution is valid), its runtimes remain reasonable given its very generalized application and widely available modern computing resources.

## 7. NAbC: Beyond 'Distance' to Generalized Entropy

In a relevant and validating digression, it is intriguing and important to note that the (two-sided) cell-level p-values NAbC provides (see Q3 and Q7 in Table B above) actually can be used to construct a competitor to commonly used distance metrics, such as norms, and it has a number of advantages over them in this setting. Some commonly used norms for measuring correlation 'distances' include the Taxi, Frobenius/Euclidean, and Chebyshev norms (collectively, the Minkowski norm), shown below in (40).

(40) 
$$||x|| = \left(\sum_{i=1}^{d} |x_i|^m\right)^{1/m}$$

where x is a distance from a presumed or baseline correlation value, d=number of observations, and m=1, 2, and  $\infty$  correspond to the Taxi, Frobenius/Euclidean, and Chebyshev norms, respectively.

All of these norms measure absolute distance from a presumed or baseline correlation/dependence value. But the range of all relevant and widely used dependence measures is bounded, either from –1 to 1 or 0 to 1, and both the relative impact and *meaning* of a given distance at the boundaries are not the same as those in the middle of the range. In other words, a shift of 0.02 from an original or presumed correlation/dependence value of, say, 0.97, means something very different than the same shift from 0.37. NAbC's p-values attribute probabilistic MEANING to these two different cases, while a norm would JD Opdyke, Chief Analytics Officer Page 74 of 92 Correlation and Beyond treat them identically, even though they very likely indicate what are very different events of very different relative magnitudes with potentially very different consequences.

Therefore, a natural, PROBABILISTIC distance measure, based directly on NAbC's cell-level p-values, is the natural log of the product of the p-values, dubbed 'LNP' in (41) below:

(41) "LNP" = 
$$\ln\left(\prod_{i=1}^{q} p\text{-value}_i\right) = \sum_{i=1}^{q} \ln\left[p\text{-value}_i\right]$$
 where  $q = p(p-1)/2$  and  $p\text{-value}_i$  is 2-sided.

Using a Pearson's correlation matrix under the (Gaussian) identity matrix, LNP shows a very strong correspondence with the entropy of the correlation matrix, defined by Felippe et al. (2021 and 2023) as (42) below:

(42) Entropy = 
$$Ent(R/p) = -\sum_{j=1}^{p} \lambda_j \ln(\lambda_j)$$

where *R* is the sample correlation matrix and  $\lambda_j$  are the p eigenvalues of the correlation matrix after it is scaled by its dimension, *R*/p. Importantly, this result (42), like NAbC, is valid for ANY positive definite measure of dependence, not just Pearson's. Graph 18 below compares LNP of Kendall's Tau matrix to the entropy of Kendall's Tau matrix in 10,000 simulations (with n=126 for half a year of daily returns) under the Gaussian identity matrix, and the Pearson's correlation between them (0.98) is virtually identical to the same comparison based on Pearson's matrix rather than Kendall's matrix (just under 0.99).<sup>71</sup>

It is important to note, however, that entropy here is limited to being calculated relative to the case of independence, which for many dependence measures corresponds only with the identity matrix.<sup>72</sup> In contrast, LNP can be calculated, and retains its meaning, in all cases, based on ANY values of the dependence matrix, not just the case of independence. Yet the correspondence of LNP to entropy under this specific case speaks to LNP's natural interpretation as a meaningful measure of deviation or distance or disorder (depending on your interpretation), and one that also is more flexible and granular than entropy as it is measured cell-by-cell, p(p-1)/2 times, as opposed to only p times for p eigenvalues. As such, LNP might be considered a type of 'generalized entropy' relative to any baseline of the dependence measure, as specified by the researcher, including as a special case perfect (in)dependence. Such entropy-related measures certainly are relevant in this setting as entropy has been used increasingly in the literature to measure, monitor, and analyze financial markets (see Meucci, 2010b, Almog and Shmueli, 2019, Chakraborti et al., 2020, and Vorobets, 2024, 2025, for several examples). So the use of LNP here warrants further investigation as a matrix-level measure that, unlike widely used distance measures such as norms, has a solid and meaningful probabilistic foundation. Its

<sup>&</sup>lt;sup>71</sup> In addition, the Pearson's correlation between LNP and the entropy of Felippe et al. (2021 and 2023), under these conditions of the Gaussian identity matrix, was the same – 0.98 – for both Spearman's and Chatterjee's (symmetric version).

<sup>&</sup>lt;sup>72</sup> Recall, of course, that a zero value for Pearson's or Kendall's or Spearman's does not imply independence, but independence does imply a zero value for these measures (the exception being Pearson's under Gaussian data, for which a zero value does indicate independence).



Graph 18: Identity Matrix Simulations for Kendall's Tau – LNP vs. Correlation Matrix Entropy

calculation applies not only beyond the independence case generally, but also to ALL positive definite measures of dependence, regardless of their values. LNP's range of application is as wide as that of NAbC's matrix-level p-value, and the two are readily calculated side-by-side as they are both based on NAbC's cell-level (two-sided) p-values for the entire matrix. These are intriguing results with possibly far-reaching implications.

## 8. NAbC: Future Research and Additional Applications

There are a number of areas where additional research can further validate and potentially increase the utility and breadth of NAbC's application.

Analytic Angles Distributions: I provide above the derivation of NAbC's fully analytic solution under the Gaussian identity matrix, but this is a narrow (albeit foundational) case. Although NAbC's general solution remains 'runtime reasonable' given its generality and objectives, expanding the range of conditions for an analytic solution for the angles distributions would dramatically speed NAbC's implementation. Deriving an "all cases" analytic solution currently appears to be a nontrivial problem, but even providing this under additional specific cases would be very useful and directly useable in NAbC's application.

<u>Competing Distributional Methods</u>: Implementing and comparing NAbC's results to those of competing, if less flexible methods, like Hansen & Archakov (2021) and the Bayesian approaches of Lan et al. (2020) and Ghosh et al. (2021), likely would be useful and insightful exercises, especially if the focus is on power studies and tests of robustness under common dependence structures in finance (e.g. spiked covariance matrices and otherwise near-singular matrices, as well as complex marginal returns distributions). The same goes for the two-sample case, where NAbC compares two sample matrices against the null hypothesis of no difference between them: comparing NAbC's results against those from some of the purportedly more generalizable competitors, like Ding et al. (2023), Bulut (2025), and Lam et al. (2025) (after covariances are converted to matrices of Pearson's correlations), and Wang et al. (2025), similarly would add to our knowledgebase.

Statistical Process Control: A full implementation of NAbC within statistical process control (SPC) monitoring frameworks would be useful to compare against potential competitors like those of Adegoke et al. (2022), ), Quessy et al. (2013), Lemyre and Quessy (2024), Ajadi et al. (2021), Bours & Steland (2020), Wang et al. (2019), Choi and Shin (2021and those reviewed in Ebadi et al. (2021). While a major focus should be on power-related metrics like average run length, special scrutiny should be placed on robustness and the (nonparametric) generalizability of NAbC vs these competitors, since these characteristics arguably are areas of weakness in the SPC literature, and where NAbC might make its most meaningful contributions.

<u>Causal Models: Recovering DAGs, and Empirical Robustification</u>: Causal models are not new (see Wright, 1921), and as we have seen above, neither are directional measures of association, which in recent times go back over a dozen years (see Zheng et al., 2012) but have direct foundations in papers from the nineteenth century (see Yule, 1897, as well as Allena and McAleerb, 2018, for a thorough analysis of Yule, 1897). Their recent use in causal frameworks already has made notable inroads (see Pascual-Marqui et al., 2024; Blömbaum et al., 2019; and MacKinnon & Lamp, 2022), and serves to validate NAbC's potential contribution in this area.

Causal model frameworks often are defined, in part or in whole, by directed acyclic graphs (DAGs), and the recovery of the 'ground truth' DAG, assuming it is rightly specified,<sup>73</sup> is one of causal modeling's fundamental tasks.<sup>74</sup> NAbC obviously is not designed to provide "all else equal" estimates of the magnitudes of treatment effects that regression approaches within causal frameworks can provide (see MacKinnon & Lamp, 2022). But it should be able to enhance covariate classification efforts for accurate

<sup>&</sup>lt;sup>73</sup> "The correct causal model is an exacting qualification, requiring a program of research with precise definition of causal effects, specification of assumptions, and sensitivity analysis for how violating assumptions affects results. Statistical analysis is useful for demonstrating associations between variables that are consistent or inconsistent with a causal model." (MacKinnon & Lamp, 2022).

<sup>&</sup>lt;sup>74</sup> Note that Czado (2025) demonstrates that vine copulas, described above as being a very flexible and effective method for *estimating* dependence structure under real-world conditions (if not for *inference* regarding all-pairwise matrices), also can be remarkably effective in the causal discovery setting. See also the innovative causal modeling approaches of Rodriguez Dominguez & Yadav (2024), and Rodriguez Dominguez (2023, 2025).

DAG recovery. For example, when using a directional dependence measure, say, Chatterjee's improved correlation (see Xia et al., 2024), we can apply NAbC twice, once with the treatment variable (X) and dependent variable (Y) and relevant covariates (V1, V2, V3) in one order in the matrix (e.g. with column and row ordering of X, V1, V2, V3, Y), and once in the reverse order in the matrix (Y, V3, V2, V1, X). The two resulting matrices will together capture all potential associations, in both directions, of all the variables. And all the cells of the two estimated dependence matrices will fully map to the relevant causal categories that make up a DAG (e.g. the confounders, colliders, mediators, independent variables, causes of X, consequences of X, causes of Y, and consequences of Y).

What NAbC could provide here is two things. First, NAbC would provide p-values associated with each of these DAG categories, for each variable, to assist in their classification. These p-values would properly take into account the entire dependence matrix, with its inherent and immutable constraints, when estimating all the pairwise relationships simultaneously. Secondly, many asymmetric dependence measures are not readily useable within regression frameworks, even when such frameworks are appropriately directional (see MacKinnon & Lamp, 2022; however, see Pascual-Margui et al., 2024, for an intriguing and innovative exception). For example, I am not aware of any regression, directional or otherwise, that allows for the use of Zhang's (2024a) combined Spearman's+Chatterjee's measure, or the asymmetric tail dependence measure of Deidda et al, (2023) when estimating (directional) covariate effects.<sup>75</sup> Yet these directional dependence measures may have more power under certain data conditions for identifying, and thus classifying, these relationships, and thus, when used by NAbC alongside existing causal frameworks, could enhance their power for accurate DAG recovery. To reemphasize, this is not a proposal to use NAbC as a standalone causal model, but rather, as a possible way that it could increase the power of an existing causal model framework in obtaining accurate DAG recovery. Of course, this begs the bigger question of whether DAGs can be used reliably within "selfreferencing open systems like capital markets" to begin with (Polakow et al., 2023). Importantly, many express strong caution, based on recent and rigorous research, regarding its application in this setting (see de Lara, 2023; Gong et al., 2024).<sup>76</sup> I propose only that NAbC can play an effective role here if the answer to this question turns out to be "yes" or "under some conditions."

<sup>&</sup>lt;sup>75</sup> However, note that Andu et al. (2021) take a very interesting approach using adaptive elastic net regression wherein Szekely's (2007) distance correlation is used to weight parameter estimates in the L1 penalty term of the regression. What's more, Pascual-Marqui et al. (2024) combine their multivariate distance-based Chatterjee correlation with the regression approach of Blömbaum et al. (2019) to extend and robustify association-based results to causal results, thus supporting the utility of using such measures in the causal modeling setting.

<sup>&</sup>lt;sup>76</sup> From Polakow et al. (2023): "The clarion call for causal reduction in the study of capital markets is intensifying. However, in self-referencing and open systems such as capital markets, the idea of unidirectional causation (if applicable) may be limiting at best, and unstable or fallacious at worst." From Gong et al. (2024): "... potential outcomes (PO) and structural causal models (SCMs) stand as the predominant frameworks. However, these frameworks face notable challenges in practically modeling counterfactuals ... we identify an inherent model capacity limitation, termed as the 'degenerative counterfactual problem', emerging from the consistency rule that is the cornerstone of both frameworks." And from De Lara (2024): "Most of the literature on causality considers the structural framework of Pearl and the potential-outcomes framework of Neyman and Rubin to be formally equivalent, and therefore interchangeably uses the do-notation and the potential-outcome subscript JD Opdyke, Chief Analytics Officer Page 78 of 92 Correlation and Beyond

Aside from DAGs, NAbC also can be directly useful to test and robustify the implementation of other causal frameworks. Consider the innovative work of Rodriguez Dominguez (2023, 2025), for example, in which each portfolio is associated with a common causal manifold, enabling future asset trajectories to be projected into its tangent space (and thus, it adroitly circumvents the time-dimension challenges that trip up many DAG-based approaches). Underlying an important part of this sophisticated approach, notably, is the trusty, association-based covariance matrix. What NAbC can provide here is the 95% confidence intervals on this matrix, under challenging, real-world financial returns data, without distortionary transformations or unrealistic assumptions, and push these upper- and lower-bound matrices through the framework to assess the (potentially highly nonlinear) effects of this estimation on the back end. The same potential for robustification can be applied when implementing the causal model of Cai et al. (2025), as they utilize the Cholesky factorization of the covariance matrix as the foundation of their algorithm, which efficiently achieves state-of-the-art performance for DAG recovery. Additionally, above I have cited Pascual-Marqui et al. (2024), who combine Chatterjee's and Szekely's measures to effectively perform directional, causal regressions. This is, in fact, a pattern: some of the best applied causal frameworks appear to be those that in no small part intelligently utilize the association-based dependence measures treated in this monograph,<sup>77</sup> and to which NAbC can be applied. So it would appear that NAbC's generality and breadth of application effectively extends to complementing and potentially enhancing the effectiveness of causal frameworks as well.

## 9. Conclusions

NAbC defines the finite sample distributions of an extremely broad class of dependence measures – all those with all-pairwise matrices that are positive definite – under challenging, real-world financial data conditions. This enables robust inference and ceteris paribus analyses in many cases where none before were possible. Motivation for NAbC's development has been the need for a method that satisfies all eight of the objectives listed below, because to date, no extant method has addressed all of these "real-world necessary" requirements simultaneously. Yet anything less than this, when modeling dependence structure in our risk and investment portfolios, fails to rise to the same level of analytical rigor as has been applied to the other parameters of these models. That is indefensible given, as is recognized in the literature, that the effects of dependence structure can be larger than many, if not all of the other parameters combined, especially when accurate models are needed most: that is, during (and just prior to) correlation breakdowns. I list again the eight objectives below for the reader's convenience:

JD Opdyke, Chief Analytics Officer

notation to write counterfactual outcomes. In this paper, we ... prove that structural counterfactual outcomes and potential outcomes do not coincide in general – not even in law." See Opdyke (2024b) for a more complete review of this literature.

<sup>&</sup>lt;sup>77</sup> Note the selection of causal drivers in Rodriguez Dominguez (2023, 2025) follows the Reichenbach Common Cause Principle (Reichenbach, 1956)), a foundational idea in probabilistic causality which emphasizes the identification of common causes through observed **correlations**.

1. NAbC remains valid under challenging, real-world data conditions, with marginal asset distributions characterized by notably different and varying degrees of serial correlation, non-stationarity, heavy-tailedness, and asymmetry.

2. NAbC can be applied to ANY positive definite dependence measure.

3. NAbC remains "estimator agnostic," that is, valid regardless of the sample-based estimator used to estimate any of the above-mentioned dependence measures.

4. NAbC provides valid confidence intervals and p-values at both the matrix level and the pairwise cell level, with analytic consistency between these two levels (i.e. the confidence intervals for all the cells define that of the entire matrix, and the same is true for the p-values; this effectively facilitates, and in many cases makes possible, granular and targeted attribution analyses).

5. NAbC provides valid confidence intervals and p-values not only for one-sample tests against matrices of fixed, assumed 'true' values, but also for two-sample tests comparing two matrices, so that we can assess inferentially whether dependence structures truly are different, for example, across different sectors or segments of our businesses.

6. NAbC provides a one-to-one quantile function, translating a matrix of all the cells' cumulative distribution function (cdf) values to a (unique) correlation/dependence measure matrix, and back again, enabling precision in reverse scenarios and stress testing, as well as informed and targeted 'what if' analyses.

7. All the above results remain valid even when selected cells in the matrix are 'frozen' for a given scenario or stress test – that is, unaffected by the scenario – thus enabling flexible, granular, and realistic scenarios.

8. NAbC remains valid not just asymptotically, i.e. for sample sizes presumed to be infinitely large, but rather, for the specific sample sizes we have in reality (for full-rank matrices with n>p), enabling inferentially reliable application in actual, real-world, non-textbook settings.

For the narrow but fundamental case of Pearson's correlation under the Gaussian identity matrix, I derive NAbC's fully analytic solution, with p-values and confidence intervals at both the cell and matrix levels (along with a measure of generalized entropy), provided in an interactive spreadsheet.

http://www.datamineit.com/JD%20Opdyke--The%20Correlation%20Matrix-Analytically%20Derived%20Inference%20Under%20the%20Gaussian%20Identity%20Matrix--02-18-24.xlsx

But way beyond Pearson's, the fully general NAbC solution presented herein satisfies all of the eight objectives listed above, simultaneously. The list of critically important, applied research that NAbC now facilitates, if not makes possible, is not only expansive, but also feasible with an ease of use and interpretability, broad range of application, scalability, and robustness not found in other more limited (spectral) methods with relatively narrow ranges of application. NAbC's utility even extends to causal modeling frameworks, further expanding its already comprehensive scope.

With NAbC, we now have a powerful, applied research tool enabling the treatment of an extremely broad class of ubiquitous dependence measures with the same level of analytical rigor as the other major parameters in our financial portfolio models. We can use NAbC in frameworks that identify, probabilistically measure and monitor, and even anticipate critically important events, such as correlation breakdowns, and mitigate and manage their effects. Correlation breakdowns are widely documented, arguably endemic characteristics of major financial markets, and their destructive potential on our attempts to estimate and forecast market behavior is difficult to overstate. Modeling efforts in this area simply cannot be effective without knowledge of, and the ability to implement and utilize, the true sampling distributions of the relevant dependence measures under real world conditions. In providing exactly these distributions, in a useable, transparent, and straightforward way, NAbC should prove to be a very useful means by which we can better understand, predict, and manage financial portfolios in our multivariate world.

## 10. References

- Abul-Magd, A., Akemann, G., and Vivo, P., (2009), "Superstatistical Generalizations of Wishart-Laguerre Ensembles of Random Matrices," *Journal of Physics A Mathematical and Theoretical*, 42(17):175207.
- Adams, R., Pennington, J., Johnson, M., Smith, J, Ovadia, Y., Patton, B., Saunderson, J., (2018), "Estimating the Spectral Density of Large Implicit Matrices" <u>https://arxiv.org/abs/1802.03451</u>.
- Adegoke, N., Ajadi, J., Mukherjee, A., and Abbasi, S., (2022), "Nonparametric Multivariate Covariance Chart for Monitoring Individual Observations," *Computers & Industrial Engineering*, Vol 167.
- Afriat, S., (1957), "Orthogonal and oblique projectors and the characterization of pairs of vector spaces," *Mathematical Proceedings of the Cambridge Philosophical Society*, 53 (4): 800–816.
- AghaKouchak, A., Easterling, D.,Hsu, K., Schubert, S., and Sorooshian, S., eds, (2013), <u>Extremes in a</u> <u>Changing Climate: Detection, Analysis and Uncertainty</u>, Ch. 6: Methods of Tail Dependence Estimation (pp.163-179), Springer Nature, Part of the book series: Water Science and Technology Library (WSTL, volume 65).
- Ajadi, j., Wong, A., Mahmood, T., and Hung, K., (2021), "A new multivariate CUSUM chart for monitoring of covariance matrix with individual observations under estimated parameter," *Quality and Reliability Engineering International*, 38(2), 834-847.

Akemann, G., Fischmann, J., and Vivo, P., (2009), "Universal Correlations and Power-Law Tails in Financial Covariance Matrices," <u>https://arxiv.org/abs/0906.5249</u>.

- Allena, D., and McAleerb, M., (2018), "Generalized Measures of Correlation for Asymmetry, Nonlinearity, and Beyond': Comment," ICAE Instituto Complutense de Análisis Económico, Working Paper 1823.
- Almog, A., and Shmueli, E., (2019), "Structural Entropy: Monitoring Correlation-Based Networks over time With Application to Financial Markets," *Scientific Reports*, 9:10832.
- Alpay, D., and Mayats-Alpay, L., (2023), "Similary Metrics, Metrics, and Conditionally Negative Definite Functions," arXiv:2307.10446v1 [math.FA].
- Andu, Y., Lee, M., and Algamal, Z., (2021), "Adaptive Elastic Net with Distance Correlation on the Group Effect and Robust of High Dimensional Stock Market Price," *Sains Malaysiana*, 50(9), 2755-2764.
- Archakov, I. and Hansen, P., (2021), "A New Parametrization of Correlation Matrices," *Econometrica*, 89(4), 1699-1715.
- Aznar, D., (2023), "Portfolio Management: A Deep Distributional RL Approach," University of Barcelona, Thesis.
- Babić, S., Ley, C., Ricci, L., and Veredas, D., (2023), "TailCoR: A new and simple metric for tail correlations that disentangles the linear and nonlinear dependencies that cause extreme comovements." *PLoS ONE*, 18(1): e0278599
- Bank for International Settlements (BIS), (2011a), Basel Committee on Banking Supervision, Working Paper 19, (1/31/11), "Messages from the academic literature on risk measurement for the trading book."
- Bank for International Settlements (BIS), (2011b), Basel Committee on Banking Supervision, "Operational Risk Supervisory Guidelines for the Advanced Measurement Approaches," June, 2011.
- Barber, R., and Kolar, M., (2018), "ROCKET: Robust Confidence Intervals via Kendall's Tau for Transelliptical Graphical Models," *The Annals of Statistics*, 46(6B), 3422-3450.
- Benjamini, Y., and Hochberg, Y., (1995), "Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing," *Journal of the Royal Statistical Society*, Series B. 57 (1): 289–300.
- Besson, O., (2025), "Covariance Matrix Estimation in the Singular Case Using Regularized Cholesky Factor," arXiv:2505.16302v1 [math.ST].

- Björck, Å. and Golub, G., (1973), "Numerical Methods for Computing Angles Between Linear Subspaces," *Mathematics of Computation*, 27 (123): 579–594.
- Blömbaum, P., Janzing, D., Washio, T., Shimizu, S., and Schölkopf, B., (2019), "Analysis of Cause-Effect Inference by Comparing Regression Errors," *Peer Journal of Computational Science*, 5:e169. doi: <u>10.7717/peerj-cs.169</u>.
- Blomqvist, N. (1950) "On a Measure of Dependence between Two Random Variables", *Annals of Mathematical Statistics*, 21(4): 593-600.
- Bongiorno, C., Challet, D., and Loeper, G., (2023), "Filtering time-dependent covariance matrices using time-independent eigenvalues," *Journal of Statistical Mechanics: Theory and Experiment*, Vol. 2023.
- Bongiorno, C., and Challet, D., (2023a), "Covariance Matrix Filtering and Portfolio Optimization: The Average Oracle vs Non-linear Shrinkage and All the Variants of DCC-NLS", arXiv:2309.17219v1 [q-fin.ST].
- Bongiorno, C., and Challet, D., (2023b), "Non-linear Shrinkage of the Price Return Covariance Matrix is Far from Optimal for Portfolio Optimization," *Finance Research Letters*, Vol. 52, 103383.
- Bouchaud, J, & Potters, M., (2015), "Financial applications of random matrix theory: a short review," <u>The</u> <u>Oxford Handbook of Random Matrix Theory</u>, Eds G. Akemann, J. Baik, P. Di Francesco.
- Bouchaud, J., (2021), "Radical Complexity," Entropy, Vol. 23.
- Bours, M., and & Steland, A., (2020), "Large-sample approximations and change testing for highdimensional covariance matrices of multivariate linear time series and factor models," *Scandinavian Journal of Statistics*, 48(2), 610-654.
- Bulut, H., (2025), "A Novel Robust Test to Compare Covariance Matrices in High-Dimensional Data," *Axioms*, 14, 247.
- Burda, Z., Jurkiewicz, J., Nowak, M., Papp, G., and Zahed, I., (2004), "Free Levy Matrices and Financial Correlations," *Physica A: Statistical Mechanics and its Applications*.
- Burda, Z., Gorlich, A., and Waclaw, B., (2006), "Spectral Properties of empirical covariance matrices for data with power-law tails," *Phys. Rev., E 74*, 041129.
- Burda, Z., Jaroz, A., Jurkiewicz, J., Nowak, M., Papp, G., and Zahed, I., (2011), "Applying Free Random Variables to Random Matrix Analysis of Financial Data Part I: A Gaussian Case," *Quantitative Finance*, Volume 11, Issue 7, 1103-1124.
- Burda, Z., and Jarosz, A., (2022), "Cleaning large-dimensional covariance matrices for correlated samples," *Phys. Rev. E*, 105, 034136.
- Cai, Y., Li, X., Sun, M., and Li, P., (2025), "Recovering Linear Causal Models with Latent Variables via Cholesky Factorization of Covariance Matrix," In: Le Thi, H.A., Le, H.M., Nguyen, Q.T. (eds)
  Advances in Data Science and Optimization of Complex Systems, ICAMCS 2024, Lecture Notes in Networks and Systems, Vol. 1311, Springer, Cham.
- Cardin, M., (2009), Multivariate Measures of Positive Definiteness," *International Journal of Contemporary Mathematical Sciences*, 4(4), 191-200.
- Carrara, C., Zambon, L., Azzimonti, D., and Corani, G., (2025), "A Novel Shrinkage Estimator of the Covariance Matrix for Hierarchical Time Series," in di Bella, E., Gioia, V., Lagazio, C., Zaccarin, S. (eds), <u>Statistics for Innovation I, SIS 2025, Italian Statistical Society Series on Advances in Statistics</u>, Springer, Cham.
- Casa, A., and Cappozzo, A., (2025), "Penalized Model-Based Clustering for Covariance Matrices," in di Bella, E., Gioia, V., Lagazio, C., Zaccarin, S. (eds), <u>Statistics for Innovation I, SIS 2025, Italian</u> <u>Statistical Society Series on Advances in Statistics</u>, Springer, Cham.
- Centofanti, F., Hubert, M., and Rousseeuw, P., (2025), "Cellwise and Casewise Robust Covariance in High Dimensions," arXiv:2505.19925v1 [state.ME].

- Ciciretti, V., and Pallotta, A., (2023), "Network Risk Parity: Graph Theory-based Portfolio Construction," *Journal of Asset Management*, 25:136–146.
- Chakraborti, A., Hrishidev, Sharma, K., and Pharasi, H., (2020), "Phase Separation and Scaling in Correlation Structures of Financial Markets," *Journal of Physics: Complexity*, 2:015002.
- Chatterjee, S., (2021), "A New Coefficient of Correlation," *Journal of the American Statistical Association*, Vol 116(536), 2009-2022.
- Chatterjee, S., (2024), "A Survey of Some Recent Developments in Measures of Association," *Probability and Stochastic Processes*, Springer Nature, Singapore.
- Chmeilowski, P., (2014), "General Covariance, the Spectrum of Riemannium and a Stress Test Calculation Formula," *Journal of Risk*, 16(6), 1-17.
- Choi, J., and Shin, D., (2021), "A self-normalization break test for correlation matrix," *Statistical Papers*, 62(5).
- Church, Christ (2012). "The asymmetric t-copula with individual degrees of freedom", Oxford, UK: University of Oxford Master Thesis, 2012.
- Cordoba, I., Varando, G., Bielza, C., and Larranaga, P., (2018), "A fast Metropolis-Hastings method for generating random correlation matrices," *IDEAL*, pp. 117-124, part of Lec Notes in Comp Sci., Vol 11314.
- Cota, R., (2019), "Shortfalls of theHierarchical Risk Parity," <u>https://www.linkedin.com/pulse/shortfalls-hierarchical-risk-parity-rafael-nicolas-fermin-cota/</u>
- Czado, C., (2025), "Vine Copula Based Structural Equation Models," *Computational Statistics and Data Analysis*, pp.453-477.
- Czado, C., and Nagler, T., (2022), "Vine Copula Based Modeling," *Annual Review of Statistics and Its Application*, pp.453-477.
- Dalio, R., (2017), Principles of Life and Work, 1<sup>st</sup> ed., Simon & Schuster.
- Dalitz, C., Arning, J., and Goebbels, S., (2024), "A Simple Bias Reduction for Chatterjee's Correlation," arXiv:2312.15496v2.
- De Lara, L., (2023), "On the (in)compatibility between potential outcomes and structural causal models and its signification in counterfactual inference," arXiv:2309.05997v3 [math.ST].
- Deidda, C., Engelke, S., and De Michele, C., (2023), "Asymmetric Dependence in Hydrological Extremes," *Water Resources Research*, Vol. 59, Issue 12.
- du Plessis, H., and van Rensburg, P., (2020), "Risk-based Portfolio Sensitivity to Covariance Estimation," *Investment Analysts Journal*, 49(3), 243-268.
- Ebadi, M., Chenouri, S., Lin, D., and Steiner, S., (2021), "Statistical monitoring of the covariance matrix in multivariate processes: A literature review," *Journal of Quality Technology*, 54(3), 269-289.
- Embrechts, P., Hofert, M., and Wang, R., (2016), "Bernoulli and Tail-Dependence Compatibility," *The Annals of Applied Probability*, Vol. 26(3), 1636-1658.
- Engle, R., Ledoit, O., and Wolf, M., (2019), "Large dynamic covariance matrices," *Journal of Business & Economic Statistics*, 37(2):363–375.
- Epozdemir, M., (2021), "Reverse Stress Testing: A critical assessment tool for risk managers and regulators," S&P Global, BLOG Aug 10, 2021, <u>https://www.spglobal.com/market-intelligence/en/news-insights/research/reverse-stress-testing-assessment-tool-risk-managers-regulators</u>
- European Banking Authority, (2013), REGULATION (EU) No 575/2013 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL, Capital Requirements Regulation (CRR), Articles 375(1), 376(3)(b), 377. https://www.eba.europa.eu/regulation-and-policy/single-rulebook/interactive-single-rulebook/12674

- Fang, Q., Jiang, Q., and Qiao, X., (2024), "Large-Scale Multiple Testing of Cross-Covariance Functions with Applications to Functional Network," arXiv:2407.19399v1 [math.ST] 28 Jul.
- Felippe, H., Viol, A., de Araujo, D. B., da Luz, M. G. E., Palhano-Fontes, F., Onias, H., Raposo, E. P., and Viswanathan, G. M., (2021), "The von Neumann entropy for the Pearson correlation matrix: A test of the entropic brain hypothesis," working paper, arXiv:2106.05379v1
- Felippe, H., Viol, A., de Araujo, D. B., da Luz, M. G. E., Palhano-Fontes, F., Onias, H., Raposo, E. P., and Viswanathan, G. M., (2023), "Threshold-free estimation of entropy from a Pearson matrix," working paper, arXiv:2106.05379v2.
- Feng, C., and Zeng, X., (2022), "The Portfolio Diversification Effect of Catastrophe Bonds and the Impact of COVID-19," working paper, https://ssrn.com/abstract=4215258
- Fernandez-Duran, J.J., and Gregorio-Dominguez, M.M., (2023), "Testing the Regular Variation Model for Multivariate Extremes with Flexible Circular and Spherical Distributions," arXiv:2309.04948v2.
- Fisher, R. A., (1915), "Frequency distribution of the values of the correlation coefficient in samples from an indefinitely large population," *Biometrika*, 10, 507-521.
- Fisher, R. A., (1921), "On the 'probable error' of a coefficient of correlation deduced from a small sample," *Metron*, 1(4), 1-32.
- Fisher, R. A., (1928), "The General Sampling Distribution of the Multiple Correlation Coefficient," Rothamsted Experimental Station, Harpenden, Herts.
- Franca, W., and Menegatto, V., (2022), "Positive definite functions on products of metric spaces by integral transforms," *Journal of Mathematical Analysis and Applications*, 514(1).
- Galeeva, R., Hoogland, J., & Eydeland, A., (2007), "Measuring Correlation Risk," publicly available manuscript.
- Garcin, M., and Nicolas, M.L.D., (2024), "Nonparametric estimator of the tail dependence coefficient: balancing bias and variance," *Statistical Papers*, 65, 4875–4913.
- Gao, M., Li, Q., (2024), "A Family of Chatterjee's Correlation Coefficients and Their Properties," arXiv:2403.17670v1 [stat.ME]
- Gebelein, H. (1941), "Das statistische problem der korrelation als variations-und eigenwertproblem und sein zusammenhang mit der ausgleichsrechnung," ZAMM-Journal of Applied Mathematics and Mechanics/Zeitschrift fur Angewandte Mathematik und Mechanik, 21 (6), 364–379.
- Genest, C., Neslehova, J., and Ghorbal, N., (2010), "Spearman's footrule and Gini's gamma: a review with complements," *Journal of Nonparametric Statistics*, 22(8), 937-954.
- Ghosh, R., Mallick, B., and Pourahmadi, M., (2021) "Bayesian Estimation of Correlation Matrices of Longitudinal Data," *Bayesian Analysis*, 16, Number 3, pp. 1039–1058.
- Gini, C., (1914), "L'Ammontare e la Composizione della Ricchezza delle Nazione," Torino: Bocca.
- Golts, M., and Jones, G., "A Sharper Angle on Optimization," ssrn.com, 1483412.
- Gong, H., Lu, C., and Zang, Y., (2024), "Distribution-consistency Structural Causal Models" arXiv:2401.15911v2 [cs.Al]
- Greenspan, A., (1999), "New Challenges for Monetary Policy," Remarks before a symposium sponsored by the Federal Reserve Federal Reserve Bank of Kansas City, Jackson Hole, Wyoming, August 27, 1999.
- Greiner, R. (1909), 'Über das fehlersystem kollektivmaßlehre', *Zeitschrift für Mathematik und Physik* **57**, 121–158,225–260,337–373.
- Gretton, A., Fukumizu, K., Teo, C., Song, L., Schölkopf, B., and Smola, A. (2007), 'A Kernel Statistical Test of Independence', *Advances in Neural Information Processing Systems*, 20.
- Grothe, O., Schnieders, J., and Segers, J., (2014), "Measuring Associatoin and Depenence Between Random Vectors," Journal of Multivariate Analysis, Vol. 123, 96-110.

JD Opdyke, Chief Analytics Officer Page 85 of 92

**Correlation and Beyond** 

Gupta, A.K. and Nagar, D.K., (2000), <u>Matrix Variate Distribution</u>, Hall/CRC, Boca Raton.

- Hamed, K., (2011), "The distribution of Kendall's tau for testing the significance of cross-correlation in persistent data," *Hydrological Sciences Journal*, 56:5, 841-853.
- Han, F., (2021), "On extensions of rank correlation coefficients to multivariate spaces," *Bernoulli News*, 28(2): 7-11.
- Hardin, J., Garcia, S., and Golan, D., (2013), "A method for generating realistic correlation matrices," *Annals of Applied Statistics*, 7(3): 1733-1762.
- Heller, R., Heller, Y., and Gorfine, M., (2013), 'A Consistent Multivariate Test of Association Based on Ranks of Distances', *Biometrika*, 100(2), 503–510.
- Hansen, P., and Luo, Y., (2024), "Robust Estimation of Realized Correlation: New Insight about Intraday Fluctuations in Market Betas," arXiv:2310.19992v1.
- Heinen, A., and Valdesogo, A., (2022), "The Kendall and Spearman rank correlations of the bivariate skew normal distribution," *Scandinavian Journal of Statistics*, 49:1669–1698.
- Heiny, J., and Yao, J., (2022), "Limiting Distributions fro Eigenvalues of Sample Correlation Matrices from Heavy-tailed Populations," arXiv:2003.03857v2 [math.PR].
- Hellton, K., (2020), "Penalized Angular Regression for Personalized Predictions," arXiv:2001.09834v2 [stat.ME].
- Higham, N., (1988), "Computing the nearest symmetric positive semi-definite matrix," *Linear algebra and its applications*, 103 (1988), 103–118.
- Higham, N., (2002), "Computing the nearest correlation matrix a problem from finance," *IMA Journal of Numerical Analysis*, 22, 329-343.
- Hirschfeld, H., (1935), "A connection between correlation and contingency," *Mathematical Proceedings* of the Cambridge Philosophical Society, 31 (4), 520–524.
- Hisakado, M. and Kaneko, T., (2023), "Deformation of Marchenko-Pastur distribution for the correlated time series," arXiv:2305.12632v1.
- Hoeffding, W. (1948). "A Non-parametric Test of Independence." *Annals of Mathematical Statistics*, 19:546–557.
- Holzmann, H., and Klar, B., (2024) "Lancaster Correlation A New Dependence Measure Linked to Maximum Correlation," arXiv:2303.17872v2 [stat.ME].
- Joarder, A., and Ali, M., (1992), "Distribution of the Correlation Matrix for a Class of Elliptical Models," Communications in Statistics – Theory and Methods, 21(7), 1953-1964.
- Johnstone, I., (2001), "On the distribution of the largest eigenvalue in principal components analysis," *The Annals of Statistics*, 29(2): 295–327, 2001.
- Jondeau, E., (2016), "Asymmetry in Tail Dependence of Equity Portfolios," *Computational Statistics & Data Analysis*, Vol 100, pp351-368.
- Jordan, Camille, (1875), "Essai sur la géométrie à dimensions," *Bulletin de la Société Mathématique de France*, 3: 103–174.
- Junker. R., Griessenberger, F., and Trutschnig, W., (2021), "Estimating scale-invariant directed dependence of bivariate distributions," *Computational Statistics & Data Analysis*, Volume 153.
- Kelly, B., Malamud, S., Pourmohammadi, M., and Trojani, F., (2024), "Universal Portfolio Shrinkage," NBER Working Paper Series, Working Paper 32004, http://www.nber.org/papers/w32004
- Ke, C., (2019), "A New Independence Measure and its Applications in High Dimensional Data Analysis," Doctoral Dissertation, University of Kentucky.
- Kendall, M. (1938), "A New Measure of Rank Correlation," *Biometrika*, 30 (1–2), 81–89.
- Kim, J., and Finger, C., (1998), "A Stress Test to Incorporate Correlation Breakdown," *Journal of Risk*, 2(3), 5-19.

- Kim, W., and Lee, Y., (2016), "A Uniformly Distributed Random Portfolio," *Quantitative Finance*, Vol. 16, No. 2, pp.297-307.
- Knyazev, A. and Argentati, M., (2002), "Principal Angles between Subspaces in an A-Based Scalar Product: Algorithms and Perturbation Estimates", *SIAM Journal on Scientific Computing*, 23 (6): 2009–2041.
- Koike, T., Lin, L., and Wang, R., (2024), "Invariant Correlation Under Marginal Transforms," *Journal of Multivariate Analysis*, 204, 105361.
- Krupskii, P., and Joe, H., (2014), "Tail-weighted Measures of Dependence," *Journal of Applied Statistics*, 42(3), 614-629.
- Kubiak, S., Weyde, T., Galkin, O., Philps, D., and Gopal, R., (2024), "Denoising Diffusion Probabilistic Model for Realistic Financial Correlation Matrices," *ICAIF '24: Proceedings of the 5th ACM International Conference on AI in Finance*, pp. 1-9. Also see https://github.com/szymkubiak/DDPMfor-Correlation-Matrices
- Kurowicka, D., (2014). "Joint Density of Correlations in the Correlation Matrix with Chordal Sparsity Patterns," *Journal of Multivariate Analysis*, 129 (C): 160–170.
- Lam, T., Dornemann, N., and Dette, H., (2025), "A New Two-Sample Test for Covariance Matrices in High Dimensions: U-Statistics Meet Leading Eigenvalues," arXiv:2506.06550v1 [math.ST].
- Lan, S., Holbrook, A., Elias, G., Fortin, N., Ombao, H., andShahbaba, B. (2020), "Flexible Bayesian Dynamic Modeling of Correlation and Covariance Matrices," *Bayesian Analysis*, 15(4), 1199–1228.
- Latif, S., and Morettin, P., (2014), "Estimation of a Spearman-Type Multivariate Measure of Local Dependence," *International Journal of Statistics and Probability*, 3(2).
- Lauria, D., Rachev, S., and Trindade, A., (2021), "Global and Tail Dependence: A Differential Geometry Approach," arXiv:2106.05865v1 [stat.AP].
- Ledoit, O., and Wolf, M., (2017), "Nonlinear shrinkage of the covariance matrix for portfolio selection: Markowitz meets goldilocks," *The Review of Financial Studies*, 30(12):4349–4388.
- Ledoit, O., and Wolf, M., (2022a), "Markowitz portfolios under transaction costs," Working paper series, Department of Economics, (420).
- Ledoit, O., and Wolf, M., (2022b), "Quadratic shrinkage for large covariance matrices," *Bernoulli*, 28(3):1519–1547.
- Lee, Y., Kim, J., Kim, W., and Fabozzi, F., (2024), "An Overview of Machine Learning for Portfolio Optimization," The Journal of Portfolio Management, 51(2), 131-148.
- Lemyre. F., and Quessy, J., (2024), "Kendall's-tau-based Inference for Gradually Changing Dependence Structures," *Statistical Papers*, Volume 65, 2033–2075.
- Lewandowski, D.; Kurowicka, D.; Joe, H. (2009). "Generating random correlation matrices based on vines and extended onion method". *Journal of Multivariate Analysis*, 100 (9): 1989–2001.
- Li, D., Cerezetti, F., and Cheruvelil, R., (2024), "Correlation breakdowns, spread positions and central counterparty margin models," Journal of Financial Market Infrastructures, 11(3), 2049-5404.
- Li, Q., (2018), "Covariance Modelling with Hypersphere Decomposition Method and Modified Hypersphere Decomposition Method," Doctoral Thesis, University of Manchester.
- Li, G., Zhang, A., Zhang, Q., Wu, D., and Zhan, C., (2022), "Pearson Correlation Coefficient-Based Performance Enhancement of Broad Learning System for Stock Price Prediction," *IEEE Transactions on Circuits and Systems—II: Express Briefs*, Vol 69(5), 2413-2417.
- Li, W. ,Yao, J., (2018), "On structure testing for component covariance matrices of a high-dimensional mixture," *Journal of the Royal Statistical Society Series B (Statistical Methodology)*, 80(2):293-318.
- Li, X., and Joe, H., (2024), "Multivariate Directional Tail-weighted Dependence Measures," *Journal of Multivariate Analysis*, Vol 203.
- JD Opdyke, Chief Analytics Officer Page 87 of 92

- Li, Y., (2025), "Large Sample Correlation Matrices with Unbounded Spectrum," *Journal of Multivariate Analysis*, 205, 105373.
- Lin, Z., and Han, F., (2023), "On Boosting the Power of Chatterjee's Rank Correlation," *Biometrika*, 110(2), 283–299.
- Lindskog, F., McNeil, A., Schmock, U., (2003), "Kendall's Tau for Elliptical Distributions," In: Bol, G., Nakhaeizadeh, G., Rachev, S.T., Ridder, T., Vollmer, KH. (eds) <u>Credit Risk. Contributions to Economics</u>, Physica-Verlag HD.
- Liu, Y., and Shang, P., (2025), "Differential Distance Correlation and Its Applications," arXiv:2507.00524 [stat.ME].
- Loretan, M., and English, W., (2000), "Evaluating Correlation Breakdowns during Periods of Market Volatility," International Finance Discussion Papers (IFDP), <u>https://www.federalreserve.gov/econres/ifdp/evaluating-correlation-breakdowns-during-periods-of-</u> market-volatility.htm and https://www.bis.org/publ/confer08k.pdf
- Lu, F., Xue, L., and Wang, Z., (2019), "Triangular Angles Parameterization for the Correlation Matrix of Bivariate Longitudinal Data," J. of the Korean Statistical Society, 49:364-388.
- MacKinnon, D., and Lamp, S., (2021), "A Unification of Mediator, Confounder, and Collider Effects," *Prev Sci.*, 22(8), 1185-1193.
- Madar, V., (2015), "Direct Formulation to Cholesky Decomposition of a General Nonsingular Correlation Matrix," *Statistics & Probability Letters*, Vol 103, pp.142-147.
- Makalic, E., Schmidt, D., (2018), "An efficient algorithm for sampling from sin(x)^k for generating random correlation matrices," arXiv: 1809.05212v2 [stat.CO].
- Maltsev, A., and Malysheva, S. (2024), "Eigenvalue Statistics of Elliptic Volatility Model with Power-law Tailed Volatility," arXiv:2402.02133v1 [math.PR].
- Manistre, J., (2008), "A Practical Concept of Tail Correlation," Society of Actuaries.
- Marchenko, A., Pastur, L., (1967), "Distribution of eigenvalues for some sets of random matrices," *Matematicheskii Sbornik*, N.S. 72 (114:4): 507–536.
- Markowitz, H., (1952), "Portfolio Selection," The Journal of Finance, 7, 77-91.
- Marti, G., (2019), "CorrGAN: Sampling Realistic Financial Correlation Matrices Using Generative Adversarial Networks," arXiv:1910.09504v1 [q-fin.ST].
- Marti, G., Goubet, V., and Nielsen, F., (2021), "cCorrGAN: Conditional Correlation GAN for Learning Empirical Conditional Distributions in the Elliptope," arXiv:2107.10606v1 [q-fin.ST]
- Marti, G. (2020) TF 2.0 DCGAN for 100×100 financial correlation matrices [Online]. Available at: <u>https://marti.ai/ml/2019/10/13/tf-dcgan-financial-correlation-matrices.html</u>. (Accessed: 17 Aug 2020)
- Martin, C. and Mahoney, M., (2018), "Implicit Self-Regularization in Deep Neural Networks: Evidence from Random Matrix Theory and Implications for Learning," *Journal of Machine Learning Research*, 22 (2021) 1-73.
- McNeil, A , Frey, R., and Embrechts, P., (2015), <u>Quantitative risk management: Concepts, techniques and</u> <u>tools</u>, Princeton University Press.

Meckes, M., (2013), "Positive Definite Metric Spaces," arXiv:1012.5863v5 [math.MG].

- Metsämuuronen, J., (2022), "Reminder of the Directional Nature of the Product-Moment Correlation Coefficient," *Academia Letters*, Article 5313.
- Meucci, A., (2010a), "The Black-Litterman Approach: Original Model and Extensions," <u>The Encyclopedia</u> of Quantitative Finance, Wiley, 2010
- Meucci, A., (2010b), "Fully Flexible Views: Theory and Practice," arXiv:1012.2848v1

Mitchell, C., Ryne, R., and Hwang, K., (2022), "Using kernel-based statistical distance to study the dynamics of charged particle beams in particle-based simulation codes," *Phys. Rev. E*, 106, 065302.

Muirhead, R., (1982), Aspects of Multivariate Statistical Theory, Wiley Interscience, Hoboken, New Jersey.

Nawroth, A., Anfuso, F. and Akesson, F., (2014), "Correlation Breakdown and the Influence of Correlations on VaR," <u>https://ssrn.com/abstract=2425515</u> or <u>http://dx.doi.org/10.2139/ssrn.2425515</u>

- Ng, F., Li, W., and Yu, P., (2014), "A Black-Litterman Approach to Correlation Stress Testing," *Quantitative Finance*, 14:9, 1643-1649.
- Niu, L., Liu, X., and Zhao, J., (2020), "Robust Estimator of the Correlation Matrix with Sparse Kronecker Structure for a High-Dimensional Matrix-variate," *Journal of Multivariate Analysis*, 177, 104598.
- Opdyke, JD, (2022), "Beating the Correlation Breakdown: Robust Inference and Flexible Scenarios and Stress Testing for Financial Portfolios," QuantMindsEdge: Alpha and Quant Investing: New Research: Applying Machine Learning Techniques to Alpha Generation Models, June 6 (available at www.DataMineit.com).
- Opdyke, JD, (2023), "Beating the Correlation Breakdown: Robust Inference and Flexible Scenarios and Stress Testing for Financial Portfolios," Columbia University, NYC–School of Professional Studies: Machine Learning for Risk Management, Invited Guest Lecture, March 20 (available at www.DataMineit.com).
- Opdyke, JD, (2024a), Keynote Presentation: "Beating the Correlation Breakdown, for Pearson's and Beyond: Robust Inference and Flexible Scenarios and Stress Testing for Financial Portfolios," QuantStrats11, NYC, March 12, 2024 (available at www.DataMineit.com).
- Opdyke, JD, (2024b), RoundTable Writeup: "Association-based vs Causal Research: the Hype, the Contrasts, and the Stronger-than-expected Complementary Overlaps," QuantStrats11, NYC, March 12, 2024 (available at www.DataMineit.com).
- Packham, N., and Woebbeking, F., (2023), "Correlation scenarios and correlation stress testing," *Journal of Economic Behavior & Organization*, Vol 205, pp.55-67.
- Pafka, S., and Kondor, I., (2004), "Estimated correlation matrices and portfolio optimization," Physica A: Statistical Mechanics and its Applications, Vol 343, 623-634.
- Palomar, D., (2025), Portfolio Optimization: Theory and Application, Cambridge University Press.
- Papenbrock, J., Schwendner, P., Jaeger, M., and Krugel, S., (2021), "Matrix Evolutions: Synthetic Correlations and Explainable Machine Learning for Constructing Robust Investment Portfolios," *Journal of Financial Data Science*, 51-69.

Parlatore, C. and Philippon, T., (2024), "Designing Stress Scenarios," NBER Working Paper No. w29901.

- Pascual-Marqui, R., Kochi, K., and Kinoshita, T., (2024), "Distance-based Chatterjee correlation: a new generalized robust measure of directed association for multivariate real and complex-valued data: arXiv:2406.16458 [stat.ME].
- Pearson, K., (1895), "VII. Note on regression and inheritance in the case of two parents," *Proceedings of the Royal Society of London*, 58: 240–242.
- Pham-Gia, T., Choulakian, V., (2014), "Distribution of the Sample Correlation Matrix and Applications," *Open Journal of Statistics*, 4(5).
- Pinheiro, J., and Bates, D., (1996), "Unconstrained parametrizations for variance-covariance matrices," *Statistics and Computing*, Volume 6, 289–296.
- Polakow, D., Gebbie, T., and Flint, E., (2023), "Epistemic Limits of Empirical Finance: Causal Reductionism and Self-Reference," arXiv:2311.16570v2 [q-fin.GN]
- Pramanik, P., (2024), "Measuring Asymmetric Tails Under Copula Distributions," *European Journal of Statistics*, 4:7.

- Puccetti, G., (2022), "Measuring linear correlation between random vectors," *Information Sciences*, Vol 607, 1328-1347.
- Quessy, J., Meriem, S., and Favre, A., (2013), "Multivariate Kendall's Tau for Change-Point Detection in Copulas," *The Canadian Journal of Statistics*, 41(1), 1-25.
- Qian, E. and Gorman, S. (2001). "Conditional Distribution in Portfolio Theory." *Financial Analysts Journal*, 44-51.
- Qin, T., and Wei-Min, H., (2024), "Epanechnikov Variational Autoencoder," arXiv:2405.12783v1 [stat.ML] 21 May 2024.
- Rebonato, R., and P. Jäckel, P., (2000) "The Most General Methodology to Create a Valid Correlation Matrix for Risk Management and Option Pricing Purposes," *Journal of Risk*, 2(2), 17-27.
- Reddi, S., Ramdas, A., Poczos, B., Singh, A., and Wasserman, L., (2015), "On Decreasing Power of Kernel and Distance based Nonparametric Hypothesis Tests in High Dimensions," arXiv:1406.2083v2 [stat.ML].
- Reichenbach, H., (1956). The Direction of Time, University of California Press.
- Rodgers, J., and Nicewander, A., (1988), The American Statistician, 42(1), 59-66.
- Rodriguez Dominguez, A., and Yadav, O., (2024), "Measuring causality with the variability of the largest eigenvalue," *Data Science in Finance and Economics*.
- Rodriguez Dominguez, A., (2023), "Portfolio optimization based on neural networks sensitivities from ass ets dynamics respect common drivers," *Machine Learning with Applications*, 11(15), 100447.
- Rodriguez Dominguez, A., (2025), "Causal Portfolio Optimization: Principles and Sensitivity-Based Solutions," arXiv:2504.05743v2 [q-fin.PM].
- Rubsamen, R., (2023), "Random Correlation Matrices Generation," https://github.com/lequant40/random-correlation-matrices-generation
- Sabato, S., Yom-Tov, E., Tsherniak, A., Rosset, S., (2007), "Analyzing systemlogs: A new view of what's important," *Proceedings, 2nd Workshop of Computing Systems ML*, pp.1–7.
- Saxena, S., Bhat, C., and Pinjari, A., (2023), "Separation-based parameterization strategies for estimation of restricted covariance matrices in multivariate model systems," *Journal of Choice Modelling*, Vol. 47, 100411.
- Schmidt, R. and Stadtmüller, U., (2006), "Non-parametric estimation of tail dependence," *Scandinavian Journal of Statistics*, 33(2):307-335.
- Sejdinovic, D., Sriperumbudur, B., Gretton, A., and Fukumizu, K., (2013) "Equivalence of Distance-Based and RKHS-Based Statistics in Hypothesis Testing," *The Annals of Statistics*, 41(5), 2263-2291.
- Sheppard, W., (1899), "On the application of the theory of error to cases of normal distribution and normal correlation," *Philosophical Transactions of the Royal Society of London (A)*, 92, 101–167.
- Siburg, K., Strothmann, C., and Weiß, G., (2024), "Comparing and quantifying tail dependence," Insurance: Mathematics and Economics, Vol 118, 95-103.
- Spearman, C., (1904), "'General Intelligence,' Objectively Determined and Measured," *The American Journal of Psychology*, 15(2), 201–292.
- Sun, F., and Huang, X., (2025), "Application of Regularized Covariance Matrices in Logistic Regression and Portfolio Optimization," Scientific Reports, 15(1).
- Szekely, G., Rizzo, M., and Bakirov, N., (2007), "Measuring and Testing Dependence by Correlation of Distances," *The Annals of Statistics*, 35(6), pp2769-2794.
- Taleb, N., (2007), <u>The Black Swan: The Impact of the Highly Improbable</u>, Random House, Inc, New York, New York.
- Taraldsen, G. (2021), "The Confidence Density for Correlation," The Indian Journal of Statistics, 2021.

- Thakkar, A., Patel, D., and Shah, P., (2021), "Pearson Correlation Coefficient-based performance enhancement of Vanilla Neural Network for Stock Trend Prediction," *Neural Computing and Applications*, 33:16985-17000.
- Tripathi, Y., Chatla, S., Chang, Y., Huang, L., and Shiefh, G., (2022), "A Nonlinear Correlation Measure with Applications to Gene Expression Data," *PLoS ONE*, 17(6): e0270270.
- Trucíos Maza, Carlos César, (2025), "Hierarchical risk clustering versus traditional risk-based portfolios: an empirical out-of-sample comparison," <u>https://ssrn.com/abstract=5247627</u> or <u>http://dx.doi.org/10.2139/ssrn.5247627</u>
- Tsay, R., and Pourahmadi, M., (2017), "Modelling structured correlation matrices," *Biometrika*, 104(1), 237-242.

Tumminello, M., Aset, T., Di Matteo, T., and Mantegna, R., (2005), "A tool for filtering information in complex systems," *Proceedings of the National Academy of Sciences*, 102(30), 10421-10426.

van den Heuvel, E., and Zhan, Z., (2022), "Myths About Linear and Monotonic Associations: Pearson's r, Spearman's ρ, and Kendall's τ," *The American Statistician*, 76:1, 44-52.

- Vanni, F., Hitaj, A., and Mastrogiacomo, E., (2024), "Enhancing Portfolio Allocation: A Random Matrix Theory Perspective," *Mathematics*, 12(9), 1389
- Veleva, E., (2017), "Generation of Correlation Matrices," *AIP Conference Proceedings* 1895, 1230008, https://doi.org/10.1063/1.5007425
- Vinod, H., (2022), "Asymmetric Dependence Measurement and Testing," arXiv:2211.16645v1 [stat.ME]. Vorobets, A., (2024), "Sequential Entropy Pooling Heuristics,"
- https://ssrn.com/abstract=3936392 or http://dx.doi.org/10.2139/ssrn.3936392
- Vorobets, A., (2025), "Portfolio Construction and Risk Management," https://ssrn.com/abstract=4807200 or http://dx.doi.org/10.2139/ssrn.4807200
- Wahba, G., (2017), "Emanuel Parzen and a Tale of Two Kernels," *Technical Report No. 1183*, University of Wisconsin, Madison.
- Wang, Z, Wu, Y., and Chu, H., (2018), "On equivalence of the LKJ distribution and the restricted Wishart distribution," arXiv:1809.04746v1.
- Wang, B., Xu, F., and Shu, L., (2019), "A Bayesian Approach to Diagnosing Covariance Matrix Shifts," *Quality and Reliability Engineering International*, 36(2).
- Wang, J., Zhu, T., and Zhang, J., (2025), "Test of the Equality of Several High-Dimensional Covariance Matrices: A Normal-Reference Approach," *MDPI: Mathematics*, 13, 295.
- Welsch, R., and Zhou, X., (2007), "Application of Robust Statistics to Asset Allocation Models," *REVSTAT–Statistical Journal*, Volume 5(1), 97–114.
- Westfall, P., and Young, S., (1993), <u>Resampling Based Multiple Testing</u>, Wiley Series in Probability and Mathematical Statistics, John Wiley & Sons, Inc., New York, New York.
- Wright, S., (1921), "Correlation and Causation," Journal of Agricultural Research, 20(7), 557-585.
- Xia, L., Cao, R., Du, J., and Chen, X., (2024), "The Improved Correlation Coefficient of Chatterjee," *Journal* of Nonparametric Statistics, pp1-17.
- Yu, P., Li, W., Ng, F., (2014), "Formulating Hypothetical Scenarios in Correlation Stress Testing via a Bayesian Framework," *The North Amer. J. of Econ. and Finance*, Vol 27, 17-33.
- Yu, S., Alesiani, F., Yu, X., Jenssen, R., and Principe, J., (2021), "Measuring Dependence with Matrix-based Entropy Functional," *The Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI-21)*.
- Yu, C., Zhu, Z., and Zhu, K., (2025), "Tensor Dynamic Conditional Correlation Model: A New Way To Pursuit 'Holy Grail of Investing'," arXiv:2502.13461v1 [q-fin.PM]
- Yule, G.U., (1897), "On the Significance of Bravais Formulæ for Regression, in the case of skew correlation," Proceedings of The Royal Society London, 477-489.

Zar, J., (1999), <u>Biostatistical Analysis</u>, 4<sup>th</sup> Ed., Prentice Hall, New Jersey.

- Zhang, Q., (2024a), "On relationships between Chatterjee's and Spearman's correlation coefficients," Communications in Statistics Theory and Methods, 54(1), 259-279.
- Zhang, Q., (2024b), "On the Extensions of the Chatterjee-Spearman Test" arXiv:2406.16859v1 [stat.ME]
- Zhang, Y., (2022), "Angle Based Dependence Measures in Metric Space," arXiv:2206.01459v1 [stat.ME].
- Zhang, Y., and Yang, S., (2023), "Kernel Angle Dependence Measures for Complex Objects," arXiv:2206.01459v2
- Zhang, W., Leng, C., and Tang, Y., (2015), "A Joint Modeling Approach for Longitudinal Studies," *Journal of the Royal Stat. Society, Series B*, 77(1), 219-238.
- Zhang, Y., Tao, J., Yin, Z., and Wang, G., (2022), "Improved Large Covariance Matrix Estimation Based on Efficient Convex Combination and Its Application in Portfolio Optimization," *Mathematics*, 10(22), 4282.
- Zhang, Y., Tao, J., Lv, Y., and Wang, G., (2023), "An Improved DCC Model Based on Large-Dimensional Covariance Matrices Estimation and Its Applications," *Symmetry*, 15, 953.
- Zhang, G., Jiang, D., and Yao, F., (2024), "Covariance Test and Universal Bootstrap by Operator Norm," arXiv:2412.20019v1 [math.ST].
- Zhangshuang, S., Gao, X., Luo, K., Bai, Y., Tao, J., and Wang, G., (2025), "Enhancing High-Dimensional Dynamic Conditional Angular Correlation Model Based on GARCH Family Models: Comparative Performance Analysis for Portfolio Optimization," *Finance Research Letters*, 106808.
- Zhao, T., Roeder, K., and Liu, H., (2014), "Positive Semidefinite Rank-based Correlation Matrix Estimation with Application to Semiparametric Graph Estimation," *Journal of Computational Graphics and Statistics*, 23(4):895–922.
- Zheng, S., Shi, N.-Z., and Zhang, Z. (2012), "Generalized Measures of Correlation for Asymmetry, Nonlinearity, and Beyond," *Journal of the American Statistical Association*, 107, 1239–1252.