

Beating the Correlation Breakdown, for Pearson's and Beyond: Robust Inference and Flexible Scenarios and Stress Testing for Financial Portfolios

JD Opdyke, Chief Analytics Officer, Partner, Sachs Capital Group Asset Management, LLC
JDOpdyke@gmail.com, August – September, 2024: LinkedIn Posts 1, 2, 3, and 4

NOTE: These posts summarize a chapter in my forthcoming monograph for Cambridge University Press.

SUMMARY:

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 a myriad of 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 (which is when we need valid inferences the most).

The following four sequential LinkedIn posts introduce a new and straightforward method – Nonparametric Angles-based Correlation (“NAbC”) – for defining the finite-sample distributions of a very wide range of dependence measures for financial portfolio analysis. These include ANY that are positive definite, such as 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), as well as measures designed to capture highly non-linear and/or cyclical dependence such as the tail dependence matrix (see Embrechts, Hofert, and Wang, 2016), 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).

No other method **simultaneously** provides 1. Finite-sample (non-asymptotic) p-values and confidence intervals for the entire pairwise matrices of ALL of these dependence measures, that are 2. analytically consistent with the individual cell-level p-values and confidence intervals of the pairwise matrices; 3. validity of 1. and 2. under real-world financial data conditions (e.g. multivariate portfolio data with marginal returns distributions differing by degree of asymmetry, heavy-tailedness, serial correlation, and (non)stationarity); 4. validity of 1. and 2. under ANY selected scenario from the pairwise matrix, where some cells are allowed to vary while others are held constant and remain unaffected by the scenario; 5. a matrix-level quantile function, allowing for one-to-one translation between unique correlation/concordance matrices and their cdf matrices; and 6. validity of 1. and 2. under any valid estimators of these dependence measures.

With NAbC, we now have a powerful, applied approach enabling us to treat an extremely broad class of widely used dependence measures with the same level of analytical rigor as the other major parameters in our (finite sample) 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. It should prove to be a very useful means by which we can better understand, predict, and manage portfolios in our multivariate world.

Beating the Correlation Breakdown, for Pearson's and Beyond: Robust Inference and Flexible Scenarios and Stress Testing for Financial Portfolios

Post 1 of 4: INTRODUCTION

JD Opdyke, Chief Analytics Officer, Partner, Sachs Capital Group Asset Management, LLC
JDOpdyke@gmail.com

NOTE: These posts summarize a chapter in my forthcoming monograph for Cambridge University Press.

Introduction

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.

Many different measures of dependence structure are widely used, including the foundational Pearson's product moment correlation matrix, rank-based measures like Kendall's Tau and Spearman's Rho, as well as measures designed to capture highly non-linear dependence such as the tail dependence matrix, Chatterjee's correlation, Lancaster's correlation, and Szekely's distance correlation and their many variants.

While dependence structure can drive portfolio results more than many other parameters in investment and risk models – sometimes even more than their combined effects – the literature provides relatively little to define the finite-sample distributions of these dependence measures under challenging, real-world data conditions. 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, and mitigating the effects of correlation breakdowns during market upheavals (which is when we need valid inferences the most).

This is the Introduction to a series of four posts that present a straightforward method– Nonparametric Angles-based Correlation (“NAbC”) – for defining the finite-sample distributions of a very wide range of dependence measures for portfolio analysis. The next post starts with a fully analytic solution for a narrow but foundational case (with a link provided to an interactive, downloadable spreadsheet), and sequentially expands NAbC's application in each post to eventually cover ANY positive definite dependence measure (including and beyond those listed above). NAbC remains highly flexible and straightforward in its implementation, yet robustly unaffected and unrestricted by the distributional challenges of real-world financial returns (see 1. in pdf below).

Motivation for NAbC's development has been its effective application for real-world financial portfolios (as opposed to textbook settings), so the solution is characterized by seven critically necessary results that no other method provides simultaneously:

1. validity under challenging, real-world data conditions, with marginal asset distributions characterized by notably varying degrees of serial correlation, non-stationarity, heavy-tailedness, and asymmetry
2. application to ANY positive definite dependence measure, including, for example, Pearson's product moment correlation, rank-based measures like Kendall's tau and Spearman's rho, the kernel-based generalization of Szekely's distance correlation, and the tail dependence matrix, among others.

Beating the Correlation Breakdown, for Pearson's and Beyond: Robust Inference and Flexible Scenarios and Stress Testing for Financial Portfolios

JD Opdyke, Chief Analytics Officer, Partner, Sachs Capital Group Asset Management, LLC
JDOpdyke@gmail.com

Post 2 of 4: Pearson's Under The Gaussian Identity Matrix

NOTE: These posts summarize a chapter in my forthcoming monograph for Cambridge University Press.

INTRODUCTION

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 these dependence measures 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 a myriad of 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 (which is when we need valid inferences the most).

This is the second in a series of four posts which introduces a new and straightforward method – Nonparametric Angles-based Correlation (“NAbC”) – for defining the finite-sample distributions of a very wide range of dependence measures for financial portfolio analysis. These include ANY that are positive definite, such as 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), as well as measures designed to capture highly non-linear dependence such as the tail dependence matrix (see Embrechts, Hofert, and Wang, 2016, and Shyamalkumar and Tao, 2020), 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).¹

This post focuses on NAbC's application to a narrow but foundational case, which is used as a baseline to greatly expand its range of application in Posts 3 and 4. The core method itself, however, remains little changed under very general conditions.

¹ Note that “positive definite” throughout these four posts refers to the dependence measure calculated on the matrix of all pairwise associations in the portfolio, that is, on a bivariate basis. Some of these dependence measures (eg Szekely's correlation) can be applied on a multivariate basis, in arbitrary dimensions, for example, to test the hypothesis of multivariate independence. But “positive definite” herein is not applied in this sense, and I explain below some of the reasons for using the dependence framework of pairwise associations.

POST 2: NAbC applied to Pearson’s under the Gaussian identity matrix.

POST 3: NAbC applied to Pearson’s under ALL correlation matrix values and ALL relevant, challenging, real-world financial returns data conditions.²

POST 4: NAbC applied to ALL matrix values and ALL positive definite measures of portfolio dependence measures, under ALL relevant, challenging, real-world financial data conditions.

PEARSON’S CORRELATION, GAUSSIAN DATA, and the IDENTITY MATRIX

We begin with Pearson’s product moment correlation matrix, the oldest and arguably most broadly used measure of dependence. 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 it remains either optimal or centrally relevant for wide-ranging purposes (e.g. 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, 2011), 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, and the data and correlation structure we presume is Gaussian data under no correlation: that is, Pearson correlation values of zero off the diagonal of the matrix as in (1).³

1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1

(1) identity matrix = for p = 4 assets

If we take two variables, such as the returns of two assets, X and Y, over a time period with n observations, we calculate Pearson’s correlation coefficient for this sample as (2):⁴

$$(2) \quad r = \frac{\sum_{i=1}^n \left(X_i - \frac{1}{n} \sum_{i=1}^n X \right) \left(Y_i - \frac{1}{n} \sum_{i=1}^n Y \right)}{\sqrt{\sum_{i=1}^n \left(X_i - \frac{1}{n} \sum_{i=1}^n X \right)^2} \sqrt{\sum_{i=1}^n \left(Y_i - \frac{1}{n} \sum_{i=1}^n Y \right)^2}} = \frac{Cov(X, Y)}{s_X s_Y}$$

² I take ‘real-world’ financial returns data to be multivariate with marginal distributions that vary notably in their degrees of heavy-tailedness, serial correlation, asymmetry, and (non-)stationarity.

³ Note, of course, that a zero value for Pearson’s correlation does not imply independence, but independence does imply a zero value for Pearson’s correlation.

⁴ Recall that Pearson’s requires that the first two moments (the mean and the variance) of the distributions of X and Y are finite.

For the corresponding matrix of all pairwise correlations, we have:

$$R = \begin{bmatrix} 1 & r_{1,2} & r_{1,3} & r_{1,4} \\ r_{2,1} & 1 & r_{2,3} & r_{2,4} \\ r_{3,1} & r_{3,2} & 1 & r_{3,4} \\ r_{4,1} & r_{4,2} & r_{4,3} & 1 \end{bmatrix}, \text{ with the usual, following characteristics:}$$

- i. Symmetry: $r_{i,j} = r_{j,i}$
- ii. Unit diagonal entries: $r_{i=j} = 1$
- iii. Bounded non-diagonal entries: $-1 \leq r_{i,j} \leq 1$
- iv. The matrix is positive definite, i.e. all eigenvalues $\lambda_i > 0$

For completeness and for reference throughout this post, we define eigenvalues here:

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

$$\det(\lambda I - R) = 0, \text{ then } \det(\lambda I - R)v = 0, \text{ where } I \text{ is the identity matrix and } \det \text{ is the determinant}$$

The eigenvalue can be thought of as the magnitude of the (portfolio) variance in the direction of the eigenvector. Also note that with iii. above, this range can be tighter under specific circumstances, such as for equicorrelation matrices where $-1/(p-1) \leq r \leq 1$, $p = \dim(r)$.

ANGLE VALUES vs CORRELATION VALUES

The key to the NAbC approach rests in its use of the ANGLE θ between the two mean-centered data vectors of X and Y , as opposed to directly and only using of the values of the correlations themselves. 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 (4):⁵

$$(4) \quad \cos(\theta) = \frac{\text{inner product}}{\text{product of norms}} = \frac{\langle \mathbf{X}, \mathbf{Y} \rangle}{\|\mathbf{X}\| \|\mathbf{Y}\|} = \frac{\sum_{i=1}^N (X_i - E(X))(Y_i - E(Y))}{\sqrt{\sum_{i=1}^N (X_i - E(X))^2} \sqrt{\sum_{i=1}^N (Y_i - E(Y))^2}} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \rho, \text{ with } 0 \leq \theta \leq \pi$$

⁵ While r typically is used to represent Pearson's calculated on a sample, ρ often is used to represent Pearson's calculated on a population.

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 (4), as long as the matrix is symmetric-positive-definite,⁶ is well established in the literature (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 (5)-(7) and in code in Table A. The steps for translating between correlations and angles, in both directions, are 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 and Wang, 2015, but note a typo in the formula in Pourahmadi and 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 hard-coded function in most statistical and mathematical software. The relevant formulae are included below for completeness.

(5) A correlation matrix R will be real, symmetric positive-definite, 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:

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

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; or in reverse, we obtain a correlation/dependence value from each cell's angle value (see Pourahmadi and Wang, 2015, as well as

⁶ Note that this is true not only for Pearson's, but also for all relevant dependence measures in this setting, as will be discussed in Posts 3 and 4.

Rapisarda et al., 2007, for a meticulous derivation 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.

(6)

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

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 $\theta_{i,j}$, we have:

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

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

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

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

TABLE A:

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

The above all is well-established and straightforward, and demonstrates, as we know intuitively, that scale does not (and should not) matter when it comes to dependence measures;⁷ 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 post-sample verification of positive definiteness, and subsequent resampling when this requirement is violated⁸ (see Makalic and Schmidt, 2018, Cordoba et al. 2018, and Papenbrock et al., 2021). This inefficiency grows very rapidly with the size of the matrix/portfolio, as shown in the ratio below in (8) (see Bohn and Hornik, 2024 and Pourahmadi and Wang, 2015).

⁷ Scale invariance is widely proved and cited for Pearson’s, Kendall’s, and Spearman’s (see Xu et al., 2013, and Schreyer et al., 2017 examples).

⁸ As shown below, this approach also much more straightforward, not to mention more generalizable, than the other, more complex sampling algorithms that have been proposed, such as the vine and extended onion algorithms of Lewandowski et al. (2009), the Metropolis-Hastings and Metropolis algorithms of Cordoba et al. (2018), and the restricted Wishart distribution approach of Wang et al. (2018).

$$(8) \quad \Pr(\text{rand "R" } \sim \text{PosDef}) = X = \frac{\prod_{j=1}^{p-1} \left[\sqrt{\pi} \Gamma\left(\frac{j+1}{2}\right) \right]^j}{2^{p(p-1)/2}} < \prod_{j=1}^{p-1} \left[\frac{\sqrt{\pi}}{2} \right]^j = \left[\frac{\sqrt{\pi}}{2} \right]^{p(p-1)/2} ; \lim_{p \rightarrow \infty} [X] = 0$$

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 than 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. 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 time, to solve the cube. The valid solution must be obtained by (always) following the rules governing shifts in the cube, each of which affects many of the individual cubes (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, simply by following steps A., B., and C., and C., B., and A., above.

Another crucial characteristic of these angles is that the distribution of each is **independent** with respect to those of the others (see Pourahmadi and Wang, 2015, Tsay and Pourahmadi, 2017, and Ghosh et al., 2020). 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 the dependence between the two variables (see Fernandez-Duren & Gregorio-Dominguez, 2023, and Zhang & Songshan, 2023, as well as Opdyke, 2024). This will be covered more extensively in subsequent posts.

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

FULLY ANALYTIC ANGLES DENSITY – EFFICIENT SAMPLE GENERATION

Once we have the matrix of angles, one for each pairwise correlation (dependence measure), 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 (9) (see Cordoba, 2018, Pourahmadi and Wang, 2015, Lewandowski et al., 2009).

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

(9)

We see directly from (6) that $\sin^k(x)$, suitably normalized in (10), satisfies this requirement (see Pourahmadi and Wang, 2015, and Makalic and Schmidt, 2018).

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

(10)

Although not mentioned in Makalic and Schmidt (2018), importantly note that $k = \# \text{columns} - \text{column\#}$ (so for the first column of a $p=10 \times 10$ matrix, $k=9$; for the second column, $k=8$, etc.).

However, we need both the cumulative distribution function (cdf) and its inverse, the quantile function, to make use of this 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 values. Yet regarding the cdf corresponding to (10) above, Makalic and 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 (2020) derived an analytic, non-recursive expression of the cdf below in (11).

(11)

$$F_x(x; k) \sim \frac{1}{2} - c_k \cdot \cos(x) \cdot {}_2F_1\left[\frac{1}{2}, \frac{1-k}{2}; \frac{3}{2}; \cos^2(x)\right] \quad \text{for } x < \frac{\pi}{2},$$

$$\sim \frac{1}{2} + c_k \cdot \cos(x) \cdot {}_2F_1\left[\frac{1}{2}, \frac{1-k}{2}; \frac{3}{2}; \cos^2(x)\right] \quad \text{for } x \geq \frac{\pi}{2}$$

where the Gaussian hypergeometric function ${}_2F_1[a, b; c; r] = \sum_n \frac{(a)_n (b)_n}{(c)_n} \cdot \frac{r^n}{n!}$

where $(h)_n = h(h+1)(h+2) \cdots (h+n-1)$, $n \geq 1$, $(h)_0 = 1$, and $|r| < 1$, $c \neq 0, -1, -2, \dots$

Interestingly, the Gaussian hypergeometric function makes many appearances in this setting,⁹ but it is admittedly cumbersome mathematically. But Opdyke (2022, 2023, and 2024) has shown that (11) can be simplified further, based on some arguably obscure hypergeometric identities:

⁹ The (Gaussian) hypergeometric function appears in derivations of the distribution of individual correlations (see Muirhead, 1982, and Taraldsen, 2021), moments of the spectral distribution under some conditions (see Adams et al. 2018, and <https://reference.wolfram.com/language/ref/MarchenkoPasturDistribution.html>), and in the definition of positive definite functions (see Franca & Menegatto, 2022).

(12)

For $c = a + 1$ and $0 < r < 1$ simultaneously, which holds in this setting, we have ${}_2F_1[a, b; c; r] = B(r; a, 1 - b)(a/r^a)$

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

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

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

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:

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

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

Let $p = \Pr(x \geq X)$. Then for $x < \frac{\pi}{2}$,

$$p = \frac{1}{2} - \left(\frac{1}{2}\right) \cdot F_{Beta}\left[\cos^2(x); \frac{1}{2}, \frac{1+k}{2}\right]$$

$$-2p = -1 + F_{Beta}\left[\cos^2(x); \frac{1}{2}, \frac{1+k}{2}\right]$$

$$1 - 2p = F_{Beta} \left[\cos^2(x); \frac{1}{2}, \frac{1+k}{2} \right]$$

$$F_{Beta}^{-1} \left(1 - 2p; \frac{1}{2}, \frac{1+k}{2} \right) = \cos^2(x)$$

$$\sqrt{F_{Beta}^{-1} \left(1 - 2p; \frac{1}{2}, \frac{1+k}{2} \right)} = \cos(x)$$

$$\arccos \left(\sqrt{F_{Beta}^{-1} \left(1 - 2p; \frac{1}{2}, \frac{1+k}{2} \right)} \right) = x$$

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

We must reflect the symmetric angle density for $p \geq 0.5$, so we have

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

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

So now from (12) above we have for the angles distribution, under the Gaussian identity matrix, for the first time together, the pdf, cdf, and quantile function:

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

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

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

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

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

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.¹⁰

The above (12) now provides a fully analytic solution,¹¹ and in fact is so straightforward as to be readily implemented in a spreadsheet, and one is provided for download via the link below and included as a file upload in this Post 2.

<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 and Schmidt (2018), the straightforward approach of inverse transform sampling CAN be used in this setting, for this narrow case, to efficiently sample the correlation matrix. And in fact, this is the most efficient way to sample it. Roman (2023) has compared Makalic and Schmidt (2018) to the above method (defined in Opdyke, 2022, 2023, and 2024) and obtained over 30% decrease in runtime.

But 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 $[1 - (\text{CDF value})]$ for the upper tail (13), 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 (14)

(13) one-sided p-value = $F_x(x; k)$ or $1 - F_x(x; k)$; two-sided p-value = 2 x one-sided p-value

(14) $F^{-1}(\alpha/2; k)$ and $F^{-1}(1 - \alpha/2; k)$ 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.

¹⁰ See Xi'an, March, 2018 (<https://stats.stackexchange.com/questions/331253/draw-n-dimensional-uniform-sample-from-a-unit-n-1-sphere-defined-by-n-1-dime/331850#331850> and <https://xianblog.wordpress.com/2018/03/08/uniform-on-the-sphere-or-not/>).

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 (<https://stats.stackexchange.com/users/7224/xian>).

¹¹ Note that we use the term 'analytic' as opposed to 'closed-form' because we are 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.

FINITE-SAMPLE DISTRIBUTION OF THE CORRELATION MATRIX

As mentioned previously, a key characteristic of the angles distributions is that they are independent vis-à-vis each other, which makes defining their multivariate distribution straightforward: it is simply the product of all the angles' pdf's. But what does this mean for the p-value and confidence intervals for the entire matrix? Given the null hypothesis of the identity matrix (under the presumption of Gaussian data here), the (2-sided) p-value of the entire matrix is simply one minus the probability of no false positives, which is the definition of controlling the family-wise error rate (FWER) of the matrix (15).

$$(15) \text{ matrix (2-sided) } pvalue = \left[1 - \prod_{i=1}^{p(p-1)/2} (1 - p-value_i) \right] \text{ where } p-value_i \text{ is the 2-sided p-value.}$$

Again, because the cell-level distributions are independent, their p-values are independent, and otherwise statistically more powerful approaches for calculating the FWER that rely on, for example, resampling methods (Westfall and Young, 1993, and Romano and Wolf, 2016), do not apply here. In other words, they provide no power gain over (15) because under independence, there is no dependence structure for them to exploit. So the straightforward calculation above in (15) is, by definition, the most powerful for FWER control.

Similarly, calculation of the confidence interval for the entire matrix (16) is essentially the same as that of the p-value, but of course it is divided in half to account for each tail, and the root of the critical values is taken, rather than the product. Otherwise, the calculations are identical to obtain the critical alphas for these 'simultaneous confidence intervals.'

$$(16) \alpha_{crit-simult-LOW} = \left(1 - [1 - \alpha/2]^{(1/[p(p-1)/2])} \right) \text{ and } \alpha_{crit-simult-HIGH} = 1 - \alpha_{crit-simult-LOW}$$

These critical alphas, when inserted as values in the cdf functions, 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 in this case is the identity matrix. Independence of the angles distributions again makes these simultaneous confidence intervals very straightforward to calculate.

Importantly, again note that because we derived the quantile (inverse cdf) function in (12) 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. 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, with visible formulae corresponding to each step of these calculations for full transparency.

P-VALUES vs ENTROPY: USING P-VALUES AS A MEASURE OF MATRIX DISTANCE/DISPERSION

Before describing how NAbC, unlike competing methods, enables granular, highly flexible scenarios for dependence measures (key result #6 of POST 1), lets take a moment to examine the meaning and implications of the cell-level p-values derived above in (12) and (13).

The (2-sided) p-value of (13) provides what can be viewed as a distance metric that has some advantages over more traditional distance metrics, such as norms. Some commonly used norms in this setting for measuring correlation 'distances' are listed below in (17).

$$(17) \quad \|x\| = \left(\sum_{i=1}^d |x_i|^m \right)^{1/m}$$

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

All of these norms measure absolute distance from a presumed or baseline correlation value. But the range of all relevant and widely used dependence measures is bounded, either from -1 to 1 or 0 to 1 , and 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 0.01 from an original or presumed correlation value of, say, 0.97 , means something very different than the same shift from 0.07 . NAbC attributes probabilistic MEANING to these two different cases, while a norm would 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 these cell-level p-values from (13) is the natural log of the product of the p-values, dubbed 'LNP' in (18) below:

$$(18) \quad \text{"LNP"} = \ln \left(\prod_{i=1}^q p\text{-value}_i \right) = \sum_{i=1}^q \ln [p\text{-value}_i] \text{ where } q = p(p-1)/2 \text{ and } p\text{-value}_i \text{ is 2-sided.}$$

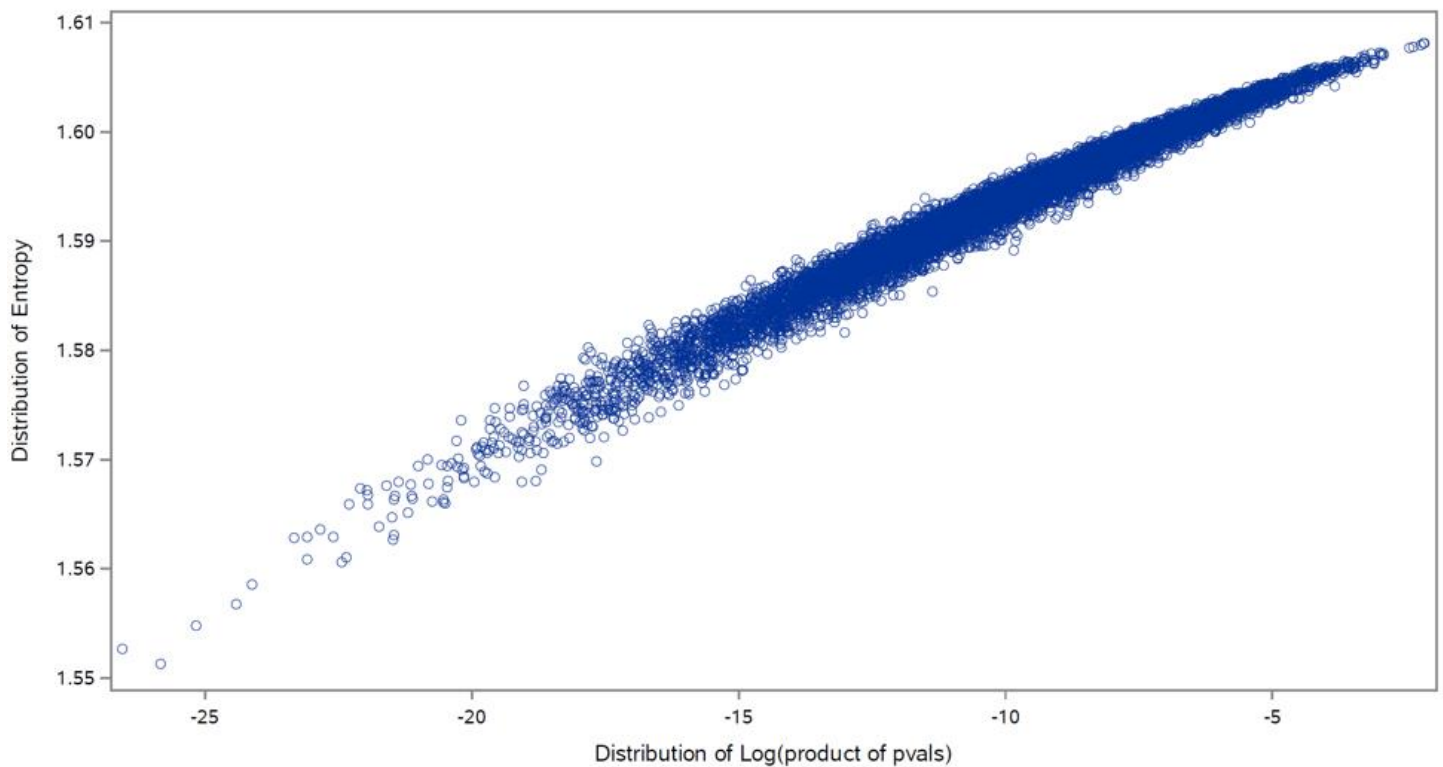
Intriguingly, LNP shows a remarkable correspondence with the entropy of the correlation matrix, defined by Felipe et al. (2021 and 2023) as (19) below:

$$(19) \quad \text{Entropy} = Ent(R/p) = - \sum_{j=1}^p \lambda_j \ln(\lambda_j)$$

where R is the sample correlation matrix and λ_j are the p eigenvalues of the correlation matrix after it is scaled by its dimension, R/p . (Note that this result (19), like NAbC, is valid for ANY positive definite measure of dependence, not just Pearson's, as will be discussed in POSTs 3 and 4).

Graph 1 compares LNP to the entropy of the correlation matrix in 10,000 simulations under the Gaussian identity matrix. The resulting Pearson's correlation between them is just shy of 0.99 .

GRAPH 1: Identity Matrix Simulations -- LNP v Entropy



What makes this result worthy of further investigation is that it indicates a broad and useful generalizability of LNP. As will be discussed in Posts 3 and 4, LNP can be calculated for ANY correlation/dependence matrix, not just the identity matrix. Entropy, on the other hand, can be calculated only with reference to the identity matrix as a baseline. Yet the correspondence of LNP to entropy under this specific case speaks to LNP's natural interpretation as a meaningful measure of deviation/distance/dispersion, 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. This topic will be treated in subsequent posts, but is mentioned here as it provides further validation of this approach under this narrow case, as well as much more general conditions.

GRANULAR, HIGHLY FLEXIBLE SCENARIOS

I have taken a very granular, 'bottom up' approach to defining the finite-sample distribution of the correlation matrix here, based on the distributions of the individual correlation cells. In addition to analytical consistency, this provides a flexibility that other approaches, such as those based on the spectrum of the dependence measure's matrix, cannot provide, because with only p eigenvalues, they simply are at the wrong level of aggregation to flexibly vary (or freeze) the $p(p-1)/2$ cells for different scenarios.¹² Correlation (dependence) matrices under a tech market bubble (2000) vs those under a

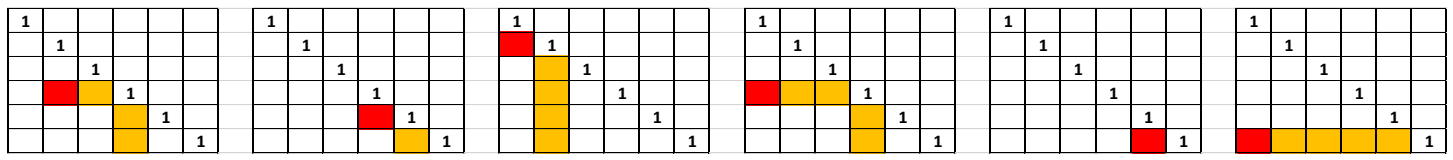
¹² Importantly, when we deviate from the identity matrix (covered in Posts 3 and 4), spectral distributions are far less robust than angles distributions. The latter are bounded, typically unimodal, smooth, not excessively asymmetric, and stable as JD Opdyke, Chief Analytics Officer Page 14 of 21 Beating the Correlation Breakdown: Post 2 of 4

housing bubble (2008) vs those under Covid (2020) will change very different individual cells, 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. In other words, while correlation ‘breakdowns’ will occur under all of these extreme conditions, the granular nature of pairwise association matrices ensures that the fundamentally different nature of these breakdowns will be captured and reflected empirically in all related analyses. The only way to flexibly and realistically model this is at the most granular level – that of the individual correlation cells.

Fortunately, when using NAbC, several results allow for this. First, 1. independence of the angles distributions allows us to vary individual cells. 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 and Wang, 2015, and Lewandowski et al., 2009). Third, 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.

Focus only on the lower triangle of the correlation matrices below in Graphs 2-4, since the upper triangle is just its reflection. Note again that using NAbC, we only perturb angles. We never perturb the correlation values directly. We must always convert to angles, perturb the angle values (in this narrow case for this Post 2, using inverse transform), 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 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 2 below.¹³

GRAPH 2: Mechanics of Matrix Multiplication



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, according to the fill order in Graph 3 below.

matrices approach singularity; in contrast, the former remain unbounded, often are multi-modal, and are far less stable as dependence matrices approach singularity, which is more the rule than the exception when portfolio sizes are not small.

¹³ Note that not all of these (orange) cells will necessarily change if values of zero are involved, but none OTHER than these (orange) cells CAN change when only the red cell changes.

GRAPH 3: Rightmost Triangle Fill Order

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 cells 1,2 and 1,3 and 2,3 and 4,3 are now in the rightmost triangle of the lower triangular matrix, in the fill order shown above.

GRAPH 4: Example of Mechanics of Matrix Multiplication Applied to Rightmost Triangle Fill Order

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

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

Changes in Corresponding Angles Cells ONLY change Same in Resorted Matrix

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

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, 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 are not targeted): they only must fill the rightmost triangle according to the order of the middle 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. In contrast, changes to cell “4,3” will affect the values of the other orange cells. Readers are encouraged to test this in the attached spreadsheet.

So we can exploit these four simultaneous conditions – 1. independence of the angles distributions; 2. (correlation) distribution invariance to row and column order; 3. the mechanics of matrix multiplication; and 4. the granular, cell-level geometry of NAbC – to obtain great flexibility in defining scenarios wherein some cells vary and some do not. 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 analyses cannot do in this setting.

The only arguable drawback of this approach is that it can be limited by the structure of measuring dependence in pairwise associations. As shown in Graph 4 above, for the $p=5$ asset matrix, there are only $p!$ (ie $5!=120$) ways to sort the rows and columns, but there are $[p(p-1)/2]!$ (ie $15!=1,307,674,368,000$) ways to sort the 15 cells. The matrix obviously cannot accommodate freely sorting the individual cells in this way because it breaks the pairwise structure of the matrix. Some scenarios, therefore, could conceivably be required to include for perturbation some few additional cells in the 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. But dealing with these potential cases appears to be well worth the price of the unmatched flexibility that this approach provides, not to mention the other advantages it maintains over more complex, strictly multivariate dependence structures. For usage with actual market data, the latter typically are much more difficult to estimate with the same levels of accuracy, let alone to manipulate for purposes of intervention or mitigation. In contrast, pairwise associations are directly identifiable, typically more easily and accurately estimated, and interventions more targeted and transparent.

CONCLUSION

In Post 1 I listed the seven characteristics of the full NAbC solution, and for completeness I list them here below:

1. validity under challenging, real-world financial data conditions, with marginal asset distributions characterized by notably different degrees of serial correlation, non-stationarity, heavy-tailedness, and asymmetry
2. application to ANY positive definite dependence measure, including, for example, Pearson's product moment correlation, rank-based measures like Kendall's tau and Spearman's rho, the kernel-based generalization of Szekely's distance correlation, and the tail dependence matrix, among others.
3. it remains "estimator agnostic," that is, valid regardless of the sample-based estimator used to estimate any of the above-mentioned dependence measures
4. it provides valid confidence intervals and p-values at both the matrix-level and the pairwise cell-level, with analytic consistency between these two levels (ie 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 attribution analyses)

5. it provides a one-to-one quantile function, translating a matrix of all the cells' cdf values to a (unique) correlation (dependence measure) matrix, and back again, enabling precision in reverse scenarios and stress testing
6. all the above results remain valid even when selected cells in the matrix are 'frozen' for a given scenario or stress test, enabling granular and realistic scenarios
7. it remains valid not just asymptotically, ie for sample sizes presumed to be infinitely large, but rather, for the specific sample sizes we have in reality, enabling reliable application in actual, imperfect, non-textbook settings

This Post 2 covers 4, 5, 6, and 7 above. The next Post 3 expands NAbC to cover 1 as well, using exactly the same angles-based framework. The utility of using the foundational, but undeniably narrow case of the Gaussian identity matrix in this Post 2 rests in establishing the framework and proving out the mechanics of how and why it works, so that we can expand its range of application to the real-world cases of challenging, financial portfolio data. Finally, in Post 4, I expand NAbC's range of application to Characteristics 2 and 3 above, not only to challenging, real-world data conditions, but also simultaneously beyond Pearson's to ALL positive definite measures of dependence.

REFERENCES

Adams, R., Pennington, J., Johnson, M., Smith, J, Ovadia, Y., Patton, B., Saunderson, J., (2018), "Estimating the Spectral Density of Large Implicit Matrices" <https://arxiv.org/abs/1802.03451>.

Askitis, D., (2017), "Asymptotic expansions of the inverse of the Beta distribution," <https://arxiv.org/abs/1611.03573>

BIS, Basel Committee on Banking Supervision, Working Paper 19, (1/31/11), "Messages from the academic literature on risk measurement for the trading book."

Chatterjee, S., (2021), "A New Coefficient of Correlation," *Journal of the American Statistical Association*, Vol 116(536), 2009-2022.

Digital Library of Mathematical Functions (DLMF), Section 8.17.ii, Hypergeometric Representations, National Institute of Standards and Technology (NIST), Handbook of Mathematical Functions, US Department of Commerce, by Cambridge University Press, Online Version 1.2.1; Release date 2024-06-15 (<https://dlmf.nist.gov/8.17#ii>).

Embrechts, P., Hofert, M., and Wang, R., (2016), "Bernoulli and Tail-Dependence Compatibility," *The Annals of Applied Probability*, Vol. 26(3), 1636-1658.

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.

- 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).
- Gao, M., and Li, Q., (2024), “A Family of Chatterjee’s Correlation Coefficients and Their Properties,” arXiv:2403.17670v1 [stat.ME].
- Ghosh, R., Mallick, B., and Pourahmadi, M., (2021) “Bayesian Estimation of Correlation Matrices of Longitudinal Data,” *Bayesian Analysis*, 16, Number 3, pp. 1039–1058.
- Holzmann, H., and Klar, B., (2024) “Lancaster Correlation - A New Dependence Measure Linked to Maximum Correlation,” arXiv:2303.17872v2 [stat.ME].
- Kendall, M. (1938), "A New Measure of Rank Correlation," *Biometrika*, 30 (1–2), 81–89.
- 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.
- Makalic, E., Schmidt, D., (2018), “An efficient algorithm for sampling from $\sin(x)^k$ for generating random correlation matrices,” arXiv: 1809.05212v2 [stat.CO].
- Meucci, A., (2010a), “The Black-Litterman Approach: Original Model and Extensions,” [The Encyclopedia of Quantitative Finance](#), Wiley, 2010
- Meucci, A., (2010b), “Fully Flexible Views: Theory and Practice,” arXiv:1012.2848v1
- Muirhead, R., (1982), [Aspects of Multivariate Statistical Theory](#), Wiley Interscience, Hoboken, New Jersey.
- Opdyke, JD, (2020), “Full Probabilistic Control for Direct & Robust, Generalized & Targeted Stressing of the Correlation Matrix (Even When Eigenvalues are Empirically Challenging),” QuantMinds/RiskMinds Americas, Sept 22-23, Boston, MA.
- 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.
- 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.
- Opdyke, JD, (2024), Keynote Address: “Beating the Correlation Breakdown, for Pearson’s and Beyond: Robust Inference and Flexible Scenarios and Stress Testing for Financial Portfolios,” QuantStrats11, NYC, March 12.
- Pafka, S., and Kondor, I., (2004), “Estimated correlation matrices and portfolio optimization,” *Physica A: Statistical Mechanics and its Applications*, Vol 343, 623-634.

- 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.
- 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.
- Pinheiro, J. and Bates, D. (1996), "Unconstrained parametrizations for variance-covariance matrices," *Statistics and Computing*, Vol. 6, 289–296.
- Pourahmadi, M., Wang, X., (2015), "Distribution of random correlation matrices: Hyperspherical parameterization of the Cholesky factor," *Statistics and Probability Letters*, 106, (C), 5-12.
- Qian, E. and Gorman, S. (2001). "Conditional Distribution in Portfolio Theory." *Financial Analysts Journal*, 44-51.
- Rapisarda, F., Brigo, D., & Mercurio, F., (2007), "Parameterizing Correlations: A Geometric Interpretation," *IMA Journal of Management Mathematics*, 18(1), 55-73.
- Rebonato, R., and Jackel, P., (2000), "The Most General Methodology for Creating a Valid Correlation Matrix for Risk Management and Option Pricing Purposes," *Journal of Risk*, 2(2)17-27.
- Romano, J., and Wolf, M., (2016), "Efficient computation of adjusted p-values for resampling-based stepdown multiple testing," *Statistics & Probability Letters*, Vol 113, 38-40.
- Rubsamen, Roman, (2023), "Random Correlation Matrices Generation," <https://github.com/lequant40/random-correlation-matrices-generation>
- Schreyer, M., Paulin, R., and Trutschnig, W., (2017), "On the exact region determined by Kendall's tau and Spearman's rho," arXiv: 1502:04620.
- 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.
- Sharma, D., and Chakrabarty, T., (2017), "Some General Results on Quantile Functions for the Generalized Beta Family," *Statistics, Optimization and Information Computing*, 5, 360-377.
- Shyamalkumar, N., and Tao, S., (2020), "On tail dependence matrices: The realization problem for parametric families," *Extremes*, Vol. 23, 245–285.
- Spearman, C., (1904), "'General Intelligence,' Objectively Determined and Measured," *The American Journal of Psychology*, 15(2), 201–292.
- Szekely, G., Rizzo, M., and Bakirov, N., (2007), "Measuring and Testing Dependence by Correlation of Distances," *The Annals of Statistics*, 35(6), pp2769-2794.
- 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.

Tsay, R., and Pourahmadi, M., (2017), "Modelling structured correlation matrices," *Biometrika*, 104(1), 237–242.

Xu, W., Hou, Y., Hung, Y., and Zou, Y., (2013), "A Comparative Analysis of Spearman's Rho and Kendall's Tau in Normal and Contaminated Normal Models," *Signal Processing*, 93, 261–276.

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.

Wang, Z, Wu, Y., and Chu, H., (2018), "On equivalence of the LKJ distribution and the restricted Wishart distribution," arXiv:1809.04746v1.

Weisstein, E., (2024a), "Beta Distribution." From *MathWorld--A Wolfram Web Resource*.
<https://mathworld.wolfram.com/BetaDistribution.html>

Weisstein, E., (2024b), "Regularized Beta Function." From *MathWorld--A Wolfram Web Resource*.
<https://mathworld.wolfram.com/RegularizedBetaFunction.html>

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), Resampling Based Multiple Testing, Wiley Series in Probability and Mathematical Statistics, John Wiley & Sons, Inc., New York, New York.

Zhang, Y., and Songshan, Y., (2023), "Kernel Angle Dependence Measures for Complex Objects," arXiv:2206.01459v2

Beating the Correlation Breakdown, for Pearson's and Beyond: Robust Inference and Flexible Scenarios and Stress Testing for Financial Portfolios

JD Opdyke, Chief Analytics Officer, Partner, Sachs Capital Group Asset Management, LLC
JDOpdyke@gmail.com

Post 3 of 4: Pearson's Under ANY Values and Real-World Financial Data Conditions

NOTE: These posts summarize a chapter in my forthcoming monograph for Cambridge University Press.

Introduction

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 these 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 a myriad of 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 (which is when we need valid inferences the most).

This is the third in a series of four posts which introduces a new and straightforward method – Nonparametric Angles-based Correlation (“NAbC”) – for defining the finite-sample distributions of a very wide range of dependence measures for financial portfolio analysis. These include ANY that are positive definite, such as 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), as well as measures designed to capture highly non-linear dependence such as the tail dependence matrix (see Embrechts, Hofert, and Wang, 2016, and Shyamalkumar and Tao, 2020), 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).¹

This Post 3 expands NAbC's application from Pearson's under the Gaussian Identity Matrix to Pearson's under ANY correlation values and challenging, real-world financial data conditions. Post 4 will expand its range of application even further, beyond Pearson's to ANY positive definite dependence measure.

¹ Note that “positive definite” throughout these four posts 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 these dependence measures (eg Szekely's correlation) can be applied on a multivariate basis, in arbitrary dimensions, for example, to test the hypothesis of multivariate independence. But “positive definite” herein is not applied in this sense, and I explain below some of the reasons for using the dependence framework of pairwise associations, which is highly flexible, and allows for more precise attribution and intervention analyses.

POST 1: NAbC introduced.

POST 2: NAbC applied to Pearson's under the Gaussian identity matrix.

POST 3: NAbC applied to Pearson's under ALL correlation matrix values and ALL relevant, challenging, real-world financial returns data conditions.²

POST 4: NAbC applied to ALL matrix values and ALL positive definite measures of portfolio dependence measures, under ALL relevant, challenging, real-world financial data conditions.

Review of Post 2: Correlations and Angles

To briefly review from Post 2, I defined and reviewed the relationship between the correlation cells in a Pearson's correlation matrix and the angles of their corresponding pairwise data vectors. There exists an angle value for every correlation value in the matrix. For a single, bivariate correlation, this can be seen directly via the widely used cosine similarity in (1),³ but the matrix analog also is well established in the literature as shown in (2.a) and (2.b) (see Rebonato & Jaeckel, 2000, Rapisarda et al., 2007, and Pourahmadi and Wang, 2015, but note a typo in the formula in Pourahmadi and Wang, 2015 corresponding to (2.b) below):

$$(1) \quad \cos(\theta) = \frac{\text{inner product}}{\text{product of norms}} = \frac{\langle \mathbf{X}, \mathbf{Y} \rangle}{\|\mathbf{X}\| \|\mathbf{Y}\|} = \frac{\sum_{i=1}^N (X_i - E(X))(Y_i - E(Y))}{\sqrt{\sum_{i=1}^N (X_i - E(X))^2} \sqrt{\sum_{i=1}^N (Y_i - E(Y))^2}} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \rho, \text{ with } 0 \leq \theta \leq \pi$$

$$(2.a) \quad R = \begin{bmatrix} 1 & r_{1,2} & r_{1,3} & \cdots & r_{1,p} \\ r_{2,1} & 1 & r_{2,3} & \cdots & r_{2,p} \\ r_{3,1} & r_{3,2} & 1 & \cdots & r_{3,p} \\ r_{4,1} & r_{4,2} & r_{4,3} & \cdots & r_{4,p} \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ r_{p,1} & r_{p,2} & r_{p,3} & \cdots & 1 \end{bmatrix},$$

(2.a). For R, a p x p correlation matrix,

² I take 'real-world' financial returns data to be multivariate with marginal distributions that vary notably from each other 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.

³ While r typically is used to represent Pearson's calculated on a sample, ρ often is used to represent Pearson's calculated on a population.

$R = BB'$ where B is the Cholesky factor (defined in Post 1) 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 $\theta_{i,j}$, we have⁴:

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

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

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

This relationship is one-to-one and bi-directional. I present below straightforward SAS/IML code translating correlations to angles (2.a) and angles to correlations (2.b) in Table A:

⁴ Note that a similar recursive relationship exists between partial correlations (Madar, 2015), although its sample-generating algorithm it is not generalizable beyond Pearson's correlations, ie to all positive definite measures of dependence, as shown in my upcoming Post 4.

TABLE A:

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

The above all is well-established and straightforward. But why are we interested in these angles in this setting? There are several very important reasons:

A. Because they are derived based on the matrix's Cholesky factor, the angles, unlike the correlations themselves, are forced on to the unit hyper-(hemi)sphere, where **positive definiteness automatically is enforced**. This is necessary for efficient sampling, as well as for direct and proper definition of the multivariate sample space.

B. Crucially, the **distributions of all of the angles are independent**, which makes sampling, and more importantly, construction of their multivariate distribution (and that of the translated correlation matrix), straightforward and useable, where it otherwise would remain intractable.

C. **The angles contain all information regarding dependence structure** (see Fernandez-Duren & Gregorio-Dominguez, 2023, and Zhang & Songshan, 2023, as well as Opdyke, 2024). On the UNIT hyper-

(hemi)sphere, the only thing we lose is scale, but scale does not and should not matter for any useful and useable measure of dependence.⁵

D. Finally, **angles distributions are more robust and** much better behaved than spectral distributions, and unlike the latter, are **at the right level of aggregation for granular scenarios** (for examples of the dramatic changes of spectral distributions under heavy-tails, see Opdyke, 2024, Burda et al., 2004, Burda et al., 2006, Akemann et al., 2009; Abul-Magd et al., 2009, Bouchaud & Potters, 2015, Martin & Mahoney, 2018), and under serial correlation (see Opdyke, 2024, and Burda et al., 2004, 2011). As discussed below, I present some empirical examples of this in graphs below under real-world financial data conditions.

Fortunately, all of the above advantages of relying on angle values hold not only for the Gaussian identity matrix, but also for any values of Pearson's matrix under ANY data conditions found in challenging, real-world financial settings.

Beyond the Gaussian Identity Matrix: Pearson's Distribution for Any Values, Any Data

Recall from Post 2 that the only requirement for the bi-directional, one-to-one relationship between correlations and angles is that the correlation matrix be symmetric positive definite,⁶ which, numerical issues aside,⁷ is always the case for Pearson's product moment correlation matrix regardless of its values.⁸ Consequently, we are not restricted to the Gaussian identity matrix if we want to use angles and their distributions to define the finite sample distribution of Pearson's matrix. Under the Gaussian identity matrix in Post 2, I first derived the straightforward, analytic distributions of the angles, presenting together their pdf's, cdf's, and quantile function (although the cdf and quantile function previously were claimed to be analytically intractable – see Makalic and Schmidt, 2018). Second, because of A. through D. above, I showed how it is straightforward to use these angles distributions to sample the correlation matrix and, more importantly, to define its finite sample distribution. This is used to obtain both cell-level and matrix-level p-values and confidence intervals that maintain analytic consistency across the two

⁵ Scale invariance is widely proved and cited for Pearson's rho, Kendall's tau, and Spearman's rho (see Xu et al., 2013, and Schreyer et al., 2017 examples).

⁶ See Pinheiro and Bates (1996), Rebonato and Jackel (2000), Rapisarda et al. (2007), Pouramadi and Wang (2015), and Cordoba et al. (2018). I discuss in Post 4 that this requirement of symmetric positive definiteness is true for any dependence measure, not just Pearson's.

⁷ Below I discuss how angles distributions are more stable and robust than spectral distributions by several criteria, including numerically, especially as dependence matrices approach singularity, which arguably is the rule rather than the exception for non-small, real-world investment portfolios.

⁸ Note that for Pearson's specifically, the first and second moments (mean and variance) of the distributions of the returns must exist.

levels.⁹ The only difference between the Gaussian identity matrix and the more general case covered in this Post 3 – any correlation values under any real world financial data – is the angles distributions themselves: all other relationships (ie angles to correlations and correlations to angles) and conditions (A. – D.) hold. So all we need are the angles distributions under general conditions to obtain the distribution of Pearson’s matrix under general correlation value and data conditions.

Angles Distributions for ALL Pearson’s Values, for ANY Real-World Financial Data Conditions

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; and Opdyke, 2024), as well as by different degrees of serial correlation (see Burda et al., 2004, 2011, and Opdyke, 2024), 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 independence (Marchenko and Pastur, 1967).¹⁰ If angles distributions are of similar complexity, deriving their general analytic form under general conditions, if possible, currently appears to be a non-trivial, unsolved problem.

However, this need not 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 in this setting, by multiple criteria, including structurally, empirically, and distributionally.

Structurally: Aggregation level becomes relevant and important here. For a given correlation matrix R there are many more angles distributions than there are spectral distributions (i.e. $p(p-1)/2$ vs p , a factor of $(p-1)/2$ more). As a matrix approaches singularity (non-positive definiteness (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 than is true for eigenvalue distributions. Consequently, the overall construction of the correlation matrix via $R = BB^T$ generally will remain much more stable than one based on an eigen-decomposition of $R = V\Lambda V^{-1}$ where V is a matrix with column eigenvectors and Λ is a diagonal matrix of the corresponding eigenvalues.

⁹ I describe below how this cell-level and matrix-level consistency is critical for attribution analyses specifically, not to mention correct inferences generally.

¹⁰ Note that some exceptions to convergence to this celebrated distribution do exist (see Li and Yao (2018), Hisakado and Kaneko (2023), and Maltsev and Malysheva (2024) for examples).

Empirically: If an angle distribution approaches degeneracy, most of its values typically will approach 0 or π . But the relevant trigonometric functions (sin, 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 that is approaching 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 they have no upper bound (in the general case), makes their computation comparatively less numerically stable as matrices approach NPD.

Distributionally: 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 spikey¹¹ and highly multimodal, and their unboundedness (in the general case) 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 matrix R under real-world financial data.

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. As discussed in more detail below, 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 the angles distributions (and correlations!) 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 correlation 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 (on both a general basis and relative to the alternative of spectral distributions), not to mention the fact that they remain at the right level of aggregation for granular analysis of the correlation 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 correlation matrix, and the known or well-estimated data generating mechanism. 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 densities of all the individual angles, but also the multivariate density of the matrix, which is

¹¹ 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.

just the product of all the individual densities due to their independence per B. above. Note that this goes in both directions: we can perform ‘look-ups’ 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). This process is described step by step below.

1. Simulate samples (say, $N=10k$) based on the specified/known or well estimated correlation matrix and the specified/known or well estimated data generating mechanism.
2. Calculate the corresponding N correlation matrices, and their Cholesky factorizations, and transform each of these into a lower triangle matrix of angles (as described above in (2.a)).
3. Fit kernel densities to each of the $p(p-1)/2$ empirical angle distributions, each having N observations.
4. Generate samples based on the densities in 3.¹²
5. Convert the samples from 4. back to a re-parameterized Cholesky factorization, and then multiply by its transpose to obtain a set of N validly sampled correlation matrices (as described above in (2.b)). Positive definiteness is enforced automatically as the Cholesky factor places us on the **unit** hyper-hemisphere.

The distribution of correlation matrices from 5. is identical to that of 2., but after the kernel densities are fit once in 3., generating samples in 4. is orders of magnitude faster than relying on direct simulations in steps 1. and 2. And of course, using 3.-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!) and subsequently, the proper transformation of the angles to correlations. This reliance on the angles, and their subsequent transformation to correlations, allows us to isolate specifically the distribution of the entire correlation matrix, for probabilistic inference, without touching any other distributional aspect of the data, which is the point of the methodology.

So this framework is essentially identical to that for the specific case of the Gaussian identity matrix derived in Post 2, the only difference being it is based on nonparametrically defined, as opposed to analytically defined, angles distributions. Before covering implementation details below, I show some examples of graphs of the angles distributions and the corresponding spectral distribution under real-world financial returns data. 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 (asym parm= $1, 0.6, 0, -0.6, -1$), non-stationarity (std dev= $3\sigma, \sigma/3, \sigma; n/3$ obs each), and serial correlation (AR1= $-0.25, 0, 0.25, 0.50, 0.75$), with a block correlation structure shown in (3) below and $n=126$ observations.¹³ The spectral distribution is compared against Marchenko-Pastur as a baseline.

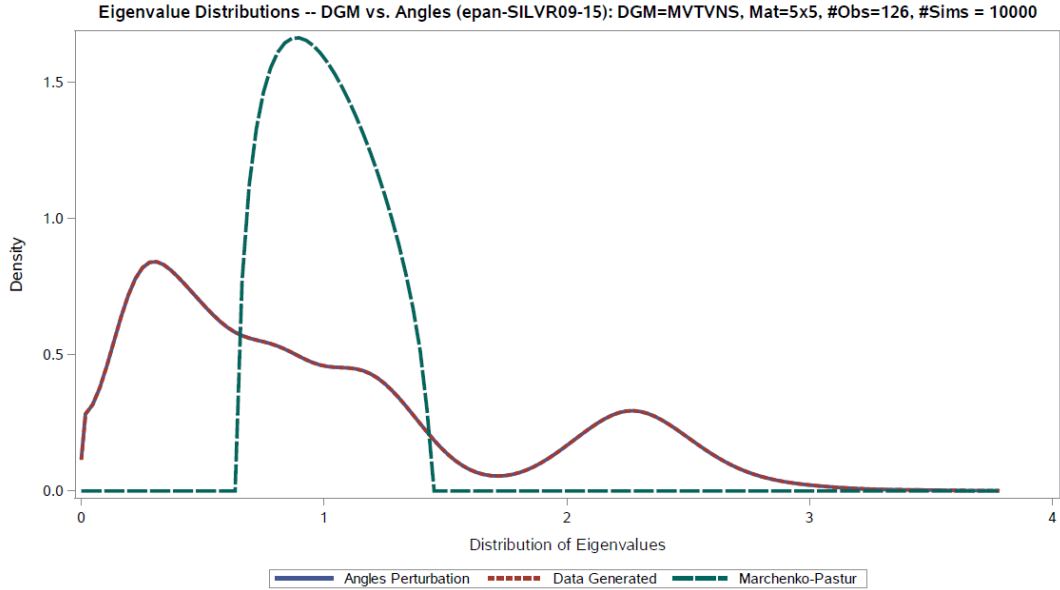
¹² Algorithms for sample generation from broadly used kernels (e.g. the Gaussian and Epanechnikov) are widely known. An example of the latter is simply the median of three uniform random variates (see Qin and Wei-Min, 2024).

¹³ Note 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 on the data (and then empirically rescale the marginal densities).

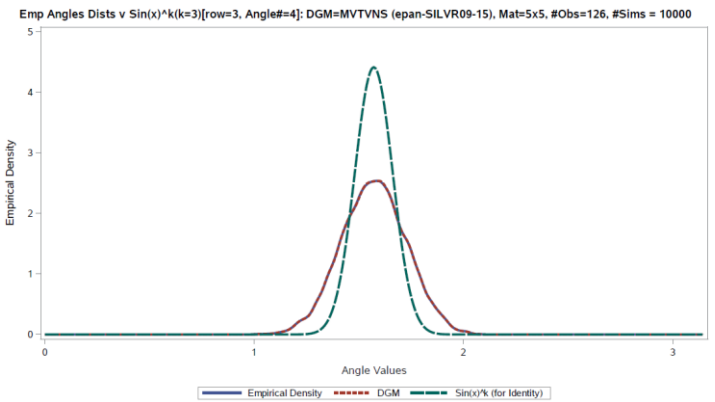
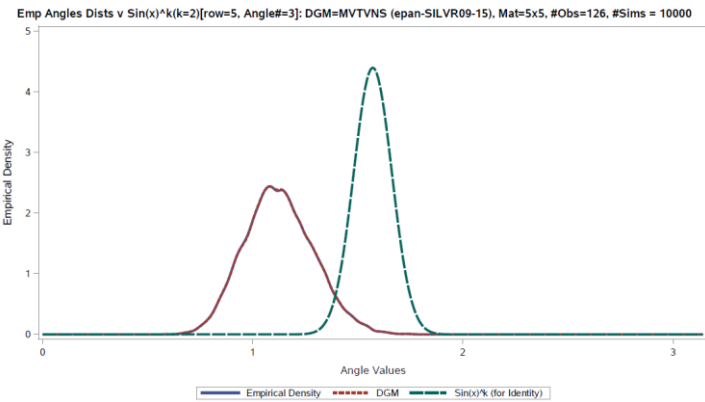
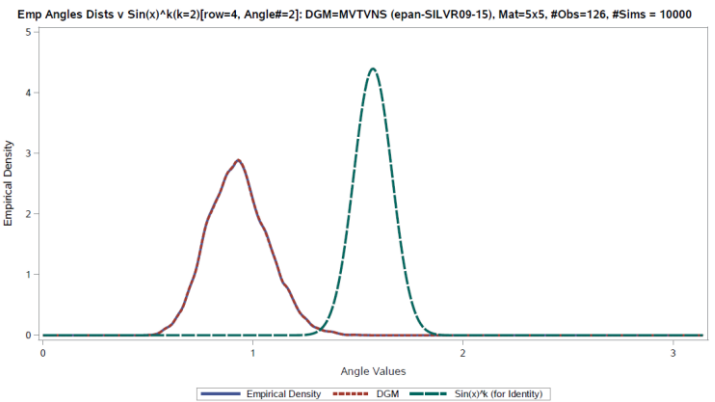
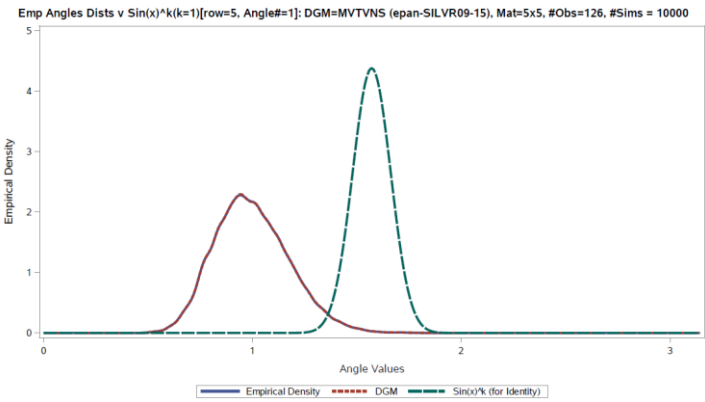
1	-0.3	-0.3	0.2	0.2
-0.3	1	-0.3	0.2	0.2
-0.3	-0.3	1	0.2	0.2
0.2	0.2	0.2	1	0.7
0.2	0.2	0.2	0.7	1

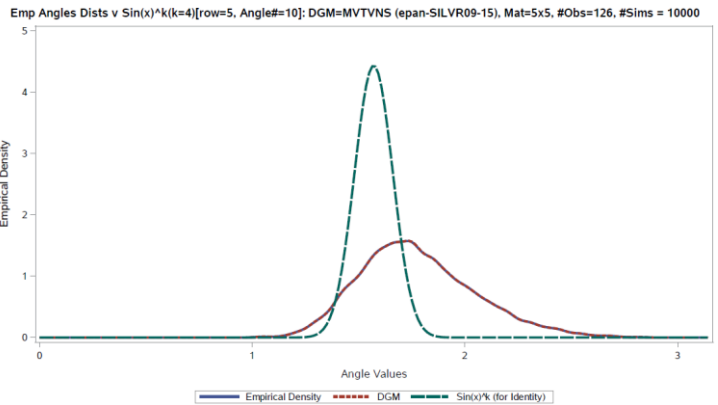
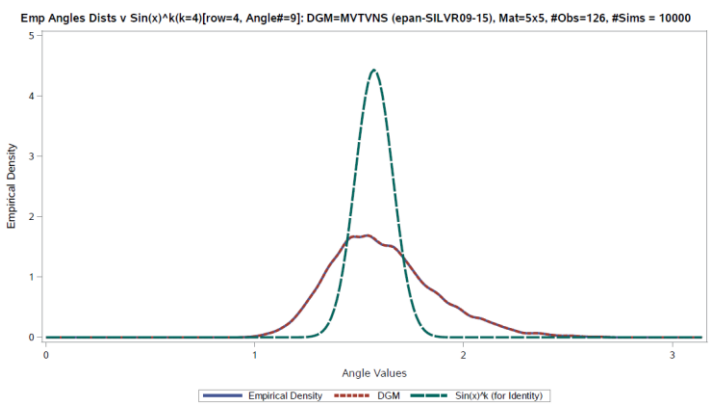
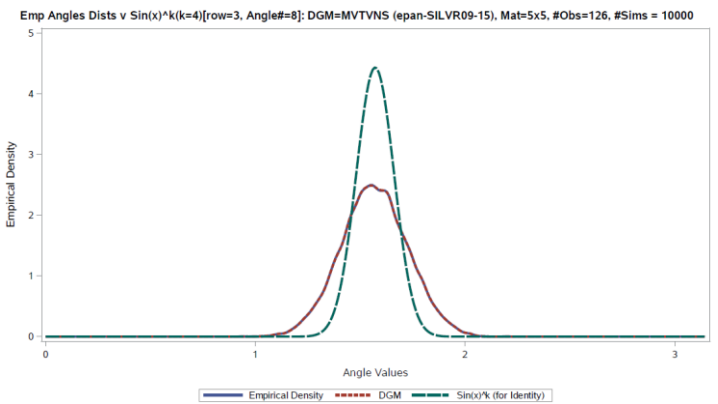
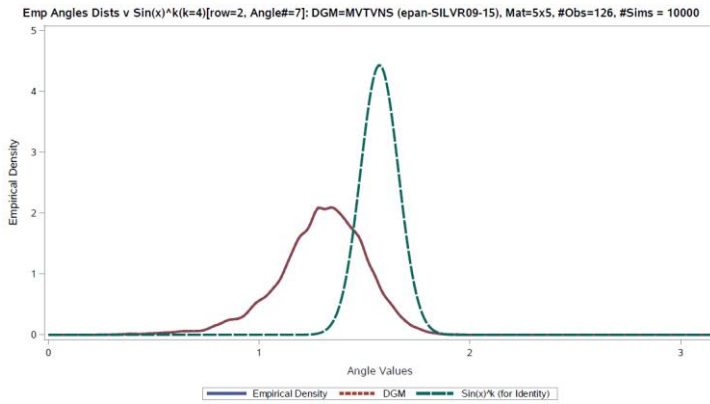
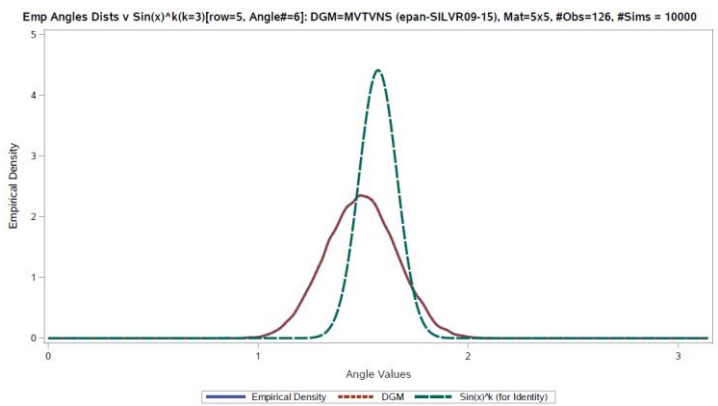
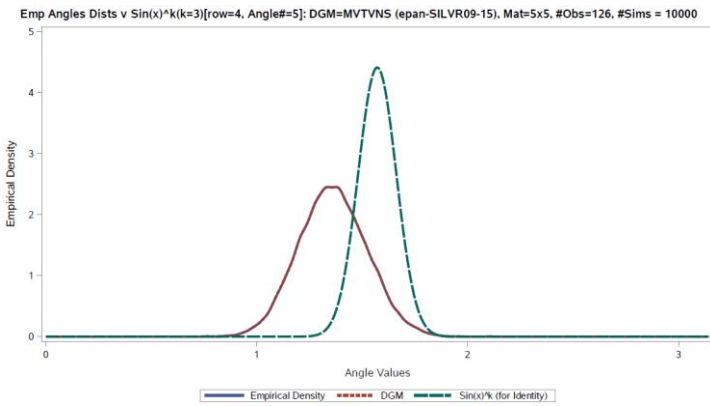
(3)

Graph 1: Spectral Distribution – Angles/Kernel Perturbation v Data Simulations v Marchenko Pastur



Graphs 2-10: Angles Distributions – Angles/Kernel Perturbation v Data Simulations v Independence





Several points are worth noting and reemphasizing from these graphs. First, the graphs of the angles distributions contain three densities: A. one based on angles perturbation (i.e. sampling from the fitted kernel) as described above in steps 3.-5., 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). The spectral distributions based on the samples generated in both A. and B. also are identical, as are a wide range of additional analyses not presented herein (e.g. various norms, VaR-based economic capital, and ‘generalized entropy’ as described below). This empirically validates that the angles-perturbation approach is an efficient and correct method for isolating and generating the density of the correlation 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 here (see Fernandez-Duren & Gregorio-Dominguez, 2023, and Zhang & Songshan, 2023, as well as Opdyke, 2024).

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 and unimodal and well bounded, N=10,000 simulations almost always suffice to provide a precise and accurate measure of their densities. Poorly behaved distributions that are very spikey, highly multi-modal, and unbounded 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 $p=100$ and $p(p-1)/2=4,950$ pairwise associations/cells.

Finally, as described above, note the multi-modal and unbounded nature of the spectral distribution for this portfolio compared to the angles distributions, where the biggest thing approaching an estimation challenge is a slight 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 a single spectral density, even if perfectly estimated, will be able to capture and reflect all the richness of dependence structure reflected here at the granular level of the pairwise cells, for any useful purposes, including cell-level attribution analyses, granular scenario and reverse scenario analyses, cell-level intervention ‘what if’ analyses, and customized stress testing, let alone precise and correct inference at either the cell level OR the matrix level. I now stop beating this drum¹⁴ and leave comparisons to spectral distributions behind to cover implementation issues below.

Nonparametric Kernel Implementation

Due to the bounded nature of the angles distributions on $\theta \in (0, \pi)$, the kernel must be appropriately reflected at the boundary (see Silverman, 1986) via: if $\theta < 0$ then $\theta \leftarrow -\theta$; if $\theta > \pi$ then $\theta \leftarrow (2\pi - \theta)$. As per the standard implementation, the kernel itself is defined as

$$f_h(\theta) = \frac{1}{N} \sum_{i=1}^N K_h(\theta - \theta_i) = \frac{1}{hN} \sum_{i=1}^N K_h([\theta - \theta_i]/h) \quad \text{with}$$

¹⁴ 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.”

Gaussian: $K(\theta) = (1/\sqrt{2\pi}) \cdot e^{-\theta^2/2}$, Epanechnikov: $K(\theta) = (3/4) \cdot (1 - \theta^2)$, $|\theta| \leq 1$.

Both the Gaussian and the Epanechnikov kernels have been tested extensively in this setting, 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

$h = 2.34 \cdot \hat{\sigma} \cdot N^{-1/5}$ 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 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. The bandwidth

that appears to perform best across wide-ranging conditions is $h = 0.9 \cdot \min(\text{IQR}/1.34, \hat{\sigma}) \cdot N^{-1/5}$.

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 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 much larger than $p=100$, a further tightening of this factor may be required, and this is readily determined by empirical testing of the aggregated metrics of interest.

This application of a nonparametric kernel for density estimation is straightforward and very well established in the literature, as are algorithms to generate samples from them (see Qin and Wei-Min, 2024). And as shown in the graphs above, when fitted to angles distributions and used as a basis for subsequent angle perturbation, it generates angles distributions, and corresponding correlation matrices and spectral distributions, that all are empirically identical to those based on direct data simulations, thus confirming the utility and appropriateness of this approach in this setting.

Analytically Consistent Cell-Level and Matrix-Level p-values and Confidence Intervals

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 cells and the entire correlation matrix are the same as those derived in Post 2 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 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 (3):

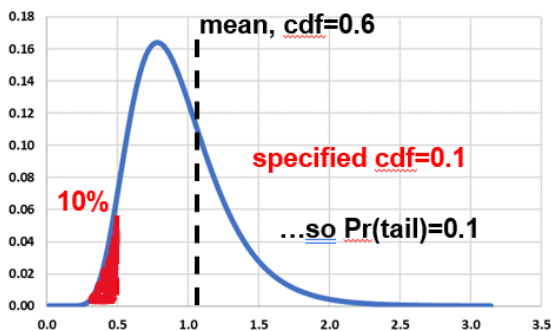
(3) one-sided p-value = $F_X(x;k)$ or $1 - F_X(x;k)$ for lower and upper tails, respectively,
 where $k = n - \text{column\#} - 2$

However, due to possible (probable) asymmetry, the two-sided p-value is different, requiring first the calculation of the empirical mean correlation matrix from the simulations in step 2. above. This mean correlation matrix is then translated into a matrix of angles, and we obtain the empirical cdf's 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. The two-sided p-values are based on the distance between the cdf's of each of the angles of the specified correlation matrix being 'tested' and those of the "mean angles," where 'distance' is the integrated density (probability), not distance of the angle size on the x-axis. Specifically, the two-sided p-values are the sum of the probability in the tails BEYOND this distance.¹⁵ Formulaically this is simply (4):

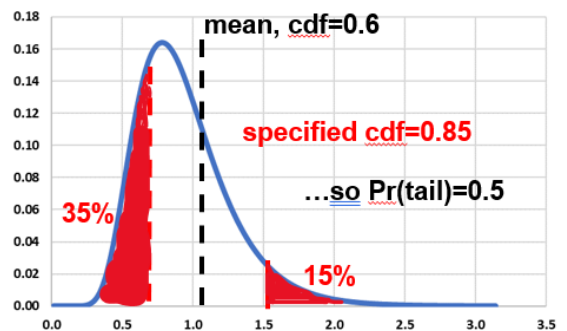
(4) two-sided p-value = $\max[0, \text{Mcdf} - d] + \max[0, 1 - (\text{Mcdf} + d)]$, where
 $d = \text{abs}(\text{Mcdf} - \text{cdf})$, Mcdf = mean angle cdf, cdf = cdf of specified angle

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 (ie the correlation matrix for which we are obtaining p-values, confidence intervals, etc.) 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 single specified (more) extreme angle cdf



Graph 12: p-value for a single specified non-extreme angle 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 probabilistically from the mean cdf=0.6, with a probabilistic distance of 0.5 for Graph 11, and 0.25 for Graph 12. Moreover, a value as extreme as the case of Graph 11 is associated with only 1/5 the probability of being observed compared to that of Graph 12 (compare the red shaded areas). This example demonstrates why asymmetry must be properly taken

¹⁵ So this 'distance' is 0.5 for Graph 1 and 0.25 for Graph 2.

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 the Scenarios section.

Cell-level confidence intervals simply are calculated as in (5), which automatically takes asymmetry into account. Asymmetry notwithstanding, this is identical to the same calculation under the (Gaussian) identity matrix.

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

$$k = n - \text{column\#} - 2$$

The above describes p-values and confidence intervals at the cell level, i.e. for every cell in the correlation matrix, individually. The p-value and confidence interval(s) at the matrix level are based directly on these cell-level calculations and remain otherwise identical to those calculated in Post 2 under the (Gaussian) identity matrix. The matrix-level p-value, again, is simply one minus the probability of no false positives, which is the definition of controlling the family-wise error rate (FWER) of the matrix (6).¹⁶

(6) matrix (2-sided) *pvalue* = $\left[1 - \prod_{i=1}^{p(p-1)/2} (1 - p\text{-value}_i) \right]$ where *p-value_i* is the 2-sided p-value.

Because the cell-level distributions are independent, their p-values are independent, and otherwise statistically more powerful approaches for calculating the FWER that rely on, for example, resampling methods (Westfall and Young, 1993, and Romano and Wolf, 2016), do not apply here. In other words, they provide no power gain over (6) because under independence, there is no dependence structure for them to exploit. So the straightforward calculation above in (6) is, by definition, the most powerful for FWER control.

Similarly, just as under the (Gaussian) identity matrix, calculation of the confidence interval for the entire matrix (7) is essentially the same as that of the p-value, but of course it is divided in half to account for each tail, and the root of the critical values is taken, rather than the product. Otherwise, the calculations are identical to obtain the critical alphas for these ‘simultaneous confidence intervals.’

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

These critical alphas, when inserted as values in the empirically-based cdf ‘look-up’ functions, 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 in this case is the specified correlation matrix being ‘tested,’

¹⁶ Note that this approach has been used in the literature for addressing very closely related problems in this setting (see Fang et al., 2024).

and is no longer strictly only the identity matrix. Independence of the angles distributions again makes these simultaneous confidence intervals very straightforward to calculate.

Importantly, again note that we can go in either direction regarding these results: we can specify a correlation 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 matrix. 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 below, but first I discuss the scenario-restricted case.

Flexible Scenarios, Reverse Scenarios, and Customized Stress Tests

NAbC is a ‘bottom up’ approach to defining the finite-sample distribution of the correlation matrix, based on the distributions of the individual correlation cells. In addition to analytic consistency, this provides a flexibility in scenario definition and scenario analytics that other approaches cannot match. 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, 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. In other words, while correlation ‘breakdowns’ will occur under all of these extreme conditions, the granular nature of pairwise association matrices ensures that the fundamentally different nature of these breakdowns will be captured and reflected empirically in all related analyses. The only way to flexibly and realistically model this is at the most granular level – that of the individual correlation cells.

Fortunately, as described in detail in Post 2, NAbC allows for this, with full inferential powers under any definable scenario within the framework of pairwise associations defined by a correlation (dependence) matrix. Without repeating this description in detail, this is made possible by exploiting four simultaneous conditions – 1. independence of the angles distributions; 2. (correlation) distribution invariance to row and column order; 3. the mechanics of matrix multiplication; and 4. the granular, cell-level geometry of NAbC, which allows arbitrarily chosen cells to vary and the rest to remain constant/unaffected by the scenario, without violating positive definiteness. 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 analyses cannot do effectively in this setting).

And while NAbC covers inference for a matrix of all pairwise associations, the same level of flexibility and sophistication exists in their estimation and simulation via, for example, vine copulas (see Czado and

Nagler, 2022).¹⁷ So rather than imposing unrealistic restrictions, the framework of all pairwise associations is greatly liberating in its analytics, whether the focus is on inference, estimation, and/or (synthetic) scenario simulation. This all will be explored further in Post 4, which expands NAbC's range of application to dependence measures beyond Pearson's.

In this Post 3, however, one difference in scenario definition and implementation, compared to that of Post 2, is relevant for these scenario analytics: when allowing only selected cells of the correlation matrix to vary for a given scenario, while holding the remaining cells constant, we must insert angle values into the 'frozen' cells that not only hold the correlations constant, but also enforce positive definiteness. Where do we get those values? From the "mean angles" matrix defined above. As the mean of N=10,000 simulations, this is a stable and robust estimator of the correlation matrix under the (simulated) null hypothesis.¹⁸ It will both hold constant the correlation values at their means, and it is itself positive definite, as it is based on a linear combination of positive definite matrices. We simply perform steps 1.-5. as usual, but after step 4. we overwrite values of the simulated angles in those cells to be held constant with the mean angle values, so that the 'frozen cells' in every simulation from 4. contain their respective (constant) mean values. The resulting scenario-restricted correlation matrices always will be positive definite, because the systematic row and column sorting for the scenario (described in detail in Post 2) in effect 'disables' systematic changes to these cells if we do not change the angle values, and their values will be 'frozen' at their mean correlation values. We did not need to do this under the (Gaussian) identity matrix because the mean values all were zero, by definition, and because simulation is not required in this case.

As mentioned above and described in detail in Post 2, this provides an unmatched degree of flexibility for scenario definition and scenario analytics: this approach works for ANY scenario restricted matrix within the framework of all pairwise associations, because the distributions of both the entire matrix and the individual cells are invariant to row and column order. Post 2 describes in detail how simply re-sorting the matrix allows for the selection of only specific cells that are meant to vary for a particular scenario, while maintaining proper inferential properties for both the individual cells and the entire matrix. Examples applying NAbC for inference on both unrestricted and scenario-restricted matrices are presented below.

¹⁷ Even though vine copulas are a sophisticated and highly flexible method to define, estimate, and simulate bivariate dependence structures, NAbC (from Post 2) replaces earlier attempts to define and sample the gaussian identity correlation matrix using these methods, as they are unnecessarily computationally demanding and complex for this purpose (see Lewandowski et al., 2009 and Kurowicka, 2014) when a fully analytic solution exists, as shown in Post 2.

¹⁸ The degree of empirical accuracy attained when using these means is based directly on the number of simulations. This can be seen in the scenario-restricted results (specifically, in the 'frozen' cells of the scenario-restricted matrices) in the next section.

NAbC Applied: Unrestricted and Scenario-Restricted p-values and Confidence Intervals

Below I apply NAbC to obtain both p-values and confidence intervals under two cases: unrestricted, and scenario-restricted. Solely for ease of replication, the data generating mechanism for these examples is simply multivariate standard normal, with $N=25k$ simulations and number of observations $n = 160$.

UNRESTRICTED CASE: Given a specified or well-estimated correlation matrix [A], and its specified or well-estimated data generating mechanism:

[A]

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

[B]

0.8				
0.7	0.8			
0.8	0.7	0.7		
0.7	0.8	0.8	0.7	

[C]

1				
0.40	1			
0.20	0.10	1		
0.03	-0.07	-0.20	1	
0.33	0.60	0.25	-0.23	1

- Q1. **Confidence Intervals:** What are the two correlation 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 and every cell of [A]?
- Q2. **Quantile Function:** What is the unique correlation matrix associated with [B], a matrix of cumulative distribution function values associated with the corresponding cells of [A]?
- Q3. **p-values:** Under the null hypothesis that observed correlation 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 and every cell of [C]?

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

[D]

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

[E]

	0.8			
	0.8	0.8	0.7	

[F]

1				
0.2	1			
-0.1	0.015	1		
0.3	-0.3	-0.1	1	
0.6	0.5	-0.2	0.302	1

- Q4. **Confidence Intervals:** What are the two correlation 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 only those cells of [D] that are relevant to the scenario (green)?
- Q5. **Quantile Function:** What is the unique correlation matrix associated with [E], a matrix of cumulative distribution function values associated with the corresponding cells of [D]?
- Q6. **p-values:** Under the null hypothesis that observed correlation 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 every (non-constant, green) cell of [F]?

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, under any data conditions. Only NAbC can simultaneously answer Q1.-Q6. above under general data conditions, as shown below.

Q1	Q2	Q3	Q4	Q5	Q6																																																																																																				
	<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.273</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.052</td><td>0.365</td><td>1</td><td></td><td></td></tr> <tr><td>0.369</td><td>-0.209</td><td>-0.060</td><td>1</td><td></td></tr> <tr><td>0.631</td><td>0.488</td><td>0.116</td><td>0.183</td><td>1</td></tr> </table>	1					0.273	1				-0.052	0.365	1			0.369	-0.209	-0.060	1		0.631	0.488	0.116	0.183	1	<p>p-value=0.1473</p> <table border="1"> <tr><td>0.0033</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.0075</td><td>0.0290</td><td></td><td></td><td></td></tr> <tr><td>0.0227</td><td>0.0297</td><td>0.0079</td><td></td><td></td></tr> <tr><td>0.0401</td><td>0.0021</td><td>0.0101</td><td>0.0049</td><td></td></tr> </table>	0.0033					0.0075	0.0290				0.0227	0.0297	0.0079			0.0401	0.0021	0.0101	0.0049			<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.1996</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.0995</td><td>0.3679</td><td>1</td><td></td><td></td></tr> <tr><td>0.2996</td><td>-0.2991</td><td>-0.0998</td><td>1</td><td></td></tr> <tr><td>0.5988</td><td>0.4304</td><td>0.0521</td><td>0.1312</td><td>1</td></tr> </table>	1					0.1996	1				-0.0995	0.3679	1			0.2996	-0.2991	-0.0998	1		0.5988	0.4304	0.0521	0.1312	1	<p>p-value=0.0526</p> <table border="1"> <tr><td></td><td></td><td></td><td></td><td></td></tr> <tr><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>0.04032</td><td></td><td></td><td></td><td></td></tr> <tr><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>0.00008</td><td>0.00184</td><td>0.01088</td><td></td><td></td></tr> </table>											0.04032										0.00008	0.00184	0.01088							
1																																																																																																									
0.273	1																																																																																																								
-0.052	0.365	1																																																																																																							
0.369	-0.209	-0.060	1																																																																																																						
0.631	0.488	0.116	0.183	1																																																																																																					
0.0033																																																																																																									
0.0075	0.0290																																																																																																								
0.0227	0.0297	0.0079																																																																																																							
0.0401	0.0021	0.0101	0.0049																																																																																																						
1																																																																																																									
0.1996	1																																																																																																								
-0.0995	0.3679	1																																																																																																							
0.2996	-0.2991	-0.0998	1																																																																																																						
0.5988	0.4304	0.0521	0.1312	1																																																																																																					
0.04032																																																																																																									
0.00008	0.00184	0.01088																																																																																																							
<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>-0.017</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.316</td><td>0.117</td><td>1</td><td></td><td></td></tr> <tr><td>0.089</td><td>-0.558</td><td>-0.214</td><td>1</td><td></td></tr> <tr><td>0.439</td><td>0.126</td><td>-0.345</td><td>-0.136</td><td>1</td></tr> </table>	1					-0.017	1				-0.316	0.117	1			0.089	-0.558	-0.214	1		0.439	0.126	-0.345	-0.136	1	<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.406</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>0.130</td><td>0.517</td><td>1</td><td></td><td></td></tr> <tr><td>0.486</td><td>0.056</td><td>0.190</td><td>1</td><td></td></tr> <tr><td>0.727</td><td>0.631</td><td>0.368</td><td>0.443</td><td>1</td></tr> </table>	1					0.406	1				0.130	0.517	1			0.486	0.056	0.190	1		0.727	0.631	0.368	0.443	1		<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.1996</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.0995</td><td>0.0827</td><td>1</td><td></td><td></td></tr> <tr><td>0.2996</td><td>-0.2991</td><td>-0.0998</td><td>1</td><td></td></tr> <tr><td>0.5988</td><td>0.3110</td><td>-0.1545</td><td>-0.0654</td><td>1</td></tr> </table>	1					0.1996	1				-0.0995	0.0827	1			0.2996	-0.2991	-0.0998	1		0.5988	0.3110	-0.1545	-0.0654	1	<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.1996</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.0995</td><td>0.5166</td><td>1</td><td></td><td></td></tr> <tr><td>0.2996</td><td>-0.2991</td><td>-0.0998</td><td>1</td><td></td></tr> <tr><td>0.5988</td><td>0.4633</td><td>0.1605</td><td>0.2680</td><td>1</td></tr> </table>	1					0.1996	1				-0.0995	0.5166	1			0.2996	-0.2991	-0.0998	1		0.5988	0.4633	0.1605	0.2680	1	
1																																																																																																									
-0.017	1																																																																																																								
-0.316	0.117	1																																																																																																							
0.089	-0.558	-0.214	1																																																																																																						
0.439	0.126	-0.345	-0.136	1																																																																																																					
1																																																																																																									
0.406	1																																																																																																								
0.130	0.517	1																																																																																																							
0.486	0.056	0.190	1																																																																																																						
0.727	0.631	0.368	0.443	1																																																																																																					
1																																																																																																									
0.1996	1																																																																																																								
-0.0995	0.0827	1																																																																																																							
0.2996	-0.2991	-0.0998	1																																																																																																						
0.5988	0.3110	-0.1545	-0.0654	1																																																																																																					
1																																																																																																									
0.1996	1																																																																																																								
-0.0995	0.5166	1																																																																																																							
0.2996	-0.2991	-0.0998	1																																																																																																						
0.5988	0.4633	0.1605	0.2680	1																																																																																																					
<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.049</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.250</td><td>0.165</td><td>1</td><td></td><td></td></tr> <tr><td>0.154</td><td>-0.497</td><td>-0.203</td><td>1</td><td></td></tr> <tr><td>0.491</td><td>0.212</td><td>-0.253</td><td>-0.089</td><td>1</td></tr> </table>	1					0.049	1				-0.250	0.165	1			0.154	-0.497	-0.203	1		0.491	0.212	-0.253	-0.089	1	<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.347</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>0.060</td><td>0.452</td><td>1</td><td></td><td></td></tr> <tr><td>0.435</td><td>-0.056</td><td>0.091</td><td>1</td><td></td></tr> <tr><td>0.693</td><td>0.569</td><td>0.265</td><td>0.341</td><td>1</td></tr> </table>	1					0.347	1				0.060	0.452	1			0.435	-0.056	0.091	1		0.693	0.569	0.265	0.341	1		<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.1996</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.0995</td><td>0.1432</td><td>1</td><td></td><td></td></tr> <tr><td>0.2996</td><td>-0.2991</td><td>-0.0998</td><td>1</td><td></td></tr> <tr><td>0.5988</td><td>0.3404</td><td>-0.1122</td><td>-0.0169</td><td>1</td></tr> </table>	1					0.1996	1				-0.0995	0.1432	1			0.2996	-0.2991	-0.0998	1		0.5988	0.3404	-0.1122	-0.0169	1	<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.1996</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.0995</td><td>0.4537</td><td>1</td><td></td><td></td></tr> <tr><td>0.2996</td><td>-0.2991</td><td>-0.0998</td><td>1</td><td></td></tr> <tr><td>0.5988</td><td>0.4492</td><td>0.1158</td><td>0.2181</td><td>1</td></tr> </table>	1					0.1996	1				-0.0995	0.4537	1			0.2996	-0.2991	-0.0998	1		0.5988	0.4492	0.1158	0.2181	1	
1																																																																																																									
0.049	1																																																																																																								
-0.250	0.165	1																																																																																																							
0.154	-0.497	-0.203	1																																																																																																						
0.491	0.212	-0.253	-0.089	1																																																																																																					
1																																																																																																									
0.347	1																																																																																																								
0.060	0.452	1																																																																																																							
0.435	-0.056	0.091	1																																																																																																						
0.693	0.569	0.265	0.341	1																																																																																																					
1																																																																																																									
0.1996	1																																																																																																								
-0.0995	0.1432	1																																																																																																							
0.2996	-0.2991	-0.0998	1																																																																																																						
0.5988	0.3404	-0.1122	-0.0169	1																																																																																																					
1																																																																																																									
0.1996	1																																																																																																								
-0.0995	0.4537	1																																																																																																							
0.2996	-0.2991	-0.0998	1																																																																																																						
0.5988	0.4492	0.1158	0.2181	1																																																																																																					

For Q1 and Q4, 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 Q6, the smaller p-values for the individual cells compared to the respective matrix-level p-values, which are larger, as expected, as they control FWER. Note also that the green cells of Q5 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

$$R = BB^T$$

NAbC Remains “Estimator Agnostic”

Another important and useful characteristic of NAbC, under both unrestricted and scenario-restricted cases, is that it remains “estimator agnostic,” that is, valid for use with any reasonable estimator of dependence structure. 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 and shifting combinations of these characteristics. 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 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. In Post 4 I address how NAbC’s estimator-agnostic nature applies beyond Pearson’s correlation, to any positive definite measure of dependence, thus adding one more flexible aspect of NAbC that further expands its already wide range of use and utility.

LNP: a Measure of Generalized Entropy

As I did in Post 2, it is worth taking an arguably minor digression here to examine further the meaning and implications of the cell-level (two-sided) p-values shown above in (4). The (two-sided) p-value provides what can be viewed as a competitor to distance metrics that has some advantages over traditional distance metrics, such as norms. Some commonly used norms in this setting for measuring correlation ‘distances’ are listed below in (8).

(8)
$$\|x\| = \left(\sum_{i=1}^d |x_i|^m \right)^{1/m}$$
 where x is a distance from a presumed or baseline correlation value, d=number of observations, and m=1, 2, and ∞ correspond to the Taxi, Frobenius/Euclidean, and Chebyshev norms, respectively.

All of these norms measure absolute distance from a presumed or baseline correlation value. But the range of all relevant and widely used dependence measures is bounded, either from -1 to 1 or 0 to 1, and 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 0.01 from an original or presumed correlation value of, say, 0.97, means something very different than the same shift from 0.07. NAbC attributes probabilistic MEANING to these two different cases, while a norm would 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 these cell-level p-values from (4) is the natural log of the product of the p-values, dubbed ‘LNP’ in (9) below:

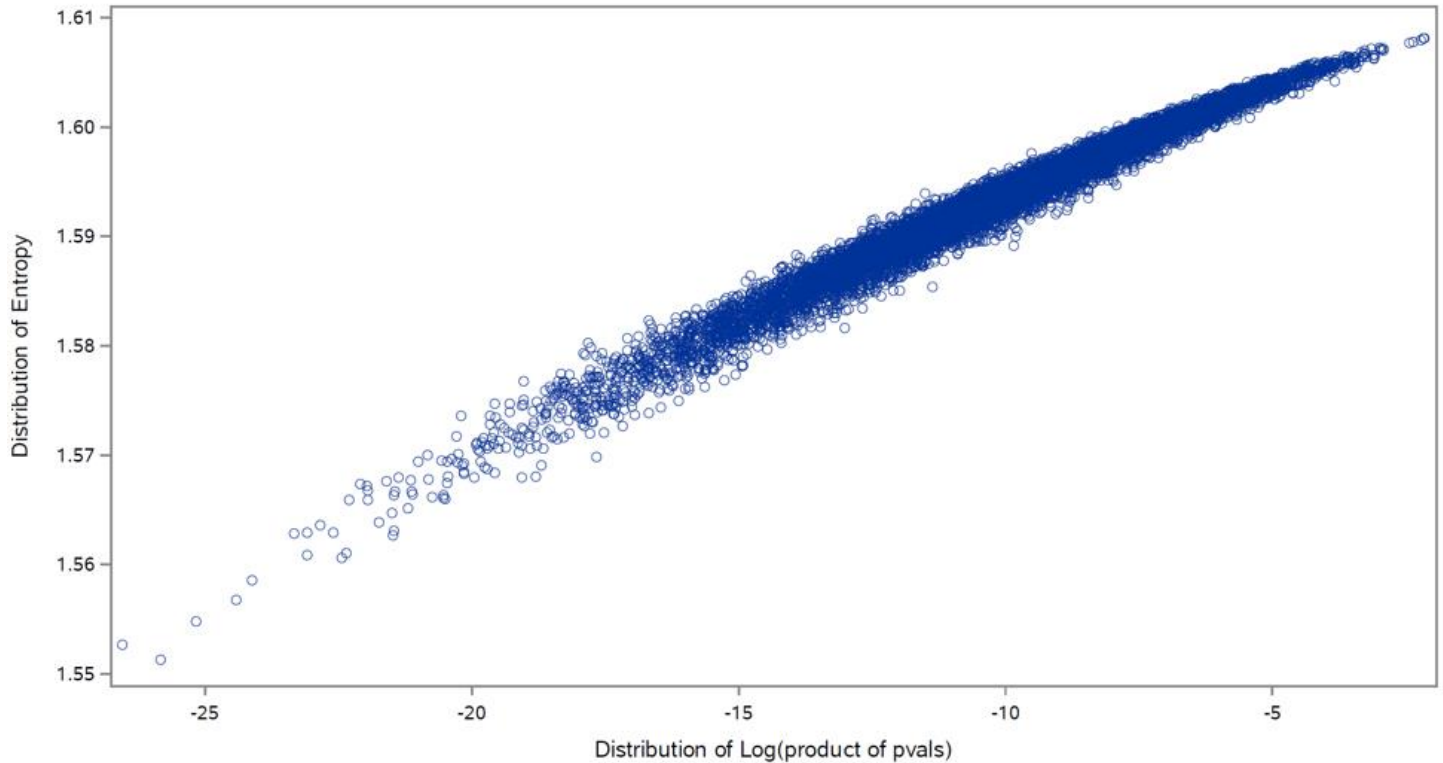
(9) "LNP" = $\ln \left(\prod_{i=1}^q p\text{-value}_i \right) = \sum_{i=1}^q \ln [p\text{-value}_i]$ where $q = p(p-1)/2$ and $p\text{-value}_i$ is 2-sided.

This was shown in Post 2, under the (Gaussian) identity matrix, to have a very strong correspondence with the entropy of the correlation matrix, defined by Felipe et al. (2021 and 2023) as (10) below:

$$(10) \quad \text{Entropy} = \text{Ent}(R/p) = -\sum_{j=1}^p \lambda_j \ln(\lambda_j)$$

where R is the sample correlation matrix and λ_j are the p eigenvalues of the correlation matrix after it is scaled by its dimension, R/p (note that this result (10), like NAbC, is valid for ANY positive definite measure of dependence, not just Pearson's, as will be discussed in Post 4). Graph 13 below compares LNP to the entropy of the correlation matrix in 10,000 simulations under the Gaussian identity matrix. The resulting Pearson's correlation between them is just shy of 0.99.

Graph 13: Identity Matrix Simulations -- LNP v Entropy



It is important to note, however, that entropy here is limited to being calculated relative to the case of independence, which for Pearson's corresponds only with the identity matrix.¹⁹ In contrast, LNP can be calculated and retains its meaning in all cases, based on ANY values of Pearson's matrix, not just the identity matrix. Yet the correspondence of LNP to entropy under the specific case of the identity matrix speaks to LNP's natural interpretation as a meaningful measure of deviation/distance/independence/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, as specified by the researcher (i.e. the specified correlation matrix), that is not necessarily perfect (in)dependence. Such measures certainly are relevant in this setting as entropy has been used increasingly in the literature to

¹⁹ Recall, of course, that a zero value for Pearson's correlation does not imply independence, but independence does imply a zero value for Pearson's correlation.

measure, monitor, and analyze financial markets (see Meucci, 2010, Almog and Shmueli, 2019, Chakraborti et al., 2020, and Vorobets, 2024a, 2024b, for several examples).

Interpretations aside, the use of LNP here warrants further investigation as a matrix-level measure that, unlike widely used measures such as various norms, has a solid and meaningful probabilistic foundation. Its calculation applies not only beyond the identity matrix for Pearson's, and the independence case generally, but also to ALL positive definite measures of dependence, regardless of their values, as discussed further in Post 4. 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.

Conclusion

In Posts 1 and 2 I listed the seven characteristics of the full NAbC solution that, taken together, are shared by no other approach, and for completeness I list them again below:

1. validity under challenging, real-world financial data conditions, with marginal asset distributions characterized by notably different degrees of serial correlation, non-stationarity, heavy-tailedness, and asymmetry
2. application to ANY positive definite dependence measure, including, for example, Pearson's product moment correlation, rank-based measures like Kendall's tau and Spearman's rho, the kernel-based generalization of Szekely's distance correlation, and the tail dependence matrix, among others.
3. it remains "estimator agnostic," that is, valid regardless of the sample-based estimator used to estimate any of the above-mentioned dependence measures
4. it 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 also effectively facilitates attribution analyses)
5. it provides a one-to-one quantile function, translating a matrix of all the cells' cdf values to a (unique) correlation (dependence measure) matrix, and back again, enabling precision in reverse scenarios and stress testing
6. all the above results remain valid even when selected cells in the matrix are 'frozen' for a given scenario or stress test, while the rest are allowed to vary, enabling granular and realistic scenarios
7. it 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, enabling reliable application in actual, imperfect, non-textbook settings

Post 2 covered 4, 5, 6, and 7 above, with NAbC providing a fully analytic solution for the finite sample distribution of Pearson's under the Gaussian identity matrix. This Post 3 expanded NAbC's range of

application to cover 1 and 3 as well, providing the solution for Pearson's under ALL matrix values and ALL real-world financial data conditions, using exactly the same angles-based framework. The only difference was the nonparametric rather than the analytic basis for defining the angles distributions, but all other components of the framework remain the same. Post 2 provided an interactive spreadsheet that implements fully analytic p-values and confidence intervals,

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

while this Post 3 provides the same answers in an example, above, under much more general conditions. Post 4 will continue to expand NAbC's range of application to characteristic 2. above, providing the finite sample distribution for cases beyond Pearson's, including ALL positive definite measures of dependence.

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.
- Akemann, G., Fischmann, J., and Vivo, P., (2009), "Universal Correlations and Power-Law Tails in Financial Covariance Matrices," <https://arxiv.org/abs/0906.5249>.
- Almog, A., and Shmueli, E., (2019), "Structural Entropy: Monitoring Correlation-Based Networks over time With Application to Financial Markets," *Scientific Reports*, 9:10832.
- Bouchaud, J., & Potters, M., (2015), "Financial applications of random matrix theory: a short review," *The Oxford Handbook of Random Matrix Theory*, Eds G. Akemann, J. Baik, P. Di Francesco.
- 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.
- 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.

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.

Czado, C., and Nagler, T., (2022), "Vine Copula Based Modeling," *Annual Review of Statistics and Its Application*, pp.453-477.

Embrechts, P., Hofert, M., and Wang, R., (2016), "Bernoulli and Tail-Dependence Compatibility," *The Annals of Applied Probability*, Vol. 26(3), 1636-1658.

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.

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.

Gao, M., and Li, Q., (2024), "A Family of Chatterjee's Correlation Coefficients and Their Properties," arXiv:2403.17670v1 [stat.ME].

Hansen, B., (2014), *Econometrics*, Ch. 20 – Nonparametric Density Estimation, p.333

Hisakado, M. and Kaneko, T., (2023), "Deformation of Marchenko-Pastur distribution for the correlated time series," arXiv:2305.12632v1.

Holzmann, H., and Klar, B., (2024) "Lancaster Correlation - A New Dependence Measure Linked to Maximum Correlation," arXiv:2303.17872v2 [stat.ME].

Johnstone, I., (2001), "On the distribution of the largest eigenvalue in principal components analysis," *The Annals of Statistics*, 29(2): 295–327, 2001.

Kendall, M. (1938), "A New Measure of Rank Correlation," *Biometrika*, 30 (1–2), 81–89.

Kurowicka, D., (2014). "Joint Density of Correlations in the Correlation Matrix with Chordal Sparsity Patterns," *Journal of Multivariate Analysis*, 129 (C): 160–170.

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, 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.

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].

Marchenko, A., Pastur, L., (1967), "Distribution of eigenvalues for some sets of random matrices," *Matematicheskii Sbornik*, N.S. 72 (114:4): 507–536.

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.

Meucci, A., (2010), "Fully Flexible Views: Theory and Practice," arXiv:1012.2848v1

Opdyke, JD, (2024), Keynote Address: "Beating the Correlation Breakdown, for Pearson's and Beyond: Robust Inference and Flexible Scenarios and Stress Testing for Financial Portfolios," QuantStrats11, NYC, March 12.

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.

Pinheiro, J. and Bates, D. (1996), "Unconstrained parametrizations for variance-covariance matrices," *Statistics and Computing*, Vol. 6, 289–296.

Pourahmadi, M., Wang, X., (2015), "Distribution of random correlation matrices: Hyperspherical parameterization of the Cholesky factor," *Statistics and Probability Letters*, 106, (C), 5-12.

Qin, T., and Wei-Min, H., (2024), "Epanechnikov Variational Autoencoder," arXiv:2405.12783v1 [stat.ML] 21 May 2024.

Rapisarda, F., Brigo, D., & Mercurio, F., (2007), "Parameterizing Correlations: A Geometric Interpretation," *IMA Journal of Management Mathematics*, 18(1), 55-73.

Rebonato, R., and Jackel, P., (2000), "The Most General Methodology for Creating a Valid Correlation Matrix for Risk Management and Option Pricing Purposes," *Journal of Risk*, 2(2)17-27.

Romano, J., and Wolf, M., (2016), "Efficient computation of adjusted p-values for resampling-based stepdown multiple testing," *Statistics & Probability Letters*, Vol 113, 38-40.

Schreyer, M., Paulin, R., and Trutschnig, W., (2017), “On the exact region determined by Kendall's tau and Spearman's rho,” arXiv: 1502:04620.

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.

Shyamalkumar, N., and Tao, S., (2020), “On tail dependence matrices: The realization problem for parametric families,” *Extremes*, Vol. 23, 245–285.

Silverman, B., (1986), Density Estimation for Statistics and Data Analysis, New York, Chapman and Hall.

Spearman, C., (1904), “‘General Intelligence,’ Objectively Determined and Measured,” *The American Journal of Psychology*, 15(2), 201–292.

Szekely, G., Rizzo, M., and Bakirov, N., (2007), “Measuring and Testing Dependence by Correlation of Distances,” *The Annals of Statistics*, 35(6), pp2769-2794.

Vorobets, A., (2024a), “Sequential Entropy Pooling Heuristics,”
<https://ssrn.com/abstract=3936392> or <http://dx.doi.org/10.2139/ssrn.3936392>

Vorobets, A., (2024b), “Portfolio Construction and Risk Management,”
<https://ssrn.com/abstract=4807200> or <http://dx.doi.org/10.2139/ssrn.4807200>

Westfall, P., and Young, S., (1993), Resampling Based Multiple Testing, Wiley Series in Probability and Mathematical Statistics, John Wiley & Sons, Inc., New York, New York.

Xu, W., Hou, Y., Hung, Y., and Zou, Y., (2013), “A Comparative Analysis of Spearman’s Rho and Kendall’s Tau in Normal and Contaminated Normal Models,” *Signal Processing*, 93, 261–276.

Zhang, Y., and Songshan, Y., (2023), “Kernel Angle Dependence Measures for Complex Objects,”
arXiv:2206.01459v2

Beating the Correlation Breakdown, for Pearson's and Beyond: Robust Inference and Flexible Scenarios and Stress Testing for Financial Portfolios

JD Opdyke, Chief Analytics Officer, Partner, Sachs Capital Group Asset Management, LLC
JDOpdyke@gmail.com

Post 4 of 4: ANY Positive Definite Measure, ALL Real-World Financial Data Conditions

NOTE: These posts summarize a chapter in my forthcoming monograph for Cambridge University Press.

Introduction

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 a myriad of 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 (which is when we need valid inferences the most).

This is the fourth in a series of four posts which introduces a new and straightforward method – Nonparametric Angles-based Correlation (“NAbC”) – for defining the finite-sample distributions of a very wide range of dependence measures for financial portfolio analysis. These include ANY that are positive definite, such as 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), as well as measures designed to capture highly non-linear dependence such as the tail dependence matrix (see Embrechts, Hofert, and Wang, 2016, and Shyamalkumar and Tao, 2020), 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).¹

¹ Note that “positive definite” throughout these four posts 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 these dependence measures (eg Szekely's correlation and variants of Chatterjee's) can be applied on a multivariate basis, in arbitrary dimensions, for example, to test the hypothesis of multivariate independence. But “positive definite” herein is not applied in this sense, 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.

This Post 4 expands NAbC’s application beyond Pearson’s to ANY positive definite dependence measure, under any values, and under all challenging, real-world financial data conditions.

POST 1: NAbC introduced.

POST 2: NAbC applied to Pearson’s under the Gaussian identity matrix (fully analytic solution).

POST 3: NAbC applied to Pearson’s under ALL correlation matrix values and ALL relevant, challenging, real-world financial returns data conditions.²

POST 4: NAbC applied to ALL matrix values and ALL positive definite measures of portfolio dependence under ALL relevant, challenging, real-world financial data conditions.

Correlations and Angles (Review of Posts 2 & 3)

To briefly review from Posts 2 & 3, I defined and reviewed the relationship between the correlation cells in a Pearson’s correlation matrix and the angles of their corresponding pairwise data vectors. There exists an angle value for every correlation value in the matrix. For a single, bivariate correlation, this can be seen directly via the widely used cosine similarity in (1),³ but the matrix analog also is well established in the literature as shown in (2.a) and (2.b) (see Pinheiro and Bates, 1996, Rebonato & Jaeckel, 2000, Rapisarda et al., 2007, and Pourahmadi and Wang, 2015, but note a typo in the formula in Pourahmadi and Wang, 2015 corresponding to (2.b) below):

$$(1) \quad \cos(\theta) = \frac{\text{inner product}}{\text{product of norms}} = \frac{\langle \mathbf{X}, \mathbf{Y} \rangle}{\|\mathbf{X}\| \|\mathbf{Y}\|} = \frac{\sum_{i=1}^N (X_i - E(X))(Y_i - E(Y))}{\sqrt{\sum_{i=1}^N (X_i - E(X))^2} \sqrt{\sum_{i=1}^N (Y_i - E(Y))^2}} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \rho, \quad \text{with } 0 \leq \theta \leq \pi$$

$$(2.a) \quad R = \begin{bmatrix} 1 & r_{1,2} & r_{1,3} & \cdots & r_{1,p} \\ r_{2,1} & 1 & r_{2,3} & \cdots & r_{2,p} \\ r_{3,1} & r_{3,2} & 1 & \cdots & r_{3,p} \\ r_{4,1} & r_{4,2} & r_{4,3} & \cdots & r_{4,p} \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ r_{p,1} & r_{p,2} & r_{p,3} & \cdots & 1 \end{bmatrix},$$

(2.a). For R, a p x p correlation matrix,

² I take ‘real-world’ financial returns data to be multivariate with marginal distributions that vary notably from each other 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.

³ While *r* typically is used to represent Pearson’s calculated on a sample, *ρ* often is used to represent Pearson’s calculated on a population.

$R = BB'$ where B is the Cholesky factor (defined in Post 2) 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 $\theta_{i,j}$, we have:⁴

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

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

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

This relationship is one-to-one and bi-directional. I present below straightforward SAS/IML code translating correlations to angles (2.a) and angles to correlations (2.b) in Table A.

The above all is well-established and straightforward. But why are we interested in these angles in this setting? There are several very important reasons:

A. Because they are derived based on the matrix's Cholesky factor, the angles, unlike the correlations themselves, are forced on to the unit hyper-(hemi)sphere, where **positive definiteness automatically is enforced**. This is necessary for efficient sampling, as well as for direct and proper definition of the multivariate sample space (see Post 2 for more detail on this).

B. Crucially, the **distributions of all of the angles are independent**, which makes sampling, and more importantly, construction of their multivariate distribution (and that of the translated correlation matrix), straightforward and useable, where it otherwise would remain intractable.

⁴ Note that a similar recursive relationship exists between partial correlations (Madar, 2015), although its sample-generating algorithm it is not generalizable beyond Pearson's correlations, ie to all positive definite measures of dependence, as shown in my upcoming Post 4.

TABLE A:

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

C. **The angles contain all information regarding dependence structure** (see Fernandez-Duren & Gregorio-Dominguez, 2023, and Zhang & Songshan, 2023, as well as Opdyke, 2024). On the UNIT hyper-(hemi)sphere, the only thing we lose is scale, but scale does not and should not matter for any useful and useable measure of dependence.⁵

D. Finally, **angles distributions are more robust and** much better behaved than spectral distributions, and unlike the latter, are **at the right level of aggregation for granular scenarios** (for examples of the dramatic changes of spectral distributions under heavy-tails, see Opdyke, 2024, Burda et al., 2004, Burda et al., 2006, Akemann et al., 2009; Abul-Magd et al., 2009, Bouchaud & Potters, 2015, Martin & Mahoney, 2018), and under serial correlation, see Opdyke, 2024, and Burda et al., 2004, 2011). As discussed in Post 3, I present some empirical examples of this in numerous graphs below under real-world financial data conditions.

⁵ Scale invariance is widely proved and cited for Pearson’s rho, Kendall’s tau, and Spearman’s rho (see Xu et al., 2013, and Schreyer et al., 2017 for examples).

Fortunately, all of the above advantages of relying on angle values hold not only for Pearson's matrix, but for ANY positive definite dependence measure, under ANY data conditions found in challenging, real-world financial settings.

Beyond Pearson's: Finite Sample Distribution of ANY Positive Definite Dependence Measure

In Post 3 I discuss in more depth why angles distributions are far more appropriate than spectral (eigenvalue) distributions for solving this particular problem, and so do not revisit this comparison here other than to reemphasize points A.-D. above. In this Post 4 I focus on the fact that the only condition required for the relationships between angles and dependence measure values, as shown in (2.a) and (2.b) above, is the symmetric positive definiteness of the dependence measure. Because this approach uses the framework of all pairwise comparisons, measuring dependence on a bi-variate basis, the requirement of symmetric positive definiteness, more precisely, is the symmetric positive definiteness of the matrix of the dependence measure calculated on every pairwise association of the all the assets in the portfolio. This distinction is important to make as many dependence measures can be calculated not only on a bi-variate basis, but also on a multivariate basis, such as Szekely's distance correlation (Szekely, Rizzo, and Bakirov, 2007) and variants of Chatterjee's correlation (see Huang et al., 2022, Gamboa et al., 2022, Fuchs, 2024, and Pascual-Marqui et al., 2024, as well as Chatterjee, 2022 for a summary of the recent literature on multivariate measures). We keep to the framework of the all-pairwise matrix here for numerous reasons: as discussed in Post 3, these include tremendous flexibility, ease and directness of application, ease, if not increased power, in estimation, and ease and transparency in intervention and what-if analyses. But the main point here is that all references to positive definiteness herein and below refer to the framework of the all-pairwise matrix.

This positive definiteness (numerical issues aside) has been long proven for the "the big three," that is, for the three most widely used dependence measures – Pearson's rho, Kendall's tau, and Spearman's rho (see Sabato et al., 2007). The values of these measures all range from -1 to 1 ,⁶ but many other measures range from 0 to 1 . These include Szekely's, Lancaster's, the Tail Dependence Matrix, Chatterjee's and its many variants (see Gao and Li, 2024) and many others. Proving that these, too, are positive definite is very straightforward, and was done by Embrechts et. al. (2016) regarding the tail dependence matrix. Recall the definition of positive definiteness (for a matrix of dimension p):

if $x'Rx > 0$ for all $x \in \mathbb{R}^p \setminus \mathbf{0}$, then R is positive definite.

Because all of the (0,1) dependence measures described above are defined by

$$0 \leq R_{i,j} \leq 1 \text{ for all } i \neq j \text{ and } R_{i,i} = 1 \text{ and } R_{i,j} = R_{j,i},$$

⁶ Of course, these are maximal bounds and many conditions exist under which actual bounds are tighter. For example, for Pearson's under the equicorrelation matrix E (all equal correlations), the lower bound is $-1/(\dim[E]-1)$ rather than -1 .

$x'Rx$ can be written in quadratic form as

$$(3) \quad x'Rx = \sum_{i=1}^p x_i^2 + 2 \sum_{i=1}^{p-1} \sum_{j=i+1}^p R_{i,j} x_i x_j$$

As long as $0 < R_{i,j} < 1$ for all $i \neq j$, that is, the coefficients on the cross terms (the second term of (3)) all remain BETWEEN 0 and 1, then

$$(4) \quad \sum_{i=1}^p x_i^2 + 2 \sum_{i=1}^{p-1} \sum_{j=i+1}^p R_{i,j} x_i x_j > 0 \quad \text{and so } x'Rx > 0, \text{ always, and so R is positive definite.}$$

In the $p = 2$ case, for example, R is positive definite if $R_{1,1} > 0$ and $(R_{1,1}R_{2,2} - R_{1,2}^2) > 0$, which is always true when $0 < R_{i,j} < 1$ for all $i \neq j$ and $R_{i,i} = 1$. For the boundary cases, if $R_{i,j} = 0$ for all $i \neq j$, R obviously remains positive definite as the first term of (3) always is greater than zero and the second term disappears; and if $R_{i,j} = 1$ for all $i \neq j$ then R is positive semi-definite, although this case of perfect multivariate dependence is only textbook relevant. In practice, empirically, positive semi-definiteness only is relevant as a boundary condition, as it relates to empirical matrices that approach singularity.

Consequently, this means that all dependence measures with values ranging from 0 to 1 are, in practice, positive definite, and that NAbC can be applied to them to define their finite sample distributions. Empirical examples of this are shown in the next section.

Operationally, implementing NAbC on these (0, 1) measures is no different from implementing it on Pearson's or Kendall's or Spearman's; the (0, 1) instead of (-1, 1) range does not even change 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 above sometimes assign an angle value slightly above $\pi/2$, even though $\pi/2$ corresponds to a measure value of zero.⁷ So this is a soft upper boundary in this case, even though the measure's range of (0,1) is not.⁸ So when NAbC generates angle θ , we continue to reflect based on

$$\text{if } \theta < 0 \text{ then } \theta \leftarrow -\theta; \text{ if } \theta > \pi \text{ then } \theta \leftarrow (2\pi - \theta)$$

since for measures with a (0,1) range, the upper bound of π will never be reached, and the lower bound of

⁷ Note that angle values (which range from zero to π on the hyper-hemisphere) decrease while dependence measure values increase, so a measure value of -1 corresponds to an angle value of π , 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).

⁸ 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 hundreds of thousands 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 made empirical matrices not positive definite.

zero remains valid and hard. So NAbC applies in exactly the same way, for all of these positive definite 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 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.

Finite Sample Distribution for ANY Dependence Measure, Under ANY Real-world Data Conditions

I present below the angles distributions for some of the dependence measures discussed above, under challenging, real-world data conditions (see Opdyke, 2024 for the application of NAbC to a large number of different data conditions). Briefly, the multivariate returns distribution of the portfolio in this case 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= $3\sigma, \sigma/3, \sigma; 1/3$ observations each), and serial correlation ($AR1=-0.25, 0, 0.25, 0.50, 0.75$), with a block correlation structure shown in (5) below and $n=126$ observations, for half a year of daily returns.⁹

(5)

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
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 angle 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(n independence) baseline (see Marchenko and Pastur, 1967). The empirical results yield both expected, and some additional interesting findings.

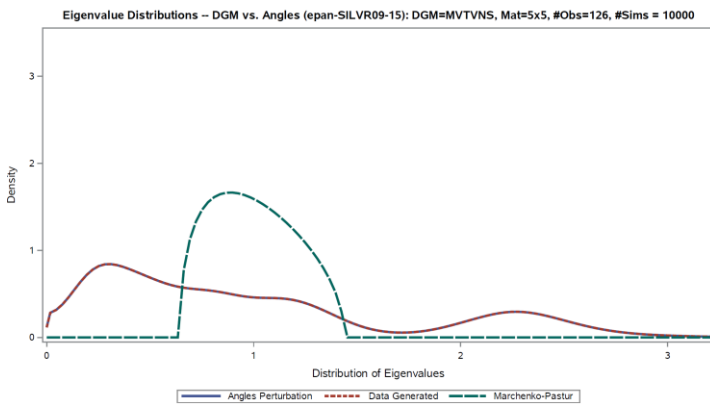
First, note that the spread, and the spread and shifts, of both the spectral and angles distributions, respectively, are larger for Pearson's than for Kendall's, which is consistent with the former's relative sensitivity to more extreme values under many conditions. The shifts and spread of both measures are much larger than those of Chatterjee,¹⁰ although this is largely due to the fact that while Chatterjee is generally more powerful under dependence that is highly nonlinear and/or highly cyclical, it is less

⁹ Note 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 on the data (and then empirically rescale the marginal densities).

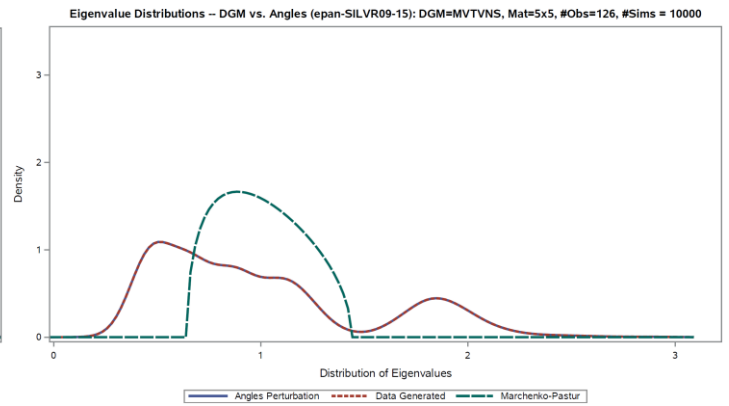
¹⁰ The symmetric version of Chatterjee's correlation coefficient is used here (see Chatterjee, 2021), with the finite sample bias correction proposed by Dalitz et. al., 2024.

Graph 1: Spectral Distribution-NAbC Angles Kernel v Data Simulations v Marchenko Pastur

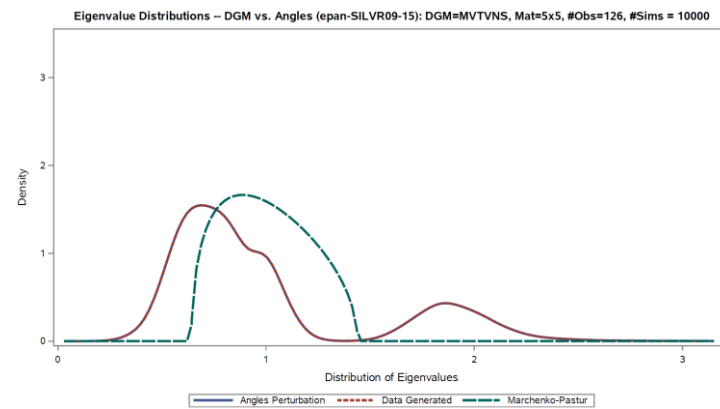
Pearson's Rho



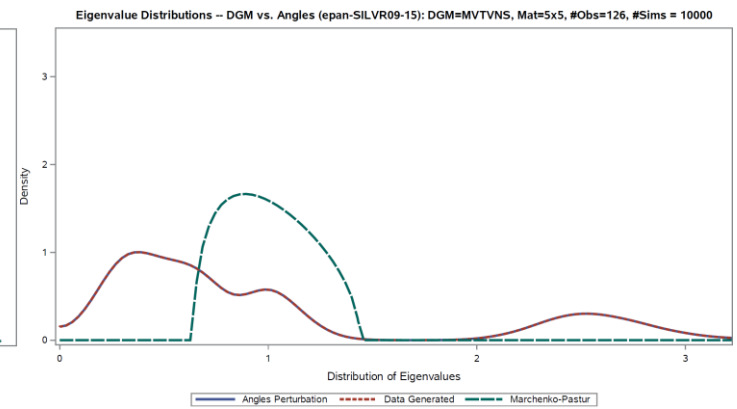
Kendall's Tau



Chatterjee's



Spearman's Rho+Chatterjee



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 (2023), which is essentially a maximum between Spearman's rho and Chatterjee's correlation, the objective being to obtain large, if not the maximum power under both types of dependence structures (i.e. strong monotonic dependence as well as highly nonlinear/cyclical 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.

Post 3 covers in detail NAbC's calculation of the dependence measure's cell level and matrix level p-values and confidence intervals, which I will not be repeat here because it is identical regardless of the dependence measure being examined (Post 3 covered only Pearson's matrix). However, I take the example from Post 3 for Pearson's matrix, which included such p-values and confidence intervals, and recreate it here using Kendall's Tau. Matrix input values necessarily are slightly different, but all other aspects of the example remain the same to demonstrate the apples-to-apples seamless application of

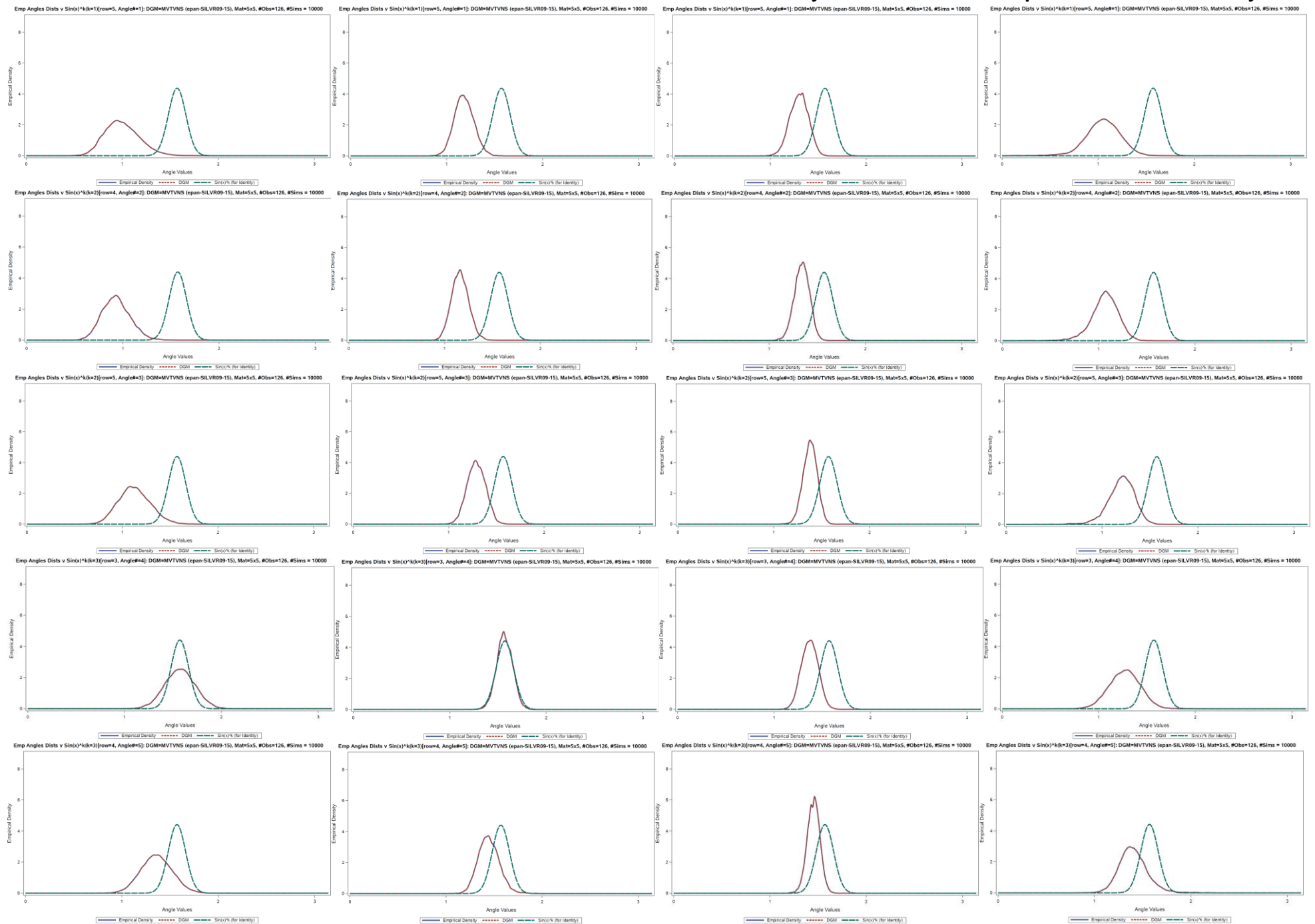
Graphs 1-5: Angles Distributions--NAbC Angles Kernel v Data Simulations v Identity Matrix

Pearson's Rho

Kendall's Tau

Chatterjee's

Spearman's Rho+Chatterjee



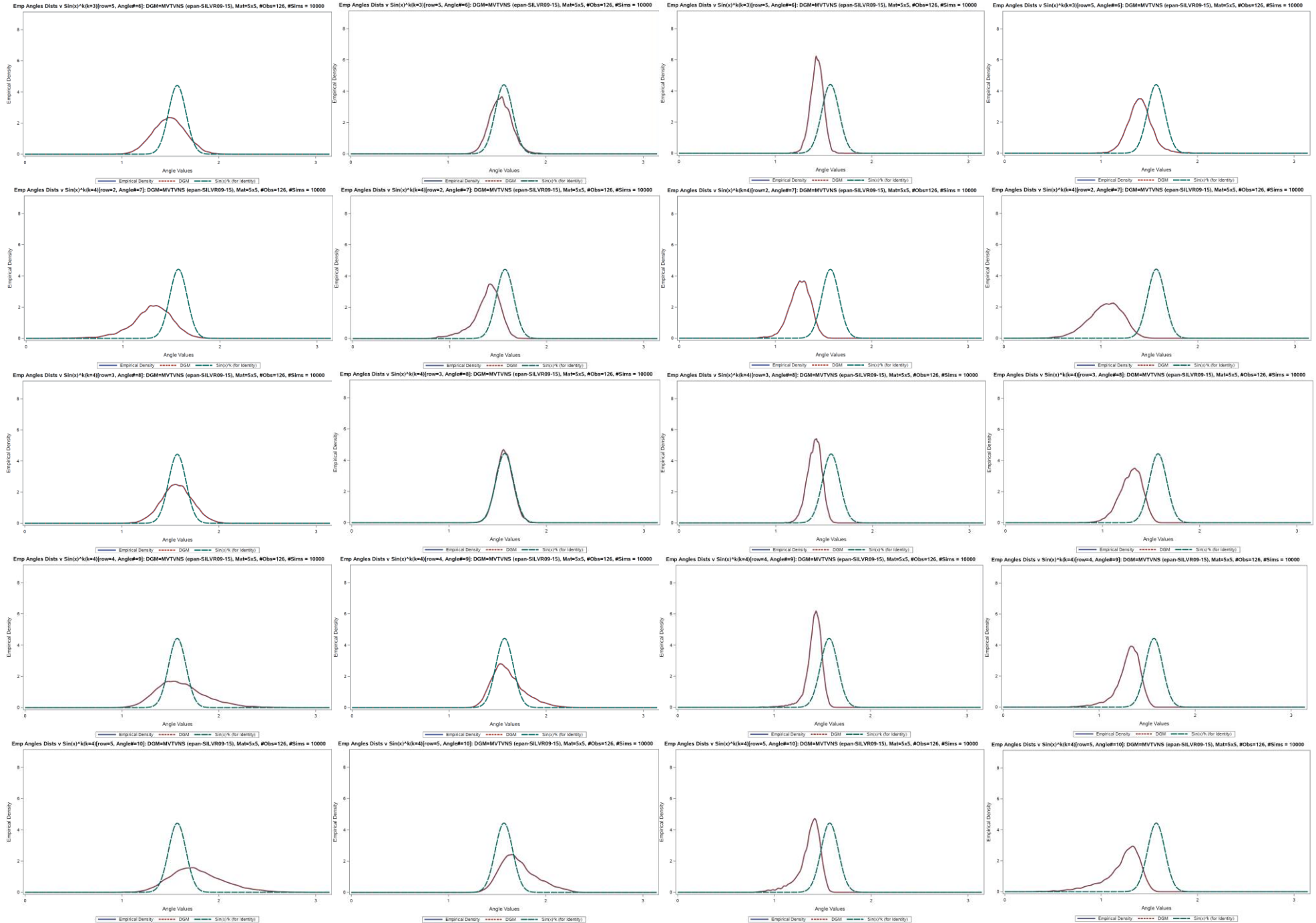
Graphs 6-10: Angles Distributions--NAbC Angles Kernel v Identity Matrix

Pearson's Rho

Kendall's Tau

Chatterjee's

Spearman's Rho+Chatterjee



NAbC Applied with Kendall's: Unrestricted + Scenario-Restricted p-values and Confidence Intervals

Below I apply NAbC to obtain both p-values and confidence intervals, for Kendall's Tau, under two cases: unrestricted, and scenario-restricted. Solely for ease of replication, the data generating mechanism for these examples is simply multivariate standard normal, with N=25k simulations and number of observations n = 160.

UNRESTRICTED CASE: Given a specified or well-estimated correlation 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 correlation 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 and every cell of [A]?
- Q2. **Quantile Function:** What is the unique correlation matrix associated with [B], a matrix of cumulative distribution function values associated with the corresponding cells of [A]?
- Q3. **p-values:** Under the null hypothesis that observed correlation 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 and every cell of [C]?

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

¹¹ Values used here for Kendall's matrix were close to those obtained when translating from the Post 3 Pearson's example using $\tau = (2/\pi) \arcsin(r)$ where r = Pearson's, which is generally valid under elliptical data (which is one of the reasons I used multivariate Gaussian data here; see McNeil et. al., 2005).

matrix [D].

[D]

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

[E]

	0.8			
	0.8	0.8	0.7	

[F]

1				
0.13	1			
-0.06	0.350	1		
0.19	-0.19	-0.06	1	
0.41	0.180	0.125	0.215	1

- Q4. **Confidence Intervals:** What are the two correlation 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 only those cells of [D] that are relevant to the scenario (green)?
- Q5. **Quantile Function:** What is the unique correlation matrix associated with [E], a matrix of cumulative distribution function values associated with the corresponding cells of [D]?
- Q6. **p-values:** Under the null hypothesis that observed correlation 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 every (non-constant, green) cell of [F]?

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, under any data conditions. Only NAbC can simultaneously answer Q1.-Q6. above under general data conditions, as shown below.

Q1	Q2	Q3	Q4	Q5	Q6																																																																																																				
	<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.1729</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.0355</td><td>0.2369</td><td>1</td><td></td><td></td></tr> <tr><td>0.2374</td><td>-0.1510</td><td>-0.0392</td><td>1</td><td></td></tr> <tr><td>0.4335</td><td>0.3159</td><td>0.0614</td><td>0.1040</td><td>1</td></tr> </table>	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	<p>p-value=0.1503</p> <table border="1"> <tr><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>0.0006</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.0090</td><td>0.0218</td><td></td><td></td><td></td></tr> <tr><td>0.0222</td><td>0.0315</td><td>0.0269</td><td></td><td></td></tr> <tr><td>0.0170</td><td>0.0157</td><td>0.0088</td><td>0.0077</td><td></td></tr> </table>						0.0006					0.0090	0.0218				0.0222	0.0315	0.0269			0.0170	0.0157	0.0088	0.0077			<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.1282</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.0636</td><td>0.2398</td><td>1</td><td></td><td></td></tr> <tr><td>0.1942</td><td>-0.1940</td><td>-0.0639</td><td>1</td><td></td></tr> <tr><td>0.4098</td><td>0.2895</td><td>0.0362</td><td>0.0867</td><td>1</td></tr> </table>	1					0.1282	1				-0.0636	0.2398	1			0.1942	-0.1940	-0.0639	1		0.4098	0.2895	0.0362	0.0867	1	<p>p-value=0.0436</p> <table border="1"> <tr><td></td><td></td><td></td><td></td><td></td></tr> <tr><td></td><td></td><td></td><td></td><td></td></tr> <tr><td></td><td>0.0047</td><td></td><td></td><td></td></tr> <tr><td></td><td></td><td></td><td></td><td></td></tr> <tr><td></td><td>0.0077</td><td>0.0148</td><td>0.0171</td><td></td></tr> </table>												0.0047										0.0077	0.0148	0.0171	
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																																																																																																					
0.0006																																																																																																									
0.0090	0.0218																																																																																																								
0.0222	0.0315	0.0269																																																																																																							
0.0170	0.0157	0.0088	0.0077																																																																																																						
1																																																																																																									
0.1282	1																																																																																																								
-0.0636	0.2398	1																																																																																																							
0.1942	-0.1940	-0.0639	1																																																																																																						
0.4098	0.2895	0.0362	0.0867	1																																																																																																					
	0.0047																																																																																																								
	0.0077	0.0148	0.0171																																																																																																						
<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>-0.0172</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.2100</td><td>0.0626</td><td>1</td><td></td><td></td></tr> <tr><td>0.0472</td><td>-0.3567</td><td>-0.1602</td><td>1</td><td></td></tr> <tr><td>0.2794</td><td>0.0926</td><td>-0.1873</td><td>-0.0830</td><td>1</td></tr> </table>	1					-0.0172	1				-0.2100	0.0626	1			0.0472	-0.3567	-0.1602	1		0.2794	0.0926	-0.1873	-0.0830	1	<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.2735</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>0.0910</td><td>0.3475</td><td>1</td><td></td><td></td></tr> <tr><td>0.3323</td><td>0.0127</td><td>0.1182</td><td>1</td><td></td></tr> <tr><td>0.5250</td><td>0.4370</td><td>0.2335</td><td>0.2789</td><td>1</td></tr> </table>	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		<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.1282</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.0636</td><td>0.0478</td><td>1</td><td></td><td></td></tr> <tr><td>0.1942</td><td>-0.1940</td><td>-0.0639</td><td>1</td><td></td></tr> <tr><td>0.4098</td><td>0.1757</td><td>-0.1144</td><td>-0.0541</td><td>1</td></tr> </table>	1					0.1282	1				-0.0636	0.0478	1			0.1942	-0.1940	-0.0639	1		0.4098	0.1757	-0.1144	-0.0541	1	<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.1282</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.0636</td><td>0.3492</td><td>1</td><td></td><td></td></tr> <tr><td>0.1942</td><td>-0.1940</td><td>-0.0639</td><td>1</td><td></td></tr> <tr><td>0.4098</td><td>0.3425</td><td>0.1150</td><td>0.1841</td><td>1</td></tr> </table>	1					0.1282	1				-0.0636	0.3492	1			0.1942	-0.1940	-0.0639	1		0.4098	0.3425	0.1150	0.1841	1	
1																																																																																																									
-0.0172	1																																																																																																								
-0.2100	0.0626	1																																																																																																							
0.0472	-0.3567	-0.1602	1																																																																																																						
0.2794	0.0926	-0.1873	-0.0830	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																																																																																																									
0.1282	1																																																																																																								
-0.0636	0.0478	1																																																																																																							
0.1942	-0.1940	-0.0639	1																																																																																																						
0.4098	0.1757	-0.1144	-0.0541	1																																																																																																					
1																																																																																																									
0.1282	1																																																																																																								
-0.0636	0.3492	1																																																																																																							
0.1942	-0.1940	-0.0639	1																																																																																																						
0.4098	0.3425	0.1150	0.1841	1																																																																																																					
<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.0250</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.1661</td><td>0.0986</td><td>1</td><td></td><td></td></tr> <tr><td>0.0926</td><td>-0.3131</td><td>-0.1427</td><td>1</td><td></td></tr> <tr><td>0.3210</td><td>0.1410</td><td>-0.1396</td><td>-0.0478</td><td>1</td></tr> </table>	1					0.0250	1				-0.1661	0.0986	1			0.0926	-0.3131	-0.1427	1		0.3210	0.1410	-0.1396	-0.0478	1	<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.2300</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>0.0424</td><td>0.3013</td><td>1</td><td></td><td></td></tr> <tr><td>0.2920</td><td>-0.0525</td><td>0.0570</td><td>1</td><td></td></tr> <tr><td>0.4929</td><td>0.3849</td><td>0.1611</td><td>0.2103</td><td>1</td></tr> </table>	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		<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.1282</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.0636</td><td>0.0904</td><td>1</td><td></td><td></td></tr> <tr><td>0.1942</td><td>-0.1940</td><td>-0.0639</td><td>1</td><td></td></tr> <tr><td>0.4098</td><td>0.2028</td><td>-0.0789</td><td>-0.0184</td><td>1</td></tr> </table>	1					0.1282	1				-0.0636	0.0904	1			0.1942	-0.1940	-0.0639	1		0.4098	0.2028	-0.0789	-0.0184	1	<table border="1"> <tr><td>1</td><td></td><td></td><td></td><td></td></tr> <tr><td>0.1282</td><td>1</td><td></td><td></td><td></td></tr> <tr><td>-0.0636</td><td>0.3006</td><td>1</td><td></td><td></td></tr> <tr><td>0.1942</td><td>-0.1940</td><td>-0.0639</td><td>1</td><td></td></tr> <tr><td>0.4098</td><td>0.3165</td><td>0.0809</td><td>0.1486</td><td>1</td></tr> </table>	1					0.1282	1				-0.0636	0.3006	1			0.1942	-0.1940	-0.0639	1		0.4098	0.3165	0.0809	0.1486	1	
1																																																																																																									
0.0250	1																																																																																																								
-0.1661	0.0986	1																																																																																																							
0.0926	-0.3131	-0.1427	1																																																																																																						
0.3210	0.1410	-0.1396	-0.0478	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																																																																																																					
1																																																																																																									
0.1282	1																																																																																																								
-0.0636	0.0904	1																																																																																																							
0.1942	-0.1940	-0.0639	1																																																																																																						
0.4098	0.2028	-0.0789	-0.0184	1																																																																																																					
1																																																																																																									
0.1282	1																																																																																																								
-0.0636	0.3006	1																																																																																																							
0.1942	-0.1940	-0.0639	1																																																																																																						
0.4098	0.3165	0.0809	0.1486	1																																																																																																					

For Q1 and Q4, 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 Q6, the smaller p-values for the individual cells compared to the respective matrix-level p-values, which are larger, as expected, as they control the family-wise error rate (FWER – see Post 3). Note also that the green cells of Q5 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, $R = BB^T$. Finally, note that the empirical values of the red cells in Q4-Q6 differ slightly from those in [D] and [F]. This is due to NAbC's conservative use of the mean of the estimated correlation matrices, rather than presuming we know the absolute 'true' values of these cells (although this is justified in some specific cases).

NAbC Remains “Estimator Agnostic”

Although this has been covered in previous Posts, it bears repeating that, regardless of the dependence measure being used, NAbC remains “estimator agnostic,” that is, valid for use with any reasonable estimator of that dependence structure. 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 and shifting combinations of these characteristics. 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 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.

LNP: A Generalized Entropy for All Positive Definite Dependence Measures

The (two-sided) p-values NAbC provides (see Q3 and Q6 above, and Post 3 for details) actually can be used to construct a competitor to commonly used distance metrics, such as norms (e.g. Taxi, Frobenius/Euclidean, and Chebyshev norms: see Post 3 for definitions), and has a number of advantages over them in this setting. Norms measure absolute distance from a presumed or baseline correlation value. But the range of all relevant and widely used dependence measures is bounded, either from -1 to 1 or 0 to 1, and 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 0.01 from an original or presumed correlation value of, say, 0.97, means something very different than the same shift from 0.07. NAbC's p-values attribute probabilistic MEANING to these two different cases, while a norm would 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 (6) below:

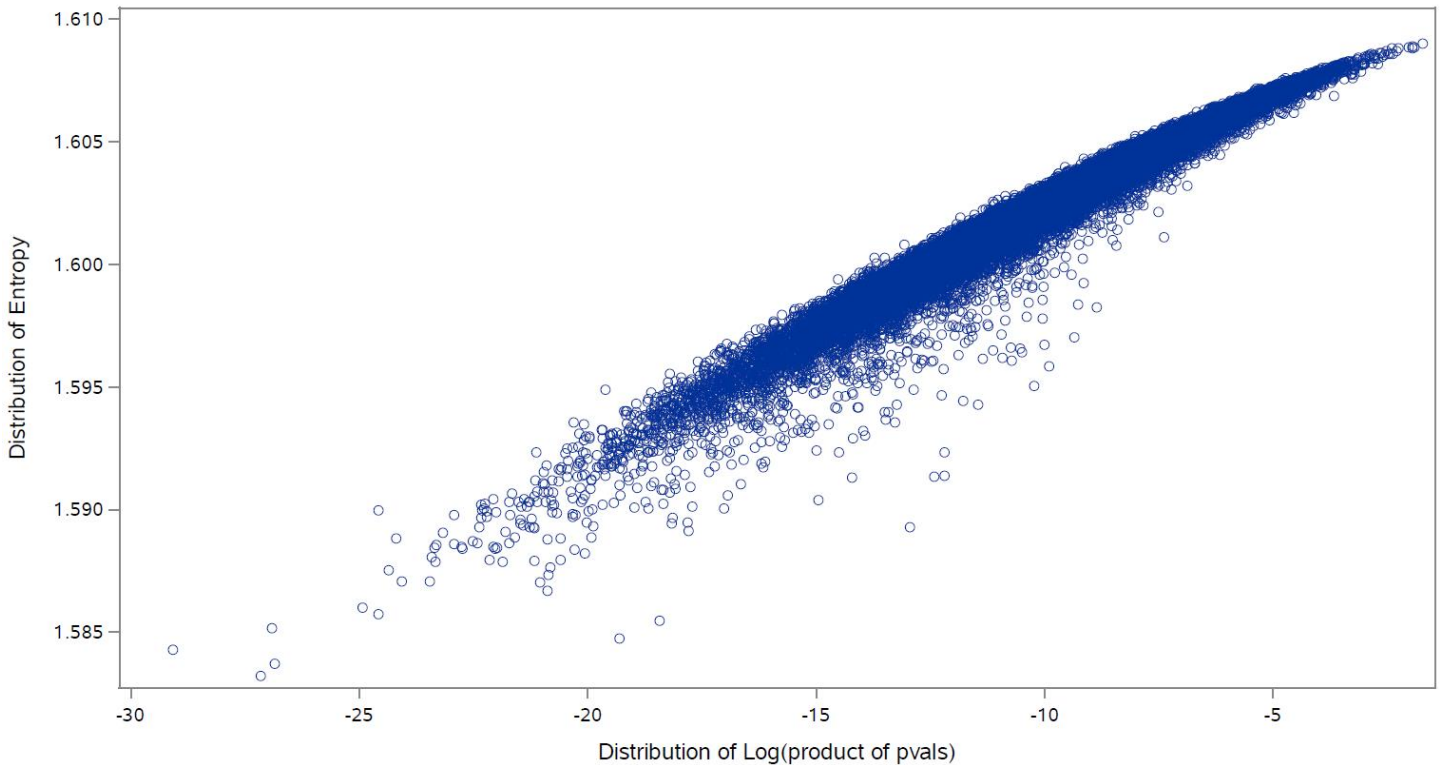
$$(6) \text{ "LNP"} = \ln \left(\prod_{i=1}^q p\text{-value}_i \right) = \sum_{i=1}^q \ln [p\text{-value}_i] \text{ where } q = p(p-1)/2 \text{ and } p\text{-value}_i \text{ is 2-sided.}$$

This was shown in Post 2, using a Pearson's correlation matrix under the (Gaussian) identity matrix, to have a very strong correspondence with the entropy of the correlation matrix, defined by Felipe et al. (2021 and 2023) as (7) below:

$$(7) \text{ Entropy} = Ent(R/p) = - \sum_{j=1}^p \lambda_j \ln(\lambda_j)$$

where R is the sample correlation matrix and λ_j are the p eigenvalues of the correlation matrix after it is scaled by its dimension, R/p. Importantly, this result (7), like NAbC, is valid for ANY positive definite measure of dependence, not just Pearson's. Graph 12 below compares LNP to the entropy of the Kendall's Tau matrix in 10,000 simulations under the Gaussian identity matrix. The resulting Pearson's correlation between them is 0.98.

Graph 12: Identity Matrix Simulations – LNP (based on Kendall's) v Entropy



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.¹² 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/distance/independence/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, as specified by the researcher (i.e. the specified dependence matrix), that is not necessarily perfect (in)dependence. Such 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, 2010, Almog and Shmueli, 2019, Chakraborti et al., 2020, and Vorobets, 2024a, 2024b, for several examples). Interpretations aside, 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 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.

Conclusion

In Posts 1,2, and 3 I listed the seven characteristics of the full NAbC solution that, taken together, are shared by no other approach, and for completeness I list them again below:

1. validity under challenging, real-world financial data conditions, with marginal asset distributions characterized by notably different degrees of serial correlation, non-stationarity, heavy-tailedness, and asymmetry
2. application to ANY positive definite dependence measure, including, for example, Pearson's product moment correlation, rank-based measures like Kendall's tau and Spearman's rho, the kernel-based generalization of Szekely's distance correlation, and the tail dependence matrix, among others.
3. it remains "estimator agnostic," that is, valid regardless of the sample-based estimator used to estimate any of the above-mentioned dependence measures
4. it 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

¹² 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.

that of the entire matrix, and the same is true for the p-values; this also effectively facilitates attribution analyses)

5. it provides a one-to-one quantile function, translating a matrix of all the cells' cdf values to a (unique) correlation (dependence measure) matrix, and back again, enabling precision in reverse scenarios and stress testing

6. all the above results remain valid even when selected cells in the matrix are 'frozen' for a given scenario or stress test, while the rest are allowed to vary, enabling granular and realistic scenarios

7. it 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, enabling reliable application in actual, imperfect, non-textbook settings

Post 2 provided, for Pearson's under the Gaussian identity matrix, an interactive spreadsheet that implements fully analytic p-values and confidence intervals

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

and combined with Post 3, both cover all but 2. in the list above. This Post 4 covers 2. above, expanding NAbC's range of application to ALL positive definite measures of dependence, with any values, under all real-world data conditions.

The objective of this work has been to provide a method that checks all of these boxes – 1. Through 7. – simultaneously, which is what is required for useful and useable portfolio analytics in real-world, non-textbook settings. 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 narrow ranges of application.

With NAbC, we now have a powerful, applied approach enabling us to treat an extremely broad class of widely used dependence measures just like the other major parameters in our (finite sample) financial portfolio models. We can use NAbC in frameworks that identify, measure and monitor, and even anticipate critically important events, such as correlation breakdowns, and mitigate and manage their effects. It should prove to be a very useful means by which we can better understand, predict, and manage portfolios in our multivariate world.

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.
- Akemann, G., Fischmann, J., and Vivo, P., (2009), "Universal Correlations and Power-Law Tails in Financial Covariance Matrices," <https://arxiv.org/abs/0906.5249>.
- Almog, A., and Shmueli, E., (2019), "Structural Entropy: Monitoring Correlation-Based Networks over time With Application to Financial Markets," *Scientific Reports*, 9:10832.
- Bouchaud, J., & Potters, M., (2015), "Financial applications of random matrix theory: a short review," *The Oxford Handbook of Random Matrix Theory*, Eds G. Akemann, J. Baik, P. Di Francesco.
- 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.
- 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., (2022), "A Survey of Some Recent Developments in Measures of Association," ArXiv preprint, arXiv:2211.04702.
- Church, Christ (2012). "The asymmetric t-copula with individual degrees of freedom", Oxford, UK: University of Oxford Master Thesis, 2012.
- Dalitz, C., Arning, J., and Goebbels, S., (2024), "A Simple Bias Reduction for Chatterjee's Correlation," arXiv:2312.15496v2.
- Embrechts, P., Hofert, M., and Wang, R., (2016), "Bernoulli and Tail-Dependence Compatibility," *The Annals of Applied Probability*, Vol. 26(3), 1636-1658.
- 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.

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.

Fuchs, S., (2024), "Quantifying Directed Dependence via Dimension Reduction," *Journal of Multivariate Analysis*, 201:105266.

Gamboa, F., Gremaud, P., Klein, T., and Lagnoux, A., (2022), "Global Sensitivity Analysis: A Novel Generation of Mighty Estimators Based on Rank Statistics," *Bernoulli*, 28(4):2345–2374.

Gao, M., Li, Q., (2024), "A Family of Chatterjee's Correlation Coefficients and Their Properties," arXiv:2403.17670v1 [stat.ME]

Holzmann, H., and Klar, B., (2024) "Lancaster Correlation - A New Dependence Measure Linked to Maximum Correlation," arXiv:2303.17872v2 [stat.ME].

Huang, Z., Deb, N., and Sen, B., (2022), "Kernel Partial Correlation Coefficient – A Measure of Conditional Dependence," *The Journal of Machine Learning Research*, 23(1):9699–9756.

Kendall, M. (1938), "A New Measure of Rank Correlation," *Biometrika*, 30 (1–2), 81–89.

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.

Madar, V., (2015), "Direct Formulation to Cholesky Decomposition of a General Nonsingular Correlation Matrix," *Statistics & Probability Letters*, Vol 103, pp.142-147.

Marchenko, A., Pastur, L., (1967), "Distribution of eigenvalues for some sets of random matrices," *Matematicheskii Sbornik*, N.S. 72 (114:4): 507–536.

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., (2005). Quantitative Risk Management: Concepts, Techniques, and Tools, Princeton, NJ: Princeton University Press.

Meucci, A., (2010), "Fully Flexible Views: Theory and Practice," arXiv:1012.2848v1

Opdyke, JD, (2024), Keynote Address: "Beating the Correlation Breakdown, for Pearson's and Beyond: Robust Inference and Flexible Scenarios and Stress Testing for Financial Portfolios," QuantStrats11, NYC, March 12.

- 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.
- Pinheiro, J. and Bates, D. (1996), "Unconstrained parametrizations for variance-covariance matrices," *Statistics and Computing*, Vol. 6, 289–296.
- Pourahmadi, M., Wang, X., (2015), "Distribution of random correlation matrices: Hyperspherical parameterization of the Cholesky factor," *Statistics and Probability Letters*, 106, (C), 5-12.
- Rapisarda, F., Brigo, D., & Mercurio, F., (2007), "Parameterizing Correlations: A Geometric Interpretation," *IMA Journal of Management Mathematics*, 18(1), 55-73.
- Rebonato, R., and Jackel, P., (2000), "The Most General Methodology for Creating a Valid Correlation Matrix for Risk Management and Option Pricing Purposes," *Journal of Risk*, 2(2)17-27.
- 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.
- 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.
- Shyamalkumar, N., and Tao, S., (2020), "On tail dependence matrices: The realization problem for parametric families," *Extremes*, Vol. 23, 245–285.
- Spearman, C., (1904), "'General Intelligence,' Objectively Determined and Measured," *The American Journal of Psychology*, 15(2), 201–292.
- Szekely, G., Rizzo, M., and Bakirov, N., (2007), "Measuring and Testing Dependence by Correlation of Distances," *The Annals of Statistics*, 35(6), pp2769-2794.
- Vorobets, A., (2024a), "Sequential Entropy Pooling Heuristics,"
<https://ssrn.com/abstract=3936392> or <http://dx.doi.org/10.2139/ssrn.3936392>
- Vorobets, A., (2024b), "Portfolio Construction and Risk Management,"
<https://ssrn.com/abstract=4807200> or <http://dx.doi.org/10.2139/ssrn.4807200>
- Zhang, Q., (2023), "On relationships between Chatterjee's and Spearman's correlation coefficients," arXiv:2302.10131v1 [stat.ME]
- Zhang, Y., and Songshan, Y., (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.