

Rehabilitating Correlations

Nick Higham
Department of Mathematics
The University of Manchester

<http://www.maths.manchester.ac.uk/~higham>

Slides available at https://bit.ly/rehab_corr

Correlation Matrix (Statistical)

Vectors $x_1, x_2, \dots, x_n \in \mathbb{R}^m$. **Covariance matrix**

$$v_{ij} = \text{cov}(x_i, x_j) = \frac{1}{n-1} (x_i - \bar{x}_i)^T (x_j - \bar{x}_j),$$

where $\bar{x}_i = \text{mean}(x_i)$. **Correlation matrix**

$$C = D^{-1/2} V D^{-1/2},$$

where $D = \text{diag}(v_{ii})$. $c_{ij} = v_{ii}^{-1/2} v_{ij} v_{jj}^{-1/2}$ is correlation between x_i and x_j .

Correlation Matrix (Mathematical)

An $n \times n$ symmetric matrix A is a **correlation matrix** if

- It has ones on the diagonal.
- All its eigenvalues are nonnegative.

Correlation Matrix (Mathematical)

An $n \times n$ symmetric matrix A is a **correlation matrix** if

- It has ones on the diagonal.
- All its eigenvalues are nonnegative.

Useful fact

For any positive semidefinite matrix, $|a_{ij}| \leq \sqrt{a_{ii}a_{jj}}$, so

$$\max_{i \neq j} |a_{ij}| \leq 1.$$

One-Parameter Family

For what values of w is

$$A_3 = \begin{bmatrix} 1 & w & w \\ w & 1 & w \\ w & w & 1 \end{bmatrix}$$

a correlation matrix?

One-Parameter Family

For what values of w is

$$A_3 = \begin{bmatrix} 1 & w & w \\ w & 1 & w \\ w & w & 1 \end{bmatrix}$$

a correlation matrix?

Answer is $-1/2 \leq w \leq 1$.

One-Parameter Family

For what values of w is

$$A_3 = \begin{bmatrix} 1 & w & w \\ w & 1 & w \\ w & w & 1 \end{bmatrix}$$

a correlation matrix?

Answer is $-1/2 \leq w \leq 1$.

For $n \times n$ matrix the condition is

$$-\frac{1}{n-1} \leq w \leq 1.$$

Spectrum of Correlation Matrix

Theorem (Schur, Horn)

A necessary and sufficient condition for a symmetric $n \times n$ A to have e'vals $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ and diagonal elements $\alpha_1 \leq \alpha_2 \leq \dots \leq \alpha_n$ (in any order along the diagonal) is that

$$\sum_{i=1}^k \lambda_i \leq \sum_{i=1}^k \alpha_i, \quad k = 1 : n,$$

with equality for $k = n$.

Spectrum of Correlation Matrix

Theorem (Schur, Horn)

A necessary and sufficient condition for a symmetric $n \times n$ A to have e'vals $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ and diagonal elements $\alpha_1 \leq \alpha_2 \leq \dots \leq \alpha_n$ (in any order along the diagonal) is that

$$\sum_{i=1}^k \lambda_i \leq \sum_{i=1}^k \alpha_i, \quad k = 1 : n,$$

with equality for $k = n$.

Conclusion: For a correlation matrix
any set of $\lambda_i \geq 0$ summing to n is possible.

Generating a Random Correlation Matrix (1)

Algorithm of **Bendel & Mickey (1978)**,
improved by **Davies & H (2000)**.

Implemented in MATLAB as `gallery('randcorr', x)`.

```
>> A = gallery('randcorr', [1 1 1])
```

```
A =
```

```
 1.0000e+00    1.6653e-16   -1.6653e-16  
 1.6653e-16    1.0000e+00   -6.9389e-17  
-1.6653e-16   -6.9389e-17    1.0000e+00
```

Generating a Random Correlation Matrix (2)

```
>> A = gallery('randcorr', [3 0 0])
```

```
A =
```

```
1.0000e+00    1.0000e+00    1.0000e+00  
1.0000e+00    1.0000e+00    1.0000e+00  
1.0000e+00    1.0000e+00    1.0000e+00
```

Generating a Random Correlation Matrix (3)

```
>> A = gallery('randcorr', [2 1 0])
```

```
A =
```

```
    1.0000e+00    8.3822e-15   -9.9938e-01  
    8.3822e-15    1.0000e+00    3.5272e-02  
   -9.9938e-01    3.5272e-02    1.0000e+00
```

```
>> eig(A)
```

```
ans =
```

```
    1.1102e-16  
    1.0000e+00  
    2.0000e+00
```

Perturbation Problem

Consider

$$\begin{bmatrix} 1 & 0.9 & ? \\ 0.9 & 1 & 0.9 \\ ? & 0.9 & 1 \end{bmatrix} .$$

Perturbation Problem

Consider

$$\begin{bmatrix} 1 & 0.9 & ? \\ 0.9 & 1 & 0.9 \\ ? & 0.9 & 1 \end{bmatrix}.$$

Is this a correlation matrix:

$$A = \begin{bmatrix} 1 & 0.9 & 0 \\ 0.9 & 1 & 0.9 \\ 0 & 0.9 & 1 \end{bmatrix}.$$

Perturbation Problem

Consider

$$\begin{bmatrix} 1 & 0.9 & ? \\ 0.9 & 1 & 0.9 \\ ? & 0.9 & 1 \end{bmatrix}.$$

Is this a correlation matrix:

$$A = \begin{bmatrix} 1 & 0.9 & 0 \\ 0.9 & 1 & 0.9 \\ 0 & 0.9 & 1 \end{bmatrix}.$$

Spectrum: $-0.2728, 1.0000, 2.2728$.

Can we find a small perturbation that makes A a correlation matrix?

Completion Problem

For what values of θ is

$$A(\theta) = \begin{bmatrix} 1 & 0.9 & \theta \\ 0.9 & 1 & 0.9 \\ \theta & 0.9 & 1 \end{bmatrix}.$$

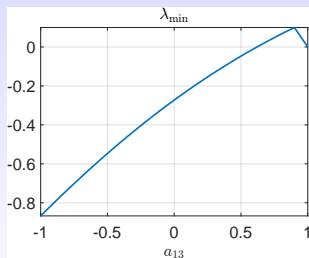
a correlation matrix?

Completion Problem

For what values of θ is

$$A(\theta) = \begin{bmatrix} 1 & 0.9 & \theta \\ 0.9 & 1 & 0.9 \\ \theta & 0.9 & 1 \end{bmatrix}.$$

a correlation matrix?



$$\lambda_{\min} \geq 0 \text{ for } a_{13} \in [0.62, 1]$$

Question from London Fund Management Company (2000)

“Given a real symmetric matrix A which is almost a correlation matrix . . .

- What is the best approximating (in Frobenius norm?) correlation matrix?
- Is it unique?
- Can we compute it?

Typically we are working with 1400×1400 at the moment, but this will probably grow to 6500×6500 .”

“Given a real symmetric matrix A which is almost a correlation matrix what is the best approximating (in Frobenius norm?) correlation matrix?”

“I am researching ways to make our company’s correlation matrix positive semi-definite.”

“Currently, I am trying to implement some real options multivariate models in a simulation framework. Therefore, I estimate correlation matrices from inconsistent data set which eventually are non psd.”

Indefinite Covariance/Correlation Matrices

A (complete) matrix may lack definiteness due to (e.g.)

- missing observations
(companies formed or closed during the period),
- asynchronous observations
(different assets sampled at different time points),
- stress testing.

What's the problem? An “indefinite correlation matrix” leads to an investment portfolio with negative risk.

How do we make the matrix (semi)definite?

How to Proceed

× Make ad hoc modifications to matrix: e.g.,

$$A \leftarrow A + \alpha I$$

$$A \leftarrow A / (1 + \alpha)$$

[In general $A \leftarrow D^{-1/2} A D^{-1/2}$ where $D = \text{diag}(A)$.]

How to Proceed

✗ Make ad hoc modifications to matrix: e.g.,

$$A \leftarrow A + \alpha I$$

$$A \leftarrow A / (1 + \alpha)$$

[In general $A \leftarrow D^{-1/2} A D^{-1/2}$ where $D = \text{diag}(A)$.]

✓ Plug the gaps in the missing data, then compute an exact correlation matrix.

How to Proceed

- ✗ Make ad hoc modifications to matrix: e.g.,
 $A \leftarrow A + \alpha I$
 $A \leftarrow A / (1 + \alpha)$
[In general $A \leftarrow D^{-1/2} A D^{-1/2}$ where $D = \text{diag}(A)$.]
- ✓ Plug the gaps in the missing data, then compute an exact correlation matrix.
- ✓ Optimally perturb the matrix.

Nearest Correlation Matrix Problem

With $\|A\|_F^2 = \sum_{i,j} a_{ij}^2$,

$$\min \{ \|A - C\|_F : C \text{ is a correlation matrix} \}.$$

- Correlation matrices are the intersection of the sets

$$S = \{ Y = Y^T \in \mathbb{R}^{n \times n} : \lambda_i(Y) \geq 0, i = 1:n \},$$

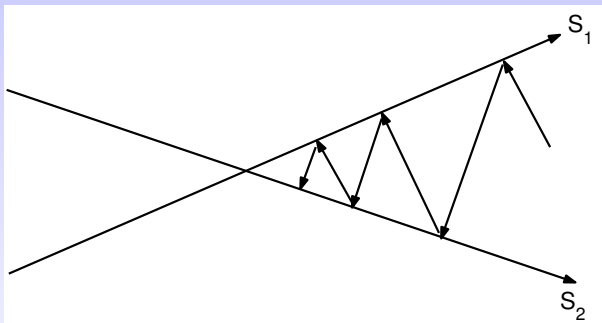
$$U = \{ Y = Y^T \in \mathbb{R}^{n \times n} : y_{ii} = 1, i = 1:n \}.$$

- Constraint set closed & convex, so unique minimizer.
- Various ad-hoc solution methods proposed in 1990s. None guaranteed to compute the solution.
- Theory and algorithm: **H** (2002).

Alternating Projections Algorithm

von Neumann (1933): for subspaces.

Dykstra (1983): corrections for closed convex sets.



- Easy to implement.
- Guaranteed convergence, at a linear rate.
- Can add further constraints/projections.

Newton Method

Qi & Sun (2006): **Newton method** based on theory of strongly semismooth matrix functions.

- Applies Newton to **dual** (unconstrained) of $\min \frac{1}{2} \|A - X\|_F^2$ problem.
- **Globally** and **quadratically** convergent.
- **H & Borsdorf (2010)**: improves to efficiency and reliability.

The basis of **G02AA** (nearest correlation matrix) in NAG Library Mark 22 (2009). See

[NAGPythonExamples/nearest_correlation_matrices](#)

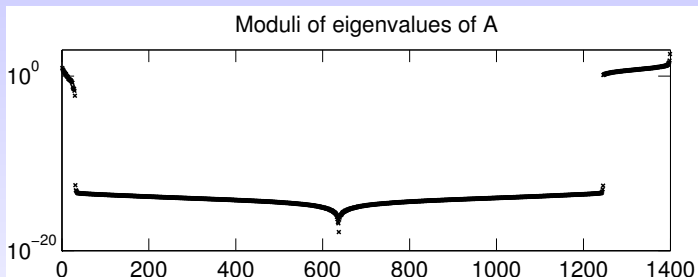
Eigenvalues of Nearest Correlation Matrix

Theorem (H, 2002)

If A has t nonpositive eigenvalues then the nearest correlation matrix has at least t zero eigenvalues.

Numerical Example from Finance

A : stock data, $n = 1399$. $a_{ii} \equiv 1$, $|a_{ij}| \leq 1$, but not psd.
 A has with 1245 nonpositive ei'vals $\Rightarrow \text{rank}(X) \leq 154$.



$$\|A - X\|_F = 20.96.$$

Alternating Projections vs Newton

- MATLAB code for alternating projections.
- NAG Library Newton code g02aa.

Matrix	n	tol	Code	Time (s)	Iters
1. Random	100	1e-10	g02aa	0.023	4
			alt proj	0.052	15
2. Random	500	1e-10	g02aa	0.48	4
			alt proj	3.01	26
3. Real-life	1399	1e-4	g02aa	6.8	5
			alt proj	100.6	68

Collection of Invalid Correlation Matrices

<https://github.com/higham/matrices-correlation-invalid> includes

- bccd16: 3250-by-3250 matrix from *EU bank data*.
- bhwi01: 5-by-5 matrix relating to *portfolio risk*.
- cor3120: 3120-by-3120 matrix from *stock data*.
- fing97: 7-by-7 matrix from *stress testing*.
- mmb13: 6-by-6 matrix from *foreign exchange trading* data supplied by the Royal Bank of Scotland.
- tec03: 4-by-4 matrix from *stress testing*.
- usgs13: 94-by-94 from national assessment of *carbon dioxide storage* resources in the USA.

Unexpected Applications of NCM

Some recent papers:

- **Applying stochastic small-scale damage functions to German winter storms** (2012)
- **Estimating variance components and predicting breeding values for eventing disciplines and grades in sport horses** (2012)
- **Characterisation of tool marks on cartridge cases by combining multiple images** (2012)
- **Experiments in reconstructing twentieth-century sea levels** (2011)
- **Media framing of Copenhagen tourism: A new approach to public opinion about tourists** (2020)

Completion Problem

- Business units 1 & 2 exposed to risks x and y .
- Only BU_1 is exposed to risk z .
- Correlations not specified between x and y in BU_2 and z in BU_1 : few experts can make judgments about cross-business risks.

Correlations		x	y	z	x	y
BU_1	x	1	0.7	0.85	0.85	0.75
	y	0.7	1	0.6	0.5	0.85
	z	0.85	0.6	1	*	*
BU_2	x	0.85	0.5	*	1	0.75
	y	0.75	0.85	*	0.75	1

Maximal Determinant Completions

- **Unique** if completions exist.
- **Maximum entropy completion** for the multivariate normal model, where maximum entropy is a principle of favouring the simplest explanations.
- **Maximum likelihood estimate** of the correlation matrix of the unknown underlying multivariate normal model.
- **Analytic centre** of the feasible region described by the pos semi defness constraints.
- Has **zeros in inverse** corresponding to the free elements of A .

A Particular Case

Theorem (Georgescu, H & Peters, 2018)

Consider the symmetric matrix

$$\begin{matrix} & n_1 & n_2 & n_3 \\ n_1 & \left[\begin{array}{ccc} A_{11} & B & C \\ B^T & A_{22} & E \\ C^T & E^T & A_{33} \end{array} \right] & & \end{matrix},$$

where E is unspecified, all the diagonal blocks are pos def, and all specified principal minors are positive. The maximal determinant completion is $E = B^T A_{11}^{-1} C$.

Example

$$A = \left[\begin{array}{cc|cc|cc} 1 & \frac{3}{5} & \frac{3}{5} & \frac{3}{5} & \frac{3}{5} & \\ \hline \frac{3}{5} & 1 & \frac{3}{5} & ? & ? & \\ \frac{3}{5} & \frac{3}{5} & 1 & ? & ? & \\ \hline \frac{3}{5} & ? & ? & 1 & \frac{3}{5} & \\ \frac{3}{5} & ? & ? & \frac{3}{5} & 1 & \end{array} \right] \Rightarrow C = \left[\begin{array}{cc|cc|cc} 1 & \frac{3}{5} & \frac{3}{5} & \frac{3}{5} & \frac{3}{5} & \\ \hline \frac{3}{5} & 1 & \frac{3}{5} & \frac{9}{25} & \frac{9}{25} & \\ \frac{3}{5} & \frac{3}{5} & 1 & \frac{9}{25} & \frac{9}{25} & \\ \hline \frac{3}{5} & \frac{9}{25} & \frac{9}{25} & 1 & \frac{3}{5} & \\ \frac{3}{5} & \frac{9}{25} & \frac{9}{25} & \frac{3}{5} & 1 & \end{array} \right].$$

$$\det(A(3/5)) = 0.087, \det(C) = 0.124.$$

$$C^{-1} = \left[\begin{array}{cc|cc|cc} \frac{29}{11} & -\frac{15}{22} & -\frac{15}{22} & -\frac{15}{22} & -\frac{15}{22} & \\ \hline -\frac{15}{22} & \frac{20}{11} & -\frac{15}{22} & 0 & 0 & \\ -\frac{15}{22} & -\frac{15}{22} & \frac{20}{11} & 0 & 0 & \\ \hline -\frac{15}{22} & 0 & 0 & \frac{20}{11} & -\frac{15}{22} & \\ -\frac{15}{22} & 0 & 0 & -\frac{15}{22} & \frac{20}{11} & \end{array} \right].$$

B
In-depth
Modelling

B
In-depth
Modelling

COMPLETING THE MATRIX

Correlation matrices arise in many applications to model the dependence between variables.

Dan Georgescu and **Nick Higham** ask what happens when we have a partially specified matrix and we wish to fill in the missing elements

In linear algebra terms, a correlation matrix is a symmetric positive semidefinite (PSD) matrix with unit diagonal. In other words, it is a symmetric matrix with ones on the diagonal whose eigenvalues are all non-negative. Eigenvalues may seem to be an unnecessary complication, but they determine whether or not a given matrix with ones on the diagonal is a correlation matrix. For example,

$$A = \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{bmatrix}$$

is not a correlation matrix: it has eigenvalues -0.414 , 1 , 2.414 . Since a correlation matrix is a scaled covariance matrix, a negative eigenvalue would suggest that one of the variables has a negative variance, which cannot happen. A matrix having positive eigenvalues is the matrix equivalent of a real number being non-negative. We also need our correlation matrices to have this property because capital models reasonably expect inputs of positive variances and simulate possible future states of the world by first calculating the square root of the correlation matrix.

Why we need correlation matrix completion
In practice, there is often incomplete or missing information for the variables and this may lead to missing values in the correlation matrix itself, hence the problem of how to complete the matrix. We show that some of these practical problems can be solved explicitly via simple formulae, and explain how to use mathematical tools to solve the more general problem where explicit solutions may not exist. (Simple¹ is, of course, a relative term.)

As a simple example, consider the 3-by-3 matrix above and suppose the entries in the top-right and bottom-left corners are omitted and those elements have to be chosen. Can we find a value for these entries that produces a valid correlation matrix? The only such value is 1, which can be seen from the requirement that a correlation matrix must have a non-negative determinant (the determinant being the product of the eigenvalues).

Correlation matrices are used in the aggregation of risk exposures required by regulations. The values of the correlations really matter because, for many insurers, aggregate diversification effects will be the most significant determinant of required capital, with a 40%-60% reduction in the capital required between risk types not being atypical for a large, well-diversified insurer or reinsurer.

Some of the correlations are known because they have been estimated from data, prescribed by regulations or assigned by expert judgment, but the other entries are not

FIGURE 3 Partially specified correlation matrix

Correlations		a	b	c	d	e
BU ₁	a	1	0.7	0.65	0.65	0.75
	b	0	1	0.6	0.6	0.65
BU ₂	a	0.65	0.6	1	0.65	0.65
	b	0.65	0.6	0.65	1	0.75
BU ₃	a	0.65	0.65	0.65	0.75	1
	b	0.65	0.65	0.75	1	

“There is often incomplete or missing information for the variables and this may lead to missing values in the correlation matrix itself, hence the problem of how to complete the matrix”

known. This could be because there is little reliable data, no specific regulations and few available experts with knowledge of both of the risks being correlated. The aim is to complete the missing entries in order to produce a valid correlation matrix to calculate a capital requirement as realistically as is possible given the absence of hard evidence.

In general there are many possible completions (our simple 3-by-3 example is unusual in having only one solution), and the choice of which to use is important because correlations determine the diversification between capital for the various risk variables.

For an example, consider the following simplified problem, which is included to illustrate the issue but

is not supposed to represent a realistic calibration. A firm is exposed to three risks, a , b and c , and is organised in two business units, BU_1 and BU_2 . Risk a represents corporate bond spread in BU_1 and BU_2 . Risk b reflects, say, equity risk. Risk c in BU_1 is an exposure to an asset class that is not traded and therefore has no market data history from which to calibrate a correlation. All the correlation coefficients are known in BU_1 and BU_2 (the upper left and bottom right of the matrix) but not between risk c in BU_1 and the risks in BU_2 . This could be because the firm operates in two markets (say the UK and Italy, corresponding to BU_1 and BU_2) and UK experts were able to advise on the correlations between risk c and risks a and b by making an expert judgment. However, no relevant judgment could be made on the correlation between risk c in BU_1 and the risks in BU_2 . Typically, missing correlations arise between risks in different business units because there are few experts who understand both businesses and can make these judgments.

For the partially specified matrix given in Figure 3, a valid correlation matrix completion must lie in the dark yellow region in Figure 4. The centre of this region is the maximum determinant completion, where x is 0.7 and y is 0.64, to two decimal places. In that sense, the maximum determinant completion is unbiased. To give an easy-to-interpret number, if the capital held for the risks in BU_1

Open Problem (1)

Rounding a Correlation Matrix

How do we round a correlation matrix to p decimal places in such a way that the rounded matrix is a correlation matrix?

Open Problem (2)

```
>> A = gallery('randcorr', [0,1 2])
```

```
A =
```

```
1.0000e+00    3.2613e-16   -2.6063e-01  
3.2613e-16    1.0000e+00   -9.6544e-01  
-2.6063e-01  -9.6544e-01    1.0000e+00
```

```
>> eig(mp(A))
```

```
ans =
```

```
2.523010255943400625391528887408137e-16  
0.9999999999999999835879781129131914  
1.9999999999999999911819193276528023
```

```
>> eig(mp(single(A)))
```

```
ans =
```




```
-1.67147076182867535347476081552851e-08  
0.9999999999999999835879779860717258  
2.00000001671470778240697367403035
```

Final Remarks

- Correlation matrices are a rich source of numerical linear algebra problems.
- Excellent methods available for **nearest correlation matrix** and **correlation matrix completion**.
- Numerical reliability and efficiency is essential.
- Partnership with **NAG** has enabled rapid inclusion of our algorithms in the NAG Library.
- More info on my blog:
<https://nhigham.com/tag/correlations/>.

Slides available at https://bit.ly/rehab_corr

References I

-  R. Borsdorf, N. J. Higham, and M. Raydan.
Computing a nearest correlation matrix with factor structure.
SIAM J. Matrix Anal. Appl., 31(5):2603–2622, 2010.
-  P. I. Davies and N. J. Higham.
Numerically stable generation of correlation matrices and their factors.
BIT, 40(4):640–651, 2000.
-  D. I. Georgescu, N. J. Higham, and G. W. Peters.
Explicit solutions to correlation matrix completion problems, with an application to risk management and insurance.
Roy. Soc. Open Sci., 5(3):1–11, 2018.

References II



N. J. Higham.

Accuracy and Stability of Numerical Algorithms.

Society for Industrial and Applied Mathematics,
Philadelphia, PA, USA, second edition, 2002.

ISBN 0-89871-521-0.

xxx+680 pp.



N. J. Higham.

Computing the nearest correlation matrix—A problem
from finance.

IMA J. Numer. Anal., 22(3):329–343, 2002.

References III



H. Qi and D. Sun.

A quadratically convergent Newton method for computing the nearest correlation matrix.

SIAM J. Matrix Anal. Appl., 28(2):360–385, 2006.