# PROBLEM:
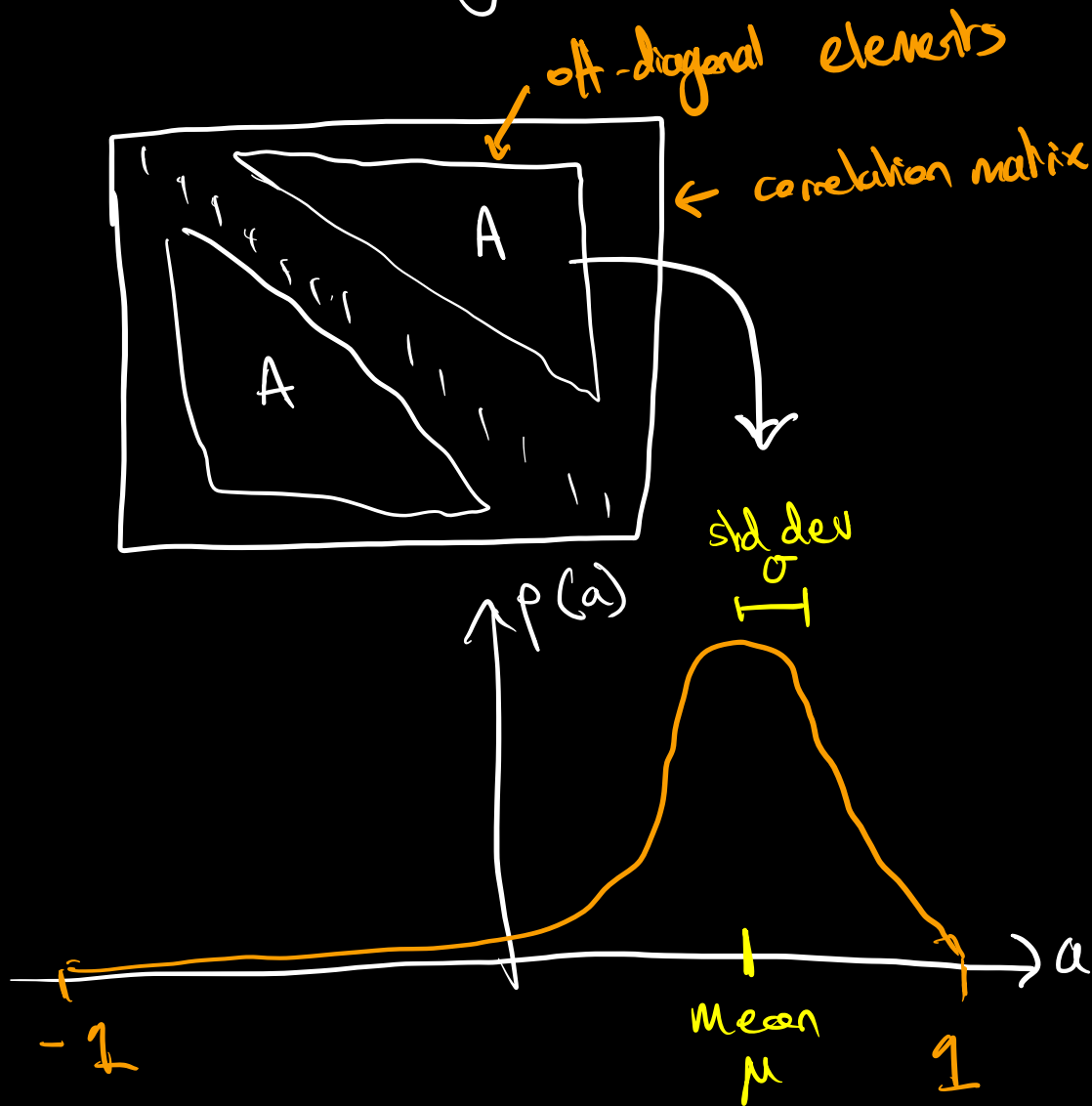
Generate random correlation matrices with particular off-diagonal distributions.

off-diagonal elements

$A$

$A$

← correlation matrix

std dev $\sigma$

$p(a)$

$a$

$-1$

Mean $\mu$

$1$

Generate correlation matrices whose off-diagonal elements have:

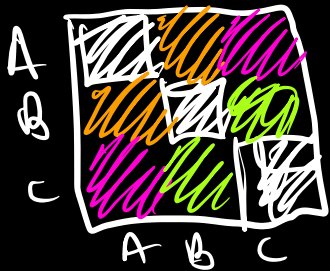* mean $\mu$
* std dev $\sigma$

# Why?

Pops up occasionally, e.g. in my neuroscience research!

# Why is it non-trivial?

You can't just sample random matrices, not all matrices are correlation matrices!

THERE IS DEPENDENCE BETWEEN THE ELEMENTS



If A + B are correlated 〰 high

AND B + C are correlated 〰 high

Then A + C MUST be quite correlated

〰 has to be large, can't just randomly sample

To understand this from another perspective, and for later use . . . .

# Geometric View on Problem:

Create a dataset, N datapoints, put them in matrix $\underline{\underline{X}}$

$$\underline{\underline{X}} = \begin{pmatrix} & & \\ & & \\ & & \end{pmatrix} \updownarrow \text{Features}$$

$\leftarrow$ Datapoints $\rightarrow$
$d_1, d_2 \ldots\ldots d_N$

$\underleftarrow{\phantom{x}}$ feature vector

$\begin{matrix} a \\ b \\ \vdots \\ z \end{matrix}$

$\Rightarrow$ create mean. centered features

$\tilde{a} = a - \mu_a \underline{1}$
$\tilde{b} = b - \mu_b \underline{1}$
$\vdots$
$\tilde{z} = z - \mu_z \underline{1}$

Feature vectors are points in some N-dimensional space



$$\sigma_a^2 = \mathbb{E}\left((a - \mu_a)^2\right)$$
$$= \frac{1}{N} \sum_i \tilde{a}_i^2 = \frac{1}{N} |\tilde{a}|^2$$

Length of the vector $\tilde{\underline{a}} = |\tilde{\underline{a}}|$ gives variance

Then correlation between two features is the angle!

$$\rho_{ab} = \frac{\mathbb{E}\left[(a-\mu_a)(b-\mu_b)\right]}{\sqrt{\mathbb{E}\left[(a-\mu_a)^2\right]\mathbb{E}\left[(b-\mu_b)^2\right]}}$$

$$= \frac{\tilde{a} \cdot \tilde{b}}{\sqrt{|\tilde{a}|^2 \, |\tilde{b}|^2}} = \cos\theta_{ab}$$

CONSTRAINT MAKES SENSE GEOMETRICALLY



$\therefore$ small $\theta_{ab}$ and small $\theta_{bc}$
$\Rightarrow$ quite small $\theta_{ac}$

# 4 Approaches to Problem

1) Onions + partial correlations

   Lewandowski, Kurowicka, & Joe 2009

   + dude on internet (amoeba)

2) Via random vectors

   me! Though discussed in Hardin, Garcia, Golan et al.
   2013 in "A Method for generating realistic correlation matrices"

3) Via neural networks

   Gautier Marti, 2019

4) Factor loadings

   dude on internet

(New 2022! Via reparameterisation, Archakov et al.
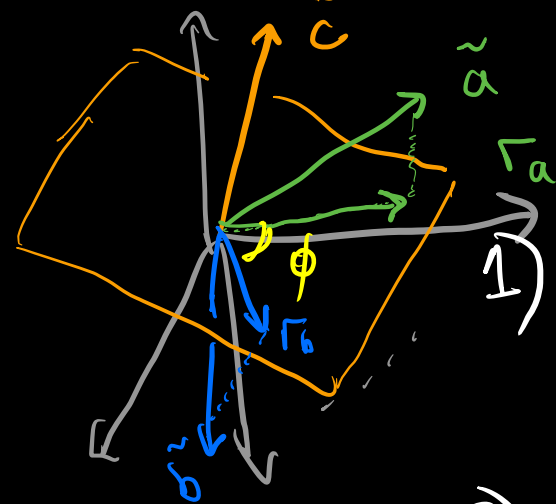   2022, A new method for generating random correlation matrices)

(New 2023! Via noise addition. HARDIN et al. 2013

   A method for generating realistic correlation matrices)

# Sampling Uniformly from set of correlation Matrices

Lewandowski, Kurowicka, & Joe, 2009  VINE METHOD

key beginning: partial correlations are independent!

Partial correlation =   correlation between

residuals of a & b after regressing with c



Geometrically

1) Projection into plane ⊥ to c

⇓

residual vector after linear regression with c

2) Angle $\phi$ of two residual vectors

= partial correlation

# ∃ a recursive formula for partial correlations

$$\rho_{xy \cdot \mathbf{z}} = \frac{\rho_{xy \cdot \mathbf{z} \backslash z_0} - \rho_{xz_0 \cdot \mathbf{z} \backslash z_0} \, \rho_{z_0 y \cdot \mathbf{z} \backslash z_0}}{\sqrt{1 - \rho^2_{xz_0 \cdot \mathbf{z} \backslash \{z_0\}}} \, \sqrt{1 - \rho^2_{z_0 y \cdot \mathbf{z} \backslash \{z_0\}}}}$$

∴ Can relate complete partial corr $\rho_{xy \cdot \mathbf{z}}$

to incomplete $\rho_{A \cdot \mathbf{z} \backslash z_0}$   $A = \begin{matrix} xy \\ xz_0 \\ yz_0 \end{matrix}$

↦ Vine = organisation of computations

Can do a recursive computation $\mathcal{O}(n^3)$ to get

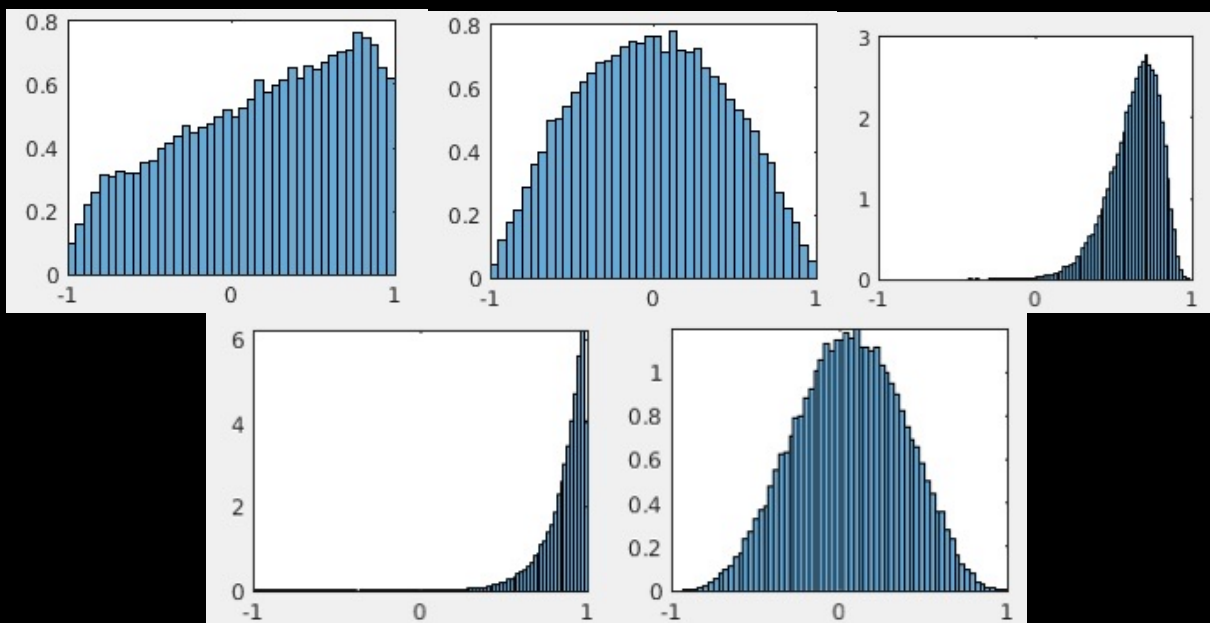from $\frac{n(n-1)}{2}$ partials to $\frac{n(n-1)}{2}$ corr s.


Paper shows how to sample partial correlations

in order to sample uniformly from correlation

matrices. OR with prob $\alpha \, |\underline{\underline{C}}|^{\eta}$

Alteration: men on internet 2014

$$\ell_{partial} \sim Beta(\alpha, \beta)$$

Vary $\alpha$ & $\beta$ manually to get desired distribution.

Some random $\alpha$ + $\beta$ values:



↑ effect of $\ell$

→ $\ell$

PRO: easy to get all kinds of distributions

CON: slow recursive formula for large dimensionality
not arbitrarily complex distributions
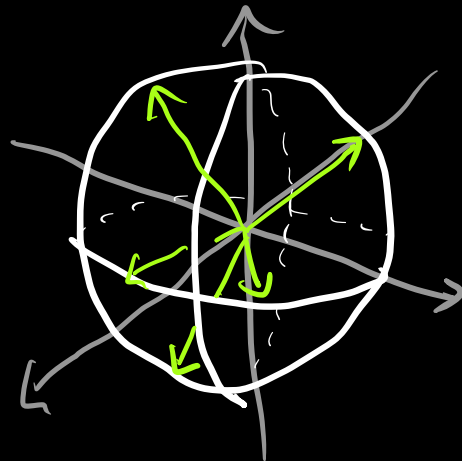probably hard to analyse

# METHOD 2 : via random vectors

## 1) Generate points on sphere

- First generate

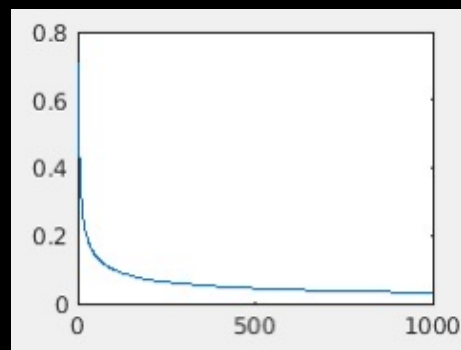$$\underline{x} \sim \mathcal{N}(\underline{0}, \underline{\underline{1}})$$

- Then normalize

$$\underline{x} = \frac{x}{|x|}$$

## 2) Get dot product matrix
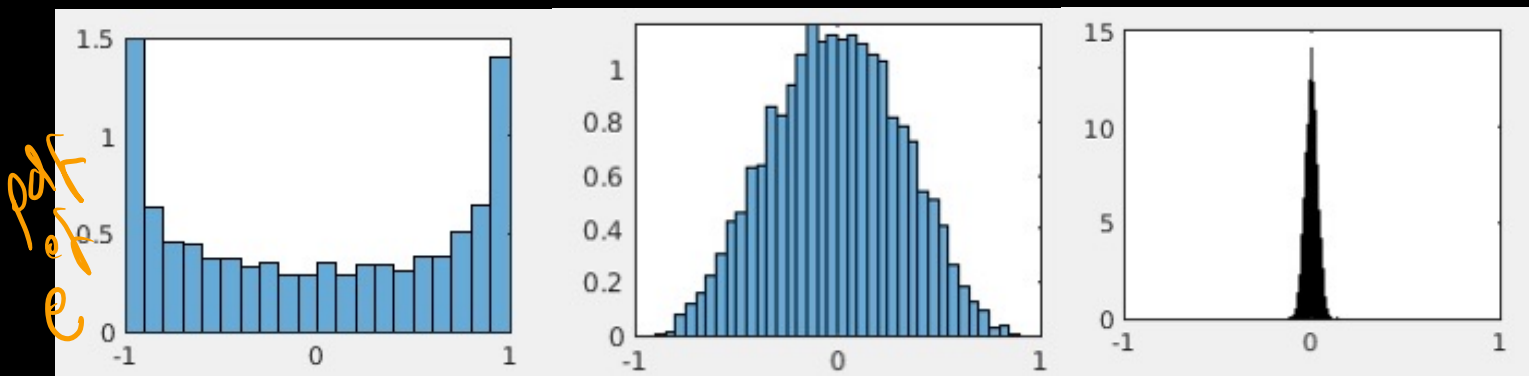
How to vary <u>dot product matrix?</u>  $\sigma(e)$

$\sigma$ : vary dimensionality $\uparrow D$

↑ variance of dot products



D

| D = 2 | D = 10 | D = 1000 |

pdf
of
e



e         e         e
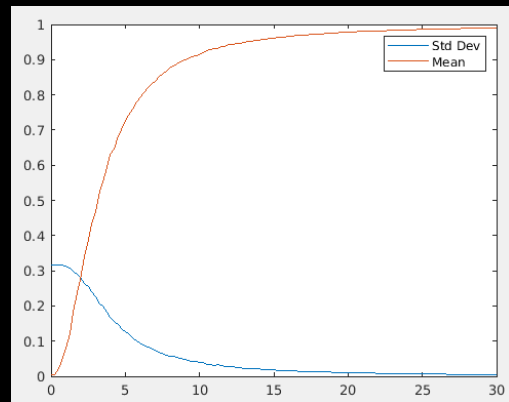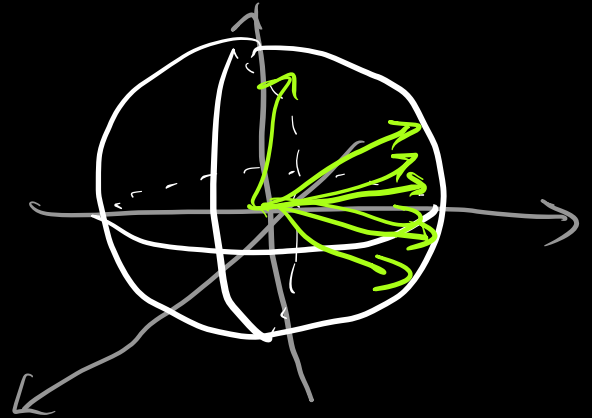
elements of dot product matrix
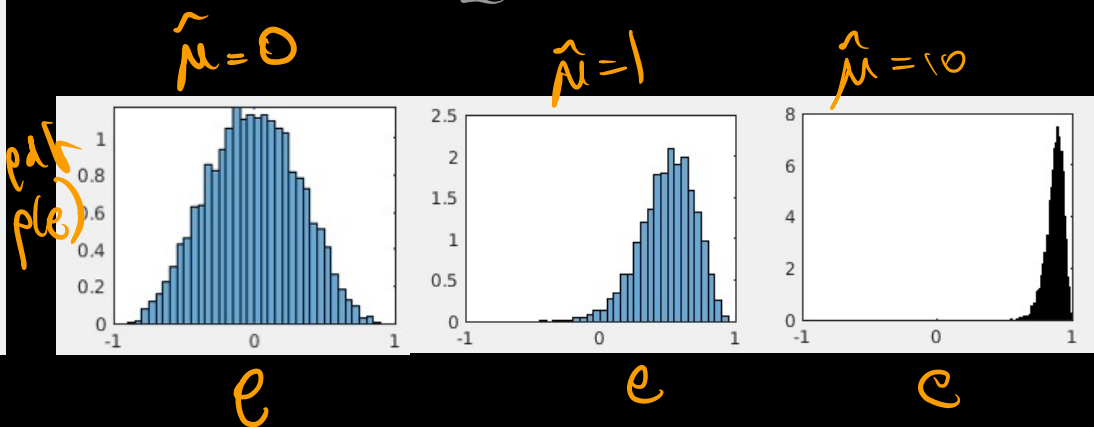
$\mu$ : shift centre

- First generate

$$\underline{x} \sim \mathcal{N}(\underline{0}, \underline{\underline{1}}) + \hat{\mu} \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

- Then normalize

$$\underline{x} = \frac{\underline{x}}{|\underline{x}|}$$





$\hat{\mu} = 0$     $\hat{\mu} = 1$     $\hat{\mu} = 10$

$p(e)$     $e$     $e$     $c$

$\hat{\mu}$

PRO :  very easy + quick

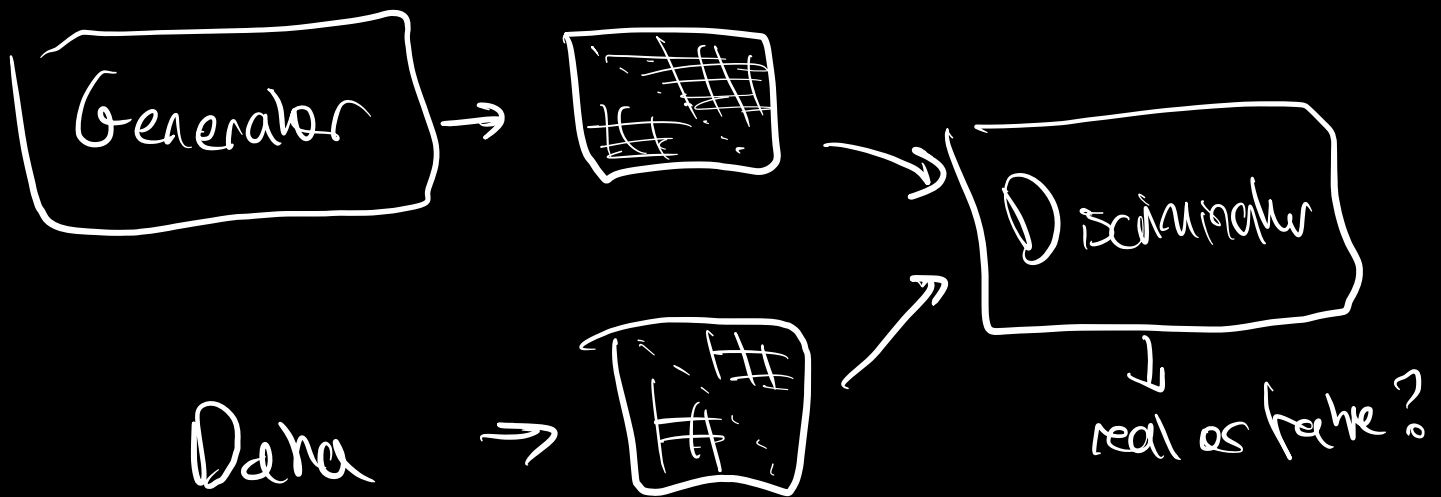probably can be analysed ? $\mu = f(\mu, 0)$ ?

CON :  have no manually tune

not arbitrarily complex

# METHOD 3

<span style="color:blue">CORRGAN , Gautier Marti, 2019</span>

Match measured correlation matrices

(e.g. from finance)

Generator $\rightarrow$ [grid] $\rightarrow$ Discriminator

Data $\rightarrow$ [grid] $\nearrow$

Discriminator $\downarrow$ real or fake?

PRO: Not just mean + variance, any structure in distribution

CON: you have to train a GAN !!

# METHOD 4: Factor loadings

$$\underline{\underline{W}} \in \mathbb{R}^{k \times d}$$

$$\underline{\underline{B}} = \underline{\underline{W}} \, \underline{\underline{W}}^T + \underline{\underline{D}}$$

<span style="color:orange">← diagonal matrix</span>

tve-definite covariance matrix

$$\underline{\underline{C}} = \underline{\underline{E}}^{-1/2} \, \underline{\underline{B}} \, \underline{\underline{E}}^{1/2}$$

correlation matrix

<span style="color:orange">$\underline{\underline{E}}$ diagonal matrix with same diagonal as $\underline{\underline{B}}$</span>

Large $k$ : random matrix, low off-diag corr

Small $k$ : v. high off-diag corr


PRO : super easy + quick
     probably analysable

CON : Only one degree of freedom (could inhedtee more)