# Covariance Matrix Cleaning
## Predictability

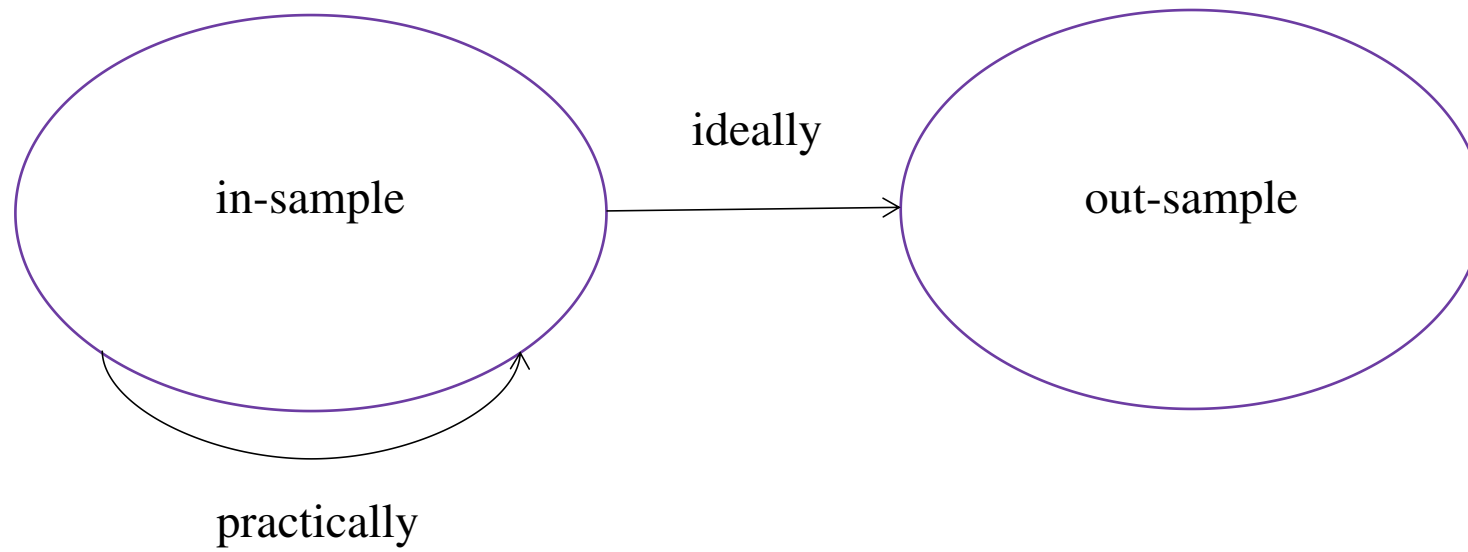Christian Bongiorno | Allocation de Portefeuille | 19-01-2023
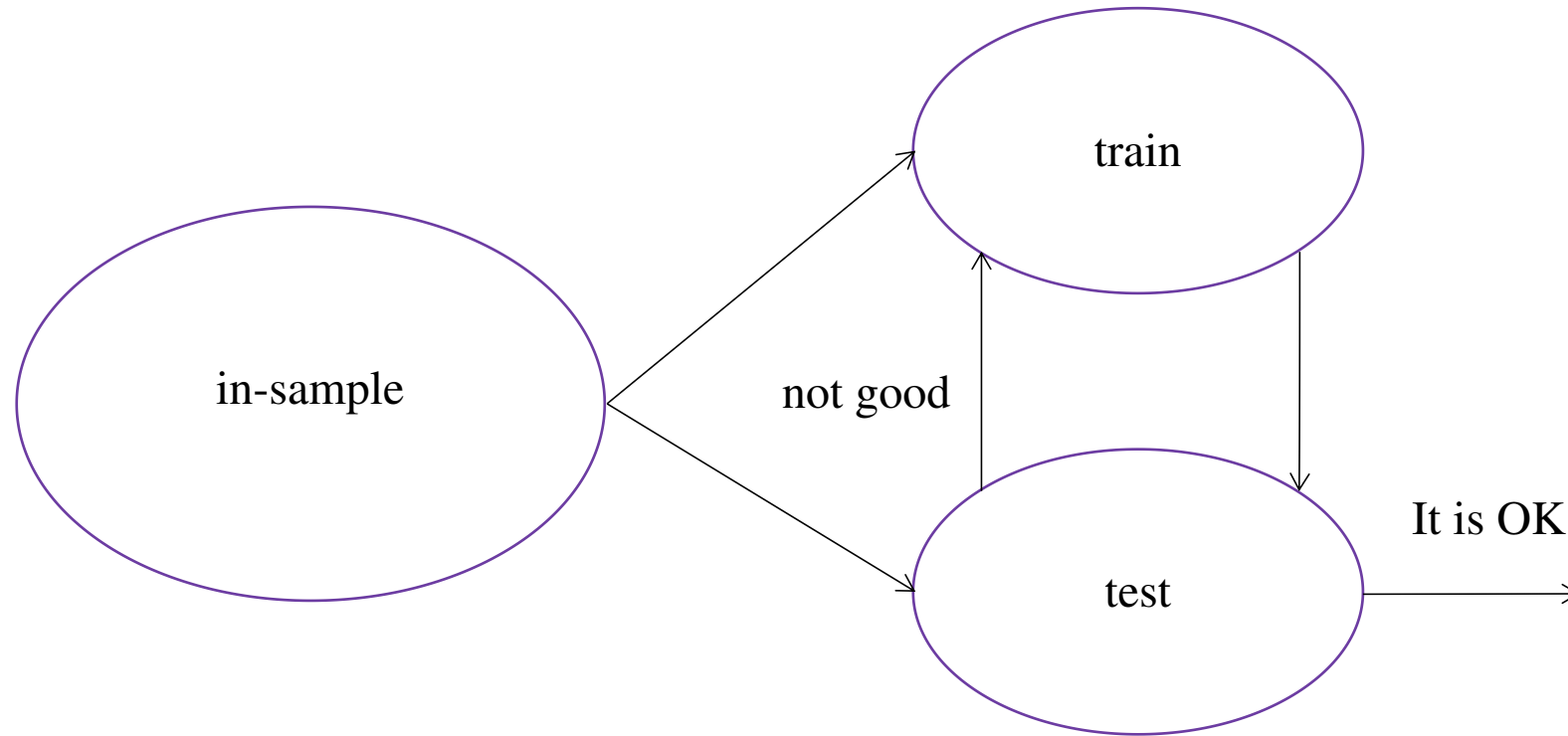
# What is Predictability?

# What is Predictability?

# Cross-Validation

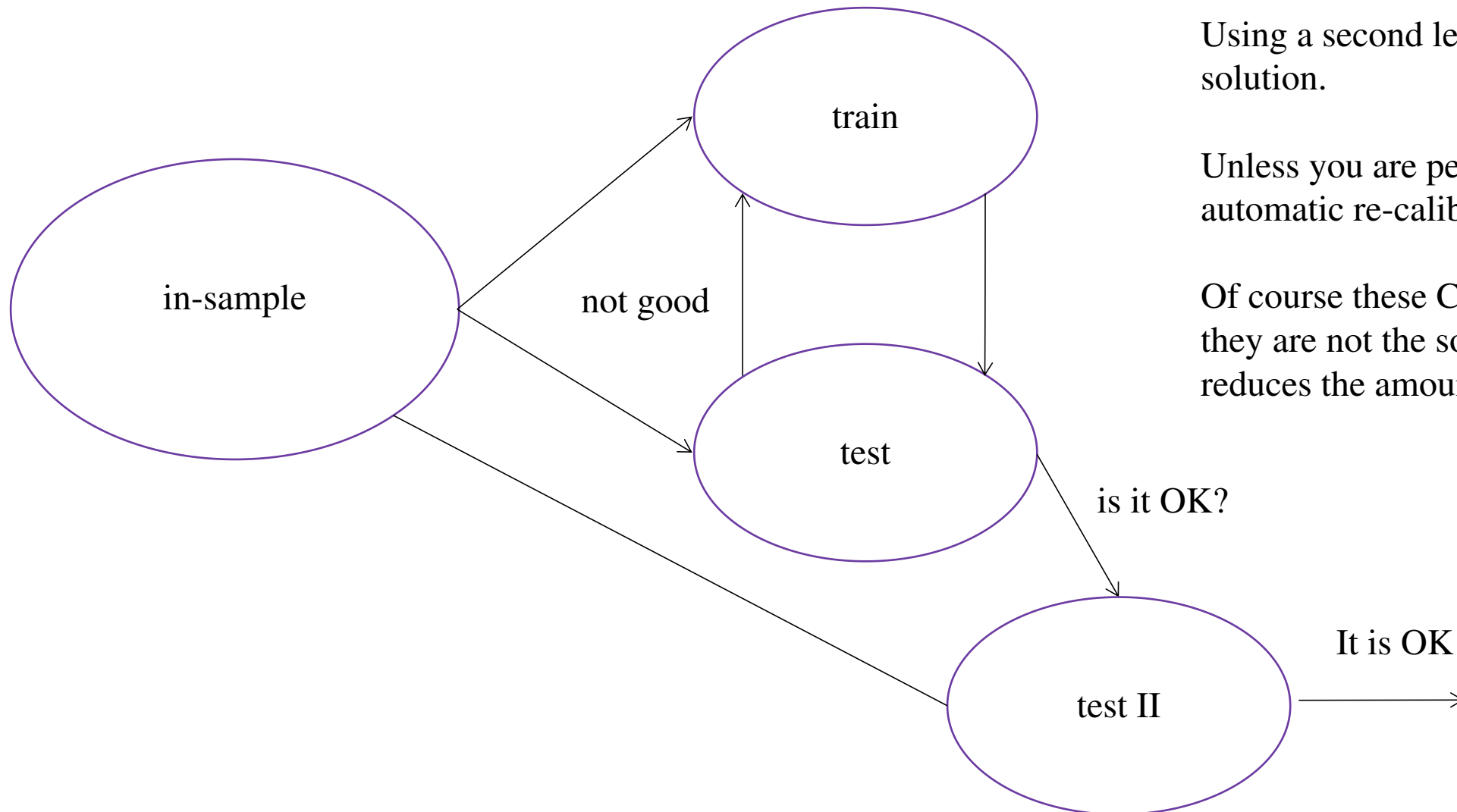# Over-fitting in Cross-Validation



If you perform many train-test adjustments you might fall again in a kind of over-fitting

# Over-fitting in Cross-Validation [Solutions?]



Using a second level test set might be the solution.

Unless you are performing again many automatic re-calibrations.

Of course these CV approaches, even if they are not the solution, for sure they reduces the amount of over-fitting.

# Cross-Validation in Finance

*We must use historical observations (in-sample), to predict to predict the next future (out-of-sample).*

Not Anchored Walk-Forward Test Procedure

1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008

**In-sample data (IS)** - Optimization

**Out-of-sample data (OOS)** - Verification (backtest)

# Do Not Mix Past and Future

## Predicting financial markets with Google Trends and not so random keywords

Damien Challet[*]

*Chaire de finance quantitative, Laboratoire de mathématiques appliquées aux systèmes,
École Centrale Paris, Grande Voie des Vignes, 92295 Châtenay-Malabry, France and
Encelade Capital SA, Parc Scientifique C, EPFL, 1015 Lausanne, Switzerland*

Ahmed Bel Hadj Ayed[†]

*Chaire de finance quantitative, Laboratoire de mathématiques appliquées aux systèmes,
École Centrale Paris, Grande Voie des Vignes, 92295 Châtenay-Malabry, France*

We discuss the claims that data from Google Trends contain enough information to predict future financial index returns. We first review the many subtle (and less subtle) biases that may affect the backtest of a trading strategy, particularly when based on such data. Expectedly, the choice of keywords is crucial: by using an industry-grade backtest system, we verify that random finance-related keywords do not to contain more exploitable predictive information than random keywords related to illnesses, classic cars and arcade games. However, other keywords applied on suitable assets yield robustly profitable strategies, thereby confirming the intuition of [24].

# Out-of-Sample Instability

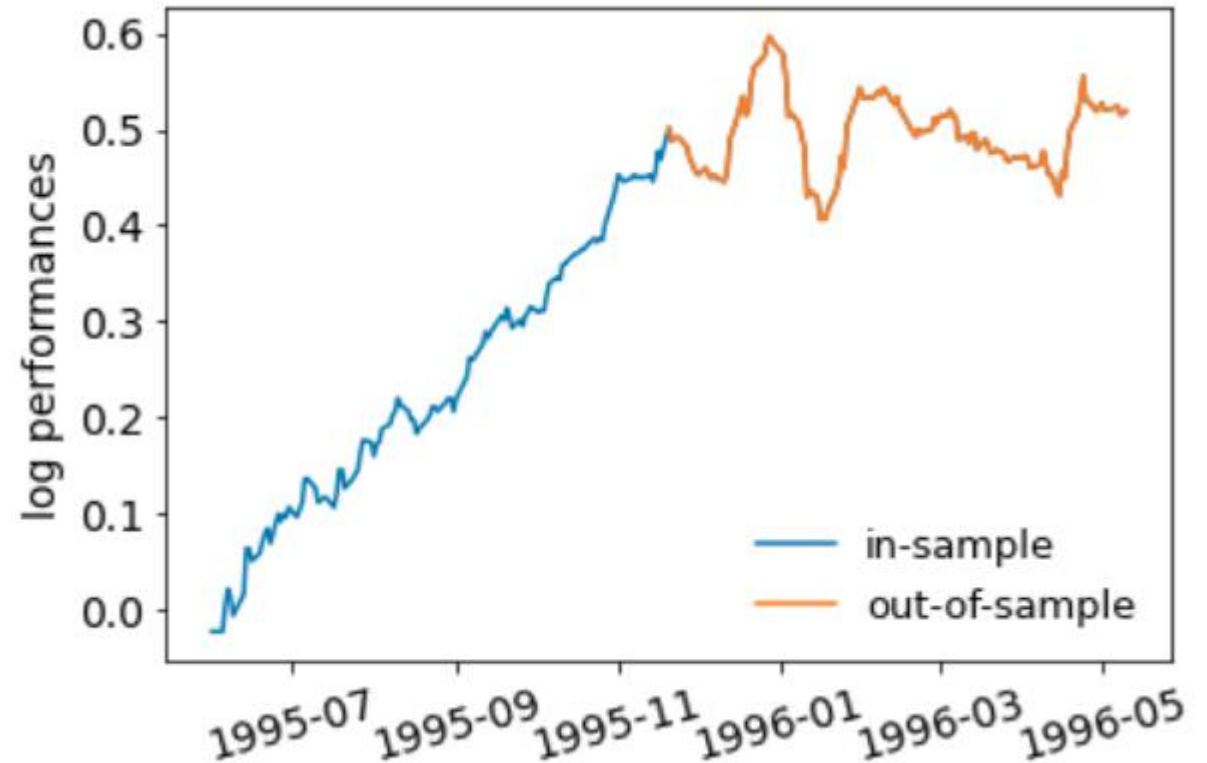*Historical values often are not persistent.*

Example 1:

**Sharpe-Ratio**

$$SR = \frac{signal}{noise} = \frac{return}{risk}$$

in-sample SR:      5.56

out-of-sample SR:            0.47

# Out-of-Sample Instability

*Example 2:*

*Correlation matrices:*

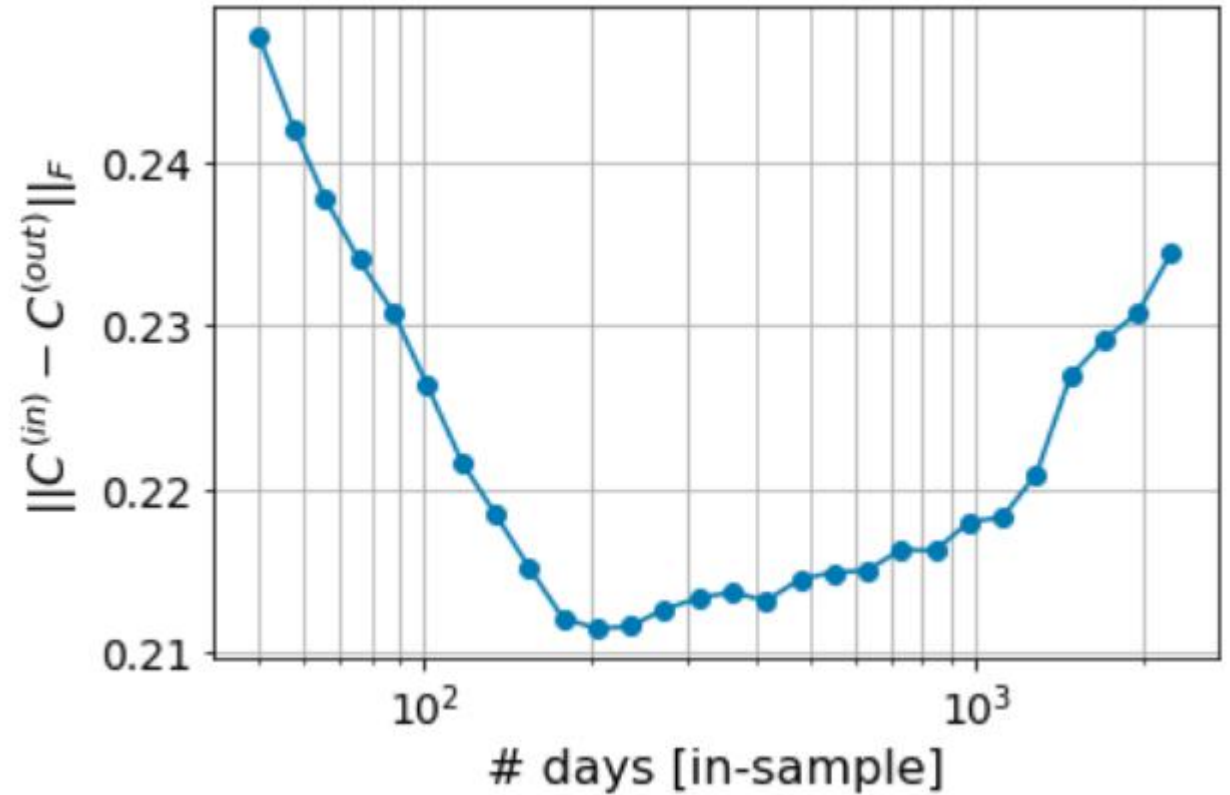$$\|C^{(\text{in})} - C^{(\text{out})}\|_F = \frac{\sqrt{\sum_{ij}\left(C_{ij}^{(\text{in})} - C_{ij}^{(\text{out})}\right)^2}}{n}$$

Too short time-series are too noisy

Too long time series messing things up for the non-stationary of the system.

**SOLUTIONS:**

1. Do not trust so much in the historical data, i,e, use robust control techniques;

2. Clean the data! [I'll tell you after what does it means]

# Sources of Instability

**1. Sample size error**:

Let us assume that you are playing with a coin that makes head 40%, but on the first set of tosses, you observed 7 out 10 head. If you do not assume sample size error you will lose everything

**2. The future is different from the past.**

# Portfolio Optimization [in-sample]



As we said before, working in-sample is easy. Let us focus on the max SR ptf.

# Portfolio Optimization [out-of-sample]



The efficient frontier in the out-of-sample shrinks and shifts to the right.

Note: This is a "lucky" case, sometimes it is y-reflexed [return inversion].

# Minimum Variance Portfolio

*Agnostic with respect to the returns*

$$w = \frac{\Sigma^{-1}\mathbb{1}}{\mathbb{1}'\,\Sigma^{-1}\mathbb{1}}$$

We reduce the information, but also the source of noise.

What do you expect?

# Minimum Variance Portfolio



out-of-sample

The Min Var portfolio has a higher SR than the Max SR portfolio.

This of course does not mean that the returns should not be accounted; however, they must be treated carefully.

In the following part we will focus only on improving the estimation of the covariance matrix.

# Spectral Decomposition of the Portfolios

*We use the Min-Var ptf to focus on the role of the covariance*

$$w = \frac{\Sigma^{-1}\mathbb{1}}{\mathbb{1}'\,\Sigma^{-1}\mathbb{1}}$$

Sum by row of the inverted covariance matrix

Sum over all elements of the inverted covariance matrix

*Spectral decomposition:*
$$\Sigma = \lambda_1 v_1 v_1' + \lambda_2 v_2 v_2' + \lambda_3 v_3 v_3' + \ldots + \lambda_n v_n v_n'$$

*Inverted matrix*
$$\Sigma^{-1} = \frac{1}{\lambda_1} v_1 v_1' + \frac{1}{\lambda_2} v_2 v_2' + \frac{1}{\lambda_3} v_3 v_3' + \ldots + \frac{1}{\lambda_n} v_n v_n'$$

# A brief note on Principal Component analysis

A first objective of principal component analysis is to determine the standardized linear combination of the original variables which has maximal variance.

Principal component analysis looks for a few linear combinations which can be used to summarize data, losing in the process as little information as possible.

# Principal Component Analysis

$$\Sigma = X\,X' = V\,\Lambda\,V' = V\,\Lambda^{\frac{1}{2}}\,\Lambda^{\frac{1}{2}}V'$$

We can define $P := V\,\Lambda^{\frac{1}{2}}$ , therefore

$$\Sigma = P\,P'$$

We can imagine the existence of a set of orthogonal axis
$$A\,A' = \mathbb{1}$$

$$\Sigma = P\,P' = P\,A\,A'P' \qquad X = P\,A$$

# The Spectrum of a Financial Correlation Matrix

A simple example with 10 stocks (2001-2003)

|       | AIG | IBM   | BAC   | AXP   | MER   | TXN   | SLB   | MOT   | RD    | OXY   |
|-------|-----|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| **AIG** | 1   | 0.413 | 0.518 | 0.543 | 0.529 | 0.341 | 0.271 | 0.231 | 0.412 | 0.294 |
| IBM   |     | 1     | 0.471 | 0.537 | 0.617 | 0.552 | 0.298 | 0.475 | 0.373 | 0.270 |
| BAC   |     |       | 1     | 0.547 | 0.591 | 0.400 | 0.258 | 0.349 | 0.370 | 0.276 |
| AXP   |     |       |       | 1     | 0.664 | 0.422 | 0.347 | 0.351 | 0.414 | 0.269 |
| MER   |     |       |       |       | 1     | 0.533 | 0.344 | 0.462 | 0.440 | 0.318 |
| TXN   |     |       |       |       |       | 1     | 0.305 | 0.582 | 0.355 | 0.245 |
| SLB   |     |       |       |       |       |       | 1     | 0.193 | 0.533 | 0.592 |
| MOT   |     |       |       |       |       |       |       | 1     | 0.258 | 0.166 |
| RD    |     |       |       |       |       |       |       |       | 1     | 0.590 |
| OXY   |     |       |       |       |       |       |       |       |       | 1     |

$$V' \, \Sigma \, V \; = \; \Lambda$$

The eingevalue spectrum of this correlation matrix is

| $\lambda_1$ | $\lambda_2$ | $\lambda_3$ | $\lambda_4$ | $\lambda_5$ | $\lambda_6$ | $\lambda_7$ | $\lambda_8$ | $\lambda_9$ | $\lambda_{10}$ |
|------|------|------|------|------|------|------|------|------|------|
| 4.72 | 1.40 | 0.96 | 0.52 | 0.48 | 0.45 | 0.42 | 0.39 | 0.35 | 0.30 |

# Eigenvalue Distribution of a Random Noise

Let consider $X \sim N(0, \mathbb{1}) \in 100 \times 300$



In this case, only using the true eigenvalues you can obtain the true population covariance

$$V^{(s)} \, \mathbb{1} \, V^{(s)\prime} = \mathbb{1} = \Sigma_{\text{true}}$$

In general, the sample covariance will be different from the population one

$$V^{(s)} \, \Lambda^{(s)} \, V^{(s)\prime} \neq \mathbb{1} = \Sigma_{\text{true}}$$

# Eigenvalue Distribution in Equity Markets



*The eigenvalue distribution of the covariance of uncorrelated random variable is call* Marchenko-Pastur distribution:

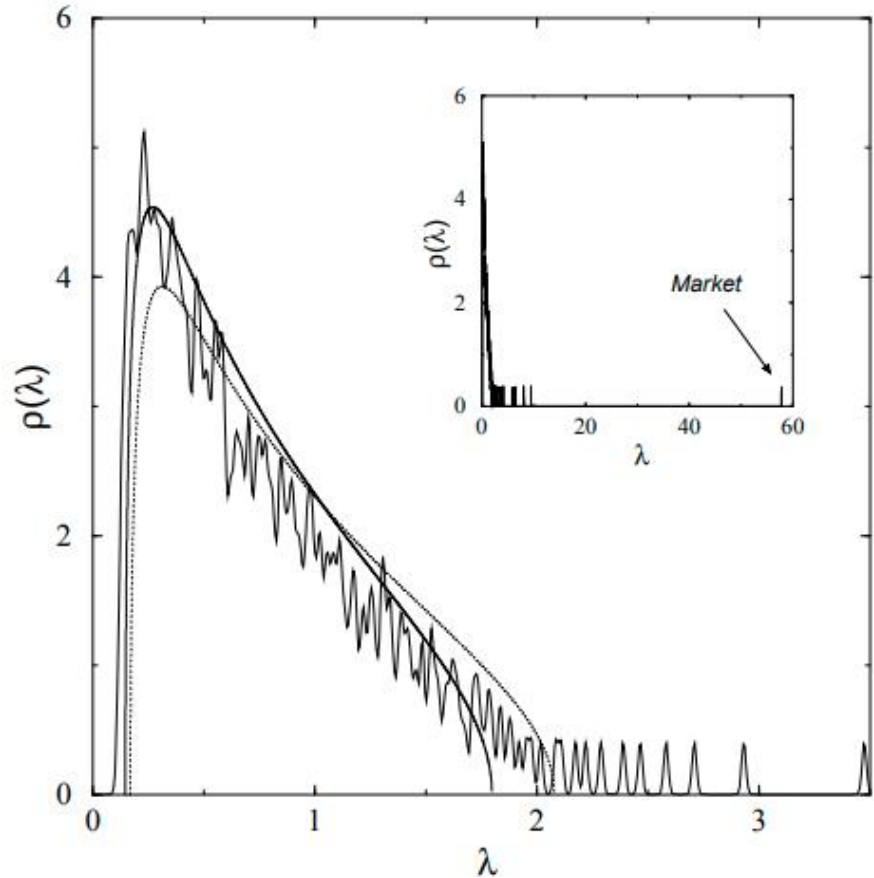$$\rho_c(\lambda) = \frac{T}{2\pi\sigma^2 N \lambda}\sqrt{(\lambda_{max} - \lambda)(\lambda - \lambda_{min})}$$

$$\lambda_{min}^{max} = \sigma^2\left(1 + \frac{N}{T} \pm 2\sqrt{\frac{N}{T}}\right)$$

In the "thermodynamic" limit of $N, T \to \infty$ with $T/N \geq 1$, the Marchenko-Pastur distribution if bounded between $(\lambda_{min}, \lambda_{max})$.

$\sigma^2$ is the variance of the elements, that is equal to one for correlation matrices.

Note the presence of outliers. We are sure that out-lier contains non-random information. Eigencomponent within the bulk of the distribution, probably contains information but they are "noise dressed"

Laloux, L., Cizeau, P., Bouchaud, J. P., & Potters, M. (1999). Noise dressing of financial correlation matrices. Physical review letters, 83(7), 1467.

# Discussion about the Eigenvalue Distribution



N=100

The correlation matrix of uncorrelated random variable has identically one eigenvalues.

We know that equities correlation matrices are dominated by the mode, therefore it is possible to discount it by setting $\sigma^2 = 1 - \lambda_1/N$

In summary, on long time-horizon, i.e. $T \gg N$ you can trust on the the eigenvalue distribution; however, due to the non-stationary you cannot trust in the correlation itself.

You should work in the regime $T \cong N$, and a correction is needed

The clipping prescribes to substitute the eigenvalue of the bulk with their average value

$$\lambda_i^{(c)} = \begin{cases} \lambda_i^{(s)} & \text{if } \lambda_i^{(s)} > \lambda_{max} \\ \langle \lambda_j^{(s)} \rangle_{\lambda_j^{(s)} < \lambda_{max}} & \text{if } \lambda_i^{(s)} < \lambda_{max} \end{cases}$$

Then, the estimator is obtained as

$$\boldsymbol{C}^{(t)} = \boldsymbol{V}\boldsymbol{\Lambda}^{(c)}\boldsymbol{V}', \quad C_{ij}^{(c)} = \frac{C_{ij}^{(t)}}{C_{ii}^{(t)}C_{jj}^{(t)}}, \qquad \Sigma_{ij}^{(c)} = C_{ij}^{(c)}\sqrt{\Sigma_{ii}\Sigma_{jj}}$$

# Oracle Estimator

Let us define an optimal target for the eigenvalues.

$$\Sigma = V \Lambda V'$$

$$V' \Sigma V = \Lambda$$

The oracle eigenvalues is a cheat. We assume to know the out-of-sample covariance.

The question is to modify only the eigenvalues and fix the eigenvectors

$$(V'_{in} \Sigma_{out} V_{in})_d = O$$

$$\Xi = V_{in} O V'_{in}$$

The oracle eigenvalues are the ones that minimize $\|\Sigma_{out} - \Xi\|_F$

Bun, J., Allez, R., Bouchaud, J. P., & Potters, M. (2016). Rotational invariant estimator for general noisy matrices. IEEE Transactions on Information Theory, 62(12), 7475-7490

# Is the MP Bulk really Random Noise?

I show you a test by randomly selecting N=100 stocks and T=400.

If the clipping is correct we should see something flat.

The Spearman correlation coefficient measures the presence of a monotonic increasing function.



There is a relationship, but probably is not linear.

# Cross-Validated Eigenvalue Shrinkage

The idea of the method is simple: Split the in-sample window into a nested in-sample* and out-of-sample*, but without preserving the temporal order



It is a random sampling procedure, if our split preserves the temporal order we can have only one split, so they propose to mix the temporal order of in-sample* and out-of-sample*, therefore they can address only in-sample sample size error, not the non-stationarity.

# The Cross-Validation Set-up

in-sample day indices

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |

| 1 | 2 | 3 | 5 | 6 | 8 | 10 |     | 4 | 7 | 9 |

in-sample* day indices

out-of-sample* day indices

The in-sample* is a sub-sampling of the in-sample pool, i.e., without replacement. The out-of-sample* is the left out.

The size of the out-of-sample* is not very  non-influential, but must be smaller than in-sample* (10%), and avoid to have too few point (like 5)

Note: No repeated days, and no missing days.

# The Eigenvalue Realizations with CV

The idea is to have an estimate of the out-of-sample variance on the in-sample eigen-directions

$$\left(V'_{in*(i)} \; \Sigma_{out*(i)} \; V_{in*(i)}\right)_d \;=\; \Lambda^{(t)}_{\text{cv}(i)}$$

That is an oracle estimator of the eigenvalues.
you should perform independent sub-sampling and estimate for each of them the oracle eigenvalues

components (N stocks)

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1.826186 | 0.450317 | 0.774918 | 0.277689 | 0.301081 | 0.284515 | 0.262410 | 0.312953 | 0.278350 | 0.132593 |
| 1 | 1.848474 | 0.499851 | 0.467592 | 0.447850 | 0.450473 | 0.386625 | 0.407486 | 0.366422 | 0.276440 | 0.173862 |
| 2 | 1.077437 | 0.432883 | 0.483675 | 0.349912 | 0.278088 | 0.233495 | 0.235383 | 0.179277 | 0.169767 | 0.119360 |
| 3 | 2.039757 | 0.485312 | 0.534770 | 0.449123 | 0.435343 | 0.264084 | 0.379235 | 0.298592 | 0.296285 | 0.178891 |
| 4 | 1.689370 | 0.466213 | 0.627736 | 0.402640 | 0.383613 | 0.242141 | 0.336123 | 0.313137 | 0.255725 | 0.117384 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 65 | 1.321976 | 0.285968 | 0.426944 | 0.332646 | 0.359031 | 0.337895 | 0.371682 | 0.262732 | 0.139557 | 0.109207 |
| 66 | 1.365696 | 0.417466 | 0.524789 | 0.301595 | 0.298802 | 0.221709 | 0.293173 | 0.296263 | 0.226661 | 0.208604 |
| 67 | 1.609150 | 0.550512 | 0.466619 | 0.377195 | 0.291001 | 0.399002 | 0.294533 | 0.262527 | 0.279050 | 0.138084 |
| 68 | 1.741880 | 0.421837 | 0.691825 | 0.398323 | 0.420704 | 0.253085 | 0.319976 | 0.245361 | 0.120431 | 0.109834 |
| 69 | 1.414913 | 0.326146 | 0.438764 | 0.257212 | 0.282921 | 0.298153 | 0.258446 | 0.286538 | 0.215420 | 0.114045 |

Independent subsamplings

Each row is the diagonal of $\Lambda^{(t)}_{\text{cv}(i)}$

I multiplied by 1000 only for visualization you shouldn't do it

# CV Shrinked Eigenvalues (1)

components (N stocks)

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1.826186 | 0.450317 | 0.774918 | 0.277689 | 0.301081 | 0.284515 | 0.262410 | 0.312953 | 0.278350 | 0.132593 |
| 1 | 1.848474 | 0.499851 | 0.467592 | 0.447850 | 0.450473 | 0.386625 | 0.407486 | 0.366422 | 0.276440 | 0.173862 |
| 2 | 1.077437 | 0.432883 | 0.483675 | 0.349912 | 0.278088 | 0.233495 | 0.235383 | 0.179277 | 0.169767 | 0.119360 |
| 3 | 2.039757 | 0.485312 | 0.534770 | 0.449123 | 0.435343 | 0.264084 | 0.379235 | 0.298592 | 0.296285 | 0.178891 |
| 4 | 1.689370 | 0.466213 | 0.627736 | 0.402640 | 0.383613 | 0.242141 | 0.336123 | 0.313137 | 0.255725 | 0.117384 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 65 | 1.321976 | 0.285968 | 0.426944 | 0.332646 | 0.359031 | 0.337895 | 0.371682 | 0.262732 | 0.139557 | 0.109207 |
| 66 | 1.365696 | 0.417466 | 0.524789 | 0.301595 | 0.298802 | 0.221709 | 0.293173 | 0.296263 | 0.226661 | 0.208604 |
| 67 | 1.609150 | 0.550512 | 0.466619 | 0.377195 | 0.291001 | 0.399002 | 0.294533 | 0.262527 | 0.279050 | 0.138084 |
| 68 | 1.741880 | 0.421837 | 0.691825 | 0.398323 | 0.420704 | 0.253085 | 0.319976 | 0.245361 | 0.120431 | 0.109834 |
| 69 | 1.414913 | 0.326146 | 0.438764 | 0.257212 | 0.282921 | 0.298153 | 0.258446 | 0.286538 | 0.215420 | 0.114045 |

Independent subsamplings

Average

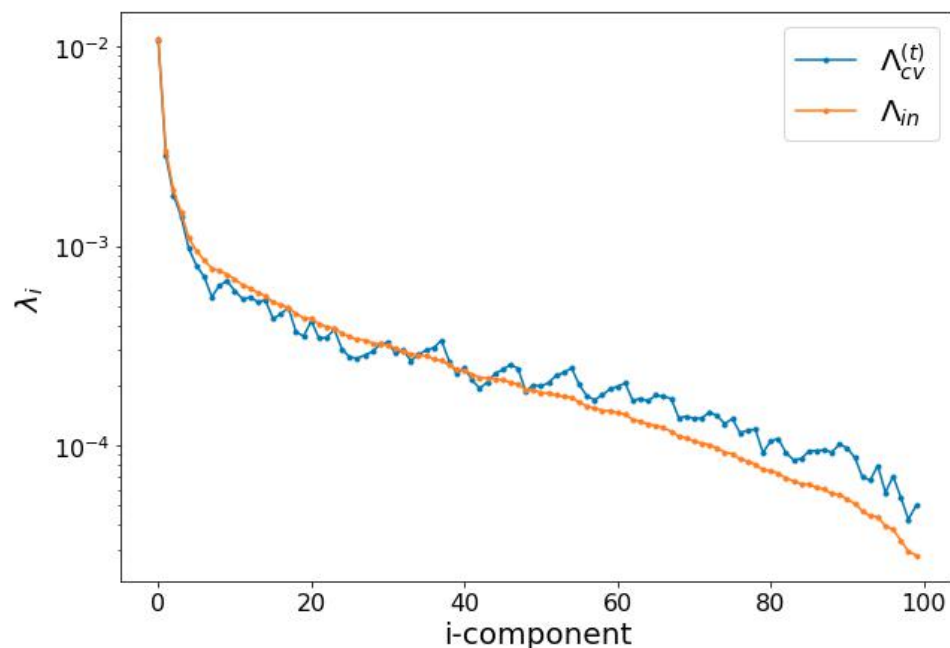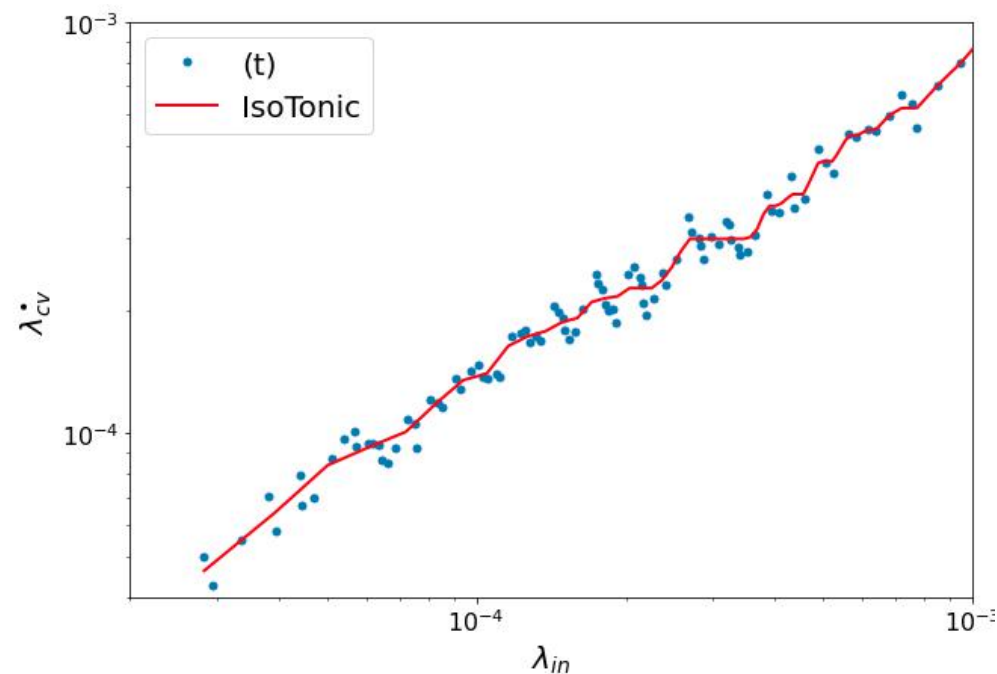|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1.472078 | 0.411935 | 0.515671 | 0.35736 | 0.369621 | 0.29016 | 0.302178 | 0.246213 | 0.214249 | 0.138083 |

$\Lambda_{\text{cv}}^{(t)}$     We need the last correction

# CV Shrinked Eigenvalues (2) [Regression]

The average eigenvalues are not monotonically decreasing

Ref [1] proposed an Isotonic regression of the observed in-sample eigenvalues on the average CV eigenvalues
[There is an error in their equation]



The isotonic regression is a technique of fitting a free-form line to a sequence of observations such that the fitted line is non-decreasing (or non-increasing) everywhere, and lies as close to the observations as possible. [Wikipedia]

[1] Reigneron, P. A., Nguyen, V., Ciliberti, S., Seager, P., & Bouchaud, J. P. (2019). The case for long-only agnostic allocation portfolios. arXiv:1906.05187.

# CV Shrinked Eigenvalues (3)



Finally, the CV estimator of the covariance matrix is defined as:

$$\boldsymbol{\Sigma}_{cv} = \boldsymbol{V}_{in}\boldsymbol{\Lambda}_{cv}^{(ISO)}\boldsymbol{V}'_{in}$$

If the sample-size error is the only source of noise $\boldsymbol{\Sigma}_{cv}$ converges to the oracle estimator

# Questions?

# Covariance Matrix Cleaning
## Predictability [TP]

Christian Bongiorno | Allocation de Portefeuille | 19-01-2023

# Project Overview

- Application of CV-Shrinkage on a selection of SP500.

- I will introduce you some performance metrics you should apply

- I will discuss some technical problem dealing with real data

- We will work together till the end of the lesson, then you can finish the assignment at home. You can work in groups of 3 (o 4), and I am expecting the notebook on Edunao within 1 week. I want a single notebook for group with the names of the students both code and pdf print with plots.

- Please comment everything. Do plots with x-labels, y-labels and legends.

- Suggestions: Use numpy, vectorialize whenever you can. Use for-cycle only if it is strictly necessary

# The Dataset

| index | A | AA | AAL | AAN | AAP | AAPL | ABBV | ABC | ABMD | ABT | ... | YUMC | Z | ZAYO | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2019-11-06 | 0.003172 | -0.031021 | -0.007081 | -0.049157 | 0.009235 | 0.000428 | 0.002928 | -0.016624 | -0.013560 | 0.009428 | ... | -0.002343 | 0.008061 | 0.001751 | 0.020! |
| 2019-11-07 | 0.001582 | 0.029692 | -0.002264 | -0.037938 | -0.002478 | 0.011475 | -0.000853 | -0.007465 | 0.020843 | 0.003243 | ... | 0.020428 | -0.005666 | -0.000583 | -0.009( |
| 2019-11-08 | 0.009571 | -0.001331 | -0.004217 | -0.030116 | -0.011168 | 0.002733 | 0.038278 | 0.023392 | 0.002234 | 0.004188 | ... | 0.007554 | 0.115921 | -0.001168 | 0.014< |
| 2019-11-11 | 0.002736 | -0.013855 | -0.005542 | -0.014708 | 0.008447 | 0.007888 | 0.005267 | -0.023868 | 0.002092 | 0.000239 | ... | -0.005030 | 0.028356 | -0.000877 | -0.004; |
| 2019-11-12 | 0.002599 | -0.001802 | -0.035605 | 0.016595 | -0.078080 | -0.000916 | 0.009065 | 0.017691 | 0.015599 | 0.006901 | ... | -0.018973 | -0.020927 | 0.001169 | 0.008< |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... ... | ... | ... | ... | |
| 2020-03-25 | 0.029414 | -0.025071 | 0.100391 | 0.132570 | 0.045156 | -0.005524 | 0.006204 | -0.058006 | 0.022165 | 0.015096 | ... | 0.035718 | 0.080084 | NaN | 0.036( |
| 2020-03-26 | 0.066168 | -0.035898 | 0.017392 | -0.002950 | 0.072035 | 0.051285 | 0.083582 | 0.110576 | 0.015403 | 0.069078 | ... | 0.013183 | 0.022066 | NaN | 0.043 |
| 2020-03-27 | -0.038863 | -0.043323 | -0.109199 | -0.001690 | -0.015775 | -0.042284 | -0.015837 | -0.028687 | -0.026699 | -0.016626 | ... | -0.056858 | -0.072186 | NaN | -0.060; |
| 2020-03-30 | 0.027059 | -0.076106 | -0.136384 | -0.047628 | 0.009370 | 0.028138 | 0.034754 | 0.074440 | 0.025492 | 0.062138 | ... | -0.005891 | -0.023468 | NaN | 0.034! |
| 2020-03-31 | -0.014554 | 0.014718 | -0.004910 | 0.010148 | -0.033509 | -0.002043 | 0.012547 | 0.004530 | -0.026043 | -0.005434 | ... | 0.007535 | -0.028464 | NaN | 0.028( |

The data are daily adjusted close-to-close **log-returns.**

Some data could miss for different reasons. Of course the methods I show you cannot handle missing data.
Even if you should not use future information in-sample, I can accept to select the set of stocks that are without NaN both in-sample and out-of-sample

# Risk Measure

- **Expected Risk**: $\sigma_P^{in} = \sqrt{252\ \boldsymbol{w} \boldsymbol{\Sigma}^{\blacksquare} \boldsymbol{w}'}$ where $\boldsymbol{\Sigma}^{\blacksquare}$ could be the filtered or the non-filtered covariance.

- **Realized Risk**: $\sigma_P^{out} = \sqrt{252\ \boldsymbol{w} \boldsymbol{\Sigma}^{out} \boldsymbol{w}'}$ of course $\boldsymbol{\Sigma}^{out}$ is the non-filtered covariance.

Note that the covariance could be defined as

$$\frac{1}{N-1} \sum_{i=1}^{N} \left( X_{ij} - \bar{X}_j \right) \left( X_{ik} - \bar{X}_k \right) \qquad \text{np.cov(X)}$$

$$\frac{1}{N} \sum_{i=1}^{N} \left( X_{ij} - \mathrm{E}(X_j) \right) \left( X_{ik} - \mathrm{E}(X_k) \right). \qquad \text{np.cov(X, bias=True)} \qquad \textbf{[USE THIS]}$$
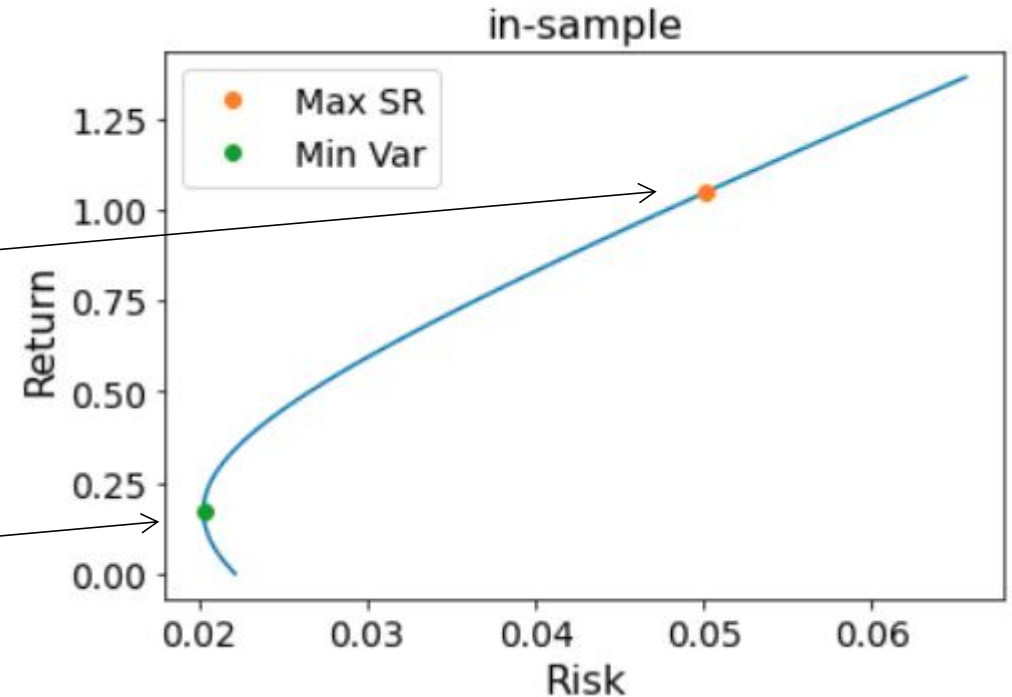
# Sharpe-Ratio

- $SR_P^{in} = = \dfrac{\langle r_P^{in} \rangle}{\sigma_P^{in}} \sqrt{252}$

- $SR_P^{out} = = \dfrac{\langle r_P^{out} \rangle}{\sigma_P^{out}} \sqrt{252}$

- To be correct the SR must be computed on the log-returns rather than on vanilla-returns. However, on short time-horizon they are approximately equivalents. Alternatively, you can computed the time-series of the daily log-returns of your portfolio, then $\sigma_P^{\blacksquare}$ will be the standard-deviation

# Analytical Solutions

$$w_{\text{MSR}} = \frac{\Sigma^{-1}\mu}{\mathbb{1}'\,\Sigma^{-1}\mu}$$

$$w_{\text{GMV}} = \frac{\Sigma^{-1}\mathbb{1}}{\mathbb{1}'\,\Sigma^{-1}\mathbb{1}}$$



in-sample

For all the other portfolios you have to solve the QP problem. However, if you do not have long-only contrain the solution is analytical.

# Portfolio Metrics

- The effective portfolio diversification is defined as

$$N_{eff} = \frac{1}{\sum_{i=1}^{N} w_i^2}$$

which represent the effective number of stocks with a significant amount of money invested.

- **Gross Leverage**

$$G = \sum_{i=1}^{1} |w_i|$$

When no short selling is allowed G=1. Portfolios with $G > 1$ have an additional intrinsic risk

due to high level of short-selling

Pantaleo, E., Tumminello, M., Lillo, F., & Mantegna, R. N. (2011). When do improved covariance matrix estimators enhance portfolio optimization? An empirical comparative study of nine estimators. Quantitative Finance, 11(7), 1067-1080.

# Assignments

1. Use the first 1000 days to obtain an efficient frontier that include the minVar and the maxSR portfolio (N=100).

      1a. Compare with the realized out-of-sample on the following 252 days.

      1b. Do the same by using in-sample the eigenvalue clipping and CV.

      1c. Comment the differences.

2. Study for different in-sample window size the minVar portfolio performances (in-sample vs out-of-sample) for the filtered and non-filtered covariance matrix.

      2a. Describe the behaviour of the portfolios with respect the previous cited metrics.

      2b. Comment the results from the point of view of the spectral decomposition.

      2c. Compare the eigenvalue distributions of the unfiltered and clipped eigenvalue , CV and the oracle estimators  for different in-sample window size and comment the results.