

---

# Penalised regressions and sparse hedging

# Summary

---

## Portfolio choice v2.0

- ▶ Classical portfolio theory (mean-variance) is polluted by **estimation errors**
- ▶ It can be shown that Minimum Variance weights are closely linked to **cross-sectional regressions**
- ▶ In order to reduce estimation risk,<sup>1</sup> it can be useful to resort to *penalised* regressions instead
- ▶ The latter can also be used solely for **prediction purposes**

---

<sup>1</sup>Or fine-tune the bias-variance tradeoff.

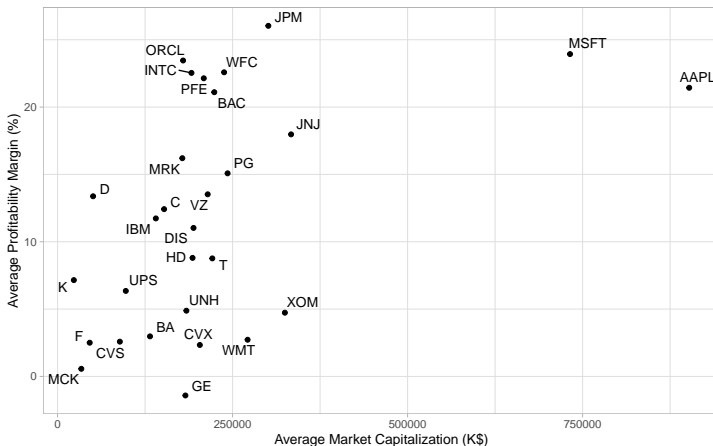
# Agenda

---

- 1 Regular regressions
- 2 Penalized regressions
- 3 Bias-variance tradeoff
- 4 Portfolio considerations
- 5 Implementation details

# For starters: illustrative example (1/4)

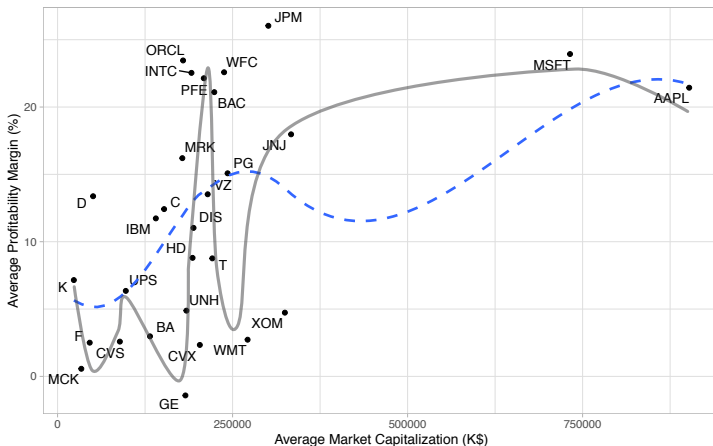
Let's say we are given the current data (plotted below)



We wish to model the dependence between **profitability** and **size**...

# For starters: illustrative example (2/4)

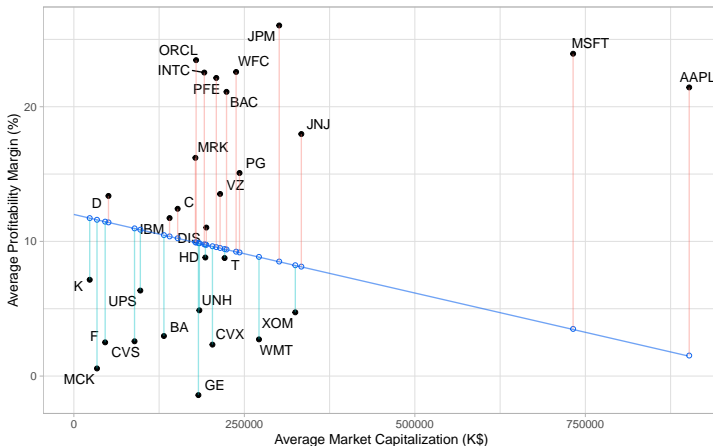
There are tons of ways to find fitting functions:



Maybe simplicity and parsimony are preferable...

# For starters: illustrative example (3/4)

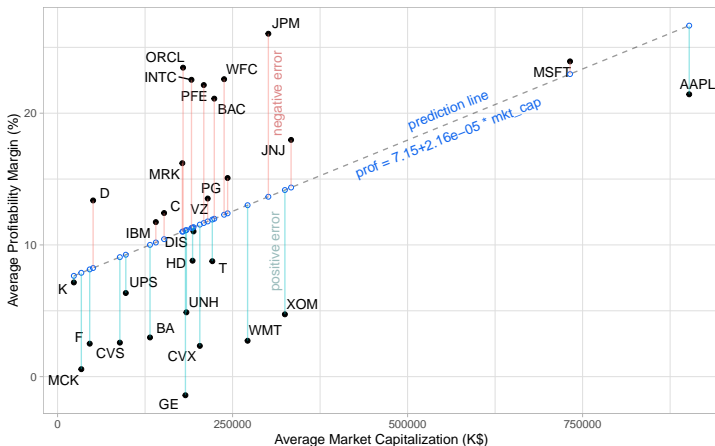
The simplest form is the linear function:  $prof = a + b * mkt\_cap$ .



But these errors are huge! You can't choose the line at random!

# For starters: illustrative example (4/4)

If we minimize the sum of squared errors:



Much better! It's the "best" possible line.

# Reminder (1/4)

## What's a (linear) regression?

You have a **dependent** variable  $Y$  and  $N$  **exogenous** variables  $X^{(n)}$  that are used to 'explain' (or forecast)  $Y$ . The model is *linear*, i.e.:

$$Y_t = \beta^{(1)} X_t^{(1)} + \dots + \beta^{(N)} X_t^{(N)} + \epsilon_t,$$

where  $t$  is the index of the observation. If there are  $T$  observations, the  $\beta^{(n)}$  are usually **estimated** (optimised, in fact) by minimising the quadratic error (sum of squared residuals):

$$\underset{\beta}{\operatorname{argmin}} \sum_{t=1}^T \left[ Y_t - \left( \beta^{(1)} X_t^{(1)} + \dots + \beta^{(N)} X_t^{(N)} \right) \right]^2.$$

In compact mode, the model is  $\mathbf{Y} = \mathbf{X}\beta + \epsilon$  and the solution to the above problem is  $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ . In terms of dimensions,  $\mathbf{X}$  is  $(T \times N)$ ,  $\beta$  is  $(N \times 1)$  so that both  $\mathbf{Y}$  and  $\epsilon$  are  $T$ -dimensional (column) vectors.



# Reminder (2/4)

## Proof of the OLS expression

The sum of squared errors is

$$\begin{aligned}\epsilon'\epsilon &= (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) \\ &= \mathbf{y}'\mathbf{y} - 2\beta'\mathbf{X}'\mathbf{y} - \beta'\mathbf{X}'\mathbf{X}\beta\end{aligned}$$

so the gradient (multivariate derivative) is

$$\frac{\partial}{\partial \beta} \epsilon'\epsilon = -2\mathbf{X}'\mathbf{y} - 2\mathbf{X}'\mathbf{X}\beta$$

and it's equal to  $\mathbf{0}$  if and only if  $\beta = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ .

Minimization would also require to check the second order condition  
→  $-2\mathbf{X}'\mathbf{X}$  as Hessian matrix.

# Reminder (3/4)

## Inference vs forecasting

The model is  $\mathbf{y} = \mathbf{X}\beta + \mathbf{e}$ .

- ▶ in **inference**, the aim is to quantify the absolute and relative importance of the variables (columns of  $\mathbf{X}$ ) in the way they explain  $\mathbf{y}$ . This is done by making assumptions on the errors  $\mathbf{e}$ , e.g. they follow a Gaussian law and

$$\mathbb{E}[\mathbf{e}|\mathbf{X}] = \mathbf{0}.$$

- ▶ in **forecasting**, the goal is simply to build the best possible model (i.e., that works out of sample). Errors are only interesting if they can be used to improve the accuracy of the model.

# Reminder (4/4)

## Where problems arise

The usual solution  $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}...$

- ▶ Obviously, the right part  $\mathbf{X}'\mathbf{Y}$  is not going to be a problem
- ▶ The left part  $(\mathbf{X}'\mathbf{X})^{-1}$  on the other hand...
- ▶ Whenever the number of features (expl. variables)  $N$  is larger than the number of observations ( $T$ ), the matrix  $\mathbf{X}'\mathbf{X}$  is **singular** (colinearity implies the matrix is not full rank).
- ▶ This happens when working with returns when the number of dates is smaller than the number of assets. This explains why daily returns are often required to compute covariance matrices (more on that later).

By the way, a great reference on regression analysis is **Econometric Analysis** by Greene (see Chap. 2→5).

# The partitioned form

## Another look when a constant is included

The compact form  $\mathbf{Y} = \alpha \mathbf{1}_T + \mathbf{X}\beta + \epsilon$  is solved using the following trick. We define  $\tilde{\mathbf{Y}} = \mathbf{Y} - \bar{Y} \mathbf{1}_T$  and  $\tilde{\mathbf{X}} = \mathbf{X} - \frac{1}{T} \mathbf{1}_T \mathbf{1}_T' \mathbf{X}$ , which are the original values, where each column has seen its mean retrieved (i.e., all columns have zero-mean). Then,

$$\hat{\beta} = (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}' \tilde{\mathbf{Y}},$$
$$\hat{\alpha} = \frac{1}{T} \mathbf{1}_T' (\mathbf{Y} - \mathbf{X} \hat{\beta})$$

where the first term in  $\hat{\beta}$  is the **inverse of the covariance matrix** of  $\mathbf{X}$  and the second term is the vector of **covariances** between  $\mathbf{X}$  and  $\mathbf{Y}$ .

It is a corollary from the Frisch-Waugh-Lovell theorem (Chap. 3 in Greene).

# Bias & Variance

## Useful for later..

Assuming a linear **generation** of data:  $\mathbf{Y} = \mathbf{X}\beta + \epsilon$  (where  $\beta$  is unknown), the estimator satisfies

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\beta + \epsilon) = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\epsilon$$

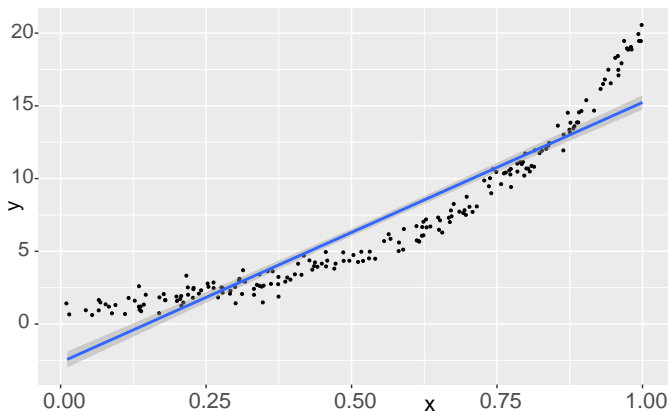
Under the **usual restriction**  $\mathbb{E}[\epsilon|\mathbf{X}] = 0$ , we have  $\mathbb{E}[\hat{\beta}] = \beta$ . Also,

$$\begin{aligned}\mathbb{V}[\hat{\beta}|\mathbf{X}] &= \mathbb{E}[(\hat{\beta} - \beta)'(\hat{\beta} - \beta)|\mathbf{X}] \\ &= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\epsilon\epsilon'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}|\mathbf{X}] \\ &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1},\end{aligned}$$

as long as  $\mathbb{E}[\epsilon\epsilon'] = \sigma^2\mathbf{I}_T$ .

# A key point

About the 'usual' restriction  $\mathbb{E}[\epsilon|\mathbf{X}] = 0$ :



The error terms must have zero mean at any location of  $X$ .

# Agenda

---

- 1 Regular regressions
- 2 Penalized regressions
- 3 Bias-variance tradeoff
- 4 Portfolio considerations
- 5 Implementation details

# Definition(1/2)

## Several families of penalised regressions

- ▶ **LASSO** (with  $\delta > 0$ ):

$$\min_{\beta} (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta), \quad \text{s.t. } \|\beta\|_1 \leq \delta$$

- ▶ **Ridge** (with  $\delta > 0$ ):

$$\min_{\beta} (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta), \quad \text{s.t. } \|\beta\|_2^2 \leq \delta$$

- ▶ **Elastic-net** (with  $\delta > 0$  and  $\alpha \in (0, 1)$ ): linear convex combination of both norms

$$\min_{\beta} (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta), \quad \text{s.t. } \alpha\|\beta\|_1 + (1 - \alpha)\|\beta\|_2^2 \leq \delta$$

We recall the norm notation:  $\|\beta\|_p = \left( \sum_{n=1}^N |\beta_n|^p \right)^{1/p}$ .



## Definition (2/2)

### Lagrange form (preferred)

- ▶ **LASSO** (with  $\lambda > 0$ ):

$$\min_{\beta} (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta) + \lambda \|\beta\|_1$$

- ▶ **Ridge** (with  $\lambda > 0$ ):

$$\min_{\beta} (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta) + \lambda \|\beta\|_2^2$$

- ▶ **Elastic-net** (with  $\lambda > 0$ ,  $\alpha \in [0, 1]$ ):

$$\min_{\beta} (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta) + \lambda \alpha \|\beta\|_1 + \lambda (1 - \alpha) \|\beta\|_2^2$$

We recall the norm notation:  $\|\beta\|_p = \left( \sum_{n=1}^N |\beta_n|^p \right)^{1/p}$ .

# Ridge: the closed form solution

## The only nice form

$$\begin{aligned}L(\beta) &= (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta) + \lambda\beta'\beta \\ &= \mathbf{Y}'\mathbf{Y} + \beta'\mathbf{X}'\mathbf{X}\beta - 2\beta'\mathbf{X}'\mathbf{Y} + \lambda\beta'\beta\end{aligned}$$

so that

$$\frac{\partial L}{\partial \beta} = 2\mathbf{X}'\mathbf{X}\beta - 2\mathbf{X}'\mathbf{Y} + 2\lambda\beta$$

and

$$\frac{\partial L}{\partial \beta} = 0 \Leftrightarrow \beta = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_N)^{-1}\mathbf{X}'\mathbf{Y}$$

where  $\mathbf{I}_N$  is the  $(N \times N)$  identity matrix. The shape is pretty close to the original OLS form. The only (**BIG**) difference is the  $\lambda\mathbf{I}_N$  inside the matrix inversion. This term ensures that the matrix is indeed **invertible**!

See Ledoit and Wolf (2004) for more details on the shrunk covariance matrix.

# Bias & Variance

## Slight changes...

Assuming a linear generation of data:  $\mathbf{Y} = \mathbf{X}\beta + \epsilon$ , the estimator satisfies

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_N)^{-1} \mathbf{X}'\mathbf{Y} = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_N)^{-1} \mathbf{X}'(\mathbf{X}\beta + \epsilon) \\ &= (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_N)^{-1} (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_N - \lambda\mathbf{I}_N)\beta + (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_N)^{-1} \mathbf{X}'\epsilon\end{aligned}$$

Under the usual restriction  $\mathbb{E}[\epsilon|\mathbf{X}]$ ,

$$\mathbb{E}[\hat{\beta}] - \beta = -\lambda(\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_N)^{-1} \beta.$$

The bias is not equal to zero. With regard to the variance, it can be shown to be equal to<sup>2</sup>

$$\mathbb{V}[\hat{\beta}] = \sigma^2(\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_N)^{-1} \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_N)^{-1}$$

and goes to zero as  $\lambda$  increases to infinity ( $\hat{\beta}$  converges to  $\mathbf{0}$ ).

<sup>2</sup>See Lecture notes on ridge regression by van Wieringen (2018)

# Example (0/3)

---

## A simple goal

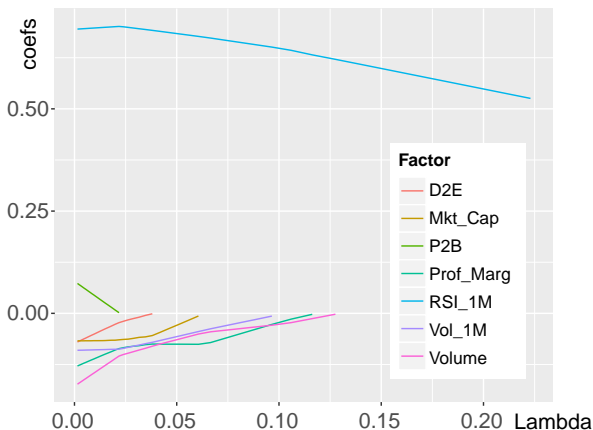
We aim to explain the return of the Apple Inc stock with the help of other **characteristics**:

$$r_t = \alpha + \beta^{Cap} Cap_t + \dots + \beta^{P2B} P2B_t + \epsilon_t,$$

where the factors (characteristics) have been normalised (because, e.g., Market Cap is out of range).

# Example (1/3)

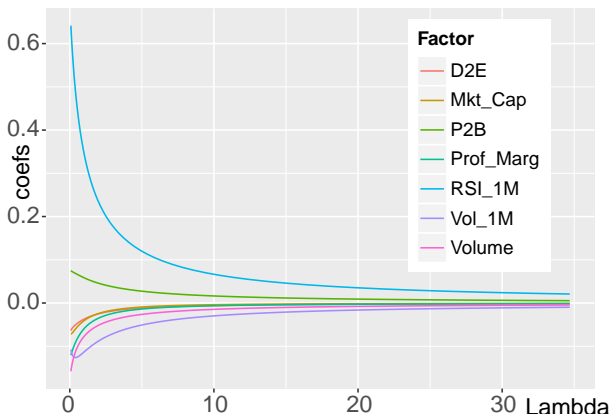
The **Lasso** first:



The coefficients converge to zero pretty quickly!  
→ the Lasso can serve as **feature selection** tool!

## Example (2/3)

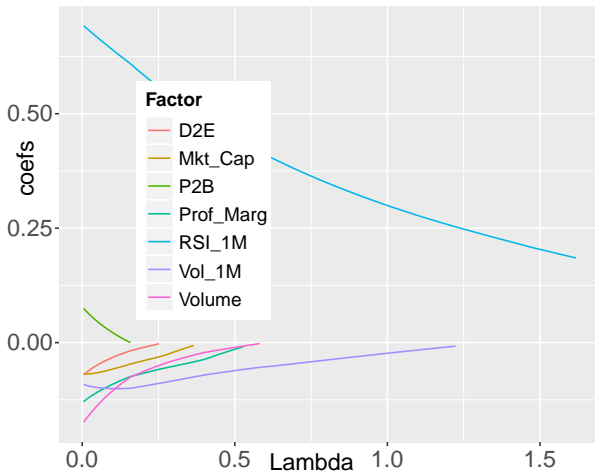
The **Ridge** next:



The coefficients converge to zero very slowly!

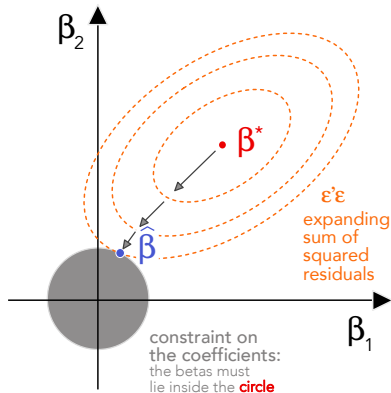
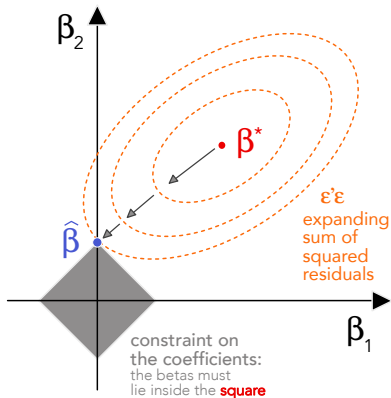
## Example (3/3)

The **elasticnet** at last:



The convergence speed/behaviour is somewhere in the middle!

# Why the LASSO is selective





# Agenda

---

- 1 Regular regressions
- 2 Penalized regressions
- 3 Bias-variance tradeoff
- 4 Portfolio considerations
- 5 Implementation details

# Definition (1/2)

## For any model

The aim is to find a model  $y = f(x) + \epsilon$  with  $\epsilon$  having zero mean,  $\sigma$  sd and  $\mathbb{E}[\epsilon|x] = 0$ . Assume a fitted model  $\hat{f}$  (regression, tree, SVM, NN, etc.). The quadratic error at point  $x$  is:

$$\begin{aligned}\mathbb{E}[(y - \hat{f}(x))^2] &= \mathbb{E}\left[\left(f(x) - \hat{f}(x) + \epsilon\right)^2\right] \\ &= \mathbb{E}\left[\left(f(x) - \hat{f}(x)\right)^2\right] + \sigma^2 \\ &= \underbrace{\left(\mathbb{E}\left[f(x) - \hat{f}(x)\right]\right)^2}_{\text{Bias}^2} + \underbrace{\mathbb{E}\left[\hat{f}(x)^2\right] - \mathbb{E}\left[\hat{f}(x)\right]^2}_{\text{Variance}} + \sigma^2\end{aligned}$$

In the above expression, the *sample*  $x$  is **random**.

## Definition (2/2)

In the above expression:

- ▶ the **bias** assesses the *average* distance between the *true*  $f$  and the fitted one  $\hat{f}$ ,
- ▶ the **variance** measures the dispersion of  $\hat{f}$  over all possible  $x$ ,
- ▶  $\sigma$  is the **irreducible error**: it is the level of approximation reached by the true  $f$ .

A very complex model is going to reduce bias but by doing so, it will move a lot to capture a maximum of points. To the other end of the spectrum, a constant model  $\hat{f}(x) = c$  has zero variance, but a rather large bias.

# Agenda

---

- 1 Regular regressions
- 2 Penalized regressions
- 3 Bias-variance tradeoff
- 4 Portfolio considerations
- 5 Implementation details

# References

---

The interested reader should have a look at:

- ▶ **On the Inverse of the Covariance Matrix in Portfolio Analysis** (Stevens, JF 1998)
- ▶ **Improving Mean Variance Optimization through Sparse Hedging Restrictions** (Goto and Xu JFQA 2015)

See also:

**The Sampling Error in Estimates of Mean-Variance Efficient Portfolio Weights** (Britten-Jones, JF 1999)

# Mean-variance portfolios

The textbook derivation (no budget constraint: simplicity)

Minimise variance for a fixed return... and hope for the best!

$$\min_{\mathbf{w}} \mathbf{w}' \Sigma \mathbf{w}, \quad \text{s.t. } \mathbf{w}' \boldsymbol{\mu} = r^*$$

with  $\boldsymbol{\mu}$  and  $\Sigma$  first and second moment of expected (excess) returns.

$$L(\mathbf{w}) = \mathbf{w}' \Sigma \mathbf{w} - \lambda(\mathbf{w}' \boldsymbol{\mu} - r^*)$$

Hence

$$\frac{\partial L}{\partial \mathbf{w}} = 2\Sigma \mathbf{w} - \lambda \boldsymbol{\mu} = 0 \Leftrightarrow \mathbf{w} = \frac{\lambda}{2} \Sigma^{-1} \boldsymbol{\mu},$$

where the  $\lambda$  scaling factor relates to risk preferences.

In practice,  $\boldsymbol{\mu}$  and  $\Sigma$  are **estimated**. Usually,  $\boldsymbol{\mu}$  is the one that is hardest to evaluate (and has the most importance!).

# MV portfolios as regression estimates!

## A surprising identity

Britten-Jones (1999) remarked that you can interpret:

$$\mathbf{w} = \underbrace{\Sigma^{-1}}_{(\mathbf{X}'\mathbf{X})^{-1}} \times \underbrace{\boldsymbol{\mu}}_{\mathbf{X}'\mathbf{1}}$$

Thus, if  $\mathbf{X}$  is a matrix of returns,  $\mathbf{w}$  is the OLS estimate of:

$$\mathbf{1} = \mathbf{X}\mathbf{w} + \mathbf{e}$$

Amazing!

# Minimum variance portfolios

## The textbook derivation (with budget constraint here)

Minimise variance and that's all!

$$\min_{\mathbf{w}} \mathbf{w}' \Sigma \mathbf{w}, \quad \text{s.t. } \mathbf{w}' \mathbf{1}_N = 1,$$

where the last constraint is simply the **budget constraint**. Then,

$$L(\mathbf{w}) = \mathbf{w}' \Sigma \mathbf{w} - \lambda(\mathbf{w}' \mathbf{1}_N - 1)$$

Hence

$$\frac{\partial L}{\partial \mathbf{w}} = 2\Sigma \mathbf{w} - \lambda \mathbf{1}_N = 0 \Leftrightarrow \mathbf{w} = \frac{\lambda}{2} \Sigma^{-1} \mathbf{1}_N,$$

where the  $\lambda$  is determined by the budget:

$$\mathbf{w} = \frac{\Sigma^{-1} \mathbf{1}_N}{\mathbf{1}_N' \Sigma^{-1} \mathbf{1}_N}.$$

IN ANY CASE,  $\Sigma^{-1}$  IS REQUIRED!



# A deep interpretation (1/3)

## A closer look at the inverse

If we **decompose** the matrix  $\Sigma$  into:

$$\Sigma = \begin{bmatrix} \sigma^2 & \mathbf{c}' \\ \mathbf{c} & \mathbf{C} \end{bmatrix},$$

classical **partitioning** results (e.g., via Schur complements) imply

$$\Sigma^{-1} = \begin{bmatrix} (\sigma^2 - \mathbf{c}'\mathbf{C}^{-1}\mathbf{c})^{-1} & -(\sigma^2 - \mathbf{c}'\mathbf{C}^{-1}\mathbf{c})^{-1}\mathbf{c}'\mathbf{C}^{-1} \\ -(\sigma^2 - \mathbf{c}'\mathbf{C}^{-1}\mathbf{c})^{-1}\mathbf{C}^{-1}\mathbf{c} & \mathbf{C}^{-1} + (\sigma^2 - \mathbf{c}'\mathbf{C}^{-1}\mathbf{c})^{-1}\mathbf{C}^{-1}\mathbf{c}\mathbf{c}'\mathbf{C}^{-1} \end{bmatrix}.$$

We are interested in the first line, which has 2 components: the factor  $(\sigma^2 - \mathbf{c}'\mathbf{C}^{-1}\mathbf{c})^{-1}$  and the line vector  $\mathbf{c}'\mathbf{C}^{-1}$ .  $\mathbf{C}$  is the covariance matrix of assets 2 to  $N$  and  $\mathbf{c}$  is the covariance between the first asset and all other assets. The first line of  $\Sigma^{-1}$  is

$$(\sigma^2 - \mathbf{c}'\mathbf{C}^{-1}\mathbf{c})^{-1} \begin{bmatrix} 1 & \underbrace{-\mathbf{c}'\mathbf{C}^{-1}}_{N-1 \text{ terms}} \end{bmatrix}.$$

# A deep interpretation (2/3)

## A special regression

We regress the returns of the first asset against those of all other assets:

$$r_{1,t} = a_1 + \sum_{n=2}^N \beta_{1|n} r_{n,t} + \epsilon_t, \text{ i.e., } \mathbf{r}_1 = a_1 \mathbf{1}_T + \mathbf{R}_{-1} \boldsymbol{\beta}_1 + \boldsymbol{\epsilon}$$

The **OLS estimator** for  $\boldsymbol{\beta}_1$  is

$$\hat{\boldsymbol{\beta}}_1 = \mathbf{C}^{-1} \mathbf{c}$$

and in addition (exercise!):

$$(1 - R^2) \sigma_{r_1}^2 = \sigma_{r_1}^2 - \mathbf{c}' \mathbf{C}^{-1} \mathbf{c} = \sigma_{\epsilon,1}^2$$

(where in fact  $\sigma_{r_1}^2 = \sigma^2$ ) and the first line of  $\boldsymbol{\Sigma}^{-1}$  is equal to

$$\frac{1}{\sigma_{\epsilon,1}^2} \times \begin{bmatrix} 1 & -\hat{\boldsymbol{\beta}}_1' \end{bmatrix}$$

# A deep interpretation (3/3)

## Remember Minimum Variance?

Each line of the inverse covariance matrix is multiplied by  $\mathbf{1}_N$ , and one line has form:

$$\frac{1}{\sigma^2(1-R^2)} \begin{bmatrix} 1 & -\beta' \end{bmatrix}$$

long position in stock 1

hedging positions in all other stocks

The red scaling factor is the inverse of  $\hat{\epsilon}'\hat{\epsilon}$ : when errors are large, positions are smaller. Note:  $\beta$  was chosen to **minimise the error/variance** pertaining to the above portfolio, which is a clear reference to 'minimum variance'.

This applies to all other lines/stocks!

# Bottom line

---

## Reduce estimation risk

- ▶ When  $T$  is not much larger than  $N$ , the **covariance matrix** is noisy, and its inverse even more. It is possible to **regularise** the estimation of these matrices upfront (see the work of Ledoit and Wolf).
- ▶ Nonetheless, the regularisation can occur at the stock level, using **hedging relationships**.
- ▶ The generalisation is straightforward: **bias vs tradeoff AND transaction costs** imply that penalised regression *could* be a good idea.

# Agenda

---

- 1 Regular regressions
- 2 Penalized regressions
- 3 Bias-variance tradeoff
- 4 Portfolio considerations
- 5 Implementation details

# The strategies (1/2)

**Sparse** portfolios.

## Pseudo-algorithm

For all dates  $t$ ,

► For all stocks  $i$ ,

1. estimate the elastic-net regression over the  $t = 1, \dots, T$  samples

$$\operatorname{argmin}_{\beta_{i|}} \left\{ \sum_{t=1}^T \left( r_{i,t} - a_i - \sum_{n \neq i}^N \beta_{i|n} r_{n,t} \right)^2 + \lambda \alpha \|\beta_{i|}\|_1 + \lambda (1 - \alpha) \|\beta_{i|}\|_2^2 \right\}$$

2. to get the weights of asset  $i$ :  $w_i = \frac{1 - \sum_{j \neq i} \beta_{i|j}}{\sigma_{\epsilon,i}^2}$

Note, when the number of asset is large, this may take too much time and the regularisation can be performed directly on the inverse covariance matrix (for Lasso and Ridge at least): precision matrix estimation via the **GLASSO** for instance.

# The strategies (2/2)

## Predictive LASSO

### Pseudo-algorithm

For all dates  $t$ ,

- ▶ For all stocks  $i$ ,
  1. estimate a predictive model

$$r_{n,t} = \alpha + \sum_{k=1}^K \beta_n^k f_{t-1}^k + \epsilon_{n,t}, \quad \text{s.t.} \quad \alpha \|\beta_n\|_1 + (1 - \alpha) \|\beta_n\|_2^2 \leq \delta$$

2. form a prediction with current factor/feature values  $f_t^k$
3. devise an allocation scheme based on this prediction

There is a large palette of potential **predictors**  $f_t^k$  (e.g., 'true' asset pricing factors), but we will simply work with (rescaled) firm attributes.

---

Thank you for your attention

Any questions?



# Proof for $R^2$

With  $\mathbf{X}$  being the concatenation of  $\mathbf{1}_T$  with returns  $\mathbf{R}_{-1}$  and with  $\mathbf{y} = \mathbf{r}_1$ ,  
 $R^2 = 1 - \frac{\epsilon'\epsilon}{T\sigma_Y^2} = 1 - \frac{\mathbf{y}'\mathbf{y} - \hat{\beta}'\mathbf{X}'\mathbf{X}\hat{\beta}}{T\sigma_Y^2} = 1 - \frac{\mathbf{y}'\mathbf{y} - \mathbf{y}'\mathbf{X}\hat{\beta}}{T\sigma_Y^2}$ , with fitted values  
 $\mathbf{X}\hat{\beta} = \hat{a}_1\mathbf{1}_T + \mathbf{R}_{-1}\mathbf{C}^{-1}\mathbf{c}$ . Hence,

$$T\sigma_{r_1}^2 R^2 = T\sigma_{r_1}^2 - \mathbf{r}_1'\mathbf{r}_1 + \hat{a}_1\mathbf{1}_T'\mathbf{r}_1 + \mathbf{r}_1'\mathbf{R}_{-1}\mathbf{C}^{-1}\mathbf{c}$$

$$T(1 - R^2)\sigma_{r_1}^2 = \mathbf{r}_1'\mathbf{r}_1 - \hat{a}_1\mathbf{1}_T'\mathbf{r}_1 - \left(\tilde{\mathbf{r}}_1 + \frac{\mathbf{1}_T\mathbf{1}_T'\mathbf{r}_1}{T}\right)' \left(\tilde{\mathbf{R}}_{-1} + \frac{\mathbf{1}_T\mathbf{1}_T'\mathbf{R}_{-1}}{T}\right) \mathbf{C}^{-1}\mathbf{c}$$

$$T(1 - R^2)\sigma_{r_1}^2 = \mathbf{r}_1'\mathbf{r}_1 - \hat{a}_1\mathbf{1}_T'\mathbf{r}_1 - T\mathbf{c}'\mathbf{C}^{-1}\mathbf{c} - \mathbf{r}_1'\frac{\mathbf{1}_T\mathbf{1}_T'}{T}\mathbf{R}_{-1}\mathbf{C}^{-1}\mathbf{c}$$

$$T(1 - R^2)\sigma_{r_1}^2 = \mathbf{r}_1'\mathbf{r}_1 - \frac{(\mathbf{1}_T'\mathbf{r}_1)^2}{T} - T\mathbf{c}'\mathbf{C}^{-1}\mathbf{c}$$

$$(1 - R^2)\sigma_{r_1}^2 = \sigma_{r_1}^2 - \mathbf{c}'\mathbf{C}^{-1}\mathbf{c}$$

where in the fourth equality we have plugged  $\hat{a}_1 = \frac{\mathbf{1}_T'}{T}(\mathbf{r}_1 - \mathbf{R}_{-1}\mathbf{C}^{-1}\mathbf{c})$ .  
(there is probably a simpler proof - see e.g. Section 3.5 in Greene)