

# Machine Learning with Python

## MTH786U/P 2023/24

### Week 3: Unstable regression problems

Nicola Perrà, Queen Mary University of London (QMUL)

# Recap of last week

Mathematical formulation of the regression problem:

Given input/output pairs  $\{(\mathbf{x}_i, y_i)\}_{i=1}^s$  find function  $f$  with

$$y_i \approx f(\mathbf{x}_i) \quad \forall i \in \{1, \dots, s\}$$



# Recap of last week

Mathematical formulation of the regression problem:

Given input/output pairs  $\{(\mathbf{x}_i, y_i)\}_{i=1}^s$  find function  $f$  with

$$y_i \approx f(\mathbf{x}_i) \quad \forall i \in \{1, \dots, s\}$$

Important to notice how each  $\mathbf{x}_i$  is a vector describing  $d$  features/variables

$$\mathbf{x}_i = (x_{i1}, \dots, x_{id})$$



# Example: linear regression

$$y_i \approx f(\mathbf{x}_i) \quad \forall i \in \{1, \dots, s\}$$



# Example: linear regression

$$y_i \approx f(\mathbf{x}_i) \quad \forall i \in \{1, \dots, s\}$$

How do we parametrise  $f$  ?



# Example: linear regression

$$y_i \approx f(\mathbf{x}_i) \quad \forall i \in \{1, \dots, s\}$$

How do we parametrise  $f$  ?

Example:

$$f(\mathbf{x}_i) = w_0 + \sum_{j=1}^d x_{ij} w_j$$



# Example: linear regression

$$y_i \approx f(\mathbf{x}_i) \quad \forall i \in \{1, \dots, s\}$$

How do we parametrise  $f$  ?

Example:

$$f(\mathbf{x}_i) = w_0 + \sum_{j=1}^d x_{ij} w_j$$

Linear transformation of vector  $\mathbf{x}_i = (x_{i1}, \dots, x_{id})$  with weights  $\mathbf{w} \in \mathbb{R}^{d+1}$



# Example: linear regression

Notation:  $f(\mathbf{x}_i) = w_0 + \sum_{j=1}^d w_j x_{ij} = \langle \mathbf{w}, \mathbf{x}_i \rangle$





# Example: linear regression

Notation:  $f(\mathbf{x}_i) = w_0 + \sum_{j=1}^d w_j x_{ij} = \langle \mathbf{w}, \mathbf{x}_i \rangle = \mathbf{w}^\top \mathbf{x}_i = \mathbf{x}_i^\top \mathbf{w}$



# Example: linear regression

Notation:  $f(\mathbf{x}_i) = w_0 + \sum_{j=1}^d w_j x_{ij} = \langle \mathbf{w}, \mathbf{x}_i \rangle = \mathbf{w}^\top \mathbf{x}_i = \mathbf{x}_i^\top \mathbf{w}$

$$\mathbf{x}_i := \begin{pmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{pmatrix} \in \mathbb{R}^{d+1}$$



# Example: linear regression

Notation:  $f(\mathbf{x}_i) = w_0 + \sum_{j=1}^d w_j x_{ij} = \langle \mathbf{w}, \mathbf{x}_i \rangle = \mathbf{w}^\top \mathbf{x}_i = \mathbf{x}_i^\top \mathbf{w}$

$$\mathbf{x}_i := \begin{pmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{pmatrix} \in \mathbb{R}^{d+1}$$

$$\mathbf{w} := \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix} \in \mathbb{R}^{d+1}$$

# Example: linear regression

Notation:  $f(\mathbf{x}_i) = w_0 + \sum_{j=1}^d w_j x_{ij} = \langle \mathbf{w}, \mathbf{x}_i \rangle = \mathbf{w}^\top \mathbf{x}_i = \mathbf{x}_i^\top \mathbf{w}$

Where this comes from?

$\mathbf{x}_i := \begin{pmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{pmatrix} \in \mathbb{R}^{d+1}$

$\mathbf{w} := \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix} \in \mathbb{R}^{d+1}$



# Example: linear regression

Notation:  $f(\mathbf{x}_i) = w_0 + \sum_{j=1}^d w_j x_{ij} = \langle \mathbf{w}, \mathbf{x}_i \rangle = \mathbf{w}^\top \mathbf{x}_i = \mathbf{x}_i^\top \mathbf{w}$

Where this comes from?

$$\mathbf{x}_i := \begin{pmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{pmatrix} \in \mathbb{R}^{d+1}$$

$$\mathbf{w} := \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix} \in \mathbb{R}^{d+1}$$

How do we choose  $w$  such that  $y_i \approx f(x_i)$  ?

# Example: linear regression

More in general?

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$



# Example: linear regression

More in general?

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$

Quality function

$$\text{MSE}(\mathbf{w}) = \frac{1}{2s} \sum_{i=1}^s \left| (\mathbf{X}\mathbf{w})_i - y_i \right|^2 = \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$



# Example: linear regression

More in general?

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$

Quality function

$$\text{MSE}(\mathbf{w}) = \frac{1}{2s} \sum_{i=1}^s \left| (\mathbf{X}\mathbf{w})_i - y_i \right|^2 = \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

The optimal solution is the minimisers of the MSE





# Example: linear regression

More in general?

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$

Quality function

$$\text{MSE}(\mathbf{w}) = \frac{1}{2s} \sum_{i=1}^s \left| (\mathbf{X}\mathbf{w})_i - y_i \right|^2 = \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

The optimal solution is the minimisers of the MSE

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \text{MSE}(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \right\}$$

# Example: linear regression

More in general?

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$

Quality function

$$\text{MSE}(\mathbf{w}) = \frac{1}{2s} \sum_{i=1}^s \left| (\mathbf{X}\mathbf{w})_i - y_i \right|^2 = \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

The optimal solution is the minimisers of the MSE

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \text{MSE}(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \right\}$$

Due to the convexity of the MSE the minimisers can be found as:

# Example: linear regression

More in general?

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$

Quality function

$$\text{MSE}(\mathbf{w}) = \frac{1}{2s} \sum_{i=1}^s \left| (\mathbf{X}\mathbf{w})_i - y_i \right|^2 = \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

The optimal solution is the minimisers of the MSE

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \text{MSE}(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \right\}$$

Due to the convexity of the MSE the minimisers can be found as:

$$\nabla \text{MSE}(\hat{\mathbf{w}}) \stackrel{!}{=} 0 \quad \Rightarrow \quad \mathbf{X}^\top \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^\top \mathbf{y}$$

# Example: linear regression

More in general?

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$

Quality function

$$\text{MSE}(\mathbf{w}) = \frac{1}{2s} \sum_{i=1}^s \left| (\mathbf{X}\mathbf{w})_i - y_i \right|^2 = \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

The optimal solution is the minimisers of the MSE

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \text{MSE}(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \right\}$$

Due to the convexity of the MSE the minimisers can be found as:

$$\nabla \text{MSE}(\hat{\mathbf{w}}) \stackrel{!}{=} 0 \quad \Rightarrow \quad \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y} \quad \Rightarrow \quad \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

# Unstable regression problems

The subject of this lecture is - yet again - the solution of regression problems of the form

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \text{MSE}(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \right\}$$



# Unstable regression problems

The subject of this lecture is - yet again - the solution of regression problems of the form

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \text{MSE}(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \right\}$$

We have established in previous lectures that the solution to this problem is

$$\mathbf{X}^\top \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^\top \mathbf{y}$$



# Unstable regression problems

The subject of this lecture is - yet again - the solution of regression problems of the form

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \text{MSE}(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \right\}$$

We have established in previous lectures that the solution to this problem is

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$



# Unstable regression problems

The subject of this lecture is - yet again - the solution of regression problems of the form

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \text{MSE}(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2s} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \right\}$$

We have established in previous lectures that the solution to this problem is

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

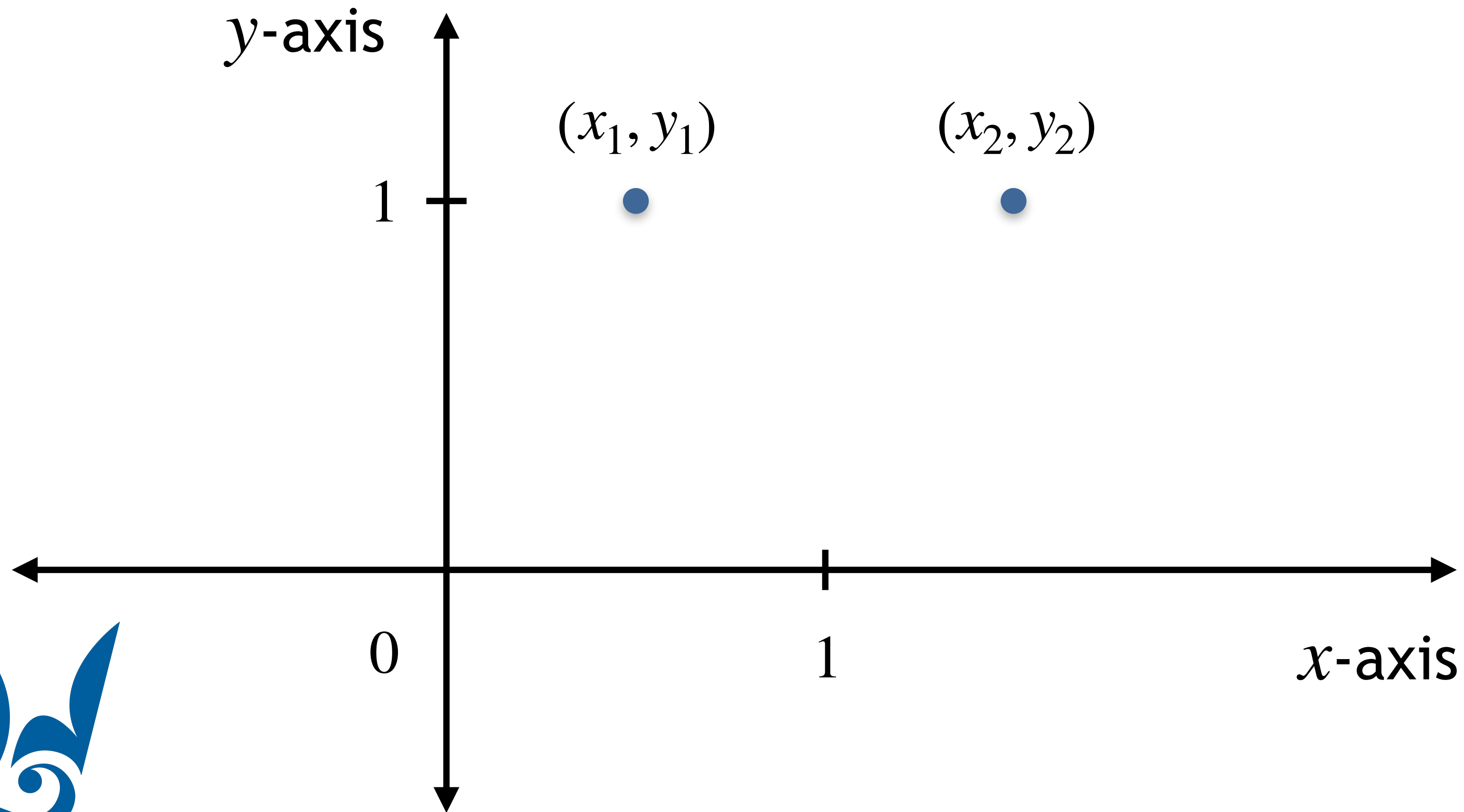
In this lecture we investigate if and when such solutions can become unstable (and what this means)





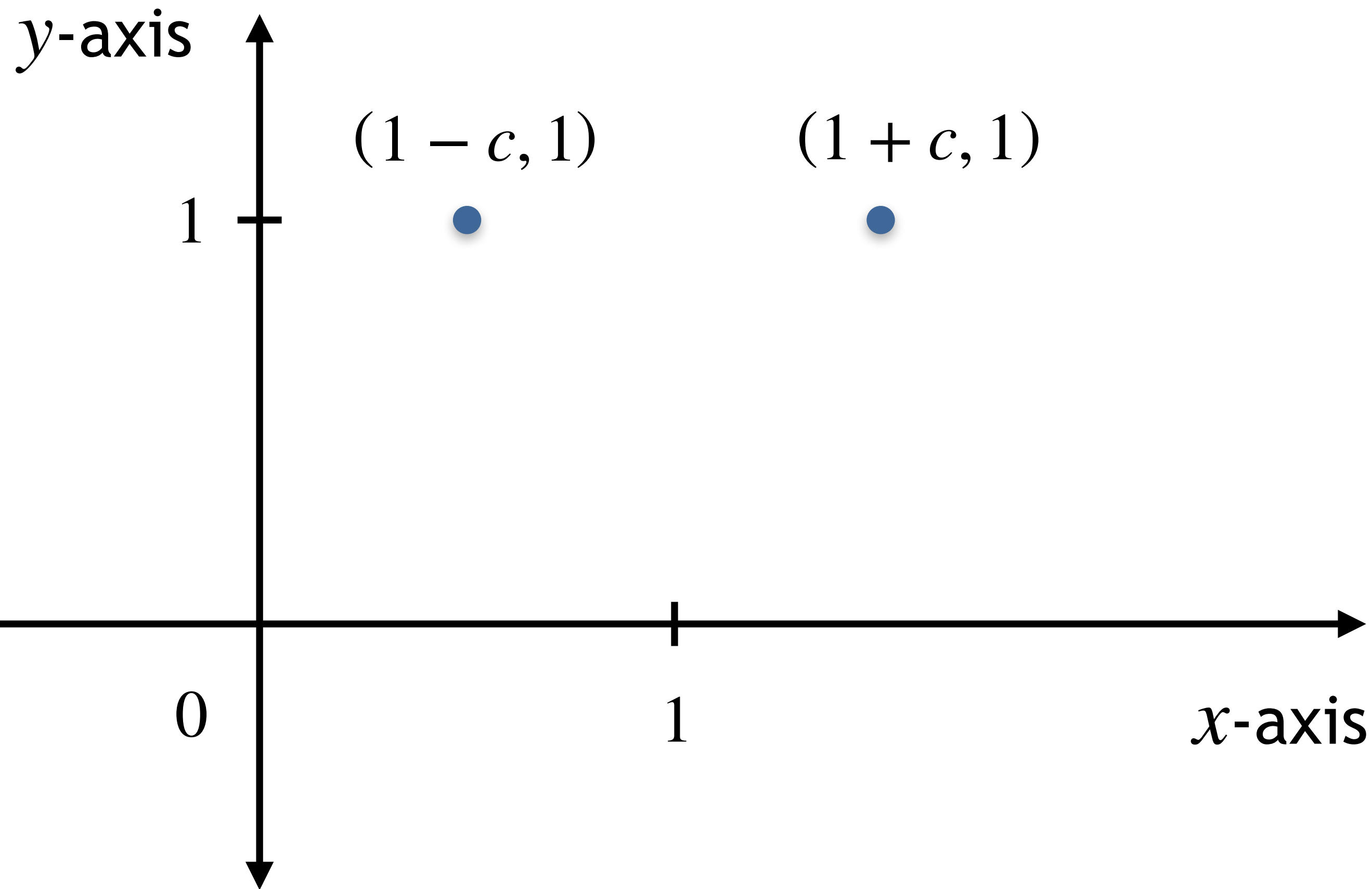
# Unstable regression problems

Let us consider the following simple regression problem:



# Unstable regression problems

Let us consider the following simple regression problem:



Data points

$$x_1 = 1 - c$$

$$y_1 = 1$$

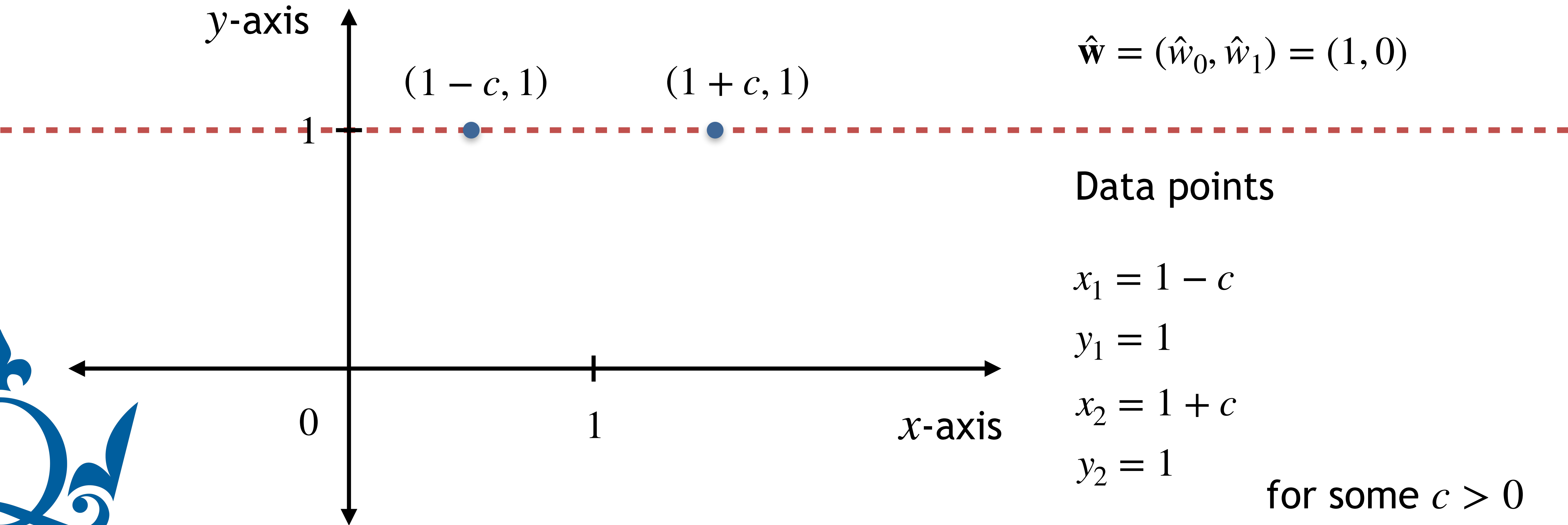
$$x_2 = 1 + c$$

$$y_2 = 1$$

for some  $c > 0$

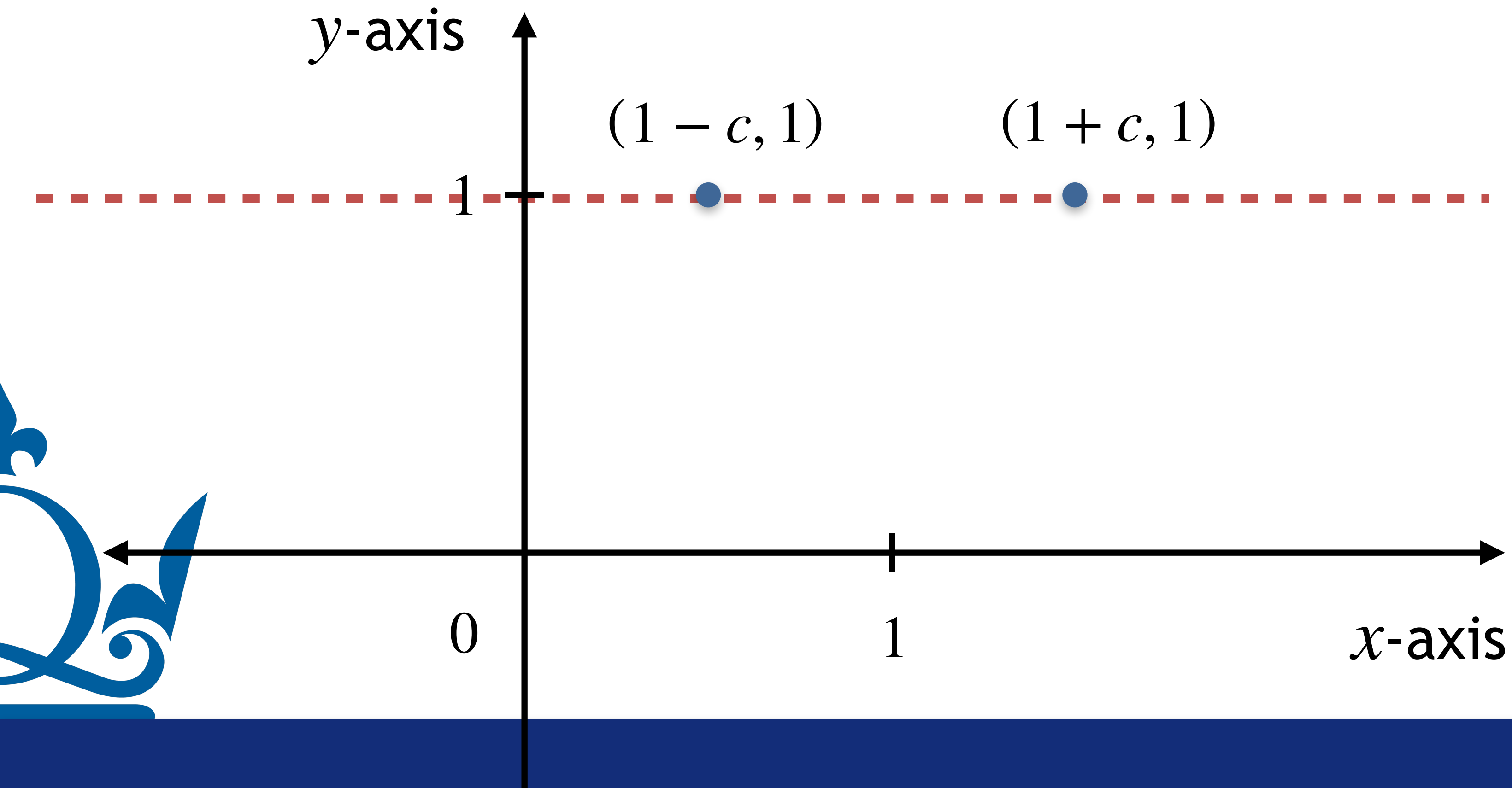
# Unstable regression problems

Let us consider the following simple regression problem:



# Why that's the solution?!

We can just solve the normal equation  $\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$



# Solving the normal equation

Normal equation  $\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$



# Solving the normal equation

Normal equation  $\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$

How do we solve it?



# Solving the normal equation

Normal equation  $\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$

How do we solve it?

Key: the unknown is...



# Solving the normal equation

Normal equation  $\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$

How do we solve it?

Key: the unknown is... the vector of weights  $\hat{\mathbf{w}}$





# Solving the normal equation

Normal equation  $\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$

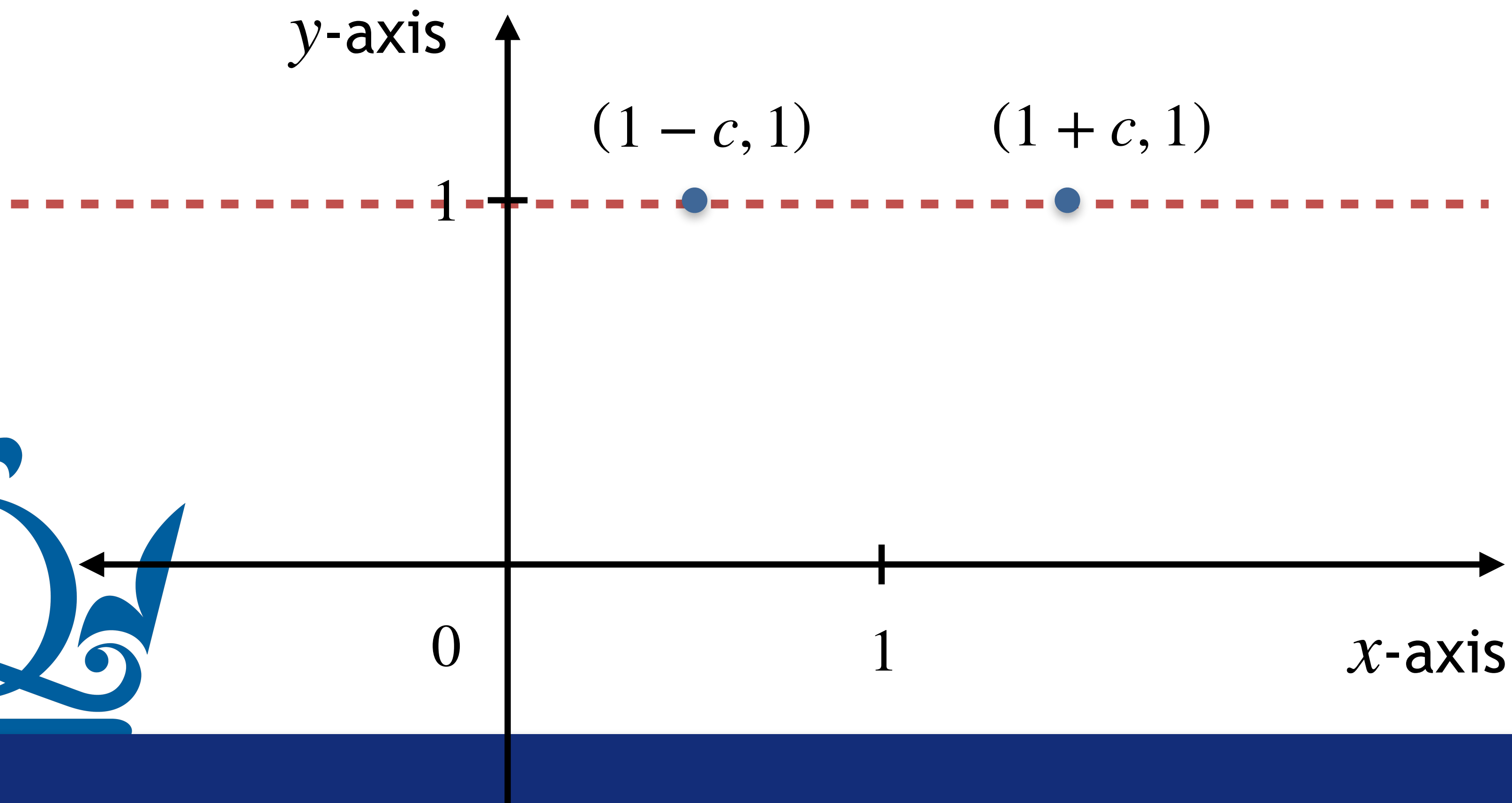
How do we solve it?

Key: the unknown is... the vector of weights  $\hat{\mathbf{w}}$

Hence we need just to evaluate  $\mathbf{X}^T \mathbf{X}$  and  $\mathbf{X}^T \mathbf{y}$  then solve the system!

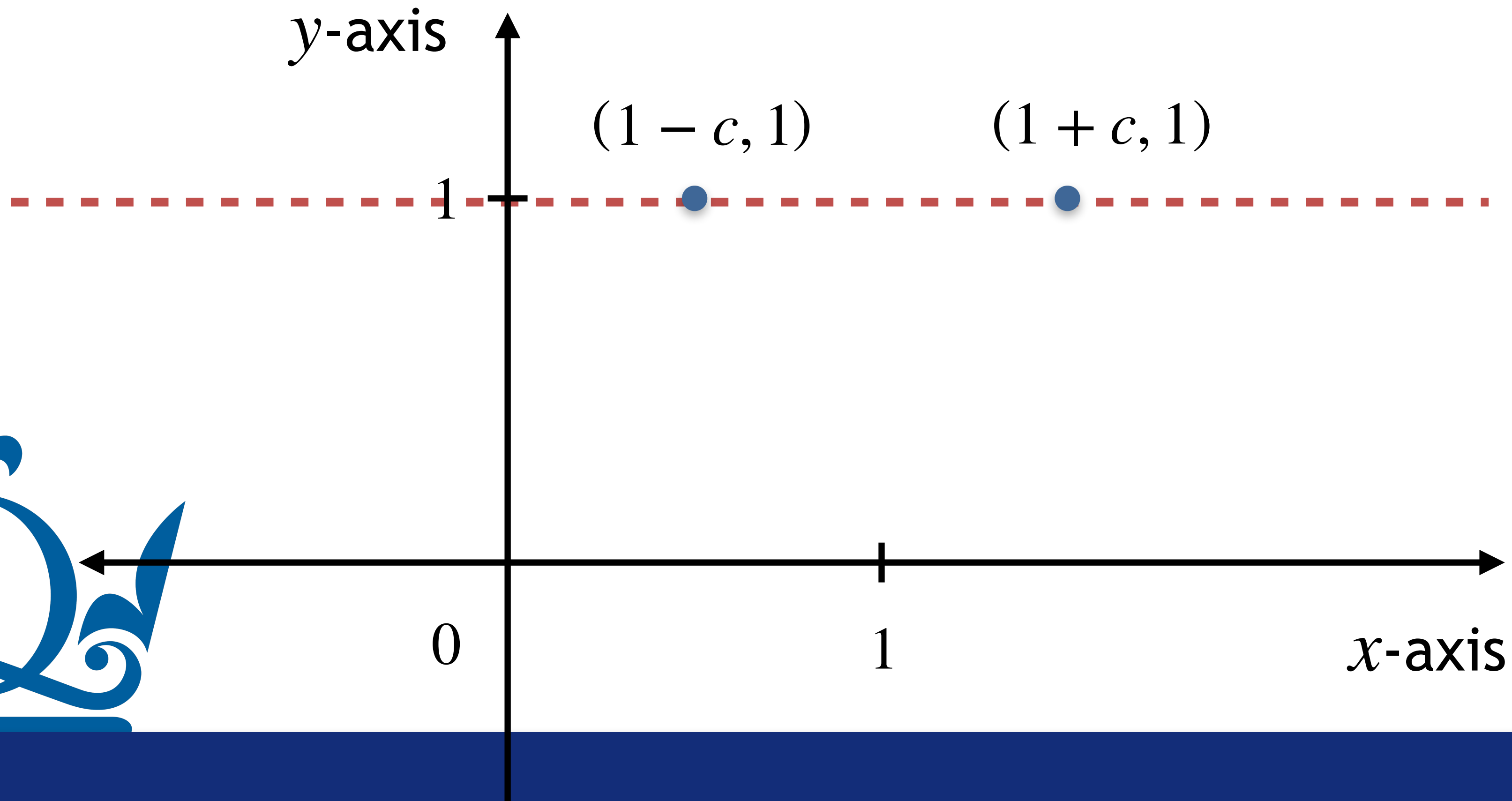


# Solving the normal equation



# Solving the normal equation

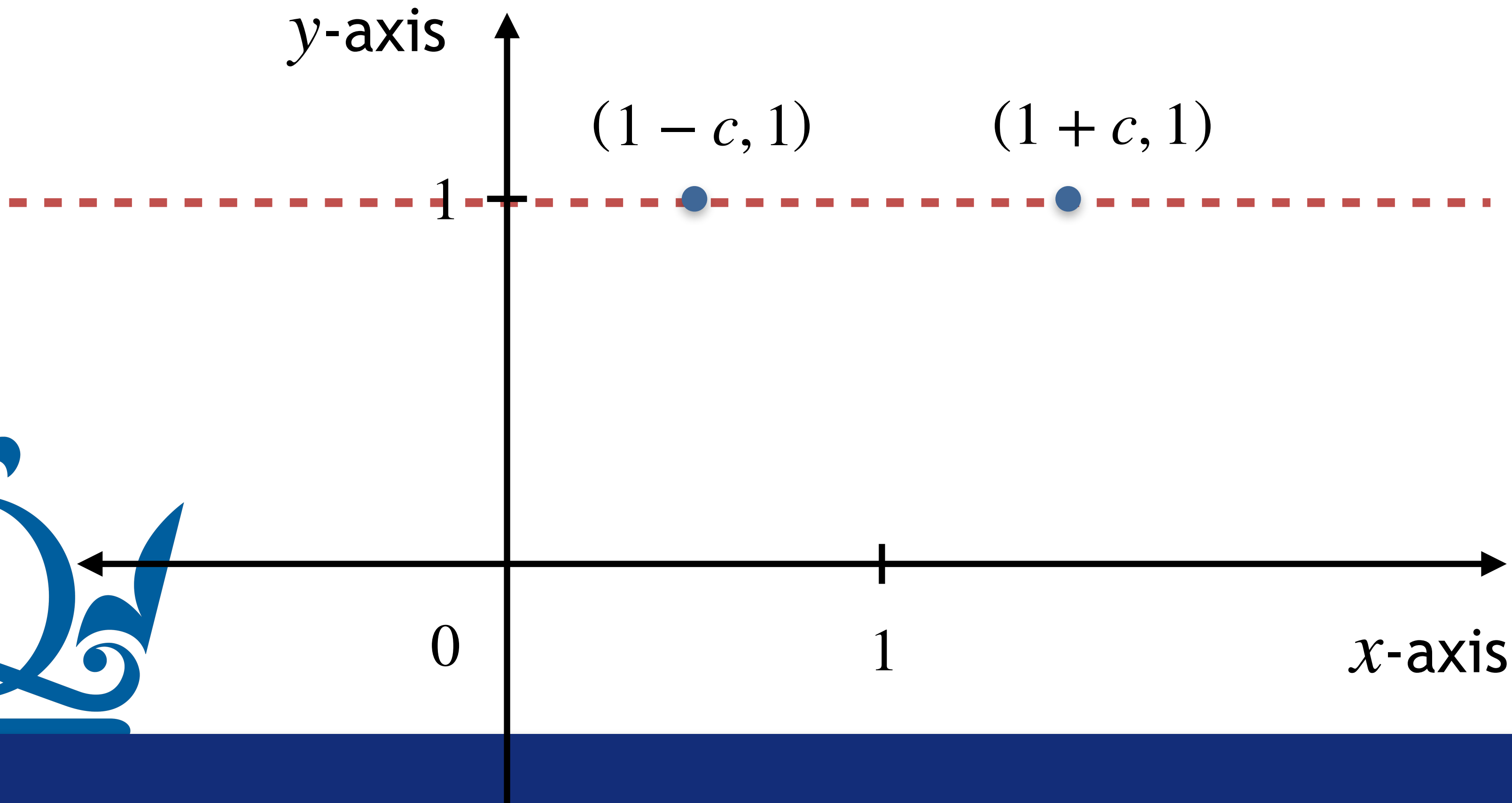
What's the data matrix  $\mathbf{X}$



# Solving the normal equation

What's the data matrix  $\mathbf{X}$

$$\mathbf{X} = \begin{pmatrix} 1 & 1 - c \\ 1 & 1 + c \end{pmatrix}$$

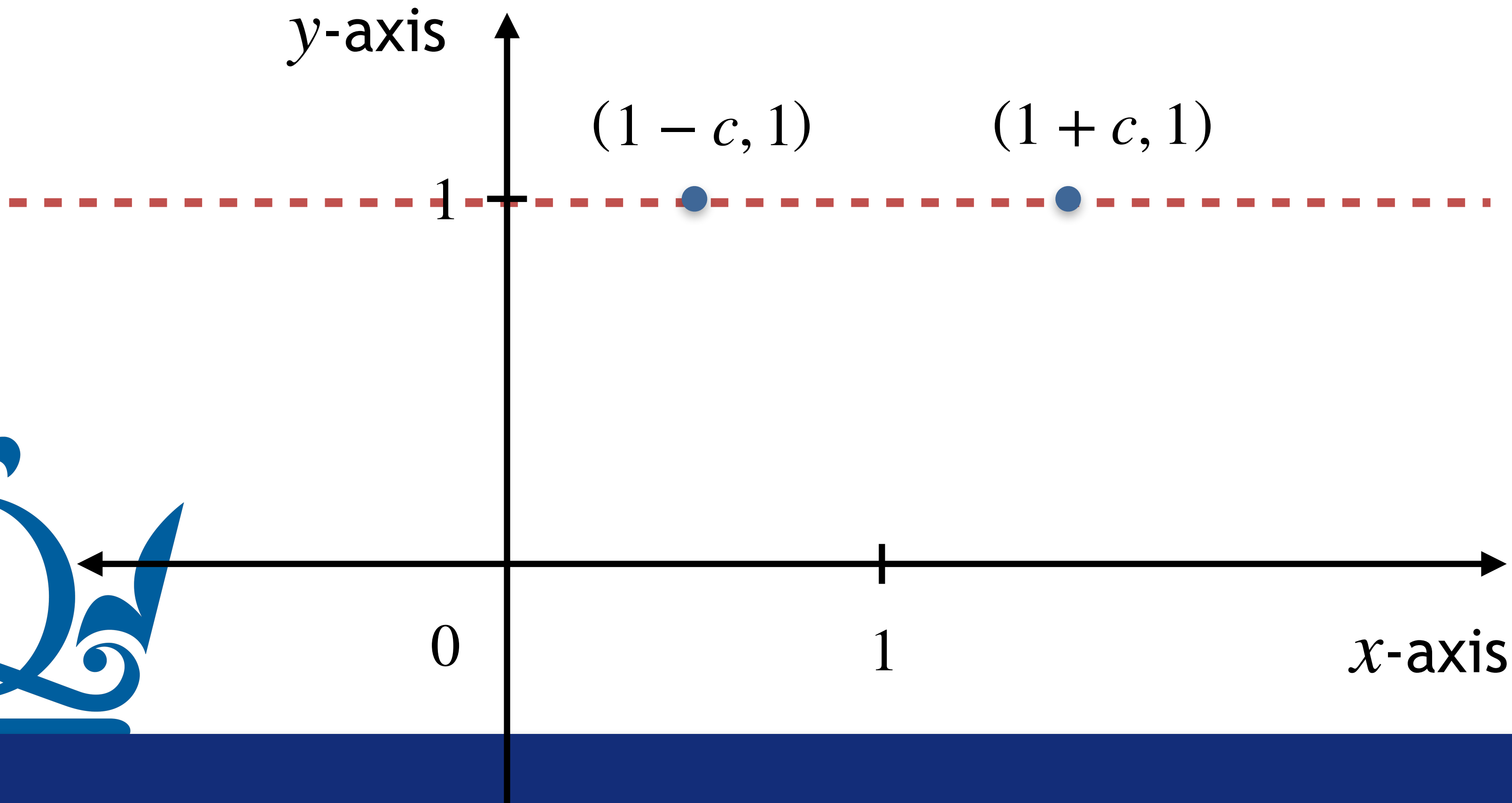


# Solving the normal equation

What's the data matrix  $\mathbf{X}$

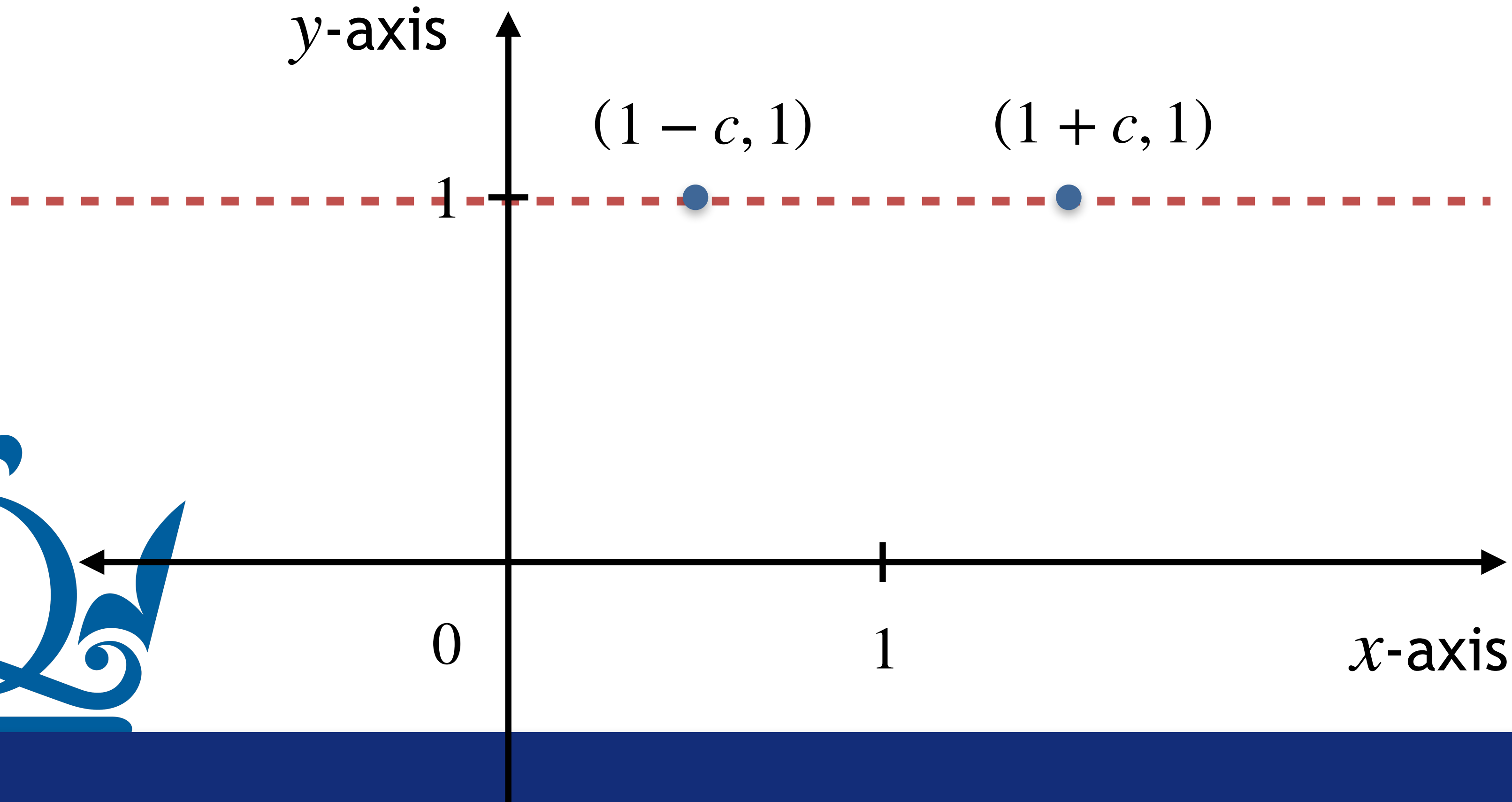
$$\mathbf{X} = \begin{pmatrix} 1 & 1 - c \\ 1 & 1 + c \end{pmatrix}$$

The transpose?



# Solving the normal equation

What's the data matrix  $\mathbf{X}$



$$\mathbf{X} = \begin{pmatrix} 1 & 1 - c \\ 1 & 1 + c \end{pmatrix}$$

The transpose?

$$\mathbf{X}^T = \begin{pmatrix} 1 & 1 \\ 1 - c & 1 + c \end{pmatrix}$$

# Solving the normal equation

What is the product  $\mathbf{X}^\top \mathbf{X}$

$$\mathbf{X}^\top = \begin{pmatrix} 1 & 1 \\ 1 - c & 1 + c \end{pmatrix}$$

$$\mathbf{X} = \begin{pmatrix} 1 & 1 - c \\ 1 & 1 + c \end{pmatrix}$$



# Solving the normal equation

What is the product  $\mathbf{X}^\top \mathbf{X}$

$$\mathbf{X}^\top = \begin{pmatrix} 1 & 1 \\ 1 - c & 1 + c \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} 1 & 1 - c \\ 1 & 1 + c \end{pmatrix}$$

$$\mathbf{X}^\top \mathbf{X} = \begin{pmatrix} 2 & 2 \\ 2 & 2 + 2c^2 \end{pmatrix}$$



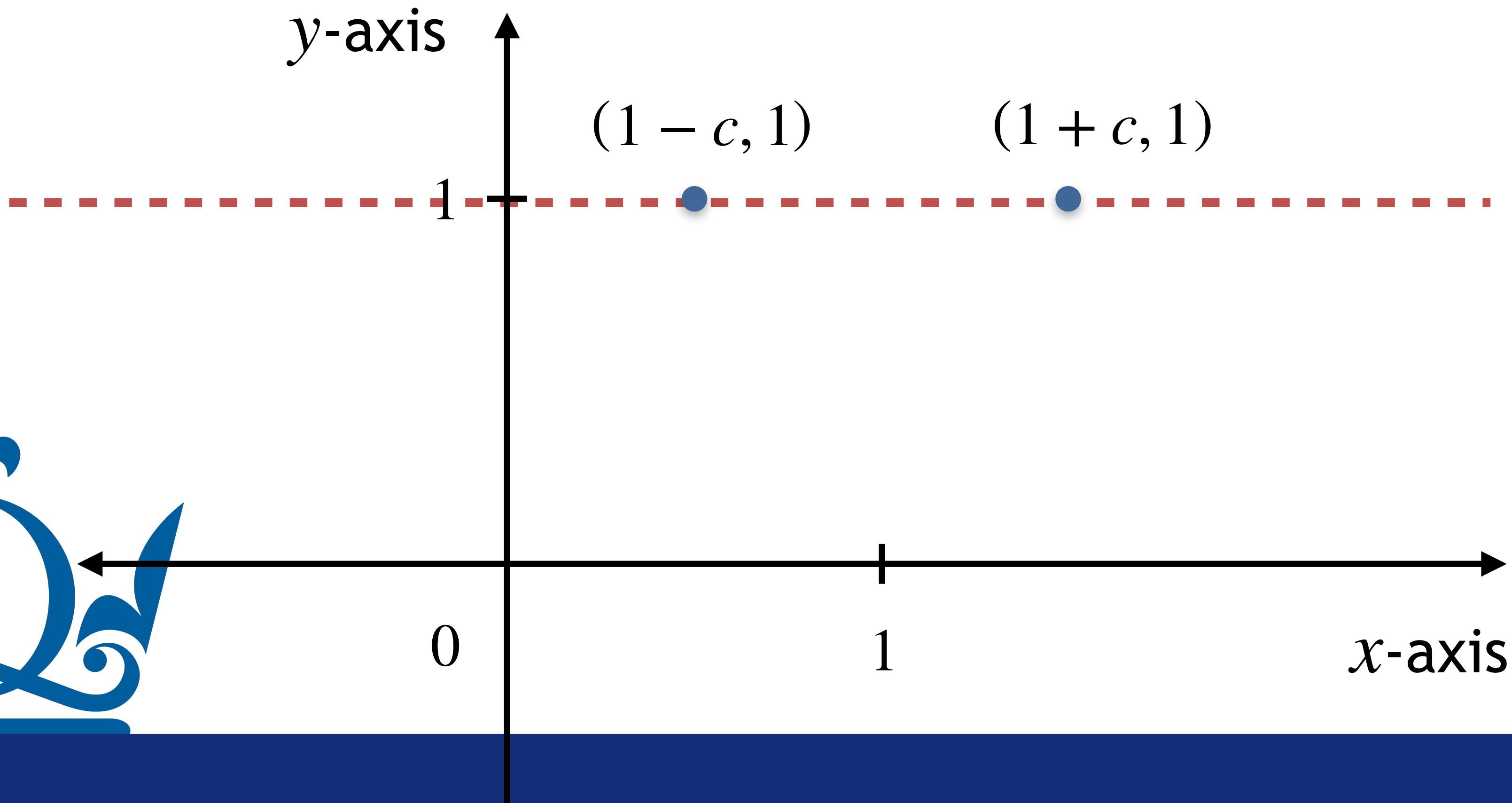


# Solving the normal equation

What is the product  $\mathbf{X}^T \mathbf{y}$

$$\mathbf{X}^T = \begin{pmatrix} 1 & 1 \\ 1-c & 1+c \end{pmatrix}$$

$$\mathbf{y} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$



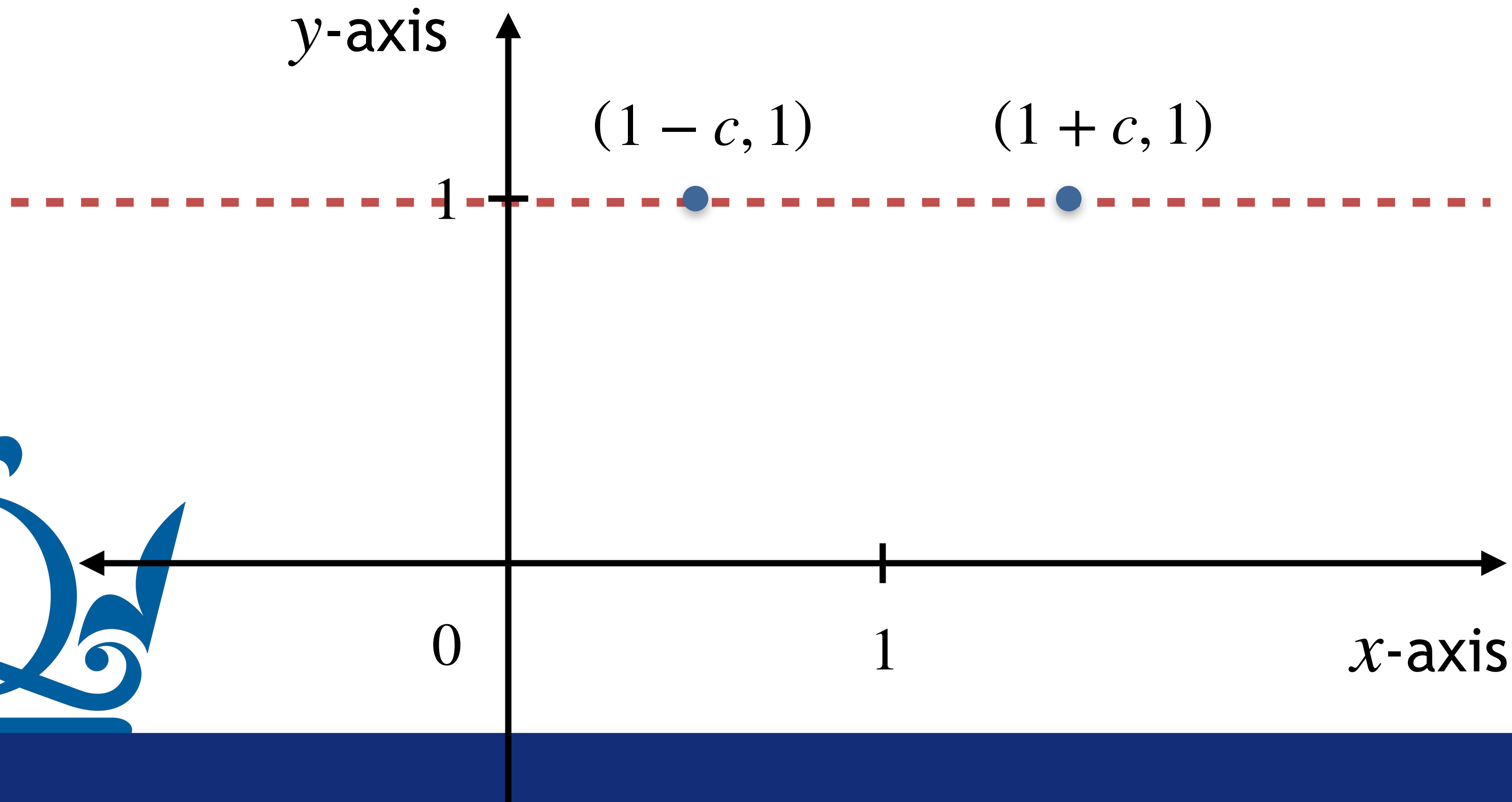
# Solving the normal equation

What is the product  $\mathbf{X}^T \mathbf{y}$

$$\mathbf{X}^T = \begin{pmatrix} 1 & 1 \\ 1-c & 1+c \end{pmatrix}$$

$$\mathbf{y} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\mathbf{X}^T \mathbf{y} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$



# Solving the normal equation

Almost done! We are solving  $\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$

$$\mathbf{X}^T \mathbf{X} = \begin{pmatrix} 2 & 2 \\ 2 & 2 + 2c^2 \end{pmatrix} \quad \mathbf{X}^T \mathbf{y} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$



# Solving the normal equation

Almost done! We are solving  $\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$

$$\mathbf{X}^T \mathbf{X} = \begin{pmatrix} 2 & 2 \\ 2 & 2 + 2c^2 \end{pmatrix} \quad \mathbf{X}^T \mathbf{y} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

Hence, we can write



# Solving the normal equation

Almost done! We are solving  $\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$

$$\mathbf{X}^T \mathbf{X} = \begin{pmatrix} 2 & 2 \\ 2 & 2 + 2c^2 \end{pmatrix} \quad \mathbf{X}^T \mathbf{y} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

Hence, we can write

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$



# Solving the normal equation

Almost done! We are solving  $\mathbf{X}^\top \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^\top \mathbf{y}$

$$\mathbf{X}^\top \mathbf{X} = \begin{pmatrix} 2 & 2 \\ 2 & 2 + 2c^2 \end{pmatrix} \quad \mathbf{X}^\top \mathbf{y} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

Hence, we can write

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

What is the solution?



# Solving the normal equation

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$



# Solving the normal equation

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Implies





# Solving the normal equation

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Implies

{



# Solving the normal equation

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Implies

$$\begin{cases} \hat{w}_1 + \hat{w}_2 = 1 \end{cases}$$



# Solving the normal equation

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Implies

$$\begin{cases} \hat{w}_1 + \hat{w}_2 = 1 \\ \hat{w}_1 + (1 + c^2)\hat{w}_2 = 1 \end{cases}$$



# Solving the normal equation

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Implies

$$\begin{cases} \hat{w}_1 + \hat{w}_2 = 1 \\ \hat{w}_1 + (1 + c^2)\hat{w}_2 = 1 \end{cases}$$



# Solving the normal equation

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Implies

$$\begin{cases} \hat{w}_1 + \hat{w}_2 = 1 \\ \hat{w}_1 + (1 + c^2)\hat{w}_2 = 1 \end{cases} \quad \begin{cases} \hat{w}_1 = 1 - \hat{w}_2 \end{cases}$$



# Solving the normal equation

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Implies

$$\begin{cases} \hat{w}_1 + \hat{w}_2 = 1 \\ \hat{w}_1 + (1 + c^2)\hat{w}_2 = 1 \end{cases}$$

$$\begin{cases} \hat{w}_1 = 1 - \hat{w}_2 \\ -\hat{w}_2 + (1 + c^2)\hat{w}_2 = 0 \end{cases}$$



# Solving the normal equation

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Implies

$$\begin{cases} \hat{w}_1 + \hat{w}_2 = 1 \\ \hat{w}_1 + (1 + c^2)\hat{w}_2 = 1 \end{cases}$$

$$\begin{cases} \hat{w}_1 = 1 - \hat{w}_2 \\ -\hat{w}_2 + (1 + c^2)\hat{w}_2 = 0 \end{cases}$$

$$\hat{w}_1 = 1$$



# Solving the normal equation

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Implies

$$\begin{cases} \hat{w}_1 + \hat{w}_2 = 1 \\ \hat{w}_1 + (1 + c^2)\hat{w}_2 = 1 \end{cases}$$

$$\begin{cases} \hat{w}_1 = 1 - \hat{w}_2 \\ -\hat{w}_2 + (1 + c^2)\hat{w}_2 = 0 \end{cases}$$

$$\hat{w}_1 = 1$$

$$\hat{w}_2 = 0$$





# Solving the normal equation

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \begin{pmatrix} \hat{w}_1 \\ \hat{w}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Implies

$$\begin{cases} \hat{w}_1 + \hat{w}_2 = 1 \\ \hat{w}_1 + (1 + c^2)\hat{w}_2 = 1 \end{cases}$$

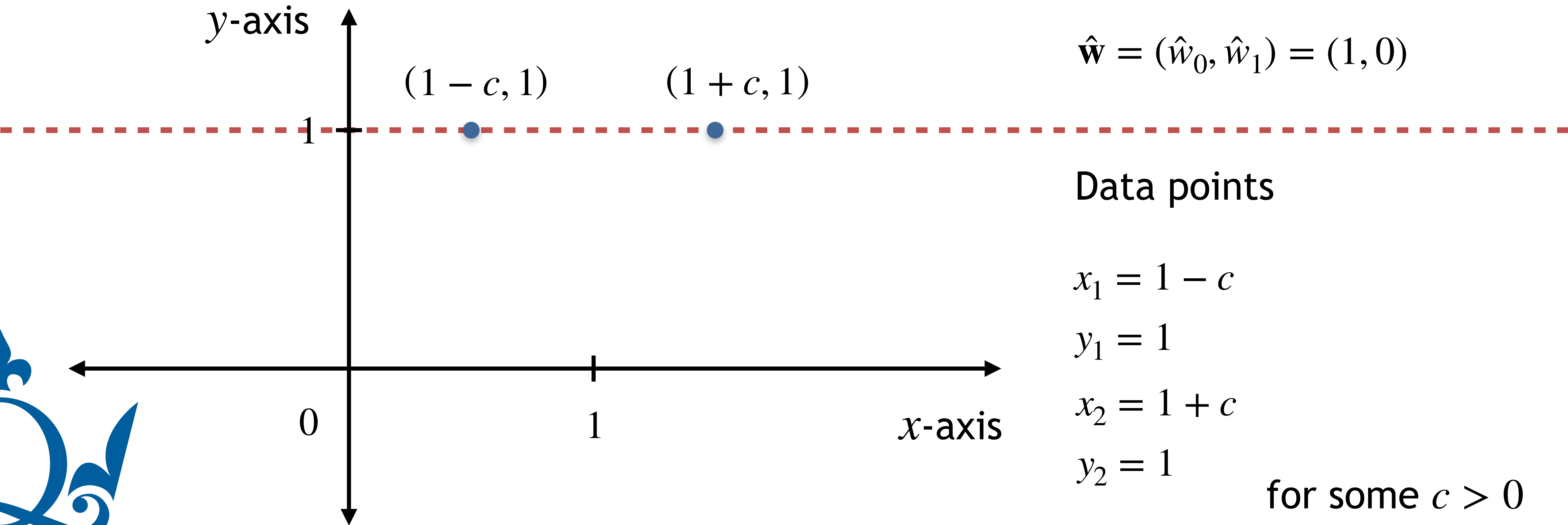
$$\begin{cases} \hat{w}_1 = 1 - \hat{w}_2 \\ -\hat{w}_2 + (1 + c^2)\hat{w}_2 = 0 \end{cases}$$

$$\begin{matrix} \hat{w}_1 = 1 \\ \hat{w}_2 = 0 \end{matrix} \rightarrow \hat{\mathbf{w}} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$



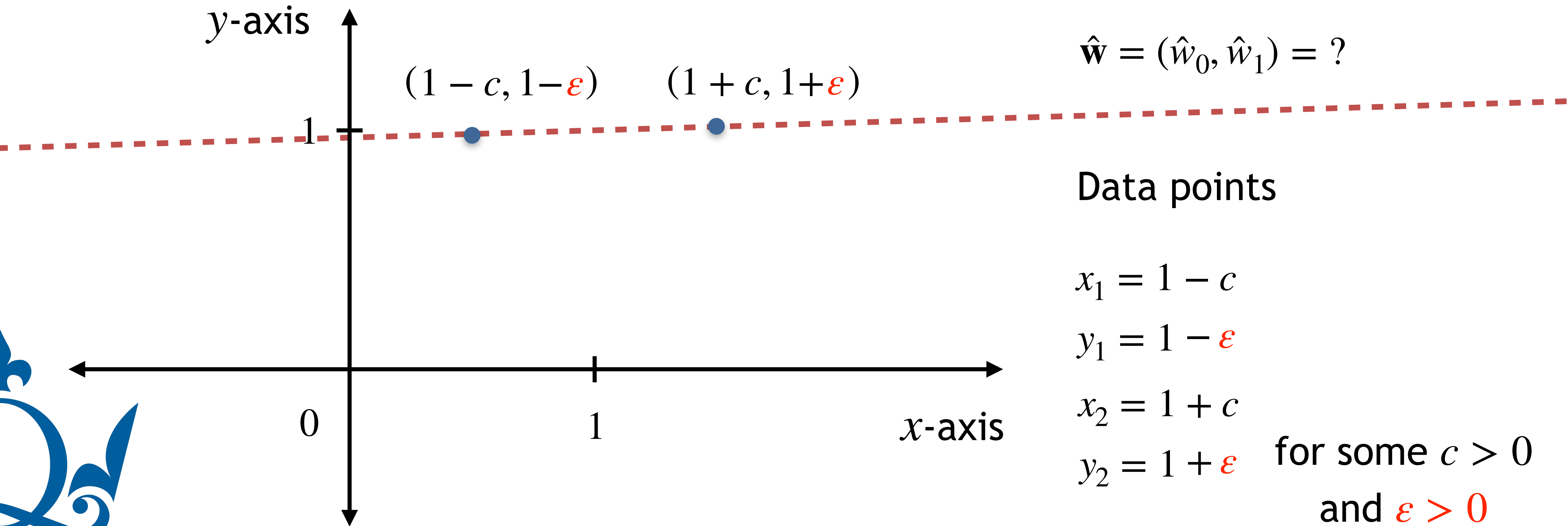
# Unstable regression problems

What if we make a small error when measuring  $y_1$  and  $y_2$ ?



# Unstable regression problems

What if we make a small error when measuring  $y_1$  and  $y_2$ ?



# Unstable regression problems

Data matrix  $\mathbf{X}$  and data vector  $\mathbf{y}$  of our problem read

$$\mathbf{X} = \begin{pmatrix} 1 & 1 - c \\ 1 & 1 + c \end{pmatrix} \quad \text{and} \quad \mathbf{y} = \begin{pmatrix} 1 - \varepsilon \\ 1 + \varepsilon \end{pmatrix};$$



# Unstable regression problems

Data matrix  $\mathbf{X}$  and data vector  $\mathbf{y}$  of our problem read

$$\mathbf{X} = \begin{pmatrix} 1 & 1 - c \\ 1 & 1 + c \end{pmatrix} \quad \text{and} \quad \mathbf{y} = \begin{pmatrix} 1 - \varepsilon \\ 1 + \varepsilon \end{pmatrix};$$

hence, we compute

$$\mathbf{X}^T \mathbf{X} = \begin{pmatrix} 2 & 2 \\ 2 & 2 + 2c^2 \end{pmatrix} \quad \text{and} \quad \mathbf{X}^T \mathbf{y} = \begin{pmatrix} 2 \\ 2 + 2c\varepsilon \end{pmatrix}$$



# Unstable regression problems

Data matrix  $\mathbf{X}$  and data vector  $\mathbf{y}$  of our problem read

$$\mathbf{X} = \begin{pmatrix} 1 & 1 - c \\ 1 & 1 + c \end{pmatrix} \quad \text{and} \quad \mathbf{y} = \begin{pmatrix} 1 - \varepsilon \\ 1 + \varepsilon \end{pmatrix};$$

hence, we compute

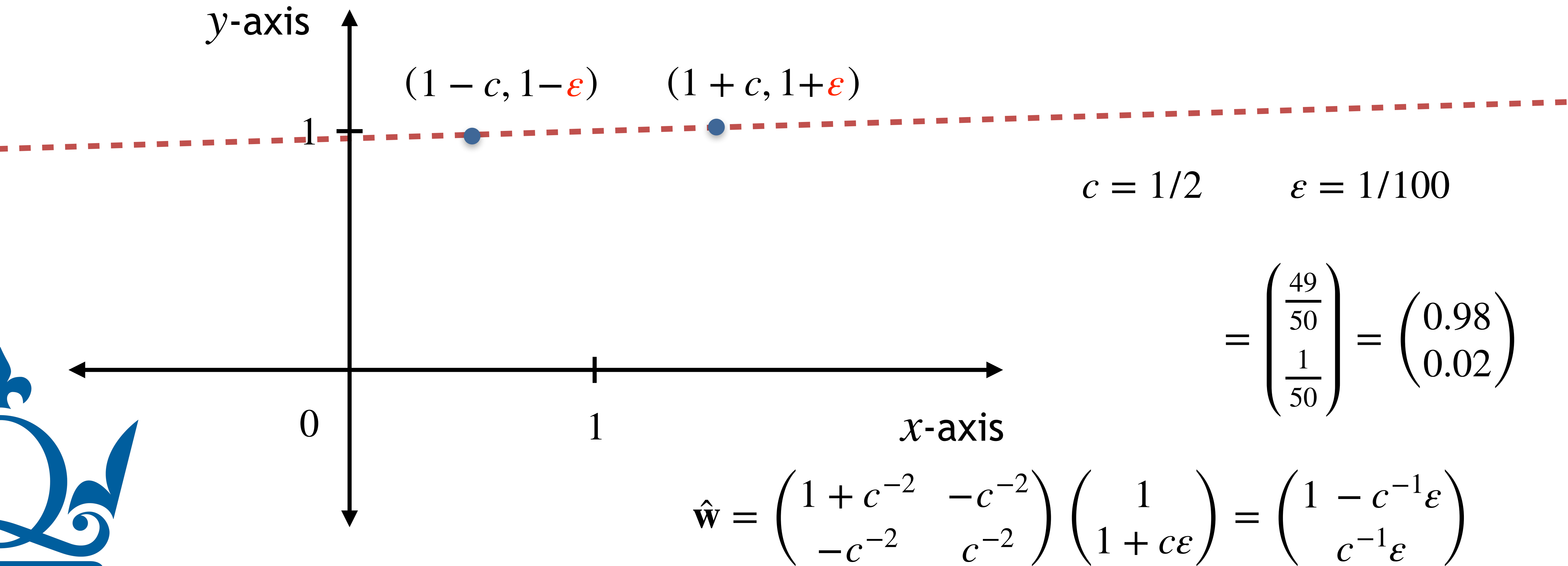
$$\mathbf{X}^T \mathbf{X} = \begin{pmatrix} 2 & 2 \\ 2 & 2 + 2c^2 \end{pmatrix} \quad \text{and} \quad \mathbf{X}^T \mathbf{y} = \begin{pmatrix} 2 \\ 2 + 2c\varepsilon \end{pmatrix}$$

and compute  $\hat{\mathbf{w}}$  via

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + c^2 \end{pmatrix} \hat{\mathbf{w}} = \begin{pmatrix} 1 \\ 1 + c\varepsilon \end{pmatrix} \implies \hat{\mathbf{w}} = \begin{pmatrix} 1 + c^{-2} & -c^{-2} \\ -c^{-2} & c^{-2} \end{pmatrix} \begin{pmatrix} 1 \\ 1 + c\varepsilon \end{pmatrix} = \begin{pmatrix} 1 - c^{-1}\varepsilon \\ c^{-1}\varepsilon \end{pmatrix}$$

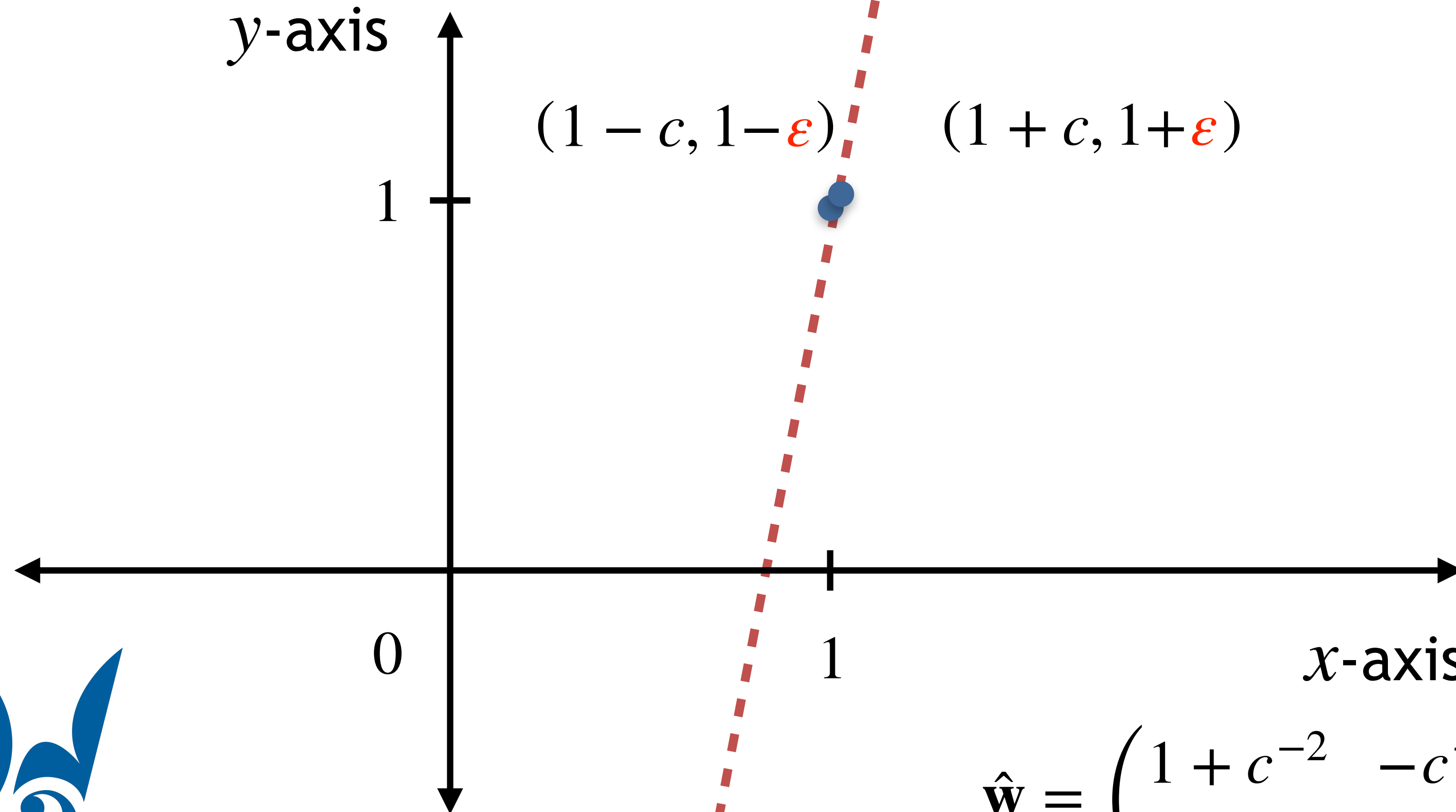
# Unstable regression problems

What if we make a small error when measuring  $y_1$  and  $y_2$ ?



# Unstable regression problems

What if we make a small error when measuring  $y_1$  and  $y_2$ ?



$$c = 1/1000 \quad \varepsilon = 1/100$$

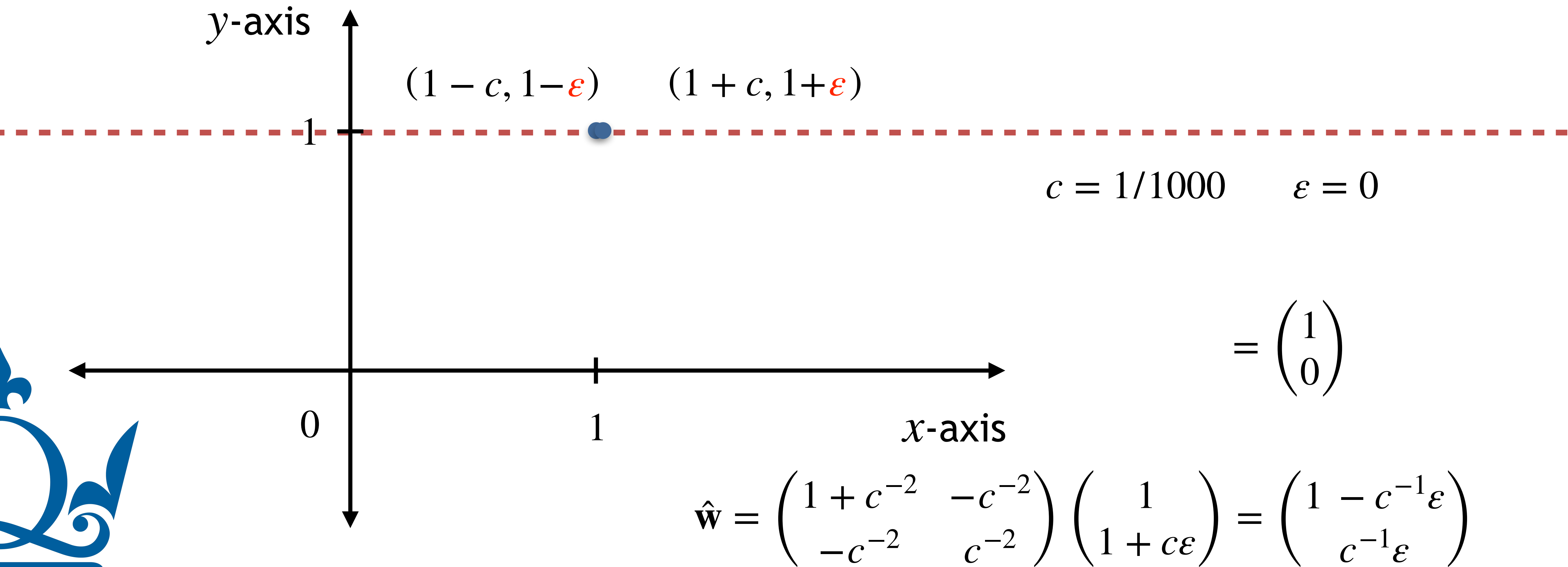
$$= \begin{pmatrix} -9 \\ 10 \end{pmatrix}$$

$$\hat{\mathbf{w}} = \begin{pmatrix} 1 + c^{-2} & -c^{-2} \\ -c^{-2} & c^{-2} \end{pmatrix} \begin{pmatrix} 1 \\ 1 + c\varepsilon \end{pmatrix} = \begin{pmatrix} 1 - c^{-1}\varepsilon \\ c^{-1}\varepsilon \end{pmatrix}$$



# Unstable regression problems

What if we make a small error when measuring  $y_1$  and  $y_2$ ?



# Unstable regression problems

Polynomial regression

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \text{MSE}(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2s} \|\Phi(X)\mathbf{w} - \mathbf{y}\|^2 \right\}$$

with

$$\Phi(X) = \begin{pmatrix} 1 & x_1 & \cdots & x_1^d \\ 1 & x_2 & \cdots & x_2^d \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_s & \cdots & x_s^d \end{pmatrix}$$



# Unstable regression problems

Polynomial regression

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \text{MSE}(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2s} \|\Phi(X)\mathbf{w} - \mathbf{y}\|^2 \right\}$$

with

$$\Phi(X) = \begin{pmatrix} 1 & x_1 & \dots & x_1^d \\ 1 & x_2 & \dots & x_2^d \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_s & \dots & x_s^d \end{pmatrix}$$

Solution:

$$\Phi(X)^\top \Phi(X) \hat{\mathbf{w}} = \Phi(X)^\top \mathbf{y}$$



# Unstable regression problems

Polynomial regression

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \text{MSE}(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2s} \|\Phi(\mathbf{X})\mathbf{w} - \mathbf{y}\|^2 \right\}$$

with

$$\Phi(\mathbf{X}) = \begin{pmatrix} 1 & x_1 & \dots & x_1^d \\ 1 & x_2 & \dots & x_2^d \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_s & \dots & x_s^d \end{pmatrix}$$

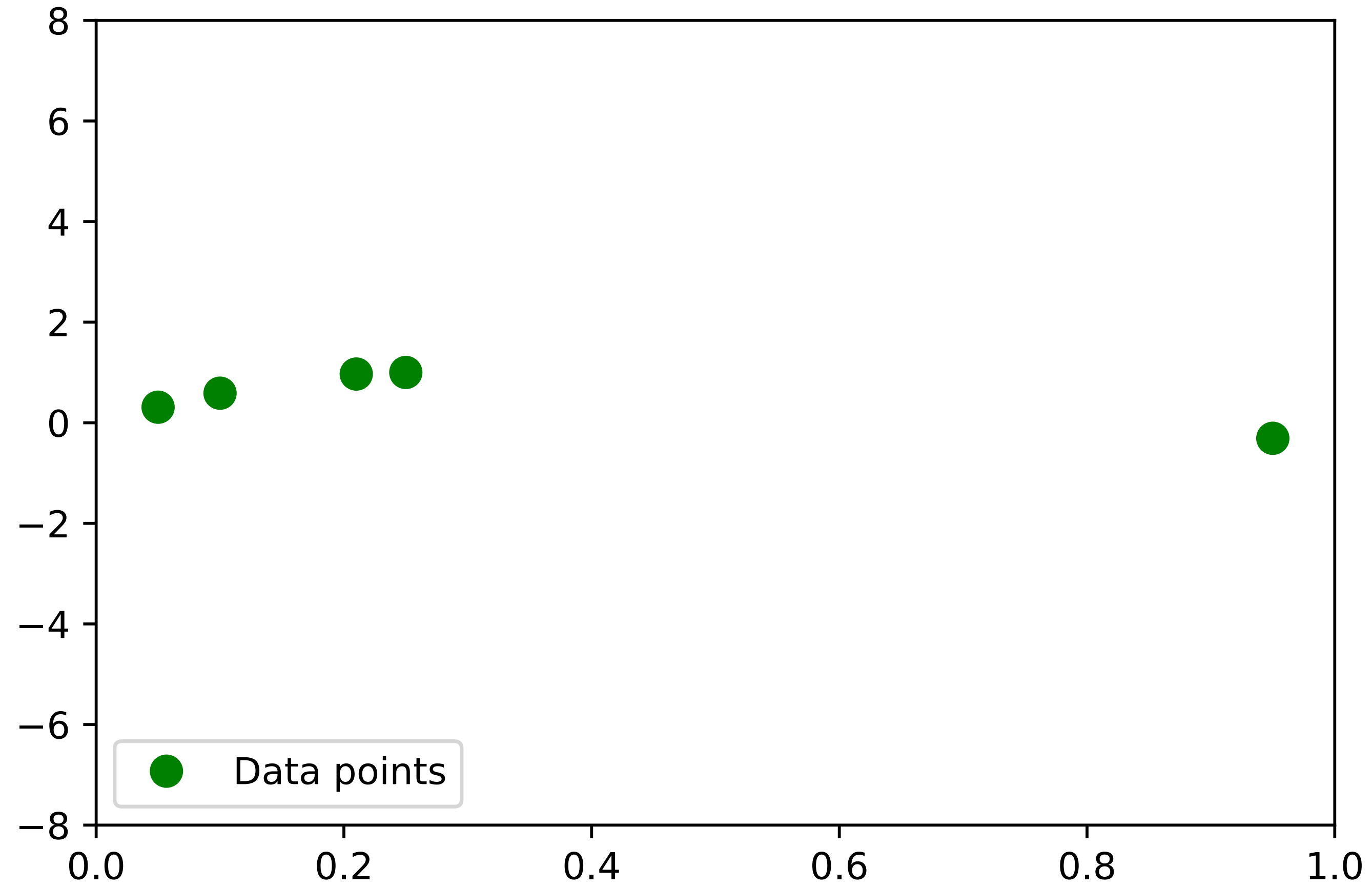
Solution:

$$\Phi(\mathbf{X})^\top \Phi(\mathbf{X}) \hat{\mathbf{w}} = \Phi(\mathbf{X})^\top \mathbf{y}$$

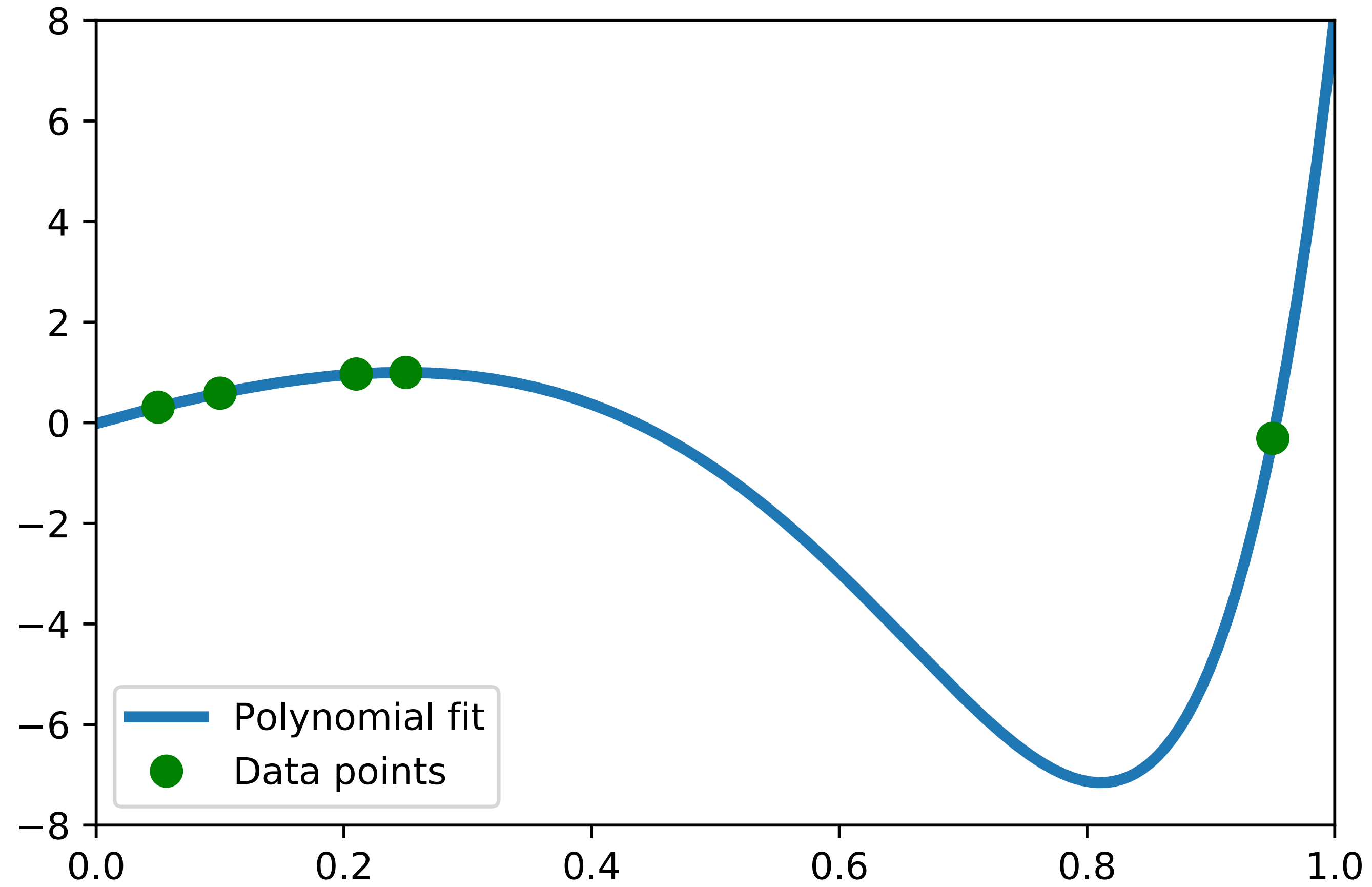
Entries of  $\Phi(\mathbf{X})^\top \Phi(\mathbf{X})$ :

$$(\Phi(\mathbf{X})^\top \Phi(\mathbf{X}))_{jk} = \sum_{i=1}^s x_i^{j+k-2}$$

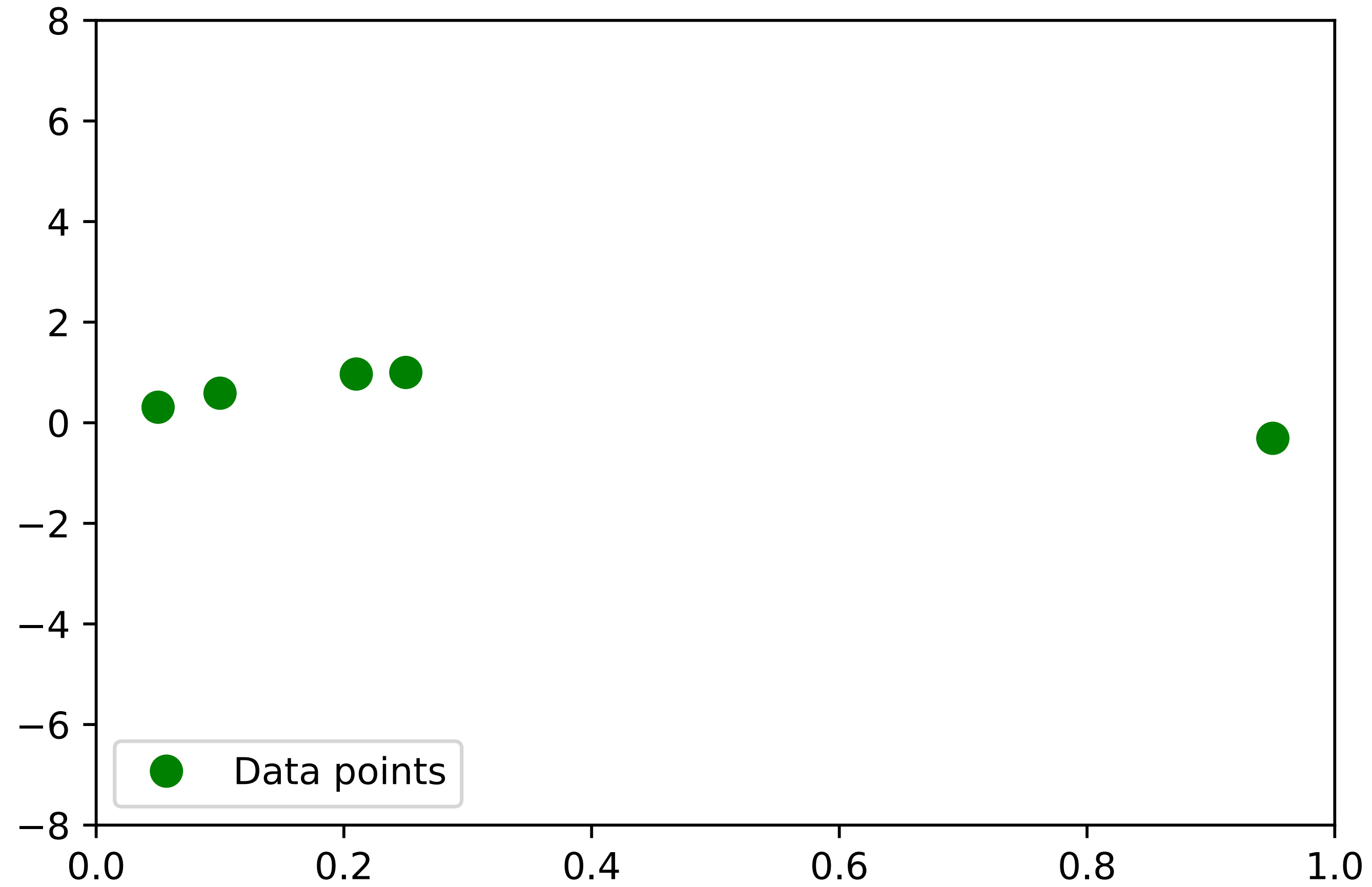
# Unstable regression problems



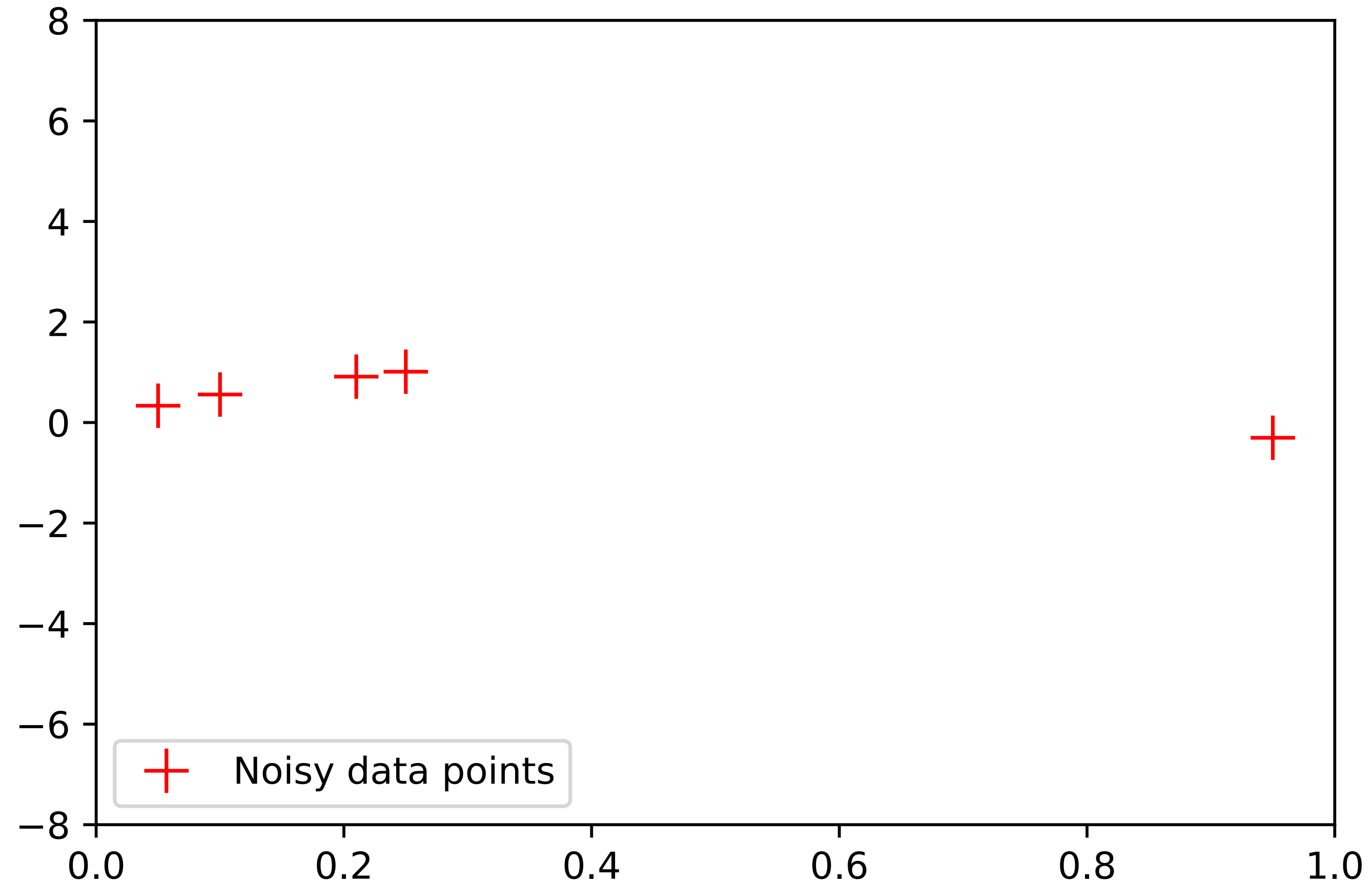
# Unstable regression problems



# Unstable regression problems

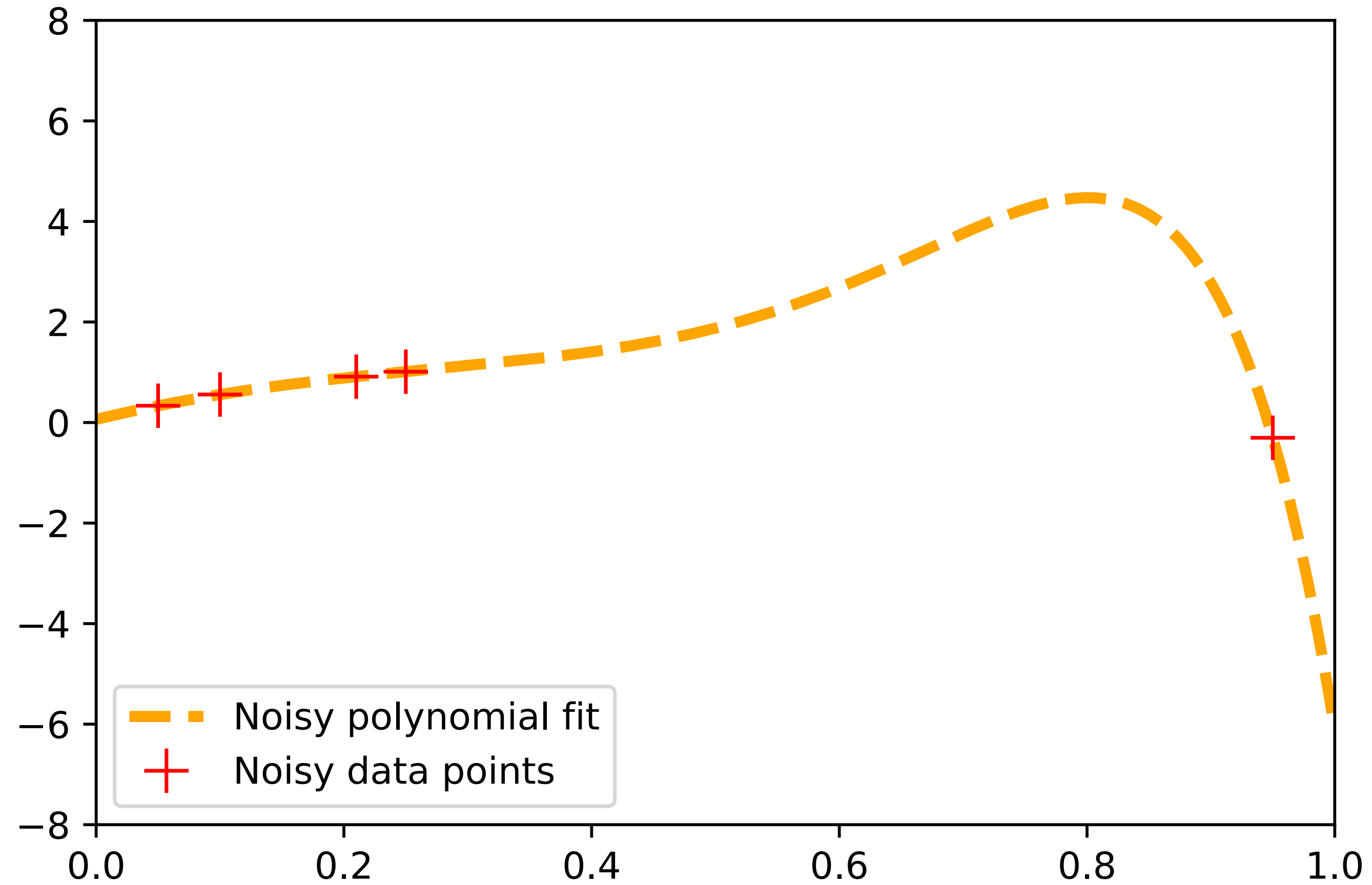


# Unstable regression problems

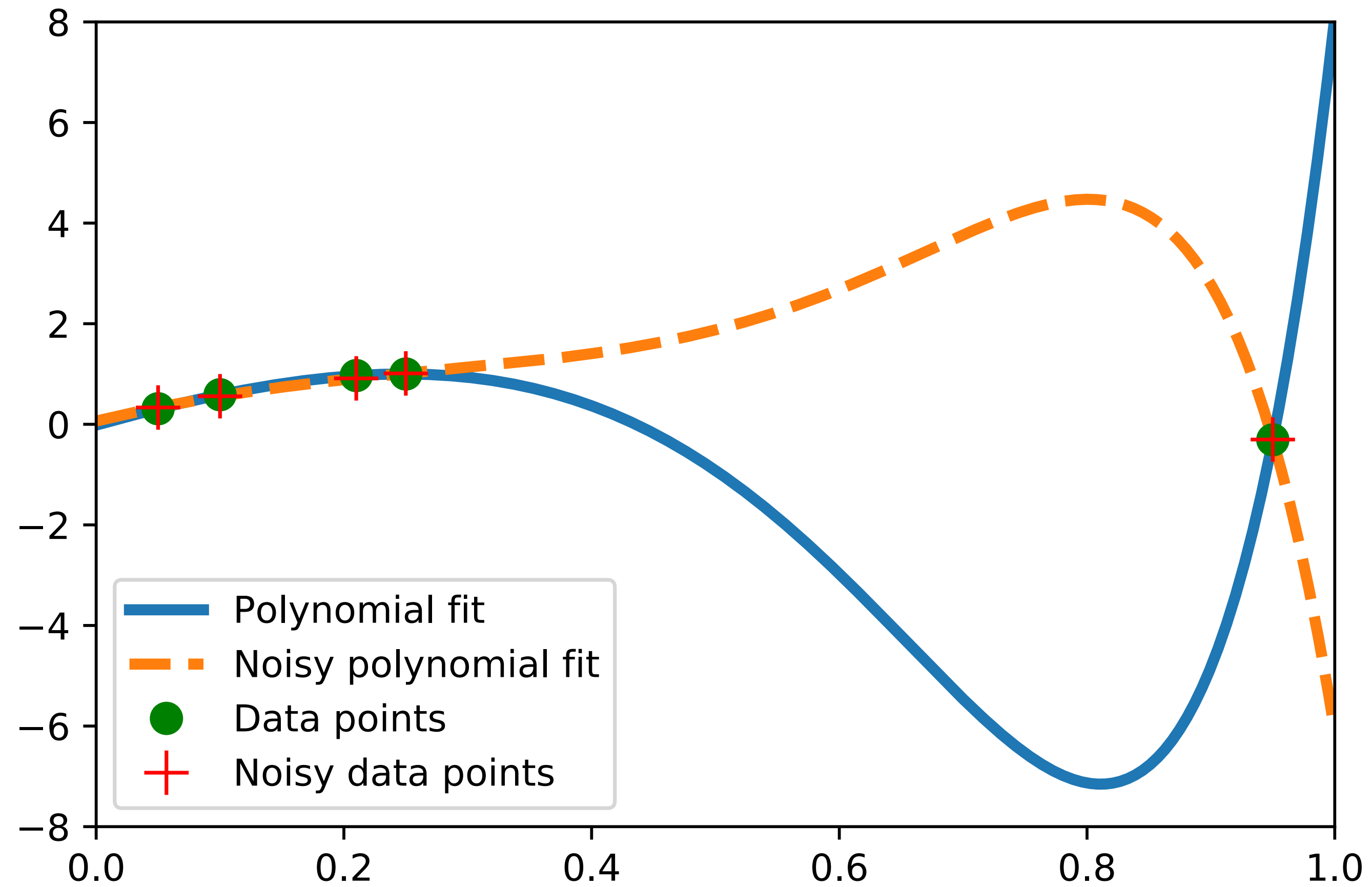




# Unstable regression problems



# Unstable regression problems



# Stability

We have discovered two phenomena:



# Stability

We have discovered two phenomena:

1. Having two slightly different outputs  $\mathbf{y}$  and  $\mathbf{y}_\delta$  with  $\|\mathbf{y} - \mathbf{y}_\delta\| \leq \delta$ , we have seen examples where  $\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\| \gg \delta$ , for

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2} \|\Phi(\mathbf{X})\mathbf{w} - \mathbf{y}\|^2 \right\} \quad \text{and} \quad \hat{\mathbf{w}}_\delta = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2} \|\Phi(\mathbf{X})\mathbf{w} - \mathbf{y}_\delta\|^2 \right\}$$



# Stability

We have discovered two phenomena:

1. Having two slightly different outputs  $\mathbf{y}$  and  $\mathbf{y}_\delta$  with  $\|\mathbf{y} - \mathbf{y}_\delta\| \leq \delta$ , we have seen examples where  $\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\| \gg \delta$ , for

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2} \|\Phi(\mathbf{X})\mathbf{w} - \mathbf{y}\|^2 \right\} \quad \text{and} \quad \hat{\mathbf{w}}_\delta = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2} \|\Phi(\mathbf{X})\mathbf{w} - \mathbf{y}_\delta\|^2 \right\}$$

2. Having two instances  $(\mathbf{X}_t, \mathbf{y}_t)$  and  $(\mathbf{X}_v, \mathbf{y}_v)$  of data samples, we have seen examples where  $\|\Phi(\mathbf{X}_v) \hat{\mathbf{w}} - \mathbf{y}_v\| \gg \|\Phi(\mathbf{X}_t) \hat{\mathbf{w}} - \mathbf{y}_t\|$ , for

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2} \|\Phi(\mathbf{X}_t) \mathbf{w} - \mathbf{y}_t\|^2 \right\}$$

# Stability

For different values of  $c$  and  $\varepsilon$  we observed

$$c > 0$$

$$\varepsilon = 0$$

$$\hat{\mathbf{w}} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$



# Stability

For different values of  $c$  and  $\varepsilon$  we observed

$$c > 0$$

$$\varepsilon = 0$$

$$\hat{\mathbf{w}} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$c = 1/2$$

$$\varepsilon = 1/100$$

$$\hat{\mathbf{w}}_{\delta} = \begin{pmatrix} \frac{49}{50} \\ \frac{1}{50} \end{pmatrix} = \begin{pmatrix} 0.98 \\ 0.02 \end{pmatrix}$$



# Stability

For different values of  $c$  and  $\varepsilon$  we observed

$$c > 0$$
$$\varepsilon = 0$$

$$\hat{\mathbf{w}} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$c = 1/2$$
$$\varepsilon = 1/100$$

$$\hat{\mathbf{w}}_{\delta} = \begin{pmatrix} \frac{49}{50} \\ \frac{1}{50} \end{pmatrix} = \begin{pmatrix} 0.98 \\ 0.02 \end{pmatrix}$$

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_{\delta}\| = \frac{1}{\sqrt{1250}}$$
$$\approx 0.0283$$



# Stability

For different values of  $c$  and  $\varepsilon$  we observed

$$c > 0$$
$$\varepsilon = 0$$

$$\hat{\mathbf{w}} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$c = 1/2$$
$$\varepsilon = 1/100$$

$$\hat{\mathbf{w}}_\delta = \begin{pmatrix} \frac{49}{50} \\ \frac{1}{50} \end{pmatrix} = \begin{pmatrix} 0.98 \\ 0.02 \end{pmatrix}$$

$$c = 1/1000$$
$$\varepsilon = 1/100$$

$$\hat{\mathbf{w}}_\delta = \begin{pmatrix} -9 \\ 10 \end{pmatrix}$$

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\| = \frac{1}{\sqrt{1250}}$$
$$\approx 0.0283$$

# Stability

For different values of  $c$  and  $\varepsilon$  we observed

$$c > 0$$
$$\varepsilon = 0$$

$$\hat{\mathbf{w}} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$c = 1/2$$
$$\varepsilon = 1/100$$

$$\hat{\mathbf{w}}_\delta = \begin{pmatrix} \frac{49}{50} \\ \frac{1}{50} \end{pmatrix} = \begin{pmatrix} 0.98 \\ 0.02 \end{pmatrix}$$

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\| = \frac{1}{\sqrt{1250}}$$
$$\approx 0.0283$$

$$c = 1/1000$$
$$\varepsilon = 1/100$$

$$\hat{\mathbf{w}}_\delta = \begin{pmatrix} -9 \\ 10 \end{pmatrix}$$

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\| = \sqrt{200}$$
$$\approx 14.14$$

# Stability

Let us characterise the error  $\|\hat{w} - \hat{w}_\delta\|$  for general problems



# Stability

Let us characterise the error  $\|\hat{w} - \hat{w}_\delta\|$  for general problems

$$\hat{w} = \arg \min_{w \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2} \|\Phi(X)w - y\|^2 \right\} \quad \text{and} \quad \hat{w}_\delta = \arg \min_{w \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2} \|\Phi(X)w - y_\delta\|^2 \right\}$$



# Stability

The solutions for

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2} \|\Phi(\mathbf{X})\mathbf{w} - \mathbf{y}\|^2 \right\} \quad \text{and} \quad \hat{\mathbf{w}}_\delta = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2} \|\Phi(\mathbf{X})\mathbf{w} - \mathbf{y}_\delta\|^2 \right\}$$

are

$$\hat{\mathbf{w}} = (\Phi(\mathbf{X})^\top \Phi(\mathbf{X}))^{-1} \Phi(\mathbf{X})^\top \mathbf{y} \quad \text{and} \quad \hat{\mathbf{w}}_\delta = (\Phi(\mathbf{X})^\top \Phi(\mathbf{X}))^{-1} \Phi(\mathbf{X})^\top \mathbf{y}_\delta$$



# Stability

The matrix  $\Phi(X)^T \Phi(X)$  is a symmetric, positive definite matrix with real entries and

$$\hat{w} = (\Phi(X)^T \Phi(X))^{-1} \Phi(X)^T y \quad \text{and} \quad \hat{w}_\delta = (\Phi(X)^T \Phi(X))^{-1} \Phi(X)^T y_\delta$$

Using the SVD we can prove that (see notes)

$$\hat{w} = V(\Sigma^T)^{-1} U^T y \quad \text{and} \quad \hat{w}_\delta = V(\Sigma^T)^{-1} U^T y_\delta$$



# Stability

Hence

$$\hat{\mathbf{w}} - \hat{\mathbf{w}}_{\delta} = \mathbf{V}(\mathbf{\Sigma}^{\top})^{-1}\mathbf{U}^{\top}(\mathbf{y} - \mathbf{y}_{\delta})$$



# Stability

Hence

$$\hat{\mathbf{w}} - \hat{\mathbf{w}}_{\delta} = \mathbf{V}(\mathbf{\Sigma}^{\top})^{-1}\mathbf{U}^{\top}(\mathbf{y} - \mathbf{y}_{\delta})$$

We can then calculate





# Stability

Hence

$$\hat{\mathbf{w}} - \hat{\mathbf{w}}_{\delta} = \mathbf{V}(\mathbf{\Sigma}^{\top})^{-1}\mathbf{U}^{\top}(\mathbf{y} - \mathbf{y}_{\delta})$$

We can then calculate

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_{\delta}\|^2 = \|\mathbf{V}(\mathbf{\Sigma}^{\top})^{-1}\mathbf{U}^{\top}(\mathbf{y} - \mathbf{y}_{\delta})\|^2$$



# Stability

Hence

$$\hat{\mathbf{w}} - \hat{\mathbf{w}}_{\delta} = \mathbf{V}(\boldsymbol{\Sigma}^{\top})^{-1}\mathbf{U}^{\top}(\mathbf{y} - \mathbf{y}_{\delta})$$

We can then calculate

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_{\delta}\|^2 = \|\mathbf{V}(\boldsymbol{\Sigma}^{\top})^{-1}\mathbf{U}^{\top}(\mathbf{y} - \mathbf{y}_{\delta})\|^2 \leq \|\mathbf{V}\|^2\|(\boldsymbol{\Sigma}^{\top})^{-1}\|^2\|\mathbf{U}^{\top}\|^2\|\mathbf{y} - \mathbf{y}_{\delta}\|^2$$

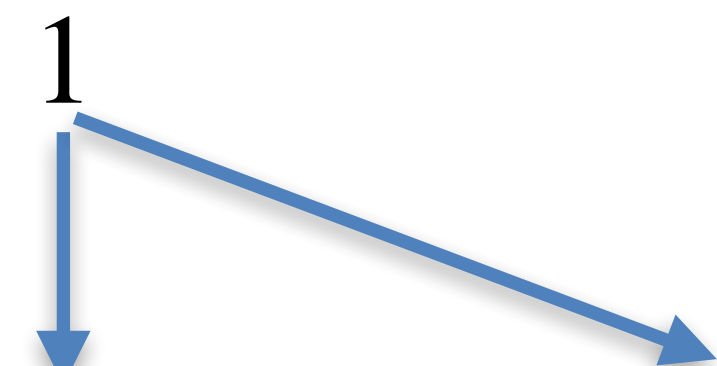


# Stability

Hence

$$\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta = \mathbf{V}(\boldsymbol{\Sigma}^\top)^{-1}\mathbf{U}^\top(\mathbf{y} - \mathbf{y}_\delta)$$

We can then calculate

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 = \|\mathbf{V}(\boldsymbol{\Sigma}^\top)^{-1}\mathbf{U}^\top(\mathbf{y} - \mathbf{y}_\delta)\|^2 \leq \overset{1}{\| \mathbf{V} \|^2 \| (\boldsymbol{\Sigma}^\top)^{-1} \|^2 \| \mathbf{U}^\top \|^2} \|\mathbf{y} - \mathbf{y}_\delta\|^2$$


# Stability

Hence

$$\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta = \mathbf{V}(\boldsymbol{\Sigma}^\top)^{-1}\mathbf{U}^\top(\mathbf{y} - \mathbf{y}_\delta)$$

We can then calculate

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 = \|\mathbf{V}(\boldsymbol{\Sigma}^\top)^{-1}\mathbf{U}^\top(\mathbf{y} - \mathbf{y}_\delta)\|^2 \leq \overset{1}{\| \mathbf{V} \|^2 \| (\boldsymbol{\Sigma}^\top)^{-1} \|^2 \| \mathbf{U}^\top \|^2} \|\mathbf{y} - \mathbf{y}_\delta\|^2$$

Cauchy-Schwarz

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|$$

# Stability

Hence

$$\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta = \mathbf{V}(\boldsymbol{\Sigma}^\top)^{-1}\mathbf{U}^\top(\mathbf{y} - \mathbf{y}_\delta)$$

We can then calculate

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 = \|\mathbf{V}(\boldsymbol{\Sigma}^\top)^{-1}\mathbf{U}^\top(\mathbf{y} - \mathbf{y}_\delta)\|^2 \leq \overset{1}{\|\mathbf{V}\|^2\|(\boldsymbol{\Sigma}^\top)^{-1}\|^2\|\mathbf{U}^\top\|^2\|\mathbf{y} - \mathbf{y}_\delta\|^2}$$

Cauchy-Schwarz

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|$$

# Stability

Hence

$$\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta = \mathbf{V}(\boldsymbol{\Sigma}^\top)^{-1}\mathbf{U}^\top(\mathbf{y} - \mathbf{y}_\delta)$$

We can then calculate

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 = \|\mathbf{V}(\boldsymbol{\Sigma}^\top)^{-1}\mathbf{U}^\top(\mathbf{y} - \mathbf{y}_\delta)\|^2 \leq \overset{1}{\|\mathbf{V}\|^2\|(\boldsymbol{\Sigma}^\top)^{-1}\|^2\|\mathbf{U}^\top\|^2\|\mathbf{y} - \mathbf{y}_\delta\|^2}$$

def. of norm + diagonal matrix

Cauchy-Schwarz

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|$$

# Stability

Hence

$$\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta = \mathbf{V}(\boldsymbol{\Sigma}^\top)^{-1}\mathbf{U}^\top(\mathbf{y} - \mathbf{y}_\delta)$$

We can then calculate

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 = \|\mathbf{V}(\boldsymbol{\Sigma}^\top)^{-1}\mathbf{U}^\top(\mathbf{y} - \mathbf{y}_\delta)\|^2 \leq \overset{1}{\|\mathbf{V}\|^2} \|\boldsymbol{\Sigma}^\top\|^{-2} \|\mathbf{U}^\top\|^2 \|\mathbf{y} - \mathbf{y}_\delta\|^2$$

$= \frac{\delta^2}{\sigma_{d+1}^2}$  def. of norm + diagonal matrix

Cauchy-Schwarz

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|$$

# Stability

We just showed how

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 \leq \frac{\delta^2}{\sigma_{d+1}^2}$$





# Stability

We just showed how

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 \leq \frac{\delta^2}{\sigma_{d+1}^2}$$

In the worst case the deviation is amplified by the smallest singular value



# Stability

From

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_{\delta}\|^2 \leq \|(\boldsymbol{\Sigma}^{\top})^{-1}\|^2 \|\mathbf{y} - \mathbf{y}_{\delta}\|^2$$



# Stability

From

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 \leq \|(\boldsymbol{\Sigma}^\top)^{-1}\|^2 \|\mathbf{y} - \mathbf{y}_\delta\|^2$$

If we assume that  $\mathbf{y} = \boldsymbol{\Phi}\mathbf{r}$  and  $\mathbf{y}_\delta = \boldsymbol{\Phi}\mathbf{p}$



# Stability

From

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 \leq \|(\boldsymbol{\Sigma}^\top)^{-1}\|^2 \|\mathbf{y} - \mathbf{y}_\delta\|^2$$

If we assume that  $\mathbf{y} = \boldsymbol{\Phi}\mathbf{r}$  and  $\mathbf{y}_\delta = \boldsymbol{\Phi}\mathbf{p}$

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 \leq \|(\boldsymbol{\Sigma}^\top)^{-1}\|^2 \|\boldsymbol{\Phi}(\mathbf{r} - \mathbf{p})\|^2 \leq \|(\boldsymbol{\Sigma}^\top)^{-1}\|^2 \|\boldsymbol{\Phi}\|^2 \|\mathbf{r} - \mathbf{p}\|^2$$



# Stability

From

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 \leq \|(\boldsymbol{\Sigma}^\top)^{-1}\|^2 \|\mathbf{y} - \mathbf{y}_\delta\|^2$$

If we assume that  $\mathbf{y} = \boldsymbol{\Phi}\mathbf{r}$  and  $\mathbf{y}_\delta = \boldsymbol{\Phi}\mathbf{p}$

$$\begin{aligned} \|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 &\leq \|(\boldsymbol{\Sigma}^\top)^{-1}\|^2 \|\boldsymbol{\Phi}(\mathbf{r} - \mathbf{p})\|^2 \leq \|(\boldsymbol{\Sigma}^\top)^{-1}\|^2 \|\boldsymbol{\Phi}\|^2 \|\mathbf{r} - \mathbf{p}\|^2 \\ &= \frac{\sigma_1^2}{\sigma_{d+1}^2} \epsilon^2 = \kappa(\boldsymbol{\Phi})^2 \epsilon^2 \end{aligned}$$



# Stability

From

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 \leq \|(\boldsymbol{\Sigma}^\top)^{-1}\|^2 \|\mathbf{y} - \mathbf{y}_\delta\|^2$$

If we assume that  $\mathbf{y} = \boldsymbol{\Phi}\mathbf{r}$  and  $\mathbf{y}_\delta = \boldsymbol{\Phi}\mathbf{p}$

$$\begin{aligned} \|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\|^2 &\leq \|(\boldsymbol{\Sigma}^\top)^{-1}\|^2 \|\boldsymbol{\Phi}(\mathbf{r} - \mathbf{p})\|^2 \leq \|(\boldsymbol{\Sigma}^\top)^{-1}\|^2 \|\boldsymbol{\Phi}\|^2 \|\mathbf{r} - \mathbf{p}\|^2 \\ &= \frac{\sigma_1^2}{\sigma_{d+1}^2} \epsilon^2 = \kappa(\boldsymbol{\Phi})^2 \epsilon^2 \end{aligned}$$

Where  $\kappa = \frac{\sigma_1}{\sigma_{d+1}}$  is the condition number that quantifies the amplification of the error in the worst case. A matrix with large kappa is called ill-conditioned

# Stability

Back to our initial example for  $d = 1$ :

$$\Phi(\mathbf{X}) = \mathbf{X} = \begin{pmatrix} 1 & 1 - c \\ 1 & 1 + c \end{pmatrix}$$



# Stability

Back to our initial example for  $d = 1$ :

$$\Phi(\mathbf{X}) = \mathbf{X} = \begin{pmatrix} 1 & 1 - c \\ 1 & 1 + c \end{pmatrix}$$

The singular values for this matrix are

$$\sigma_1 = \sqrt{2 + c^2 + \sqrt{c^4 + 4}} \quad \text{and} \quad \sigma_2 = \sqrt{2 + c^2 - \sqrt{c^4 + 4}}$$





# Stability

Back to our initial example for  $d = 1$ :

$$\Phi(\mathbf{X}) = \mathbf{X} = \begin{pmatrix} 1 & 1 - c \\ 1 & 1 + c \end{pmatrix}$$

The singular values for this matrix are

$$\sigma_1 = \sqrt{2 + c^2 + \sqrt{c^4 + 4}} \quad \text{and} \quad \sigma_2 = \sqrt{2 + c^2 - \sqrt{c^4 + 4}}$$

$\Rightarrow$

$$\kappa = \frac{\sigma_1}{\sigma_2} = \frac{2 + c^2 + \sqrt{c^4 + 4}}{2c}$$



# Stability

$$\kappa = \frac{2 + c^2 + \sqrt{c^4 + 4}}{2c}$$

$$\hat{\mathbf{w}} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$



# Stability

$$\kappa = \frac{2 + c^2 + \sqrt{c^4 + 4}}{2c}$$

$$\hat{\mathbf{w}} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$c = 1/2$$

$$\varepsilon = 1/100$$

$$\hat{\mathbf{w}}_\delta = \begin{pmatrix} \frac{49}{50} \\ \frac{1}{50} \end{pmatrix} = \begin{pmatrix} 0.98 \\ 0.02 \end{pmatrix}$$

$$\kappa \approx 4.3$$

# Stability

$$\kappa = \frac{2 + c^2 + \sqrt{c^4 + 4}}{2c}$$

$$\hat{\mathbf{w}} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$c = 1/2$$

$$\varepsilon = 1/100$$

$$\hat{\mathbf{w}}_\delta = \begin{pmatrix} \frac{49}{50} \\ \frac{1}{50} \end{pmatrix} = \begin{pmatrix} 0.98 \\ 0.02 \end{pmatrix}$$

$$\kappa \approx 4.3$$

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\| \leq \kappa \varepsilon \approx 0.04$$

# Stability

$$\kappa = \frac{2 + c^2 + \sqrt{c^4 + 4}}{2c}$$

$$\hat{\mathbf{w}} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$c = 1/2$$
$$\varepsilon = 1/100$$

$$\hat{\mathbf{w}}_\delta = \begin{pmatrix} \frac{49}{50} \\ \frac{1}{50} \end{pmatrix} = \begin{pmatrix} 0.98 \\ 0.02 \end{pmatrix}$$

$$\kappa \approx 4.3$$

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\| \leq \kappa \varepsilon \approx 0.04$$

$$c = 1/1000$$
$$\varepsilon = 1/100$$

$$\hat{\mathbf{w}}_\delta = \begin{pmatrix} -9 \\ 10 \end{pmatrix}$$

$$\kappa \approx 2000$$

# Stability

$$\kappa = \frac{2 + c^2 + \sqrt{c^4 + 4}}{2c}$$

$$\hat{\mathbf{w}} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$c = 1/2$$
$$\varepsilon = 1/100$$

$$\hat{\mathbf{w}}_\delta = \begin{pmatrix} \frac{49}{50} \\ \frac{1}{50} \end{pmatrix} = \begin{pmatrix} 0.98 \\ 0.02 \end{pmatrix}$$

$$\kappa \approx 4.3$$

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\| \leq \kappa \varepsilon \approx 0.04$$

$$c = 1/1000$$
$$\varepsilon = 1/100$$

$$\hat{\mathbf{w}}_\delta = \begin{pmatrix} -9 \\ 10 \end{pmatrix}$$

$$\kappa \approx 2000$$

$$\|\hat{\mathbf{w}} - \hat{\mathbf{w}}_\delta\| \leq \kappa \varepsilon \approx 20$$



# REGULARISATION METHODS

# Ill-conditioned problems

What can we do in order to be less sensitive towards measurement errors?

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$$





# Ill-conditioned problems

What can we do in order to be less sensitive towards measurement errors?

$$(\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad \alpha > 0$$

We can add a multiple of the identity matrix  $\rightarrow$  shift of the singular values!



# Ill-conditioned problems

What can we do in order to be less sensitive towards measurement errors?

$$(\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad \alpha > 0$$

We can add a multiple of the identity matrix  $\rightarrow$  shift of the singular values!

$$\kappa = \sqrt{\frac{\sigma_1^2}{\sigma_{d+1}^2}} \quad \Rightarrow \quad \kappa = \sqrt{\frac{\sigma_1^2 + \alpha}{\sigma_{d+1}^2 + \alpha}}$$



# Ill-conditioned problems

What can we do in order to be less sensitive towards measurement errors?

$$(\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad \alpha > 0$$

We can add a multiple of the identity matrix  $\rightarrow$  shift of the singular values!

$$\kappa = \sqrt{\frac{\sigma_1^2}{\sigma_{d+1}^2}} \quad \Rightarrow \quad \kappa = \sqrt{\frac{\sigma_1^2 + \alpha}{\sigma_{d+1}^2 + \alpha}}$$

Example:

$$\begin{aligned} \sigma_1 &= \sqrt{2} \\ \sigma_{d+1} &= 1/\sqrt{2000000} \\ \alpha &= 1 \end{aligned}$$

# Ill-conditioned problems

What can we do in order to be less sensitive towards measurement errors?

$$(\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad \alpha > 0$$

We can add a multiple of the identity matrix  $\rightarrow$  shift of the singular values!

$$\kappa = \sqrt{\frac{\sigma_1^2}{\sigma_{d+1}^2}} \quad \Rightarrow \quad \kappa = \sqrt{\frac{\sigma_1^2 + \alpha}{\sigma_{d+1}^2 + \alpha}}$$

Example:

$$\begin{aligned} \sigma_1 &= \sqrt{2} \\ \sigma_{d+1} &= 1/\sqrt{2000000} \\ \alpha &= 1 \end{aligned} \quad \Rightarrow \quad \kappa \approx \sqrt{3} \ll 2000$$



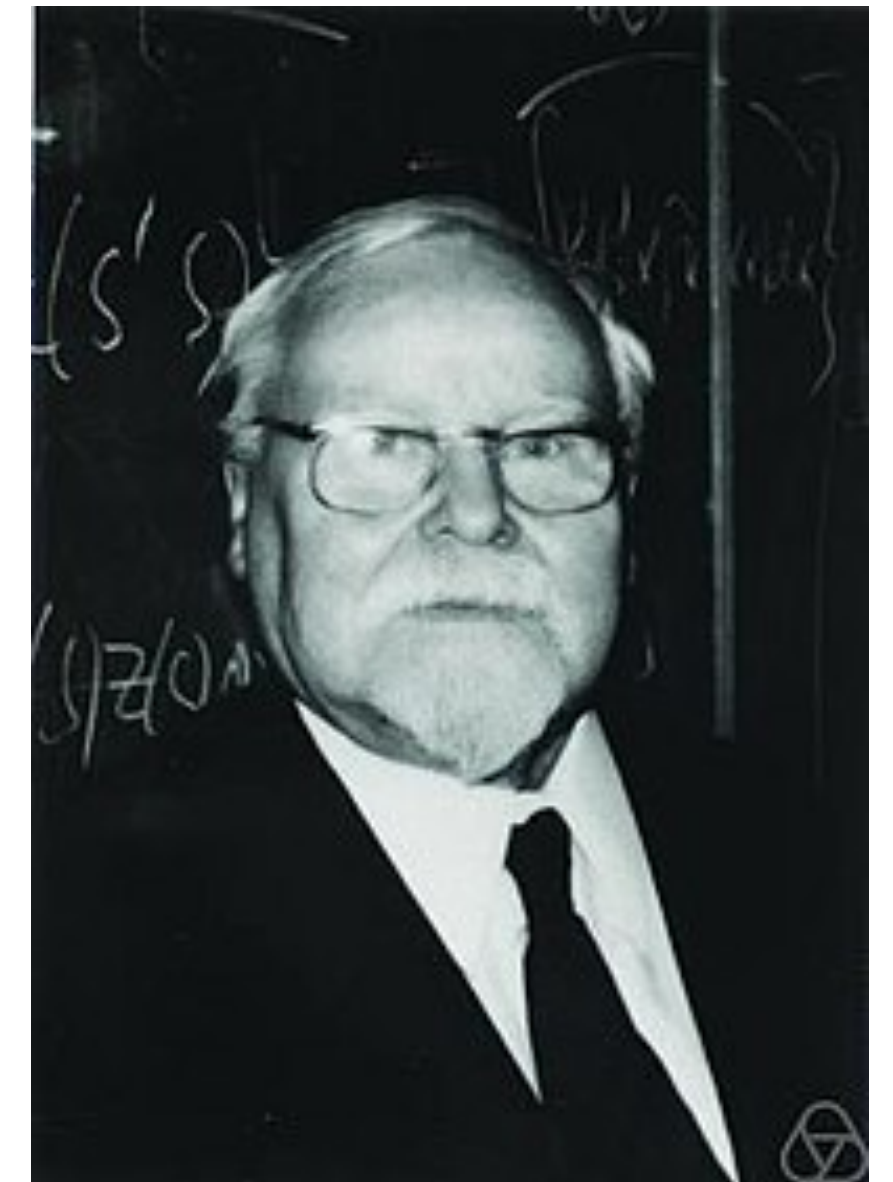
# REGULARISATION

# Ridge regression / Tikhonov regularisation

The minimisation problem

$$\hat{w} = \arg \min_w \left\{ \frac{1}{2} \|Xw - y\|^2 + \frac{\alpha}{2} \|w\|^2 \right\}$$

is also known as *Tikhonov regularisation*  
or *ridge regression*



Andrey Tikhonov, 1906 - 1993



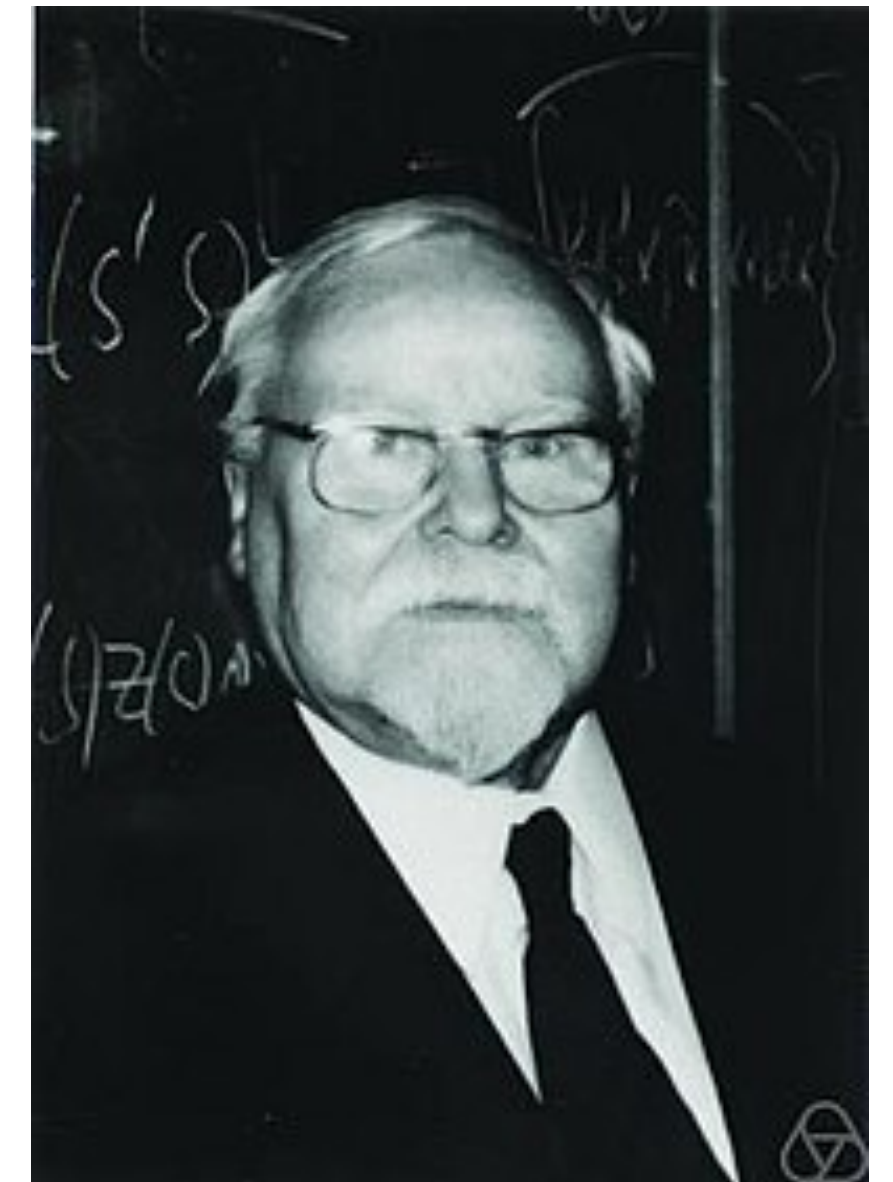
# Ridge regression / Tikhonov regularisation

The minimisation problem

$$\hat{w} = \arg \min_w \left\{ \frac{1}{2} \|Xw - y\|^2 + \frac{\alpha}{2} \|w\|^2 \right\}$$

Standard regression term

is also known as *Tikhonov regularisation*  
or *ridge regression*



Andrey Tikhonov, 1906 - 1993



# Ridge regression / Tikhonov regularisation

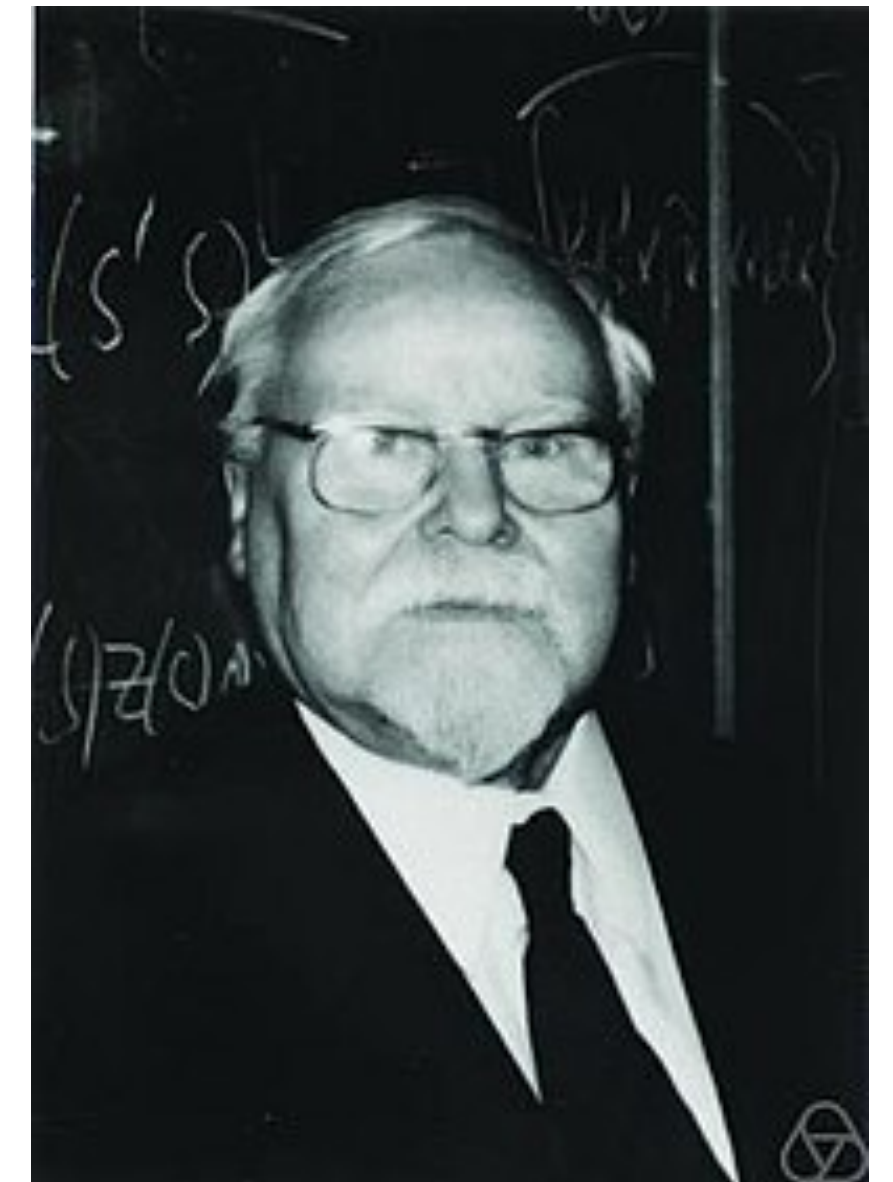
The minimisation problem

$$\hat{w} = \arg \min_w \left\{ \frac{1}{2} \|Xw - y\|^2 + \frac{\alpha}{2} \|w\|^2 \right\}$$

Standard regression term

Regularisation term

is also known as *Tikhonov regularisation*  
or *ridge regression*



Andrey Tikhonov, 1906 - 1993



# Ridge regression / Tikhonov regularisation

The minimisation problem

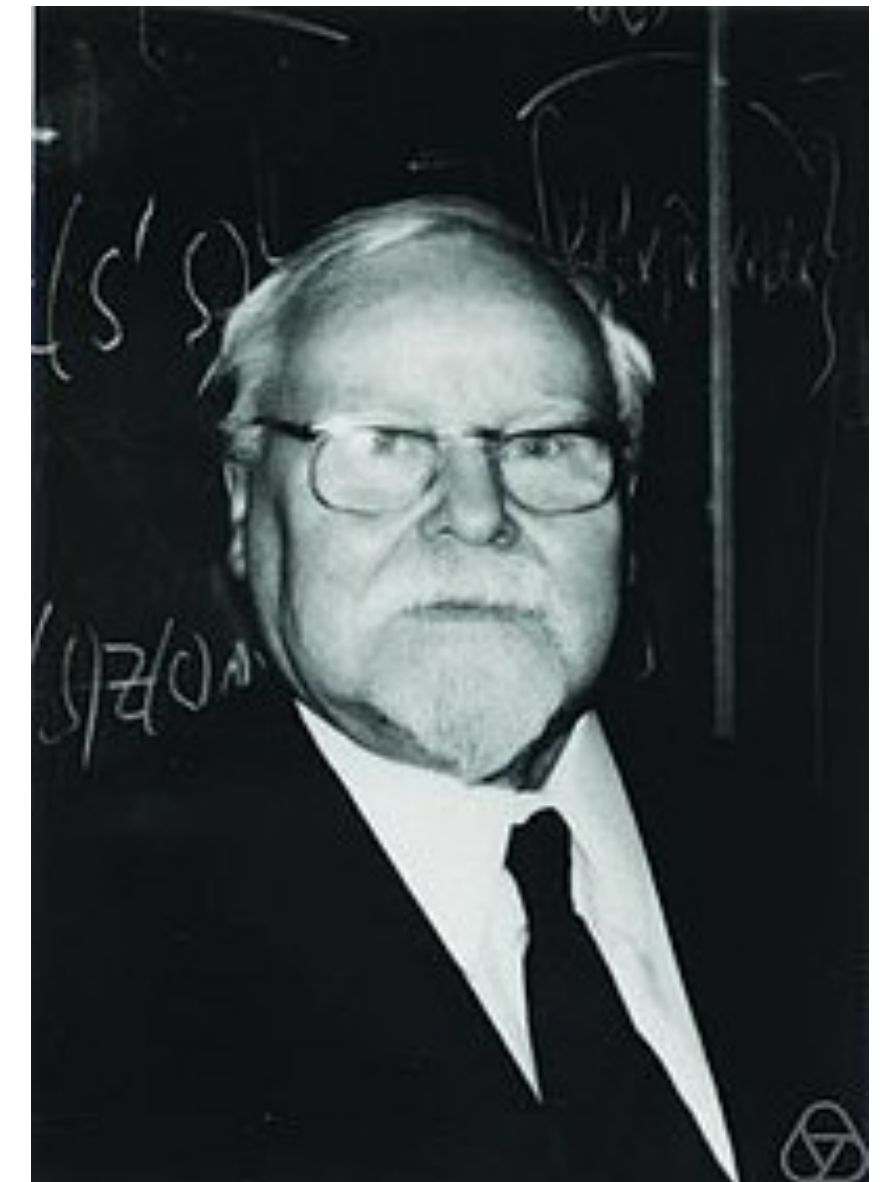
$$\hat{w} = \arg \min_w \left\{ \frac{1}{2} \|Xw - y\|^2 + \frac{\alpha}{2} \|w\|^2 \right\}$$

Standard regression term

Regularisation term

Regularisation parameter

is also known as *Tikhonov regularisation*  
or *ridge regression*



Andrey Tikhonov, 1906 - 1993

# Variational regularisation

A more general form of the previous problem is variational regularisation

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \{L(\mathbf{w}) + R(\mathbf{w})\}$$



# Variational regularisation

A more general form of the previous problem is variational regularisation

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \{L(\mathbf{w}) + R(\mathbf{w})\}$$

Data term/  
Regression term



# Variational regularisation

A more general form of the previous problem is variational regularisation

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \{L(\mathbf{w}) + R(\mathbf{w})\}$$

Data term/  
Regression term

Regularisation  
term



# Variational regularisation

A more general form of the previous problem is variational regularisation

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \{L(\mathbf{w}) + R(\mathbf{w})\}$$

Data term/  
Regression term

Regularisation  
term

Previous example:

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

$$R(\mathbf{w}) = \frac{\alpha}{2} \|\mathbf{w}\|^2$$



# Ill-conditioned problems

Note that

$$(\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

is equivalent to

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2 \right\}$$



# Ill-conditioned problems

Note that

$$(\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

is equivalent to

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2 \right\}$$

Proof sketch:



# Ill-conditioned problems

Note that

$$(\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

is equivalent to

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2 \right\}$$

Proof sketch: 1. Compute gradient of  $E(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2$ , set it to zero

and show that this coincides with  $(\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$





# Ill-conditioned problems

Note that

$$(\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

is equivalent to

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2 \right\}$$

Proof sketch: 1. Compute gradient of  $E(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2$ , set it to zero

and show that this coincides with  $(\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$

2. Show that  $E(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2$  is convex



# Good exercise!

1. Compute gradient of  $E(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2$ , set it to zero



# Good exercise!

1. Compute gradient of  $E(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2$ , set it to zero

$$E(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2 = \frac{1}{2} \sum_i \left( \sum_j X_{ij} w_j - y_i \right)^2 + \frac{\alpha}{2} \sum_i w_i^2$$



# Good exercise!

1. Compute gradient of  $E(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2$ , set it to zero

$$E(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2 = \frac{1}{2} \sum_i \left( \sum_j X_{ij} w_j - y_i \right)^2 + \frac{\alpha}{2} \sum_i w_i^2$$

$$\partial_{w_p} E(\mathbf{w}) = \frac{1}{2} \sum_i \partial_{w_p} \left( \sum_j X_{ij} w_j - y_i \right)^2 + \frac{\alpha}{2} \sum_i \partial_{w_p} w_i^2$$



# Good exercise!

1. Compute gradient of  $E(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2$ , set it to zero

$$E(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2 = \frac{1}{2} \sum_i \left( \sum_j X_{ij} w_j - y_i \right)^2 + \frac{\alpha}{2} \sum_i w_i^2$$

$$\partial_{w_p} E(\mathbf{w}) = \frac{1}{2} \sum_i \partial_{w_p} \left( \sum_j X_{ij} w_j - y_i \right)^2 + \frac{\alpha}{2} \sum_i \partial_{w_p} w_i^2$$

$$\partial_{w_p} E(\mathbf{w}) = \sum_i \left( \sum_j X_{ij} w_j - y_i \right) X_{ip} + \alpha w_p$$



# Good exercise!

$$\partial_{w_p} E(\mathbf{w}) = \sum_i \left( \sum_j X_{ij} w_j - y_i \right) X_{ip} + \alpha w_p$$



# Good exercise!

$$\partial_{w_p} E(\mathbf{w}) = \sum_i \left( \sum_j X_{ij} w_j - y_i \right) X_{ip} + \alpha w_p$$

$$\partial_{w_p} E(\mathbf{w}) = \sum_i \left( \sum_j X_{ij} w_j - y_i \right) X_{pi}^\top + \alpha w_p$$



# Good exercise!

$$\partial_{w_p} E(\mathbf{w}) = \sum_i \left( \sum_j X_{ij} w_j - y_i \right) X_{ip} + \alpha w_p$$

$$\partial_{w_p} E(\mathbf{w}) = \sum_i \left( \sum_j X_{ij} w_j - y_i \right) X_{pi}^\top + \alpha w_p$$

$$\partial_{w_p} E(\mathbf{w}) = \sum_i \sum_j \left( X_{pi}^\top X_{ij} w_j - X_{pi}^\top y_i \right) + \alpha w_p$$





# Good exercise!

$$\partial_{w_p} E(\mathbf{w}) = \sum_i \left( \sum_j X_{ij} w_j - y_i \right) X_{ip} + \alpha w_p$$

$$\partial_{w_p} E(\mathbf{w}) = \sum_i \left( \sum_j X_{ij} w_j - y_i \right) X_{pi}^\top + \alpha w_p$$

$$\partial_{w_p} E(\mathbf{w}) = \sum_i \sum_j \left( X_{pi}^\top X_{ij} w_j - X_{pi}^\top y_i \right) + \alpha w_p$$

$$\partial_{w_p} E(\mathbf{w}) = \left( \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y} + \alpha \mathbf{w} \right)_p$$



# Good exercise!

$$\partial_{w_p} E(\mathbf{w}) = (\mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y} + \alpha \mathbf{w})_p$$



# Good exercise!

$$\partial_{w_p} E(\mathbf{w}) = (\mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y} + \alpha \mathbf{w})_p$$

$$\nabla E(\mathbf{w}) = \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y} + \alpha \mathbf{w}$$



# Good exercise!

$$\partial_{w_p} E(\mathbf{w}) = (\mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y} + \alpha \mathbf{w})_p$$

$$\nabla E(\mathbf{w}) = \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y} + \alpha \mathbf{w}$$

$$\nabla E(\mathbf{w}) = \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y} + \alpha \mathbf{w} = 0$$



# Good exercise!

$$\partial_{w_p} E(\mathbf{w}) = (\mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y} + \alpha \mathbf{w})_p$$

$$\nabla E(\mathbf{w}) = \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y} + \alpha \mathbf{w}$$

$$\nabla E(\mathbf{w}) = \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y} + \alpha \mathbf{w} = 0$$

$$\mathbf{X}^\top \mathbf{X} \mathbf{w} + \alpha \mathbf{w} = \mathbf{X}^\top \mathbf{y}$$



# Maximum-Likelihood-Estimation

Given some data and some model  $f(\mathbf{w})$  we can define the **maximum likelihood**

$$\rho(\mathbf{DATA} | \mathbf{w})$$



# Maximum-Likelihood-Estimation

Given some data and some model  $f(\mathbf{w})$  we can define the **maximum likelihood**

$$\rho(\mathbf{DATA} | \mathbf{w})$$

Last week we saw that finding the minimiser of the negative log likelihood is equivalent to minimise the MSE



# Maximum-Likelihood-Estimation

This is very general method, to estimate (fit) some parameter of a model we can use the so-called **Maximum-Likelihood-Estimation** approach





# Maximum-Likelihood-Estimation

This is very general method, to estimate (fit) some parameter of a model we can use the so-called **Maximum-Likelihood-Estimation** approach

$$\hat{\mathbf{w}}_{MLE} = \arg \max_{\mathbf{w}} \rho(\mathbf{DATA} | \mathbf{w})$$



# Maximum-Likelihood-Estimation

This is very general method, to estimate (fit) some parameter of a model we can use the so-called **Maximum-Likelihood-Estimation** approach

$$\hat{\mathbf{w}}_{MLE} = \arg \max_{\mathbf{w}} \rho(\mathbf{DATA} | \mathbf{w})$$

Assuming that the samples are i.i.d.



# Maximum-Likelihood-Estimation

This is very general method, to estimate (fit) some parameter of a model we can use the so-called **Maximum-Likelihood-Estimation** approach

$$\hat{\mathbf{w}}_{MLE} = \arg \max_{\mathbf{w}} \rho(\mathbf{DATA} | \mathbf{w})$$

Assuming that the samples are i.i.d.

$$\hat{\mathbf{w}}_{MLE} = \arg \max_{\mathbf{w}} \prod_i \rho(\mathbf{DATA}_i | \mathbf{w})$$



# Maximum-Likelihood-Estimation

This is very general method, to estimate (fit) some parameter of a model we can use the so-called **Maximum-Likelihood-Estimation** approach

$$\hat{\mathbf{w}}_{MLE} = \arg \max_{\mathbf{w}} \rho(\mathbf{DATA} | \mathbf{w})$$

Assuming that the samples are i.i.d.

$$\hat{\mathbf{w}}_{MLE} = \arg \max_{\mathbf{w}} \prod_i \rho(\mathbf{DATA}_i | \mathbf{w}) = \arg \max_{\mathbf{w}} \log \prod_i \rho(\mathbf{DATA}_i | \mathbf{w})$$



# Maximum-Likelihood-Estimation

This is very general method, to estimate (fit) some parameter of a model we can use the so-called **Maximum-Likelihood-Estimation** approach

$$\hat{\mathbf{w}}_{MLE} = \arg \max_{\mathbf{w}} \rho(\mathbf{DATA} | \mathbf{w})$$

Assuming that the samples are i.i.d.

$$\begin{aligned} \hat{\mathbf{w}}_{MLE} &= \arg \max_{\mathbf{w}} \prod_i \rho(\mathbf{DATA}_i | \mathbf{w}) = \arg \max_{\mathbf{w}} \log \prod_i \rho(\mathbf{DATA}_i | \mathbf{w}) \\ &= \arg \max_{\mathbf{w}} \sum_i \rho(\mathbf{DATA}_i | \mathbf{w}) \end{aligned}$$

# Maximum-A-Posteriori Estimation

Another approach is instead to maximise another quantity the **A-Posteriori** probability

$$\rho(\mathbf{w} \mid \mathbf{DATA})$$



# Maximum-A-Posteriori Estimation

Another approach is instead to maximise another quantity the **A-Posteriori** probability

$$\rho(\mathbf{w} \mid \mathbf{DATA})$$

This is the conditional probability of the model given the data



# Maximum-A-Posteriori Estimation

The optimal model is that the maximise this quantity





# Maximum-A-Posteriori Estimation

The optimal model is that the maximise this quantity

$$\hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \rho(\mathbf{w} | \mathbf{DATA})$$



# Maximum-A-Posteriori Estimation

The optimal model is that the maximise this quantity

$$\hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \rho(\mathbf{w} | \mathbf{DATA})$$

We know however that



# Maximum-A-Posteriori Estimation

The optimal model is that the maximise this quantity

$$\hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \rho(\mathbf{w} | \mathbf{DATA})$$

We know however that

$$\rho(\mathbf{w} | \mathbf{DATA}) = \frac{\rho(\mathbf{DATA} | \mathbf{w})\rho(\mathbf{w})}{\rho(\mathbf{DATA})}$$



# Maximum-A-Posteriori Estimation

The optimal model is that the maximise this quantity

$$\hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \rho(\mathbf{w} | \mathbf{DATA})$$

We know however that

$$\rho(\mathbf{w} | \mathbf{DATA}) = \frac{\rho(\mathbf{DATA} | \mathbf{w}) \rho(\mathbf{w})}{\rho(\mathbf{DATA})}$$

Prior distribution



# Maximum-A-Posteriori Estimation

The optimal model is that the maximise this quantity

$$\hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \rho(\mathbf{w} | \mathbf{DATA})$$

Prior distribution

We know however that

$$\rho(\mathbf{w} | \mathbf{DATA}) = \frac{\rho(\mathbf{DATA} | \mathbf{w}) \rho(\mathbf{w})}{\rho(\mathbf{DATA})}$$

$$\hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \frac{\rho(\mathbf{DATA} | \mathbf{w}) \rho(\mathbf{w})}{\rho(\mathbf{DATA})}$$



# Maximum-A-Posteriori Estimation

The optimal model is that the maximise this quantity

$$\hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \rho(\mathbf{w} | \mathbf{DATA})$$

Prior distribution

We know however that

$$\rho(\mathbf{w} | \mathbf{DATA}) = \frac{\rho(\mathbf{DATA} | \mathbf{w}) \rho(\mathbf{w})}{\rho(\mathbf{DATA})}$$

$$\hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \frac{\rho(\mathbf{DATA} | \mathbf{w}) \rho(\mathbf{w})}{\rho(\mathbf{DATA})} = \arg \max_{\mathbf{w}} \rho(\mathbf{DATA} | \mathbf{w}) \rho(\mathbf{w})$$



# Maximum-A-Posteriori Estimation

The optimal model is that the maximise this quantity

$$\hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \rho(\mathbf{w} | \mathbf{DATA})$$

Prior distribution

We know however that

$$\rho(\mathbf{w} | \mathbf{DATA}) = \frac{\rho(\mathbf{DATA} | \mathbf{w}) \rho(\mathbf{w})}{\rho(\mathbf{DATA})}$$

$$\hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \frac{\rho(\mathbf{DATA} | \mathbf{w}) \rho(\mathbf{w})}{\rho(\mathbf{DATA})} = \arg \max_{\mathbf{w}} \rho(\mathbf{DATA} | \mathbf{w}) \rho(\mathbf{w})$$

This is very similar to the MLE expression but for the prior distribution!

# MLE VS MAP

$$\hat{\mathbf{w}}_{MLE} = \arg \max_{\mathbf{w}} \rho(\mathbf{DATA} | \mathbf{w})$$

$$\hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \rho(\mathbf{w} | \mathbf{DATA}) = \arg \max_{\mathbf{w}} \rho(\mathbf{DATA} | \mathbf{w})\rho(\mathbf{w})$$





# MLE VS MAP

$$\hat{\mathbf{w}}_{MLE} = \arg \max_{\mathbf{w}} \rho(\mathbf{DATA} | \mathbf{w})$$

$$\hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \rho(\mathbf{w} | \mathbf{DATA}) = \arg \max_{\mathbf{w}} \rho(\mathbf{DATA} | \mathbf{w})\rho(\mathbf{w})$$

If the prior is “flat” the two are the equivalent, otherwise they differ due to the assumptions on the prior distribution of the parameters of the models



# MLE for linear regression

Recall: maximum likelihood estimator for least-squares linear regression

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{y}, \mathbf{X} | \mathbf{w})) \right\}$$



# MLE for linear regression

Recall: maximum likelihood estimator for least-squares linear regression

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{y}, \mathbf{X} | \mathbf{w})) \right\} \\ &= \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{X} | \mathbf{w})\rho(\mathbf{y} | \mathbf{X}, \mathbf{w})) \right\} \quad \text{Factoring of likelihood}\end{aligned}$$



# MLE for linear regression

Recall: maximum likelihood estimator for least-squares linear regression

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{y}, \mathbf{X} | \mathbf{w})) \right\} \\ &= \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{X} | \mathbf{w})\rho(\mathbf{y} | \mathbf{X}, \mathbf{w})) \right\} && \text{Factoring of likelihood} \\ &= \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{X})\rho(\mathbf{y} | \mathbf{X}, \mathbf{w})) \right\} && \mathbf{X} \text{ does not depend on } \mathbf{w}\end{aligned}$$



# MLE for linear regression

Recall: maximum likelihood estimator for least-squares linear regression

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{y}, \mathbf{X} | \mathbf{w})) \right\} \\ &= \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{X} | \mathbf{w})\rho(\mathbf{y} | \mathbf{X}, \mathbf{w})) \right\} && \text{Factoring of likelihood} \\ &= \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{X})\rho(\mathbf{y} | \mathbf{X}, \mathbf{w})) \right\} && \mathbf{X} \text{ does not depend on } \mathbf{w} \\ &= \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{y} | \mathbf{X}, \mathbf{w})) \right\} && \rho(\mathbf{X}) \text{ does not depend on } \mathbf{w}\end{aligned}$$



# MLE for linear regression

Recall: maximum likelihood estimator for least-squares linear regression

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{y}, \mathbf{X} | \mathbf{w})) \right\} \\ &= \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{X} | \mathbf{w})\rho(\mathbf{y} | \mathbf{X}, \mathbf{w})) \right\} && \text{Factoring of likelihood} \\ &= \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{X})\rho(\mathbf{y} | \mathbf{X}, \mathbf{w})) \right\} && \mathbf{X} \text{ does not depend on } \mathbf{w} \\ &= \arg \min_{\mathbf{w}} \left\{ -\log (\rho(\mathbf{y} | \mathbf{X}, \mathbf{w})) \right\} && \rho(\mathbf{X}) \text{ does not depend on } \mathbf{w} \\ &= \arg \min_{\mathbf{w}} \left\{ -\log \left( \prod_{i=1}^s \rho(\mathbf{y}_i | \mathbf{x}_i, \mathbf{w}) \right) \right\} && \text{samples are iid}\end{aligned}$$



# MLE for linear regression

Recall: maximum likelihood estimator for least-squares linear regression

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\{ -\log \left( \prod_{i=1}^s \rho(\mathbf{y}_i | \mathbf{x}_i, \mathbf{w}) \right) \right\}$$



# MLE for linear regression

Recall: maximum likelihood estimator for least-squares linear regression

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \left\{ -\log \left( \prod_{i=1}^s \rho(\mathbf{y}_i | \mathbf{x}_i, \mathbf{w}) \right) \right\} \\ &= \arg \min_{\mathbf{w}} \left\{ -\log \left( \prod_{i=1}^s \mathcal{N}(\mathbf{y}_i | \langle \mathbf{x}_i, \mathbf{w} \rangle, \sigma^2) \right) \right\} \quad \mathcal{N} = \text{probability} \\ &\quad \text{density function for} \\ &\quad \text{normal distribution}\end{aligned}$$





# MLE for linear regression

Recall: maximum likelihood estimator for least-squares linear regression

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \left\{ -\log \left( \prod_{i=1}^s \rho(\mathbf{y}_i | \mathbf{x}_i, \mathbf{w}) \right) \right\} \\ &= \arg \min_{\mathbf{w}} \left\{ -\log \left( \prod_{i=1}^s \mathcal{N}(\mathbf{y}_i | \langle \mathbf{x}_i, \mathbf{w} \rangle, \sigma^2) \right) \right\} \quad \mathcal{N} = \text{probability} \\ & \quad \text{density function for} \\ & \quad \text{normal distribution} \\ &= \arg \min_{\mathbf{w}} \left\{ -\log \left( \prod_{i=1}^s \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{(\mathbf{y}_i - \langle \mathbf{x}_i, \mathbf{w} \rangle)^2}{2\sigma^2} \right) \right) \right\}\end{aligned}$$

# MLE for linear regression

Recall: maximum likelihood estimator for least-squares linear regression

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\{ -\log \left( \prod_{i=1}^s \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{(y_i - \langle \mathbf{x}_i, \mathbf{w} \rangle)^2}{2\sigma^2} \right) \right) \right\}$$



# MLE for linear regression

Recall: maximum likelihood estimator for least-squares linear regression

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \left\{ -\log \left( \prod_{i=1}^s \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{(\mathbf{y}_i - \langle \mathbf{x}_i, \mathbf{w} \rangle)^2}{2\sigma^2} \right) \right) \right\} \\ &= \arg \min_{\mathbf{w}} \left\{ s \log \left( \sqrt{2\pi\sigma^2} \right) + \frac{1}{2\sigma^2} \sum_{i=1}^s \left( \mathbf{y}_i - \langle \mathbf{x}_i, \mathbf{w} \rangle \right)^2 \right\}\end{aligned}$$



# MLE for linear regression

Recall: maximum likelihood estimator for least-squares linear regression

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \left\{ -\log \left( \prod_{i=1}^s \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{(\mathbf{y}_i - \langle \mathbf{x}_i, \mathbf{w} \rangle)^2}{2\sigma^2} \right) \right) \right\} \\ &= \arg \min_{\mathbf{w}} \left\{ s \log \left( \sqrt{2\pi\sigma^2} \right) + \frac{1}{2\sigma^2} \sum_{i=1}^s \left( \mathbf{y}_i - \langle \mathbf{x}_i, \mathbf{w} \rangle \right)^2 \right\} \\ &= \arg \min_{\mathbf{w}} \left\{ \frac{1}{2\sigma^2} \sum_{i=1}^s \left( \mathbf{y}_i - \langle \mathbf{x}_i, \mathbf{w} \rangle \right)^2 \right\}\end{aligned}$$



# MLE for linear regression

Recall: maximum likelihood estimator for least-squares linear regression

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \left\{ -\log \left( \prod_{i=1}^s \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{(\mathbf{y}_i - \langle \mathbf{x}_i, \mathbf{w} \rangle)^2}{2\sigma^2} \right) \right) \right\} \\ &= \arg \min_{\mathbf{w}} \left\{ s \log \left( \sqrt{2\pi\sigma^2} \right) + \frac{1}{2\sigma^2} \sum_{i=1}^s \left( \mathbf{y}_i - \langle \mathbf{x}_i, \mathbf{w} \rangle \right)^2 \right\} \\ &= \arg \min_{\mathbf{w}} \left\{ \frac{1}{2\sigma^2} \sum_{i=1}^s \left( \mathbf{y}_i - \langle \mathbf{x}_i, \mathbf{w} \rangle \right)^2 \right\} \quad \text{This is the MSE!}\end{aligned}$$

# MAP for ridge regression

Maximum a-posteriori estimate for ridge regression:

$$\hat{w} = \arg \min_w \left\{ -\log (\rho(w | X, y)) \right\}$$



# MAP for ridge regression

Maximum a-posteriori estimate for ridge regression:

$$\begin{aligned}\hat{w} &= \arg \min_w \left\{ -\log (\rho(w | X, y)) \right\} \\ &= \arg \min_w \left\{ -\log \left( \frac{\rho(y, X | w)\rho(w)}{\rho(y, X)} \right) \right\} \quad \text{Bayes' Rule/Theorem}\end{aligned}$$



# MAP for ridge regression

Maximum a-posteriori estimate for ridge regression:

$$\begin{aligned}\hat{w} &= \arg \min_w \left\{ -\log (\rho(w | X, y)) \right\} \\ &= \arg \min_w \left\{ -\log \left( \frac{\rho(y, X | w)\rho(w)}{\rho(y, X)} \right) \right\} && \text{Bayes' Rule/Theorem} \\ &= \arg \min_w \left\{ -\log (\rho(y, X | w)\rho(w)) \right\}\end{aligned}$$





# MAP for ridge regression

Maximum a-posteriori estimate for ridge regression:

$$\begin{aligned}\hat{w} &= \arg \min_w \left\{ -\log (\rho(w | X, y)) \right\} \\ &= \arg \min_w \left\{ -\log \left( \frac{\rho(y, X | w)\rho(w)}{\rho(y, X)} \right) \right\} && \text{Bayes' Rule/Theorem} \\ &= \arg \min_w \left\{ -\log (\rho(y, X | w)\rho(w)) \right\} \\ &= \arg \min_w \left\{ -\log (\rho(y | X, w)\rho(w)) \right\}\end{aligned}$$



# MAP for ridge regression

Maximum a-posteriori estimate for ridge regression:

$$\hat{w} = \arg \min_w \left\{ -\log (\rho(y | X, w)\rho(w)) \right\}$$



# MAP for ridge regression

Maximum a-posteriori estimate for ridge regression:

$$\begin{aligned}\hat{w} &= \arg \min_w \left\{ -\log (\rho(y | X, w)\rho(w)) \right\} \\ &= \arg \min_w \left\{ -\log \left( \rho(w) \prod_{i=1}^s \rho(y_i | x_i, w) \right) \right\}\end{aligned}$$



# MAP for ridge regression

Maximum a-posteriori estimate for ridge regression:

$$\begin{aligned}\hat{w} &= \arg \min_w \left\{ -\log (\rho(y | X, w)\rho(w)) \right\} \\ &= \arg \min_w \left\{ -\log \left( \rho(w) \prod_{i=1}^s \rho(y_i | x_i, w) \right) \right\} \\ &= \arg \min_w \left\{ -\log \left( \prod_{j=1}^{d+1} \mathcal{N} \left( w_j \mid 0, \frac{\sigma^2}{\alpha} \right) \prod_{i=1}^s \mathcal{N} \left( y_i \mid \langle x_i, w \rangle, \sigma^2 \right) \right) \right\}\end{aligned}$$

# MAP for ridge regression

Maximum a-posteriori estimate for ridge regression:

$$\hat{w} = \arg \min_w \left\{ -\log \left( \prod_{j=1}^{d+1} \mathcal{N} \left( w_j \mid 0, \frac{\sigma^2}{\alpha} \right) \prod_{i=1}^s \mathcal{N} \left( y_i \mid \langle x_i, w \rangle, \sigma^2 \right) \right) \right\}$$



# MAP for ridge regression

Maximum a-posteriori estimate for ridge regression:

$$\begin{aligned}\hat{w} &= \arg \min_w \left\{ -\log \left( \prod_{j=1}^{d+1} \mathcal{N} \left( w_j \mid 0, \frac{\sigma^2}{\alpha} \right) \prod_{i=1}^s \mathcal{N} \left( y_i \mid \langle x_i, w \rangle, \sigma^2 \right) \right) \right\} \\ &= \arg \min_w \left\{ -\log \left( \prod_{j=1}^{d+1} \sqrt{\frac{\alpha}{2\pi\sigma^2}} \exp \left( -\frac{\alpha}{2\sigma^2} |w_j|^2 \right) \prod_{i=1}^s \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{(y_i - \langle x_i, w \rangle)^2}{2\sigma^2} \right) \right) \right\}\end{aligned}$$

# MAP for ridge regression

Maximum a-posteriori estimate for ridge regression:

$$\begin{aligned}\hat{w} &= \arg \min_w \left\{ -\log \left( \prod_{j=1}^{d+1} \mathcal{N} \left( w_j \mid 0, \frac{\sigma^2}{\alpha} \right) \prod_{i=1}^s \mathcal{N} \left( y_i \mid \langle x_i, w \rangle, \sigma^2 \right) \right) \right\} \\ &= \arg \min_w \left\{ -\log \left( \prod_{j=1}^{d+1} \sqrt{\frac{\alpha}{2\pi\sigma^2}} \exp \left( -\frac{\alpha}{2\sigma^2} |w_j|^2 \right) \prod_{i=1}^s \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{(y_i - \langle x_i, w \rangle)^2}{2\sigma^2} \right) \right) \right\} \\ &= \arg \min_w \left\{ \frac{1}{2} \sum_{i=1}^s (y_i - \langle x_i, w \rangle)^2 + \frac{\alpha}{2} \|w\|^2 \right\}\end{aligned}$$

# MAP for ridge regression

Maximum a-posteriori estimate for ridge regression:

$$\begin{aligned}\hat{w} &= \arg \min_w \left\{ -\log \left( \prod_{j=1}^{d+1} \mathcal{N} \left( w_j \mid 0, \frac{\sigma^2}{\alpha} \right) \prod_{i=1}^s \mathcal{N} \left( y_i \mid \langle x_i, w \rangle, \sigma^2 \right) \right) \right\} \\ &= \arg \min_w \left\{ -\log \left( \prod_{j=1}^{d+1} \sqrt{\frac{\alpha}{2\pi\sigma^2}} \exp \left( -\frac{\alpha}{2\sigma^2} |w_j|^2 \right) \prod_{i=1}^s \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{(y_i - \langle x_i, w \rangle)^2}{2\sigma^2} \right) \right) \right\} \\ &= \arg \min_w \left\{ \frac{1}{2} \sum_{i=1}^s (y_i - \langle x_i, w \rangle)^2 + \frac{\alpha}{2} \|w\|^2 \right\}\end{aligned}$$

This is the MSE for the ridge regression