



# MTH786: Machine Learning with Python

Week 2: Linear and polynomial regressions. Convexity. Regularisation.

# Linear regression

# What is a regression?

---

The main task of **regression analysis** consist of finding a mapping  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  between  $s$  known input arguments  $\{\mathbf{x}^{(j)}\}_{j=1}^s$ , with  $\mathbf{x}^{(j)} \in \mathbb{R}^n$  for all  $1 \leq j \leq s$ , onto a collection of  $s$  output elements  $\{\mathbf{y}^{(j)}\}_{j=1}^s$ , with  $\mathbf{y}^{(j)} \in \mathbb{R}^m$  for all  $1 \leq j \leq s$  such that

$$f\left(\mathbf{x}^{(j)}\right) \approx \mathbf{y}^{(j)}, \quad \forall j \in \{1, \dots, s\}.$$

## How to find a suitable function $f(x)$ ?

1. Consider a class of functions  $\mathcal{F}$  parametrised by a set of parameters  $\vec{w}$ .
2. Use the training data provided to find an optimal value of parameters  $\vec{w}$ .
3. Validate the result using validation data.

**What does optimal mean?  $\longrightarrow$  smallest error**



# Example: linear regression

---

How a linear function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  could be parametrised?

$$f_{\mathbf{W}}(\mathbf{x}) = \mathbf{w}^{(0)} + \sum_{j=1}^n \mathbf{w}^{(j)} x_j,$$

where  $\mathbf{w}^{(j)} \in \mathbb{R}^m$  for all  $0 \leq j \leq n$  are vector-valued unknown weights.

**Remark:** unless other agreement is discussed we use superscript for object numbering and subscript for vector indexing. Vectors and matrices are **bold**.

**Note:** The sum in above goes over vector coordinates not vector labels, i.e.

$$f_{\mathbf{W}}(\mathbf{x}^{(j)}) = \mathbf{w}^{(0)} + \sum_{k=1}^n \mathbf{w}^{(k)} x_k^{(j)}, \quad \left(f_{\mathbf{W}}(\mathbf{x}^{(j)})\right)_{\ell} = w_{\ell}^{(0)} + \sum_{k=1}^n w_{\ell}^{(k)} x_k^{(j)}.$$

In the case of our main interest  $m = 1$  and one can write

$$f_{\mathbf{W}}(\mathbf{x}) = \left\langle \mathring{\mathbf{x}}, \mathbf{W} \right\rangle = \mathring{\mathbf{x}}^T \mathbf{W}, \text{ where } \mathring{\mathbf{x}} = (1, x_1, \dots, x_n)^T, \mathbf{W} = (w_0, w_1, \dots, w_n)^T.$$

# Example: linear regression

$$f_{\mathbf{W}}(\mathbf{x}) = \langle \mathring{\mathbf{x}}, \mathbf{W} \rangle = \mathring{\mathbf{x}}^T \mathbf{W}, \text{ where } \mathring{\mathbf{x}} = (1, x_1, \dots, x_n)^T, \mathbf{W} = (w_0, w_1, \dots, w_n)^T.$$

**Can we fit the data?**

**Question:** for the training data  $\{(\mathbf{x}^{(j)}, \mathbf{y}^{(j)})\}_{j=1}^s$  find  $f_{\mathbf{W}}$  such that  $f_{\mathbf{W}}(\mathbf{x}^{(j)}) = \mathbf{y}^{(j)}$

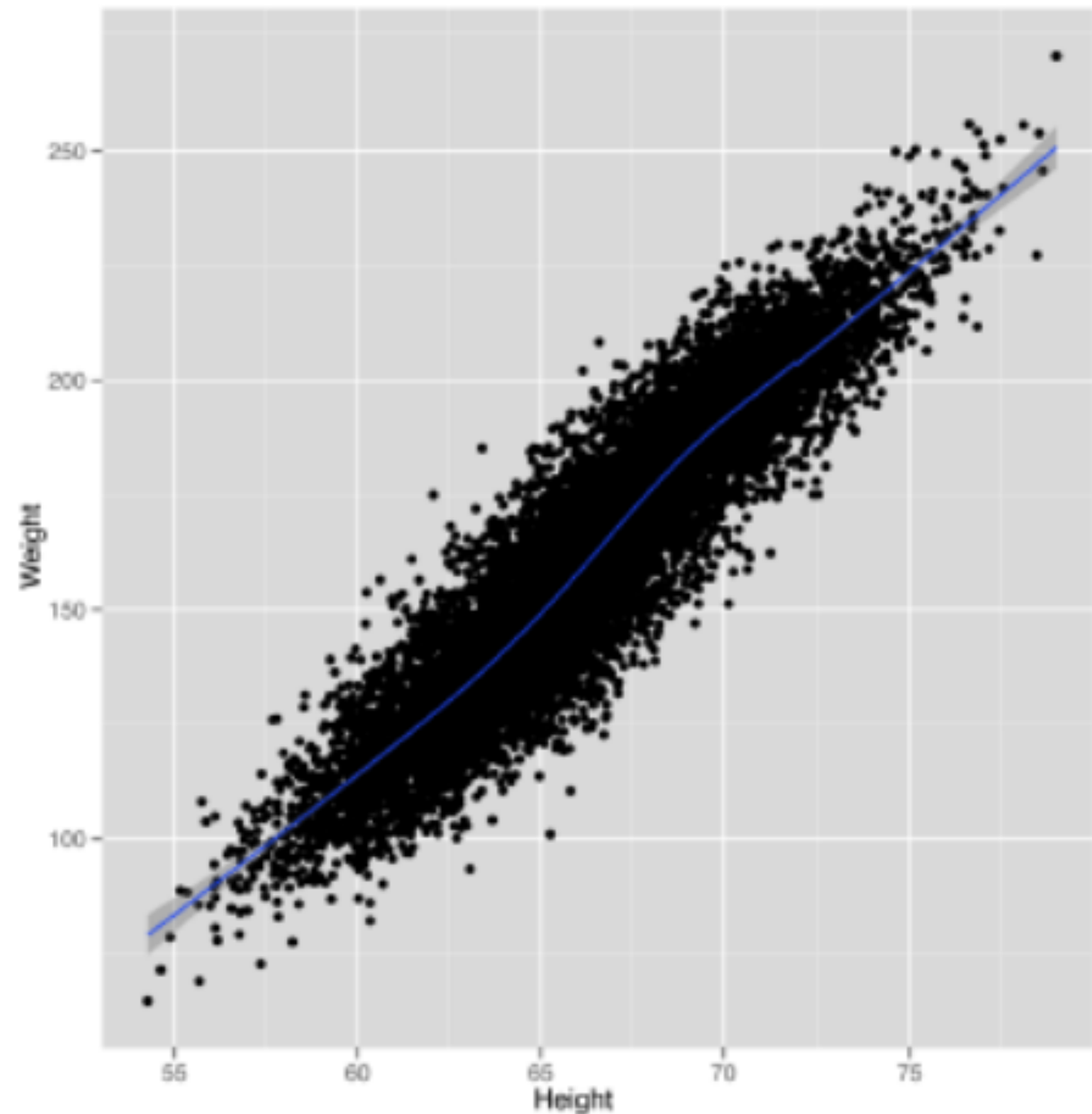
**Example:** let  $s = 3$  and  $n = 2$ , then

$$\begin{pmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \\ 1 & x_1^{(3)} & x_2^{(3)} \end{pmatrix} \begin{pmatrix} w_1^{(0)} \\ w_1^{(1)} \\ w_1^{(2)} \end{pmatrix} = \begin{pmatrix} y_1^{(1)} \\ y_1^{(2)} \\ y_1^{(3)} \end{pmatrix}$$

**Is the system uniquely solvable? —————> Yes, for a "good" input**  
**What would happen in general case? —————>  $sm$  eqs,  $(n+1)m$  vars**  
**Only if  $s = n + 1$**

# Example: linear regression

Only if  $s = n + 1$ . Is it realistic?



**Example:** Prediction of a person's weight from its height by using linear regression.

**Dimensionality:**  $n = 1$  and  $s \gg 1$ .

**Question:** how to measure an error of approximation when exact fit is not possible?

**Idea:** we need to measure how far is the prediction out of expected value in average.

Remind us the *variance*. For  $n = 1$

$$\text{MSE}(\mathbf{W}) = \frac{1}{2s} \sum_{j=1}^s \left| f_{\mathbf{W}}(x^{(j)}) - y^{(j)} \right|^2$$

Mean-squared error

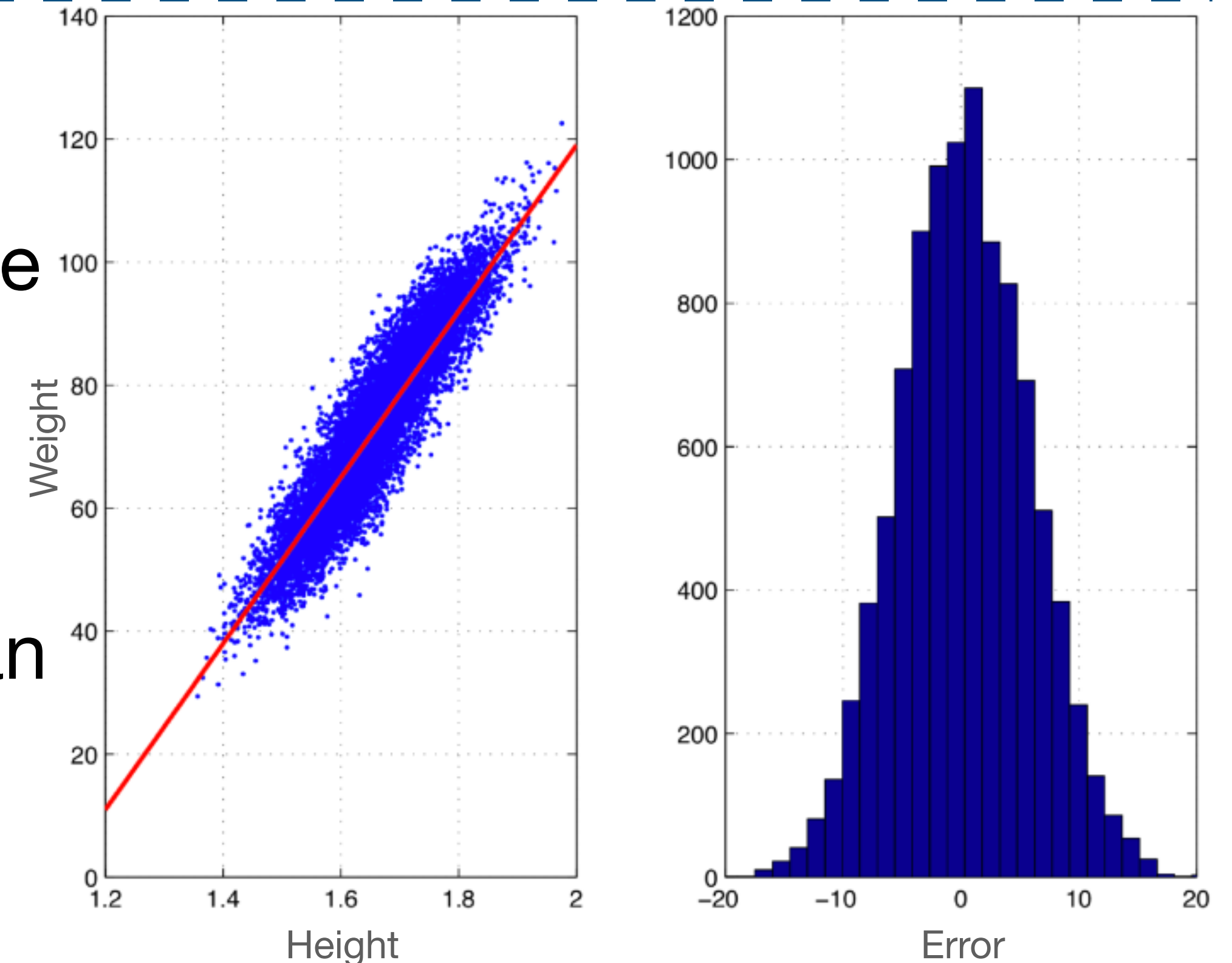
# Statistical motivation of MSE

**Question:** how did we come up with MSE?

**Motivation:** let us consider the distribution of errors  $\varepsilon^{(j)} := f_{\mathbf{W}}(\mathbf{x}^{(j)}) - \mathbf{y}^{(j)}$ . If we believe the model has been chosen right then  $\varepsilon^{(j)}$  are r.v. Moreover,  $\{\varepsilon^{(1)}, \dots, \varepsilon^{(s)}\}$  are independent.

**Observation:** the distribution is Gaussian.

**Assumption:**  $\{\varepsilon^{(1)}, \dots, \varepsilon^{(s)}\}$  are i.i.d. Gaussian random variables with mean 0 and variance  $\sigma^2$ , i.e.



$$\rho(\varepsilon^{(j)} | 0, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\varepsilon^{(j)})^2}{2\sigma^2}}$$

**Problem:** find optimal values of  $\mathbf{W}$  to maximise the likelihood of sampling  $\mathbf{y}^{(j)}$ .

[m.poplavskyi@qmul.ac.uk](mailto:m.poplavskyi@qmul.ac.uk)



# Statistical motivation of MSE

**Remark:** below we work with  $m = 1$ . For  $m > 1$  the method still works.

Likelihood

$$\begin{aligned}\rho(\varepsilon^{(1)}, \varepsilon^{(2)}, \dots, \varepsilon^{(s)} | 0, \sigma^2) &= \prod_{j=1}^s (2\pi\sigma^2)^{-\frac{1}{2}} e^{-\frac{(\varepsilon^{(j)})^2}{2\sigma^2}} = (2\pi\sigma^2)^{-\frac{s}{2}} \prod_{j=1}^s e^{-\frac{(\mathbf{y}^{(j)} - f(\mathbf{x}^{(j)}))^2}{2\sigma^2}} \\ &= \rho(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(s)} | f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(s)}), \sigma^2)\end{aligned}$$

Maximisation

$$\begin{aligned}\hat{W} &= \arg \min_{W \in \mathbb{R}^{n+1}} \{-\log(\rho(\mathbf{Y} | f_W(\mathbf{X}), \sigma^2))\} = \arg \min_{W \in \mathbb{R}^{n+1}} \left\{ -\log \left( \prod_{j=1}^s \rho(\mathbf{y}^{(j)} | f_W(\mathbf{x}^{(j)}), \sigma^2) \right) \right\} = \arg \min_{W \in \mathbb{R}^{n+1}} \left\{ -\sum_{j=1}^s \log \rho(\mathbf{y}^{(j)} | f_W(\mathbf{x}^{(j)}), \sigma^2) \right\} \\ &= \arg \min_{W \in \mathbb{R}^{n+1}} \left\{ \frac{1}{2\sigma^2} \sum_{j=1}^s (\mathbf{y}^{(j)} - f_W(\mathbf{x}^{(j)}))^2 + \frac{s}{2} \log(2\pi\sigma^2) \right\} = \arg \min_{W \in \mathbb{R}^{n+1}} \left\{ \frac{2s}{\sigma^2} \text{MSE}(W) + \frac{s}{2} \log(2\pi\sigma^2) \right\} = \arg \min_{W \in \mathbb{R}^{n+1}} \{\text{MSE}(W)\}\end{aligned}$$

**Question 1:** How to find a minimiser?

**Question 2:** Does the minimiser exist?

**Question 3:** Is the minimiser unique?



# Example: 1+1 dimension

**Problem:** For a given set of input-output pairs  $\left\{ \left( x^{(j)}, y^{(j)} \right) \right\}_{j=1}^s$  find a linear function  $f_{\mathbf{W}} : \mathbb{R} \rightarrow \mathbb{R}$  such that  $\text{MSE}(\mathbf{W})$  is minimised over all possible  $\mathbf{W}$ .

## Zero approximation

Let  $f_{\mathbf{W}}(x) = w^{(0)}$  for all  $x \in \mathbb{R}$ . Then

$$\text{MSE}(w^{(0)}) = \frac{1}{2s} \sum_{j=1}^s (w^{(0)} - y^{(j)})^2$$

$$\frac{\partial \text{MSE}(w^{(0)})}{\partial w^{(0)}} = \frac{1}{s} \sum_{j=1}^s (w^{(0)} - y^{(j)}) \stackrel{?}{=} 0$$

$$\hat{w}_0 = \frac{y^{(1)} + \dots + y^{(s)}}{s} =: \bar{y}$$

## Real linear approximation

Let  $f_{\mathbf{W}}(x) = w^{(0)} + w^{(1)}x$ ,  $x \in \mathbb{R}$ . Then

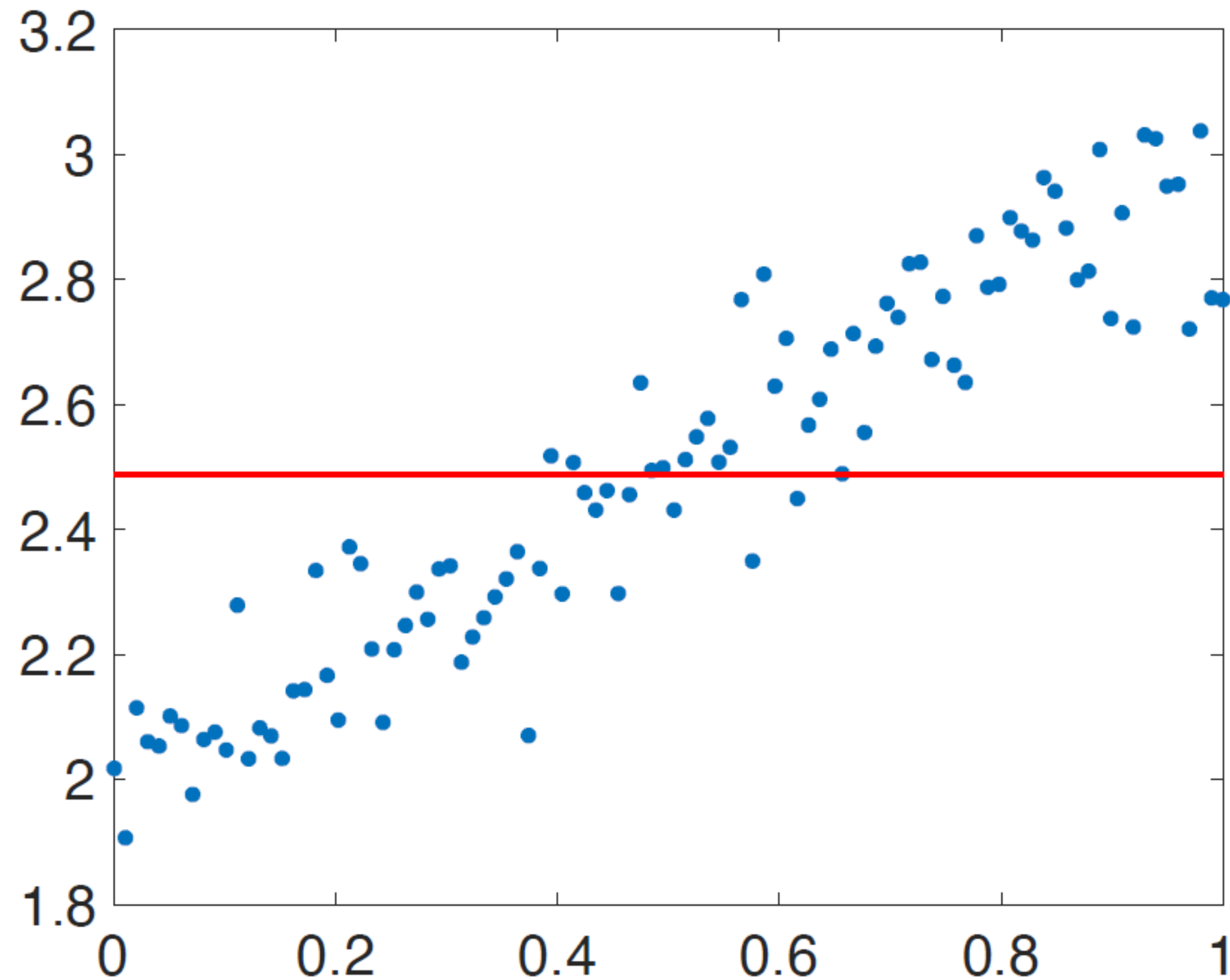
$$\text{MSE}(\mathbf{W}) = \frac{1}{2s} \sum_{j=1}^s (w^{(0)} + w^{(1)}x^{(j)} - y^{(j)})^2$$

$$\frac{\partial \text{MSE}(\mathbf{W})}{\partial w^{(0)}} = \frac{1}{s} \sum_{j=1}^s (w^{(0)} + w^{(1)}x^{(j)} - y^{(j)}) \stackrel{?}{=} 0$$

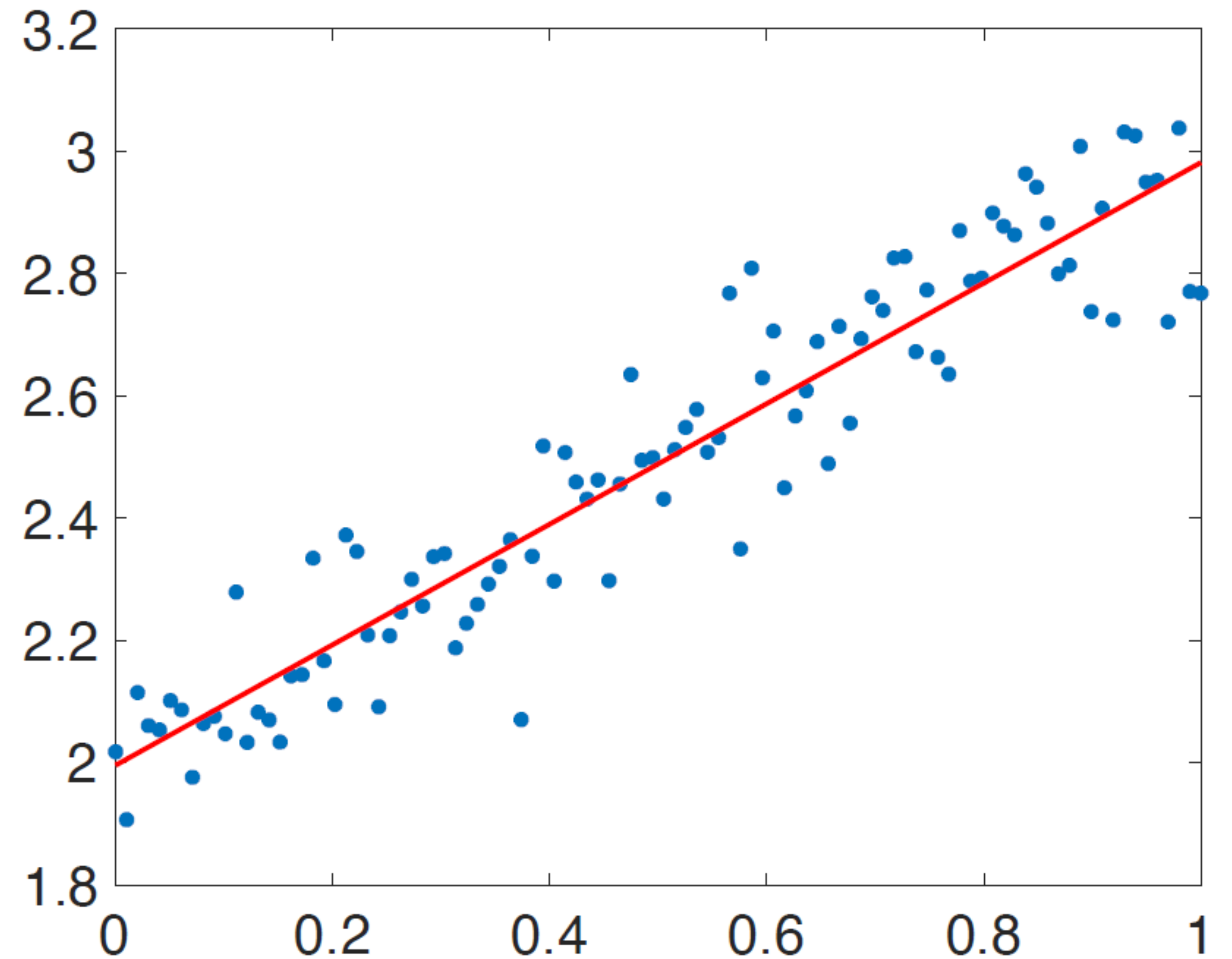
$$\frac{\partial \text{MSE}(\mathbf{W})}{\partial w^{(1)}} = \frac{1}{s} \sum_{j=1}^s x^{(j)} (w^{(0)} + w^{(1)}x^{(j)} - y^{(j)}) \stackrel{?}{=} 0$$

$$\hat{w}^{(0)} = \frac{\bar{y}\bar{x}^2 - \bar{x}\bar{xy}}{\bar{x}^2 - \bar{x}^2}, \quad \hat{w}^{(1)} = \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x}^2 - \bar{x}^2}$$

# Example: 1+1 dimension



$$\hat{w}^{(0)} = 2.4489$$



$$\begin{aligned}\hat{w}^{(0)} &= 1.9962 \\ \hat{w}^{(1)} &= 0.9854\end{aligned}$$

# Example: $n + 1$ dimensions

---

**Problem:** For a given set of input-output pairs  $\left\{ \left( \mathbf{x}^{(j)}, y^{(j)} \right) \right\}_{j=1}^s$  find a linear function  $f_{\mathbf{W}} : \mathbb{R}^n \rightarrow \mathbb{R}$  such that  $\text{MSE}(\mathbf{W})$  is minimised over all possible  $\mathbf{W}$ .

Let  $f_{\mathbf{W}}(\mathbf{x}) = w^{(0)} + \sum_{j=1}^n w^{(j)} \mathbf{x}_j = \mathbf{\hat{x}}^T \mathbf{W}$ , where  $\mathbf{W} = (w^{(0)}, \dots, w^{(n)})^T$ .

$$\text{MSE}(\mathbf{W}) = \frac{1}{2s} \sum_{j=1}^s \left( [\mathbf{XW}]_j - y^{(j)} \right)^2 = \frac{1}{2s} \left\| \mathbf{XW} - \mathbf{Y} \right\|^2.$$

$$\nabla \text{MSE}(\hat{\mathbf{W}}) = 0 \Rightarrow \mathbf{X}^T \mathbf{X} \hat{\mathbf{W}} = \hat{\mathbf{X}}^T \mathbf{Y}. \text{ (Exercise!)}$$

$$\hat{\mathbf{W}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

# Discussion of further development

---

**Question 1:** Are there any other cost-functions except of MSE ( $\mathbf{W}$ )?

**Answer 1:** Indeed there is a wide variety of cost-functions, for example in the case of  $m = 1$  one can consider mean absolute error (MAE)

$$\text{MAE}(\mathbf{W}) = \frac{1}{2s} \sum_{j=1}^s \left| f(\mathbf{x}^{(j)}) - y^{(j)} \right|.$$

+: robust to outliers (see Assignment 2) -: non-differentiable function

**Question 2:** How the same approach would work when allowing nonlinearity?

**Answer 2:** We can use the same approach for nonlinear functions  $f(\mathbf{W})$ .

**Question 3:** Whether the minimiser we have found is indeed a minimiser?

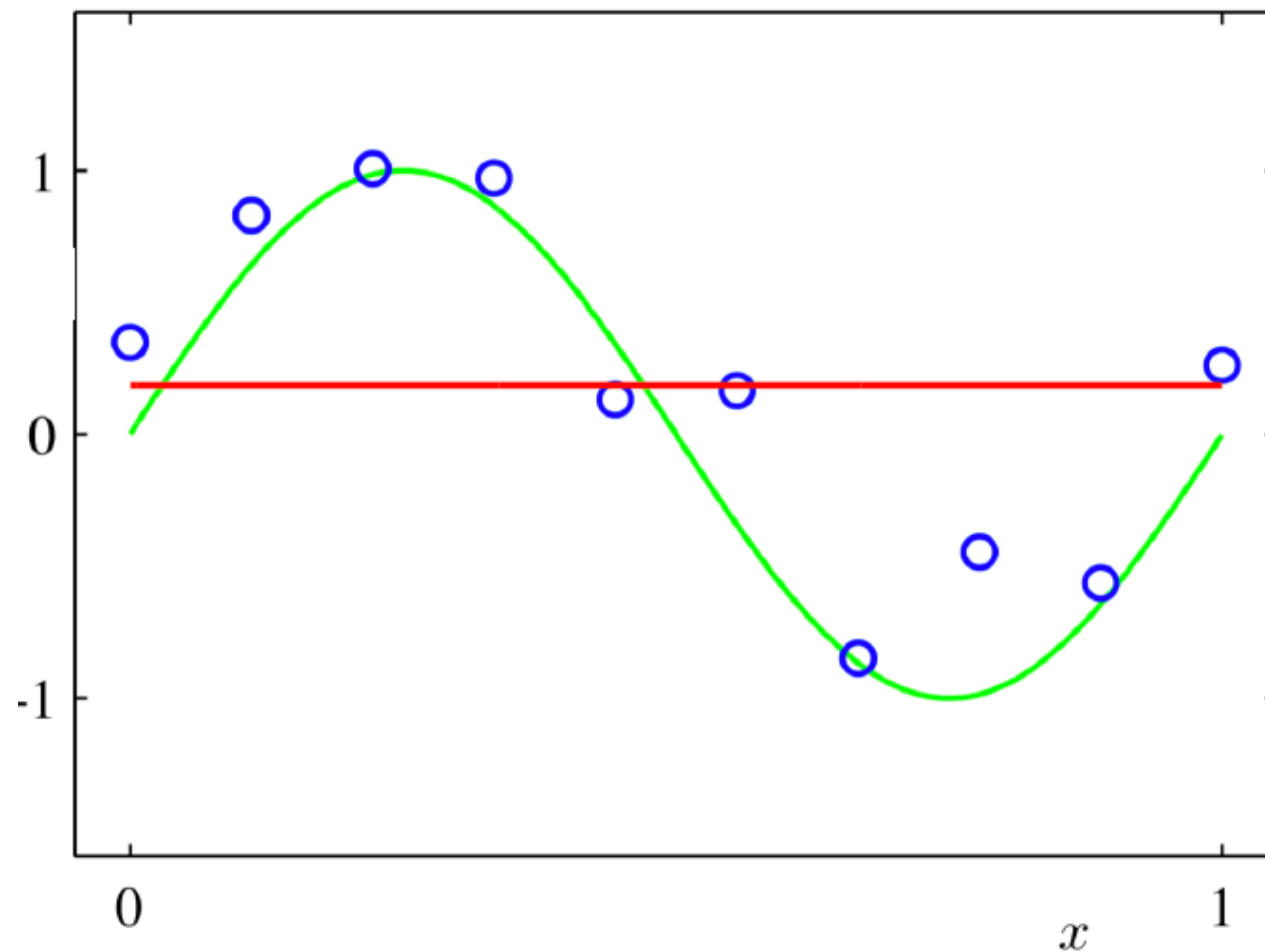
**Answer 3:** It can be checked via evaluation of higher derivatives.



# Polynomial regression

# Polynomial regression

In this section we work in 1 + 1 dimensions only. Data:  $\{x^{(j)}, y^{(j)}\}_{j=1}^s$ .



The green curve is the plot of a function. Blue dots are points from the curve perturbed by some random forces. Red line is a line of linear regression. **It does not fit well.**

**How to adapt a function  $f(x)$ ?**

The next simplest class of functions are polynomials, i.e.

$$f_{\mathbf{W}}(x) = w^{(0)} + w^{(1)}x + \dots + w^{(d)}x^d,$$

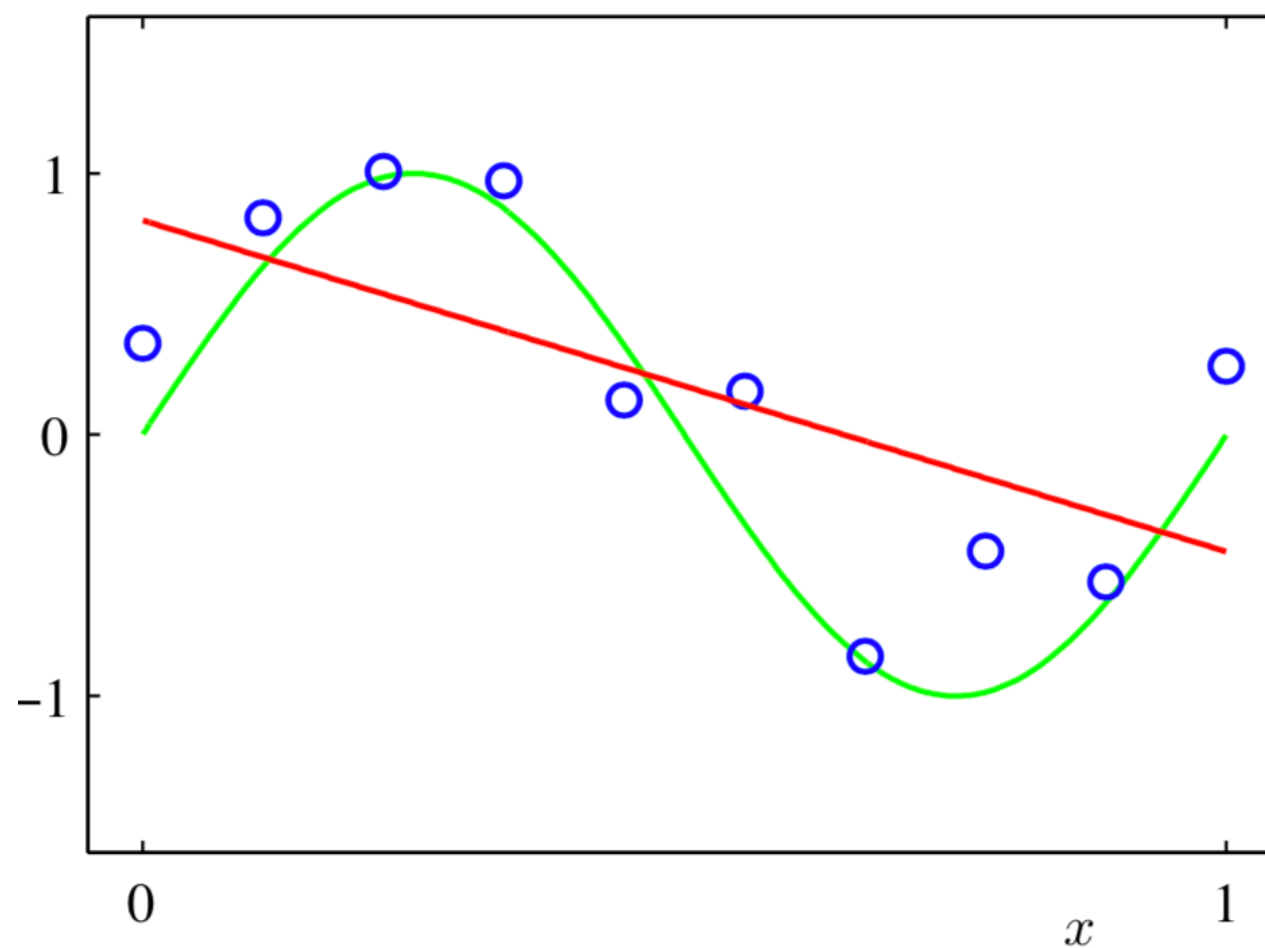
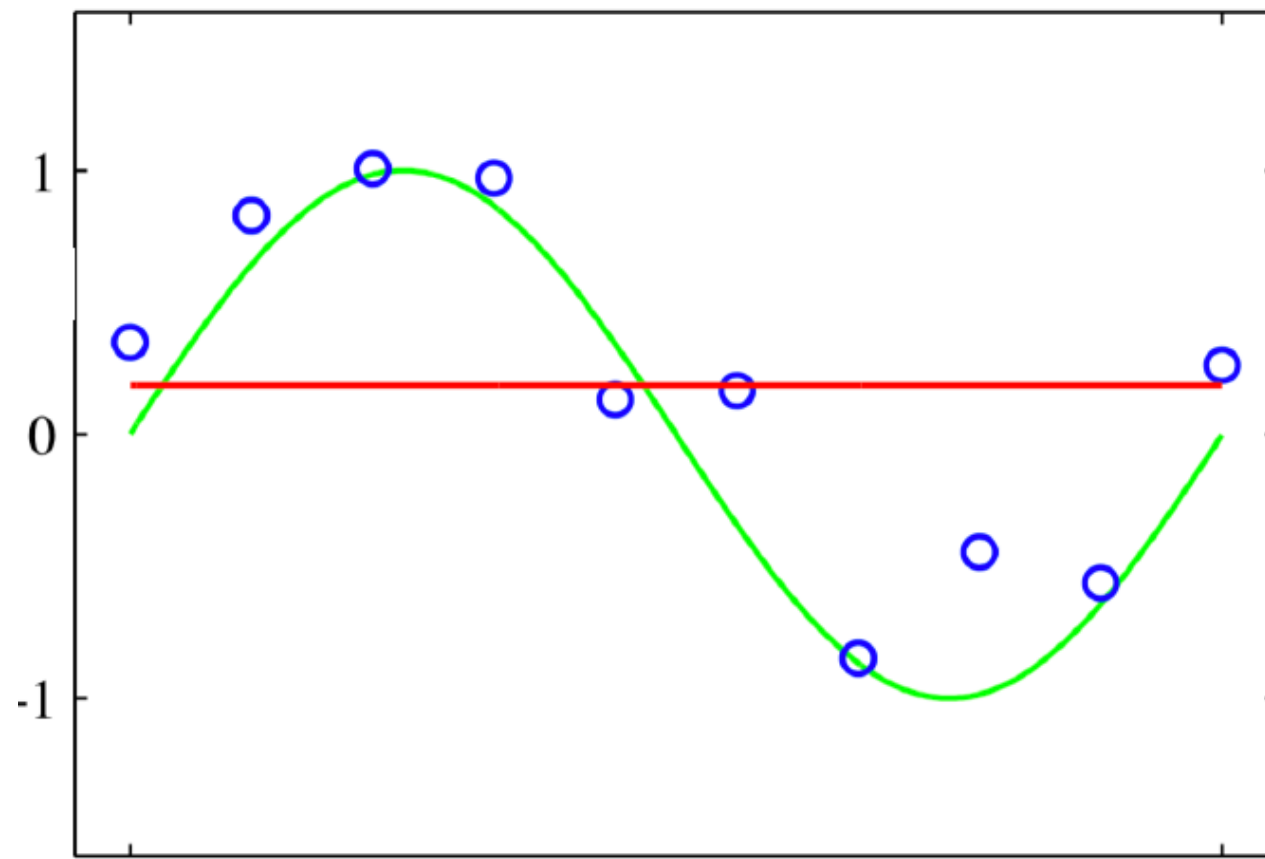
for some unknown vector  $\mathbf{W} = (w^{(0)}, \dots, w^{(d)})^T$ . Let

$$\phi(x) \stackrel{\text{def}}{=} (1, x, \dots, x^d)^T, \text{ then } f_{\mathbf{W}}(x) = \langle \phi(x), \mathbf{W} \rangle.$$



Augmented/Extended feature vectors

# Polynomial regression



How to adapt a function  $f(x)$ ?

For a new problem  $f_{\mathbf{W}}(\mathbf{x}) = \langle \phi(x), \mathbf{W} \rangle$ . Modified MSE-problem is

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W} \in \mathbb{R}^{d+1}} \left\{ \frac{1}{2s} \left\| \Phi(\mathbf{X}) \mathbf{W} - \mathbf{Y} \right\|^2 \right\},$$

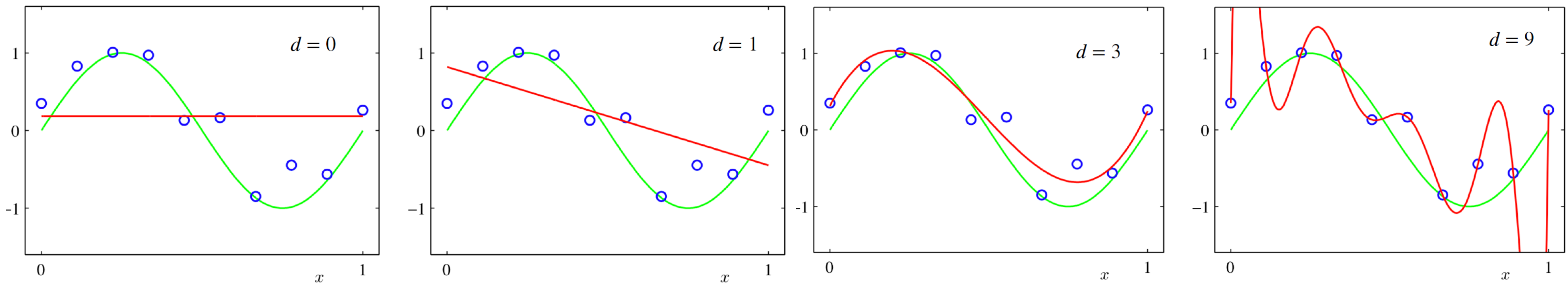
where  $\Phi(\mathbf{X}) = \begin{pmatrix} \phi(x^{(1)}) \\ \vdots \\ \phi(x^{(s)}) \end{pmatrix}$ . Solution to the problem is:

$$\hat{\mathbf{W}} = \left( \Phi(\mathbf{X})^T \Phi(\mathbf{X}) \right)^{-1} \Phi(\mathbf{X})^T \mathbf{Y}$$

# Under- and over- fitting

The models we consider can be too limited or too rich:

- We say that the function is **underfitting** the data if we can't find a function that is a good fit to our data.
- We say that the function is **overfitting** the data if we find a function that fits the data too well.



Both are issues, and difficult to address in practice, as we do not know what part of the data is signal and what is noise.

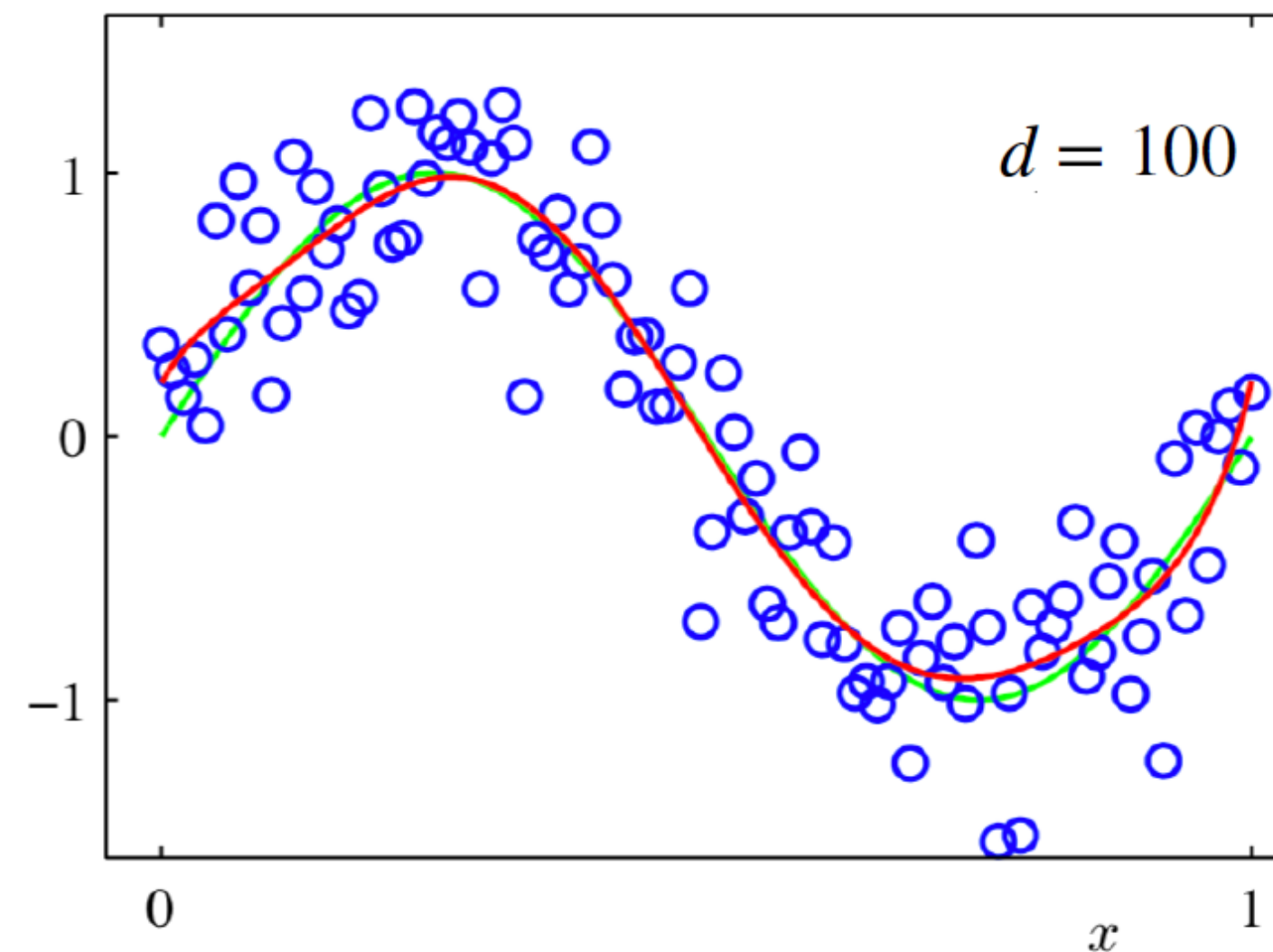
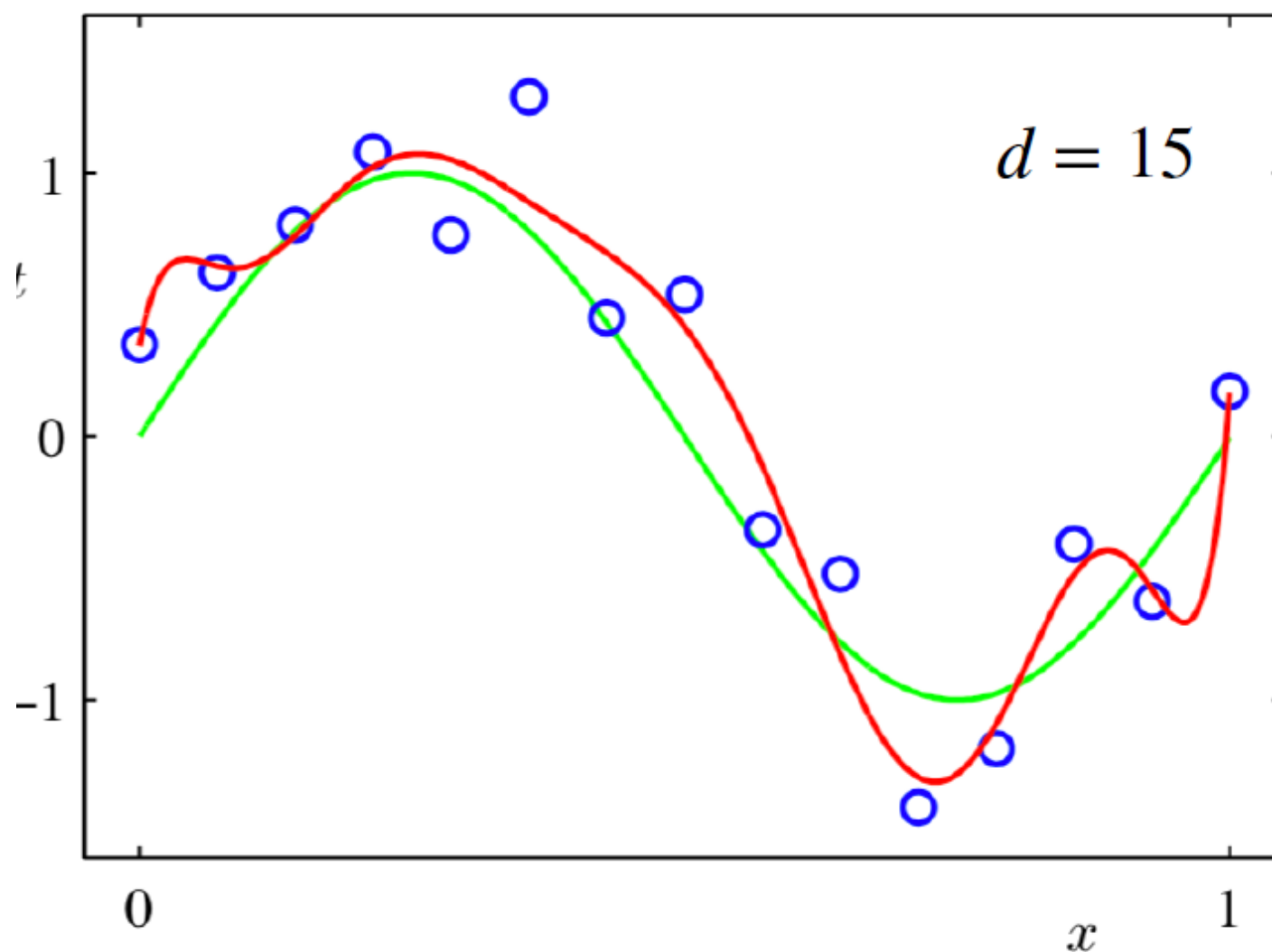


# Under- and over- fitting

The models we consider can be too limited or too rich:

- We say that the function is **underfitting** the data if we can't find a function that is a good fit to our data.
- We say that the function is **overfitting** the data if we find a function that fits the data too well.

## How to tackle the problem of overfitting?



- Increase the number of samples
- Use regularisation

# Other types of regressions

Instead of a **power basis** one may consider other basis functions, for example:

- Trigonometric functions

Let us consider functions  $a_k(x) = \sin(\pi kx)$  and  $b_k(x) = \cos(\pi kx)$  and corresponding augmented vectors

$$\phi(x) = (b_0(x), a_1(x), b_1(x), \dots, a_d(x), b_d(x))$$

- Radial basis functions

One can also consider bell functions either  $\psi(x) = e^{-\frac{x^2}{2\sigma^2}}$  or

$$\psi(x) = \left(1 - \frac{1}{\alpha} |x|\right) 1_{|x| \leq \alpha} \text{ and } \phi(x) = (\psi(x - \mu_1), \dots, \psi(x - \mu_d))$$

# Minimisers and convexity

# Minimisation problem

Previously we have seen that the problem of linear/polynomial regression is:

$$\text{MSE}(\mathbf{W}) = \frac{1}{2s} \sum_{j=1}^s \left\| f(\mathbf{x}^{(j)}) - \mathbf{y}^{(j)} \right\|^2 \rightarrow \min$$

And in case of linear regression this can be solved via:

$$\nabla \text{MSE}(\hat{\mathbf{W}}) = 0 \Leftrightarrow \mathbf{X}^T \mathbf{X} \hat{\mathbf{W}} = \mathbf{X}^T \mathbf{Y}$$

While in case of polynomial regression this can be solved via:

$$\nabla \text{MSE}(\hat{\mathbf{W}}) = 0 \Leftrightarrow \Phi(\mathbf{X})^T \Phi(\mathbf{X}) \hat{\mathbf{W}} = \Phi(\mathbf{X})^T \mathbf{Y}$$

**Can we be sure that this is the minimiser?**



# 3 questions to answer

---

There are 3 questions/concerns to raise:

1. Why solving of the equation  $\nabla \text{MSE}(\hat{W}) = 0$  is equivalent to the problem  $\hat{W} = \arg \min_{W \in \mathbb{R}^{(d+1) \times m}} \text{MSE}(W)$ .

2. Why does a solution to minimisation problem exist?

3. Why is a solution to minimisation problem unique?

**Question 2: Linear algebra**

**Question 3: Convexity**

# Existence of a minimiser

Solving  $\nabla \text{MSE}(\mathbf{W}) = 0$  in the case of linear/polynomial regression is equivalent to:

$$\Phi(\mathbf{X})^T \Phi(\mathbf{X}) \mathbf{W} = \Phi(\mathbf{X})^T \mathbf{Y} \quad (1)$$

**Lemma:** For any matrix  $\mathbf{M} \in \mathbb{R}^{n \times m}$  we have

$$\ker(\mathbf{M}^T \mathbf{M}) = \ker \mathbf{M} \qquad \text{ran}(\mathbf{M}^T \mathbf{M}) = \text{ran} \mathbf{M}$$

*Proof of the lemma: see lecture notes (non-examinable)*

*Application of the lemma:* the right hand side of (1) belongs to  $\ker \Phi(\mathbf{X})^T$ . By the lemma this belongs to  $\ker(\Phi(\mathbf{X})^T \Phi(\mathbf{X}))$ , i.e. there exist a vector  $\hat{\mathbf{W}}$  solving (1).  
**The minimisation problem is thus solvable.**

# Convexity

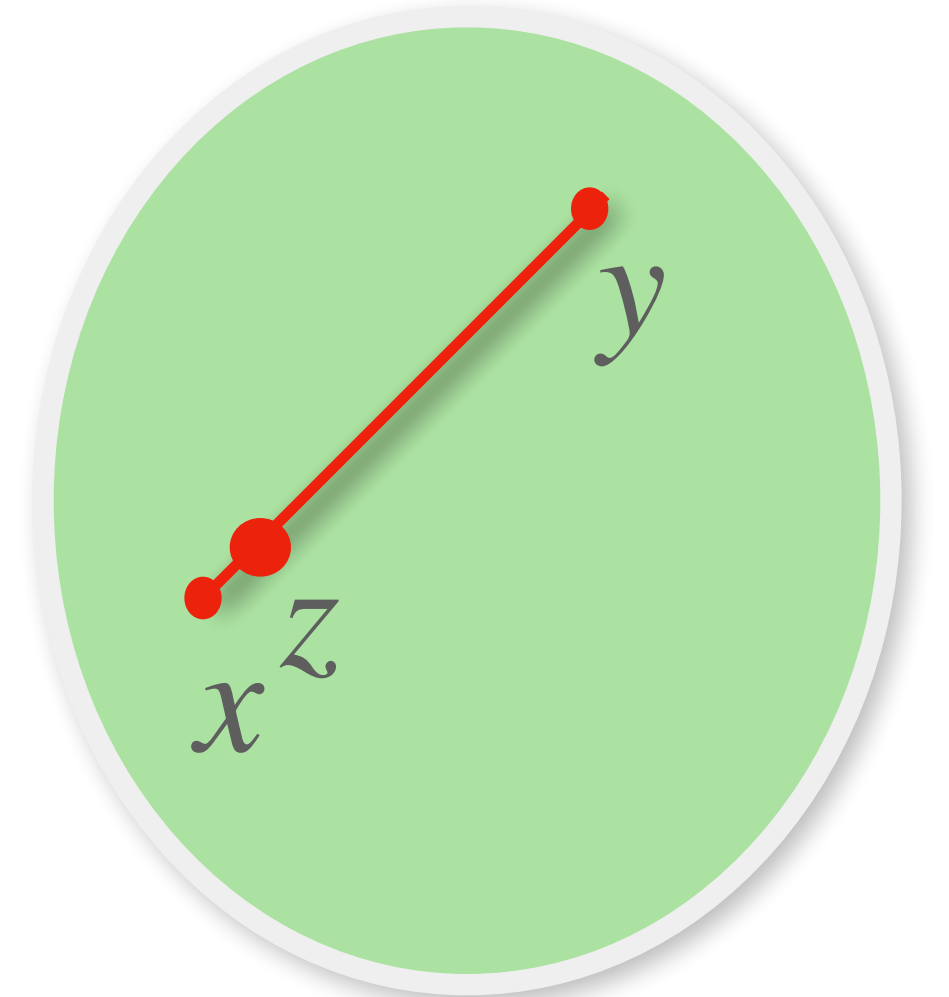
In this section we will show that the cost function is a convex one and then will use the theory of convex functions to conclude that the gradient of cost function is equal to zero only at the point of minimum.

**Definition.** A subset  $C$  of a vector space is said to be **convex set** if for any  $\lambda \in [0,1]$  and  $\mathbf{x}, \mathbf{y} \in C$  the element

$$\mathbf{z}_\lambda := \lambda \mathbf{x} + (1 - \lambda) \mathbf{y} \in C.$$

**Definition.** A function  $f$  that is defined on a convex set  $C$  is said to be a **convex function** if for any  $\lambda \in [0,1]$  and  $\mathbf{x}, \mathbf{y} \in C$

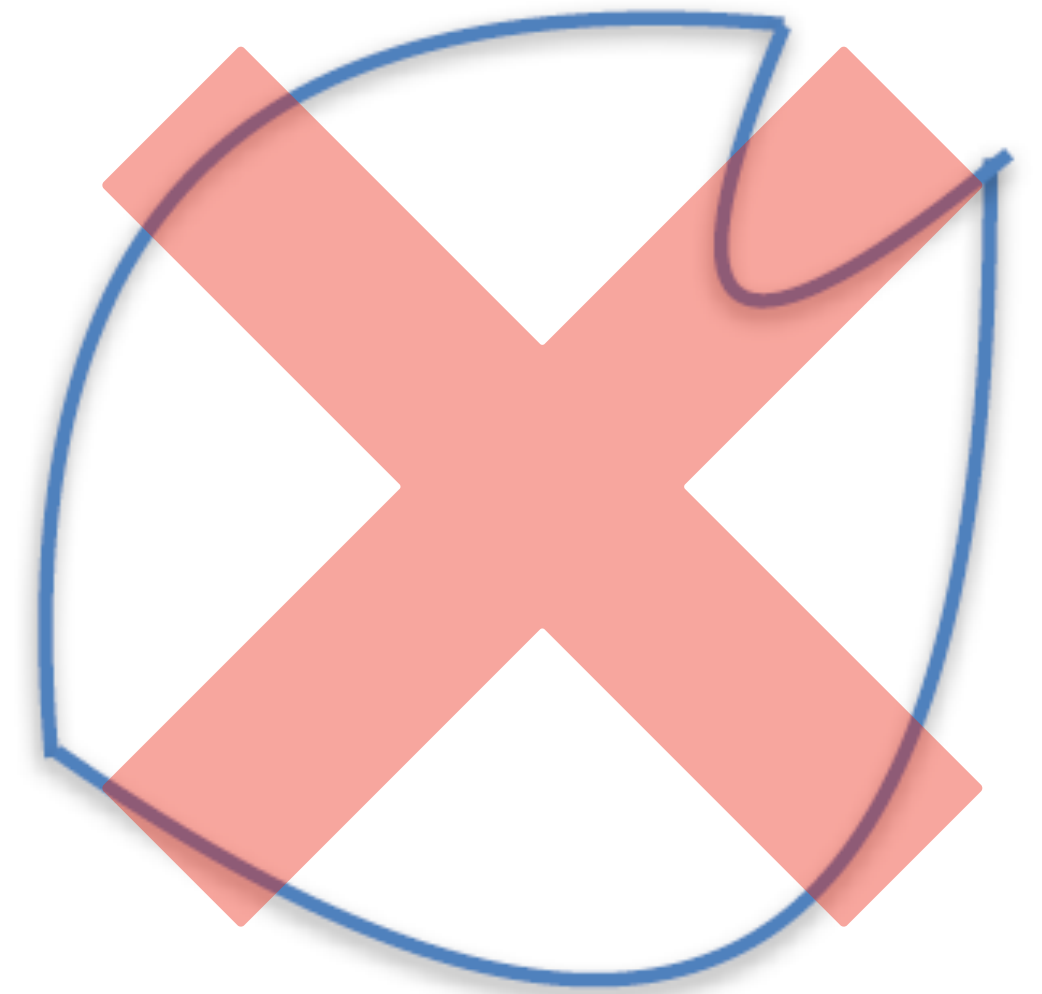
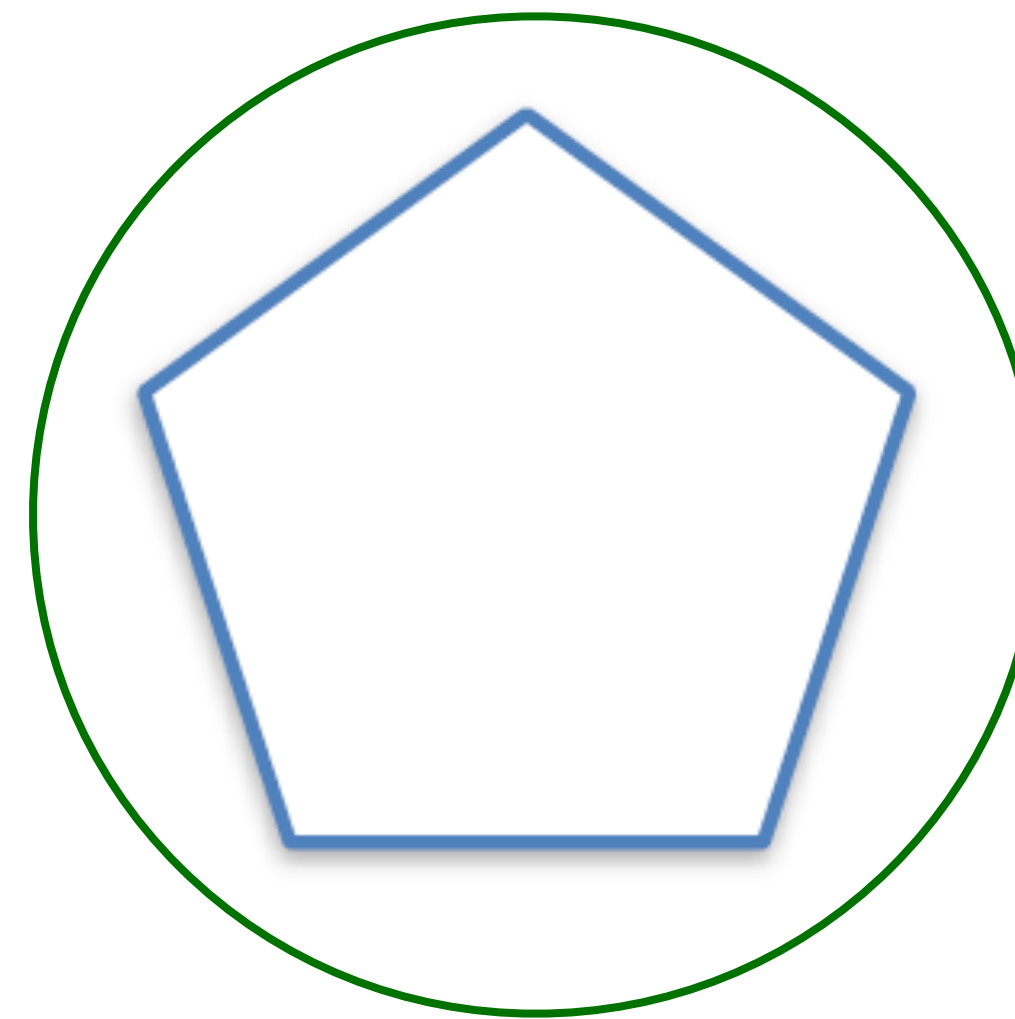
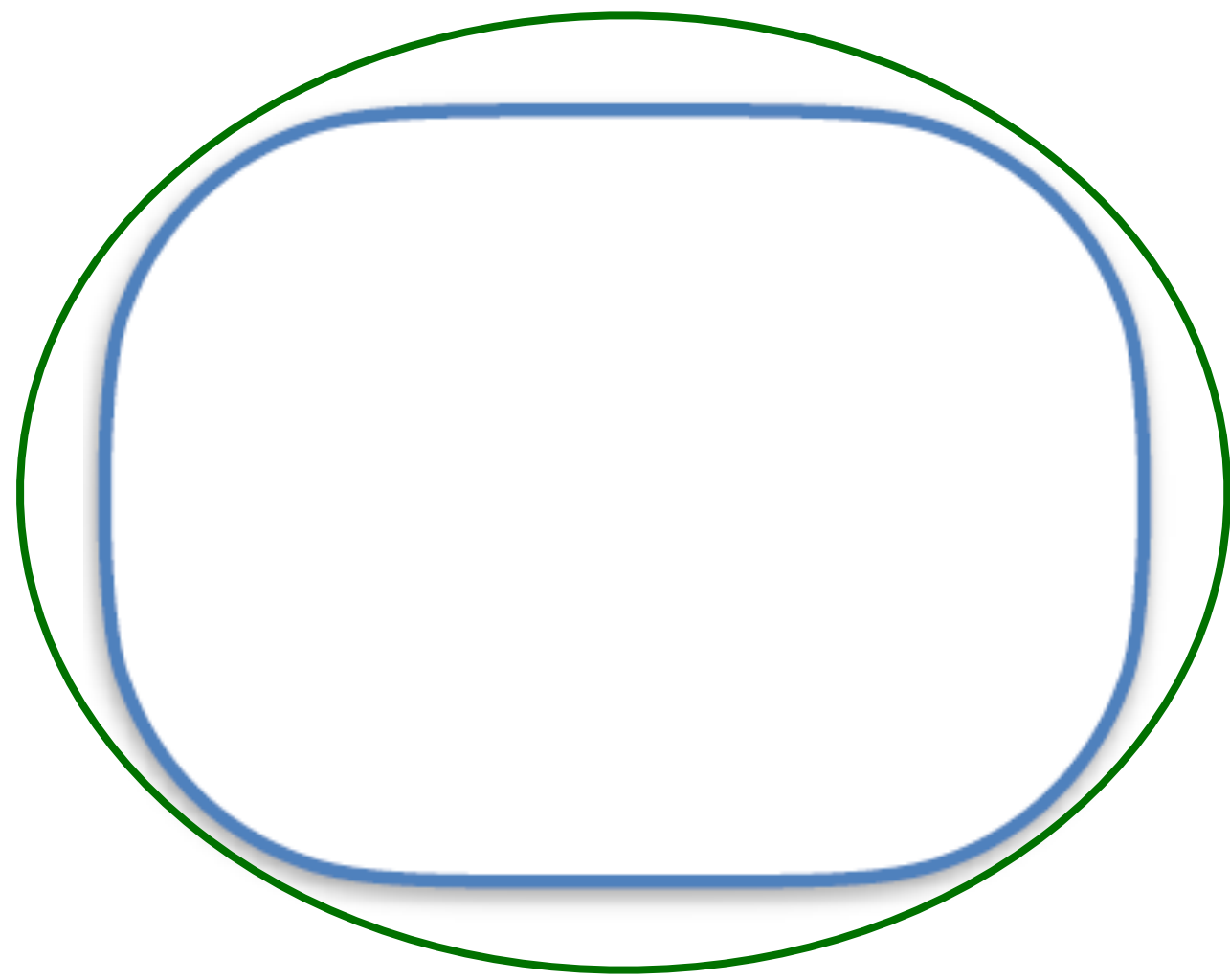
$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}).$$



# Convexity quiz

---

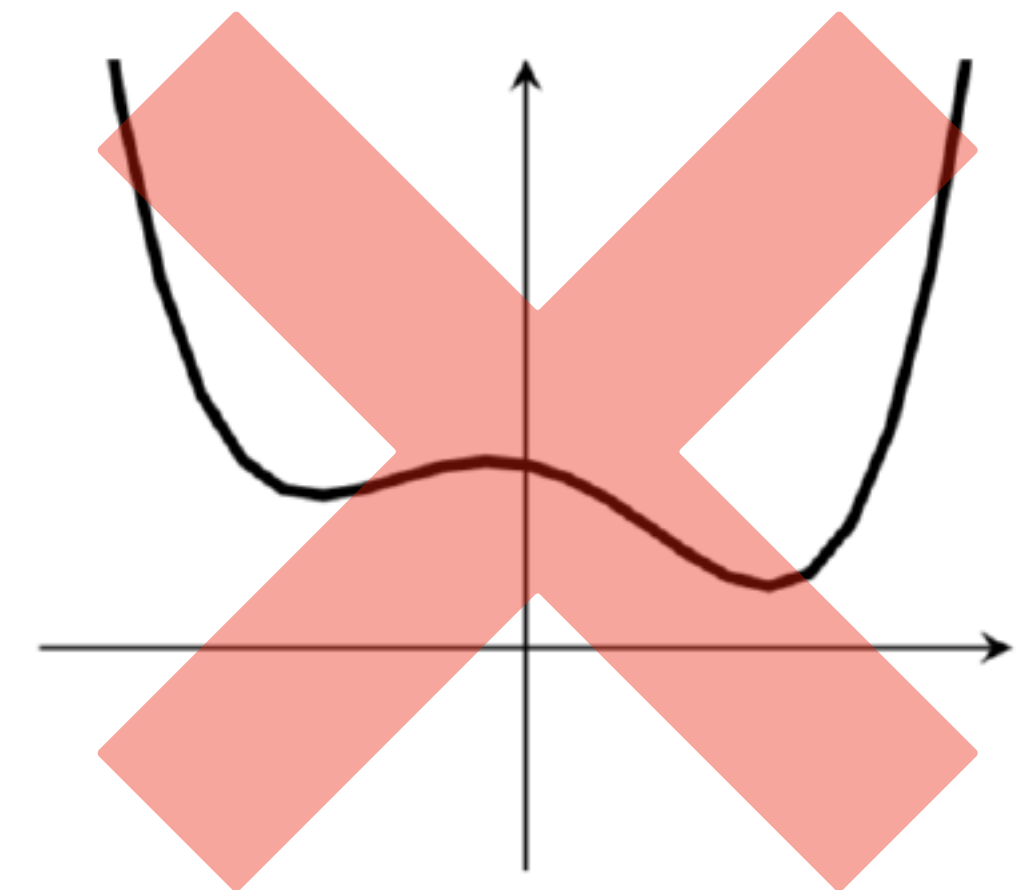
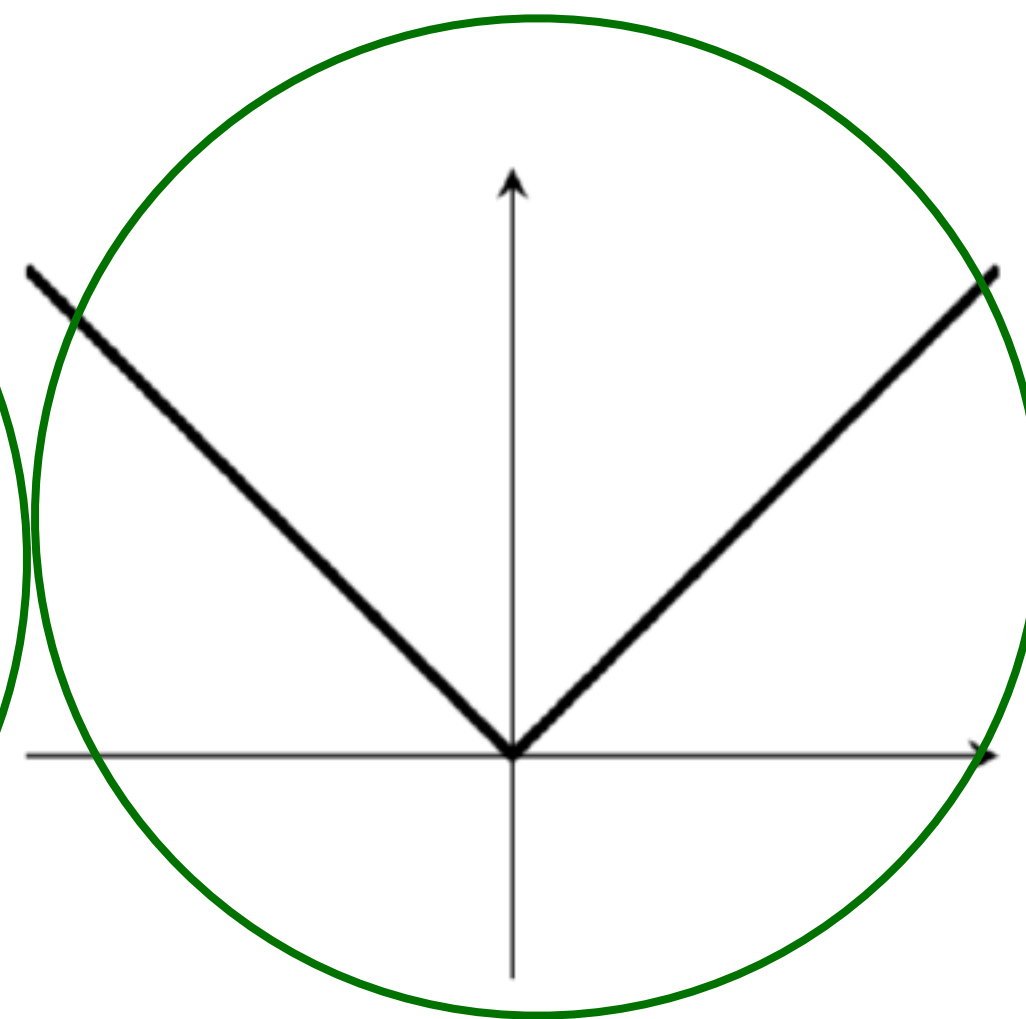
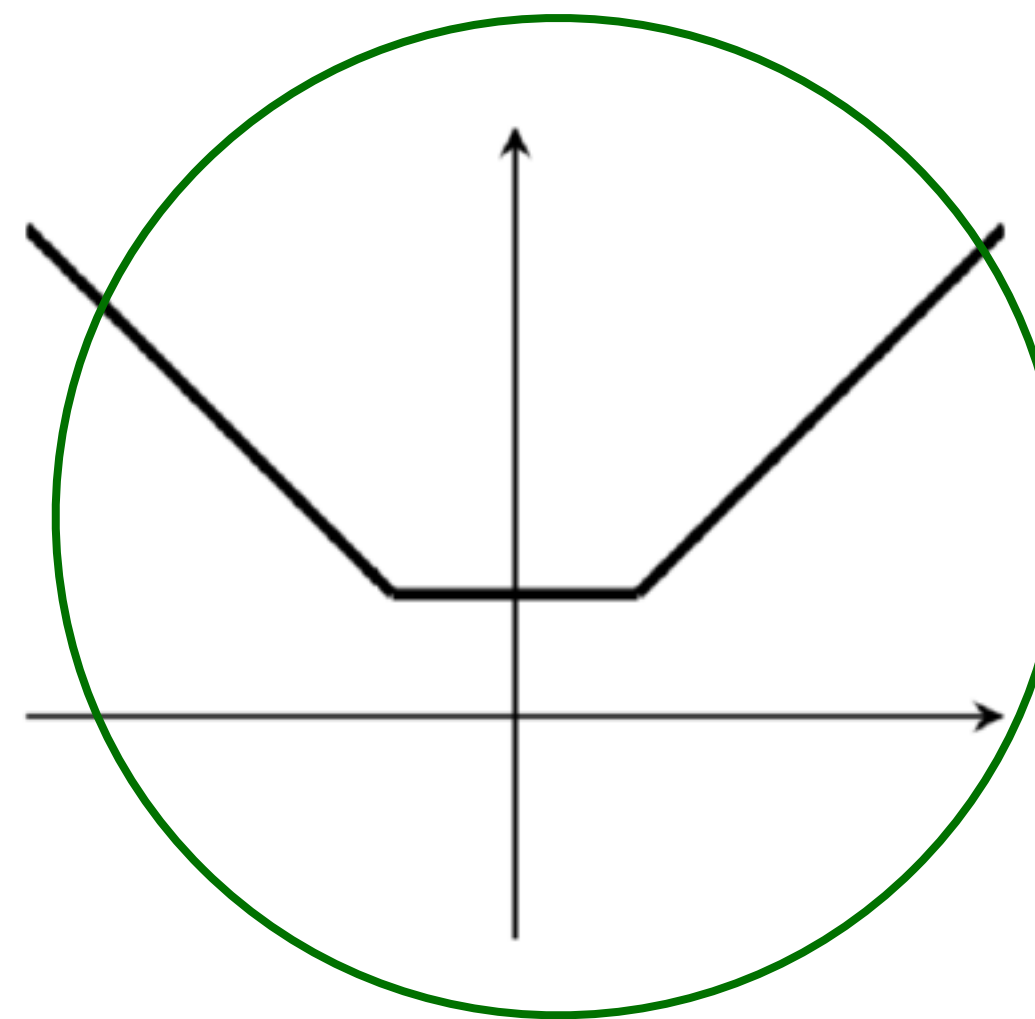
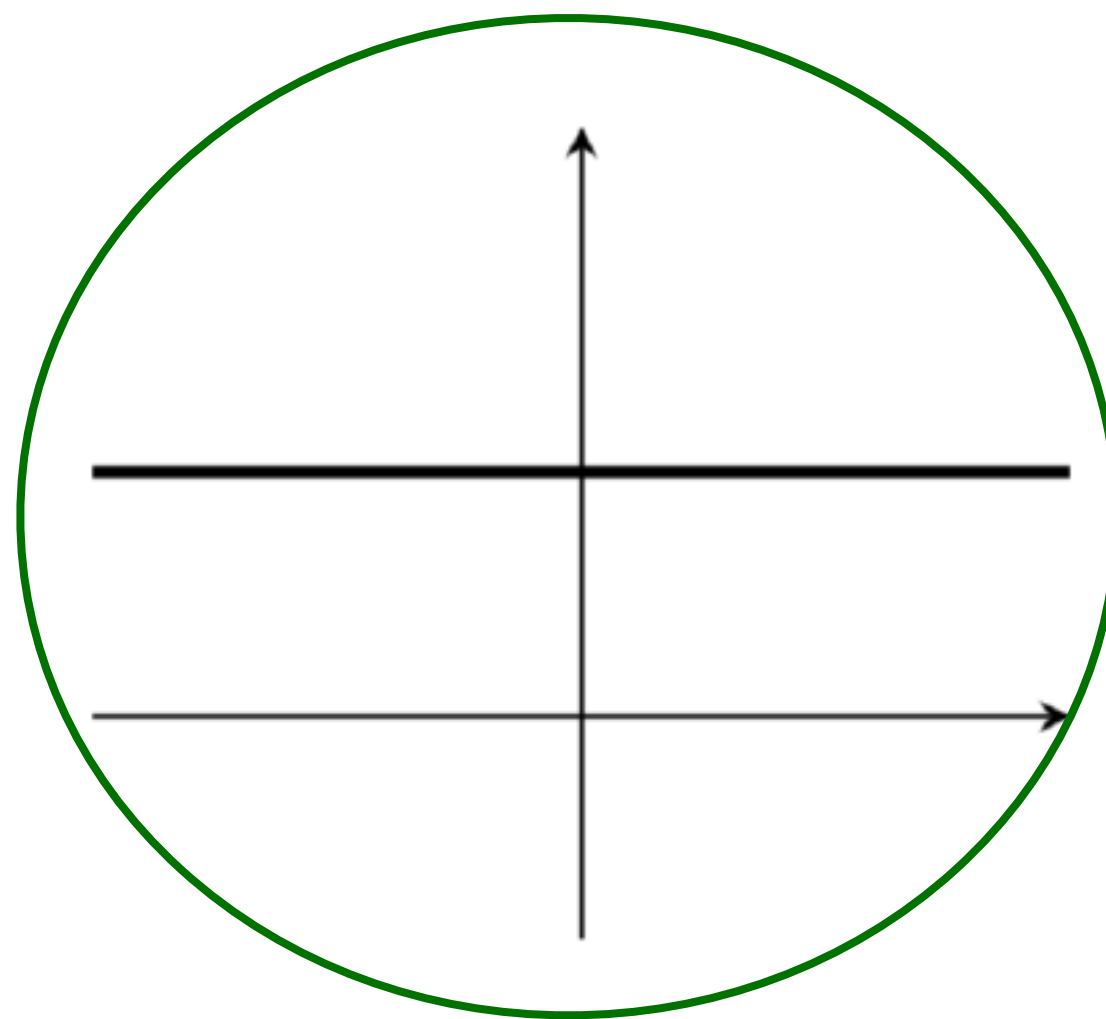
Which of these sets are convex?





# Convexity quiz

Which of these functions are convex?



# How the convexity can help?

Let  $f: \mathcal{C} \rightarrow \mathbb{R}$  be a convex function defined on a convex set  $\mathcal{C}$ .

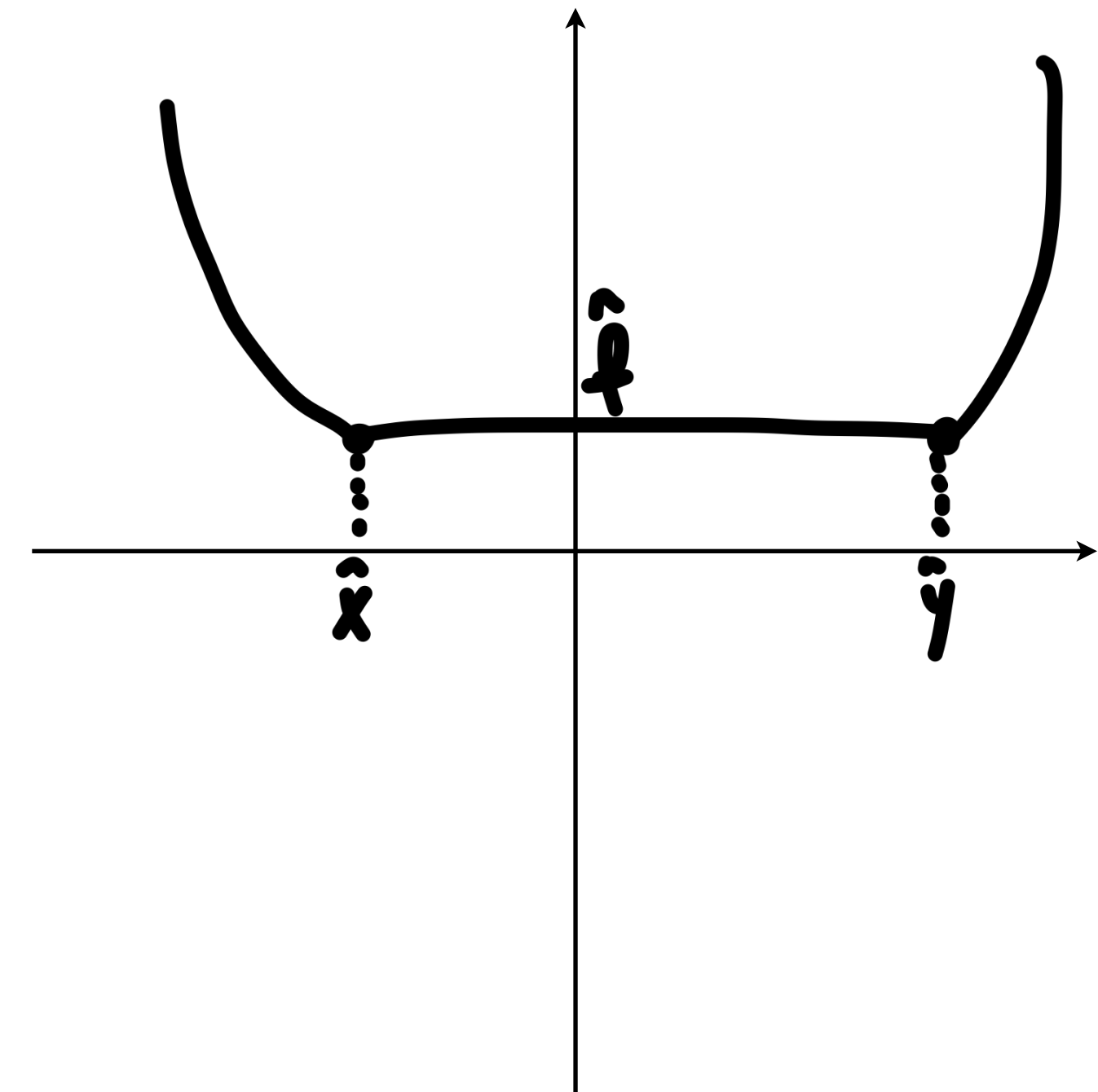
**Question: Could a convex function have two minimisers?**

Let  $\mathbf{x}, \mathbf{y}$  be two different minima,  $f(\mathbf{x}) = f(\mathbf{y}) = \hat{f} = \min f$ . Then  $\forall \lambda \in (0,1)$

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) = \lambda \hat{f} + (1 - \lambda) \hat{f} = \hat{f} \Rightarrow$$

$$\forall \lambda \in [0,1] \quad f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) = \hat{f}$$

**Answer: Convex function can have more than one minimiser if it is a constant on some interval.**



# How the convexity can help?

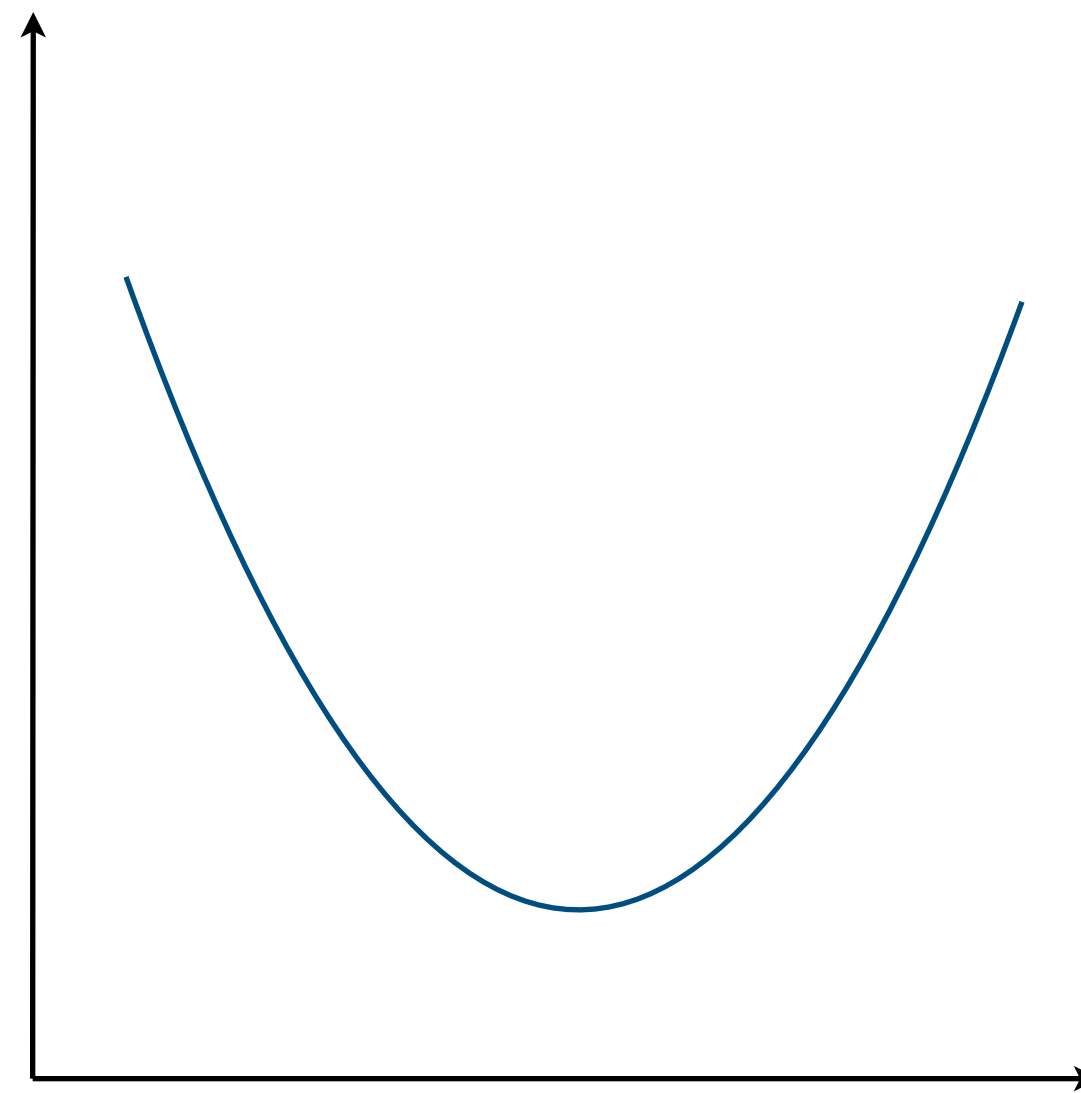
**Lemma:** Let  $f: \mathcal{C} \subset \mathbb{R}^n \rightarrow \mathbb{R}$  is a convex function that is bounded from below and also differentiable. Then an argument  $\hat{\mathbf{w}}$  satisfies  $\nabla f(\hat{\mathbf{w}}) = 0$  if and only if  $\hat{\mathbf{w}}$  is a global minimiser.

## Motivation:

We need to prove 2 results

If  $\hat{\mathbf{w}}$  is a minimiser  $\Rightarrow \nabla f(\hat{\mathbf{w}}) = 0$ .

If  $\nabla f(\hat{\mathbf{w}}) = 0 \Rightarrow \hat{\mathbf{w}}$  is a minimiser.



# How the convexity can help?

**Lemma:** Let  $f: \mathcal{C} \subset \mathbb{R}^n \rightarrow \mathbb{R}$  is a convex function that is bounded from below and also differentiable. Then an argument  $\hat{\mathbf{w}}$  satisfies  $\nabla f(\hat{\mathbf{w}}) = 0$  if and only if  $\hat{\mathbf{w}}$  is a global minimiser.

## Proof 1D: minimiser $\Rightarrow$ tangent

Let  $\hat{w}$  be a minimiser. What can we say about  $f'(\hat{w})$ ?

$$f'(\hat{w}) = \lim_{w \rightarrow \hat{w}^+} \frac{\overbrace{f(w) - f(\hat{w})}^{\geq 0}}{\underbrace{w - \hat{w}}_{>0}} \geq 0 \quad f'(\hat{w}) = \lim_{w \rightarrow \hat{w}^-} \frac{\overbrace{f(w) - f(\hat{w})}^{\geq 0}}{\underbrace{w - \hat{w}}_{<0}} \leq 0 \Rightarrow f'(\hat{w}) = 0$$



# How the convexity can help?

**Lemma:** Let  $f: \mathcal{C} \subset \mathbb{R}^n \rightarrow \mathbb{R}$  is a convex function that is bounded from below and also differentiable. Then an argument  $\hat{\mathbf{w}}$  satisfies  $\nabla f(\hat{\mathbf{w}}) = 0$  if and only if  $\hat{\mathbf{w}}$  is a global minimiser.

## Proof 1D: tangent $\Rightarrow$ minimiser

Let  $\hat{w}$  be an argument such that  $f'(\hat{w}) = 0$ . For any  $\lambda \in [0, 1]$  and  $w \in \mathcal{C}$

$$f(\lambda w + (1 - \lambda) \hat{w}) \leq \lambda f(w) + (1 - \lambda) f(\hat{w})$$

$$f(\hat{w} + \lambda (w - \hat{w})) - f(\hat{w}) \leq \lambda (f(w) - f(\hat{w}))$$

$$f(w) - f(\hat{w}) \geq \frac{f(\hat{w} + \lambda (w - \hat{w})) - f(\hat{w})}{\lambda} \rightarrow_{\lambda \rightarrow 0} f'(\hat{w}) (w - \hat{w}) = 0$$

# How the convexity can help?

---

## Conclusions:

- Solving a regression problem is equivalent to a minimisation of cost function.
- If the cost function is convex, then the above is equivalent to finding a point of zero gradient.
- The solution to normal equation exists.
- This solution is unique except of very special cases

**Question: What is left?**

**Answer: We need to show the convexity of a cost function**