


# Theme: Curve Fitting

**Mark van der Wilk**

Department of Computing  
Imperial College London

@markvanderwilk  
m.vdwilk@imperial.ac.uk

October 3, 2022

# Overview

What is regression?

Regression as Minimising a Loss

A Statistical View on Regression

Conclusion

# Curve Fitting (Regression) Examples

We will be considering *curve fitting* or *supervised learning*.

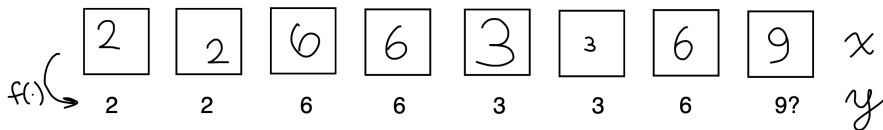
- ▶ Given a dataset of  $N$  examples of inputs and outputs...
- ▶ predict what the output will be for a new input.

# Curve Fitting (Regression) Examples

We will be considering *curve fitting* or *supervised learning*.

- ▶ Given a dataset of  $N$  examples of inputs and outputs...
- ▶ predict what the output will be for a new input.

**Image classification.** Inputs  $\in \mathbb{R}^D$ , outputs  $\in \mathbb{N}$ :

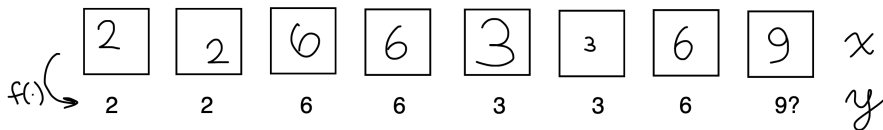


# Curve Fitting (Regression) Examples

We will be considering *curve fitting* or *supervised learning*.

- ▶ Given a dataset of  $N$  examples of inputs and outputs...
- ▶ predict what the output will be for a new input.

**Image classification.** Inputs  $\in \mathbb{R}^D$ , outputs  $\in \mathbb{N}$ :

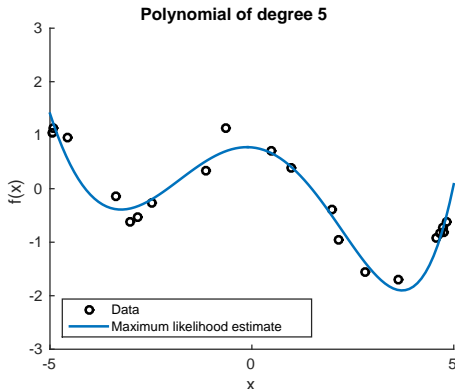


**Translation.** Inputs  $\in \bigcup_{\ell=1}^{\infty} \mathbb{N}^{\ell}$ , outputs  $\in \bigcup_{k=1}^{\infty} \mathbb{N}^k$ :

Wiskunde is belangrijk.	→	Mathematics is important.
Dutch		English

# Regression Example

Curve fitting in 1D. Inputs  $\in \mathbb{R}$ , outputs  $\in \mathbb{R}$ :



# Curve fitting

*“All the impressive achievements of deep learning  
amount to just curve fitting.”*

— Judea Pearl

# Overview

What is regression?

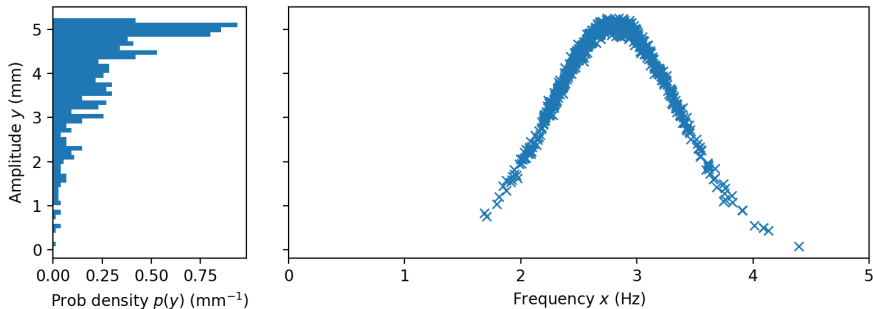
Regression as Minimising a Loss

A Statistical View on Regression

Conclusion



# Regression Example



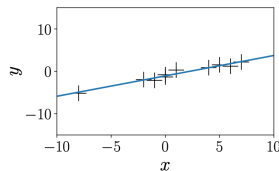
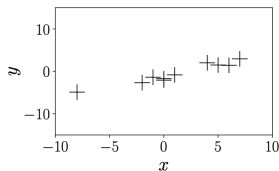
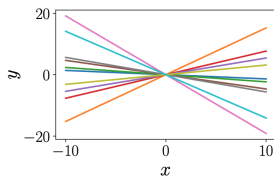
For some observed  $x$ , the world is generating data from  $\pi(y|x)$ .

We can choose two possible goals for regression:

- ▶ Loss view: Find a function  $f(x)$  that goes “near” outputs  $y$ .
- ▶ Stats view: Match a statistical model  $p(y|x, \theta)$  to  $\pi(y|x)$ .

# Loss view: Good and bad functions

We now have many functions that we can choose from:

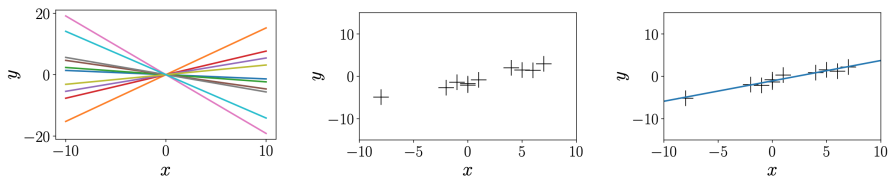


Left: example functions. Middle: Training set. Right: A good fit.

Source: Mathematics for Machine Learning book.

# Loss view: Good and bad functions

We now have many functions that we can choose from:



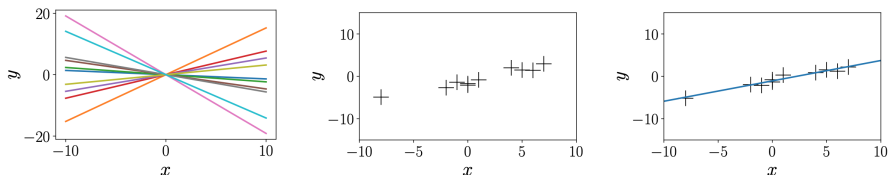
Left: example functions. Middle: Training set. Right: A good fit.

Source: Mathematics for Machine Learning book.

Little question: Which function do we pick?

# Loss view: Good and bad functions

We now have many functions that we can choose from:



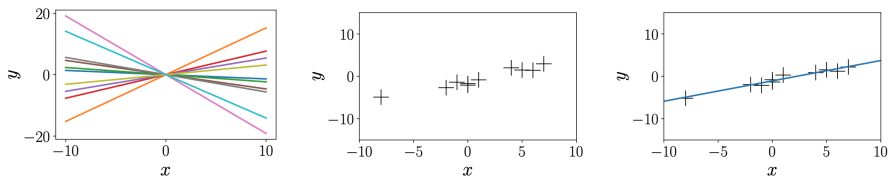
Left: example functions. Middle: Training set. Right: A good fit. Source: Mathematics for Machine Learning book.

## Little question: Which function do we pick?

- Need to define what good and bad functions are. Good functions have  $f(x_i, \theta^*) \approx y_i$ .

# Loss view: Good and bad functions

We now have many functions that we can choose from:



Left: example functions. Middle: Training set. Right: A good fit.

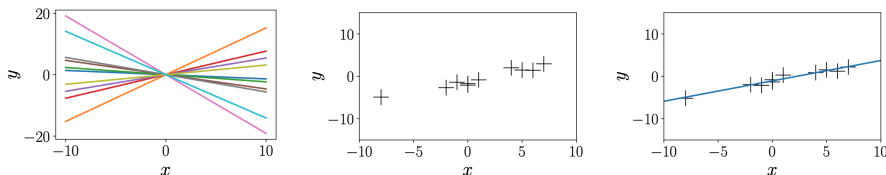
Source: Mathematics for Machine Learning book.

## Little question: Which function do we pick?

- ▶ Need to define what good and bad functions are. Good functions have  $f(x_i, \theta^*) \approx y_i$ .
- ▶ Define a **loss function**, e.g.,  $L(\theta) = \sum_{i=1}^N (y_i - f(x_i, \theta))^2$

# Loss view: Good and bad functions

We now have many functions that we can choose from:



Left: example functions. Middle: Training set. Right: A good fit.

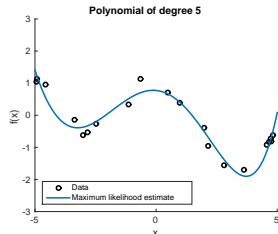
Source: Mathematics for Machine Learning book.

## Little question: Which function do we pick?

- ▶ Need to define what good and bad functions are. Good functions have  $f(x_i, \theta^*) \approx y_i$ .
- ▶ Define a **loss function**, e.g.,  $L(\theta) = \sum_{i=1}^N (y_i - f(x_i, \theta))^2$
- ▶ Choose a good function, i.e.  $\theta^* = \operatorname{argmin}_{\theta} L(\theta)$

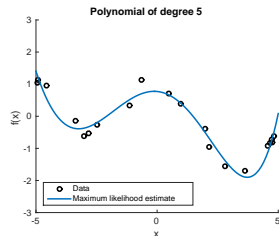
# Curve Fitting: Representing functions

Little question: How do we represent functions?



# Curve Fitting: Representing functions

Little question: How do we represent functions?

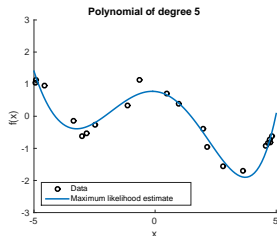


- We need a *collection* of functions from which to pick a good one.



# Curve Fitting: Representing functions

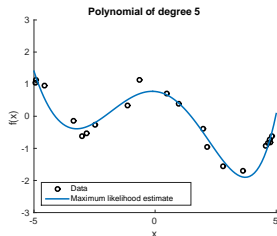
Little question: How do we represent functions?



- ▶ We need a *collection* of functions from which to pick a good one.
- ▶ **Parameterise** a set of functions, i.e. take some numbers  $\theta$  that map to a function.

# Curve Fitting: Representing functions

Little question: How do we represent functions?



- ▶ We need a *collection* of functions from which to pick a good one.
- ▶ **Parameterise** a set of functions, i.e. take some numbers  $\theta$  that map to a function.

For example, linear or polynomial functions:

$$f_{\theta}(x) = a \cdot x + b, \quad \theta = \begin{bmatrix} a \\ b \end{bmatrix}, \quad (1)$$

$$f_{\theta}(x) = a \cdot x^3 + b \cdot x^2 + c \cdot x + d, \quad \theta = \begin{bmatrix} a & b & c & d \end{bmatrix}^T. \quad (2)$$

# Overview

What is regression?

Regression as Minimising a Loss

A Statistical View on Regression

Conclusion

# Maximum Likelihood Estimation

Revision from 50008: Probability & Statistics

- ▶ Model is a probability distribution on data:  $p(y|\theta)$
- ▶ For an observed dataset (fixed), we can evaluate the probability assigned to it for different  $\theta$
- ▶ This defines the likelihood  $\ell(\theta) = p(y|\theta)$

# Maximum Likelihood Estimation

Revision from 50008: Probability & Statistics

- ▶ Model is a probability distribution on data:  $p(y|\theta)$
- ▶ For an observed dataset (fixed), we can evaluate the probability assigned to it for different  $\theta$
- ▶ This defines the likelihood  $\ell(\theta) = p(y|\theta)$

Maximum likelihood does:

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \ell(\theta) = \underset{\theta}{\operatorname{argmax}} \log \ell(\theta) \quad (3)$$

# Likelihood for Linear Regression

Assume:

- ▶ Gaussian deviations from the function:

$$p(y_n|x_n, \theta) = \mathcal{N}(y_n; f_{\theta}(x_n), \sigma^2) \quad (4)$$

- ▶ Independent deviations between datapoints. So denoting  $y \in \mathbb{R}^N$ ,  $x \in \mathbb{R}^N$  for  $N$  datapoints, we get the likelihood:

$$p(y|x, \theta) = \prod_{n=1}^N \mathcal{N}(y_n; f_{\theta}(x_n), \sigma^2) \quad (5)$$

- ▶ You will show that this is equivalent to the loss view (exercises).

# Overview

What is regression?

Regression as Minimising a Loss

A Statistical View on Regression

Conclusion

# Curve Fitting Summary

- ▶ Training data, e.g.,  $N$  pairs  $(x_i, y_i)$  of inputs  $x_i$  and observations  $y_i$
- ▶ **Parameterise** functions as  $f(\mathbf{x}_i, \boldsymbol{\theta})$
- ▶ **Training the model** means finding parameters  $\boldsymbol{\theta}^*$ , such that
  - ▶  $f(\mathbf{x}_i, \boldsymbol{\theta}^*) \approx y_i$  (loss is minimised)
  - ▶  $p(y_n|x_n, \boldsymbol{\theta}) \approx \pi(y_n|x_n, \boldsymbol{\theta})$  (max likelihood)



# Curve Fitting Summary

- ▶ Training data, e.g.,  $N$  pairs  $(x_i, y_i)$  of inputs  $x_i$  and observations  $y_i$
- ▶ **Parameterise** functions as  $f(x_i, \theta)$
- ▶ **Training the model** means finding parameters  $\theta^*$ , such that
  - ▶  $f(x_i, \theta^*) \approx y_i$  (loss is minimised)
  - ▶  $p(y_n|x_n, \theta) \approx \pi(y_n|x_n, \theta)$  (max likelihood)
- ▶ Not discussed: How to find  $\theta^*$

