



UNIVERSITY  
OF AMSTERDAM

# Machine Learning for Physics and Astronomy

Juan Rojo

VU Amsterdam & Theory group, Nikhef

***Natuur- en Sterrenkunde BSc (Joint Degree), Honours Track***  
***Lecture 2, 14/09/2021***

# Lecture outline

- ✿ Course logistics: Recap
- ✿ Why machine learning is difficult: Recap
- ✿ Deterministic optimisation algorithms: gradient descent
- ✿ Stochastic optimisation algorithms: genetic algorithms

# **Logistics Recap**

# Course schedule

Week	Monday	Tuesday	Wednesday	Thursday	
36 6/9 10/9		11-13 HC: Lecture dr. Juan Rojo Chacon	SP G2.13	9-11 LC: Practical session dr. Juan Rojo Chacon	SP B0.202
37 13/9 17/9		11-13 HC: Lecture dr. Juan Rojo Chacon	SP G2.13	9-11 LC: Practical session dr. Juan Rojo Chacon	SP B0.202
38 20/9 24/9		11-13 HC: Lecture dr. Juan Rojo Chacon	SP G2.13	9-11 LC: Practical session dr. Juan Rojo Chacon	SP B0.202
39 27/9 1/10		11-13 HC: Lecture dr. Juan Rojo Chacon	SP G2.13	9-11 LC: Practical session dr. Juan Rojo Chacon	SP B0.202
40 4/10 8/10		11-13 HC: Lecture dr. Juan Rojo Chacon	SP G2.13	9-11 LC: Practical session dr. Juan Rojo Chacon	SP B0.202
41 11/10 15/10		11-13 HC: Lecture dr. Juan Rojo Chacon	SP G2.13	9-11 LC: Practical session dr. Juan Rojo Chacon	SP B0.202
42 18/10 22/10		11-13 HC: Lecture dr. Juan Rojo Chacon	SP G2.13	9-11 LC: Practical session dr. Juan Rojo Chacon	SP B0.202

⌚ **7 lectures** (Tuesdays 11am) and **7 tutorials** (Thursdays 9am)

⌚ All course activities will be **on campus**. Recordings of the zoom lectures of the **2020-2021 academic year** are available on Canvas.

⌚ Additional discussion sessions can be scheduled upon request

# Course logistics

- 💡 All the **course material** (slides, recordings, notebooks for the tutorials, other resources) will be made available via the **Canvas page of the course**, which will also be used for the **main announcements** as well as for the submission of **homework assignments**
- 💡 This material will also be available via the course **GitHub repository**, which will always contain the most updated versions of the tutorial's codes

**<https://github.com/LHCfitNikhef/ML4PA>**

note that this repo can be directly linked to e.g. **Google's Collab** to run the codes on the cloud rather than locally

- 💡 Furthermore, we have created a **Discord server** to streamline the communications between the instructors and the students, and to facilitate discussions during the lectures. You should all have received an invitation to join the course channel, else get in touch

# Course evaluation

- Students should write a short report (around 6 pages) about a specific **application of Machine Learning algorithms** to a physics or astronomy problem that they find interesting. This report and the subsequent presentation count up to **60%** of the course grade.

You need to submit this report via Canvas by **Monday 1st November**

*the recorded presentations of the 2020/2021 cohort are available on Canvas for inspiration*

- In addition, after each tutorial we will propose some short **homework assignments** that need to be completed and submitted via Canvas. These homework assignments count up to **40%** of the final course grade

Homework assignments must be handed in **within one week** following the tutorial (see Canvas for specific deadlines)

# Tutorials

the hands-on tutorials will allow you to familiarise with the machine learning concepts presented in the lectures by means of **practical examples**

- 💡 The **Python notebooks** with the course tutorials can be either run locally (if you have an updated python installation) or on the cloud by uploading the notebook in **Google Collab**
- 💡 Please note that some tutorials require the download of heavy datasets, make sure this is done before the tutorial
- 💡 Due to time constraints we can only offer limited assistance with **software installation** - please make sure you can install and execute the notebooks beforehand, and get back to us via Discord if there is any trouble
- 💡 The tutorial sessions will be coordinated by Tommaso Giani and Ryan van Mastrigt, the course TAs. Please contact them for any **questions related to the course tutorials**

The tutorials will include both **basic, intermediate, and advanced** topics concerning ML algorithms - it is up to you to decide how far you want to push!

# Guest lectures 2020/2021

In the previous year we had several guest lecturers presenting state-of-the-art ML applications for physics and astronomy problems within our community

- 🎧 Lecture 4: **Dr. Atul Chhotray** (applications to Astronomy)
- 🎧 Lecture 5: **Dr. Christoph Weniger** (applications to Dark Matter)
- 🎧 Lecture 6: **Dr. Sascha Caron** (applications for High-Energy Physics)
- 🎧 Lecture 7: **Dr. Tristan Bereau** (applications to Condensed Matter)

Slides and recordings of these guest lectures are also **available on Canvas**:  
students are encouraged to watch them as extra study resources

# Lecture 1 Recap

# Supervised learning

problems in **Supervised Machine Learning** are defined by the following ingredients:

(1) **Input dataset:**  $\mathcal{D} = (X, Y)$

array of **independent**  
variables

array of **dependent**  
variables

$$X = (x_1, x_2, \dots, x_N)$$

$$Y = (y_1, y_2, \dots, y_N)$$

$$x_i = (x_{i,1}, x_{i,2}, \dots, x_{i,p})$$

*each independent variable contains  $p$  features*

# Supervised learning

problems in **Supervised Machine Learning** are defined by the following ingredients:

(1) **Input dataset:**  $\mathcal{D} = (X, Y)$

array of **independent**  
variables

$$X = (x_1, x_2, \dots, x_n)$$

$$x_i = (x_{i,1}, x_{i,2}, \dots, x_{i,p})$$

array of **dependent**  
variables

$$Y = (y_1, y_2, \dots, y_n)$$

each **independent variable contains  $p$  features**

(2) **Model:**

$$f(X, \theta)$$

mapping between dependent  
and independent variables

$$f : X \rightarrow Y$$

model parameters

$$\theta = (\theta_1, \theta_2, \dots, \theta_m)$$

the more complex the problem, the more flexible the model

# Supervised learning

problems in **Supervised Machine Learning** are defined by the following ingredients:

(1) **Input dataset:**  $\mathcal{D} = (X, Y)$

(2) **Model:**  $f(X, \theta)$

(3) **Cost function:**  $C(Y; f(X; \theta))$

The cost function measures how well the model (for a specific choice of its parameters) is able to **describe the input dataset**

*example of cost function for single dependent variable: sum of residuals squared*

$$C(Y; f(X; \theta)) = \frac{1}{n} \sum (y_i - f(x_i, \theta))^2$$

# Supervised learning

problems in **Supervised Machine Learning** are defined by the following ingredients:

(1) **Input dataset:**  $\mathcal{D} = (X, Y)$

(2) **Model:**  $f(X, \theta)$

(3) **Cost function:**  $C(Y; f(X; \theta))$

The cost function measures how well the model (for a specific choice of its parameters) is able to describe the input dataset

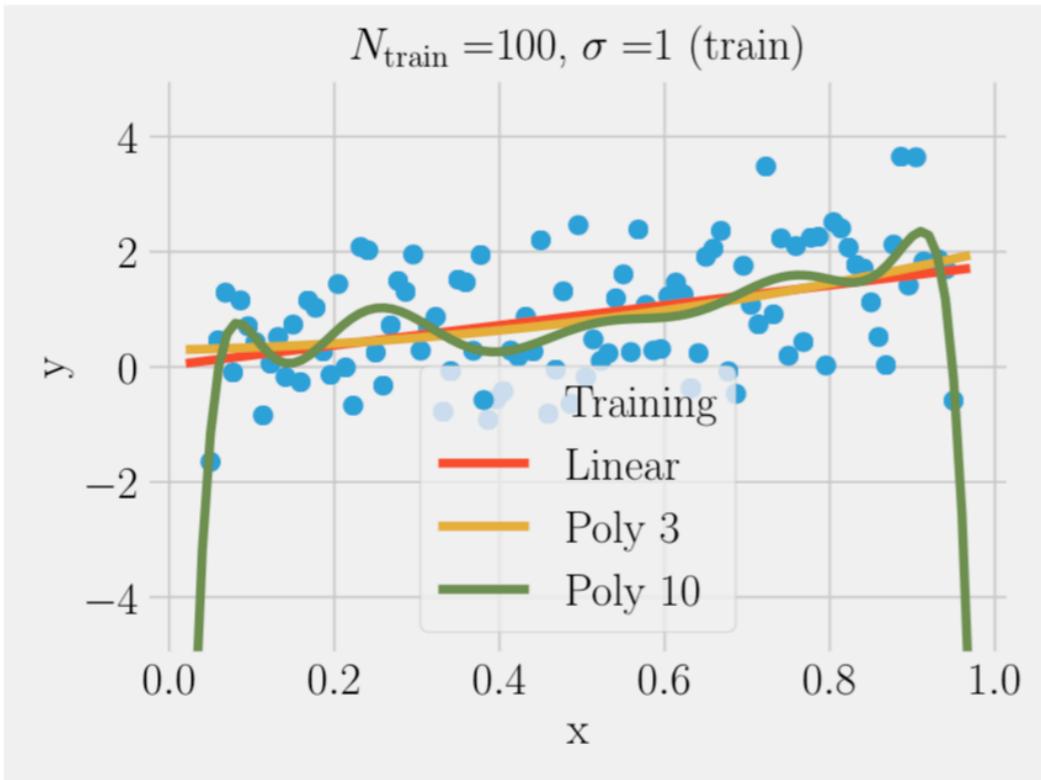
Fitting the model means determining the values of its parameters which **minimise the cost function**

$$\frac{\partial C(Y; f(X; \theta))}{\partial \theta_i} \Bigg|_{\theta=\theta_{\text{opt}}} = 0$$

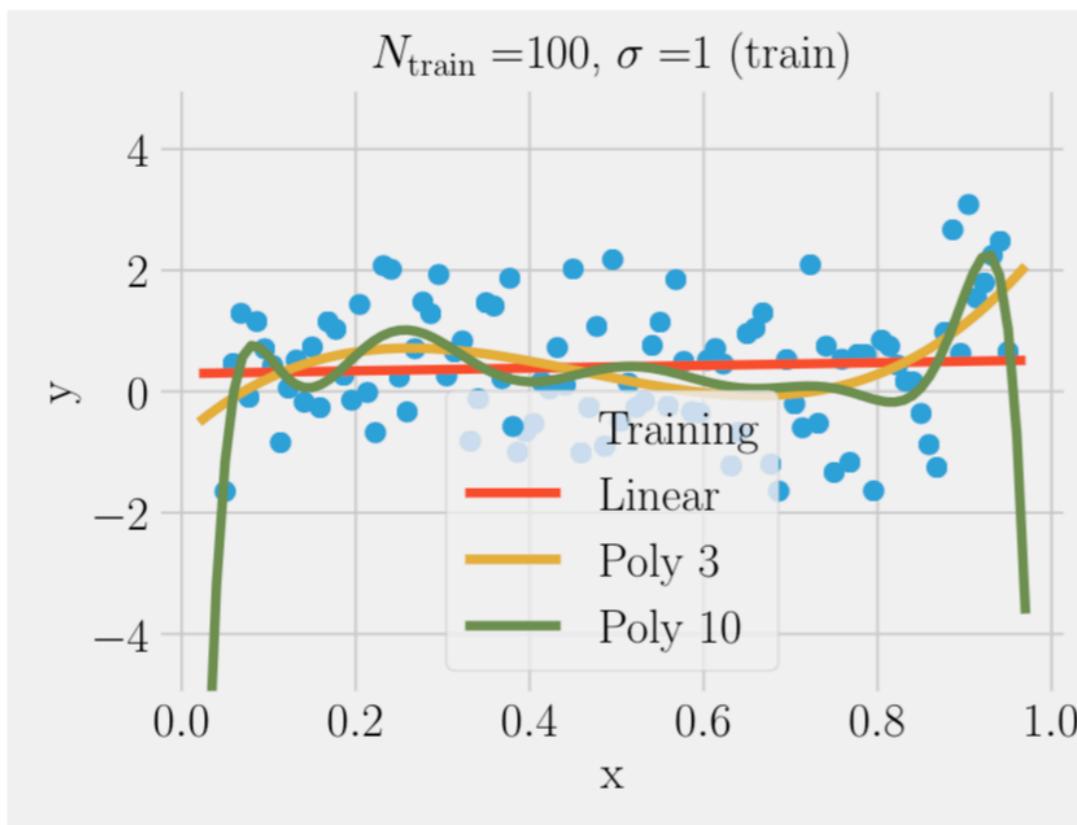
*easy right?*

# Model fitting

$$f(x) = 2x, \quad x \in [0,1]$$



$$f(x) = 2x - 10x^5 + 15x^{10}, \quad x \in [0,1]$$



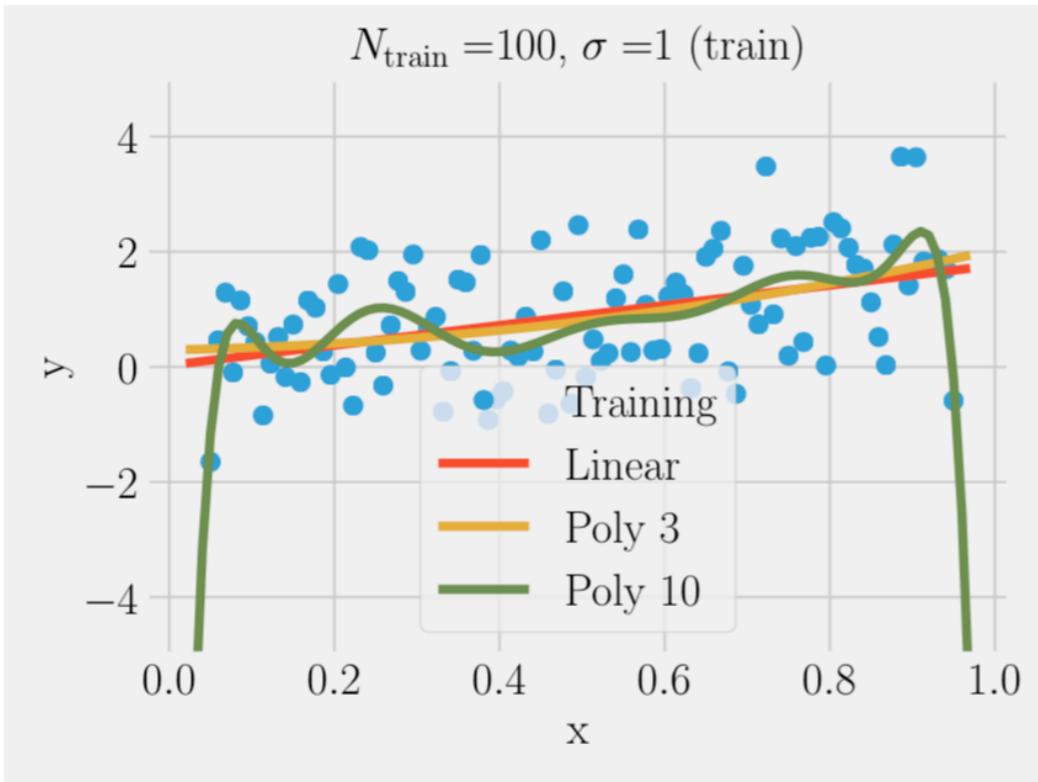
***What was the main take-home message of this exercise?***

***How do you expect these issues could affect realistic machine learning applications***

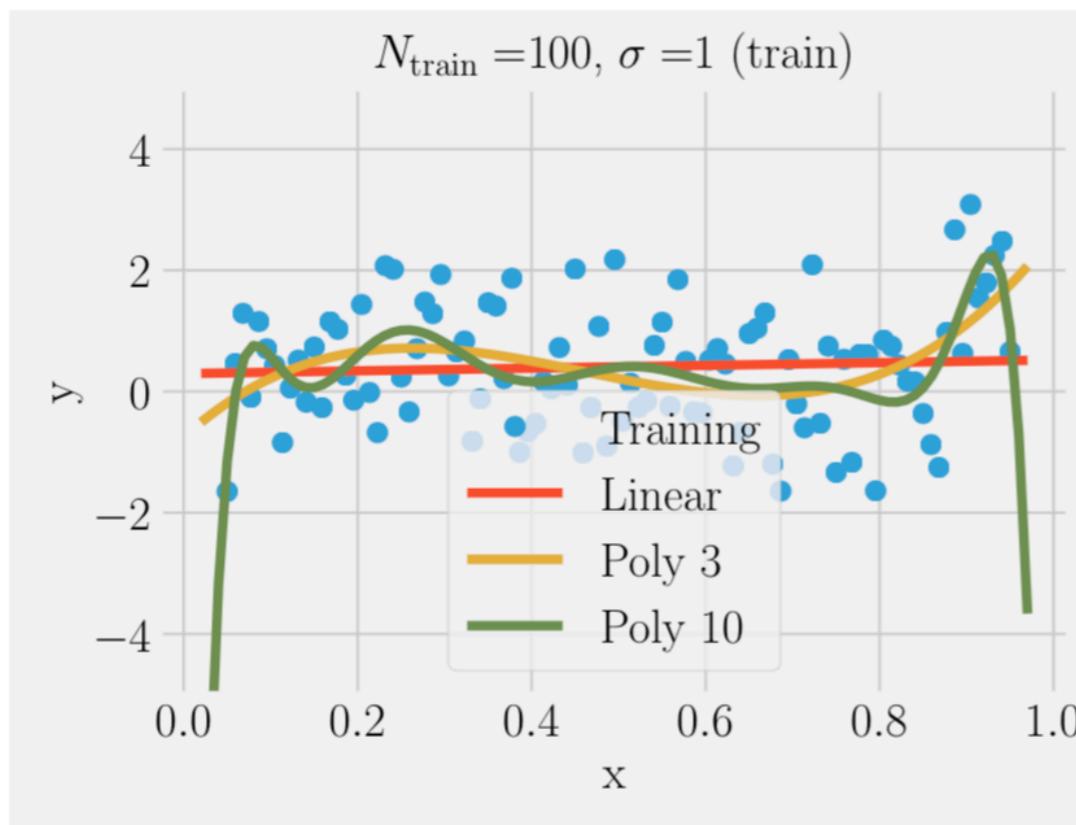
***e.g. self-driving cars***

# Model fitting

$$f(x) = 2x, \quad x \in [0,1]$$



$$f(x) = 2x - 10x^5 + 15x^{10}, \quad x \in [0,1]$$



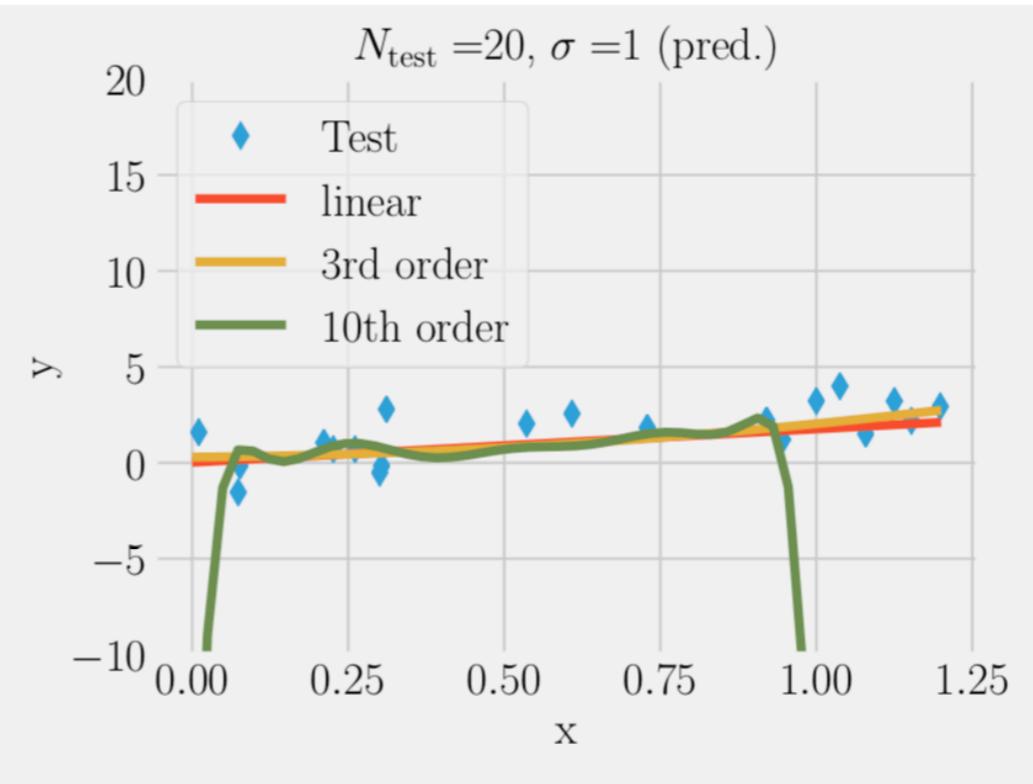
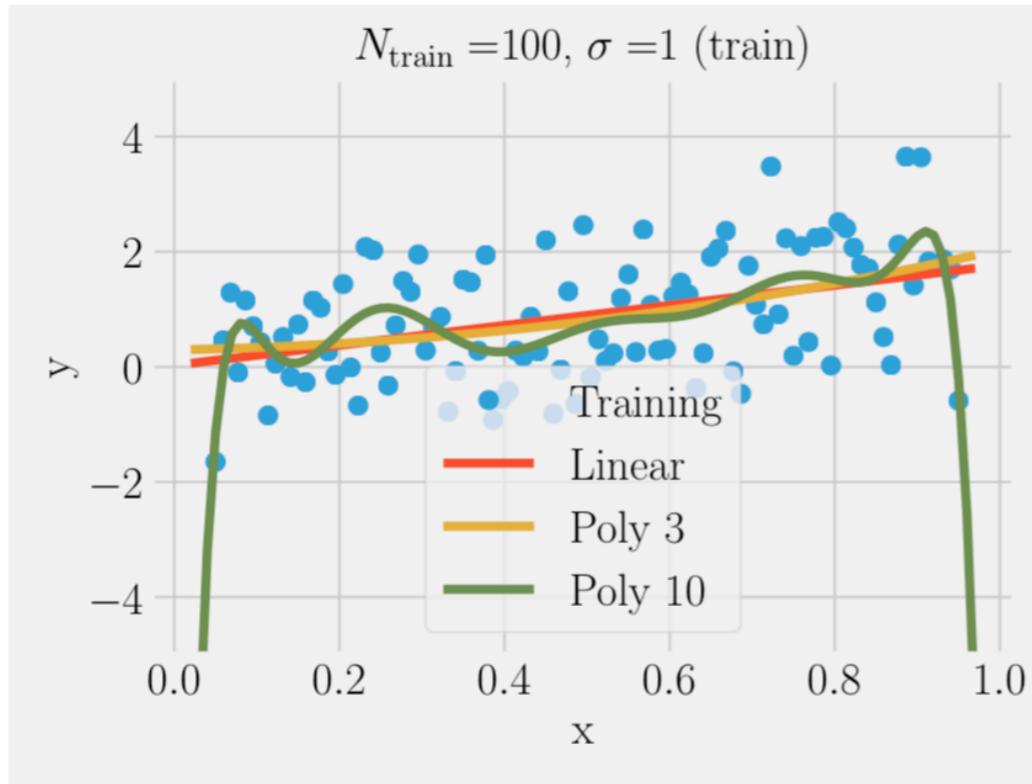
***What is the main take-home message of this exercise?***

***How do you expect these issues could affect realistic machine learning applications***

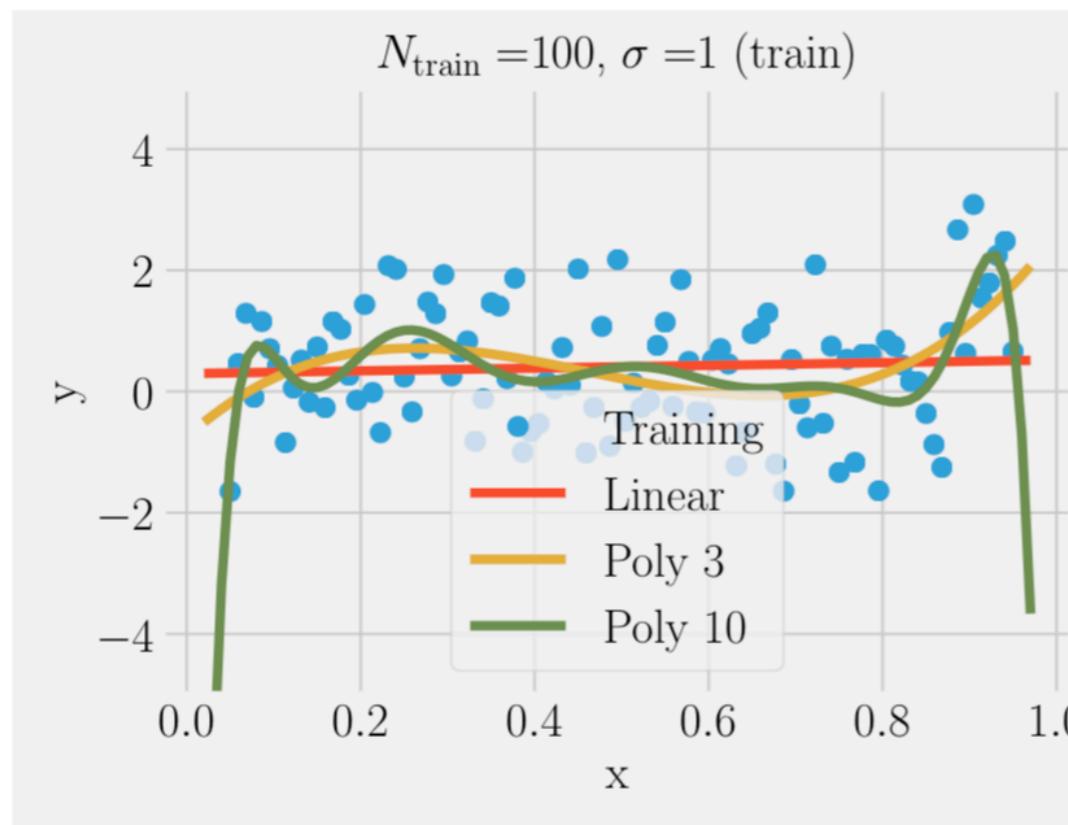
***Fitting is not Predicting!***

# Model fitting

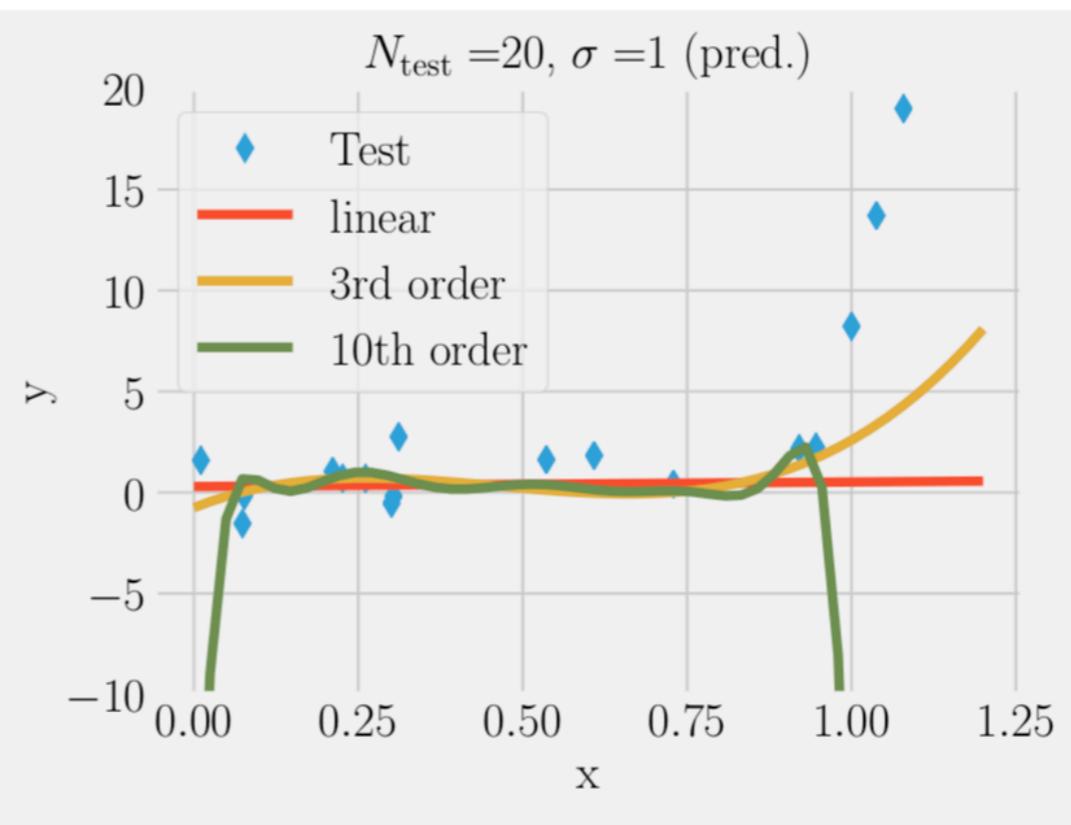
$$f(x) = 2x, \quad x \in [0,1]$$



$$f(x) = 2x - 10x^5 + 15x^{10}, \quad x \in [0,1]$$



*in the presence of noise, models with less complexity can exhibit improved predictive power*



# Why Machine Learning is difficult

- **Fitting existing data** is conceptually different from **making predictions about new data**
- Increasing the model complexity can improve the description of the training data but **reduce the predictive power** of the model due to overfitting, unless a suitable regularisation strategy is implemented
- For complex and/or small datasets, simple models can be better at prediction than complex models. The **“right” amount of complexity** cannot in general be determined from first principles
- It is difficult to **generalise** beyond the situations encountered in the training set: the model cannot learn what it has not seen
- Many problems that are **approachable in principle** can become **unfeasible in practice**, e.g. due to computational limitations, lack of convergence, instability .....

# **Supervised Learning: Deterministic Optimisation Strategies**

# Optimisation strategies

problems in **Supervised Machine Learning** are defined by the following ingredients:

**(1) Input dataset:**  $\mathcal{D} = (X, Y)$

**(2) Model:**  $f(X, \theta)$

**(3) Cost function:**  $C(Y; f(X; \theta))$

the model parameters are then found by **minimising the cost function**

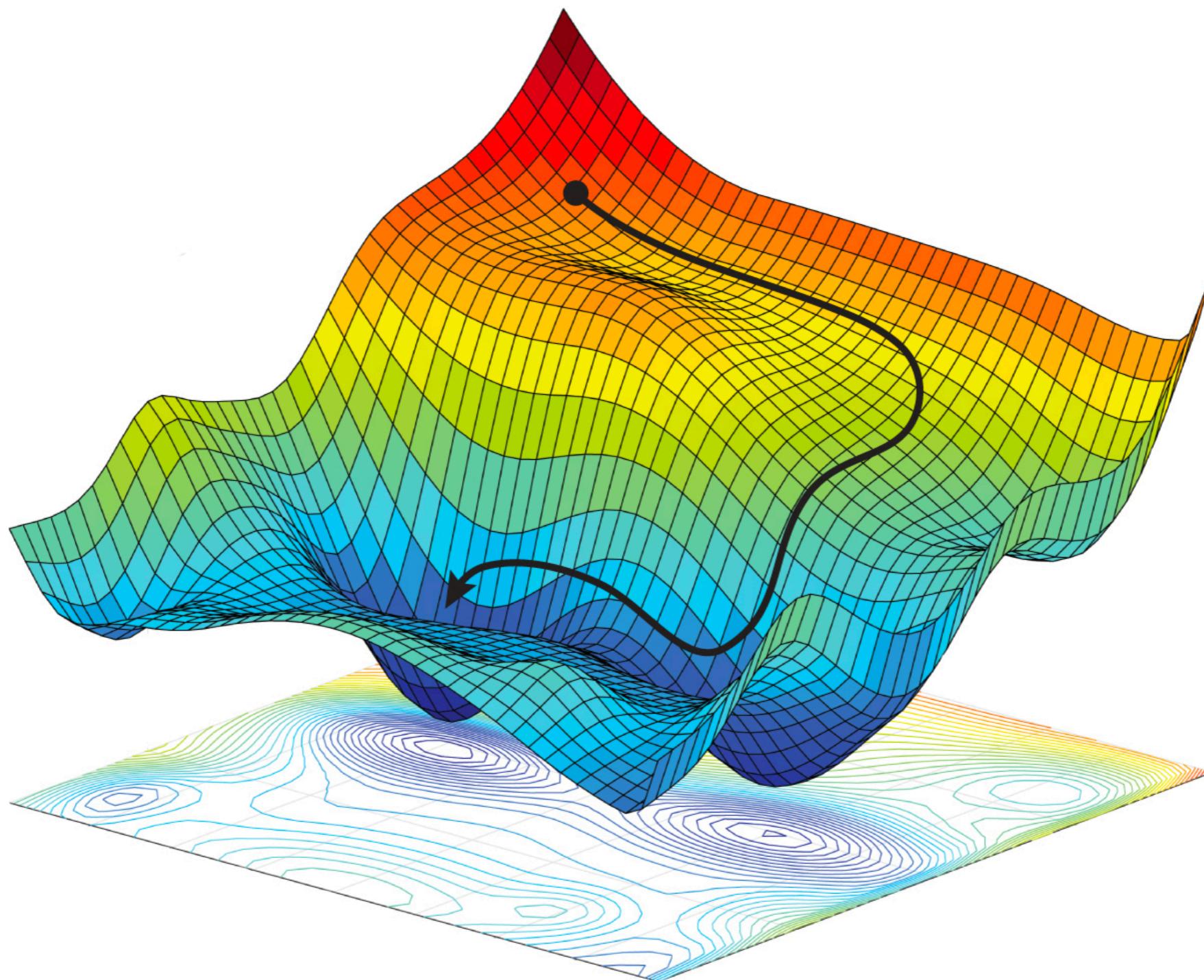
for simple models (e.g. polynomial regression), the solution of this minimisation problem can be found analytically. This is not possible with complex models in realistic applications. In the context of ML studies, there exist two main classes of **minimisers (optimisers)** that are frequently deployed:

💡 **Gradient Descent** and its generalisations

💡 **Genetic Algorithms** and its generalisations

# Gradient Descent

basic idea: iteratively adjust the model parameters in the direction where the **gradient of the cost function** is large and negative (**steepest descent direction**)



# Gradient descent

basic idea: iteratively adjust the model parameters in the direction where the **gradient of the cost function** is large and negative (**steepest descent direction**)

Our goal is thus to **minimise an error (cost) function** that can usually be expressed as

$$E(\theta) = \sum_{i=1}^n e_i(x_i; \theta)$$

e.g. in polynomial regression we had that

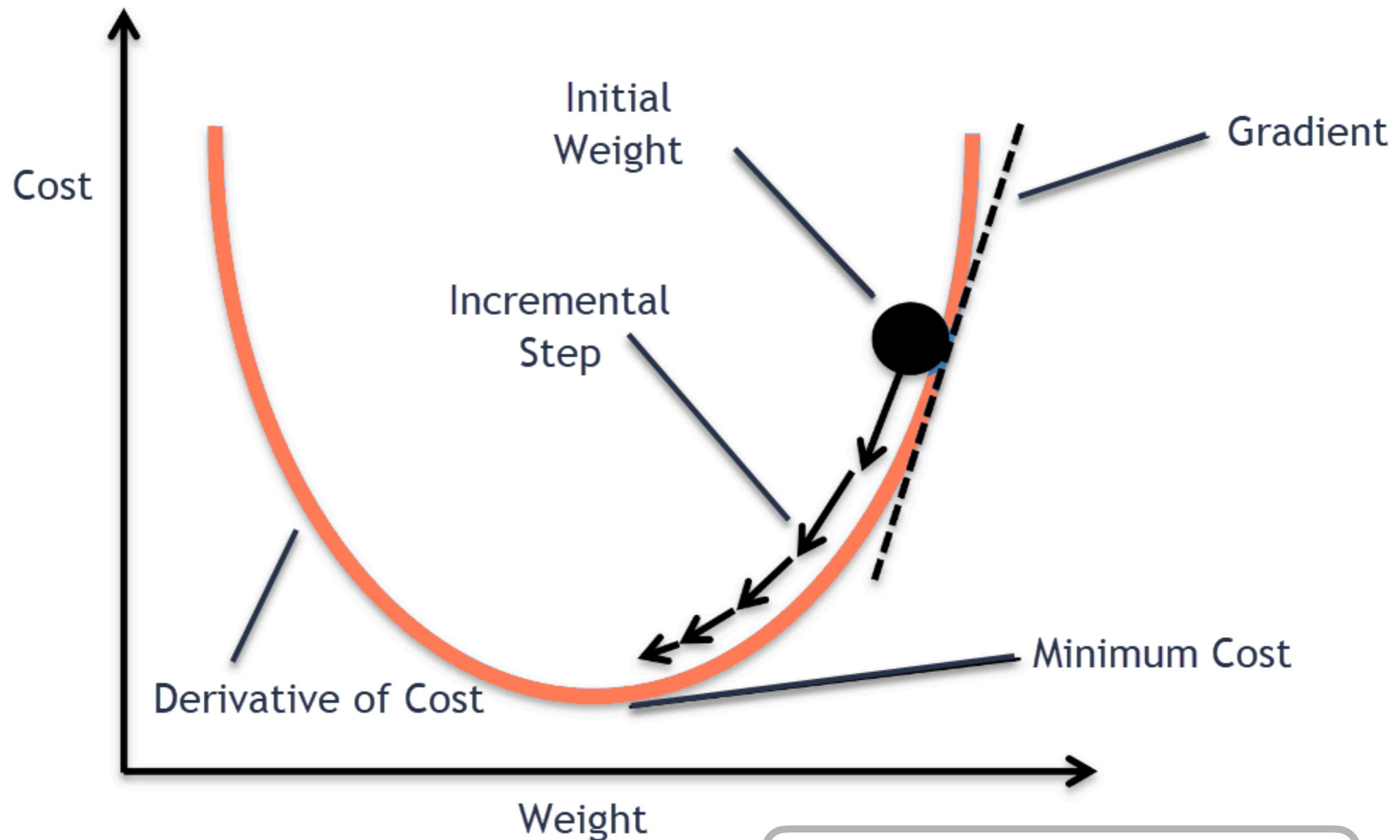
$$E(\theta) = \sum_{i=1}^n (y_i - f_\alpha(x_i; \theta_\alpha))^2 \quad \text{so} \quad e_i = (y_i - f_\alpha(x_i; \theta_\alpha))^2$$

and we will see that in **logistic regression** (classification problems)

$$E(\theta) = \sum_{i=1}^n (-y_i \ln \sigma(\mathbf{x}_i^T \theta) - (1 - y_i) \ln [1 - \sigma(\mathbf{x}_i^T \theta)])$$

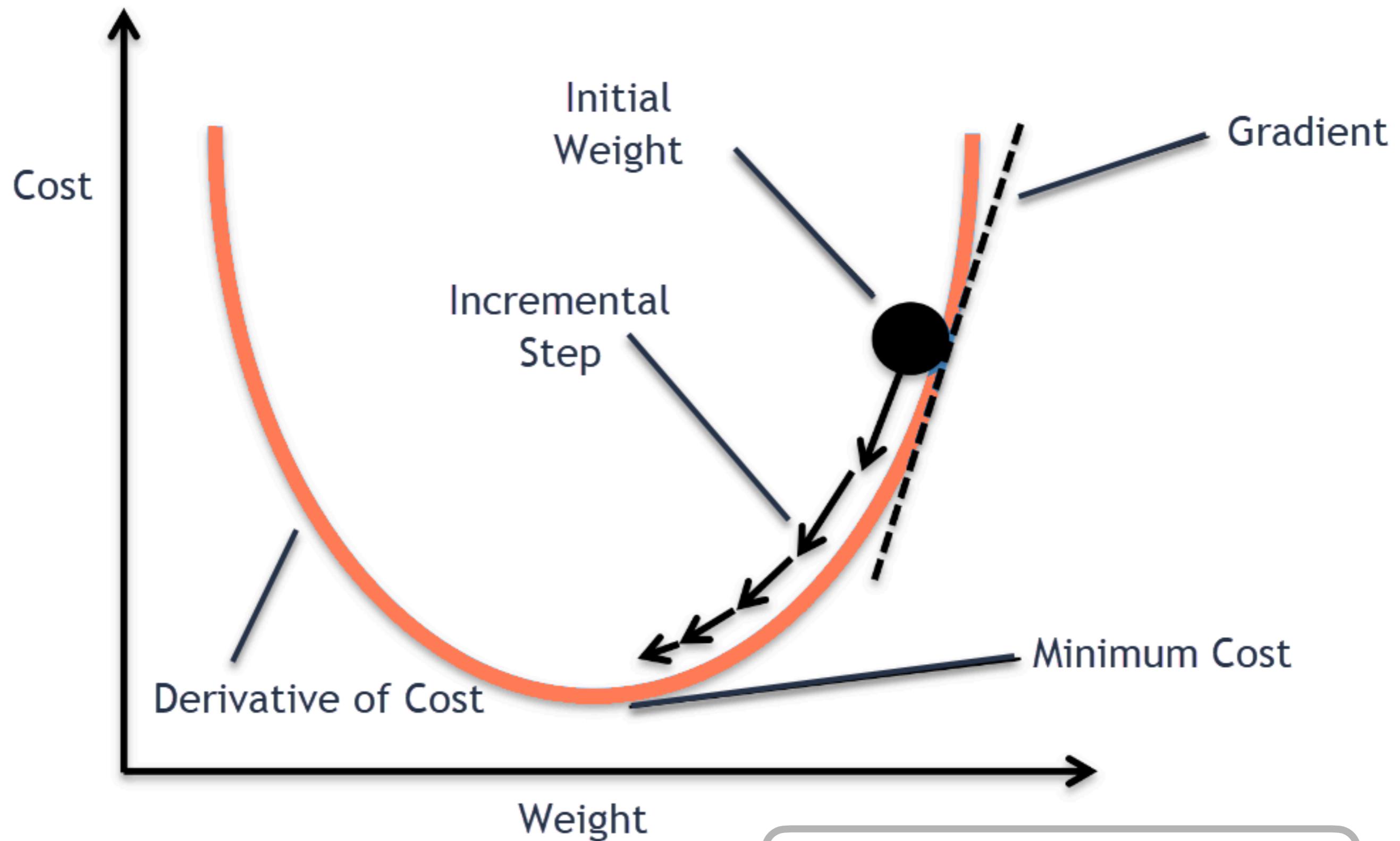
*aka the cross-entropy, relevant for categorisation*

# Gradient descent



*What do you think is key benefit of GD?*

# Gradient descent



*Only local information (gradient) required!*

# Gradient descent

starting from a suitable initial condition, GD **iteratively updates** the model parameters:

$$\mathbf{v}_t = \eta_t \nabla_{\theta} E(\boldsymbol{\theta}_t), \quad \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mathbf{v}_t$$

*learning rate*      *gradient of cost function*      *update of model parameters between iterations  $t$  and  $t+1$*

the learning rate determines the size of the step in the direction of the gradient

*What could go wrong when choosing the learning rate?*

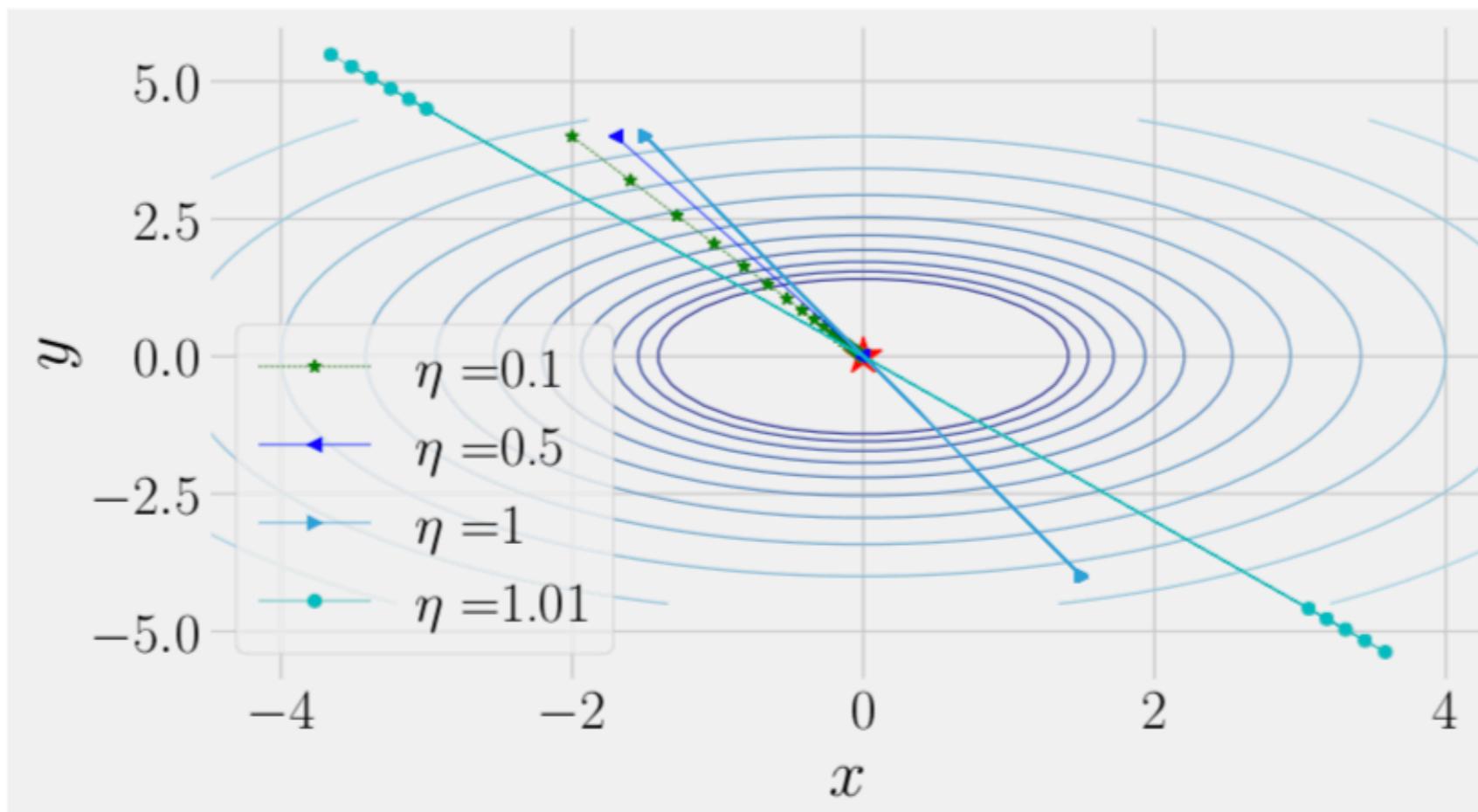
# Gradient descent

starting from a suitable initial condition, GD **iteratively updates** the model parameters:

$$\mathbf{v}_t = \eta_t \nabla_{\theta} E(\boldsymbol{\theta}_t), \quad \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mathbf{v}_t$$

*learning rate*      *gradient of cost function*      *update of model parameters between iterations  $t$  and  $t+1$*

the learning rate determines the size of the step in the direction of the gradient



*small  $\eta$ : guaranteed to find local minimum, at price of large number of iterations*

*large  $\eta$ : can overshoot minimum, problems of convergence*

*nb evaluating gradient can be very cpu-intensive!*

# Gradient descent

starting from a suitable initial condition, GD **iteratively updates** the model parameters:

$$\mathbf{v}_t = \eta_t \nabla_{\theta} E(\boldsymbol{\theta}_t), \quad \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mathbf{v}_t$$

*learning rate*      *gradient of cost function*      *update of model parameters between iterations  $t$  and  $t+1$*

compare GD with **Newton's method**: first Taylor-expand of cost function  
for a small change of the model parameters

$$E(\boldsymbol{\theta} + \mathbf{v}) \simeq E(\boldsymbol{\theta}) + \nabla_{\theta} E(\boldsymbol{\theta}) \cdot \mathbf{v} + \frac{1}{2} \mathbf{v}^T H(\boldsymbol{\theta}) \mathbf{v}$$

$$H_{ij}(\boldsymbol{\theta}) = \left\{ \frac{\partial^2 E(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} \right\}$$

*Hessian matrix (2nd derivatives)*

# Gradient descent

starting from a suitable initial condition, GD **iteratively updates** the model parameters:

$$\mathbf{v}_t = \eta_t \nabla_{\theta} E(\boldsymbol{\theta}_t), \quad \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mathbf{v}_t$$

*learning rate*      *gradient of cost function*      *update of model parameters between iterations  $t$  and  $t+1$*

compare GD with **Newton's method**: first Taylor-expand of cost function  
for a small change of the model parameters

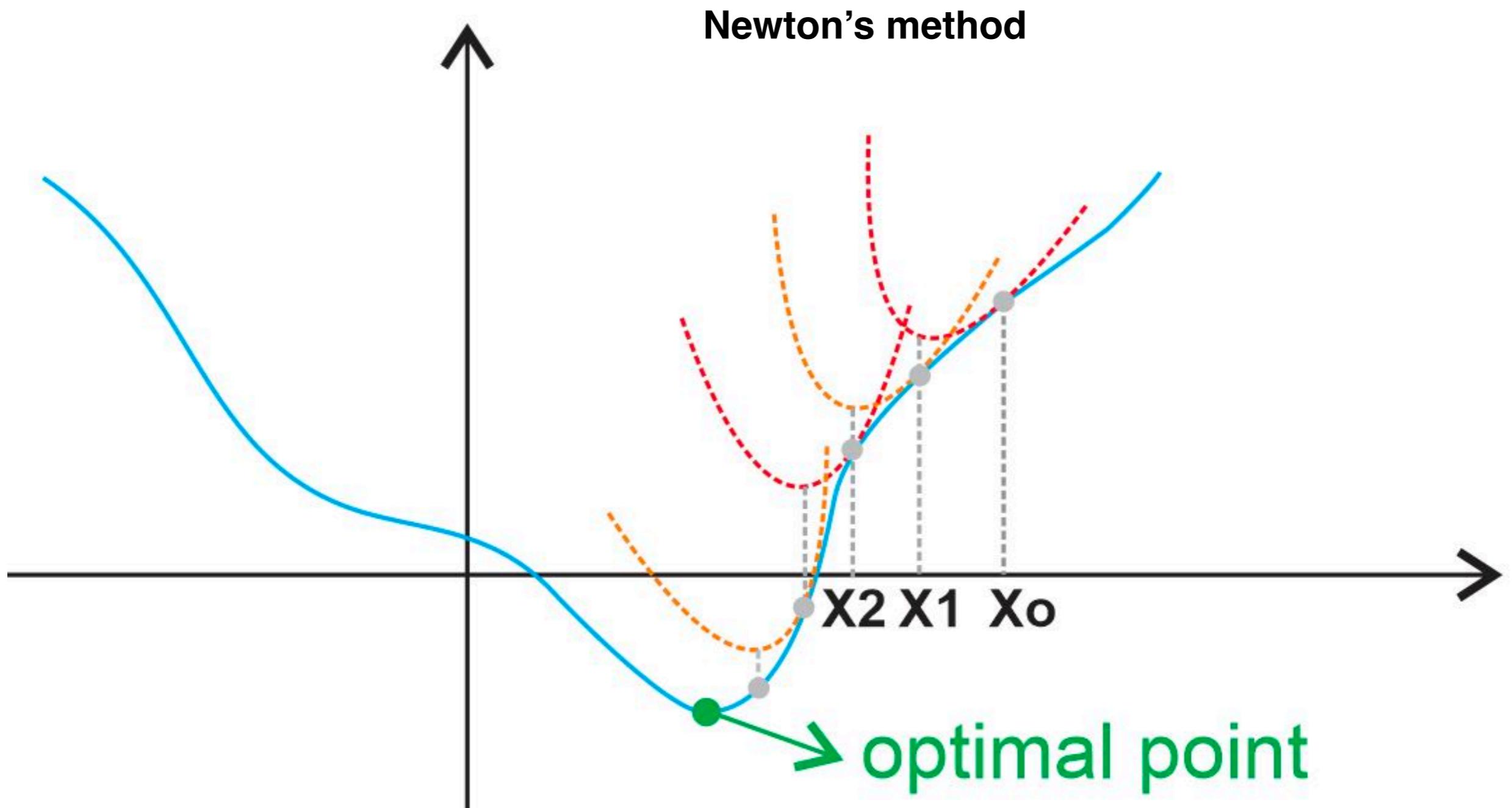
$$E(\boldsymbol{\theta} + \mathbf{v}) \simeq E(\boldsymbol{\theta}) + \nabla_{\theta} E(\boldsymbol{\theta}) \cdot \mathbf{v} + \frac{1}{2} \mathbf{v}^T H(\boldsymbol{\theta}) \mathbf{v}$$

and then differentiate with respect to the step requiring that the **linear term vanishes**

$$\left. \frac{\partial E(\boldsymbol{\theta} + \mathbf{v})}{\partial \mathbf{v}} \right|_{\mathbf{v}_{\text{opt}}} = 0 \longrightarrow \nabla_{\theta} E(\boldsymbol{\theta}) + H(\boldsymbol{\theta}) \mathbf{v}_{\text{opt}} = 0$$

$$\mathbf{v}_t = H^{-1}(\boldsymbol{\theta}_t) \cdot \nabla_{\theta} E(\boldsymbol{\theta}_t), \quad \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mathbf{v}_t$$

# Gradient descent



# Gradient descent

*Gradient Descent*  $\longrightarrow \mathbf{v}_t = \eta_t \nabla_{\theta} E(\theta_t), \quad \theta_{t+1} = \theta_t - \mathbf{v}_t$

*Newton's method*  $\longrightarrow \mathbf{v}_t = H^{-1}(\theta_t) \cdot \nabla_{\theta} E(\theta_t), \quad \theta_{t+1} = \theta_t - \mathbf{v}_t$

*Why Newton's method not suitable for machine learning problems?*

# Gradient descent

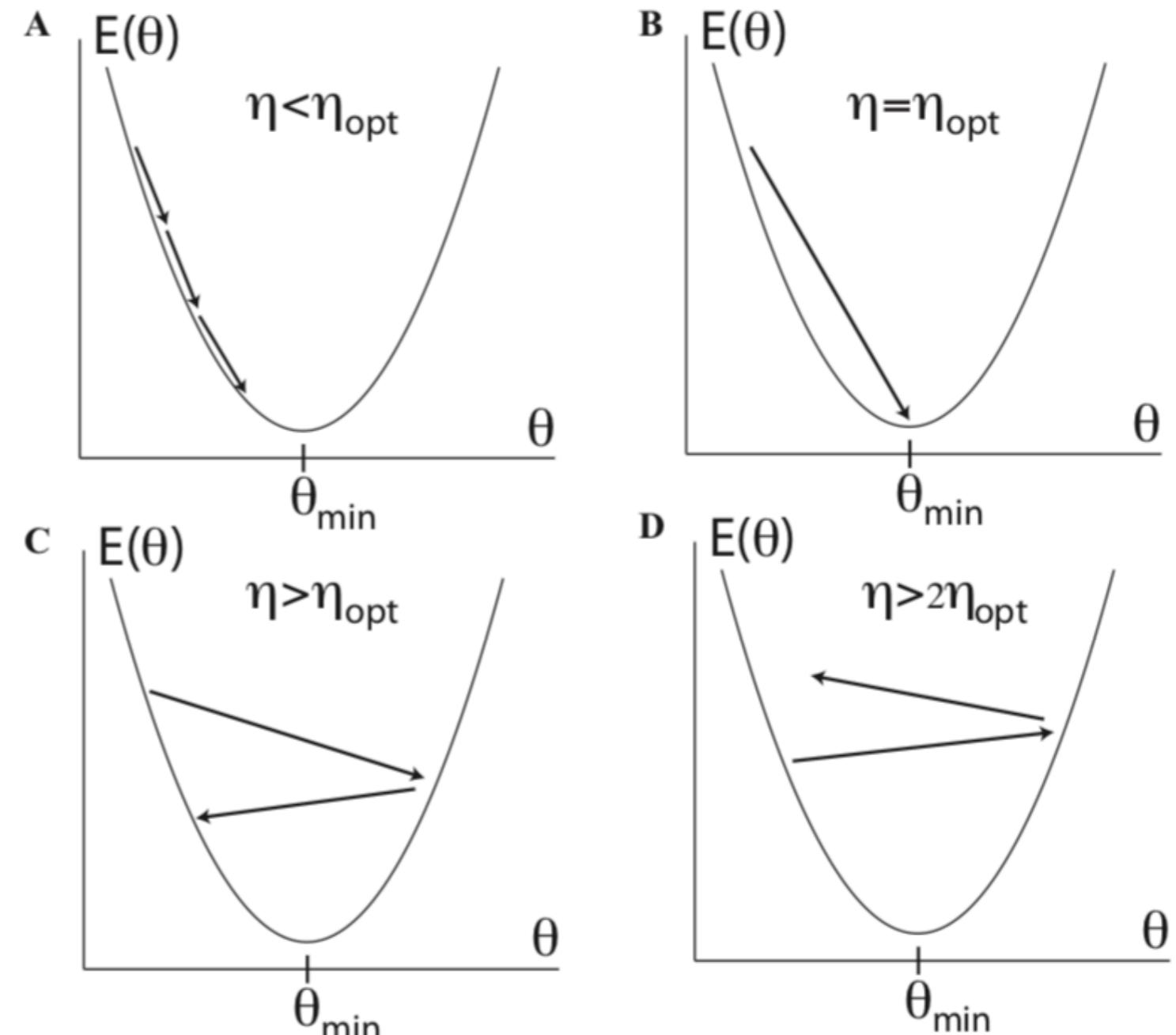
Gradient Descent  $\longrightarrow \mathbf{v}_t = \eta_t \nabla_{\theta} E(\theta_t), \quad \theta_{t+1} = \theta_t - \mathbf{v}_t$

Newton's method  $\longrightarrow \mathbf{v}_t = H^{-1}(\theta_t) \cdot \nabla_{\theta} E(\theta_t), \quad \theta_{t+1} = \theta_t - \mathbf{v}_t$

- Newton's method not practical for ML applications: it involves **evaluation and inversion of Hessian matrices with  $M^2$  entries**, with  $M$  = number of model parameters

- But it provides useful ideas about how to improve GD, for example by **adapting the learning rate to the local curvature** of region of the parameter space

*suggest improvements of GD!*



# Advanced gradient descent

the simplest implementation of GD, unsurprisingly, has several limitations

- It converges to **local, rather than global**, minima of the cost function

*poor performance for realistic applications*

# Advanced gradient descent

the simplest implementation of GD, unsurprisingly, has several limitations

- It converges to **local, rather than global**, minima of the cost function

*poor performance for realistic applications*

- Evaluating gradients is **computationally expensive** for large datasets

$$E(\boldsymbol{\theta}) = \sum_{i=1}^n (y_i - f_\alpha(x_i; \boldsymbol{\theta}_\alpha))^2$$

*involves sum over all n data points*

# Advanced gradient descent

the simplest implementation of GD, unsurprisingly, has several limitations

- It converges to **local, rather than global**, minima of the cost function

*poor performance for realistic applications*

- Evaluating gradients is **computationally expensive** for large datasets

$$E(\theta) = \sum_{i=1}^n (y_i - f_\alpha(x_i; \theta_\alpha))^2$$

*involves sum over all n data points*

- Very sensitive to choice of **learning rates**

*Ideally one would like an adaptive learning rate*

# Advanced gradient descent

the simplest implementation of GD, unsurprisingly, has several limitations

- It converges to **local, rather than global**, minima of the cost function

*poor performance for realistic applications*

- Evaluating gradients is **computationally expensive** for large datasets

$$E(\theta) = \sum_{i=1}^n (y_i - f_\alpha(x_i; \theta_\alpha))^2$$

*involves sum over all n data points*

- Very sensitive to choice of **learning rates**

*Ideally one would like an adaptive learning rate*

- Treats **uniformly all directions** in the parameter space

*and why this is a problem??*

# Advanced gradient descent

the simplest implementation of GD, unsurprisingly, has several limitations

- It converges to **local, rather than global**, minima of the cost function

*poor performance for realistic applications*

- Evaluating gradients is **computationally expensive** for large datasets

$$E(\boldsymbol{\theta}) = \sum_{i=1}^n (y_i - f_\alpha(x_i; \boldsymbol{\theta}_\alpha))^2$$

*involves sum over all n data points*

- Very sensitive to choice of **learning rates**

*Ideally one would like an adaptive learning rate*

- Treats **uniformly all directions** in the parameter space

*we would like to use info also on curvature (as in Newton's method)*

**Generalised Gradient Descent methods** have been developed to address these shortcomings and are at the basis of **modern deep learning methods**

# Stochastic gradient descent

**Stochasticity** can be added to GD by approximating the gradient on a subset of the training data, called a **mini-batch**

$$\nabla_{\theta} E(\theta) = \sum_{i=1}^n \nabla_{\theta} e_i(x_i; \theta) \simeq \nabla_{\theta} E^{\text{MB}}(\theta) \equiv \sum_{i \in B_k} \nabla_{\theta} e_i(x_i; \theta)$$

$k = 1, \dots, n/K \leftarrow \# \text{ points per batches}$

↑  
 $\# \text{ batches}$

We then **cycle over all mini-batches**, updating the model parameters at each step  $k$

$$\mathbf{v}_t = \eta_t \nabla_{\theta}^{\text{MB}} E(\theta_t), \quad \theta_{t+1} = \theta_t - \mathbf{v}_t$$

A full iteration over all  $n$  data points (over the  $n/K$  batches) is called an **epoch**

benefits of SGD: stochasticity prevents getting stuck in local minima, the calculation of the gradient is speed up & stochasticity acts as natural regulariser

# Adding momentum to SGD

SGD can be used with a **momentum term** that provides some **memory** on the direction in which one is moving in the parameter space

$$\mathbf{v}_t = \gamma \mathbf{v}_{t-1} + \eta_t \nabla_{\theta}^{\text{MB}} E(\theta_t), \quad \theta_{t+1} = \theta_t - \mathbf{v}_t$$



*momentum parameter (  $0 < \gamma < 1$  )*

$$\Delta\theta_{t+1} = \gamma \Delta\theta_t - \eta_t \nabla_{\theta} E(\theta_t), \quad \Delta\theta_t = \theta_t - \theta_{t-1}$$

*proposed parameter  
change in iteration t+1*

*proposed parameter  
change in iteration t*

momentum in SGD helps to **gain speed in directions with persistent but small gradients**, while suppressing oscillations in high-curvature directions.

# Adaptive Gradient Descent

**Adaptive GD** varies the learning rate to reflect the **local curvature** of the parameter space without the need to evaluate the Hessian matrix

**RMSprop**: keep track also of the **second moment of gradient**, similar as how the momentum term is a running average of previous gradients

$$\theta_{t+1} = \theta_t - \eta_t \frac{\mathbf{g}_t}{\sqrt{s_t + \epsilon}}$$

*regulator to avoid numerical divergences*

$$\mathbf{g}_t = \nabla_{\theta} E(\theta)$$

*gradient at iteration t*

$$s_t = \mathbb{E}[\mathbf{g}_t^2]$$

*2nd moment of gradient  
(averaged over iterations)*

$$s_t = \beta s_{t-1} + (1 - \beta) \mathbf{g}_t^2$$

*controls averaging time of 2nd  
moment of gradient*

$$\beta = 1 \rightarrow s_t = s_{t-1}$$

$$\beta = 0 \rightarrow s_t = \mathbf{g}_t^2$$

*effective learning rate reduced in directions where gradient is consistently large*

# Adaptive Gradient Descent

**Adaptive GD** varies the learning rate to reflect the **local curvature** of the parameter space without the need to evaluate the Hessian matrix

**ADAM:** adaptively change the parameters from information on **running averages** of first and second moments of the gradient

$$\mathbf{g}_t = \nabla_{\theta} E(\theta)$$

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} - (1 - \beta_1) \mathbf{g}_t \quad \mathbf{s}_t = \beta_2 \mathbf{s}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2$$

$$\widehat{\mathbf{m}}_t = \frac{\mathbf{m}_t}{1 - (\beta_1)^t}$$

*sets lifetime  
of first moment*

$$\widehat{\mathbf{s}}_t = \frac{\mathbf{s}_t}{1 - (\beta_2)^t}$$

*sets lifetime  
of second moment*

$$\theta_{t+1} = \theta_t - \eta_t \frac{\widehat{\mathbf{m}}_t}{\sqrt{\widehat{\mathbf{s}}_t + \epsilon}}$$

# Adaptive Gradient Descent

**Adaptive GD** varies the learning rate to reflect the **local curvature** of the parameter space without the need to evaluate the Hessian matrix

**ADAM** has two main advantages: (i) adapting step size to cut off large gradient directions (prevent oscillations) and (ii) measuring gradients in a natural length scale

to see this, express the ADAM update rules in terms of the **variance in parameter space**

$$\sigma_t = \hat{s}_t - (\hat{m}_t)^2$$

and we can see that the update rule for one of the parameters reads now

$$\Delta\theta_{t+1} = -\eta_t \frac{\hat{m}_t}{\sqrt{\sigma_t^2 + \hat{m}_t^2} + \epsilon}$$

*small fluctuations of gradient*      *large fluctuations of gradient*

$\Delta\theta_{t+1} \rightarrow -\eta_t$        $\Delta\theta_{t+1} = -\eta_t \hat{m}_t / \sigma_t^2$

*learning rate promotional to signal-to-noise ratio*

# Tutorial 2 (2a & 2b)

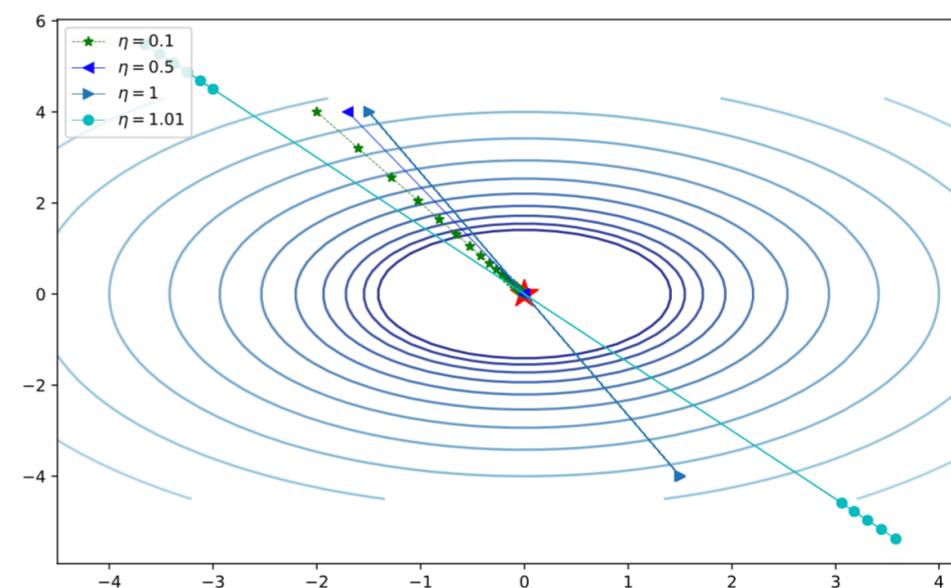
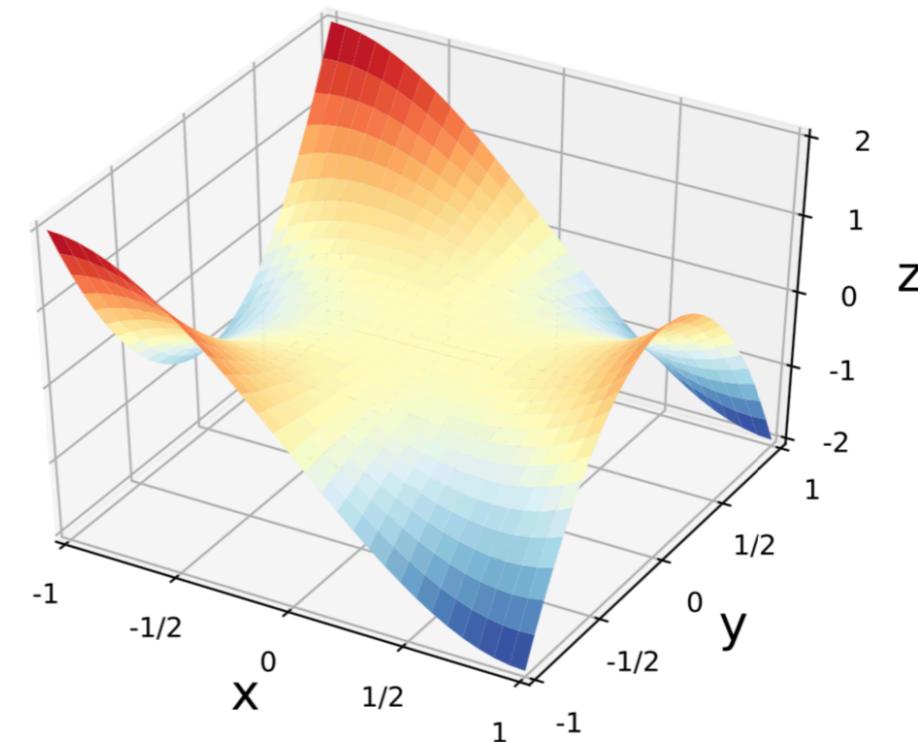
- Study different optimisation algorithms based on GD for simple **1D** (easy) and **2D** functions (more difficult)

- Visualise the **path of GD algorithms in parameter space**

- Determine which choices of the parameters allow reaching the minima quicker

- Given new 2D functions, show that you can **tune GD to find all the minima**, and compare with the analytical results

*challenge: find the most efficient optimiser settings for a given function!*



# Tutorial 2 (2a & 2b)

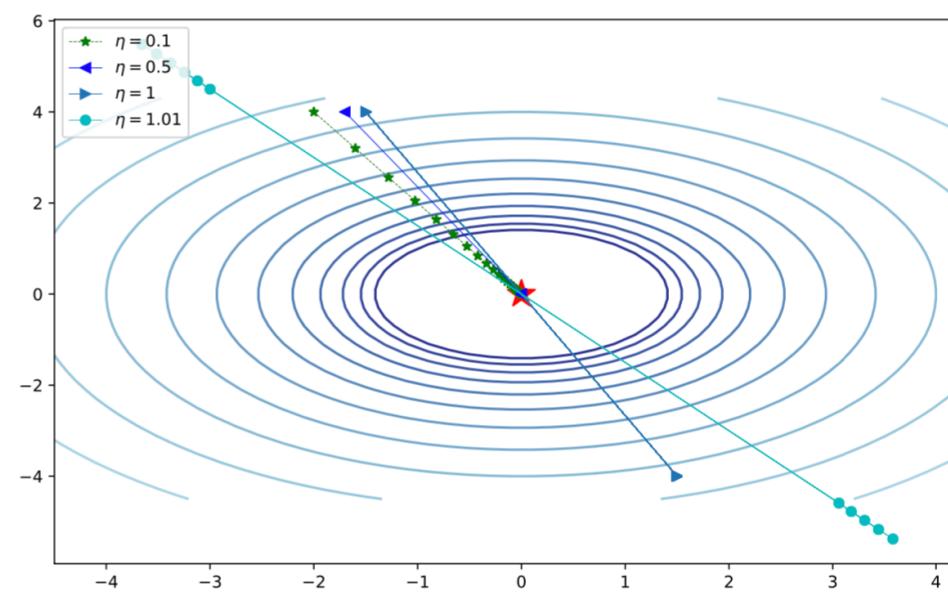
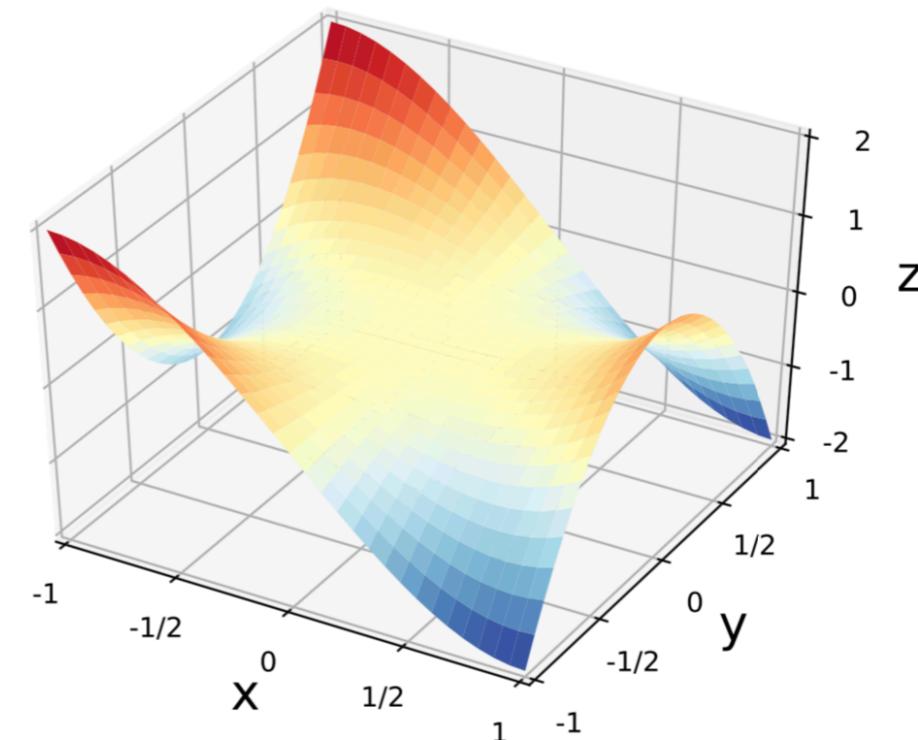
- Study different optimisation algorithms based on GD for simple **1D** (easy) and **2D** functions (more difficult)

- Visualise the **path of GD algorithms in parameter space**

- Determine which choices of the parameters allow reaching the minima quicker

- Given new 2D functions, show that you can **tune GD to find all the minima**, and compare with the analytical results

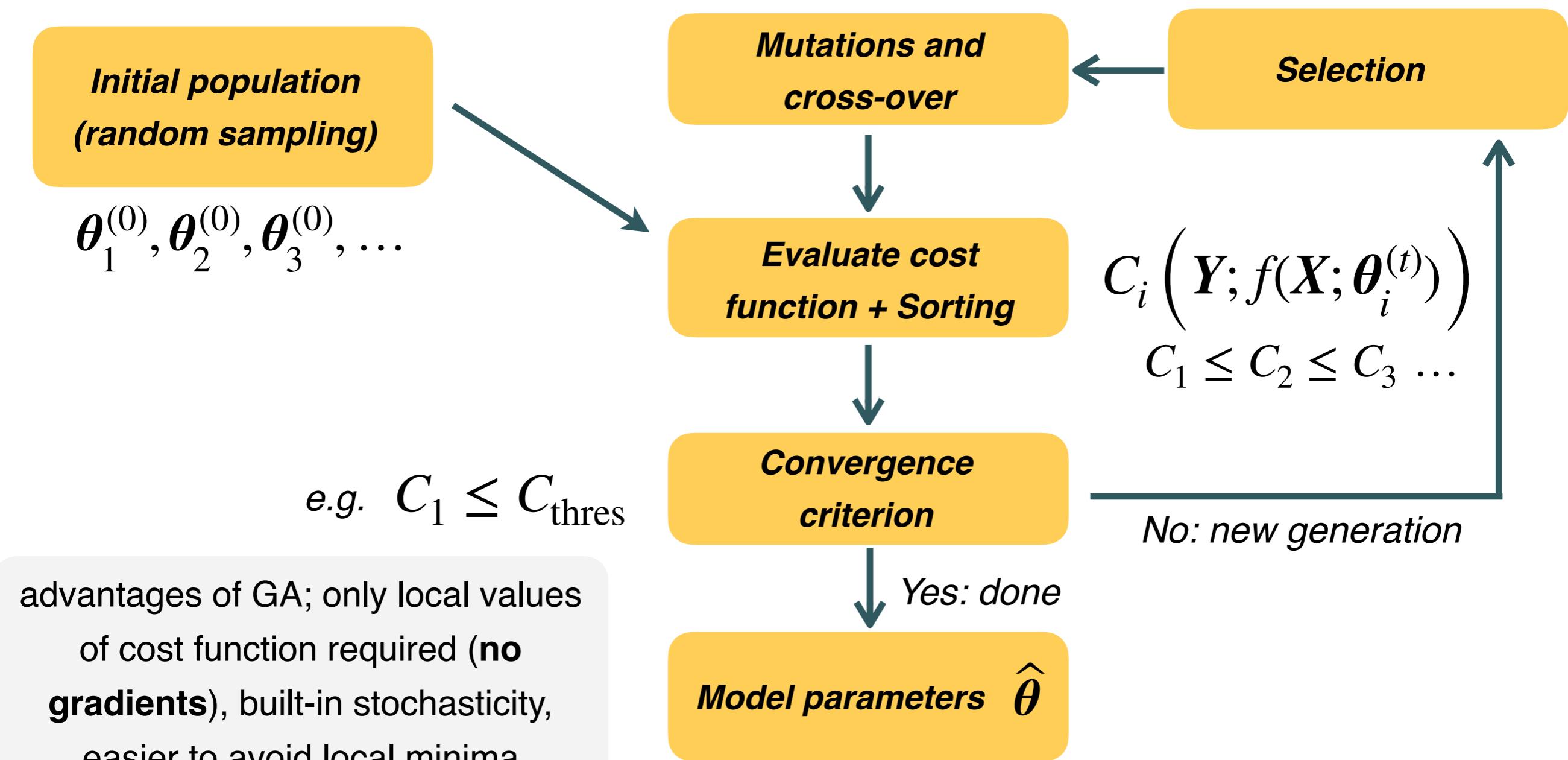
*Question: how could one go beyond hyperparameter setting by trial and error?*



# **Supervised Learning: Stochastic Optimisation Strategies**

# Genetic Algorithms

GAs combine **stochastic and deterministic ingredients** to explore the model parameter space and minimise the cost function



# Genetic Algorithms

GAs combine **stochastic and deterministic ingredients** to explore the model parameter space and minimise the cost function

**mutations**     $\theta_{i,j}^{(t+1)} = \theta_{i,j}^{(t)} \left( 1 + \eta_j^{(t)} \right), \quad j = 1 \dots, m, i = 1, \dots, M$

*learning rates*  
*(depends in generation  
and parameter)*

# model  
parameters

# mutants

# *crossover*

$$\theta_i^{(t)} = \left( \theta_{i,1}^{(t)}, \theta_{i,2}^{(t)}, \dots, \theta_{i,m}^{(t)} \right) \rightarrow \left( \theta_{i,1}^{(t)}, \theta_{i,2}^{(t)}, \dots, \theta_{i,p-1}^{(t)}, \theta_{k,p}^{(t)}, \dots, \theta_{k,m}^{(t)} \right)$$

$$\theta_k^{(t)} = \left( \theta_{k,1}^{(t)}, \theta_{k,2}^{(t)}, \dots, \theta_{k,m}^{(t)} \right) \rightarrow \left( \theta_{k,1}^{(t)}, \theta_{k,2}^{(t)}, \dots, \theta_{k,p-1}^{(t)}, \theta_{i,p}^{(t)}, \dots, \theta_{i,m}^{(t)} \right)$$

***children inherit part of the DNA of each parent***

# Genetic Algorithms

GAs combine **stochastic and deterministic ingredients** to explore the model parameter space and minimise the cost function

$$\text{mutations} \quad \theta_{i,j}^{(t+1)} = \theta_{i,j}^{(t)} \left( 1 + \eta_j^{(t)} \right), \quad j = 1 \dots, m, i = 1, \dots, M$$

*learning rates*  
*(depends in generation  
and parameter)*

*# model  
parameters*

*# mutants*

**crossover**

Chromosome1	11011 00100110110
Chromosome2	11011 11000011110
Offspring1	11011 11000011110
Offspring2	11011 00100110110

*children inherit part of the DNA of each parent*

# Genetic Algorithms

we can illustrate how **GAs work in practice** in a graphical way with this example

# The CMA-ES algorithm

Improve the efficiency of simple GAs by incorporating global (in addition to local) information on the parameter space

One example is the **Covariance matrix Adaptation - Evolutionary Strategy** (CMA-ES)

- Start by randomly initialising the model parameters using a Gaussian distribution

$$\theta^{(0)} \sim \mathcal{N}(0, \mathbf{C}^{(0)})$$

- At every generation, we generate  **$M$  mutants** according to the following distribution

$$a_k^{(t)} \sim \theta^{(t-1)} + \sigma^{(t-1)} \mathcal{N} \left( 0, \mathbf{C}^{(t-1)} \right), \quad k = 1, \dots, M$$


 A diagram illustrating the components of the equation. A red arrow points upwards from the label "step size" to the term  $\sigma^{(t-1)}$ . Another red arrow points from the label "covariance matrix of search distribution" to the term  $\mathbf{C}^{(t-1)}$ . A third red arrow points from the label "tuned during fit" to the term  $\theta^{(t-1)}$ .

***step size***      ***covariance matrix of search distribution***      ***tuned during fit***

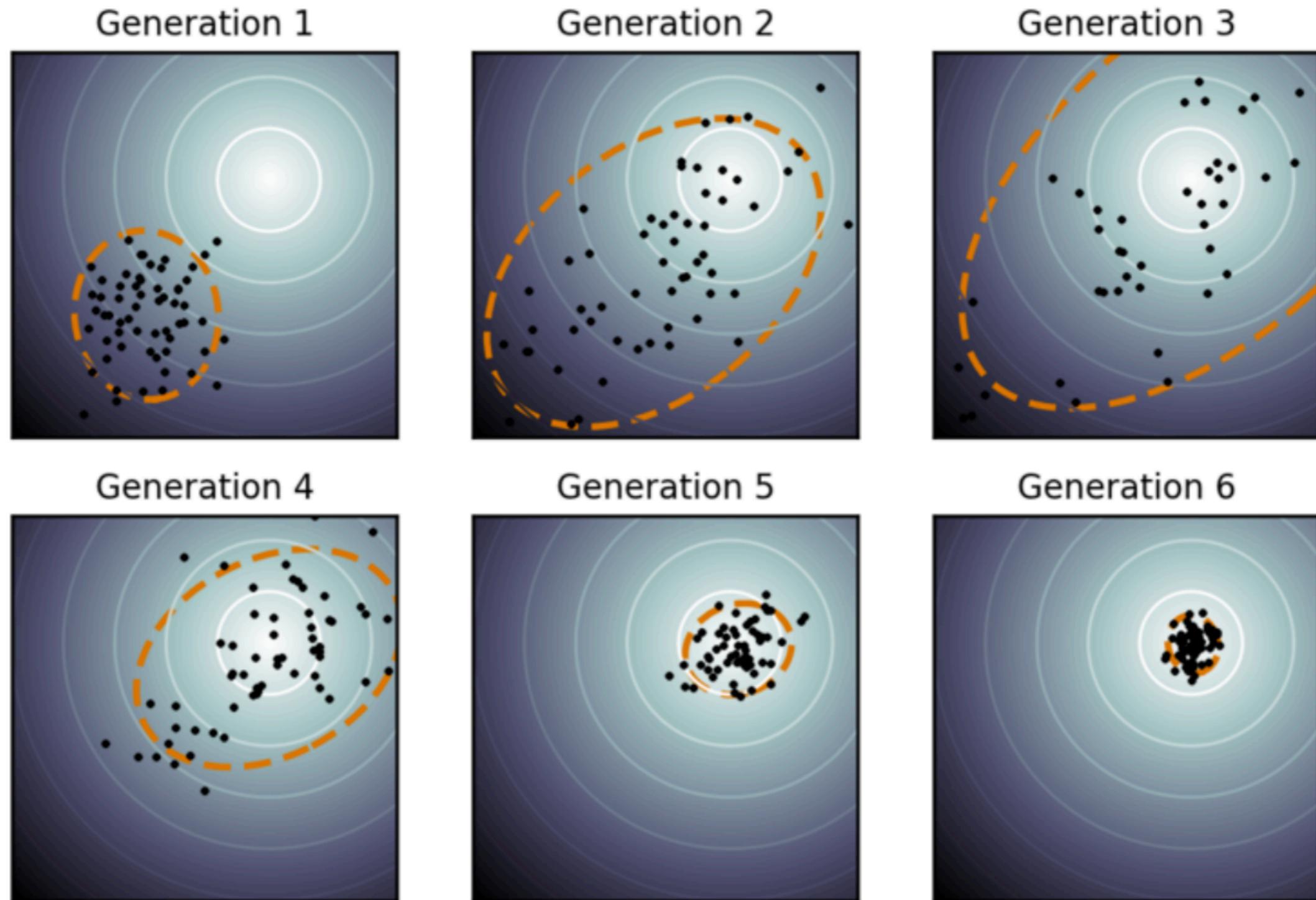
- ⌚ The **new search centre** is computed as weighted average over fraction of best mutants

$$\theta^{(t)} = \theta^{(t-1)} + \sum_{k=1}^{\mu} W_k (a_k - \theta^{(t-1)})$$

**adaptation** as the parameter space is explored is key feature of CMA-ES

# *Parameters of the algorithms*

# The CMA-ES algorithm



**main advantages:** (i) information on gradients never required, only local values of cost function, (ii) information on whole set of mutants (not only best ones) used to tune step size and search covariance matrix

# The CMA-ES algorithm applied



# The CMA-ES algorithm applied



# Optimisation algorithms: summary

- Machine learning problems have associated **high-dimensionality parameter spaces**
- Choice of optimisation algorithm and its hyperparameters crucial for efficient solution of realistic machine learning problems
- Stochastic gradient descent methods are at the core of deep-learning algorithms
- Evolutionary algorithms benefit from built-in stochasticity and can be competitive for problems with very complex space of minima (where multiple quasi-degenerate minima are present)