



MASTER OF SCIENCE  
IN ENGINEERING

# Week 2: Learning and Optimisation

TSM\_DeLearn

**Illustrated with Binary Classification of MNIST Digits**

Jean Hennebert  
Martin Melchior

# Overview

Recap from Last Week, Plan for This Week

MNIST Dataset

Data Preparation

Model and Cost

Gradient Descent

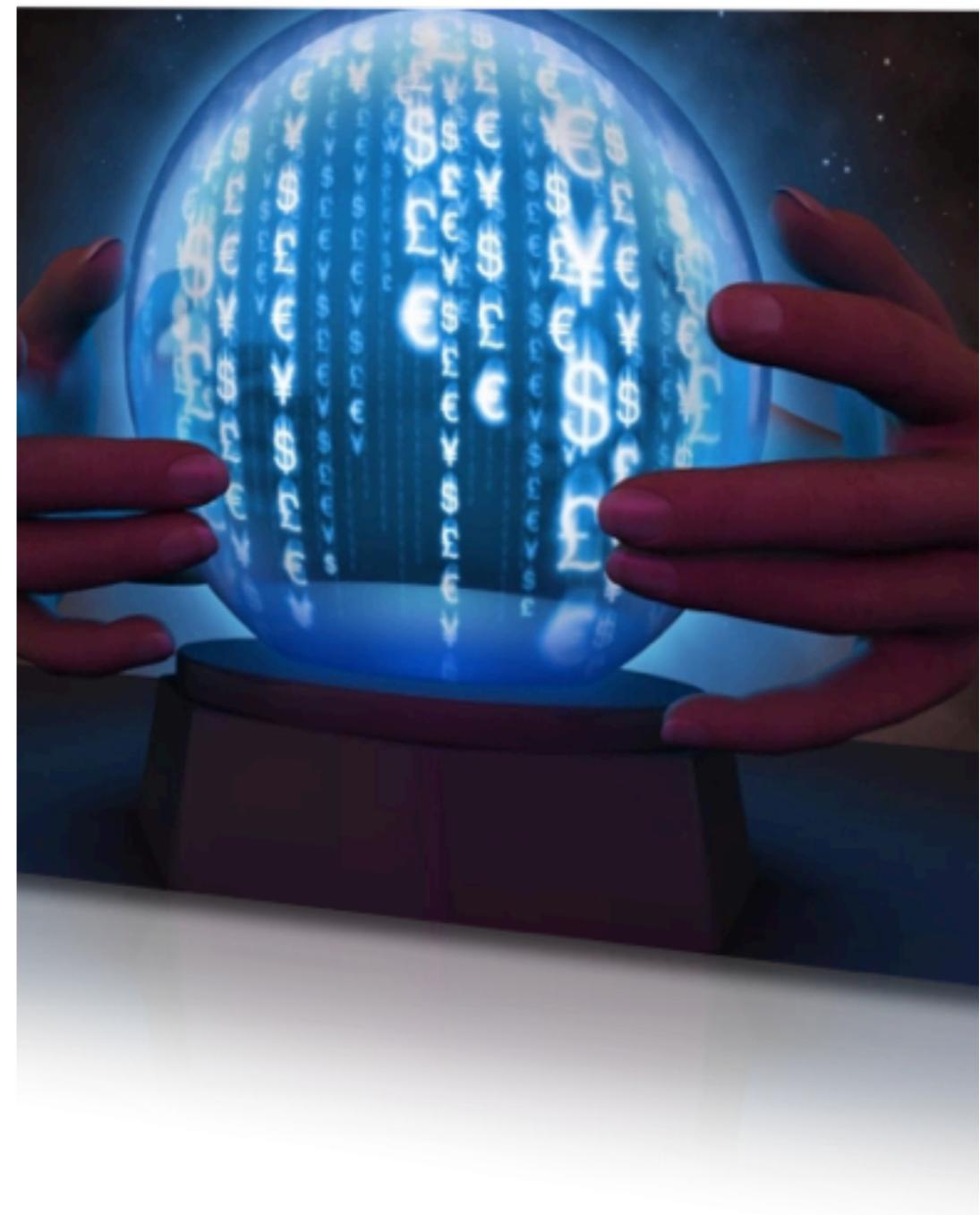
Stochastic Gradient Descent

# Recap from Last Week

From Machine Learning to  
Deep Learning

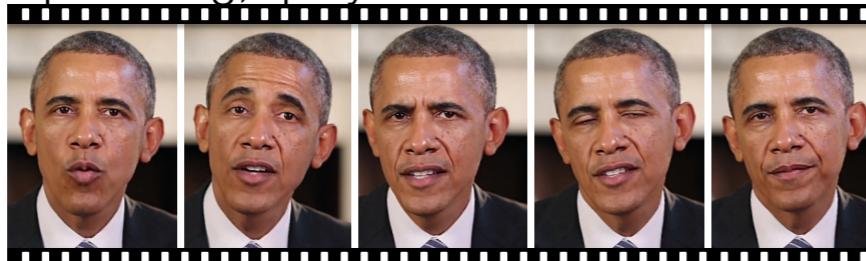
Successful Applications of  
Deep Learning

Perceptron, Learning Rule

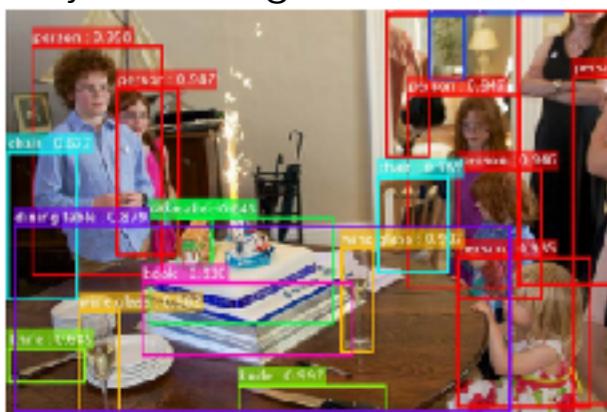


# Example Application Domains

Lip reading, lip sync from audio



Object Recognition



Robot Control

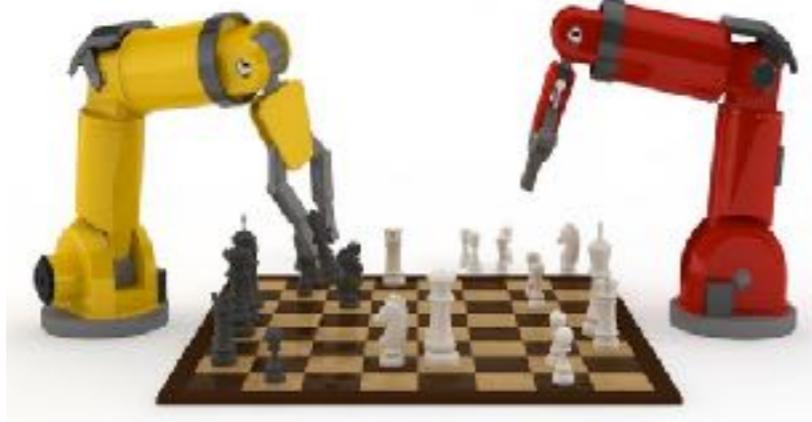
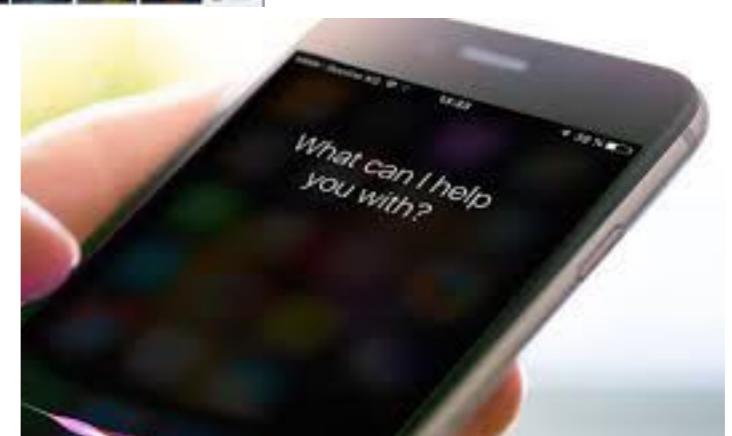


Image Tagging



Question / Answering

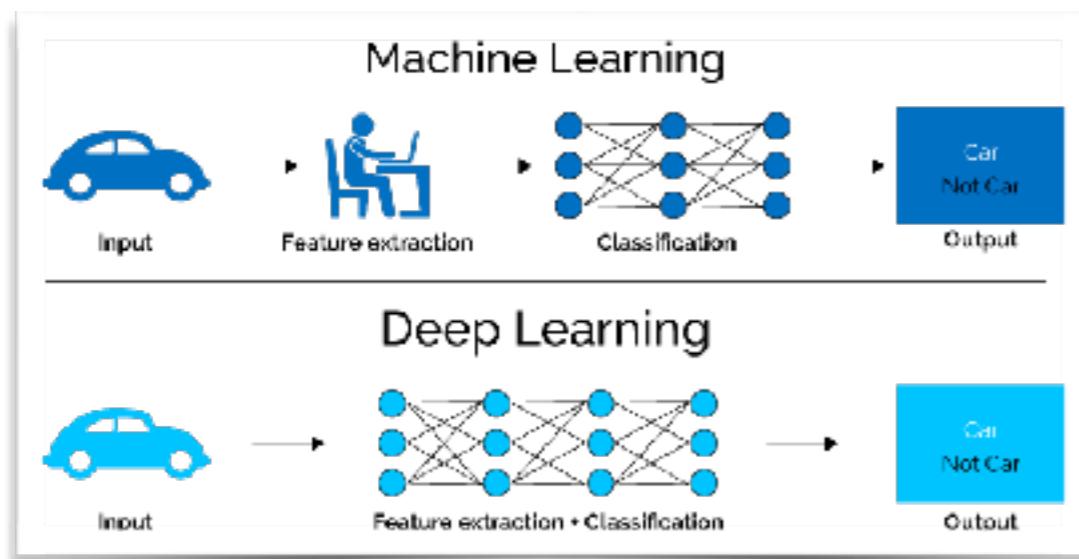


Self-Driving Cars

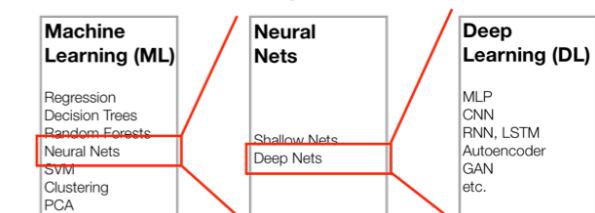


Early Detection of Glaucoma

# From Machine Learning to Deep Learning



Deep Learning can be considered as a sub-discipline of machine learning. But generally, it requires less feature engineering than standard machine learning. The machines finds the features automatically as part of the learning.



## Larger quantities of data

text, audio, images, videos, ...

## New algorithms

DBN, RBM, CNN ...

Deep  
Learning

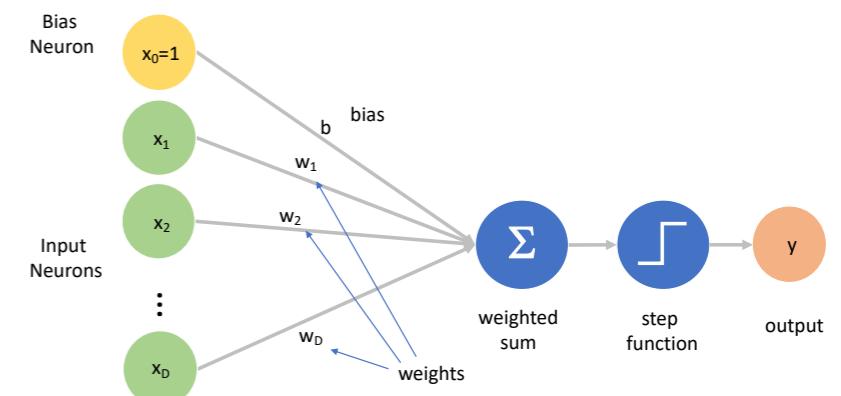
## Better computer performance

GPU, distributed computing ...



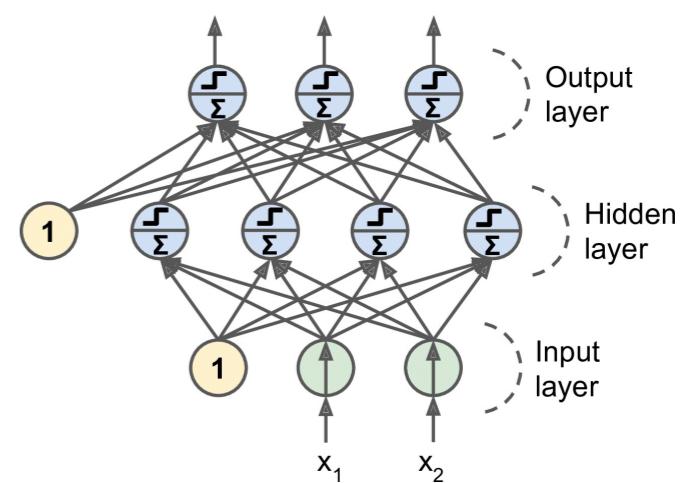
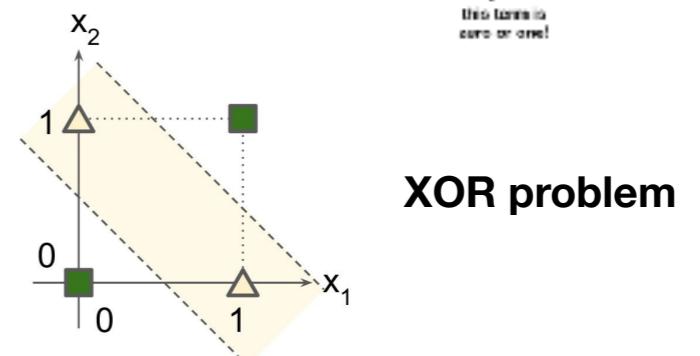
# Perceptron, MLP

- Artificial Neuron (Perceptron, LTU)
  - Rosenblatt's Perceptron with weights and bias parameters
- Perceptron Learning Algorithm
  - Defines a learning scheme for a single perceptron or single layer perceptron.
  - Convergence only for linearly separable problems.
- Not all Problems are Linearly Separable
  - XOR, etc.
- Multi-Layer Perceptron (MLP)
  - Networks with perceptrons as elementary building blocks can represent a much larger family of functions.
- However: How can general MLP's be learned?



Learning algorithm:

1. Initialise parameters zero or small random number
  2. Iterate by updating weights according to
    - A. Pick example  $(\mathbf{x}^{(i)}, y^{(i)})$
    - B. Compute predicted value:  $\hat{y}^{(i)} = H(\mathbf{w} \cdot \mathbf{x}^{(i)} + b)$
    - C. Parameter update rule: only update if  $\hat{y}^{(i)} \neq y^{(i)}$  :
 
$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \alpha \cdot (\hat{y}^{(i)} - y^{(i)}) \mathbf{x}^{(i)} \\ b &\leftarrow b - \alpha \cdot (\hat{y}^{(i)} - y^{(i)}) \end{aligned}$$
- this term is zero or one!



# Plan for this week

- Formulate Learning as an Optimisation Problem.
- Use optimisation algorithms to solve it:  
Gradient Descent
- Illustrate the procedure with a toy model:  
MNIST digits classification

# MNIST Dataset



# Original and Lightweight MNIST Dataset

Public dataset widely used for testing and illustration

- **Original MNIST**

- 70'000 images
- 28x28-pixel
- Available for download from [mldata.org](http://mldata.org)  
(scikit-learn, eras, tensor flow provide utilities to download it)

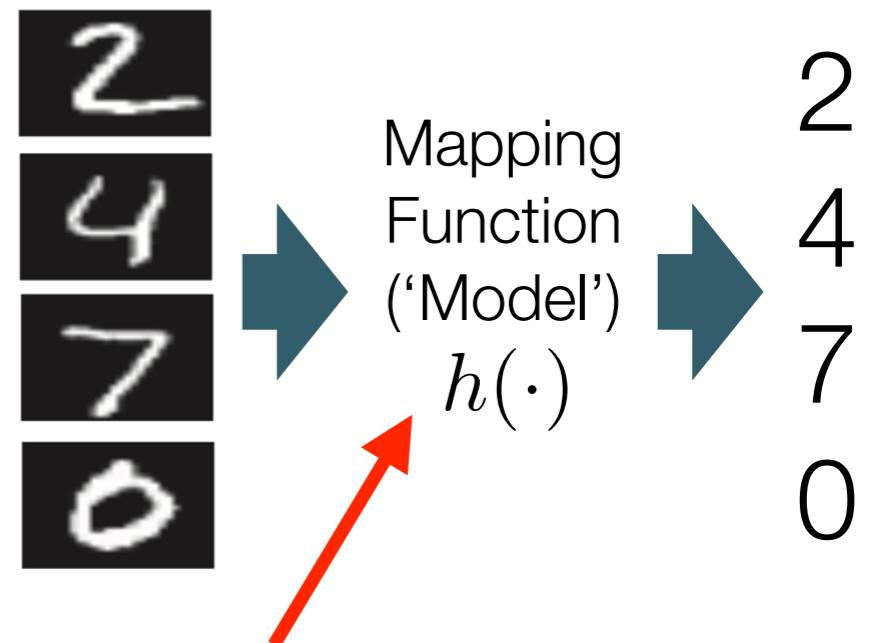
- **Lightweight MNIST**

- 1'797 images
- 8x8-pixel
- Installed as part of scikit-learn library

763423618397846730416523528784659374246192757139773948549394949405542  
20678980340814793494524743646497220913464919238929944066810693103  
151135674425612868256841432621857668652537552916359545729651641961014  
041282874659086317375583029015171769620062610213418816890387032987825  
491421500017390021398700415142238699480231210206034486931933863883136  
193043778299315737764343377158128087180811289932221270002998174259647  
4292518989767452747960449051439280851051132330191409581135097125761753  
915703643189262731004854772034572639488577859168036750224952736719269  
592436043416390930063285368479937752087665019881360012335379117119970  
949315281703692429119308146485716738791293691342519083446716328096781  
084377457447232507551269012752330042405311959292596797557865219180992  
675071053714190432263882135602771592292089638905748519668518320293503  
939052245374682930984293281391882500463248156675246134779603198329814  
093372209995144380831589574035051802078212946832471167800495463432625  
179196386868149249789557948053141152665612794403061731911734209558236  
540973123845338131421761254417350023431492752027759639022631015651647  
943612912424710461467691195319028458798398196192248621133246035794258  
746596706767139088251407028004453004632548439201819627244241702843869  
395225581188019257601040153642666276940555812860392861355000826925190  
425514143217884308679884810145150818574167598461093294406969478012701  
95714805613808672414347998697603511437978330550825290517426654199512  
642229800270575567965816172603263518479218494234775760628052479232023  
672997447257494192614235437115576002646169375765889754789145961395134  
567321602346049800830749873379842874592608457873404451893224327428145  
471003261248281874247461313803418866003948149178012037915332891509376  
092440417318828546973587852115471072163116799122831235743443412665887  
708493771287273060036200951269110399815938529388469487402564793700418  
454018938313294593821092300513151263886912002063414211299604205841939  
081743322933979131337310697641524763244573711305533163053706569002152  
157181068181819451145930213259432749828549406581849853902820606104811  
36311997271944181312593484476561014370197853737516963081900369215677  
1632838214250523936536663003636031643529655814008610681104012601326872  
358144267052462969323797263650014431782216140321614116575140921430925  
478554299317508016570280239076644714818202597191350099875251601541010  
78739693650385107111577133189924514003809338982066577812362362672148  
21939421156552905417511062948264527806612933339900964013593219783280  
702982741521126641127104408744671520574718569627371683222684255894332  
112731362611224655365156249670241245785168913358702930791095486205447  
264878524512110203148274990624736919137421354134673872833106648116548

# MNIST Classification Problem

- General problem to solve:
  - Identify digit shown in the image.
  - Classification problem with 10 classes.
- Simplified problem considered first (for illustration purpose):
  - Identify whether the digit ‘5’ is shown in the image.
  - Binary classification problem.



Different type of models  
lead to different varying  
prediction precision  
<http://yann.lecun.com/exdb/mnist/>

Is there a general  
scheme that can  
be used to train  
the model from  
the data?

# MNIST Image (reduced version of scikit-learn)

## jupyter notebook (using the Lightweight MNIST)

```
from sklearn.datasets import load_digits
digits = load_digits()
x = digits.data
y = digits.target
print("Image Data Shape" , x.shape)
print("Label Data Shape", y.shape)

image = x[0,:].reshape((8,8))
print(image)

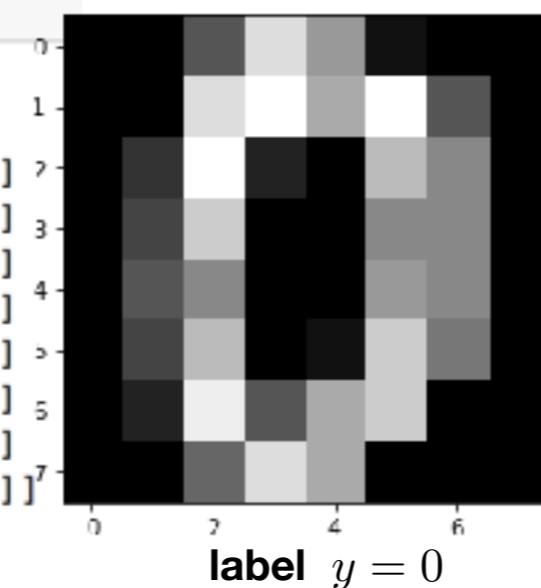
import matplotlib.pyplot as plt
%matplotlib inline
import numpy as np
plt.imshow(image, cmap=plt.cm.gray)
```

Loads all 1797 images and labels

A single image

```
Image Data Shape (1797, 64)
Label Data Shape (1797,)
[[ 0.  0.  5.  13.  9.  1.  0.  0.] 0
 [ 0.  0.  13.  15.  10.  15.  5.  0.] 1
 [ 0.  3.  15.  2.  0.  11.  8.  0.] 2
 [ 0.  4.  12.  0.  0.  8.  8.  0.] 3
 [ 0.  5.  8.  0.  0.  9.  8.  0.] 4
 [ 0.  4.  11.  0.  1.  12.  7.  0.] 5
 [ 0.  2.  14.  5.  10.  12.  0.  0.] 6
 [ 0.  0.  6.  13.  10.  0.  0.  0.]] 7
```

grayscale values  $0 \leq x_k \leq 16$



Build a model for mapping the array with  $8 \times 8 = 64$  input values (greyscale value per pixel) to the labels

$$\hat{y} = h(\mathbf{x})$$

Use a learning algorithm that also works for not linearly separable problems and for multiplayer networks.

# Data Preparation

Training and Test Sets  
Data Normalisation



# Training and Test Set

- Split data into two subsets
  - Training Set : Used for learning the task.
  - Test Set : Used for testing how well the learned model performs.
  - Split ratio: Depends on available data and the specific task to be trained and tested.

Typical ratios:

	<b>Small Datasets</b>	<b>Large Datasets</b>
Train	70-80%	99%
Test	20-30%	1%

- Make sure that training and test sets have the same characteristics
  - Randomly shuffle the dataset before splitting - unless you are sure that the dataset is already shuffled.
- Keep training set and test set strictly separate!
  - Don't use any information contained in the test set to adjust parameters that are used during learning or that define the model.

# Example: MNIST Dataset

Several libraries have built-in functionality to do the split,  
e.g. scikit-learn

**jupyter notebook (using the reduced version from scikit-learn)**

```
from sklearn.datasets import load_digits
from sklearn.model_selection import train_test_split

digits = load_digits()
x = digits.data
y = digits.target
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=0)

print(len(x),len(x_train),len(x_test))

1797 1347 450
```

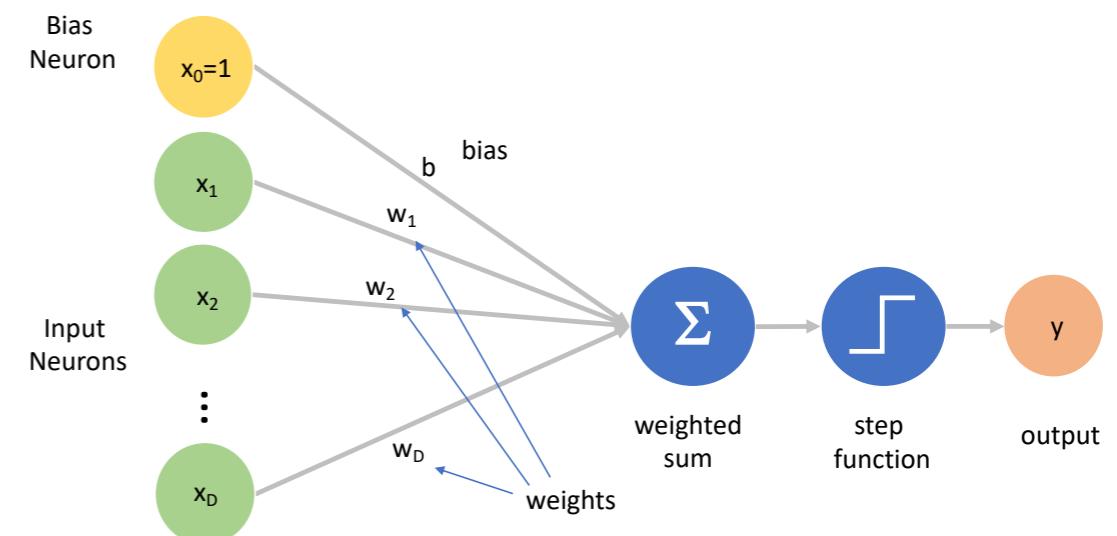
# Data Normalisation, Feature Scaling

- Data contains features on different scales — numeric values may range over different orders of magnitude.
- For numerical stability, it's indispensable to bring the numeric values on similar scales. It improves the convergence properties of the learning algorithm.
- Normalisation (also known as feature scaling) is the process that brings the different inputs to a similar range and importance.
- In general, it applies to both, input data and output data.
- Different schemes available:
  - z-normalisation
  - min/max rescaling or normalisation

# Why is Input Normalisation Important?

Weighted Sum:  $\mathbf{z} = \mathbf{W} \cdot \mathbf{x} + \mathbf{b}$

Step function (such as the Heaviside function) is not dependent on scale, i.e. all the information on the scale should be absorbed when computing  $\mathbf{z}$ .



Two options:

(A) Adjust the weights to compensate for the scale observed in the inputs



scale of the weights  
 $w$  such that  $x \cdot w \sim 1$

(B) Normalise the inputs to make them scale-independent — optimisation algorithm can concentrate on finding the weights without taking care of the scale



scale of the weights  $w$  can be initialised  
independently while still obtaining  $x \cdot w \sim 1$

It turns out that with option (B) the optimisation algorithm works better<sup>(\*)</sup> and the rule for initialising the weights is simpler.

# z-normalisation

$$x_k'^{(i)} = \frac{x_k^{(i)} - \mu_k}{\sigma_k}$$

**Shifting** and **rescaling** the data so that a zero mean and a unit-variance is obtained.

where

$$\begin{aligned}\mu_k &= \frac{i}{N} \sum_{k=1}^N x_k^{(i)} \\ \sigma_k^2 &= \frac{i}{N} \sum_{k=1}^N (x_k^{(i)} - \mu_k)^2\end{aligned}$$

Mean and standard deviation for the  $k$ -th component.

Remarks:

- The parameters  $\mu_k, \sigma_k$  should be **computed on the training set** only. Computing them on the whole data set (including the test set) is considered as cheating.
- Normalisation of **image data**: Image data are normalised typically not per pixel. Mean and standard deviation are computed over all the pixels.

# Min-Max Rescaling and Normalisation

## Min-Max Rescaling

$$x_k'^{(i)} = \frac{x_k^{(i)} - \min_j(x_k^{(j)})}{\max_j(x_k^{(j)}) - \min_j(x_k^{(j)})}$$

Features are rescaled in the [0,1] range.

## Min-Max Normalisation

$$x_k'^{(i)} = 2 \cdot \frac{x_k^{(i)} - \min_j(x_k^{(j)})}{\max_j(x_k^{(j)}) - \min_j(x_k^{(j)})} - 1$$

Features are rescaled and centred to [-1,1] range.

The same important remarks as on previous slide also holds here:

- image data: min/max computed over all pixels
- min/max computed on the basis of training set

# Scaling and Centering

In general, normalisation consists of two steps :

- **Scaling:** Bring the data to the same scale.
  - Improves the convergence speed and accuracy of the learning algorithms.
- **Centering:** Balance the data around 0.
  - Improves the robustness of the learning algorithms. This is particularly important for deep learning networks.
  - Probably also improves the convergence speed and accuracy of the learning algorithms.

Sometimes, scaling (e.g. min-max-rescaling) is sufficient.

**Remark:** Normalisation of input data is closely related to batch normalisation (discussed in the scope of **regularisation** in week 5).

# Normalisation for MNIST

MNIST data range given by integer greyscale values:

- Original MNIST:  
0,...,255
- Lightweight MNIST:  
0,...,16

## Min-Max-Scaling

We only deal with a shallow network. Hence, we don't care about centering.

```
import numpy as np
xmin = np.min(x_train)
xmax = np.max(x_train)
print(xmin, xmax)

x_train = x_train / xmax
x_test = x_test / xmax

0.0 16.0

print(x_train[0,:].reshape((8,8)))
```

[[ 0. 0.1875 0.8125 1. 0.5625 0. 0. 0. 0. ]  
 [ 0. 0.625 0.9375 0.8125 0.9375 0.125 0. 0. 0. ]  
 [ 0. 0.9375 0.25 0.25 1. 0.0625 0. 0. 0. ]  
 [ 0. 0. 0. 0.3125 1. 0.125 0. 0. 0. ]  
 [ 0. 0. 0.0625 0.875 0.8125 0. 0. 0. 0. ]  
 [ 0. 0. 0.625 1. 0.3125 0. 0. 0. 0. ]  
 [ 0. 0.25 1. 0.8125 0.5 0.625 0.5625 0.0625 0.0625]  
 [ 0. 0.125 1. 1. 0.875 0.75 0.5625 0.5625 0.0625]]

Computes the min and max over all the pixels.

Several libraries have built-in functionality to do normalisation: see e.g. <http://scikit-learn.org/stable/modules/preprocessing.html#preprocessing-normalization>.

# Notations Time ...

The notations can become quite intricate. We follow quite closely the notations of Andrew Ng on Coursera. You can find a summary in DL-notations.pdf on Moodle.

$m$	number of samples in the input dataset
$n_x$	number of input features, dimension of the input feature vector
$\mathbf{x}$	input feature vector of dimension $n_x$
$x_k$	$k$ -th component of the input feature vector ( $1 \leq k \leq n_x$ )
$\mathbf{x}^{(i)}$	input feature vector of the $i$ -th training sample ( $1 \leq i \leq m$ )
$x_k^{(i)}$	$k$ -th component of the input feature vector of the $i$ -th training sample
$y$	scalar output variable, also called target output or label
$\mathbf{y}$	output vector (or target output vector) of dimension $n_y$ vector
$y_k$	$k$ -th component of the output vector (or target output vector) ( $1 \leq k \leq n_y$ )
$\hat{y}$	predicted output, as computed by the mapping function
$\hat{\mathbf{y}}$	predicted output vector, as computed by the mapping function
$(\mathbf{x}, \mathbf{y})$	input sample (pair of input feature vector and corresponding label vector)
$(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$	$i$ -th input sample of the input dataset ( $1 \leq i \leq m$ )

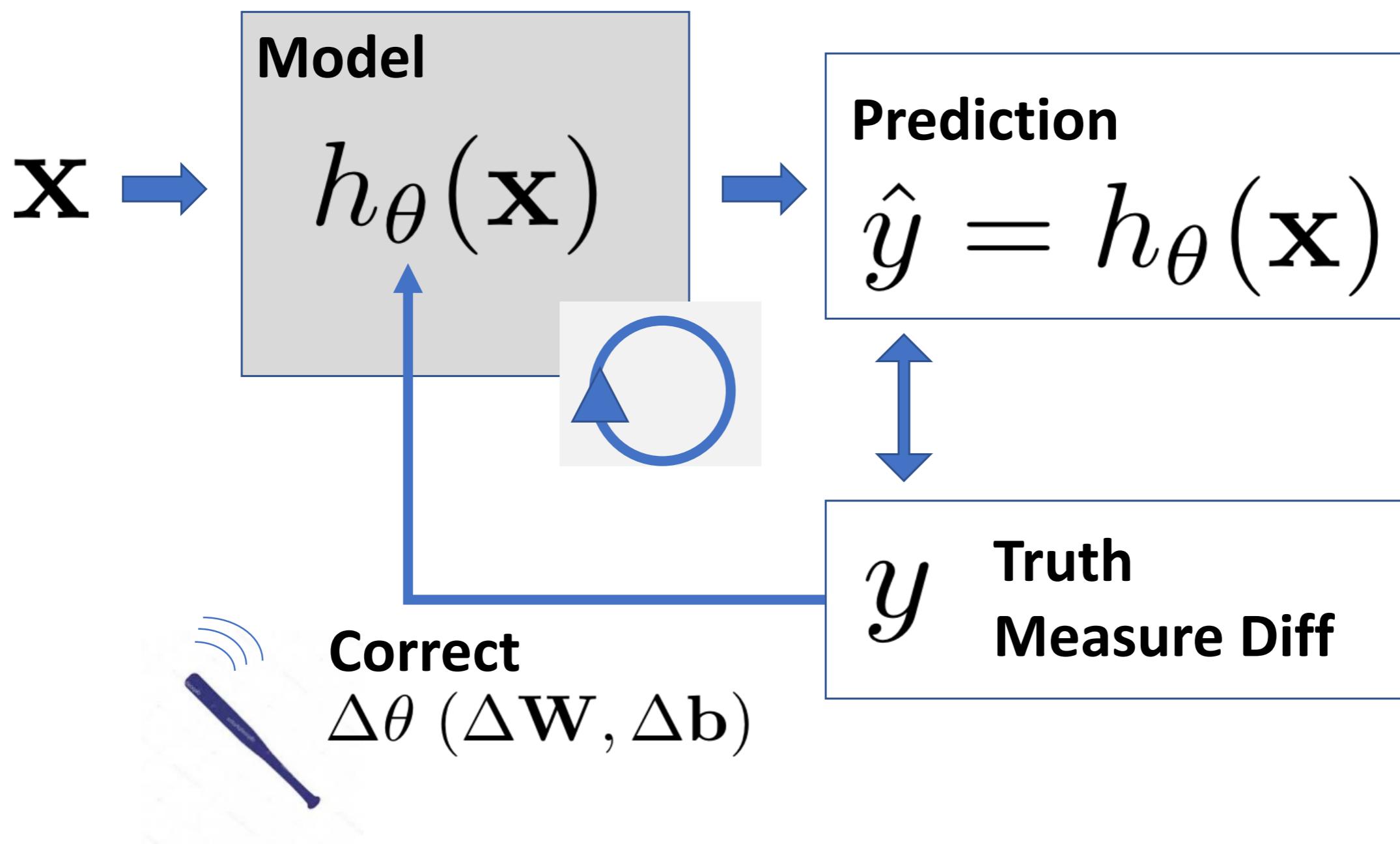
**to be continued – UFFF !\*\$^^...**

# Model and Cost

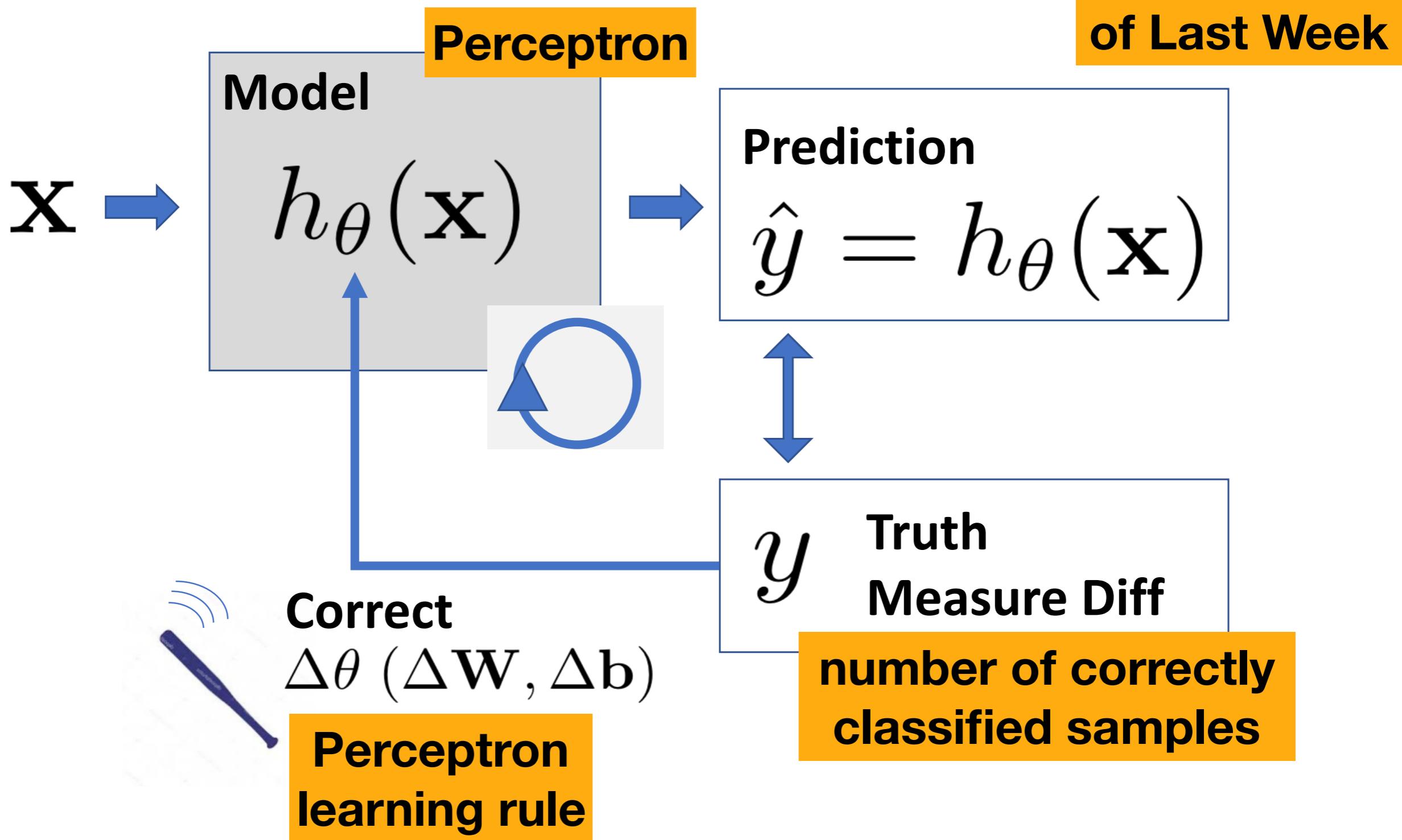
Generalised Perceptron  
Model  
Cost Function



# Training a Model

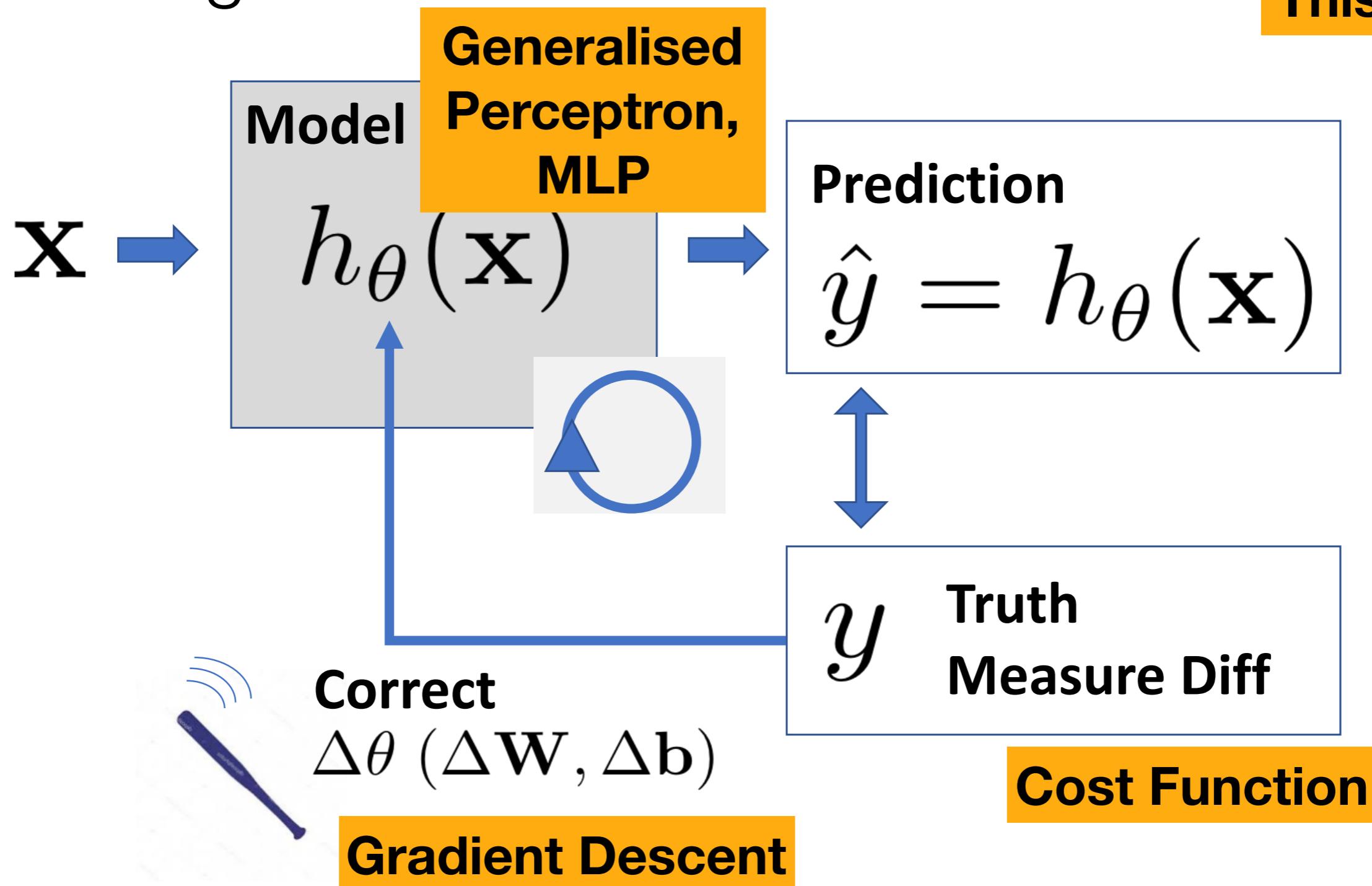


# Training a Model



# Training a Model

This Week

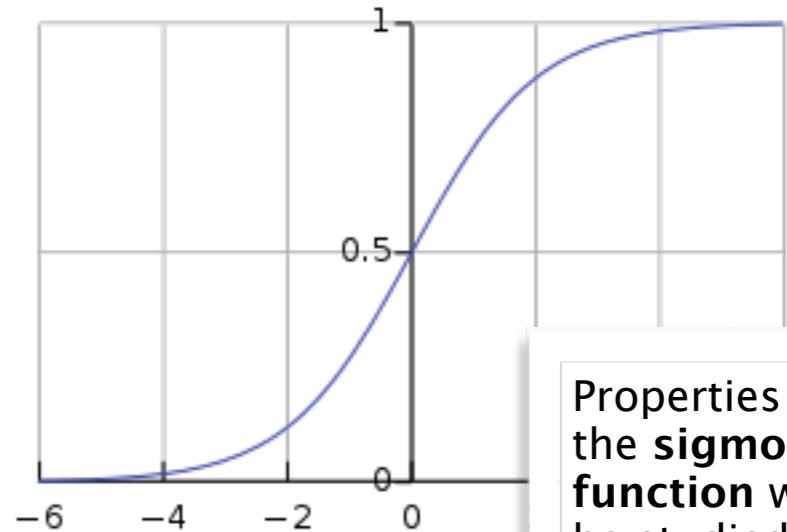


# “Generalised” Perceptron

Generalised from Rosenblatt’s Perceptron.

- 1) **Smooth activation function: Sigmoid Function**  $\sigma(z) = \frac{1}{1+e^{-z}}$   
(replaces the ‘hard-limit’ Heaviside-function)

- Allows for a probabilistic interpretation:  
Answer is given in form of a probability  
of seeing the class or not (and not just a ‘yes/no’).
- Differentiable so that optimisation techniques  
from calculus can be adopted.



Properties of the **sigmoid function** will be studied in the exercises.

- 2) **Gradient Descent Optimisation Algorithm used for Learning**

- Applicable also for training datasets that are not linearly separable and for MLPs

# Model for Simplified (Binary) MNIST Classification

$$\hat{y} = h_{\theta}(\mathbf{x}) = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$$

$\sigma(z) = \frac{1}{1 + e^{-z}}$  sigmoid activation function

$\theta = (\mathbf{w}, b)$  model parameters: weights and bias

## Output

- No longer a ‘yes’ (1) or ‘no’ (0) but a numeric value  $\hat{y} \in ]0, 1[$
- Class label can be assigned by the rule:

$$\text{yes} : \hat{y} = h_{\theta}(\mathbf{x}) \geq 0.5$$

$$\text{no} : \hat{y} = h_{\theta}(\mathbf{x}) < 0.5$$

# How to Choose the Parameters?

- Choose parameters such that the predicted values  $\hat{y}$  are in some notion of distance ‘close’ to the true outcomes  $y$ .
- The notion of distance is expressed in terms of a cost function. Choose the parameters such that the cost gets small / smaller / smallest.
- Different cost functions lead to different solutions.
- We will discuss two cost functions:
  - (1) Mean Quadratic Distance between  $y$  and  $\hat{y}$ .
  - (2) Cross-Entropy (distance between probability distributions)

# Mean Square Error Cost Function

$$J_{MSE}(\theta) = \frac{1}{2m} \sum_{i=1}^m \left( h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

This is for later math simplifications

This is to normalise the cost regarding the size of training set — so that the cost magnitude of the cost is independent of the sample size  $m$ .

The square makes sure that a positive definite distance measure is obtained — the cost is positive and 0 only if all the prediction exactly match the truth.

By minimising  $J_{MSE}(\theta)$  with respect to the parameters  $\theta$  we adjust the model to ‘best’ represent the mappings  $\mathbf{x} \rightarrow y$  seen in the training data.

Actually, there is another cost function better suited for classification problems, the cross entropy cost which introduce later.

# Solving the Minimisation Problem

- There is no closed-form solution to the minimisation problem with the model  $h_\theta(\mathbf{x})$  and the cost function  $J_{MSE}(\theta)$  as given by the perceptron.
- In general, *no closed-form solution* exists for these kind of optimisation problems — a closed-form solution only exists for very specific cases (e.g. linear regression).
- In general, *numeric methods* are needed to find a solution for the parameters  $\theta$ . Typically, these work iteratively where a solution is iteratively approached.
- *Gradient Descent* is a family of optimisation algorithms widely used in practice.

# Imagine Different Minimisation Strategies



Activity

- In groups of 2, address the following points through a short discussion:
  - Think about different strategies for how to find the minimum?
  - Imagine that you move on a landscape and try to reach the lowest point.
  - Will your strategies lead to the lowest point?

# Batch Gradient Descent

Characteristics  
Principles  
Mathematical Formulation



# Geometric Intuition about Gradient Descent



The **gradient**, a mathematical construct, gives the direction steepest ascent!

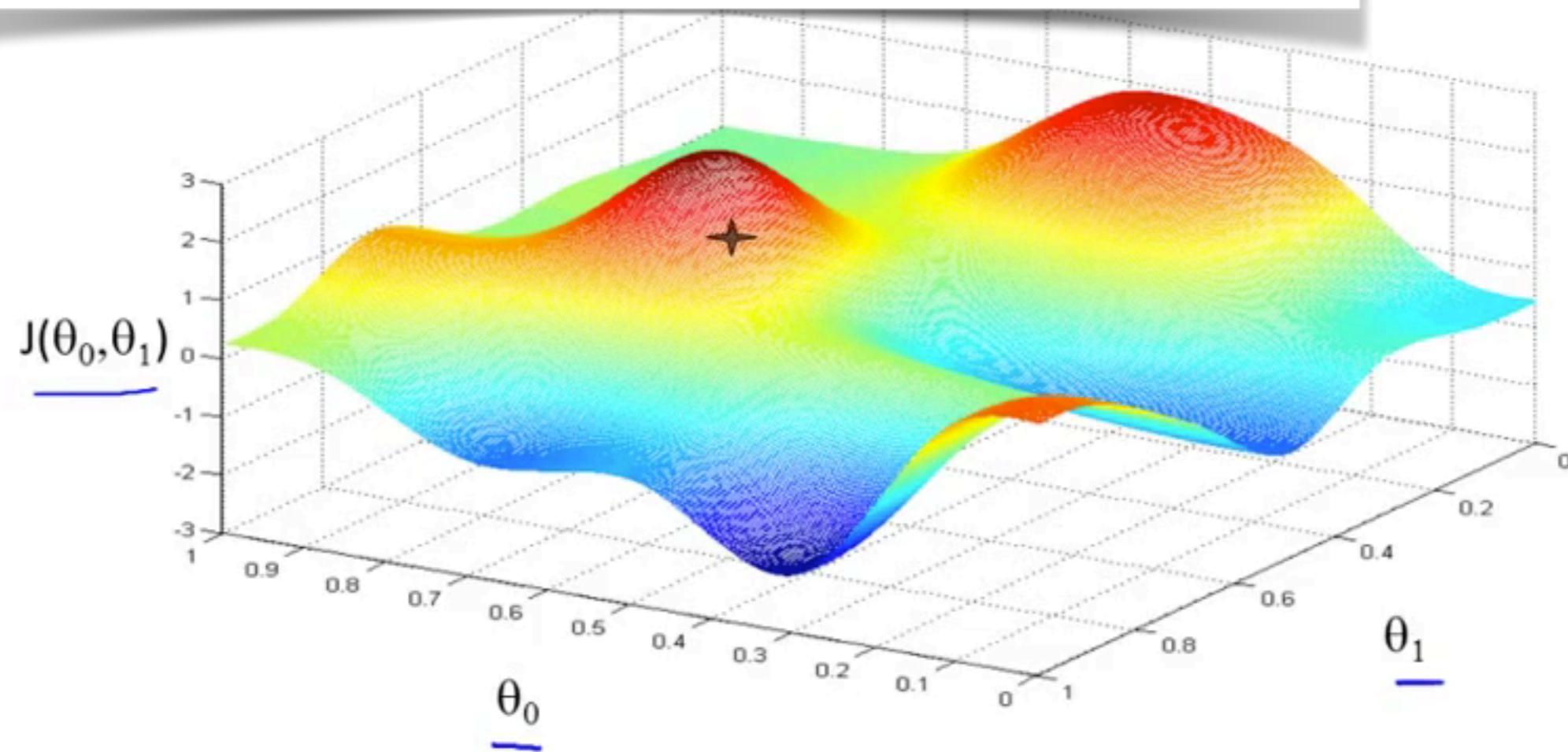
- Cost function defines a manifold in the parameter space.
- Imagine you are a hiker that walks on this manifold and you want to always go downhill — but you cannot see but only feel directions with the tip of your feet.
- Iteratively do the following until you reach **a** minimum:
  - Locally tip with your feet to identify the direction where it is going down steepest.
  - Make a step in this direction proportional to the slope.

# Basic Principle of Gradient Descent

- 1) Start with some initial value for the parameter vector (for example random or 0):  $\theta_0$
- 2) Iteratively update the parameter vector by
  - (i) Computing the gradient of the cost function  $J(\theta)$  at the last position reached ( $\theta_t$ ) :  $\nabla_{\theta} J(\theta_t)$
  - (ii) Stepping in the negative direction of the gradient according to
$$\theta_{t+1} = \theta_t - \alpha \cdot \nabla_{\theta} J(\theta_t) \quad (\alpha: \text{learning rate})$$
- (3) Stop when change in parameter vector small.

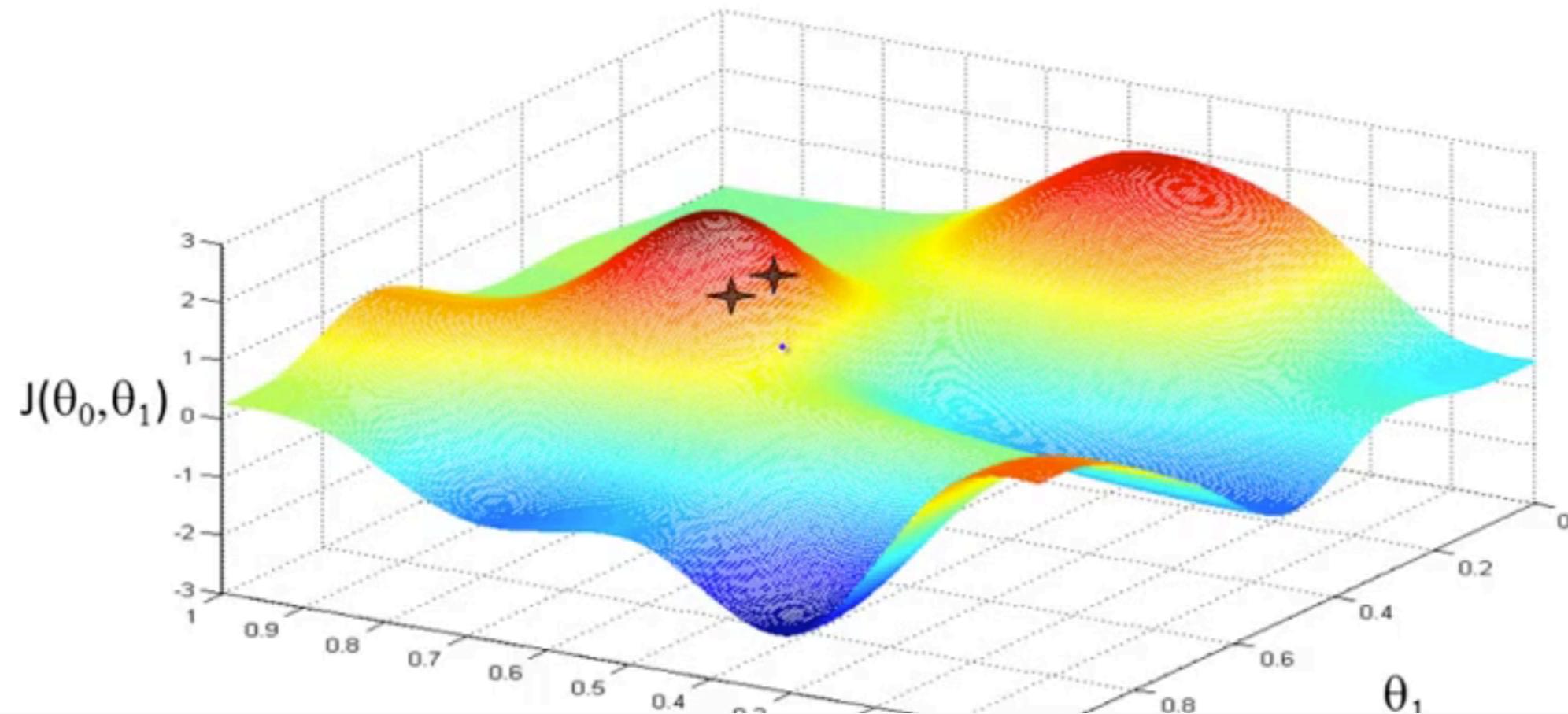
# Gradient Descent Intuition

From the initial position, take steps in the opposite direction of the gradient, “going downhill” until a minimum is reached.



Source: Andrew Ng - Machine Learning class - Stanford

# Gradient Descent Intuition



Source: Andrew Ng - Machine Learning class - Stanford

The path to the (local) minimum depends on the initial conditions. A slight difference in the initial parameter values, may lead to another minimum, a local minimum which is above the global minimum reached in the previous run.

In the case of a single unit perceptron and a sigmoid activation function, the  $J()$  function is convex, i.e. “bowl shaped” with a unique minimum, so we don’t need to care by getting stuck in a local minimum in this case.

Ng

# Characteristics of Gradient Descent

- Learning principle does not depend on the type of model  $h_\theta(\mathbf{x})$  as long as gradients can be computed (e.g. also works for deep neural networks).
- Learning principle works ‘locally’ and finds local wells – not designed to find global minima.
- In a well (local minimum), the gradient is zero. There, it cannot distinguish between local minima, local maxima or saddle points (referred to as critical points).
- Learning schema is iterative – iteratively approaches a critical point. There it fluctuates around the critical point.
- Generally, the learning schema is not guaranteed to converge (e.g. if learning rate is chosen too large, see later).

# How to Compute Gradients?

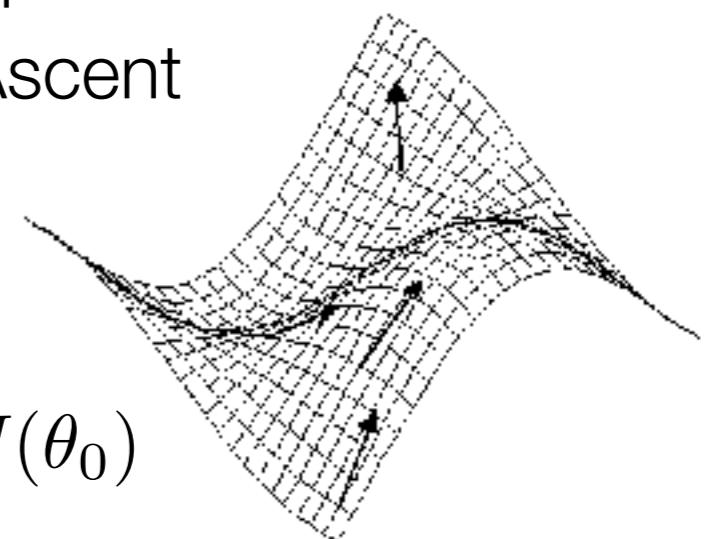
**Lots of Calculus  
Needed !**

**Partial Derivative** of Function  $J(\theta) = J(\theta_1, \dots, \theta_n)$  w.r.t.  $\theta_k$ :

$$\frac{\partial J}{\partial \theta_k} = \lim_{\epsilon \rightarrow 0} \frac{1}{\Delta \theta_k} (J(\theta_1, \dots, \theta_k + \Delta \theta_k, \dots, \theta_n) - J(\theta_1, \dots, \theta_k, \dots, \theta_n))$$

**Gradient:**  $\nabla_{\theta} J = \frac{\partial J}{\partial \theta} = \begin{pmatrix} \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{pmatrix}$

Direction of  
Steepest Ascent



Linearisation (Taylor):  $J(\theta_0 + \Delta\theta) \approx J(\theta_0) + \Delta\theta \cdot \nabla_{\theta} J(\theta_0)$

Largest Decrease in direction of  $\Delta\theta = -\alpha \nabla_{\theta} J(\theta_0)$  :

$$J(\theta_0 - \alpha \nabla_{\theta} J(\theta_0)) \approx J(\theta_0) - \alpha \|\nabla_{\theta} J(\theta_0)\|^2 \quad (\alpha \text{ sufficiently small})$$

# Gradient of the Mean Square Cost (Generalised Perceptron)

**Lots of Calculus  
Needed !**

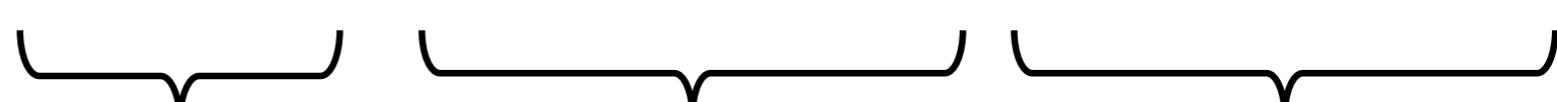
$$\begin{aligned}
 \frac{\partial J_{\text{MSE}}(\theta)}{\partial \theta_k} &= \frac{\partial}{\partial \theta_k} \frac{1}{2m} \sum_{i=1}^m (h_\theta(\mathbf{x}^{(i)}) - y^{(i)})^2 \\
 &= \frac{1}{2m} \sum_{i=1}^m \frac{\partial}{\partial \theta_k} (h_\theta(\mathbf{x}^{(i)}) - y^{(i)})^2 \\
 &= \frac{1}{m} \sum_{i=1}^m (h_\theta(\mathbf{x}^{(i)}) - y^{(i)}) \frac{\partial h_\theta(\mathbf{x}^{(i)})}{\partial \theta_k} \\
 &= \frac{1}{m} \sum_{i=1}^m (h_\theta(\mathbf{x}^{(i)}) - y^{(i)}) \sigma'(\mathbf{w} \cdot \mathbf{x}^{(i)} + b) \frac{\partial (\mathbf{w} \cdot \mathbf{x}^{(i)} + b)}{\partial \theta_k}
 \end{aligned}$$

In the third and the fourth line we have used the **chain rule of calculus**. By using  $\sigma'(z) = \sigma(z)(1 - \sigma(z))$  and  $\hat{y}^{(i)} = h_\theta(\mathbf{x}^{(i)}) = \sigma(\mathbf{w} \cdot \mathbf{x}^{(i)} + b)$  we obtain

$$\begin{aligned}
 \nabla_{\mathbf{w}} J_{\text{MSE}}(\mathbf{w}, b) &= \frac{1}{m} \sum_{i=1}^m \hat{y}^{(i)} (1 - \hat{y}^{(i)}) (\hat{y}^{(i)} - y^{(i)}) \mathbf{x}^{(i)} \\
 \nabla_b J_{\text{MSE}}(\mathbf{w}, b) &= \frac{1}{m} \sum_{i=1}^m \hat{y}^{(i)} (1 - \hat{y}^{(i)}) (\hat{y}^{(i)} - y^{(i)})
 \end{aligned}$$

# Gradient of the Mean Square Cost (Generalised Perceptron)

The gradient of the cost function is the sum over all training samples where each sample contributes in proportion to the difference between true label and its predicted value.

$$\begin{aligned}\nabla_{\mathbf{w}} J_{\text{MSE}}(\mathbf{w}, b) &= \frac{1}{m} \sum_{i=1}^m \hat{y}^{(i)} (1 - \hat{y}^{(i)}) (\hat{y}^{(i)} - y^{(i)}) \mathbf{x}^{(i)} \\ \nabla_b J_{\text{MSE}}(\mathbf{w}, b) &= \frac{1}{m} \sum_{i=1}^m \hat{y}^{(i)} (1 - \hat{y}^{(i)}) (\hat{y}^{(i)} - y^{(i)})\end{aligned}$$


Average  
over the  
training set.

Measure for the  
uncertainty in the  
prediction: It is small  
if  $\hat{y}$  is close to zero or  
close to 1 and large if  
 $p$  is around 0.5, i.e. if  
the prediction is clear.

Error Signal –  
contributes more for  
samples with  
mismatch. Similar to  
Perceptron Learning  
Rule but  $\hat{y}$  is  
different.

# General Gradient Descent Update-Rules

In vector notation:

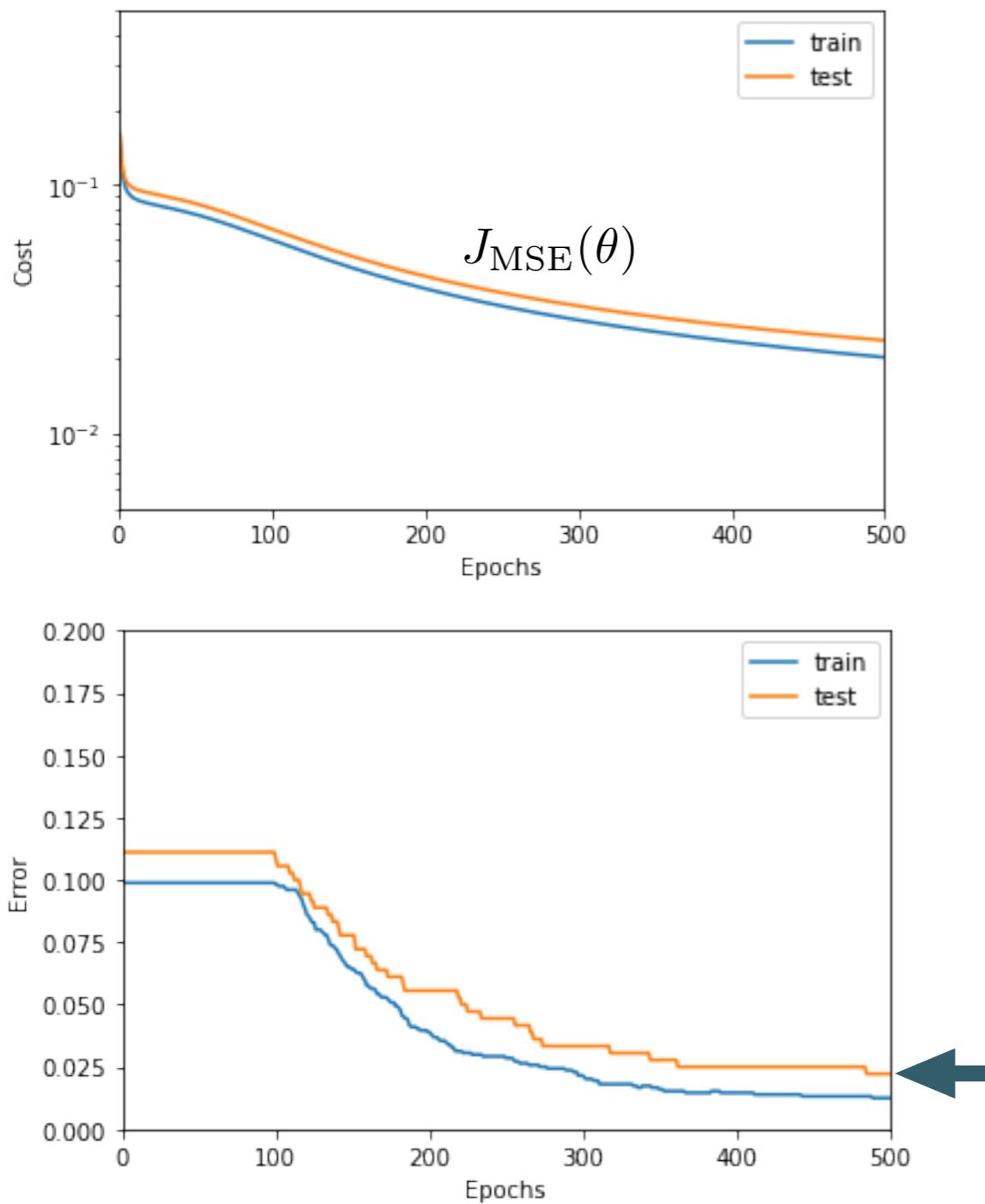
$$\theta \leftarrow \theta - \alpha \nabla_{\theta} J(\theta)$$

In coordinates:

$$\theta_k \leftarrow \theta_k - \alpha \frac{\partial J(\theta)}{\partial \theta_k}$$

$\alpha$  : learning rate.

# Results for MNIST Lightweight



- Inputs rescaled to  $[0,1]$
- Weights and bias initialised to zero
- Learning rate: 0.5
- 500 Iterations ('epochs')
- Training cost monotonically decreasing.
- Error / Performance measure:

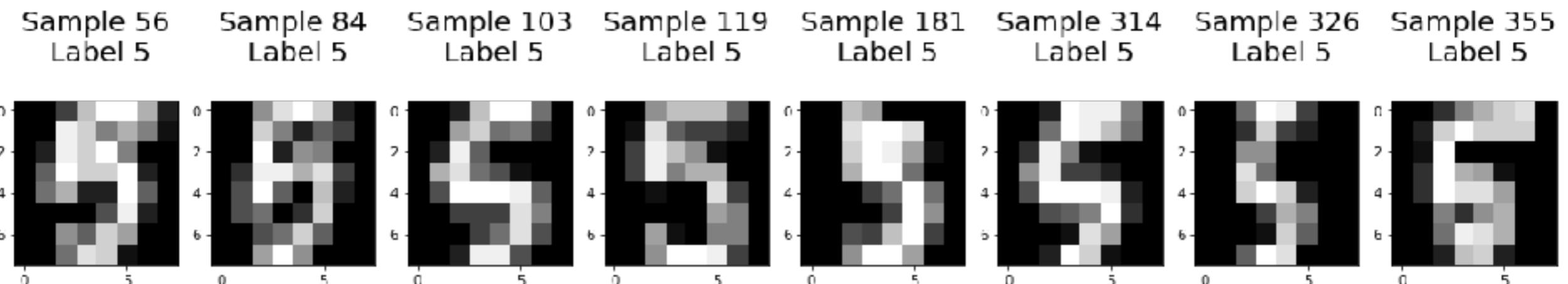
$$\text{error} = \frac{|\{\hat{y} \neq y\}|}{m_{\text{test}}}$$

with  $\hat{y} = \text{round}(h_{\theta}(\mathbf{x}))$

- After 500 epochs:
  - Training error: 1.2%
  - Test error: 2.2%  
(8 samples misclassified)

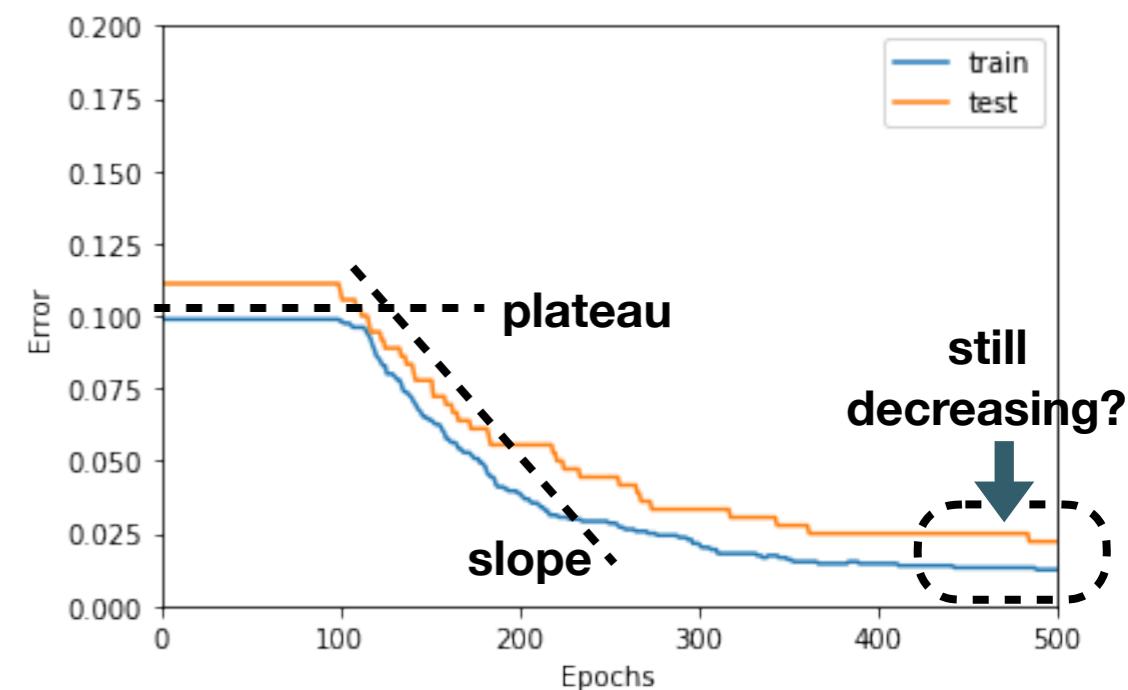
# Results for MNIST Lightweight

## Mis-classified images of the test set:



# Possible Improvements

- Try different learning rates
  - Learning rate needs to be tuned to the problem at hand.
- Run more epochs
  - Have we arrived at the bottom of the valley?
- More care when initialising parameters and preparing the input data
  - Do the initial weights properly match the distribution of the input data?
- Other cost function
  - Are there other cost functions that better capture the differences between predictions and labels?

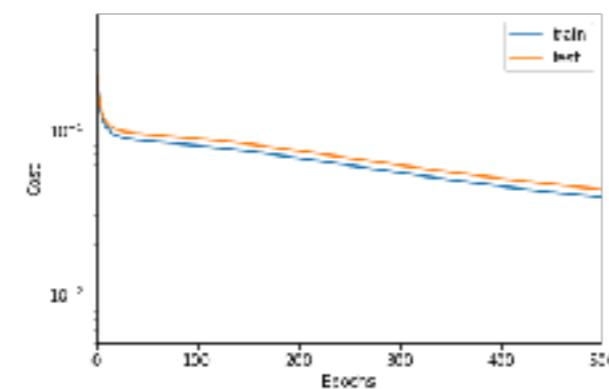


# Effect of Different Learning Rates

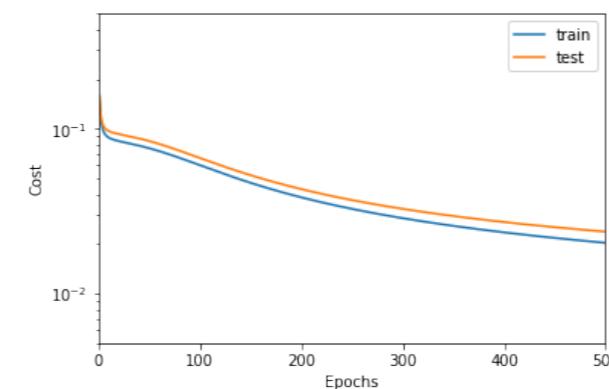
Running batch gradient descent over 500 epochs with different learning rates...

Learning rates here atypically large — usually much smaller rates are chosen.

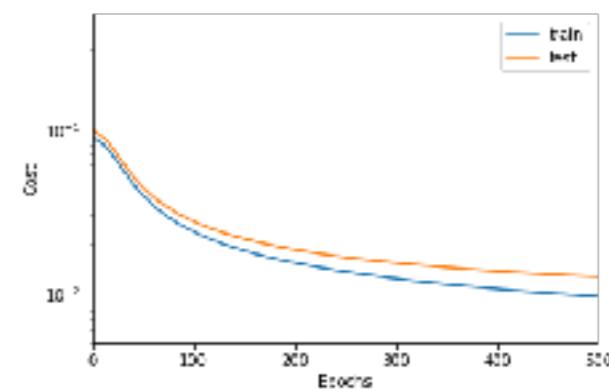
**learning rate 0.2**



**learning rate 0.5**

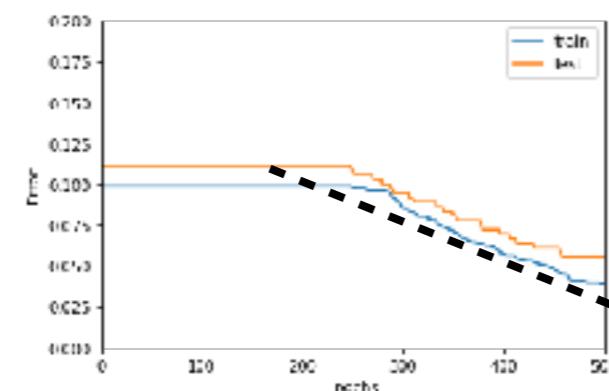


**learning rate 2.0**

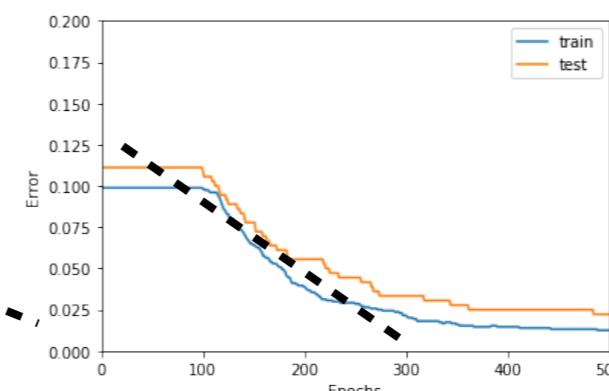


Error rate starts at 10% — corresponds to error rate of dummy predictor.

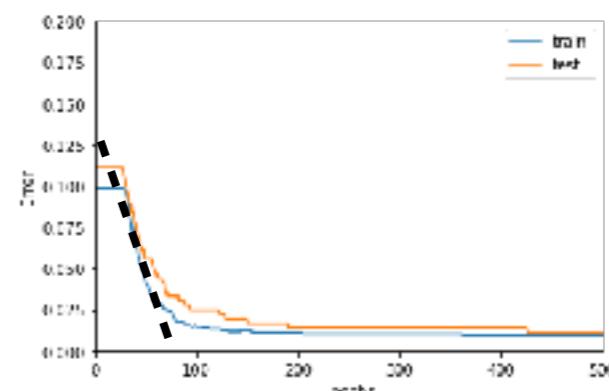
Dummy Predictor:  
None of the digits is a '5'.



slow decrease in cost and error rates  
more epochs would be needed here !



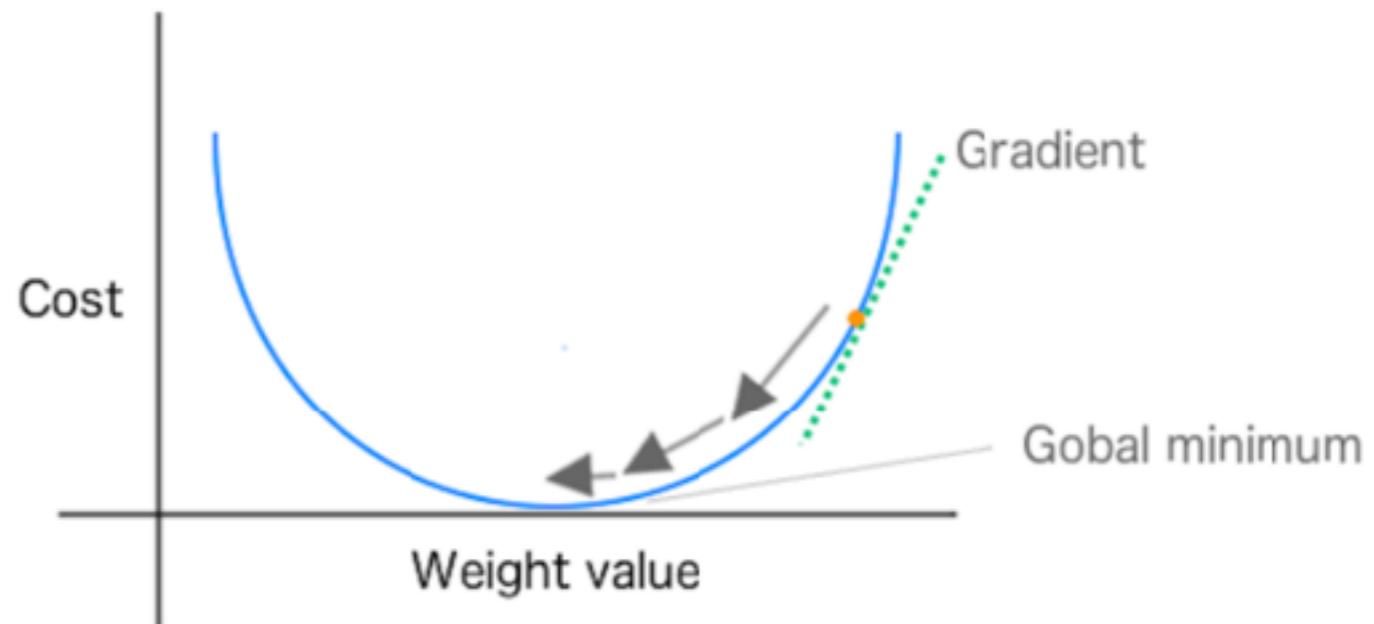
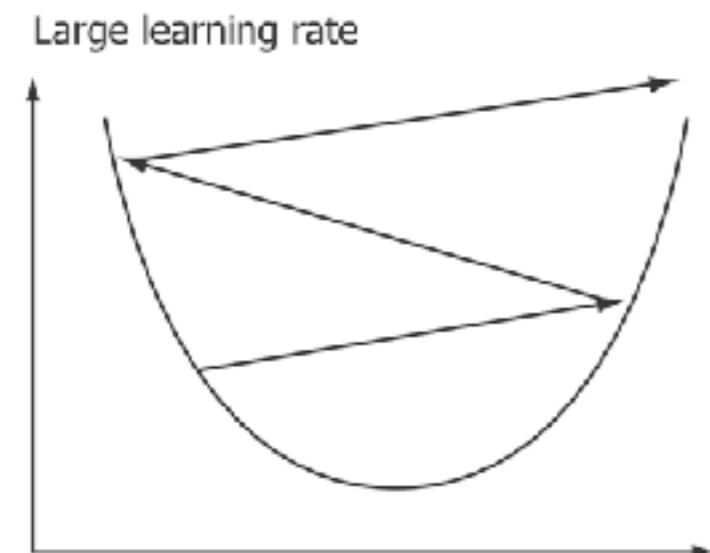
faster decrease in cost and error rates



Further reading: <https://www.jeremyjordan.me/nn-learning-rate/>

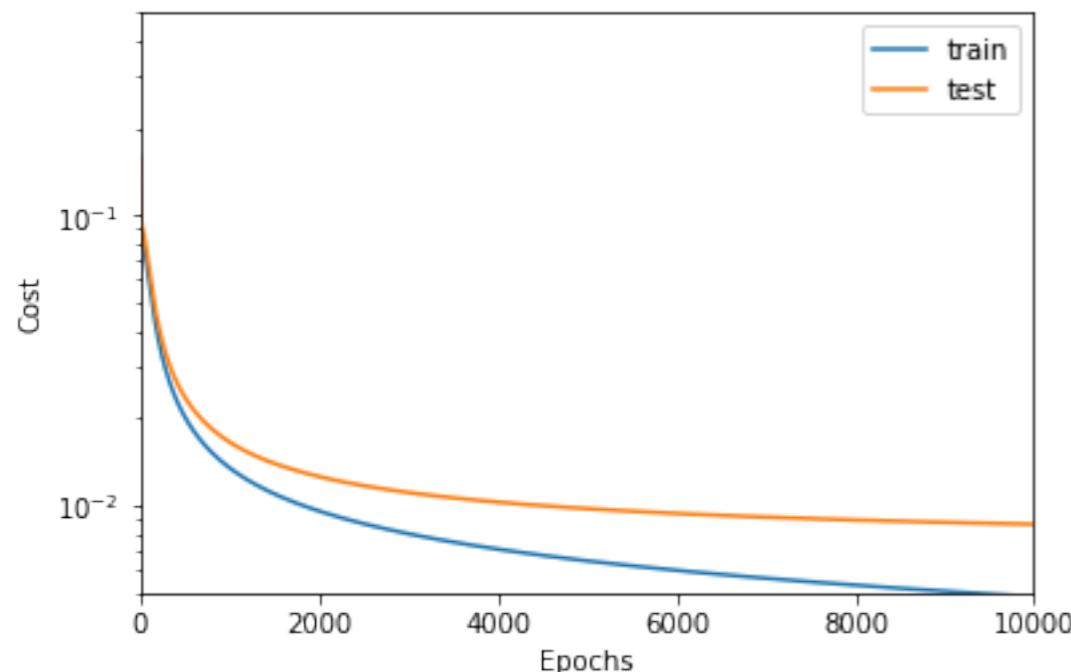
# Learning Rate $\alpha$

- Determines the learning speed
  - If too large: Oscillation around the minimum or even divergence
  - If too small: Slow convergence



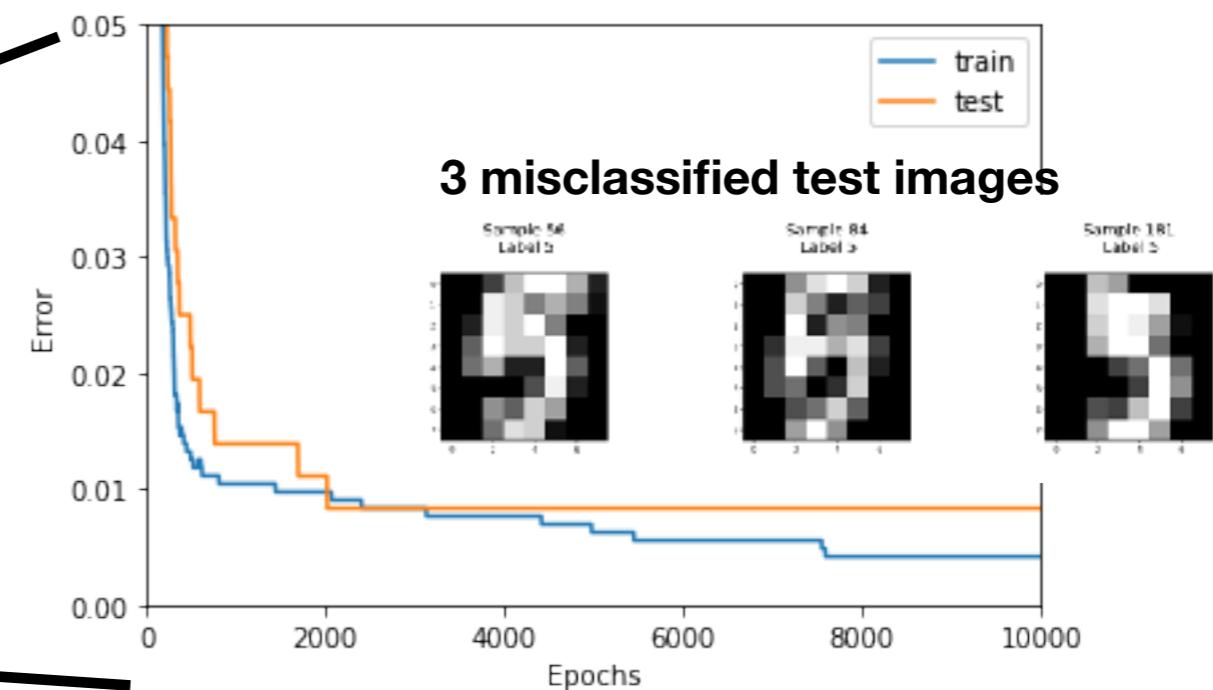
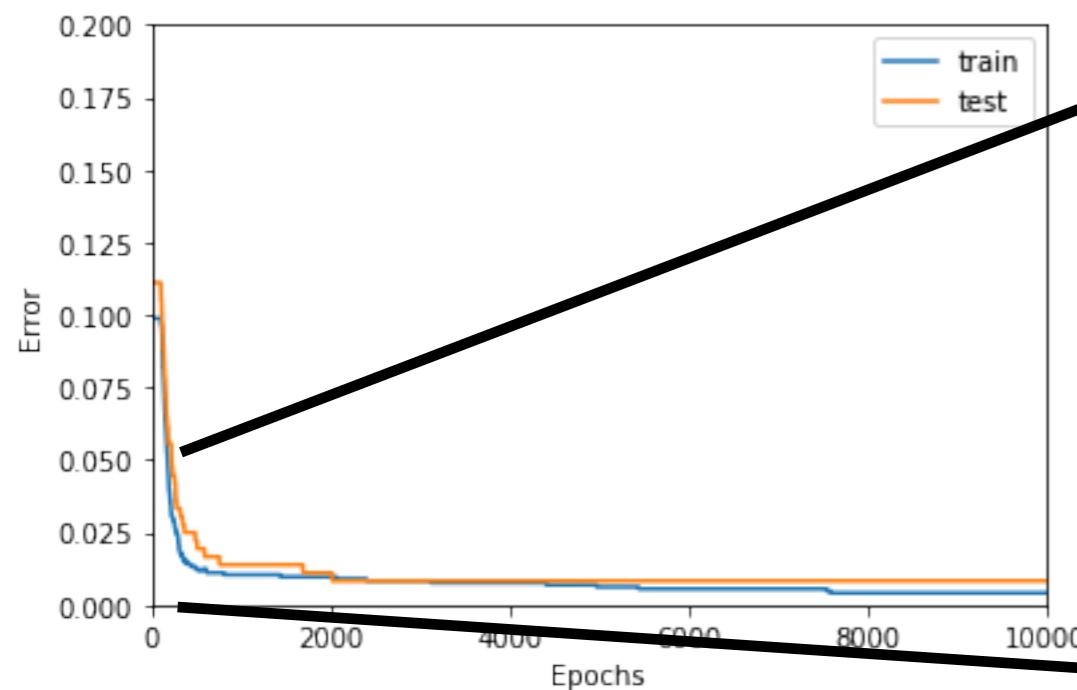
- Needs to be tuned to a given problem
  - An optimal value for a given problem may not work for another problem.
  - Thanks to data normalisation it's largely independent of the scale of the original data.

# Effect of More Epochs



After 2'000 epochs the test performance (test error) cannot be further reduced.

After 2'000 epochs the model gets overfitted. This important topic will be further covered in weeks 3 and 5).



# Issues with MSE Cost and Alternatives

- Recall that the **gradient of the MSE cost** contains in each summand a factor  $\hat{y}^{(i)}(1 - \hat{y}^{(i)}) \leq 1/4$   
As the model output  $\hat{y}^{(i)}$  gets closer to either 0 or 1 this expression can get very small.
- In result, the change in parameters can get very small and training can get stalled or stuck.
- For classification tasks, the **cross entropy cost function** (defined below) is better suited and the gradient has nicer properties.
- To understand the cross entropy cost function, a probabilistic consideration is needed — or stated differently, by using the cross entropy cost function, a probabilistic perspective is adopted.

# Probabilistic Interpretation of the Model Output

- The quantity  $h_\theta(\mathbf{x}) \in [0, 1]$  can be interpreted as the probability for ‘observing’  $y = 1$  given  $\mathbf{x}$ . Similarly,  $1 - h_\theta(\mathbf{x})$  is then interpreted as the probability of ‘observing’  $y = 0$  given  $\mathbf{x}$ .
- In maths terms:

$$\begin{aligned} p(y = 1|\mathbf{x}) &= h_\theta(\mathbf{x}) \\ p(y = 0|\mathbf{x}) &= 1 - h_\theta(\mathbf{x}) \end{aligned}$$

- Since  $p(y|\mathbf{x})$  depends on  $\theta$  we sometimes indicate that explicitly by writing:

$$p(y|\mathbf{x}, \theta) \quad \text{or} \quad p_\theta(y|\mathbf{x})$$

# Cross-Entropy Loss Function

- Suited for classification problems.
- Based on model predicting the probability for observing class  $y$  given  $\mathbf{x}$  denoted by  $p_\theta(y|\mathbf{x})$ .
- Cross-Entropy Loss defined by

$$L_{\text{CE}}((\mathbf{x}, y), \theta) = -\log(p_\theta(y|\mathbf{x}))$$

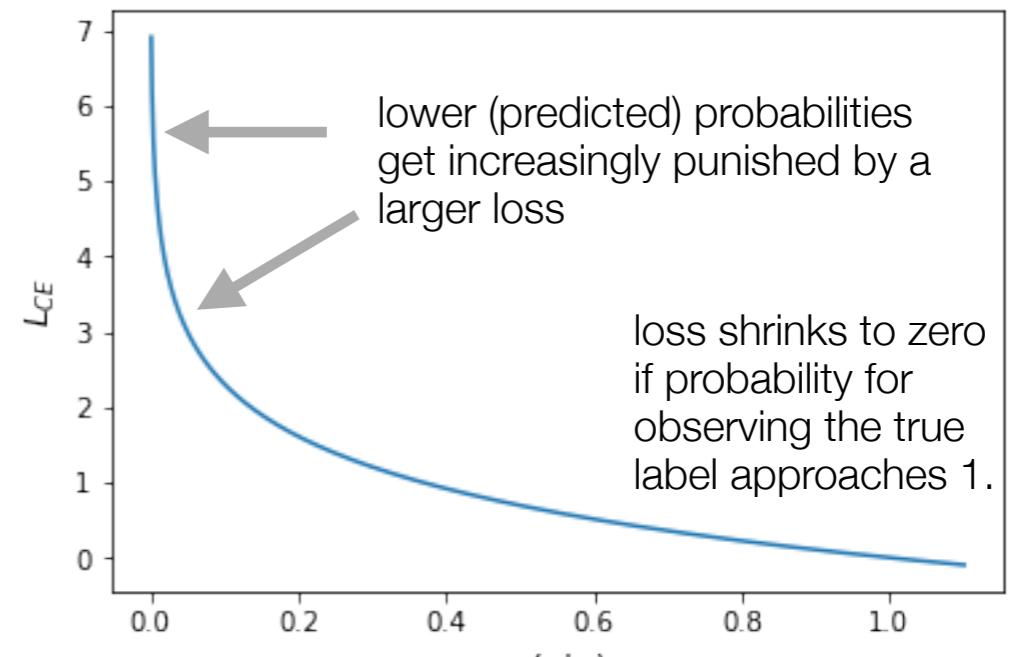
where  $y$  is the true label for  $\mathbf{x}$ .

- Example: Binary classification

$$p_\theta(y = 1|\mathbf{x}) = h_\theta(\mathbf{x})$$

$$p_\theta(y = 0|\mathbf{x}) = 1 - h_\theta(\mathbf{x})$$

$$L_{\text{CE}}((\mathbf{x}, y), \theta) = -y \log(h_\theta(\mathbf{x})) - (1 - y) \log(1 - h_\theta(\mathbf{x}))$$



probability predicted by the model

# Cross-Entropy Cost Function

- Cross entropy cost function : Averaged cross entropy loss

$$J_{\text{CE}}(\theta) = -\frac{1}{m} \sum_{i=1}^m \log \left( p(y^{(i)} | \mathbf{x}^{(i)}, \theta) \right)$$

- For the binary classification problem:

$$J_{\text{CE}}(\theta) = -\frac{1}{m} \sum_{i=1}^m \left( y^{(i)} \log(h_\theta(\mathbf{x}^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(\mathbf{x}^{(i)})) \right)$$

- In the given single layer perceptron problem the cross-entropy cost function is a *convex function*. As a result, the function has a single local minimum (a single critical point) which is the global minimum.  
⇒ As long as we choose the learning rate sufficiently small, we will always get to the global minimum with gradient descent.

# From Maximum Likelihood to Cross Entropy

- Maximum Likelihood is a widely used principle for determining the parameters of a model: Choose the parameters  $\theta$  such that the probability for observing the given dataset

$$\left\{ (\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, m_{\text{train}} \right\}$$

is maximised.

- In mathematical terms:

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta \in \mathbb{R}^n} \prod_{i=1}^m p(y^{(i)} | \mathbf{x}^{(i)}, \theta) \\ &= \arg \max_{\theta \in \mathbb{R}^n} \log \left( \prod_{i=1}^m p(y^{(i)} | \mathbf{x}^{(i)}, \theta) \right) \\ &= \arg \max_{\theta \in \mathbb{R}^n} \left( \sum_{i=1}^m \log p(y^{(i)} | \mathbf{x}^{(i)}, \theta) \right) \\ &= \arg \min_{\theta \in \mathbb{R}^n} \left( - \sum_{i=1}^m \log p(y^{(i)} | \mathbf{x}^{(i)}, \theta) \right)\end{aligned}$$

Cross-  
Entropy  
Cost

# Gradient of the Cross-Entropy Loss (Generalised Perceptron)

The calculation of the update rule for the cross-entropy cost function is more complicated. We compute it here for the binary classification case. We first state a result we have computed for the MSE cost:

$$\nabla h_\theta(\mathbf{x}) = h_\theta(\mathbf{x}) \cdot (1 - h_\theta(\mathbf{x})) \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \quad (*)$$

Now we can compute the gradient of the CE-loss:

$$\begin{aligned} -\nabla L_{\text{CE}}(\theta) &= \nabla (y \log(h_\theta(\mathbf{x})) + (1 - y) \log(1 - h_\theta(\mathbf{x}))) \\ &= \frac{y}{h_\theta(\mathbf{x})} \nabla h_\theta(\mathbf{x}) - \frac{1-y}{1-h_\theta(\mathbf{x})} \nabla h_\theta(\mathbf{x}) \\ &= \left( \frac{y}{h_\theta(\mathbf{x})} - \frac{1-y}{1-h_\theta(\mathbf{x})} \right) \nabla h_\theta(\mathbf{x}) \\ &\stackrel{(*)}{=} (y(1 - h_\theta(\mathbf{x})) - (1 - y)h_\theta(\mathbf{x})) \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \\ &= (y - h_\theta(\mathbf{x})) \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \end{aligned}$$

# Gradient of the Cross-Entropy Cost (Generalised Perceptron)

The gradient of the cost function is the gradient per sample ('loss') summed over all training samples:

$$\nabla J_{\text{CE}}(\theta) = \frac{1}{m} \sum_{i=1}^m \left( h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) \begin{pmatrix} \mathbf{x}^{(i)} \\ 1 \end{pmatrix}$$

Predicted Probability  
(for seeing label  $y=1$ )

Label

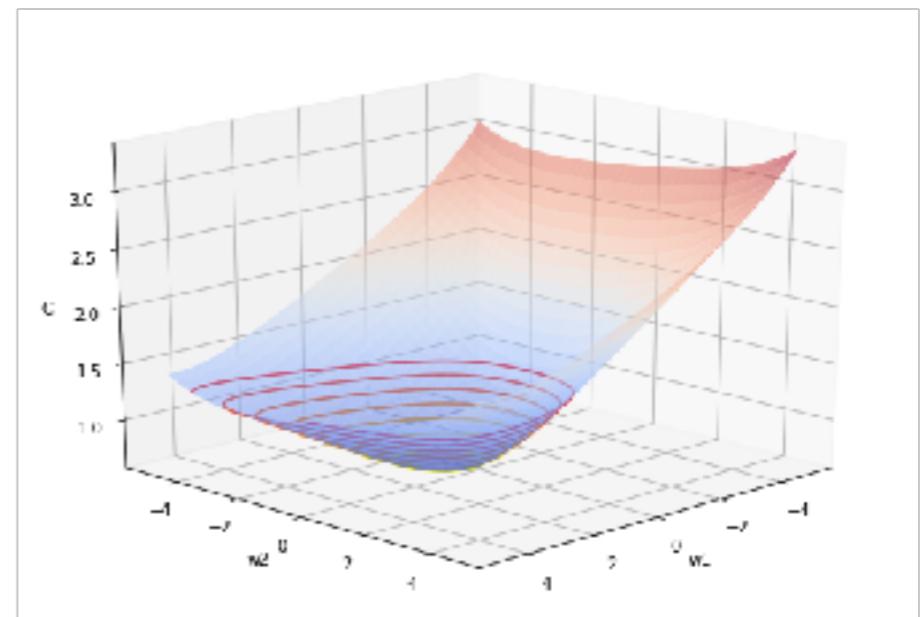
**Error Signal:** Contributes more  
for samples with a mismatch.

Very similar to the perceptron update rule! What are the differences?

# Update Rules for Generalised Perceptron based on Cross Entropy Cost

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(\mathbf{x}^{(i)}) - y^{(i)}) \mathbf{x}^{(i)}$$
$$b \leftarrow b - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(\mathbf{x}^{(i)}) - y^{(i)})$$

The cross entropy cost function for the generalised perceptron is a convex function, i.e. bowl-shaped. Therefore, the above scheme is guaranteed to find the global minimum - as long as the learning rate is kept sufficiently small.



# Note on Entropy and Cross-Entropy

- Probability distributions encode information: For a distribution  $p(x)$  it is measured with **Shannon's Entropy** as

$$H(p) = - \sum_x p(x) \ln p(x) = E_{X \sim p} [-\ln p(X)]$$

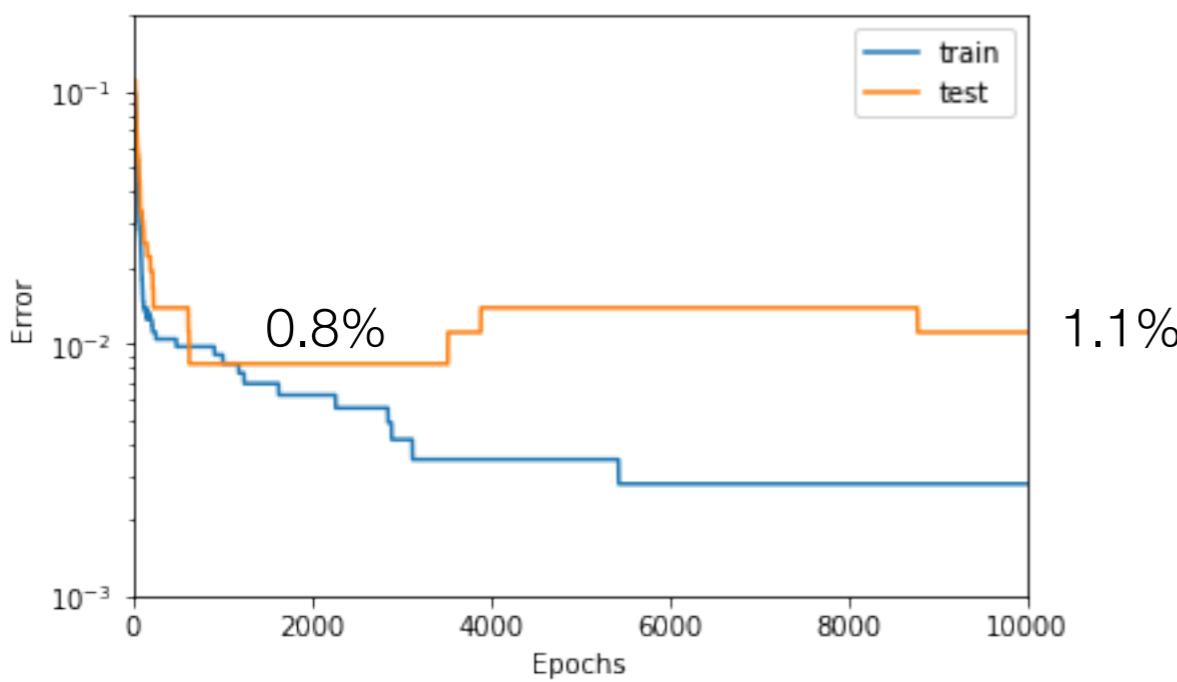
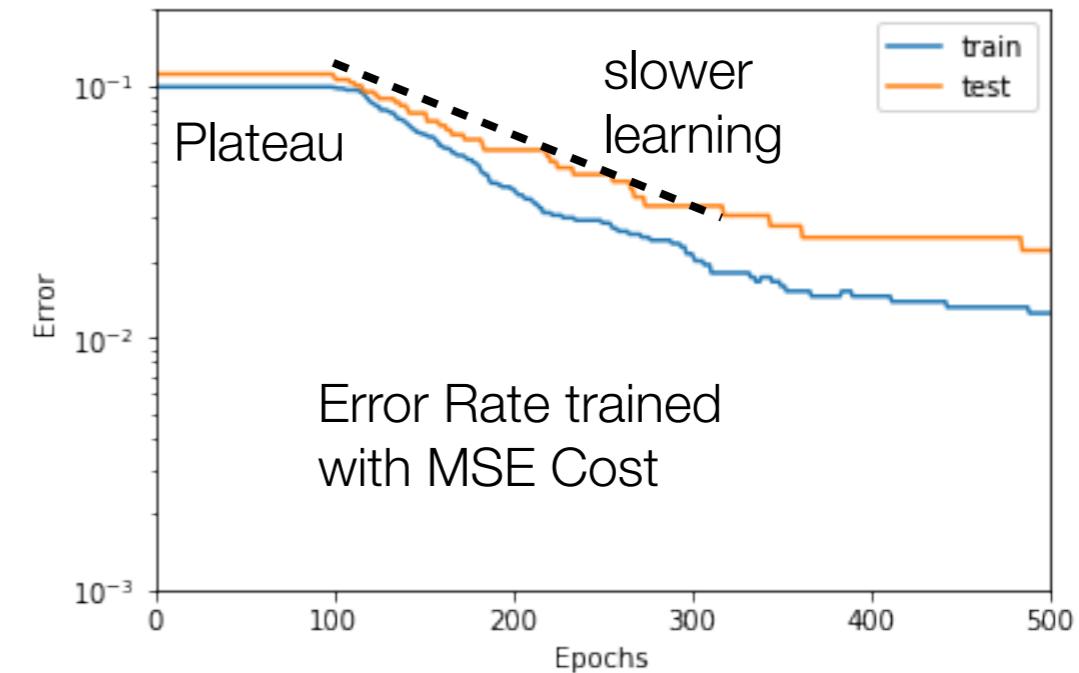
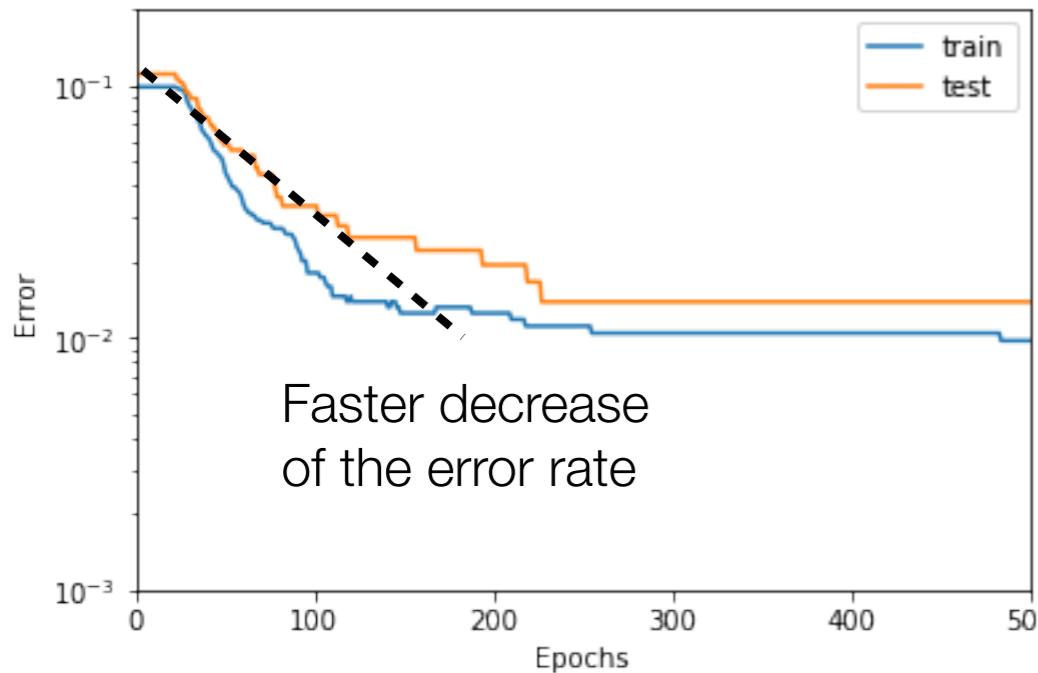
## Basic Idea of Information Theory

The more one knows about a topic, the less new information one is apt to get about it. If an event is very probable, it is no surprise when it happens and thus provides little new information. Inversely, if the event was improbable, it is much more informative that the event happened. Therefore, the information content is an increasing function of the inverse of the probability of the event ( $1/p$ ). Now, if more events may happen, entropy measures the average information content you can expect to get if one of the events actually happens. This implies that casting a die has more entropy than tossing a coin because each outcome of the die has smaller probability than each outcome of the coin. (**Source: Wikipedia**)

- **Cross-Entropy** measures the information of one distribution w.r.t. to another

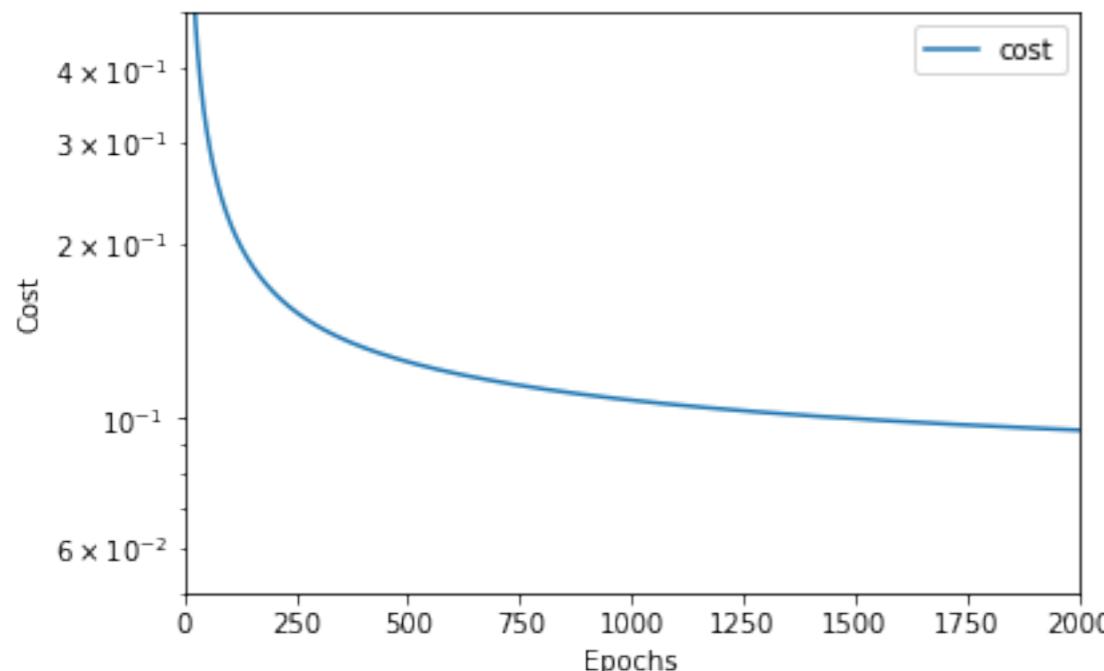
$$H(p, q) = - \sum_x p(x) \ln q(x) = E_{X \sim p} [-\ln q(X)]$$

# Results for MNIST Lightweight



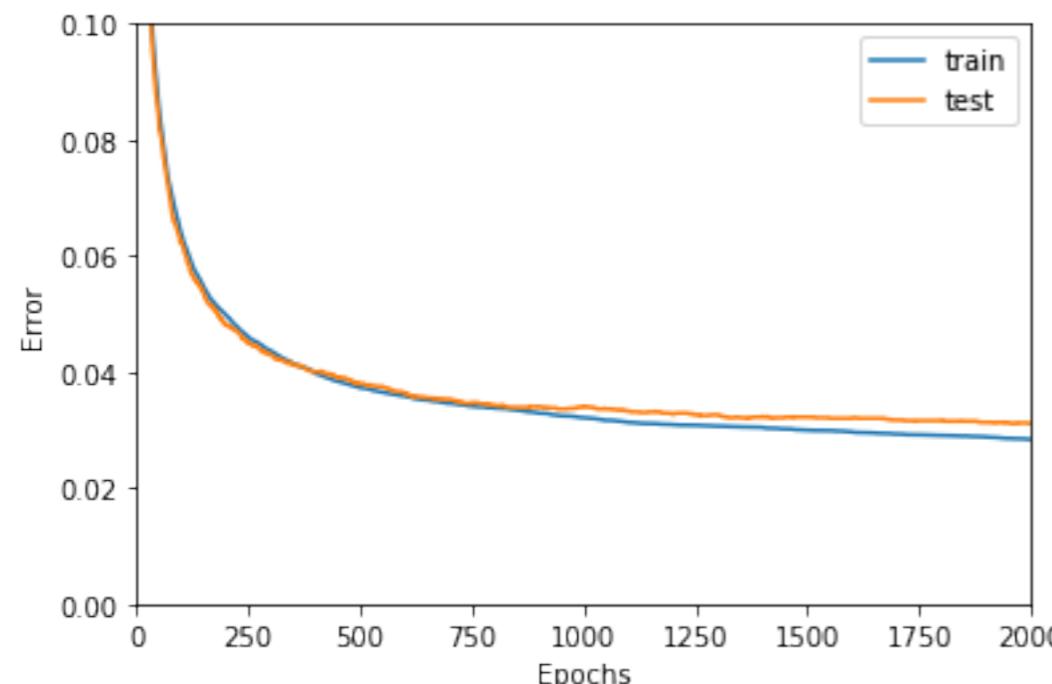
- Same settings as with MSE cost:
  - Inputs rescaled to  $[0, 1]$
  - Weights and bias initialised to zero
  - Learning rate: 0.5
- CE Cost performs better:
  - After 500 epochs: 1.4% (MSE: 2.2%)
  - After 2000 epochs: 0.8% (MSE: 1.1%)
  - Misclassifications after 500 epochs: 5

# Results for Original MNIST



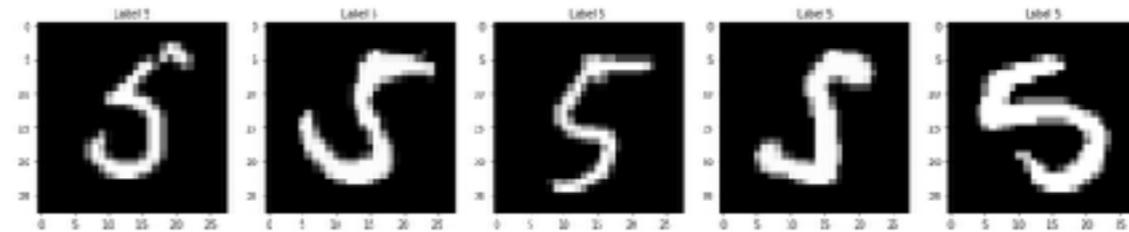
```
from sklearn.datasets import fetch_mldata
from sklearn.model_selection import train_test_split
mnist = fetch_mldata('MNIST original')
x, y = mnist['data'], np.array(mnist['target'], dtype='int')

x_train0, x_test0, y_train, y_test = train_test_split(x, y, test_size=10000, random_state=0)
xmax = np.max(x_train0)
xmin = np.min(x_train0)
x_train = x_train0 / xmax
x_test = x_test0 / xmax
```

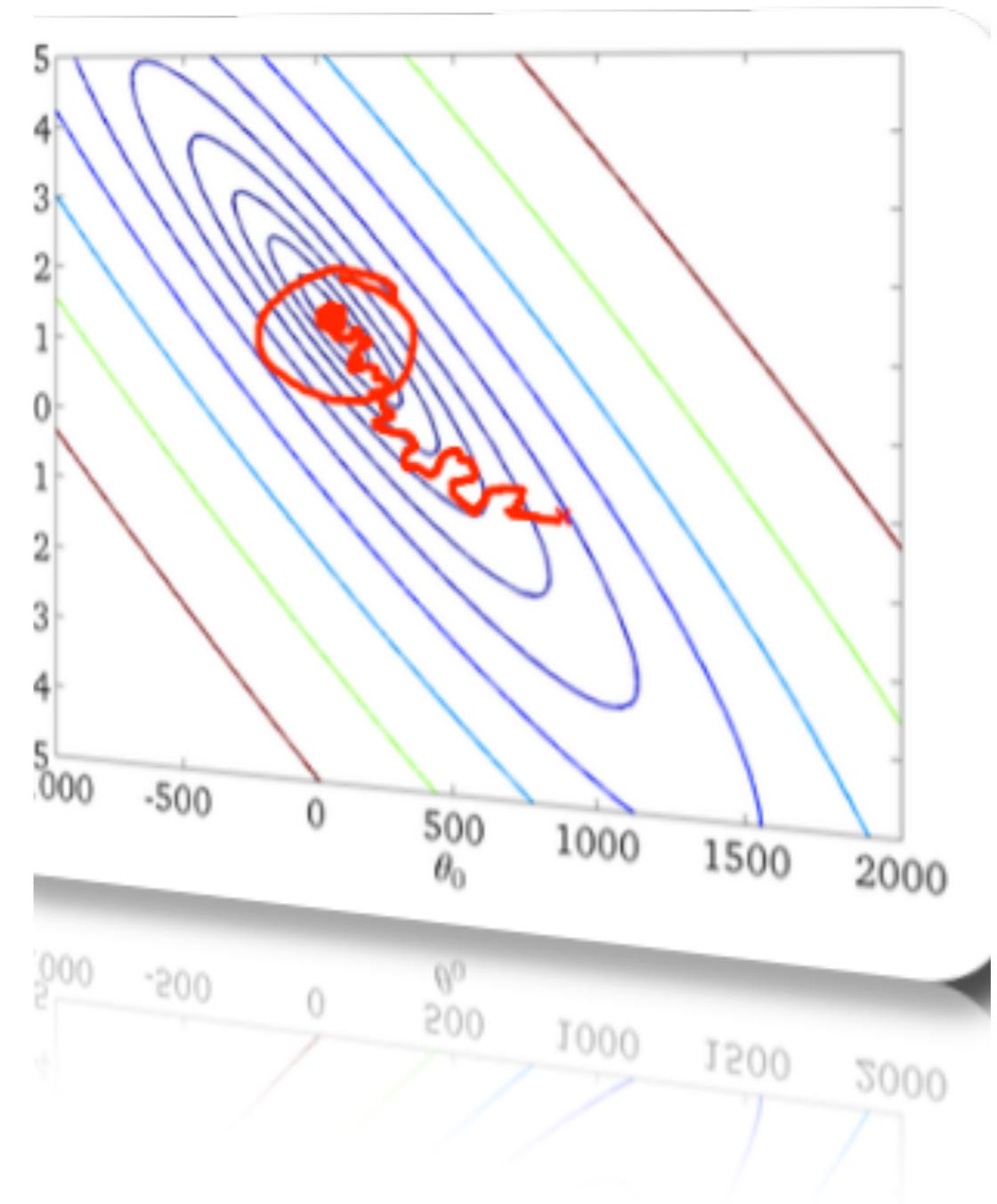


- Learning rate: 1.0
- 2000 epochs
- Training error: 2.84%
- Test error: 3.11% (311 / 10'000)

**Some images not identified as '5':**



# Stochastic Gradient Descent



# Motivation for Stochastic Gradient Descent

- So far, we have computed the gradient of the cost function defined on all the training set. This procedure is referred to as ***Batch Gradient Descent***.
- This can be costly: It involves the evaluation of the model and its derivatives for all the samples in the training set.
- This is not an issue for small models like the single (layer) perceptron — but for deep neural networks and if the training dataset is large this becomes a problem.
- Idea of ***Stochastic Gradient Descent*** or ***Mini-Batch Gradient Descent***: Use less training samples to define a single update for moving in parameter space.

# Idea behind Stochastic Gradient Descent and Mini-Batch Gradient Descent

- Cost functions ( $J_{\text{CE}}(\theta)$ ,  $J_{\text{MSE}}(\theta)$ ) are expressed as an arithmetic mean over per sample contributions ('loss').

$$J_{\text{CE}}(\theta) = -\frac{1}{m} \sum_{i=1}^m \log \left( p(y^{(i)} | \mathbf{x}^{(i)}, \theta) \right)$$

- The update rule for the parameters can be defined by including the contributions from less samples (here binary classification example).

$$\begin{pmatrix} \Delta w \\ \Delta b \end{pmatrix} = -\alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(\mathbf{x}^{(i)}) - y^{(i)}) \begin{pmatrix} \mathbf{x}^{(i)} \\ 1 \end{pmatrix}$$

- **Batch Gradient Descent:** Summation over all the training set.
- **Mini-Batch Gradient Descent:** Summation only over subset of the training samples
- **Stochastic Gradient Descent:** No summation , just a single sample
- However: No guarantee to move in parameter space in a direction that leads to smaller values of the cost function.

# Stochastic Gradient Descent (SGD)

1. Start with some initial  $\theta$ 's (for example random or 0)  
Here:  $\theta = (b, \mathbf{w})$

2. Select one training sample randomly  $(\mathbf{x}^{(i)}, y^{(i)})$   
and perform the update of the  $\theta$ 's with this example

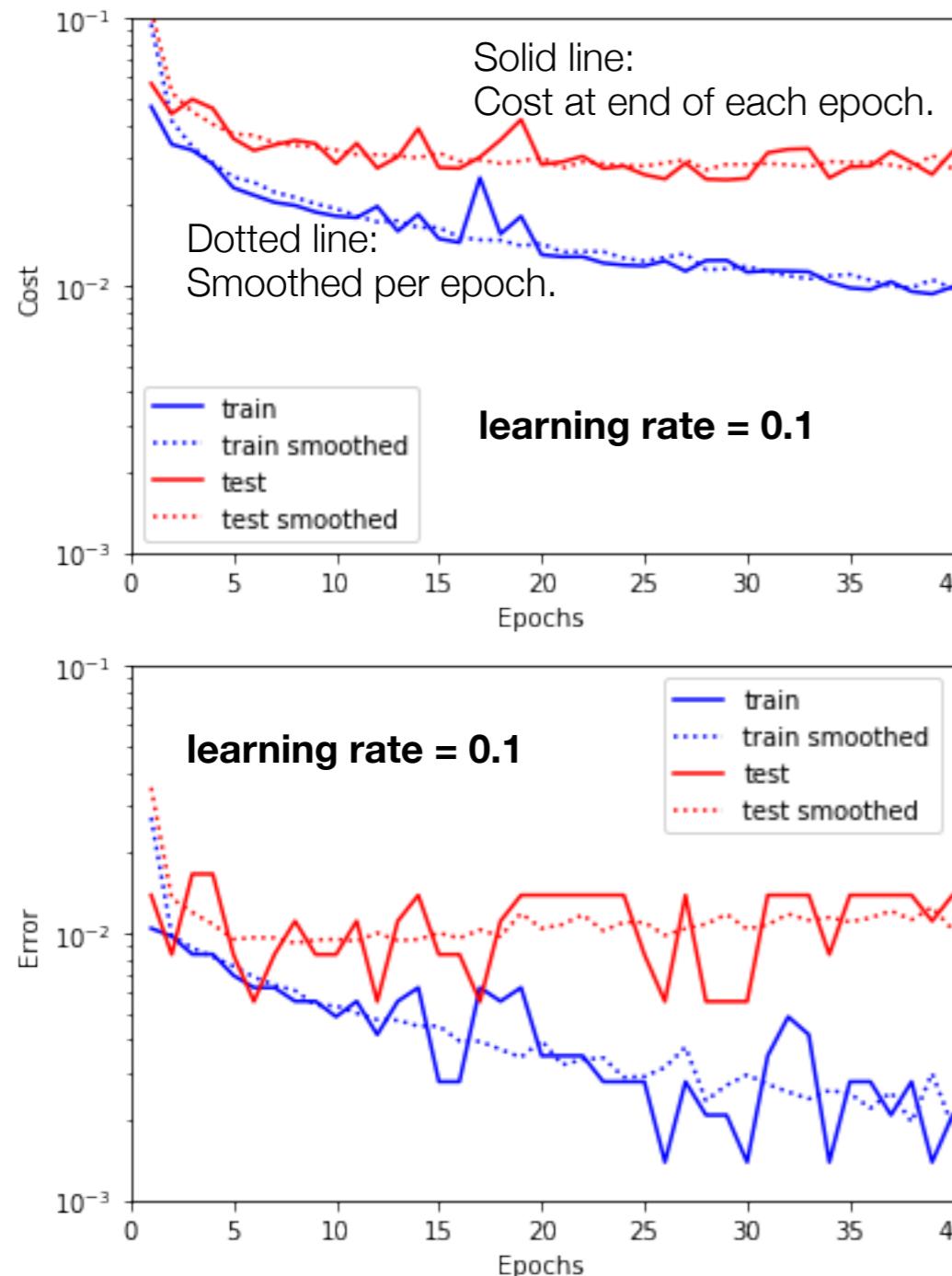
$$\begin{aligned}\mathbf{w} &\leftarrow \mathbf{w} - \alpha (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}) \mathbf{x}^{(i)} \\ b &\leftarrow b - \alpha (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})\end{aligned}$$

3. Loop in step 2 until convergence.

**Convergence** is trickier to observe as we may select favorable and less favorable points in the training set.

In practice we may compute an “averaged” cost function over the past  $U$  updates. Typically  $U \ll T$  and  $U \ll T$  if the total number of updates  $T$  is really large.

# Results for Lightweight MNIST (SGD)



Cost is not monotonically decreasing.

Values for the cost and error rates (both training and test data) fluctuate.

Reduce learning rate!

40 epochs have been computed to reach an test error rate of ~1.1%. Fewer steps needed than with BGD.

one epoch = one pass through the training data  
= 1347 update steps

(for comparison with batch gradient descent)

Random motion in parameter space leads to uncertainty in error rate and cost.

# Characteristics of SGD

## General Characteristics

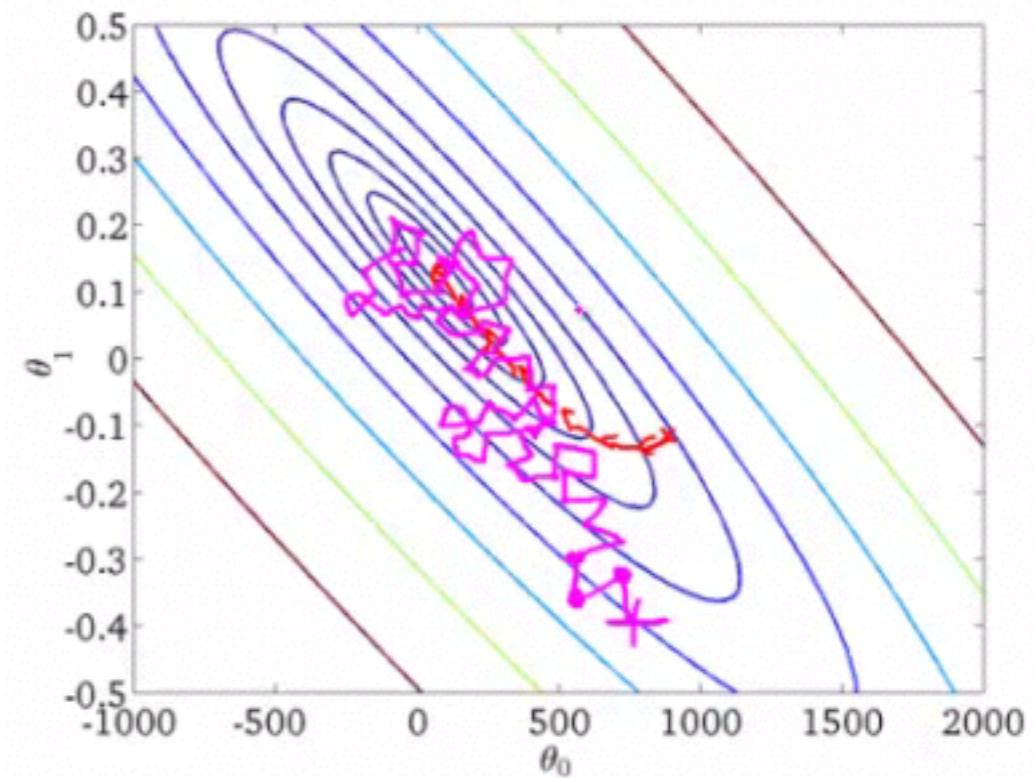
- Tends to move in direction of the global minimum, but not always.
- Never converges like batch gradient descent does, but ends up wandering around close to the global minimum — not a problem in practice as long as we are close to the minimum. Regularising effect!
- As batch gradient descent, the learning principle is generalisable to many other “hypothesis families”.

## Advantages

- Faster since parameters are updated for each training sample.
- It can handle very large sets of data - great for learning on huge datasets that do not fit in memory (**‘out-of-core learning’**).
- It allows for incremental learning (**‘online learning’**), i.e. on-the-fly adjustment of the model parameters on new incoming data.

## Disadvantage

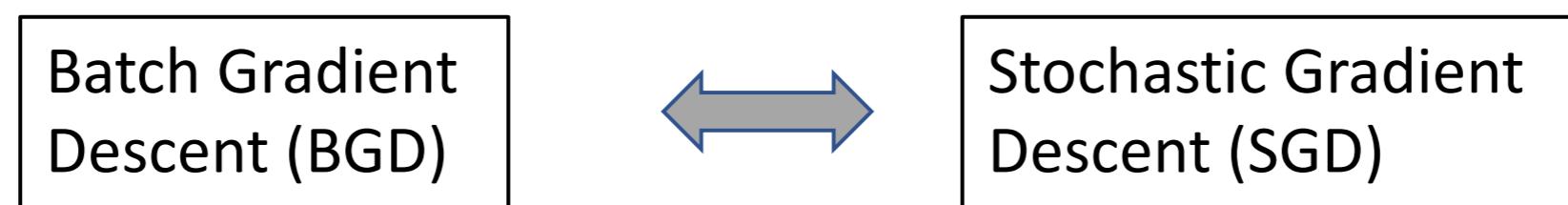
- Not easily parallelised.



Online learning is great for systems that receive data as a continuous flow!

# Mini-Batch Gradient Descent (MBGD)

**Find a good compromise** between



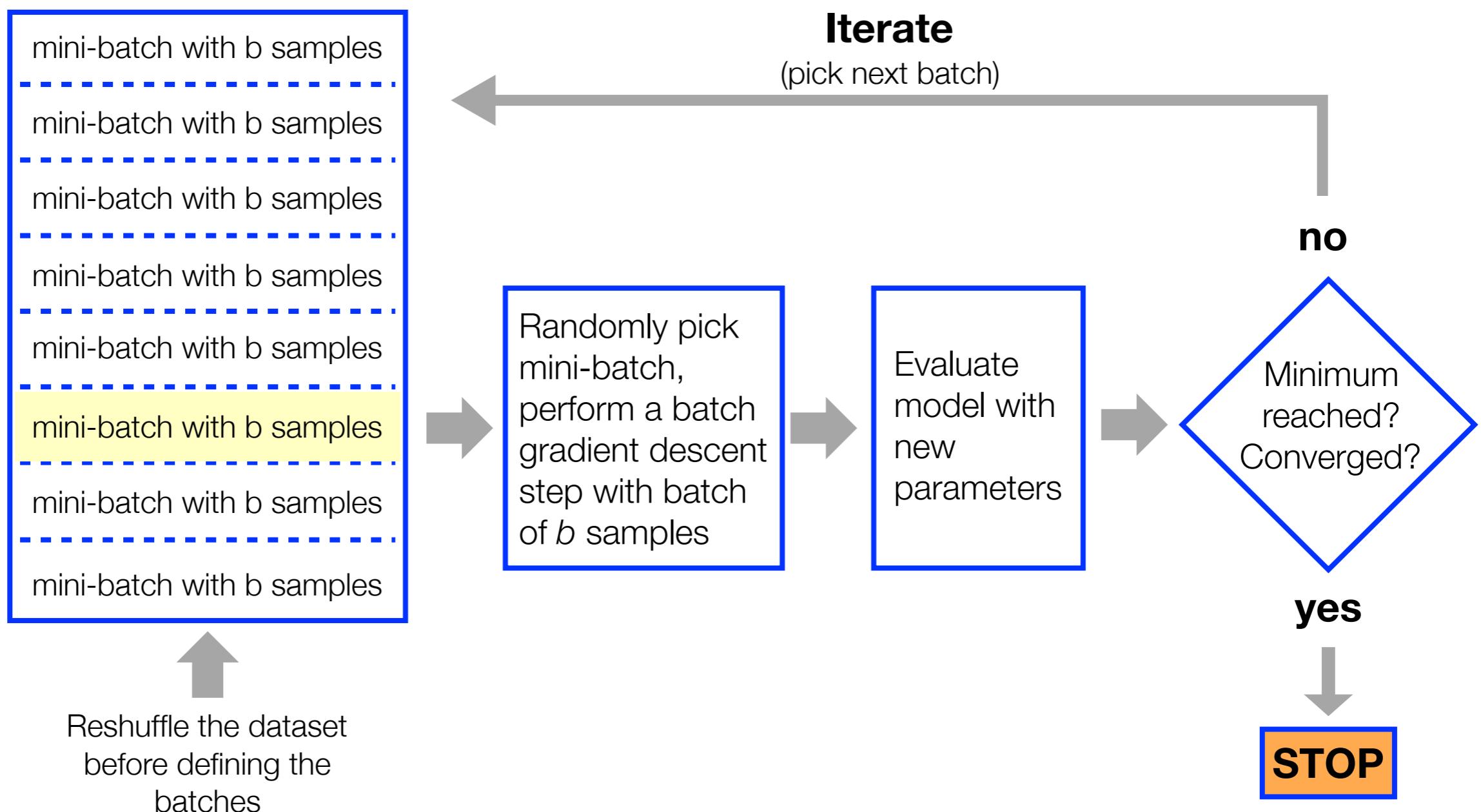
## Advantage

- Level of noise seen in the learning curve of stochastic gradient descent is reduced
- Yet allows for vectorised implementations, potentially more efficient (pipelining)
- Straightforward distribution of computational load by computing mini-batches on different machines or cores - i.e. adopting a map/reduce pattern:
  - map = distribute the mini-batches (e.g. each with 50 samples)
  - reduce = update parameters with the results of the mini-batches

## Disadvantage

- Batch size needs to be optimised

# Mini-Batch Gradient Descent Principle



# Mini-Batch Gradient Descent (MBGD)

1. Start with some initial  $\theta$ 's (for example random or 0)  
Here:  $\theta = (b, \mathbf{w})$
2. Randomly select mini-batch of  $b$  training sample with indices  $i_1, \dots, i_b$ , i.e.

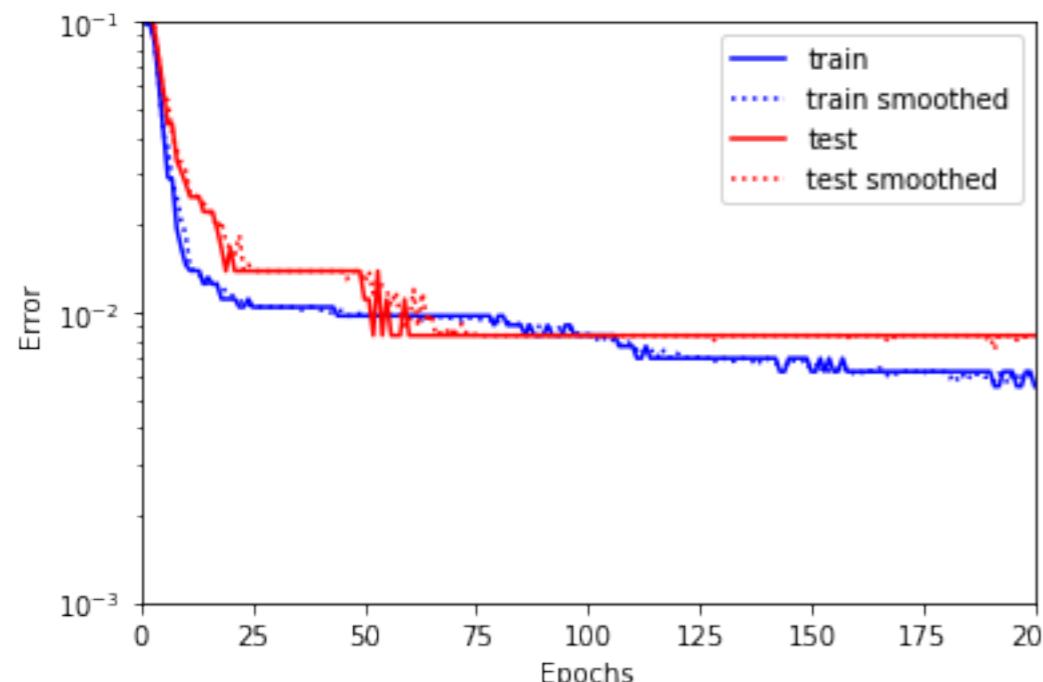
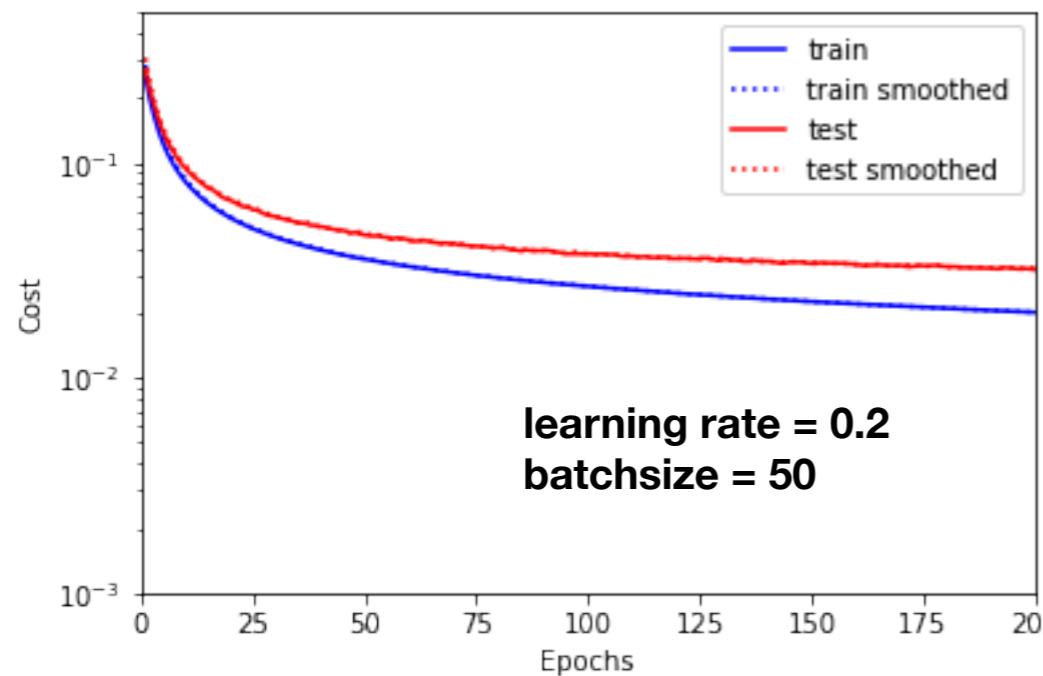
$$(\mathbf{x}^{(i_1)}, y^{(i_1)}), \dots, (\mathbf{x}^{(i_b)}, y^{(i_b)})$$

and perform the update of the  $\theta$ 's with this example

$$\begin{aligned}\mathbf{w} &\leftarrow \mathbf{w} - \alpha \sum_{j=1}^b (h_\theta(\mathbf{x}^{(i_j)}) - y^{(i_j)}) \mathbf{x}^{(i_j)} \\ b &\leftarrow b - \alpha \sum_{j=1}^b (h_\theta(\mathbf{x}^{(i_j)}) - y^{(i_j)})\end{aligned}$$

3. Loop in step 2 until convergence.

# Results for Lightweight MNIST



Cost is much smoother than with SGD,  
typically monotonically decreasing.

Values for the cost and error rates (both  
training and test data) fluctuate.

Reduce learning rate!

Test error rate of ~0.8% reached after  
~60 epochs.

As compared with SGD much less  
wiggling in the error rate due to  
random motions in parameter space.

# Characteristics of MBGD

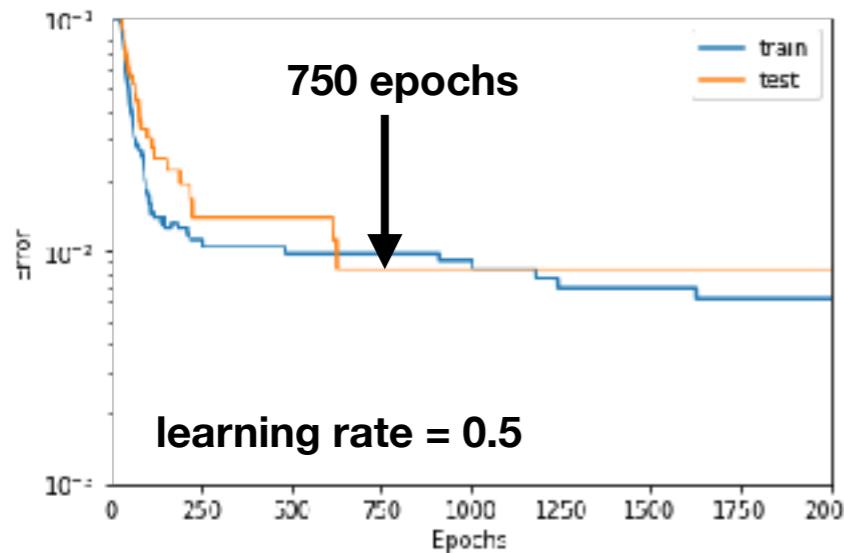
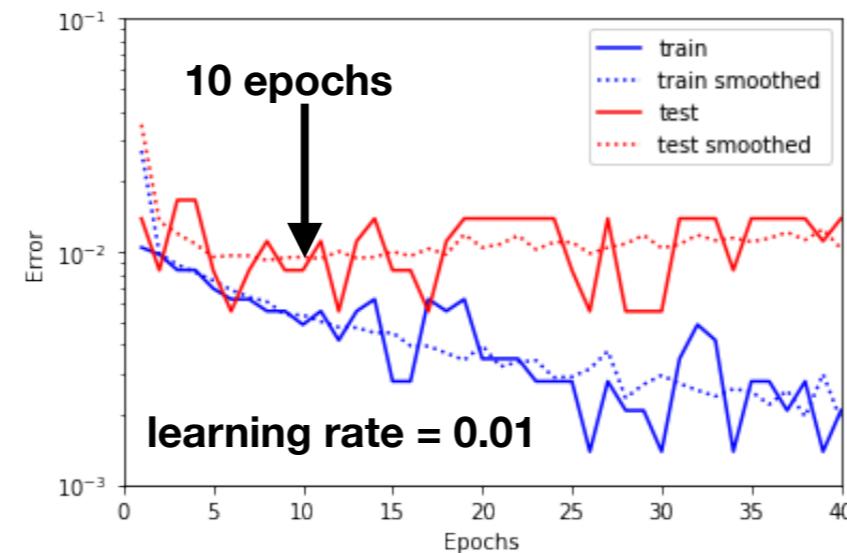
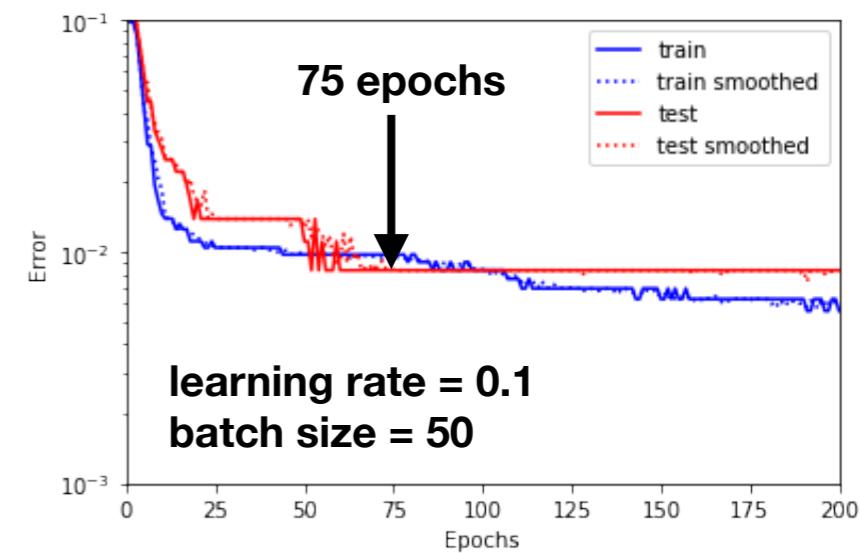
## General Characteristics

- Tends to move in direction of the global minimum, but not always.
- Ends up wandering closely around to the global minimum.  
Regularising effect!
- As batch gradient descent, the learning principle is generalisable to many other “hypothesis families”.

## Advantages

- Faster since parameters are updated for each mini-batch.
- It can handle very large sets of data - great for learning on huge datasets that do not fit in memory (**‘out-of-core learning’**).
- It allows for incremental learning (**‘online learning’**), i.e. on-the-fly adjustment of the model parameters on new incoming data.
- Easily parallelised on GPU and in HPC.

# Comparison of Gradient Descent Schemes

**BGD****SGD****MBGD**

- Smooth, strictly decreasing curves.
- Cost curves strictly decreasing. Can reach minimum.
- Slow for large  $m$ . Many epochs needed.
- Learning rate can be large( $r$ ).
- No out-of-core support - all data in RAM.
- Easy to parallelise.

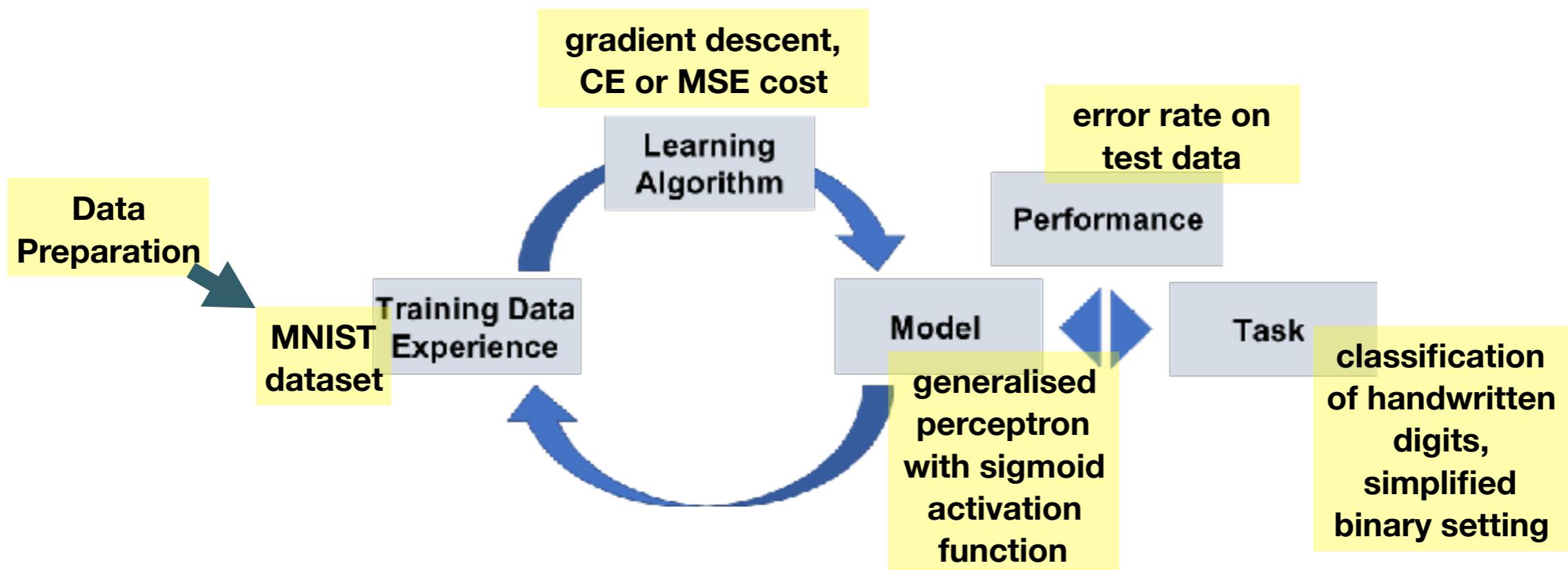
- Wiggling curves, need smoothing.
- Cost curve not necessarily decreasing. Wiggles around minimum.
- Fast for large  $m$ . Only few epochs needed.
- Learning rate must be small(er),
- Out-of-core support - not all data to be kept in RAM of a single machine.
- Not easy to parallelise.

- Slightly wiggling curves.
- Cost curve typically, but not necessarily decreasing. Wiggles slightly around minimum.
- Fast for large  $m$ . Much less epochs than BGD but more than SGD needed.
- Medium learning rate.
- Out-of-core support - not all data to be kept in RAM of a single machine.
- Easy to parallelise.

# Characteristics of Gradient Descent

- Can handle very large sets of data – particularly, the SGD and MBG.
- Allows for incremental learning, i.e. on-the-fly adaptation of the model on new incoming data.
- Learning principle is generalisable to many other forms of models  $h_\theta(\mathbf{x})$  such as deep neural networks.
- It works ‘locally’ and finds local critical points (local minima, maxima, saddle points). It is not designed to find global minima in general.

# Summary - Ingredients of ML Algorithms



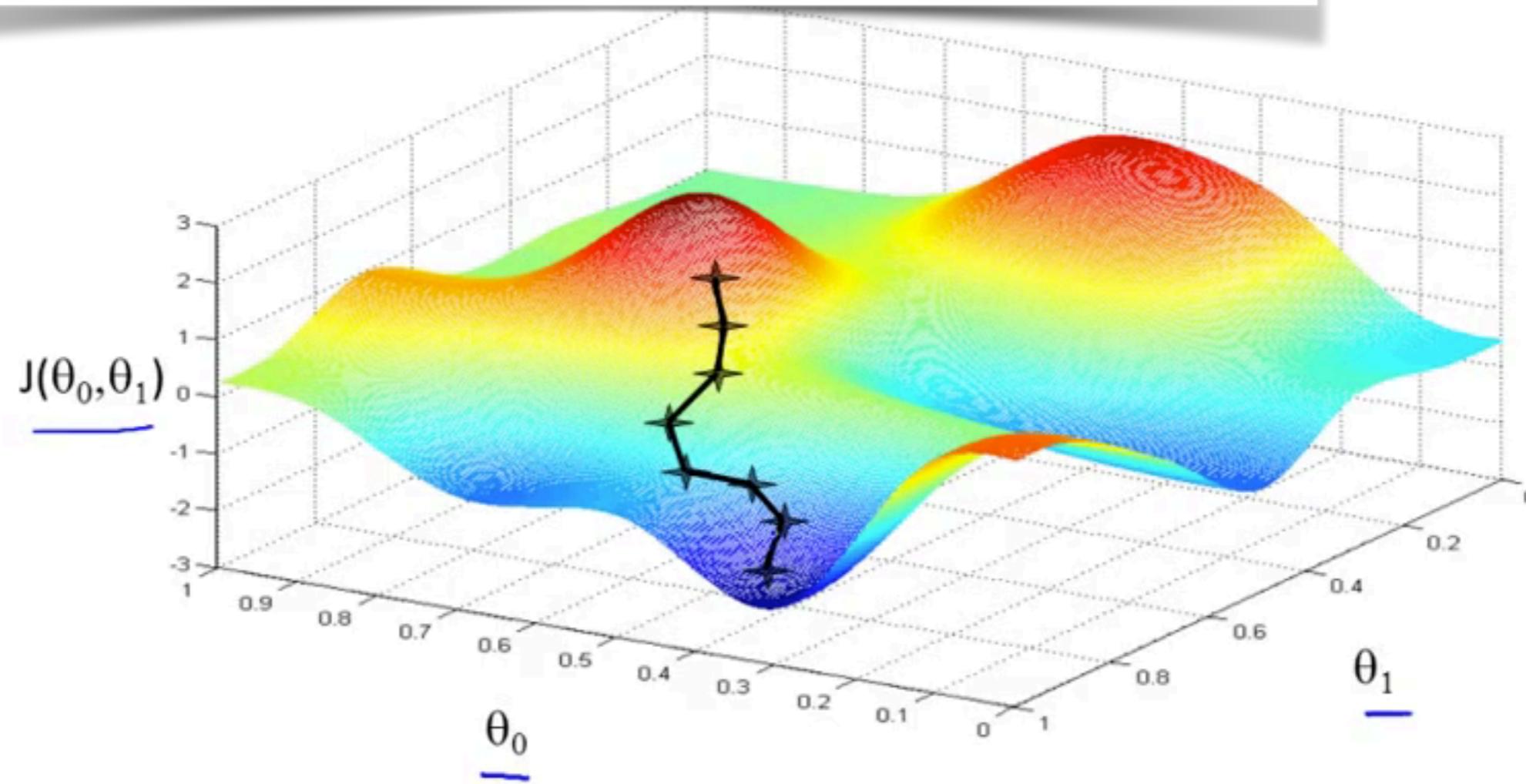
# Questions / Outlook

- How to evaluate the performance ('convergence') of the different schemes? When to stop?
- How to best choose
  - the learning rate?
  - the mini-batch size?  
... and how to evaluate that systematically?
- How can we improve the performance?

# Backup

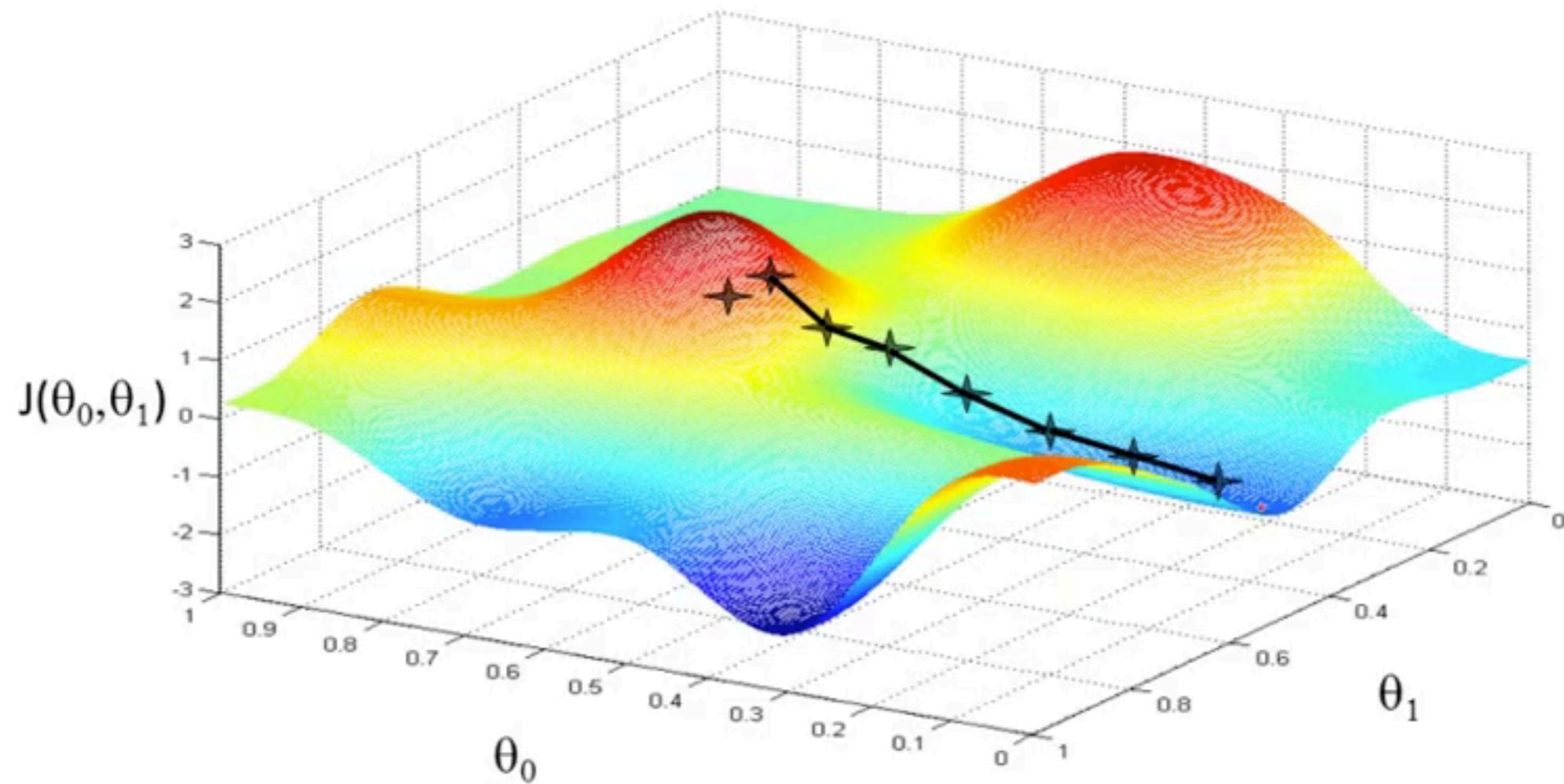
# Gradient Descent Intuition

From the initial position, take steps in the opposite direction of the gradient, “going downhill” until a minimum is reached.



Source: Andrew Ng - Machine Learning class - Stanford

# Gradient Descent Intuition



The path to the (local) minimum depends on the initial conditions. A slight difference in the initial parameter values, may lead to another minimum, a local minimum which is above the global minimum reached in the previous run.

In the case of a single unit perceptron and a sigmoid activation function, the  $J()$  function is convex, i.e. “bowl shaped” with a unique minimum, so we don’t need to care by getting stuck in a local minimum in this case.