

# Introduction to Machine Learning

Tu Bui

Central for Vision Speech and Signal Processing (CVSSP)

University of Surrey

Guildford, Surrey, United Kingdom

# Contents

- What is machine learning (ML)
- Linear regression
- Optimisation
- Classification
- ML in practice

# Machine Learning (ML)

- ML appears frequently in our daily life



Web search



Auto tagging



Product recommendation

# Machine Learning (ML)



- Definition

“Field of study that gives computers the ability to learn from and make predictions on **data** without explicitly programmed”

Arthur Samuel, 1959

- Learn from data
- NOT following static program instructions
- Humans are learning machines

# Why ML?

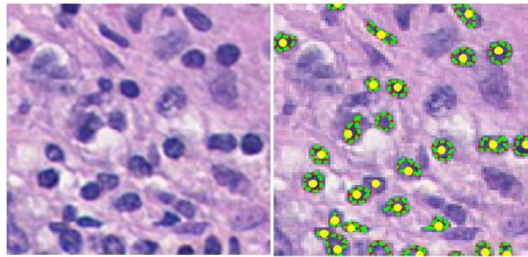
- We cannot program everything
- Some tasks are difficult to define algorithmically
  - Especially in computer vision, e.g. recognise objects
- Cost effective



# ML applications

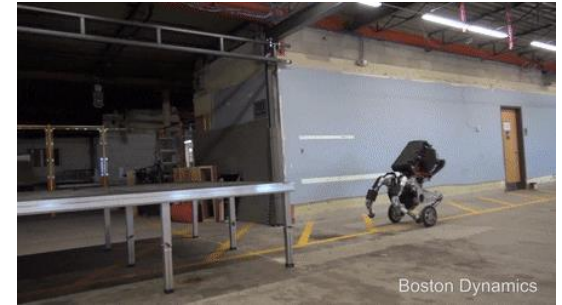


Market analysis

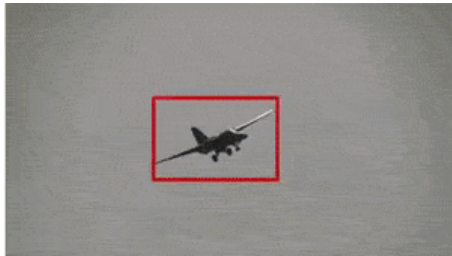


Medical imaging

<http://www.robots.ox.ac.uk>



Robot



Tracking/surveillance



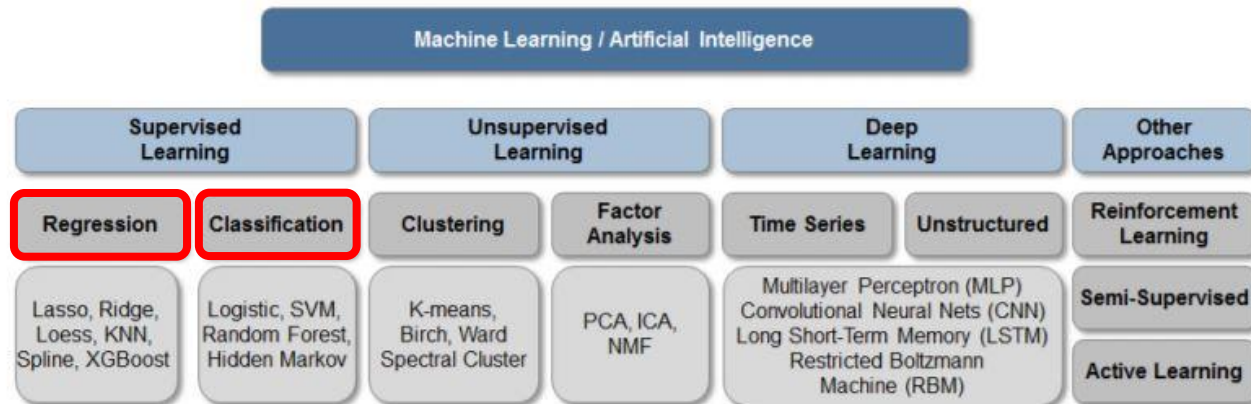
Self-driving car



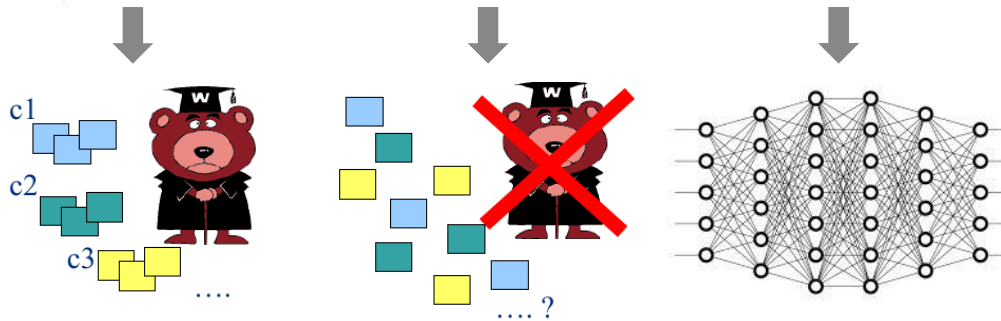
3D reconstruction/analysis

# ML types

Figure 39: Classification of Machine Learning techniques



Source: J.P.Morgan Macro QDS



# Supervised ML: regression vs. classification

- Regression
  - **“Continuous”** valued output
- Example: house price prediction



- Classification
  - **Discrete** set of output labels/classes
- Example: gender prediction





# Quiz: regression or classification?

1. You are helping a bank to write two programs that help decide whether the bank should provide credit card to new customers.
  - a) Your 1<sup>st</sup> program will output credit score for each customer based on their credit history. 1a. regression
  - b) Your 2<sup>nd</sup> program will rely on the output of the 1<sup>st</sup> program plus personal ratings of the bank manager(s) to decide if a customer will get a credit card or not. 1b. classification
2. You write a program to filter spam emails in your inbox. If an email is marked as spam by your program, it will be moved to the Junk folder. 2. classification
3. You are running a company that sells umbrella. You want to predict how many umbrellas will be sold in the next 3 months based on last 10 year sale data. 3. regression

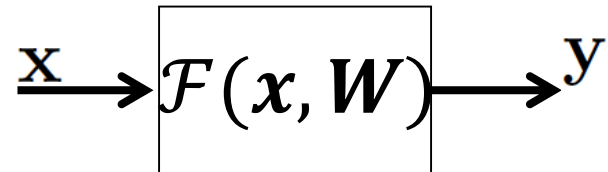
# Supervised learning: problem formulation

- Mathematically, the machine is realising a hypothesis function to approximate the output:

$$\mathcal{F}: \mathbf{x} \mapsto \mathbf{y}$$

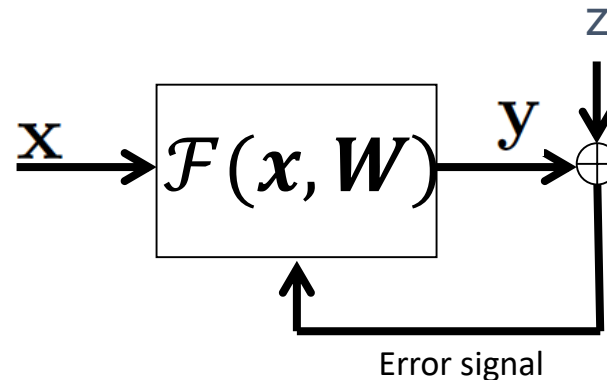
$$\mathbf{y} = \mathcal{F}(\mathbf{x}, \mathbf{W})$$

- $\mathbf{x}$  ..... D dimensional input
- $\mathbf{y}$  ..... p dimensional output
- $\mathbf{W}$  ..... parameters



# Machine training

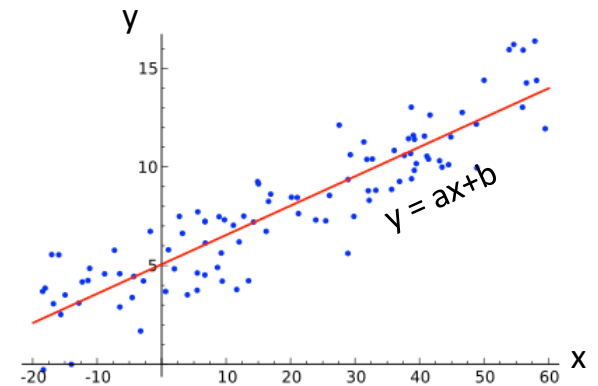
- Pre-requisites
  - Training set  $\mathbf{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$
  - Groundtruth target values  $\mathbf{Z} = \{\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)}\}$
  - Form of hypothesis function  $\mathcal{F}$
  - Cost function (objective, loss or error measure)
  - Optimisation (procedure to update  $W$ )



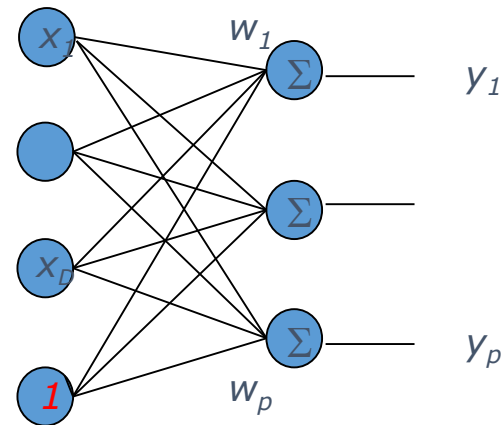
# Basic linear machine

- The simplest form of  $\mathcal{F}$  is a linear function

$$\mathbf{y} = \mathcal{F}(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x}$$



$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_p \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1D} & \textcolor{red}{b}_1 \\ w_{21} & w_{22} & \dots & w_{2D} & \textcolor{red}{b}_2 \\ & & \dots & & \\ & & \dots & & \\ & & \dots & & \\ w_{p1} & w_{p2} & \dots & w_{pD} & \textcolor{red}{b}_p \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \\ \textcolor{red}{1} \end{bmatrix} := \begin{bmatrix} w_1^T \\ w_2^T \\ \vdots \\ w_p^T \end{bmatrix} \mathbf{x}$$

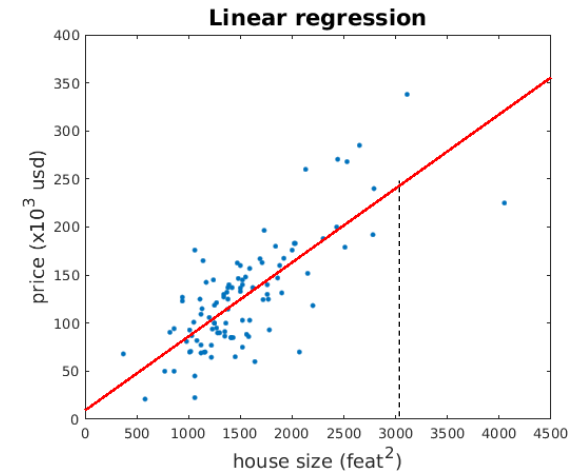
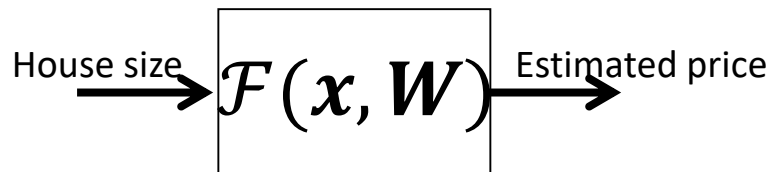


# Linear regression: example

- Pre-requisites

- Training set  $\mathbf{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$
- Groundtruth target values  $\mathbf{Z} = \{\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)}\}$
- Form of hypothesis function  $\mathcal{F}$
- Cost function (loss, objective or error measure)
- Optimisation (procedure to update  $\mathbf{W}$ )

$$\mathbf{y} = \mathcal{F}(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x}$$

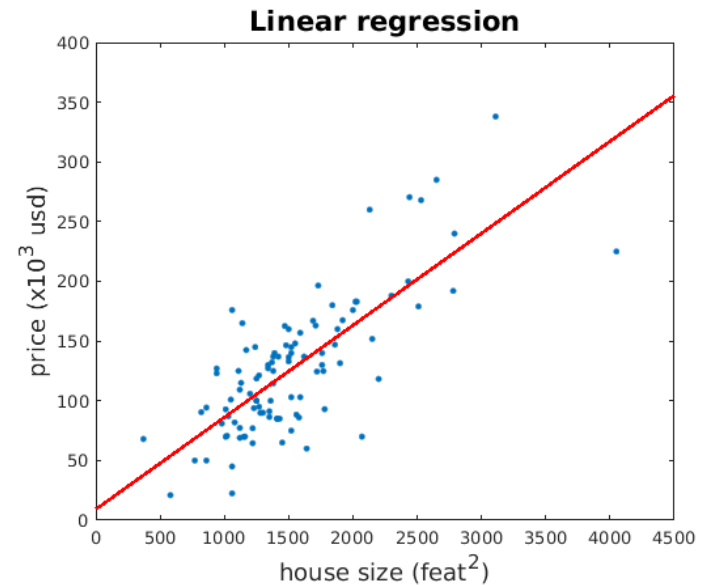


$\mathbf{X}$	House size (feet <sup>2</sup> )	Price (\$ x10 <sup>3</sup> )	$\mathbf{Z}$
$\mathbf{x}^{(1)}$	1240	145	$\mathbf{z}^{(1)}$
$\mathbf{x}^{(2)}$	370	68	$\mathbf{z}^{(2)}$
$\mathbf{x}^{(3)}$	1130	115	$\mathbf{z}^{(3)}$
$\mathbf{x}^{(4)}$	1120	69	$\mathbf{z}^{(4)}$
$\mathbf{x}^{(5)}$	1710	163	$\mathbf{z}^{(5)}$
...	...	...	...
$\mathbf{x}^{(N)}$	860	50	$\mathbf{z}^{(N)}$

# Cost/Loss function

$$\begin{aligned}\mathcal{L}(W) &= \frac{1}{2N} \sum_{i=1}^N (\mathbf{y}^{(i)} - \mathbf{z}^{(i)})^2 \\ &= \frac{1}{2N} \sum_{i=1}^N (\mathcal{F}(\mathbf{x}^{(i)}, W) - \mathbf{z}^{(i)})^2\end{aligned}$$

$\text{Min}_W \mathcal{L}(W)$

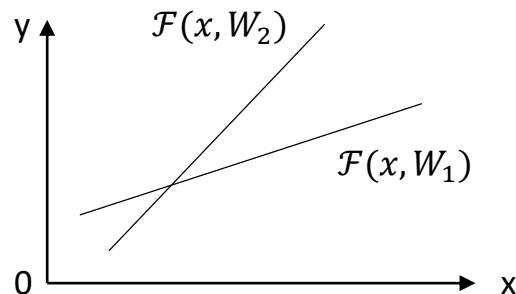


# Hypothesis function vs. cost function

## Hypothesis function

$$\begin{aligned}\mathcal{F}(\mathbf{x}, W) &= W\mathbf{x} = [a \ b] \begin{bmatrix} x \\ 1 \end{bmatrix} \\ &= ax + b\end{aligned}$$

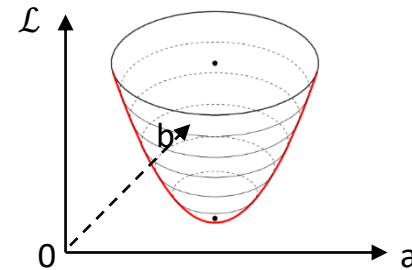
- Is a function of  $\mathbf{x}$ , given parameter  $W$



## Cost function

$$\begin{aligned}\mathcal{L}(W) &= \frac{1}{2N} \sum_{i=1}^N (\mathcal{F}(\mathbf{x}^{(i)}, W) - \mathbf{z}^{(i)})^2 \\ &= \frac{1}{2N} \sum_{i=1}^N (W\mathbf{x}^{(i)} - \mathbf{z}^{(i)})^2\end{aligned}$$

- Is a function of  $W$

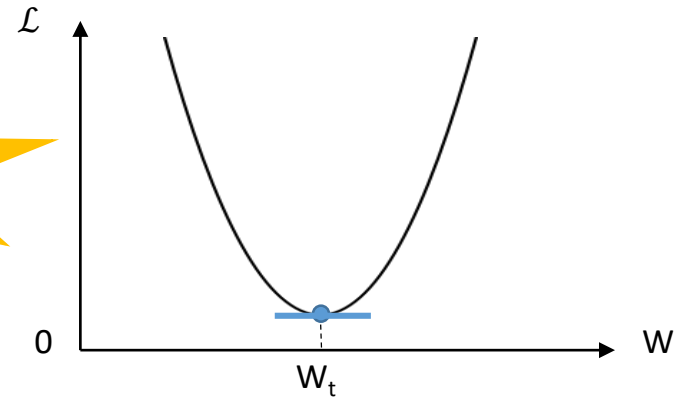


Todo: 2D map

# Optimisation method #1

$$\mathcal{F}(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x}$$

$$\mathcal{L}(\mathbf{W}) = \frac{1}{2N} \sum_{i=1}^N (\mathbf{W}\mathbf{x}^{(i)} - \mathbf{z}^{(i)})^2$$



Solve for  $\mathbf{W}$   $\boxed{\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = 0}$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \frac{1}{N} \sum_{i=1}^N (\mathbf{W}\mathbf{x}^{(i)} - \mathbf{z}^{(i)})\mathbf{x}^{(i)} = \frac{1}{N} (\mathbf{W}\mathbf{X} - \mathbf{Z})\mathbf{X}^T := 0$$

$$\boxed{\mathbf{W} = \mathbf{Z}\mathbf{X}^T(\mathbf{X}\mathbf{X}^T)^{-1}}$$



# Optimisation method #1

## Pros

- Straight forward, no loop.
- No additional parameters.

## Cons

- Need to compute:  
 $(\mathbf{X}\mathbf{X}^T)^{-1}$ 
  - slow when feature dimension is large
  - $\mathbf{X}\mathbf{X}^T$  may not invertible

⇒ Gradient descent

# Optimisation method #2: Gradient descent

Training procedure:

1. Initialisation: assign a random value to  $W$

$$W := W_0$$

2. Compute gradient of the loss function

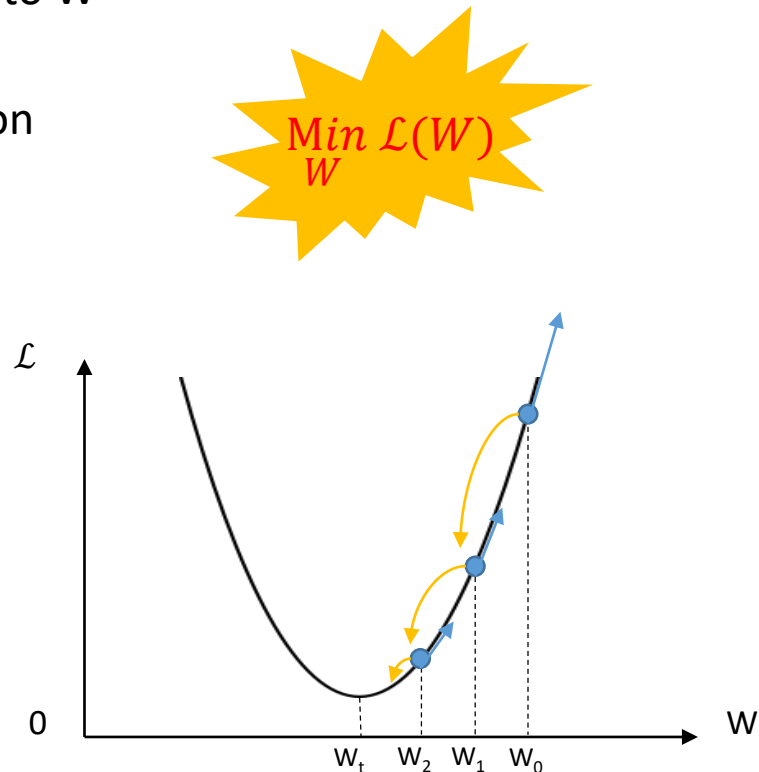
$$\frac{\partial \mathcal{L}}{\partial W}$$

3. Update  $W$

$$W := W - k \frac{\partial \mathcal{L}}{\partial W}$$

where  $k$  is the learning rate

4. Repeat step 2 & 3.



# Stopping criteria

When should we terminate the updating iterations?

- After fixed number of iterations:  $n = 1k, 1M$ ?
- When improvement in the loss is small enough?

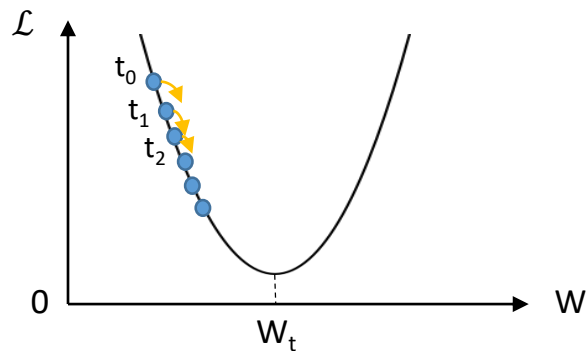
$$\Delta\mathcal{L}_t = |\mathcal{L}(W_t) - \mathcal{L}(W_{t-1})| < \epsilon$$

- Both?

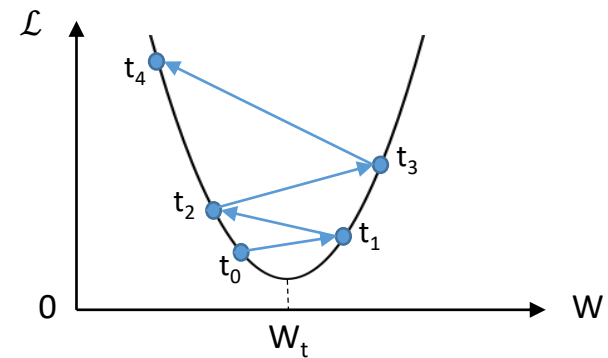
# Effect of learning rate

$$W := W - k \frac{\partial \mathcal{L}}{\partial W}$$

K too small: gradient descent can be slow



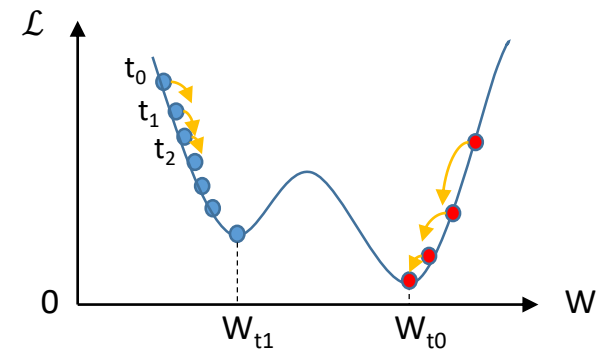
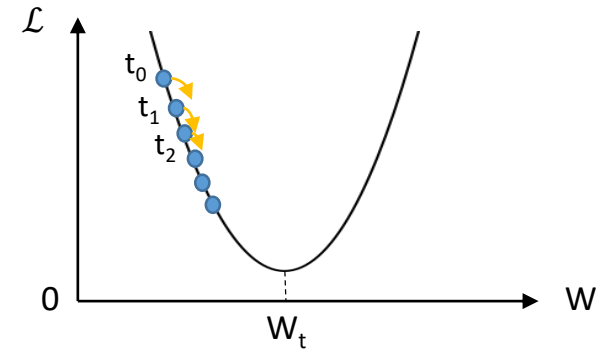
K too large: the training may not converge



**Quiz:** should we vary learning rate during the training of our model?

# Local optimal in gradient descent

- Linear regression:  $\mathcal{L}(W)$  is always convex
  - $\mathcal{L}(W) = \frac{1}{2N} \sum_{i=1}^N (W \mathbf{x}^{(i)} - \mathbf{z}^{(i)})^2$
- What if  $\mathcal{L}(W)$  non-convex?
  - Gradient descent may converge to local optimal.
- Possible solution:
  - Different parameter ( $W$ ) initialisations
  - Stochastic gradient descent



# Gradient descent types

- **Batch gradient descent:** each update iteration is derived from the whole training set

$$\mathcal{L}(W) = \frac{1}{2N} \sum (\mathcal{F}(x^{(1:N)}, W) - z^{(1:N)})^2 \quad W := W - k \frac{\partial \mathcal{L}}{\partial W}$$

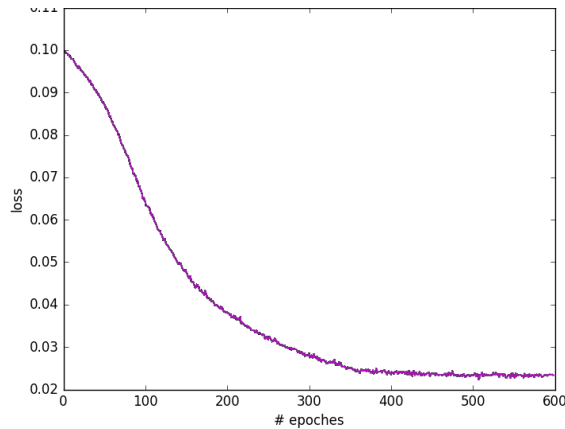
- **Stochastic gradient descent:** each update is derived from one sample in the training set

$$\mathcal{L}(W) = \frac{1}{2} (\mathcal{F}(x^{(i)}, W) - z^{(i)})^2 \quad W := W - k \frac{\partial \mathcal{L}}{\partial W}$$

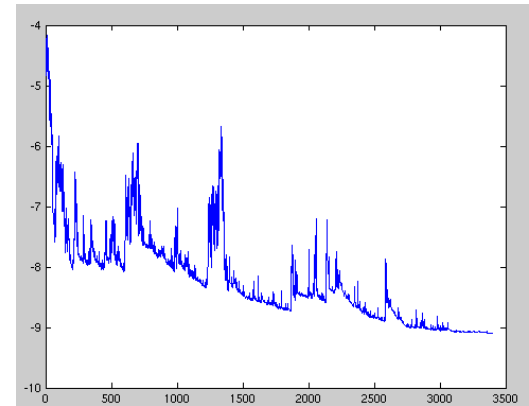
- **Mini-batch gradient descent:** perform update for every mini-batch of  $p$  training samples at a time

$$\mathcal{L}(W) = \frac{1}{2p} \sum (\mathcal{F}(x^{(i:i+p)}, W) - z^{(i:i+p)})^2 \quad W := W - k \frac{\partial \mathcal{L}}{\partial W}$$

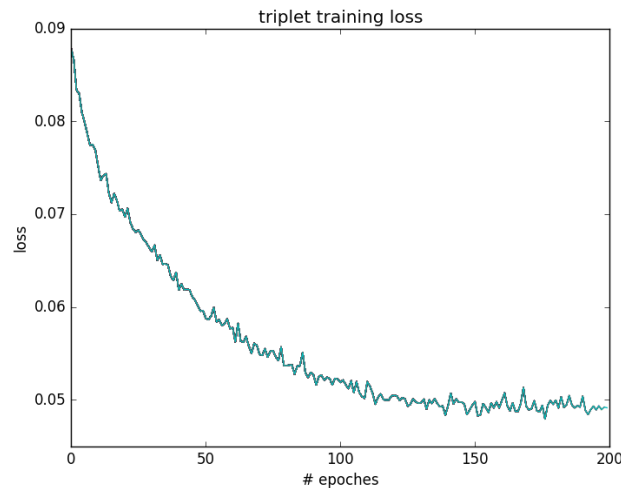
# Gradient descent types



Batch gradient descent



Stochastic gradient descent\*



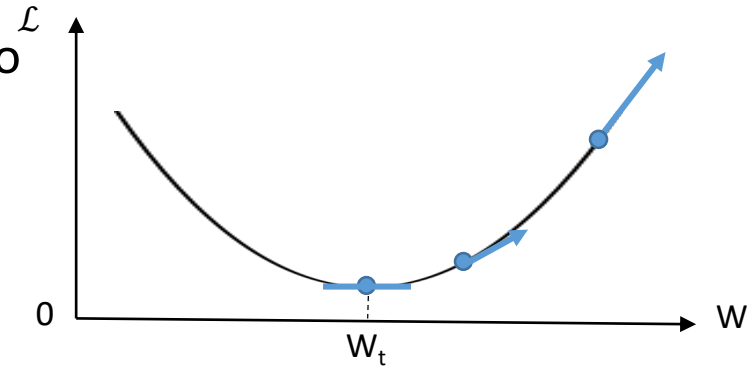
Mini-batch  
gradient descent

# Gradient descent with momentum

- Problem: update is getting slow when close to the optima (gradient becomes small).

$$W := W - k \frac{\partial \mathcal{L}}{\partial W}$$

especially when  $W$  is multi-dimensional.



- Solution: give “velocity” and “momentum” to the update

$$v_t = mv_{t-1} + k \frac{\partial \mathcal{L}}{\partial W}$$
$$W := W - v_t$$

where  $m$  is momentum,  $k$  is learning rate.  
In practice,  $m = 0.9$  (usually).



# Other optimisation methods

- Nesterov accelerated gradient

$$v_t = mv_{t-1} + k \frac{\partial}{\partial W} \mathcal{F}(W - mv_{t-1})$$

$$W := W - v_t$$

compute gradient at “future”  $W$ .

- Adagrad
  - Assign different learning rates to each component of  $W$ .
  - Perform larger updates to infrequent and smaller updates to frequent components.
- Adadelta, AdaMax, Adam, Nadam, RMSprop ... all using adaptive learning rate methods.

# Multivariate regression

House size (feet <sup>2</sup> )	#bedrooms	House age (years)	Distance from city centre (km)	Price (\$ x10 <sup>3</sup> )
1240	4	25	1.5	145
370	1	40	20.1	68
1130	3	5	13.0	115
1120	2	60	100.5	69
1710	4	13	30.7	163
...	...	...	...	...
860	2	8	46.4	50

$$y = ax + b$$

$$W = [a \ b] \quad x = \begin{bmatrix} x \\ 1 \end{bmatrix}$$

$$y = Wx$$



$$y = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + w_b$$

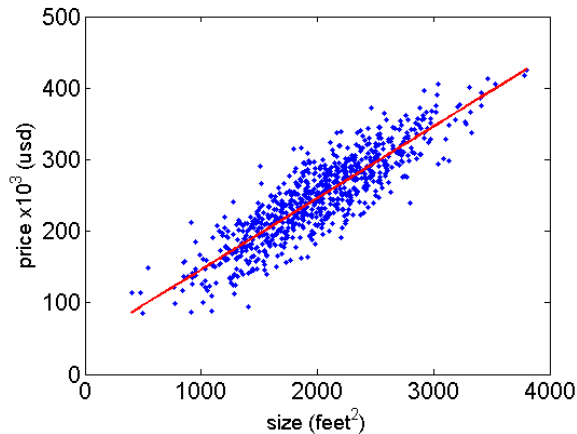
$$W = [w_1 \ w_2 \ w_3 \ w_4 \ w_b] \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ 1 \end{bmatrix}$$

$$y = Wx$$

Univariate regression

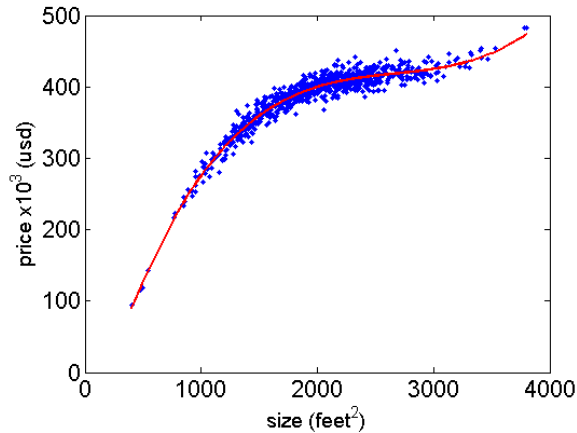
Multivariate regression

# Polynomial regression



$$y = ax + b \quad W = [a \ b], x = \begin{bmatrix} x \\ 1 \end{bmatrix}$$

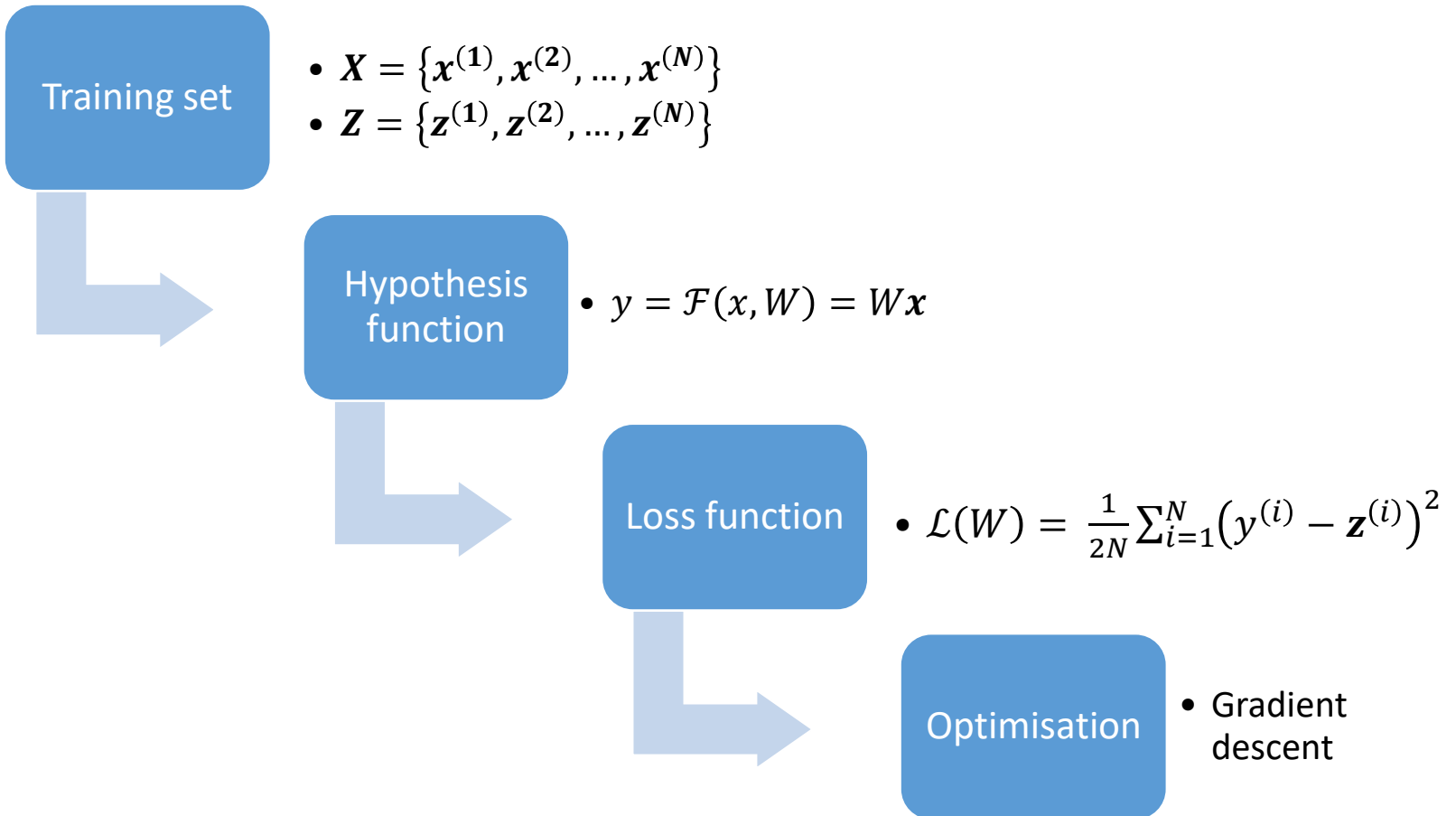
$$y = Wx$$



$$y = ax^3 + bx^2 + cx + d \quad W = [a \ b \ c \ d], x = \begin{bmatrix} x^3 \\ x^2 \\ x \\ 1 \end{bmatrix}$$

$$y = ax + b\sqrt{x} + c \quad W = [a \ b \ c], x = \begin{bmatrix} x \\ \sqrt{x} \\ 1 \end{bmatrix}$$

# Regression summary



# Classification

- Output: set of **discrete** values (label/class/category).
- Basic example: movie rating



John: ★ ★ ★ ★

Mary: ★ ★

Should I watch it or not?

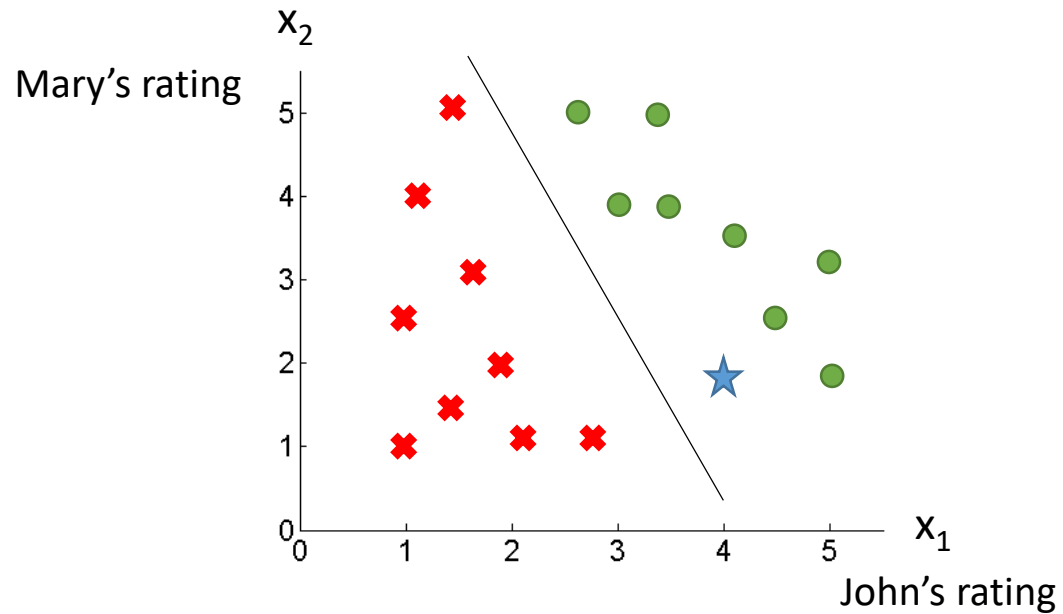
Spiderman homecoming

Binary  
Classification\*

	Rates (1-5★)		Do I like it?
	John	Mary	
Captain America: civil war	5	3	Yes
The hobbit III	4	3.5	Yes
Lalaland	1.5	5	No
....	...	...	...

\* or logistic regression

# Binary classification



# Problem formulation

- Training set

$$\mathbf{x} = \begin{bmatrix} \text{John's rating} \\ \text{Mary's rating} \\ 1 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}, \quad \mathbf{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$$

$$\mathbf{z} = \begin{cases} 1 & \text{if I like} \\ 0 & \text{if I don't like} \end{cases}, \quad \mathbf{Z} = \{\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)}\}$$

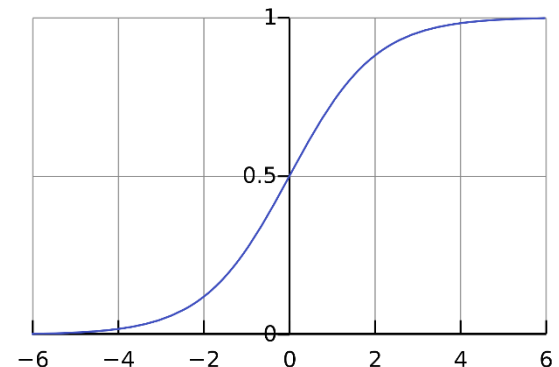
# Hypothesis function

- $t = w_1x_1 + w_2x_2 + w_b$
- $W = [w_1 \ w_2 \ w_b] \rightarrow t = Wx$
- Hypothesis function

$$\mathcal{F}(x, W) = t = Wx \quad \text{✗}$$

$$\mathcal{F}(x, W) = g(t) \quad \text{✓}$$

Sigmoid\*  $\rightarrow g(t) = \frac{1}{1 + e^{-t}}$



\* Or logistic function



# Decision making

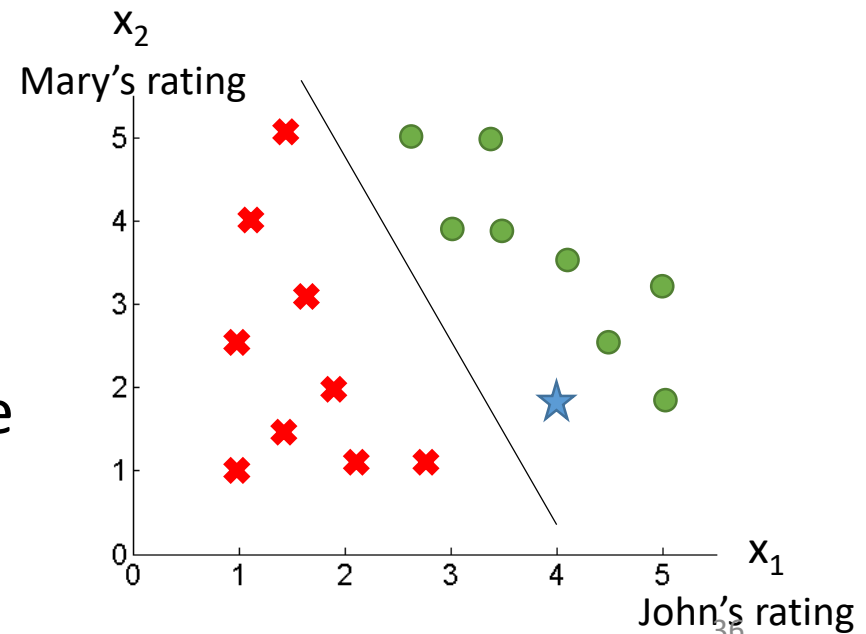
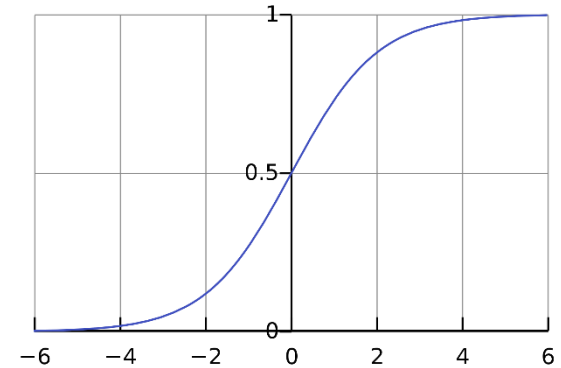
- $\mathcal{F}(x, W) = g(t) > 0.5$   
 $t = W\mathbf{x} > 0$

I like the movie

- $\mathcal{F}(x, W) < 0.5$   
 $t = W\mathbf{x} < 0$

I DON'T like the movie

- $\mathcal{F}(x, W) = 0.5$   
 $W\mathbf{x} = 0$   
 $w_1x_1 + w_2x_2 + w_b = 0$   
50-50 chance liking the movie

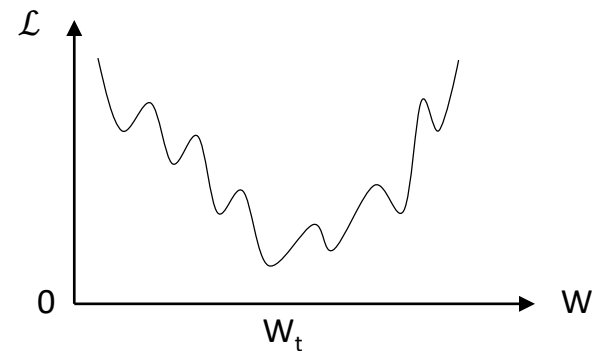


# Loss function

- If regression loss

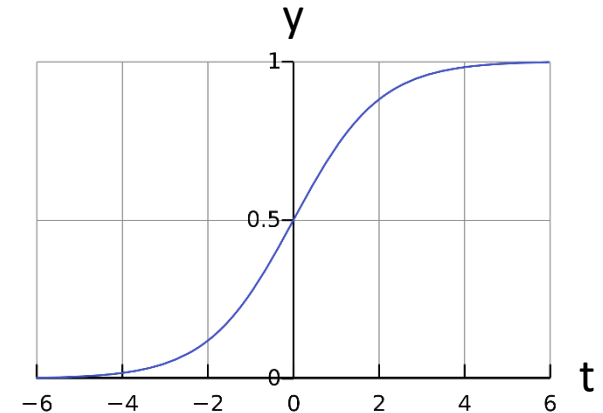
$$\begin{aligned}\mathcal{L}(W) &= \frac{1}{2N} \sum_{i=1}^N (\mathcal{F}(\mathbf{x}^{(i)}, W) - \mathbf{z}^{(i)})^2 \\ &= \frac{1}{2N} \sum_{i=1}^N \left( (1 + e^{-Wx})^{-1} - \mathbf{z}^{(i)} \right)^2\end{aligned}$$

→ non-convex



# Loss function

- $\mathcal{F}(\mathbf{x}, W) = g(t) = \frac{1}{1+e^{-Wx}}$   
represent probability at  $y = 1$
- $\mathcal{F}(\mathbf{x}, W) = P(y = 1|\mathbf{x}, W)$
- Probability  $y = 0$  is  $1 - \mathcal{F}(\mathbf{x}, W)$
- Goal: if groundtruth  $z=1$ , maximise  $\mathcal{F}(\mathbf{x}, W)$ ;  
if  $z=0$ , maximise  $1 - \mathcal{F}(\mathbf{x}, W)$
- $loss = \begin{cases} -\mathcal{F}(\mathbf{x}, W) & \text{if } z = 1 \\ -(1 - \mathcal{F}(\mathbf{x}, W)) & \text{if } z = 0 \end{cases}$

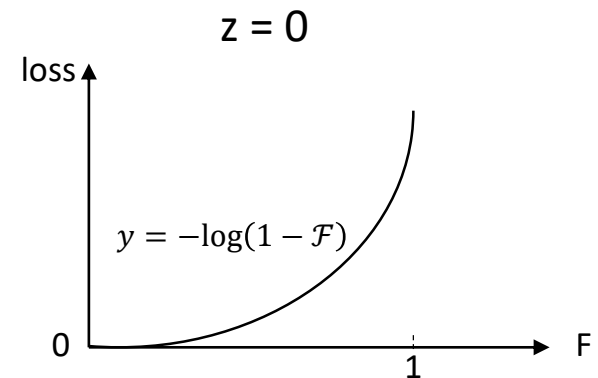
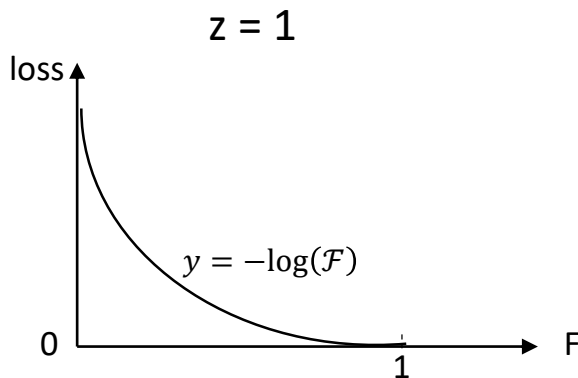


# Loss function

- $loss = \begin{cases} -\mathcal{F}(x, W) & \text{if } z = 1 \\ -(1 - \mathcal{F}(x, W)) & \text{if } z = 0 \end{cases}$  → non-convex

- $loss = \begin{cases} -\log(\mathcal{F}(x, W)) & \text{if } z = 1 \\ -\log(1 - \mathcal{F}(x, W)) & \text{if } z = 0 \end{cases}$  → convex

Negative  
log loss



# Loss function

$$\begin{aligned}\mathcal{L}(W) &= \begin{cases} -\log(\mathcal{F}(x, W)) & \text{if } z = 1 \\ -\log(1 - \mathcal{F}(x, W)) & \text{if } z = 0 \end{cases} \\ &= -z\log(\mathcal{F}(x, W)) - (1 - z)\log(1 - \mathcal{F}(x, W))\end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial \mathcal{F}} \frac{\partial \mathcal{F}}{\partial W}$$

- $\frac{\partial \mathcal{L}}{\partial \mathcal{F}} = \frac{-z}{\mathcal{F}} + \frac{1-z}{1-\mathcal{F}}$
- $\frac{\partial \mathcal{F}}{\partial W} = x \left( \frac{1}{1+e^{-Wx}} - \frac{1}{(1+e^{-Wx})^2} \right) = x\mathcal{F}(1 - \mathcal{F})$

$$\boxed{\frac{\partial \mathcal{L}}{\partial W} = x(\mathcal{F} - z)}$$

# Loss function: matrix form

- Per sample

- $\mathbf{y} = \mathcal{F}(\mathbf{x}, W) = \frac{1}{1+e^{-W\mathbf{x}}}$
- $\mathcal{L} = -\mathbf{z}\log(\mathcal{F}) - (1 - \mathbf{z})\log(1 - \mathcal{F})$
- $\frac{\partial \mathcal{L}}{\partial W} = \mathbf{x}(\mathcal{F} - \mathbf{z})$

- Over the whole training set

- $Y = \mathcal{F}(X, W) = \frac{1}{1+e^{-WX}}$
- $\mathcal{L} = \frac{-1}{N} [\log(\mathcal{F})\mathbf{Z}^T + \log(1 - \mathcal{F})(1 - \mathbf{Z})^T]$
- $\frac{\partial \mathcal{L}}{\partial W} = \frac{1}{N} (\mathcal{F} - \mathbf{Z})X^T \quad \rightarrow \text{convex}$

# Optimisation

- $\frac{\partial \mathcal{L}}{\partial W} = \frac{1}{N} (\mathcal{F} - \mathbf{Z}) \mathbf{X}^T$

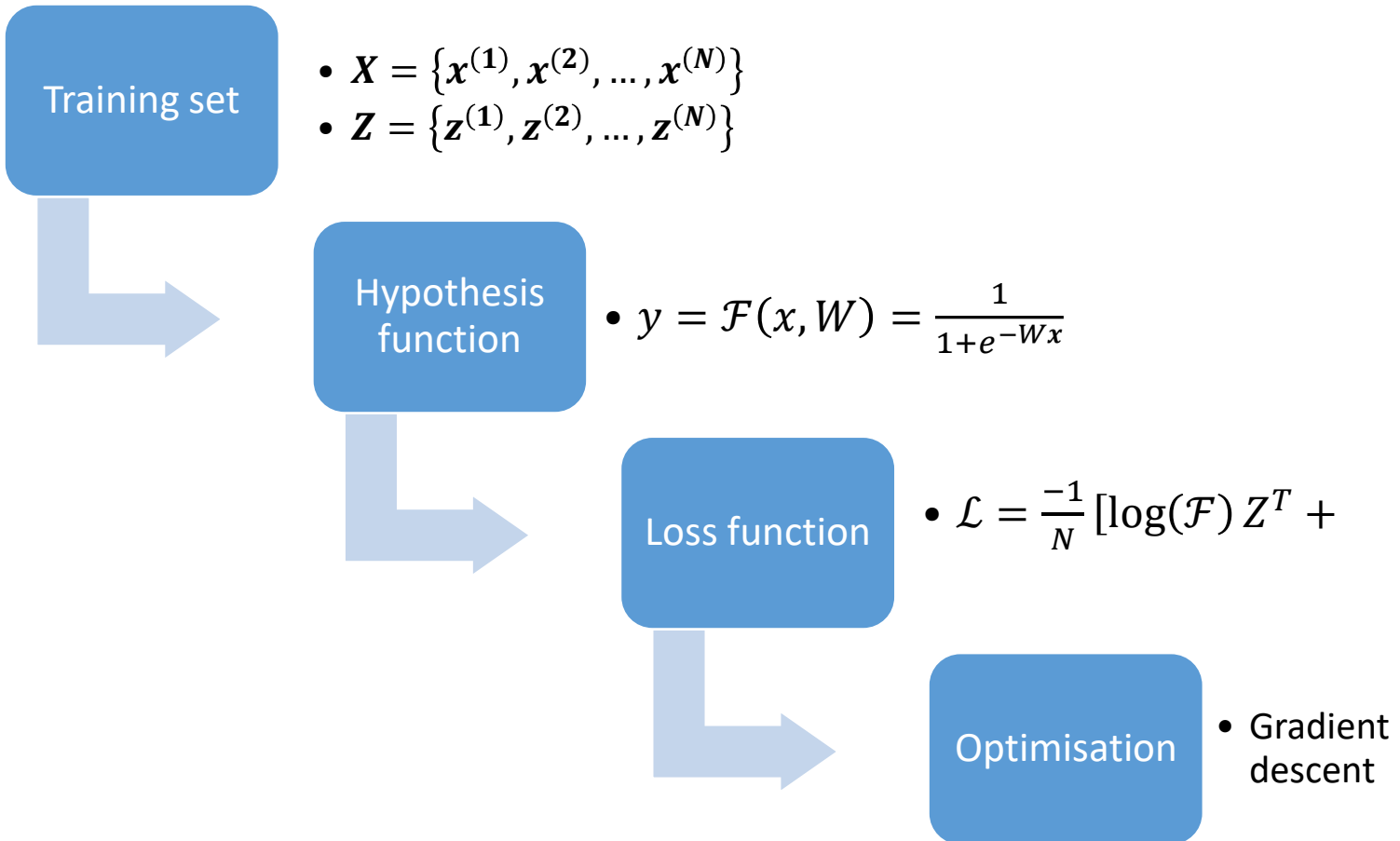
- Gradient descent

$$W := W - k \frac{\partial \mathcal{L}}{\partial W}$$

- Gradient descent with momentum

$$v_t = m v_{t-1} + k \frac{\partial \mathcal{L}}{\partial W}$$
$$W := W - v_t$$

# Binary classification summary

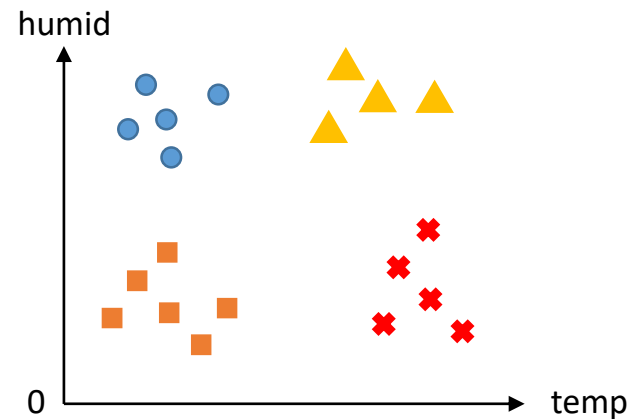




# Multi-class classification

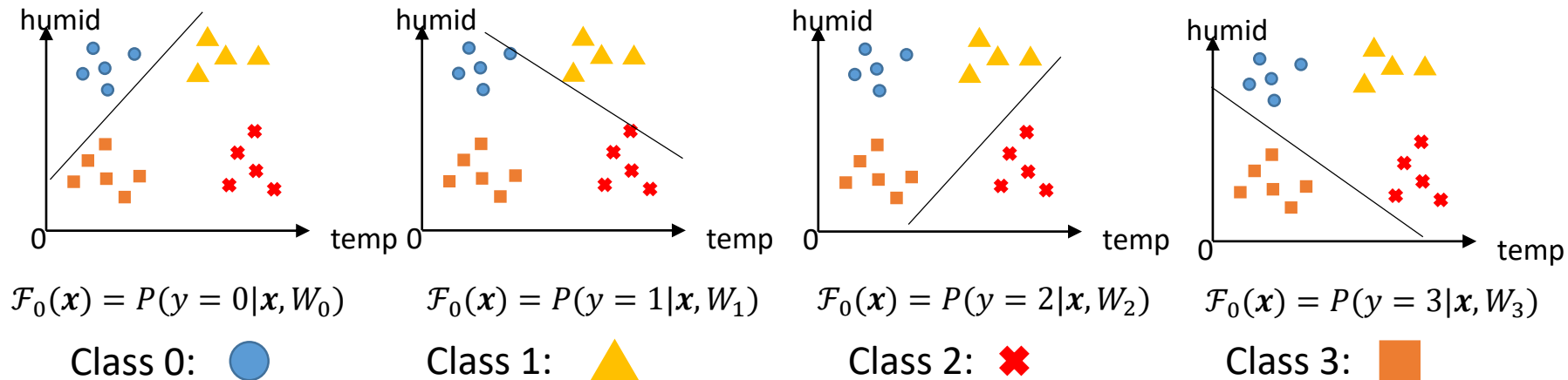
- The weather classification example:

- Input:  $x = \begin{bmatrix} \text{temperature} \\ \text{humidity} \\ 1 \end{bmatrix}$
- Output: sunny (0), cloudy (1), rain (2) or snow (3).
- $y \in \{0, 1, 2, 3\}$



# Method #1: one-vs-all

- For n-class classification, train n binary classifiers  $\mathcal{F}_i$ , each distinguish one class from the rest.

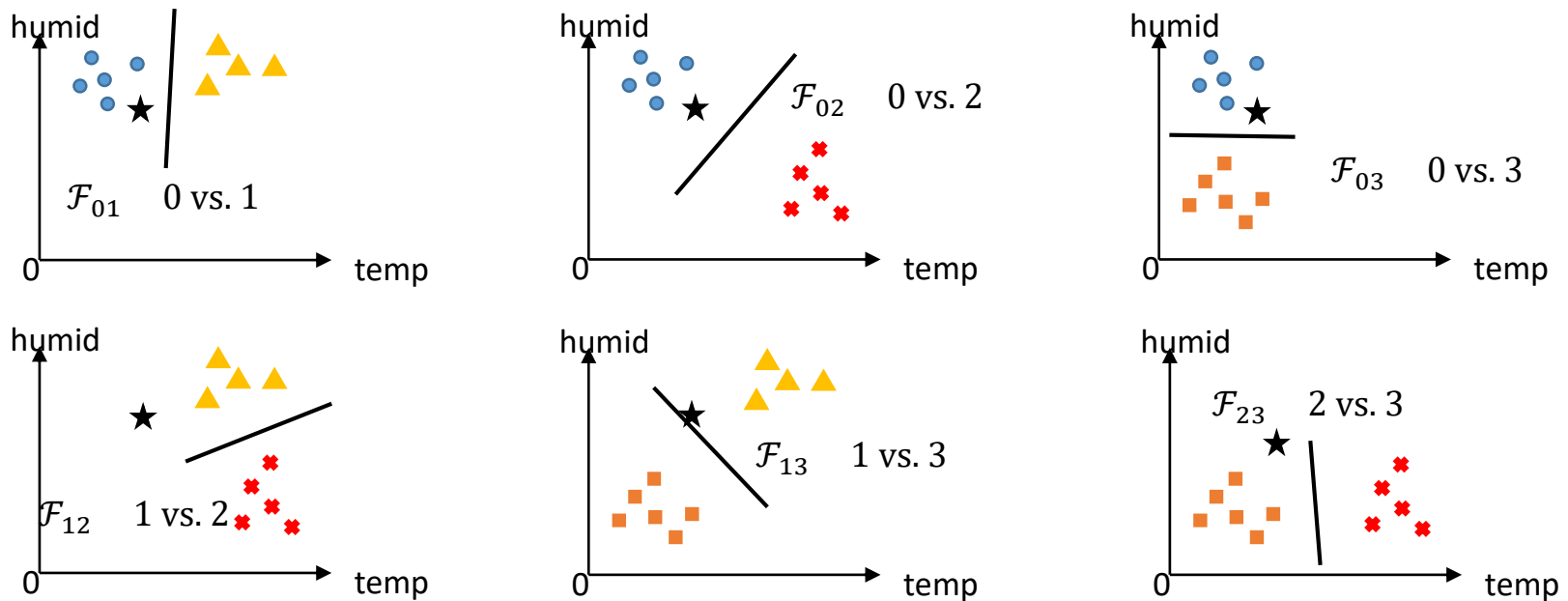


- Unknown sample  $\mathbf{x}$  is classified into class  $i$  :

$$\max_i \mathcal{F}_i(\mathbf{x}, W_i)$$

# Method #2: one-vs-one

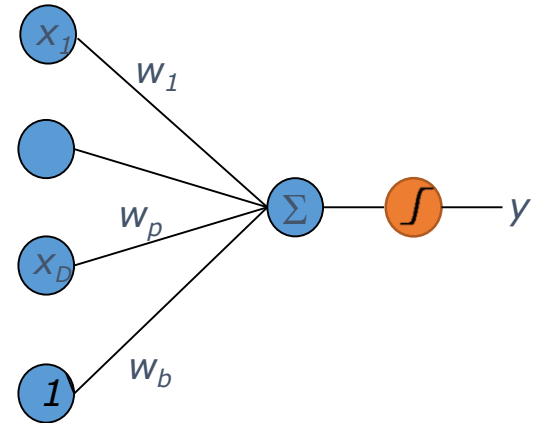
- For  $n$ -class classification, train  $n(n - 1)/2$  binary classifiers  $\mathcal{F}_{ij}$ , each distinguishes class  $i$  and class  $j$ .



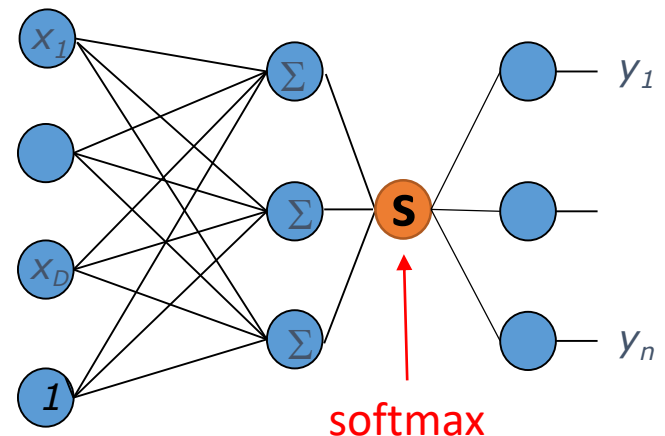
- Unknown sample  $\mathbf{x}$  belongs to class  $i$  if  $i$  gets the highest #predictions in  $n(n - 1)/2$  classifiers.

# Method #3: softmax

- Binary classification
  - $x \in \mathbb{R}^{D+1}; W \in \mathbb{R}^{D+1}; y, z \in \mathbb{R}$
  - $y = \mathcal{F}(x, W) = g(W\mathbf{x})$
- Multi-class classification
  - $x \in \mathbb{R}^{D+1}; W \in \mathbb{R}^{n \times (D+1)}; y, z \in \mathbb{R}^n$
  - n number of classes
  - $y = \mathcal{F}(x, W) = h(W\mathbf{x})$



$$z = \begin{bmatrix} \text{sunny} \\ \text{cloudy} \\ \text{rain} \\ \text{snow} \end{bmatrix} \quad \text{e.g. rain} \quad z = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$



# Classification with softmax

- Softmax function

$$t = \begin{bmatrix} t_1 \\ t_2 \\ \dots \\ t_n \end{bmatrix} \rightarrow \boxed{\mathcal{h}(\mathbf{t})} \rightarrow y = \begin{bmatrix} \frac{e^{t_1}}{\sum_i e^{t_i}} \\ \frac{e^{t_2}}{\sum_i e^{t_i}} \\ \dots \\ \frac{e^{t_n}}{\sum_i e^{t_i}} \end{bmatrix} = \begin{bmatrix} P(y = 1|x, W) \\ P(y = 2|x, W) \\ \dots \\ P(y = n|x, W) \end{bmatrix}$$

- Loss function: negative log loss for class  $c$

$$\mathcal{L}_c = -\log(P(y = c|x, W))$$

# Machine learning in practice

- Preprocessing data: feature normalisation

House size (feet <sup>2</sup> )	#bedrooms	House age (years)	Distance from city centre (km)
1240	4	25	1.5
370	1	40	20.1
1130	3	5	13.0
1120	2	60	100.5
1710	4	13	30.7
...	...	...	...
860	2	8	46.4

$$\mathbf{x} = \begin{bmatrix} \text{house size} \\ \text{bedrooms} \\ \text{age} \\ \text{distance} \\ 1 \end{bmatrix}$$

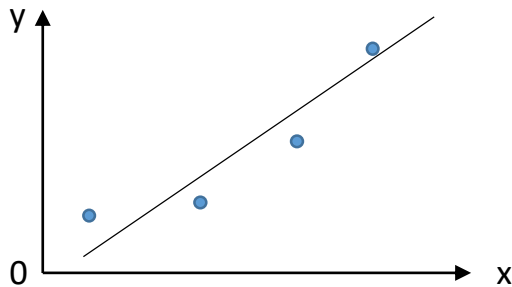
$$x_1 \in 200\text{-}5000 \quad x_2 \in 1\text{-}5 \quad x_3 \in 1\text{-}100 \quad x_4 \in 1\text{-}200$$

Normalise each feature

$$x := \frac{x - \bar{x}}{\max x - \min x} \in [-0.5, 0.5]$$

# Machine learning in practice

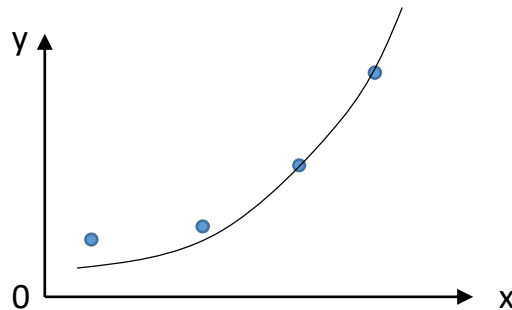
- The overfitting problem
  - Work well on training data, poor performance on test data.
  - Incorrect selection of the hypothesis function can be a source of the problem
  - Often occurs when there are many features, but not enough training samples -> more parameters to learn.



$$y = ax + b$$

$$W = [a \ b]$$

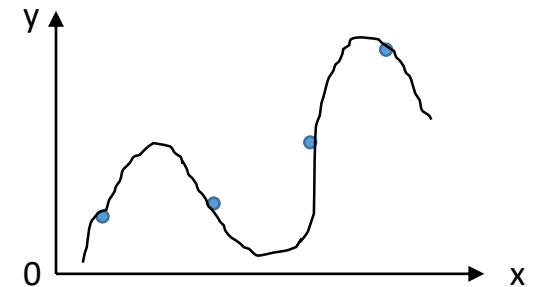
Under-fitting



$$y = ax^2 + bx + c$$

$$W = [a \ b \ c]$$

good



$$y = a_5x^5 + a_4x^4 + \dots + a_0$$

$$W = [a_5 \ a_4 \ \dots \ a_0]$$

over-fitting

# Machine learning in practice

- Combat over-fitting using regularisation
  - Keep all features but reduce the magnitude of  $W$ .
  - Implemented by adding  $|W|$  into the loss function

regression

$$\begin{aligned}\mathcal{L}(W) &= \frac{1}{2N} \sum_{i=1}^N (\mathcal{F}(\mathbf{x}^{(i)}, W) - \mathbf{z}^{(i)})^2 + \frac{1}{2} |W|^2 \\ &= \frac{1}{2N} |WX - Z|^2 + \frac{1}{2} |W|^2\end{aligned}$$

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial W} &= \frac{1}{N} \sum \mathbf{x}(\mathcal{F} - \mathbf{z}) + W \\ &= \frac{1}{N} (WX - Z)X^T + W\end{aligned}$$

classification

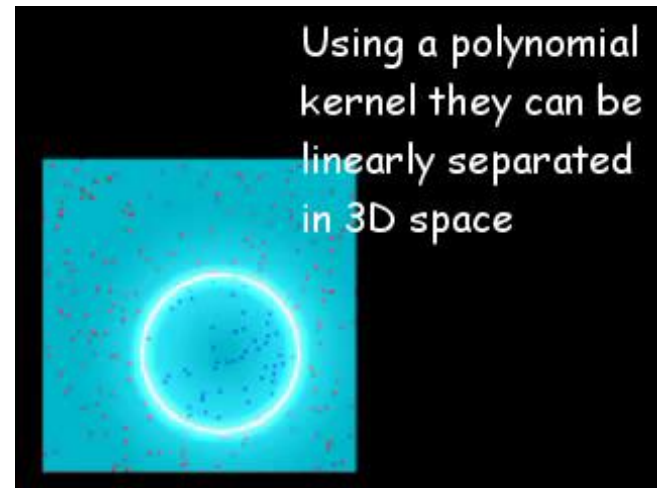
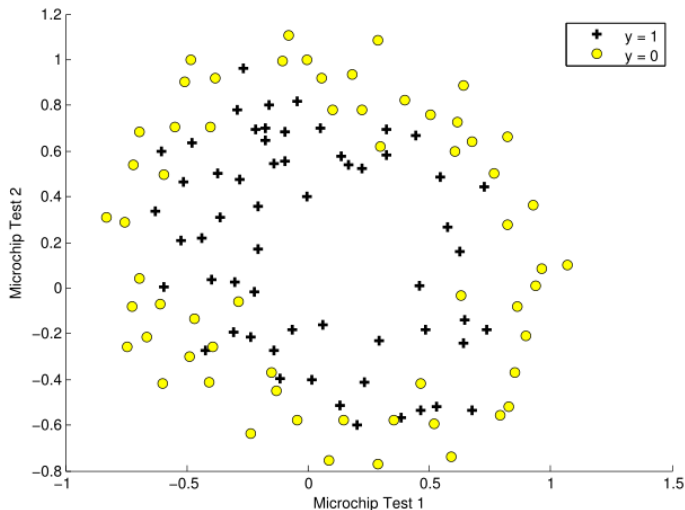
$$\begin{aligned}\mathcal{L}(W) &= \frac{-1}{N} \sum [\mathbf{z} \log(\mathcal{F}) + (1 - \mathbf{z}) \log(1 - \mathcal{F})] + \frac{1}{2} |W|^2 \\ &= \frac{-1}{N} [\log(\mathcal{F}) \mathbf{Z}^T + \log(1 - \mathcal{F})(1 - \mathbf{Z})^T] + \frac{1}{2} |W|^2\end{aligned}$$

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial W} &= \frac{1}{N} \sum \mathbf{x}(\mathcal{F} - \mathbf{z}) + W \\ &= \frac{1}{N} (WX - Z)X^T + W\end{aligned}$$



# Advanced ML techniques

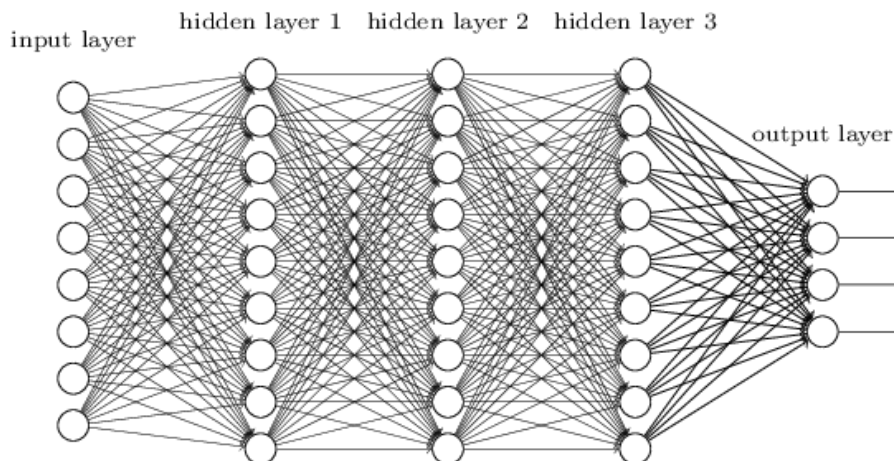
- Support vector machine (SVM)



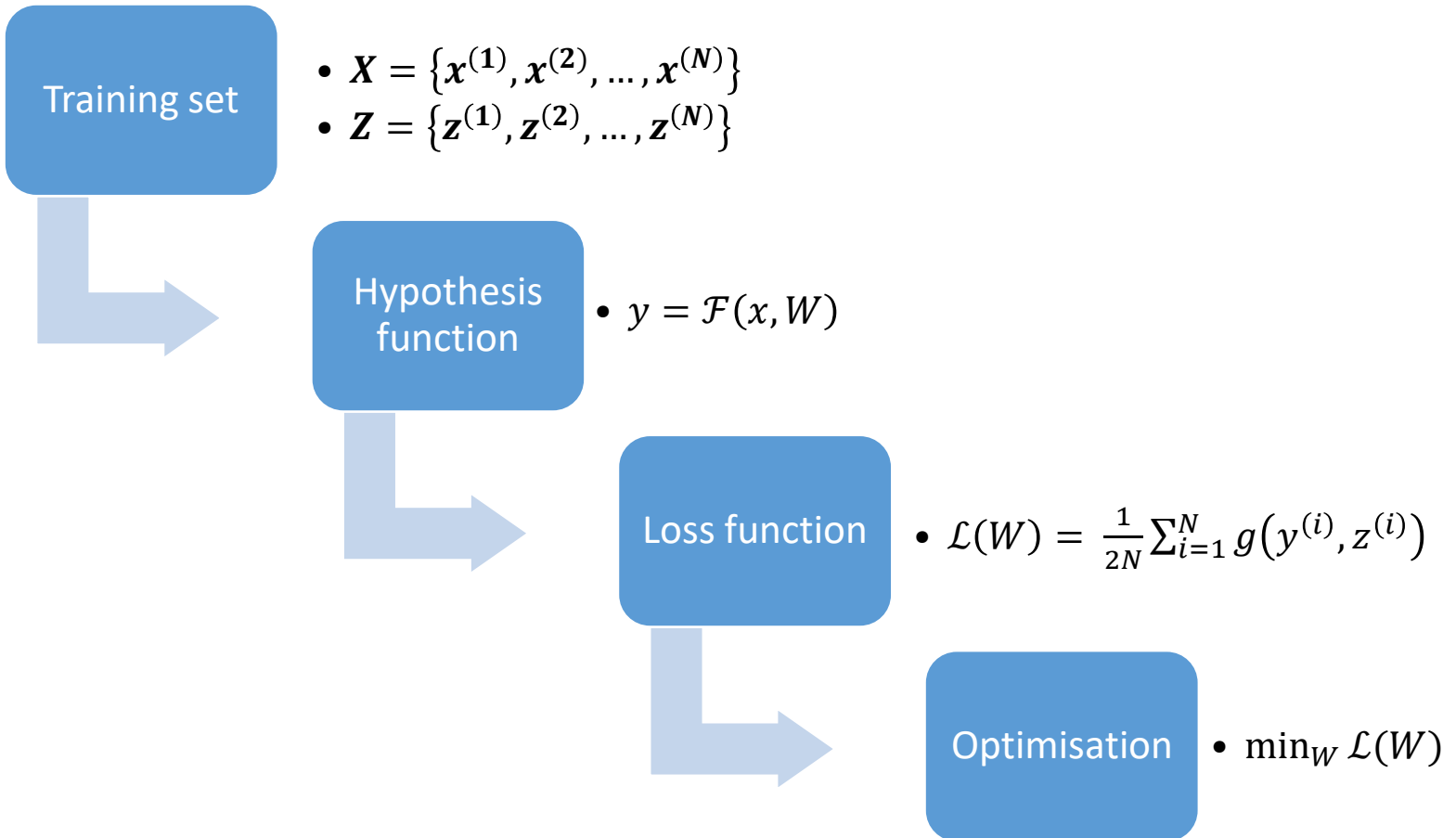
Courtesy: youtube

# Advanced ML techniques

- Deep neural networks (DNN)
  - Multi-layer neurons, each consists of a linear operation (convolution) followed by a non-linear one (activation).
  - Huge number of parameters (millions-D).
  - Compute  $\partial \mathcal{L} / \partial W$  based on chain rule and backpropagation
$$\frac{\partial}{\partial W} f \left( g \left( h \left( \dots k(W) \right) \right) \right) = \frac{\partial f}{\partial g} \frac{\partial g}{\partial h} \dots \frac{\partial k}{\partial W}$$
  - Currently state-of-the-art in many applications



# Conclusion: supervised learning procedure



# Acknowledgement

This presentations use materials from the following sources:

- Machine Learning Intro – UDRC'17 by Josef Kittler – University of Surrey
- Deep Learning Tutorials – UDRC'17 by Muhammad Awais – University of Surrey
- Machine Learning – coursera.org by Andrew Ng – Stanford University
- Neural Network reviews – MLSS'14 by Quoc Le – Google Brain

Thank you.

