



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 3, 14/09/2020

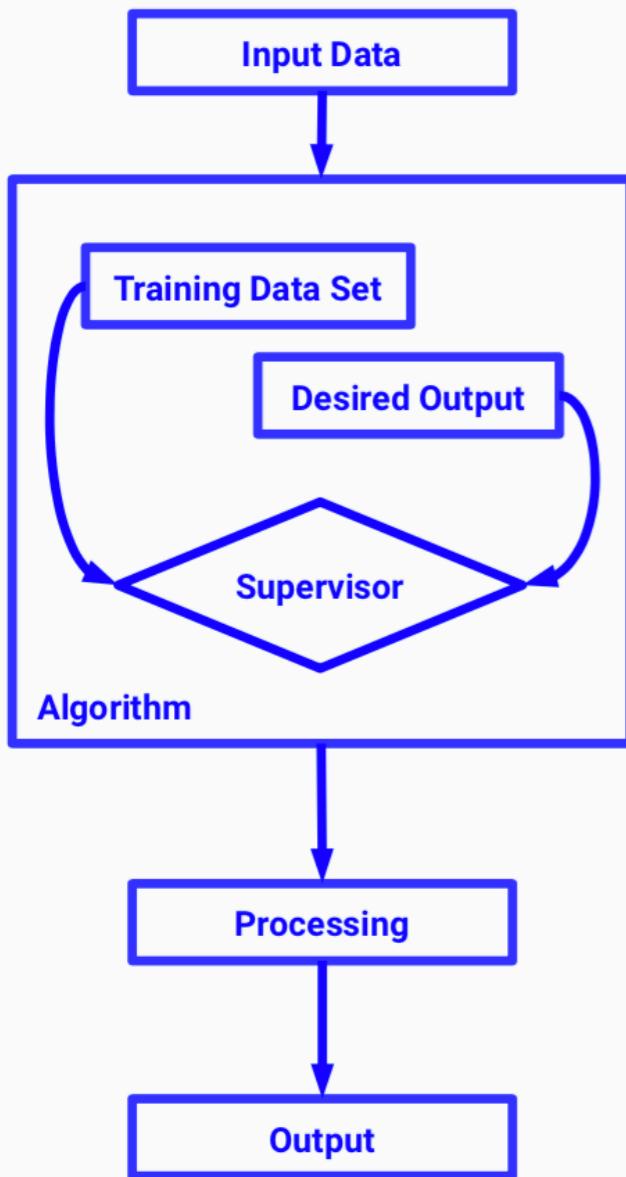
Today's lecture

- ➊ Supervised Learning for classification
- ➋ Linear classifiers and least squares
- ➌ Decision Theory
- ➍ Logistic regression and neural nets for classification
- ➎ Decision Trees and Random Forests
- ➏ Support Vector Machines

Classification with Linear Models

Supervised Learning recap

Supervised learning



Unsupervised learning



Reinforcement learning



this lecture

lecture 4

lecture 5

Supervised Learning recap

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

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

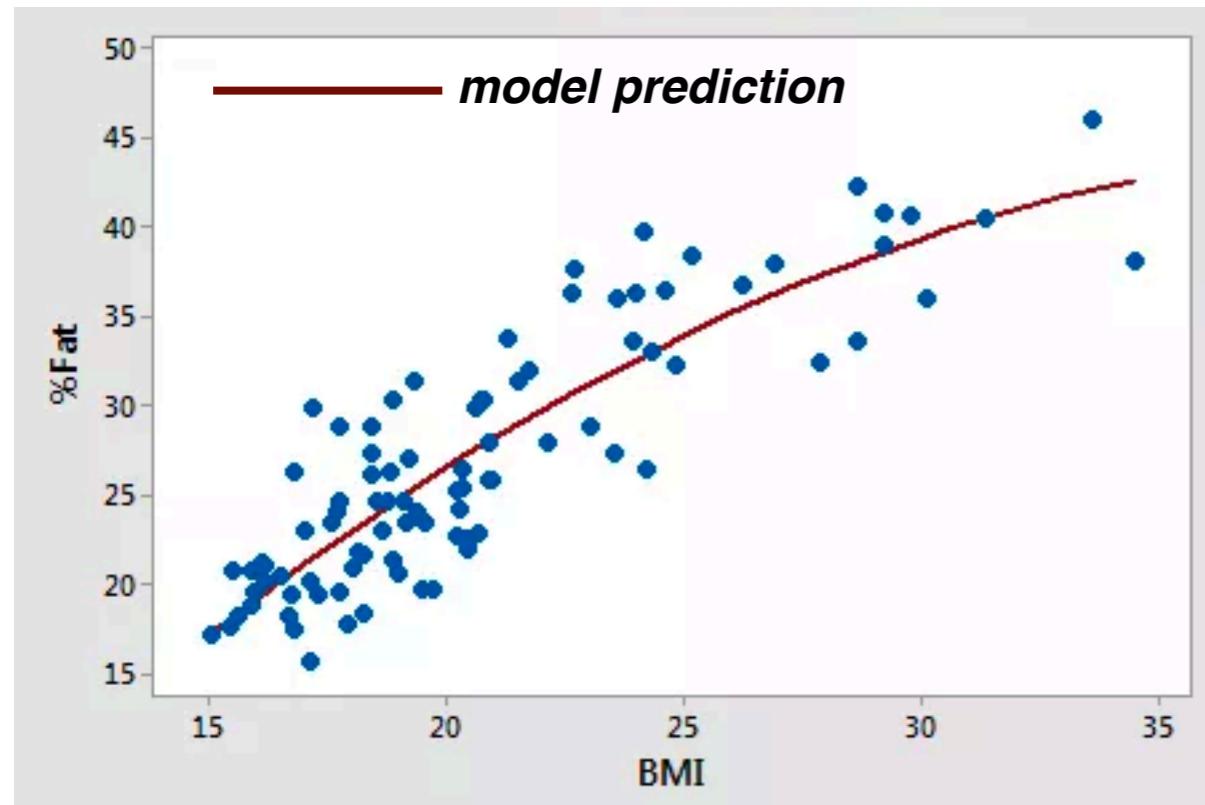
Supervised learning

We denote as **supervised learning** the ML task of **learning a function** that maps a **vector of inputs** to a **vector of outputs** from a finite set of training example

note that some assumptions will be needed: a function is an **infinite-dimensional object** but learning takes place from a **finite number of examples**

main property of supervised learning: the **training samples are labeled**

continuous outputs:
regression



discrete outputs:
classification



Classification tasks

The goal is to **predict a class label** from a pre-defined list of possibilities

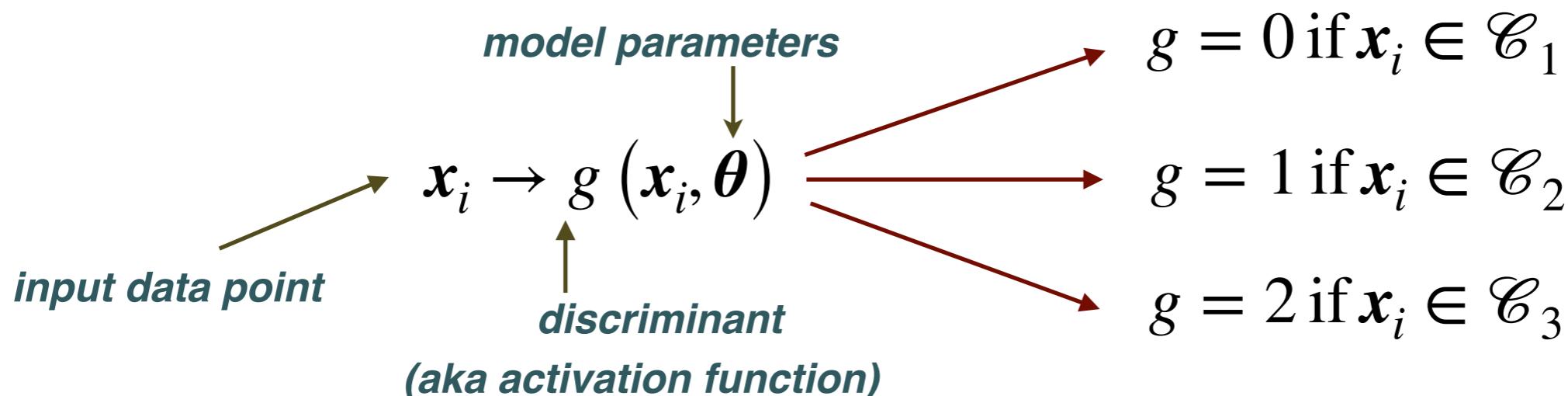
the simplest type of problem is **binary classification** (yes/no problems)

e.g. should I put this email in the spam folder?

but in general one considers **multiclass classification** (> 2 categories)

e.g. which type of bird is the one I just photographed?

In the context of ML applications there exist a large number of approaches to classifications tasks. The most basic one is based on assembling a **discriminant function** that maps each input data point to its specific class



Classification tasks

relevant for Machine Learning applications where outcomes are discrete variables, eg. **categories in classification problems**

“noise”



$$y_i = 0$$

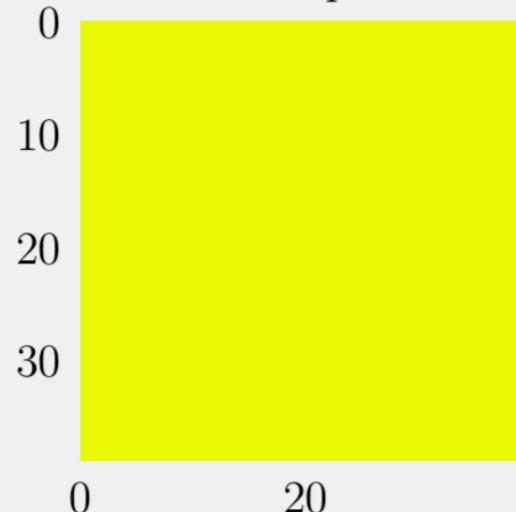
“signal”



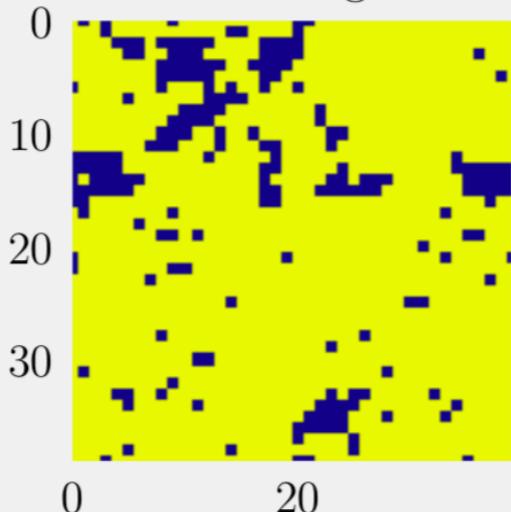
$$y_i = 1$$

*can we tell apart
cats from dogs?*

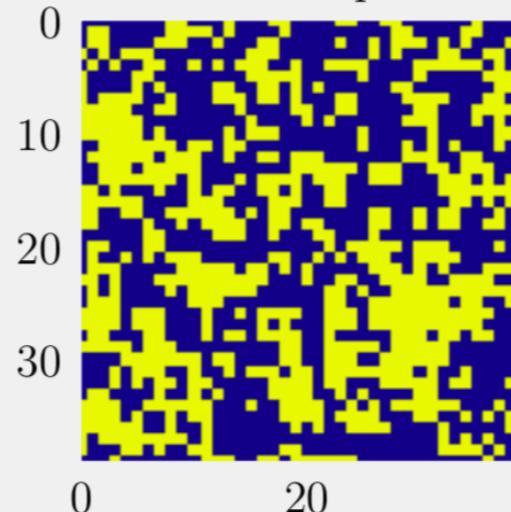
ordered phase



critical region



disordered phase



*can we identify the phase
(ordered/disordered) of
spin configurations in 2D Ising?*

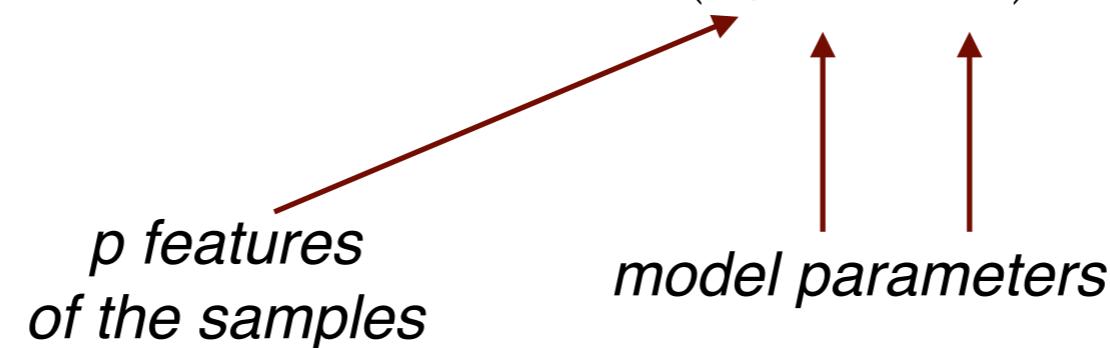
Linear classifiers

In this class of problems, the **dependent variables** y_i are discrete and take values $m=0, \dots, M-1$, so the index m also labels the M categories

The goal is to **classify the n input samples**, each composed by **p features**, into the **M possible categories** of the problem

A simple classifier is the **perceptron**: a **linear classifier** that categorises examples from a linear combination of the features

$$\sigma(s_i) = \text{sign}(s_i) = \text{sign}(\mathbf{x}_i^T \boldsymbol{\theta} + b_0)$$



a perceptron is a **hard classifier** where each sample is assigned to a category with 100% probability

more complex models can be used as classifiers, including NNs

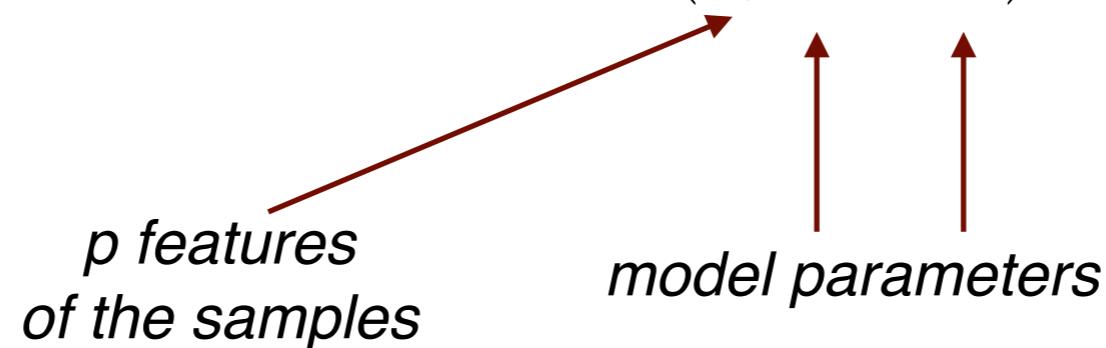
Linear classifiers

In this class of problems, the **dependent variables** y_i are discrete and take values $m=0, \dots, M-1$, so the index m also labels the M categories

The goal is to **classify the n input samples**, each composed by **p features**, into the **M possible categories** of the problem

A simple classifier is the **perceptron**: a **linear classifier** that categorises examples from a linear combination of the features

$$\sigma(s_i) = \text{sign}(s_i) = \text{sign}(\mathbf{x}_i^T \boldsymbol{\theta} + b_0)$$



a perceptron is a **hard classifier** where each sample is assigned to a category with 100% probability

more complex models can be used as classifiers, including NNs

simple linear models are great to gain physical intuition about how ML algorithms work

Linear classifiers

the perceptron (linear classifier) only has **two possible outcomes**

$$x_i \rightarrow g(x_i, \theta)$$
$$\begin{aligned} g &= 1 \text{ if } x_i \in \mathcal{C}_1, (x_i^T \theta + b_0) \geq 0 \\ g &= -1 \text{ if } x_i \in \mathcal{C}_1, (x_i^T \theta + b_0) < 0 \end{aligned}$$

the **decision boundary** of the classifier is then defined by the condition

$$(x_i^T \theta + b_0) = 0$$

the properties of the decision boundary of this problem can be easily evaluated

For two points that belong to this boundary one has

$$\begin{aligned} (x_1^T \theta + b_0) &= 0 \\ (x_2^T \theta + b_0) &= 0 \end{aligned} \quad \begin{array}{l} \xrightarrow{\hspace{1cm}} \\ \xrightarrow{\hspace{1cm}} \end{array} \quad (x_1^T - x_2^T) \theta$$

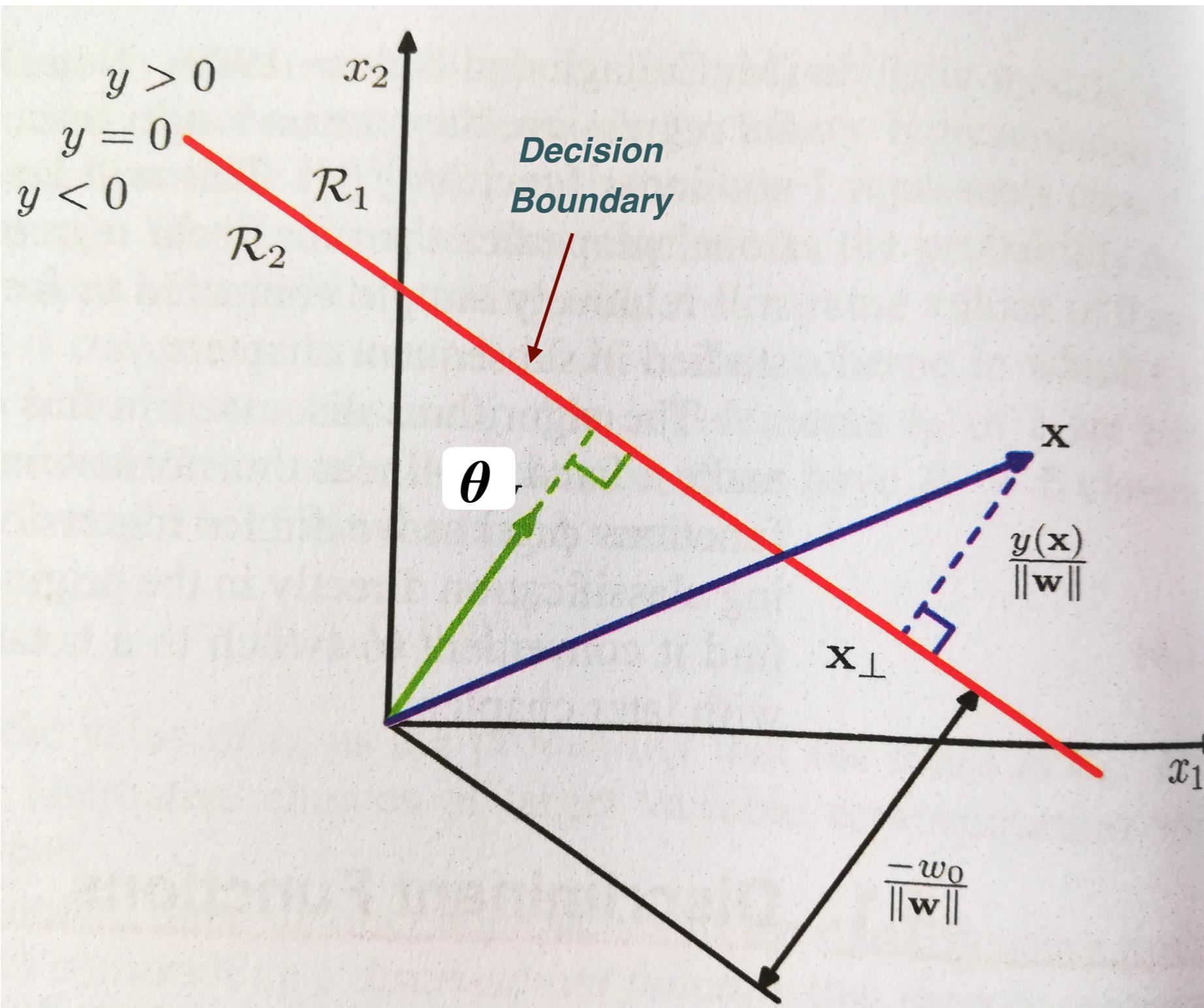
*the model parameter vector
is perpendicular to the
decision boundary!*

several other geometric properties can be evaluated

recall that both x and θ live in the p -dimensional feature space

Linear classifiers

example of linear classification for 2d datasets



Training the model

as in the case of regression problems, also for classifications we need to introduce a **figure of merit** (error function) to be optimised

the simplest option is based again on **least squares**. Let's consider a binary linear classifier

*the sign function is added
back after the training*

$$g(x_i) = (x_i^T \theta + b_0)$$

recall that x_i and θ have p features each

and one is given a **training dataset** with N instances

$$\{x_i^T, t_i\}, \quad i = 1, \dots, N, \quad t_i \in (-1, 1)$$

in this case the error function depending on the **model parameters** is

$$E_{\text{tr}}(\theta, b_0) = \frac{1}{N} \sum_{i=1}^N \left((x_i^T \theta + b_0) - t_i \right)^2$$

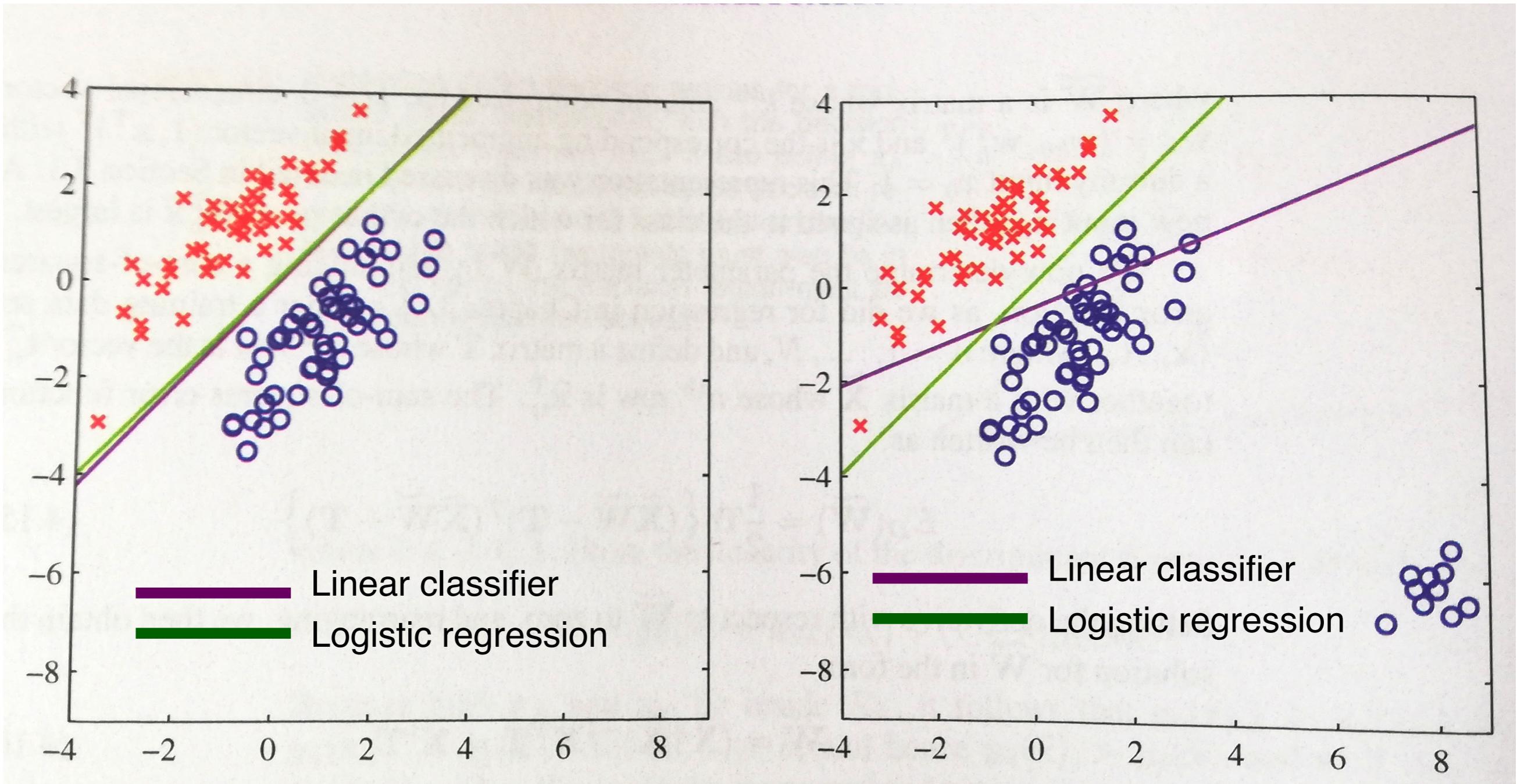
Which is amenable to analytic optimisation, for example

*linear system of equations
for the model params*

$$\frac{\partial}{\partial b_0} E_{\text{tr}} = \frac{1}{N} \sum_{i=1}^N \left((x_i^T \theta + b_0) - t_i \right) = 0 \rightarrow b_0 = \frac{1}{N} \sum_{i=1}^N (t_i - x_i^T \theta) = 0$$

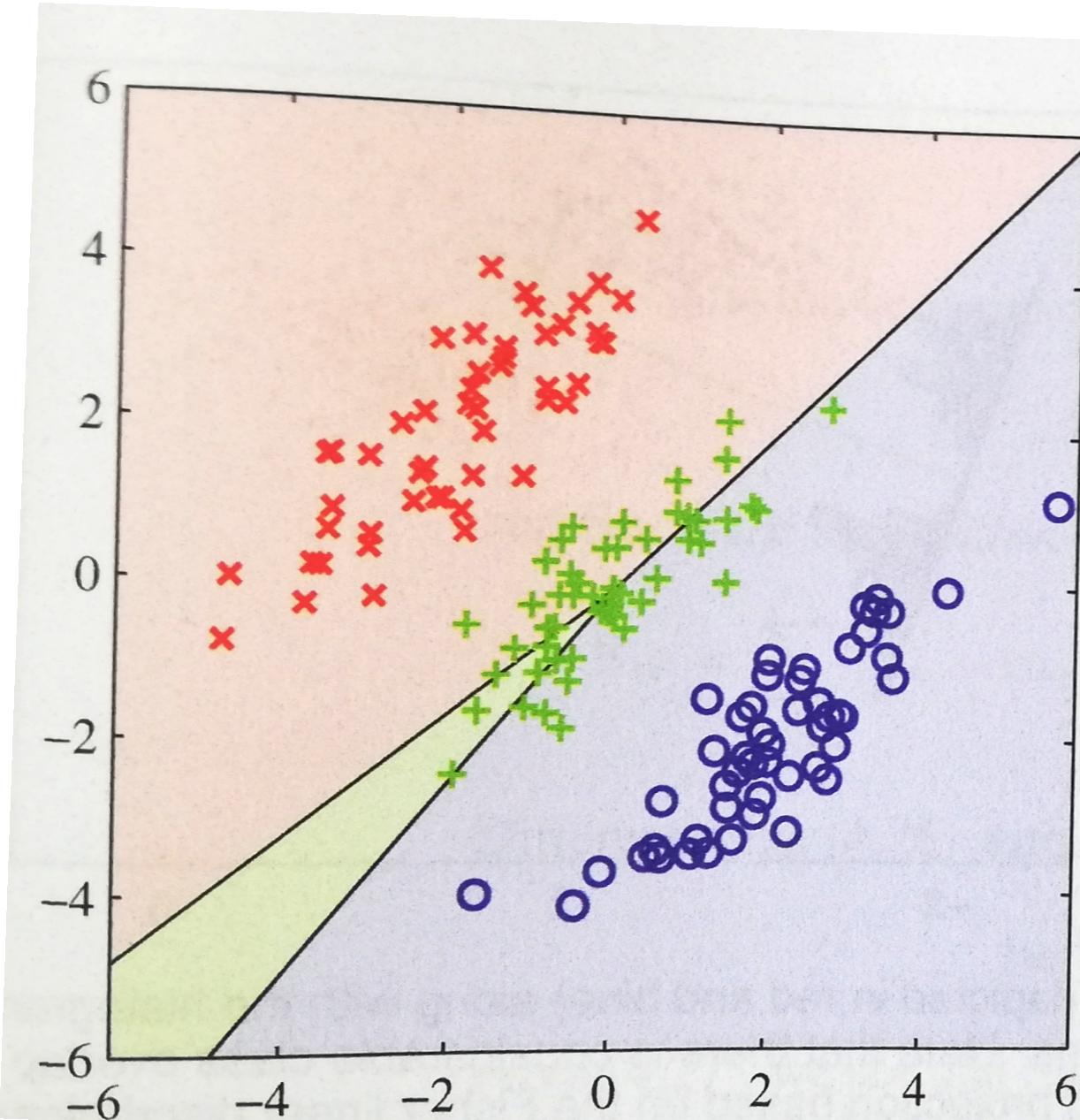
Limitations of the linear classifier

- As in regression problems, a linear classifier trained with least squares lacks robustness with respect to the presence of outliers

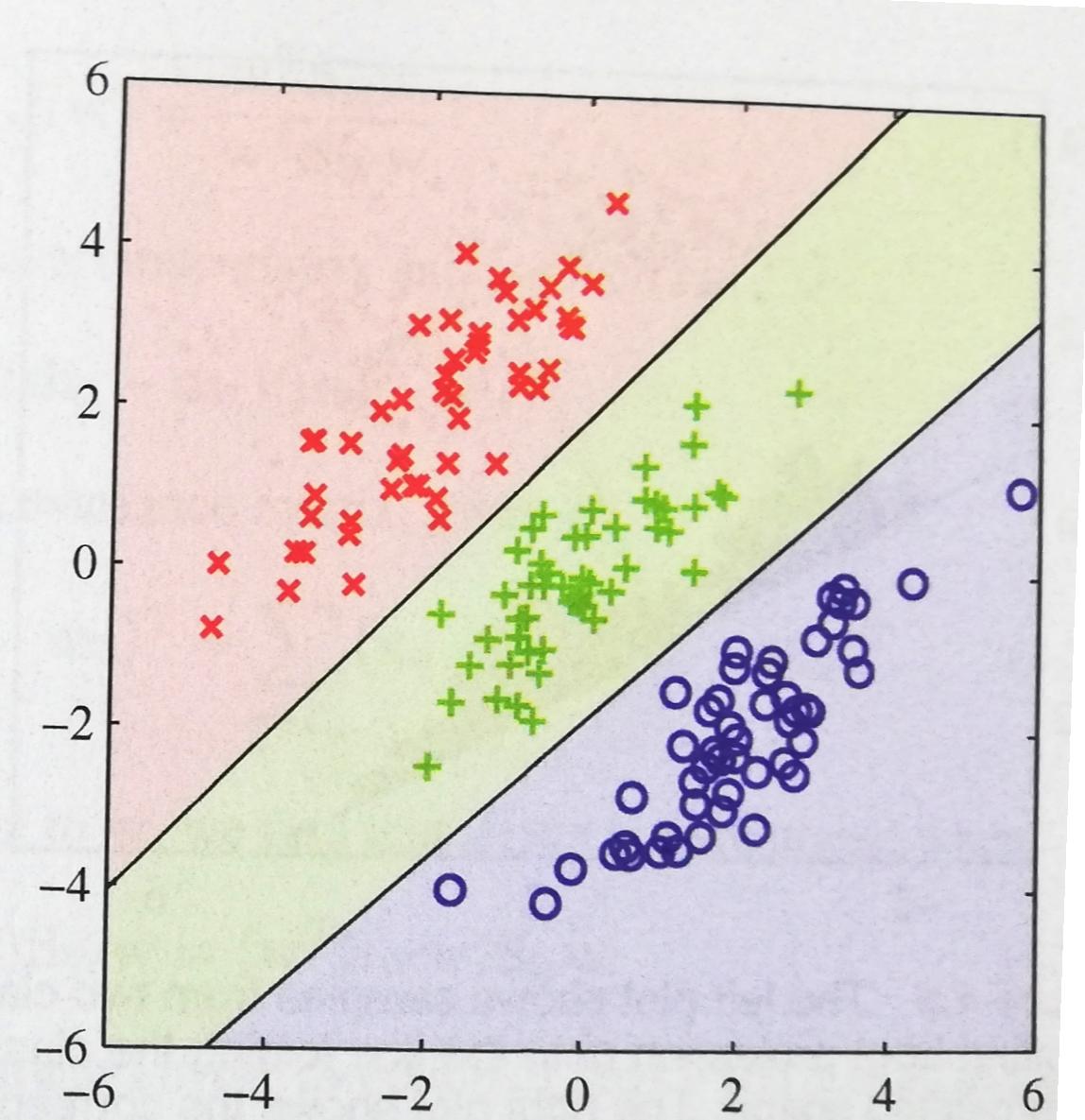


Limitations of the linear classifier

- Also the categorisation for multiclass problems can be problematic even when the **decision boundaries** should be easy to find



Linear classifier



Logistic regression

Logistic Regression

Logistic regression

in many cases a **soft classifier**, that **outputs the probability** of a given category, is advantageous over a hard classifier. Also, non-linear models exhibit many positive features for classification problems

In **logistic regression** the probability that a data point x_i belongs to a category y_i is given by

$$P(y_i = 1 | x_i, \theta) = \frac{1}{1 + e^{-x_i^T \theta}} = \sigma(x_i^T \theta)$$

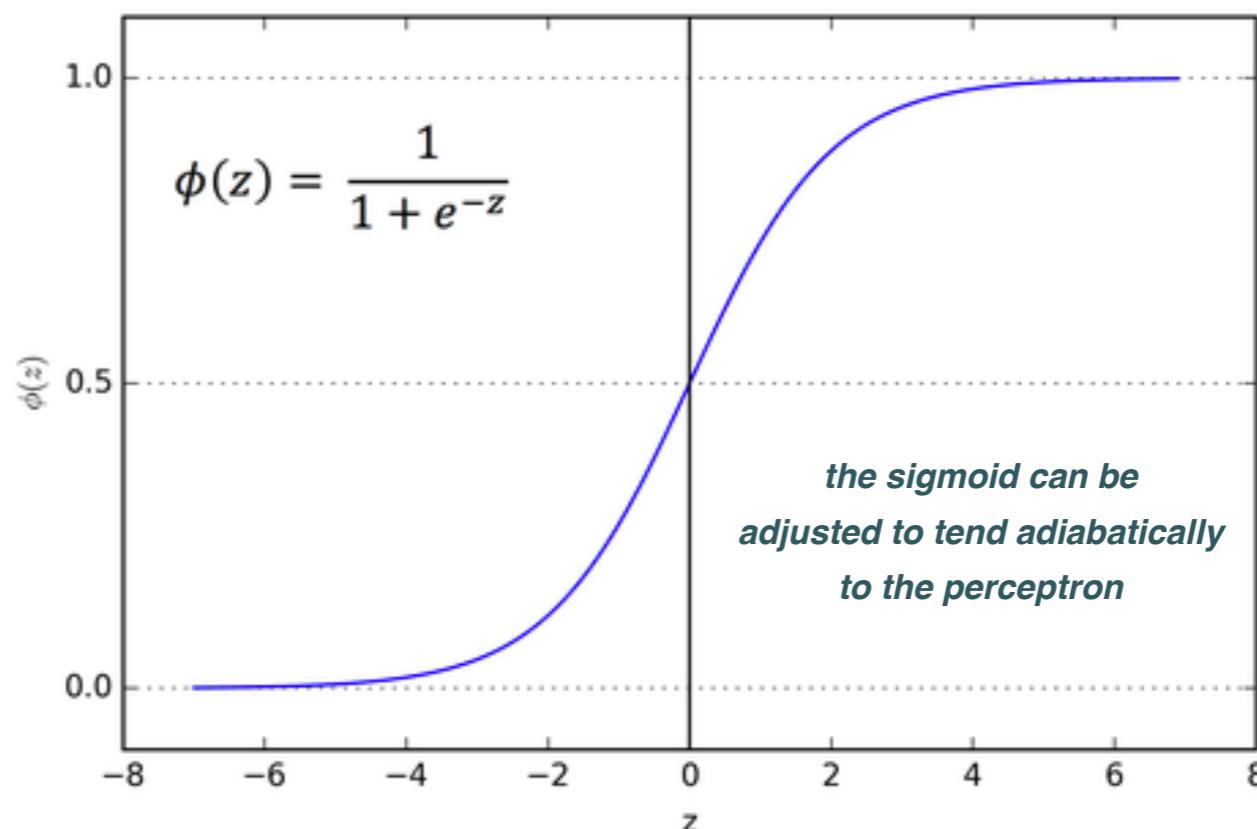
*non-linear
sigmoid function*

$$P(y_i = 0 | x_i, \theta) = 1 - P(y_i = 1 | x_i, \theta) = 1 - \sigma(x_i^T \theta) = \sigma(-x_i^T \theta)$$

where we have used the logistic (or sigmoid) function:

$$\sigma(s) = \frac{1}{1 + e^{-s}}$$

$$\sigma(-s) = 1 - \sigma(s)$$



Logistic regression

cost function for logistic regression from **Maximum Likelihood Estimation (MLE)**:
choose parameters that maximise the probability of seeing the observed data

$$P(y_i | \mathbf{x}_i, \boldsymbol{\theta}) = (P(y_i = 1 | \mathbf{x}_i, \boldsymbol{\theta}))^{y_i} \times (P(y_i = 0 | \mathbf{x}_i, \boldsymbol{\theta}))^{1-y_i}$$

``*probability interpolation*'': '*recovers limits when $y_i=0, y_i=1$*

$$P(y_i | \mathbf{x}_i, \boldsymbol{\theta}) = (\sigma(\mathbf{x}_i^T \boldsymbol{\theta}))^{y_i} \times (1 - \sigma(\mathbf{x}_i^T \boldsymbol{\theta}))^{1-y_i}$$

assuming that all observations are Bernoulli **independent**, the total likelihood is

$$\mathcal{L}(\boldsymbol{\theta} | \mathcal{D}) = \prod_{i=1}^n P(y_i | \mathbf{x}_i, \boldsymbol{\theta}) = \prod_{i=1}^n (\sigma(\mathbf{x}_i^T \boldsymbol{\theta}))^{y_i} \times (1 - \sigma(\mathbf{x}_i^T \boldsymbol{\theta}))^{1-y_i}$$

*note that this is soft classification: the probability is not only
0 or 1 but any value in between is possible*

Logistic regression

cost function for logistic regression derived from Maximum Likelihood Estimation (MLE):
choose parameters that maximise the probability of seeing the observed data

$$\mathcal{L}(\boldsymbol{\theta} | \mathcal{D}) = \prod_{i=1}^n P(y_i | \mathbf{x}_i, \boldsymbol{\theta}) = \prod_{i=1}^n (\sigma(\mathbf{x}_i^T \boldsymbol{\theta}))^{y_i} \times (1 - \sigma(\mathbf{x}_i^T \boldsymbol{\theta}))^{1-y_i}$$

Since the cost function is the negative log-likelihood, we find that for logistic regression

$$E(\boldsymbol{\theta}) = -\log \mathcal{L} = \sum_{i=1}^n - (y_i \log \sigma(\mathbf{x}_i^T \boldsymbol{\theta}) + (1 - y_i) \log(1 - \sigma(\mathbf{x}_i^T \boldsymbol{\theta})))$$

which is known as the cross-entropy function

The parameters of the model are determined by minimising the cross-entropy

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \left\{ \sum_{i=1}^n (-y_i \log \sigma(\mathbf{x}_i^T \boldsymbol{\theta}) - (1 - y_i) \log(1 - \sigma(\mathbf{x}_i^T \boldsymbol{\theta}))) \right\}$$

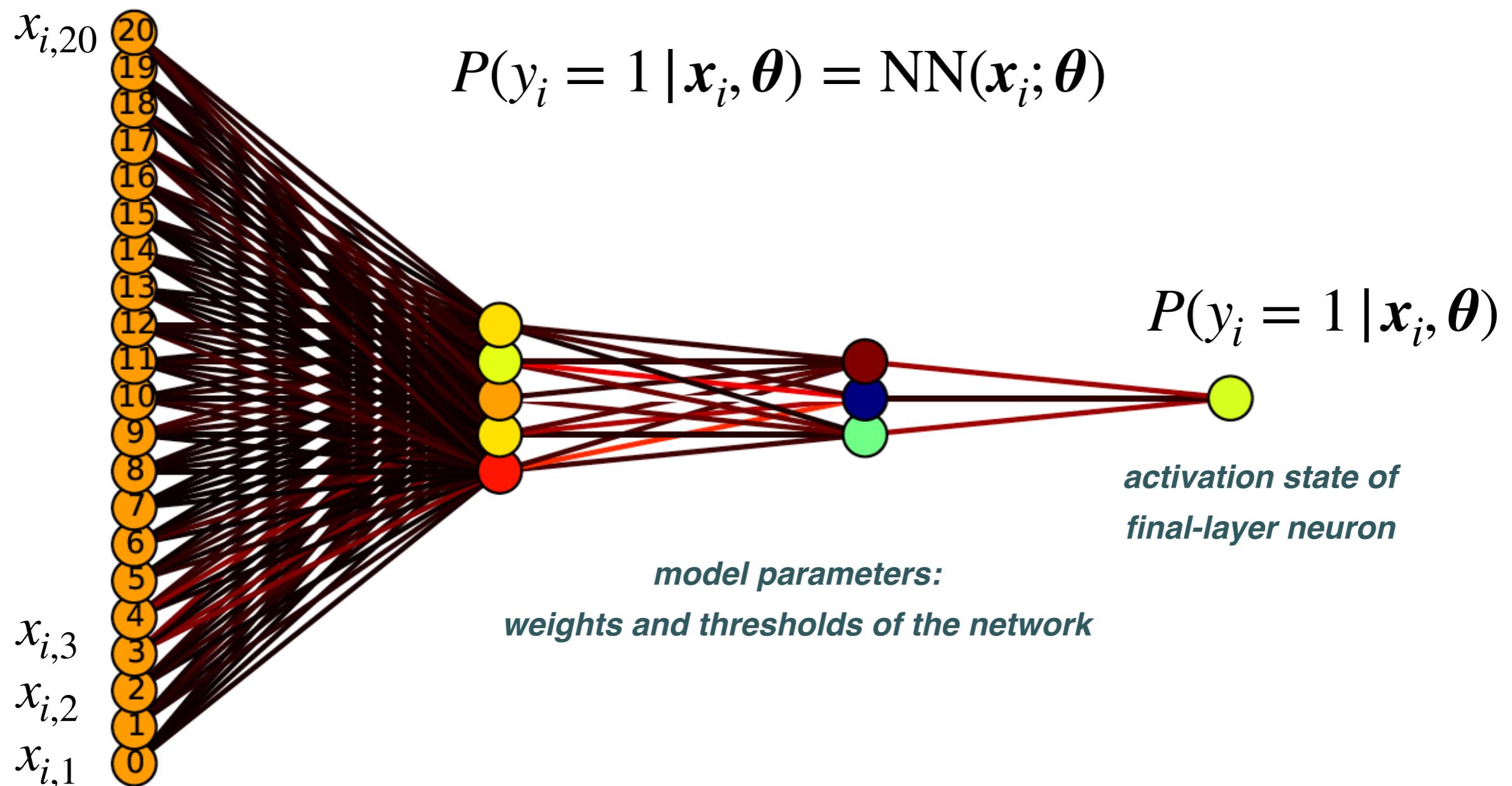
***note that no analytic solution is possible, and numerical methods are required
one can use more complex models to parametrise the probabilities, such as NNs***

Neural Networks for Classification

instead of modelling the classification probability by a **simple logistic function**

$$P(y_i = 1 | \mathbf{x}_i, \theta) = \frac{1}{1 + e^{-\mathbf{x}_i^T \theta}}$$

one can adopt more complex models, such as deep neural networks



Decision Theory

Decision Theory

we have discussed several ML models that provide a **probability distribution** for a given outcome given some inputs. However, we still need to **take decisions** based on this info

assume that we have the usual **training dataset**, either for regression or classification

$$(\mathbf{x}_1, \dots, \mathbf{x}_N) , \quad t = (t_1, \dots, t_N)$$

depending on our application, we can try to achieve different goals

- Given a new input vector \mathbf{x} , make a prediction for t
- Given a new input vector \mathbf{x} , **make a decision** based on the expectations for t
- Construct the joint probability distribution $p(\mathbf{x}, t)$ associated to the training dataset: this is known as the **inference problem**

the goal of Decision Theory is to help us making **optimal decisions** given our knowledge of the involved probabilities

Decision Theory

for example, our goal could be to **minimise the misclassification rate**

e.g. in clinical diagnosis that require an invasive treatment

we will have trained our classifier on the input examples, and end up with a map of the input space into **decision regions** (assigned to a specific class) separated by **decision boundaries**

the probability of incurring in a **classification mistake** will be given by

$$p(\text{mistake}) = p(x \in \mathcal{R}_1, \mathcal{C}_2) + p(x \in \mathcal{R}_2, \mathcal{C}_1) = \int_{\mathcal{R}_1} p(x, \mathcal{C}_2) dx + \int_{\mathcal{R}_2} p(x, \mathcal{C}_1) dx$$

probability that an input in decision region R_1 is classified as class C_2

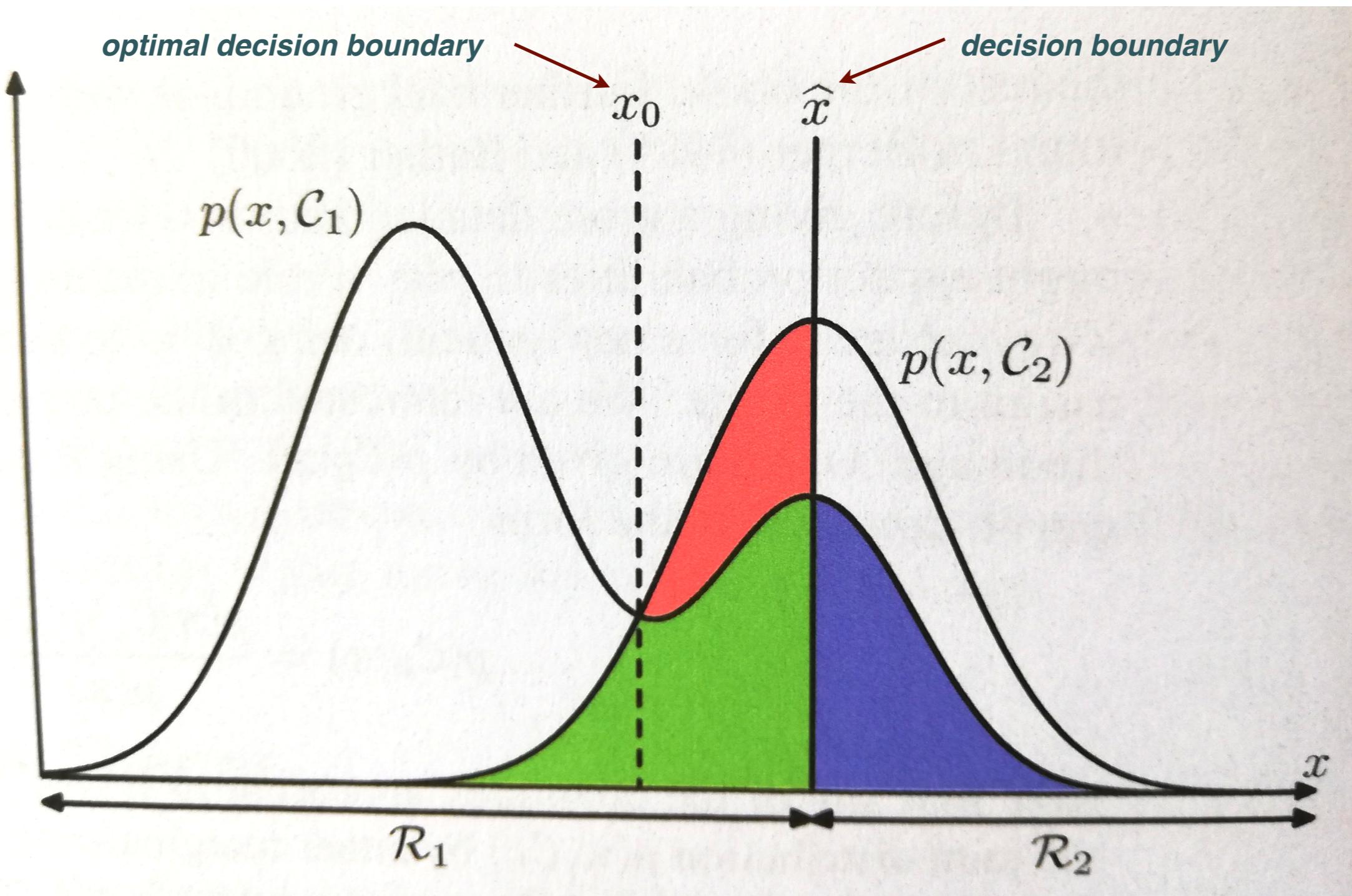
note that we still need to decide the **rule** that assigns a given point x to each category.

From the product rule of probabilities we have that

$$p(x, \mathcal{C}_k) = p(\mathcal{C}_k | x)p(x)$$

The misclassification rate is minimised by the decision boundary defined when each point x is assigned to the class for which its posterior probability $p(\mathcal{C}_k | x)$ is the largest

Decision Theory



red, green, blue regions: **misclassification probability**

green+blue constant as decision boundary varied: optimal choice is that one that minimised red region

Metrics for binary classification

there are many **useful metrics** that are used in ML binary classification problems

example classification problem: 34 training samples, of which

True Positives (TP) e.g. 8	False Positives (FP) e.g. 2
False Negatives (FN) e.g. 4	True Negatives (TN) e.g. 20

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} = \frac{8 + 20}{8 + 20 + 2 + 4} = 0.823$$

however accuracy is not the right metric for classification problems defined by a large disparity between classes, e.g. when $TN \gg TP$

$$\text{Accuracy} \simeq \frac{\text{TN}}{\text{TN} + \text{FN}} + \dots$$

Metrics for binary classification

there are other **useful metrics** that are used in ML binary classification problems

example classification problem: 34 training samples, of which

True Positives (TP) e.g. 8	False Positives (FP) e.g. 2
False Negatives (FN) e.g. 4	True Negatives (TN) e.g. 20

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} = \frac{8 + 20}{8 + 20 + 2 + 4} = 0.823$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} = 0.80$$

proportion of correct positive classifications

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} = 0.67$$

proportion of correct actual positive classifications

ROC curve

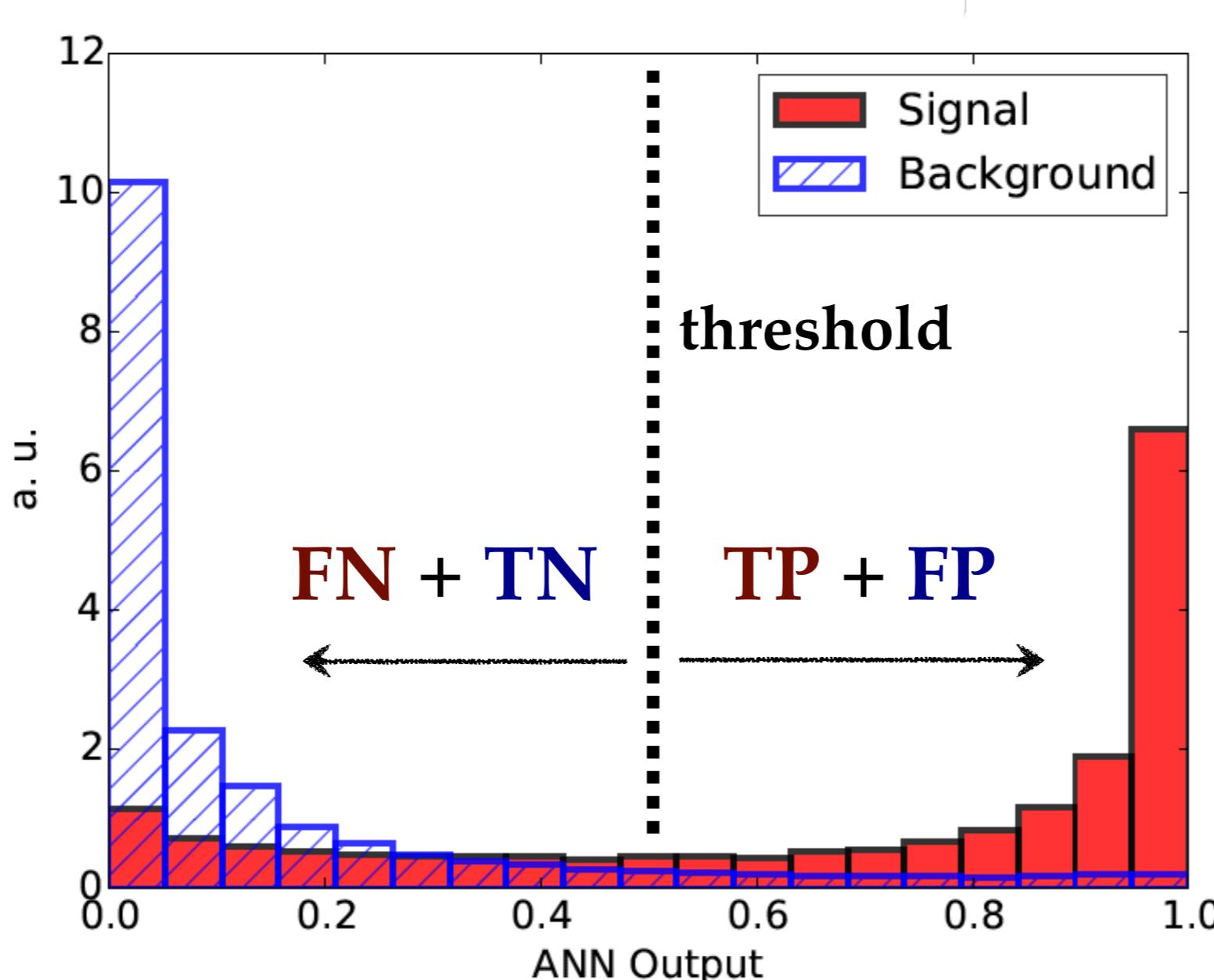
In most ML classification problems there is a **threshold** that can be varied to decide at what output of the model a given data point is assigned to each category:

how should we select this threshold for classification?

ROC curve

In most ML classification problems there is a **threshold** that can be varied to decide at what output of the model a given data point is assigned to each category:

- A conservative threshold will **maximise TP**, but also FP might be large
- An aggressive threshold **reduces FP** but then FN might be large.



ROC curve

In most ML classification problems there is a **threshold** that can be varied to decide at what output of the model a given data point is assigned to each category:

- A conservative threshold will **maximise TP**, but also FP might be large
- An aggressive threshold **reduces FP** but then FN might be large.

effect of threshold is quantified by the **Receiver Operating Characteristic (ROC)** curve

$$\text{Recall} = \text{True Positive Rate} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad \xleftarrow{\text{proportion of correctly classified positives}}$$

vs

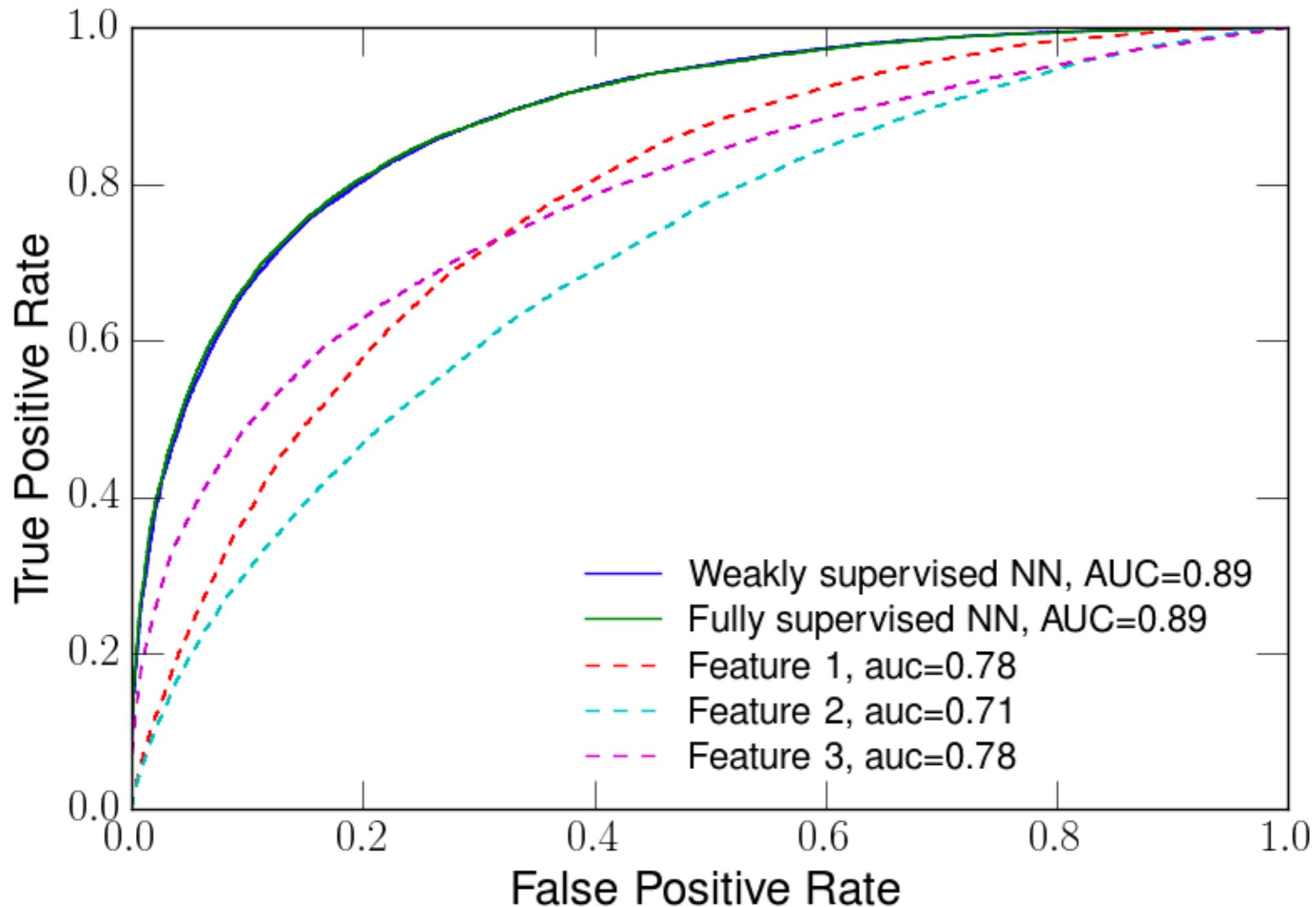
$$\text{False Positive Rate} = \frac{\text{FP}}{\text{FP} + \text{TN}} \quad \xleftarrow{\text{incorrectly classified negatives}}$$

a good classifier should maximise TPR while minimising FPR

ROC curve

$$\text{Recall} = \text{True Positive Rate} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{False Positive Rate} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

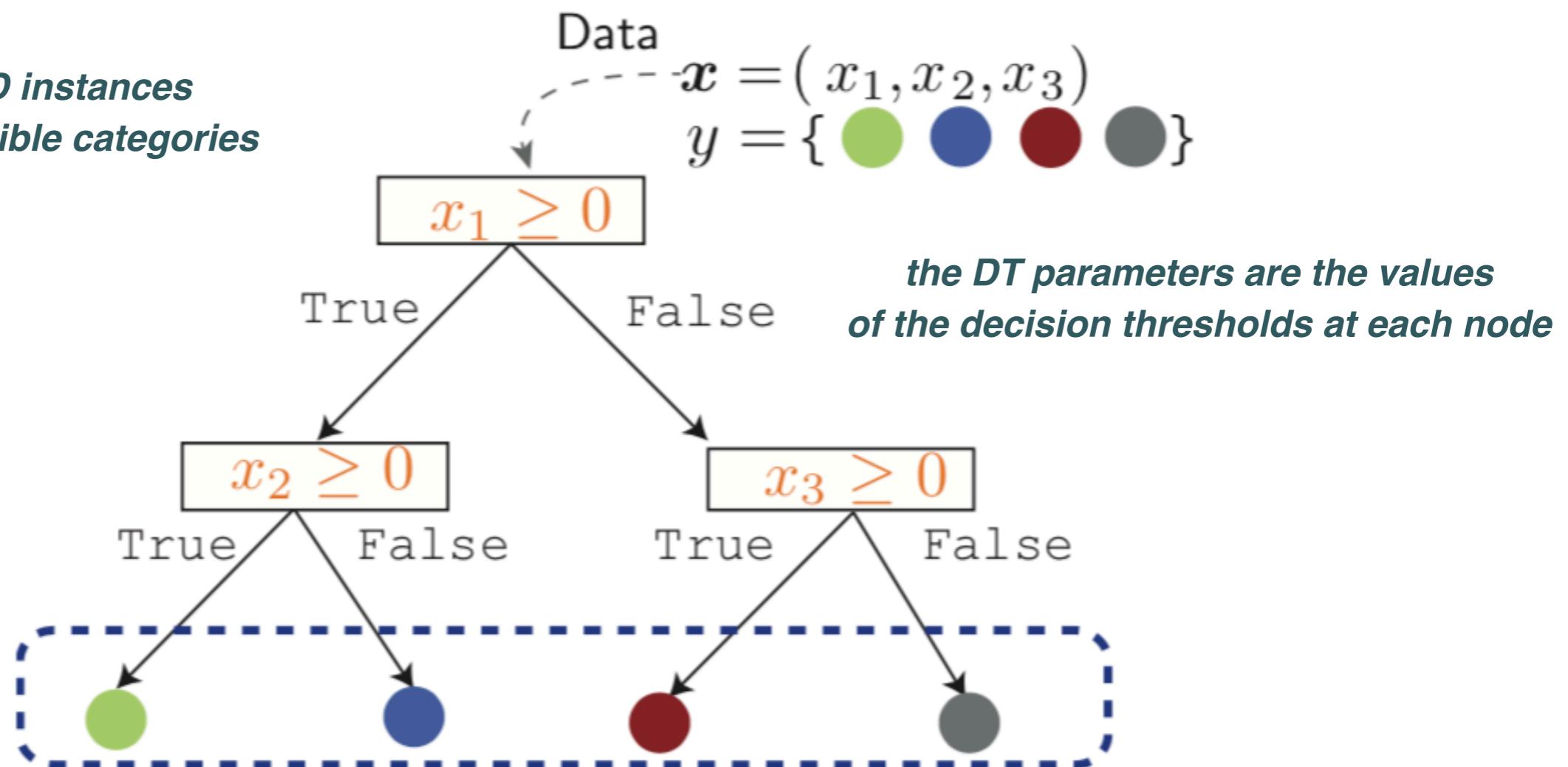


Decision Trees and Random Forests

Decision Trees

Decision Trees are classifiers that work by partitioning the input space into **hypercubes** and then assign a model (e.g. a constant) to each region

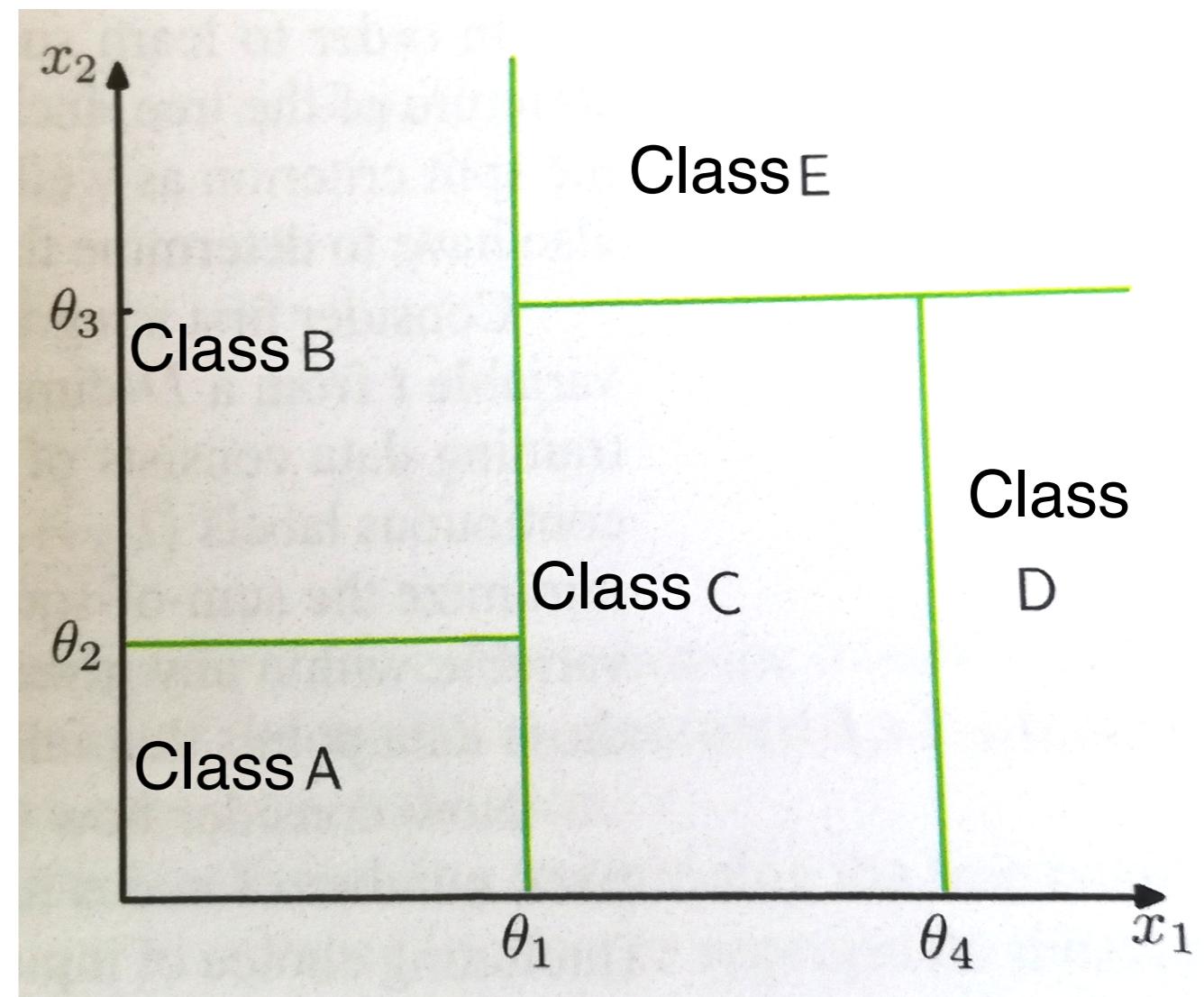
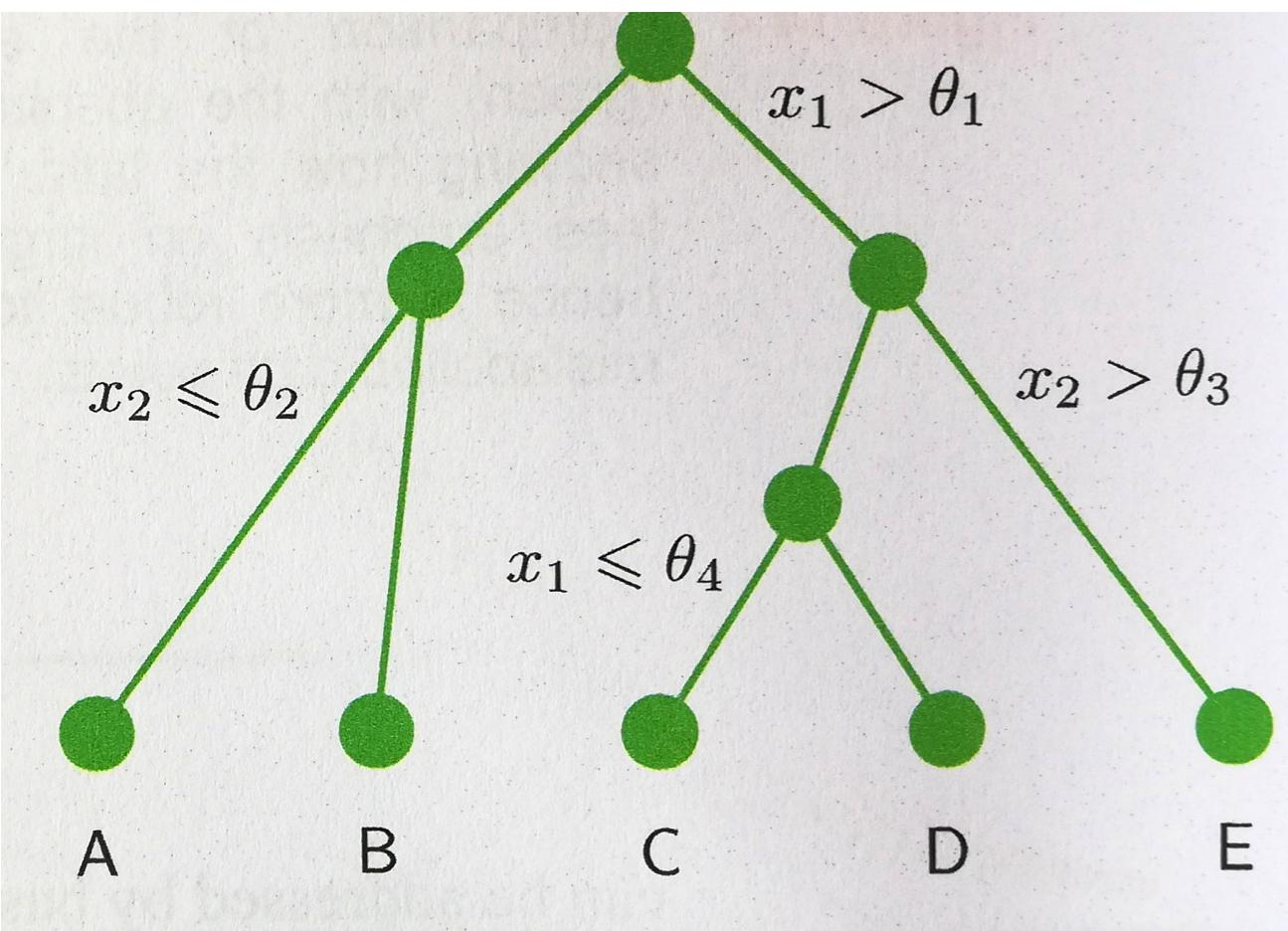
*classify 3D instances
into four possible categories*



In ML, a **decision tree** is an algorithm which uses a series of questions to hierarchically partition the data, where each branch of the tree splits the data into smaller subsets

Decision Trees

Decision trees have the key property that their output is **human-interpretable**: sequence of binary decisions applied to individual input variables



to train a decision tree we need to determine which **variable is chosen at each node** for the split criterion and the value of the **threshold for the split**

Random Forests

individual trees have often **high variance** and are weak classifiers:
we can improve by incorporating them in an ensemble method

→ We need an ensemble of **randomised decision trees** (minimised correlations)

- ➊ (1) train each decision tree on a different bootstrapped dataset: **bagged decision tree**
- ➋ (2) use different random subset of features at each split: **random forest**
*reduces correlations between trees that arise
when only few features are strongly predictive*

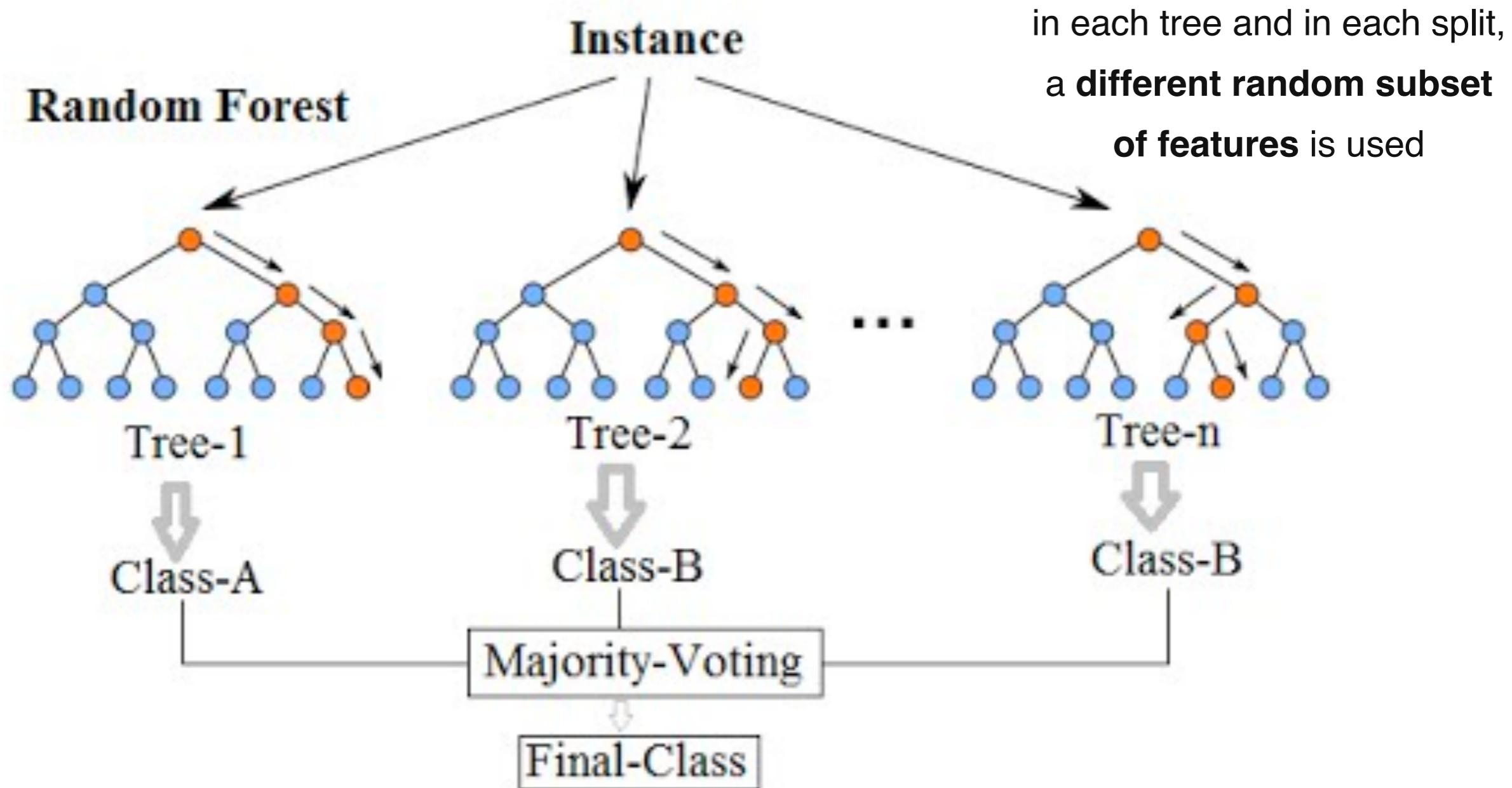
as other ML learning classifiers also Random Forests require some **regularisation**, for example a maximum depth of the tree, to control complexity and prevent overfitting

typically, a classification problem with p features, in RFs only $p^{1/2}$ features are used in each split

Random forests have other attractive features, for example, they can be used to **rank the importance of variables** in a regression or classification problem in a natural way

Random Forests

each decision tree in the ensemble is built upon a
random bootstrap sample of the original data



the final categorisation is assigned e.g. from **majority voting**

Support Vector Machines

Support Vector Machines

SVM are ML algorithms popular for classification, regression, and novelty detection

A peculiar feature of SVM is that the model parameters are found by a solution of a convex optimisation problem: any local solution is also a **global optimum**

The advantages of support vector machines include:

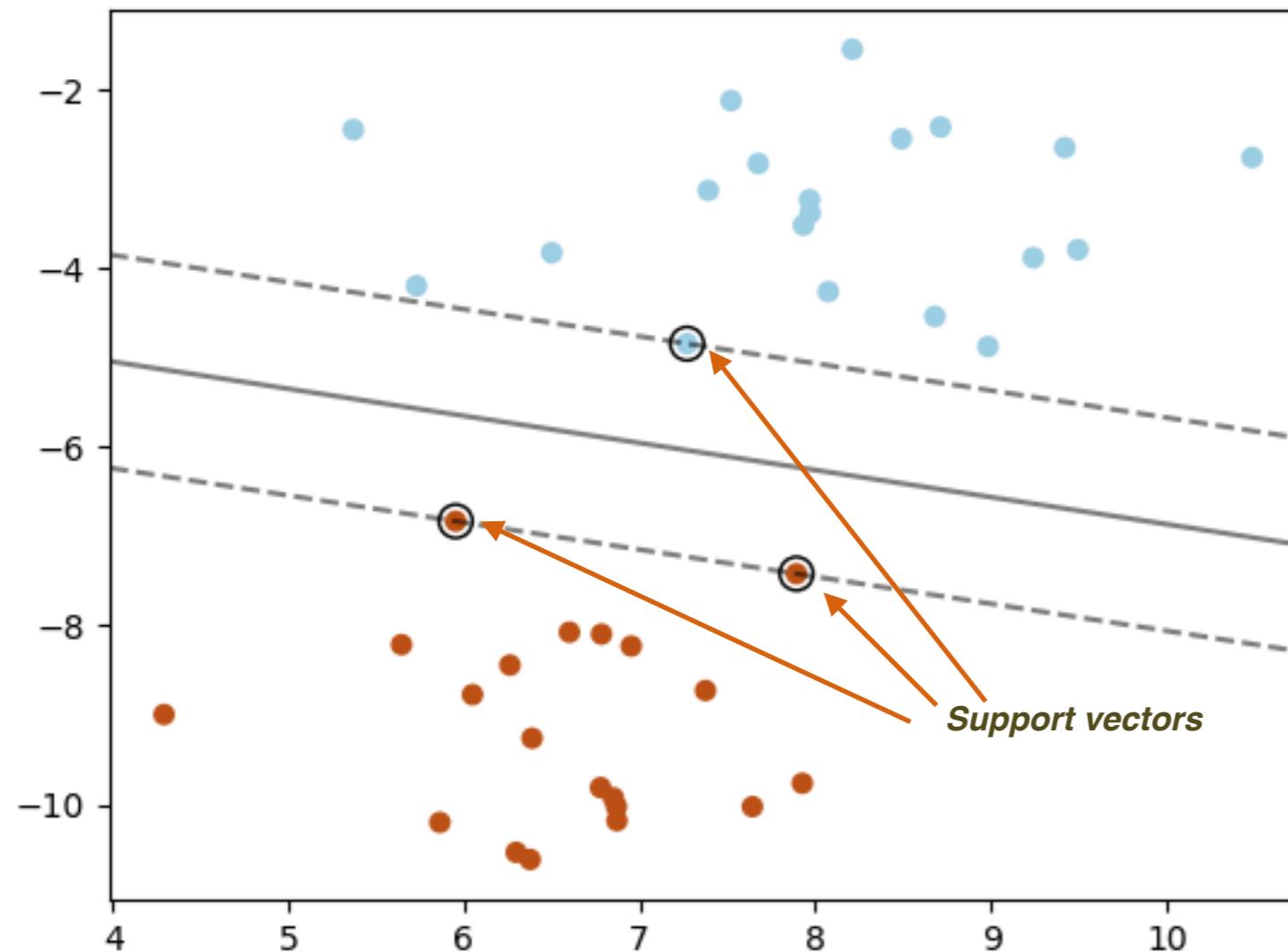
- ✿ Effective in **high dimensions**.
- ✿ Effective when number of dimensions is greater than number of samples.
- ✿ Based on a subset of training points in the decision function (**support vectors**)
- ✿ Different **kernels functions** can be used for the figure of merit

In ML, kernel methods operate in a high-dimensional, implicit feature space without computing the coordinates in data space

Support Vector Machines

SVM are ML algorithms popular for classification, regression, and novelty detection

A peculiar feature of SVM is that the model parameters are found by a solution of a convex optimisation problem: any local solution is also a **global optimum**



Support Vector Machines

As other **classification algorithms**, the problem that SVM try to solve is the following. Starting from a training dataset

$$\{\mathbf{x}_i^T, t_i\}, \quad i = 1, \dots, N, \quad t_i \in (-1, 1)$$

our goal is to find the **model parameters** such that the prediction from

$$\text{sign}(\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}) + b)$$

is correct for most of the samples. Here $\boldsymbol{\phi}$ is the a so-called **feature-space transformation**, where the *feature space* is where our data lives

Example of a feature-space transformation

\mathbf{x}	$\boldsymbol{\phi}(\mathbf{x})$
<i>Height</i>	<i>Height/Weight</i>
<i>Weight</i>	<i>Weight+Height</i>
<i>Age</i>	<i>Age * Weight</i>

Support Vector Machines

the problem then that SVM are aiming to solve is

$$\min_{\omega, b, \xi} \left(\frac{1}{2} \theta^T \theta + C \sum_{i=1}^n \xi_i \right)$$

margin between the two categories

subject to the **additional constraint** that

$$Class\ label\ of\ sample\ i \longrightarrow t_i \times (\theta^T \phi(x_i) + b) \geq 1 - \xi_i \longleftarrow tolerance\ (regulator)$$

which essentially aims to **maximise the margin between the two categories** while including a penalty when a sample is misclassified or within the margin boundary, where the term proportional to C plays the role of a **regulator**

SVM are an example of a **constrained optimisation problem**, where one tries to minimise a quadratic function subject to **linear constraints**. They are most efficiently solved by means of the Lagrange multiplier method

Support Vector Machines



<https://www.youtube.com/watch?v=5zRmhOUjjGY>