

Machine Learning

Rudolf Mayer
mayer@ifs.tuwien.ac.at

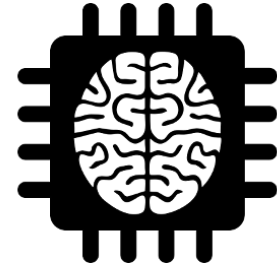
October 11th, 2017

- Machine Learning: Intro
- Perceptron
- k-NN
- Performance evaluation (intro)
- Types of data & data preparation

- Machine Learning: Intro
- Perceptron
- k-NN
- Performance evaluation (intro)
- Types of data & data preparation

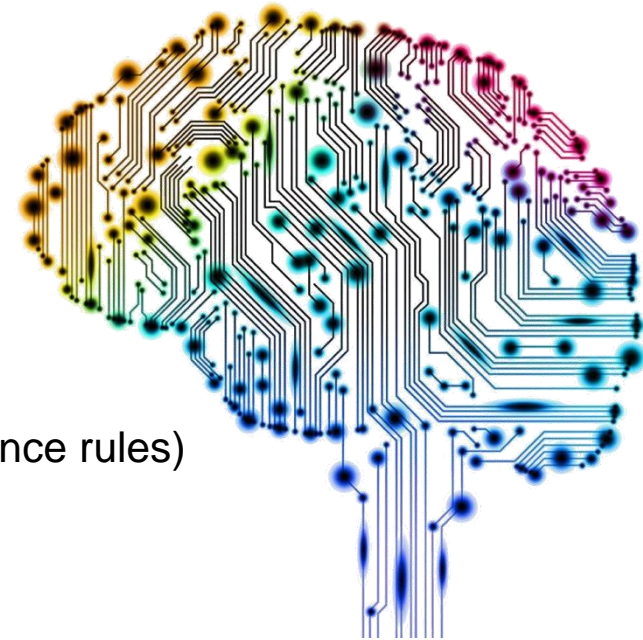
Defining the area

- Some related terms
 - Data Mining
 - Predictive analytics
 - Knowledge discovery
 - Data Science
 - Statistics
 - Cluster Analysis
 - Artificial Intelligence
 - Reasoning / deduction
 - Regression
 - Classification
 - Supervised/ unsupervised learning
 - ...



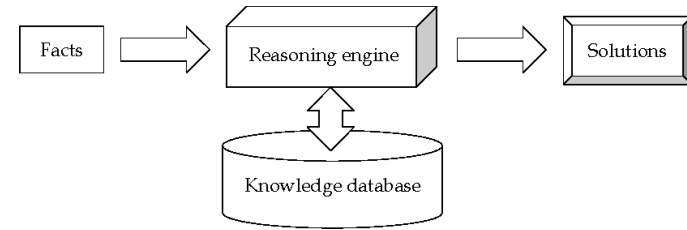
- Deals with intelligence of machines
 - *System that perceives its environment & takes actions that maximise its chances of success*

- Areas/problems of AI
 - Deduction, reasoning, problem solving
 - Often rule-based; manually created
 - Knowledge representation (reasoning)
 - Reasoning: generation of new knowledge (inference rules)
 - **Machine learning**
 - Planning / scheduling
 - Natural language processing
 - ...



- System that uses rules to make deductions or choices

- E.g. domain-specific expert system



- Two components

- Knowledge base: facts & rules (if \rightarrow then style)

- Knowledge representation (language)
 - Example rule: *IF (hot & smoke) THEN fire*

- Inference engine: applies rules to deduce new facts

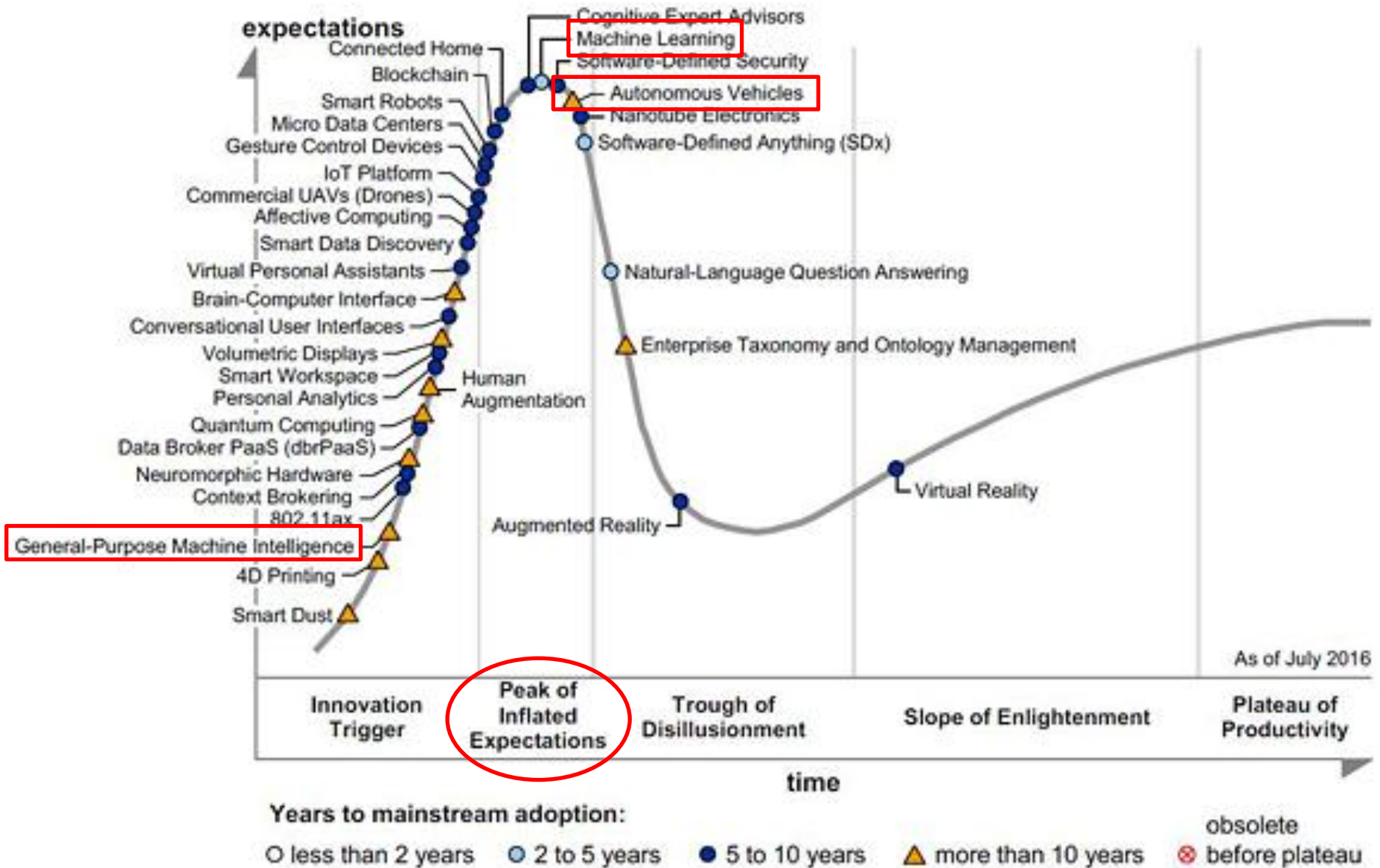
- Forward chaining: assert new facts
 - Backward chaining: start with goal \rightarrow determine which facts need to asserted

- Rules often manually specified (by expert)

- Expensive, incomplete

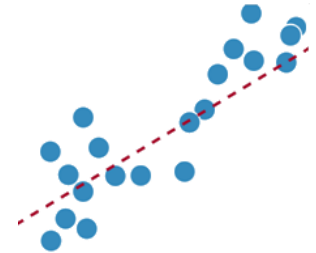
- Studies computer algorithms that can ***learn*** from data and make ***predictions*** on data
- Automatic methods – no human assistance in learning (i.e. no specification of rules!)
 - Human assistance generally required for
 - Defining problem
 - Gathering and assessing data
- Types of Machine Learning
 - Supervised
 - Unsupervised
 - Reinforcement

Machine learning – why?



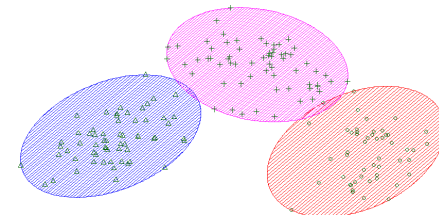
- Statistics

- Collection, interpretation and presentation of data
- Descriptive: summarise data (mean, standard deviation, probability density, correlation analysis)
- Predictive modelling: e.g. regression

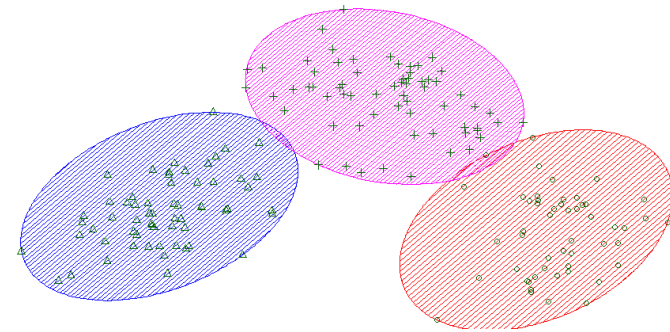


- Data Mining

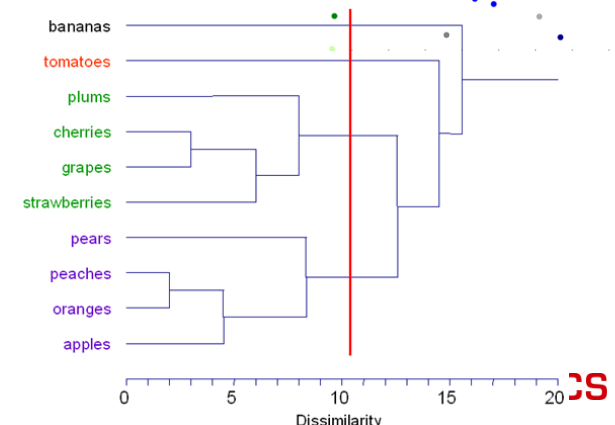
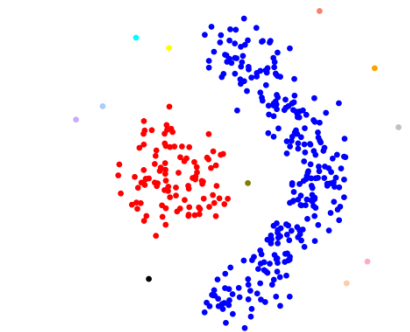
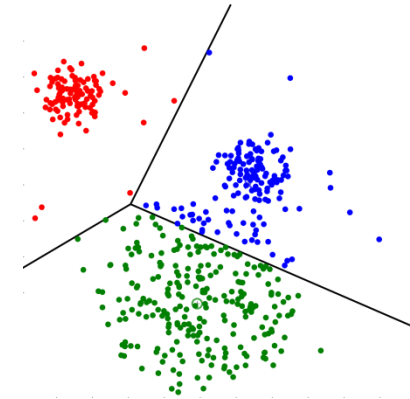
- Discovering patterns in large datasets
 - Cluster analysis, outlier detection, association rule mining



- Unsupervised learning
 - Data not *labelled*
 - No information on which and how many classes or other structures, ...
 - Goal:
 - Find structures (e.g. clustering)
 - Association rules learning
 - Find outliers (anomalies)
 - Most often associated with *Data Mining*
- Covered in other lectures
 - Business Intelligence, Selbstorganisierende Systeme, ..



- Explorative method
- Find groups of similar objects
- Applications: e.g. market segmentation
- Several popular algorithms
 - K-means clustering (centroid based)
 - DBSCAN (density-based)
 - Linkage (hierarchical)
 - single, complete, average, Ward, ..



- Animal data set
 - Describes animal by some characteristics
 - Instances: cow, duck, cat, bee, sparrow, ...
- Characteristics
 - Size (tiny, small, medium, big)
 - Number of legs (2, 4, 6, 8)
 - Feathers (yes/no)
 - Eggs (yes/no)

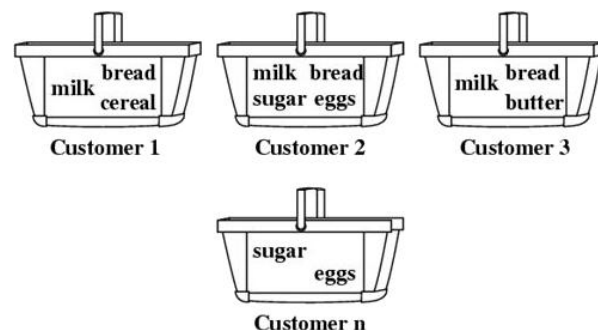
	<i>size</i>	<i>legs</i>	<i>feathers</i>	<i>eggs</i>
duck	small	2	yes	yes
dog	medium	4	no	no
spider	tiny	8	no	yes
ladybird	tiny	6	no	yes
cow	large	4	no	no
bee	tiny	6	no	no
sparrow	small	2	yes	yes

- Animal data set
 - Describes animal by some characteristics

	<i>size</i>	<i>legs</i>	<i>feathers</i>	<i>eggs</i>
duck	small	2	yes	yes
dog	medium	4	no	no
spider	tiny	8	no	yes
ladybird	tiny	6	no	yes
cow	large	4	no	no
bee	tiny	6	no	no
sparrow	small	2	yes	yes

- Goal: find groups of related animals:
 - Mammals: cat, cow
 - Birds: duck, sparrow
 - Insects: bee, ladybird
 - Invertebrate: spider

- Discover relations between variables
 - E.g. market basket analysis: which items are frequently bought together
 - ➔ Useful for marketing
 - Identify rules that are
 - Frequent (“support”)
 - Reliable (“confidence”)
 - Anecdotal example: beer & diapers
 - Related: collaborative filtering / recommender systems



- Supervised learning
 - Data labelled with actual **output** variable
 - Goal: correctly label **unknown** data
- Sometimes equivalently used with “machine learning”
 - In some definitions, machine learning means both unsupervised and supervised learning
 - In other definitions, unsupervised learning mostly equals/is a subset of data mining; and rather seen as part of statistics
- *Types of supervised learning?*
 - Regression & Classification

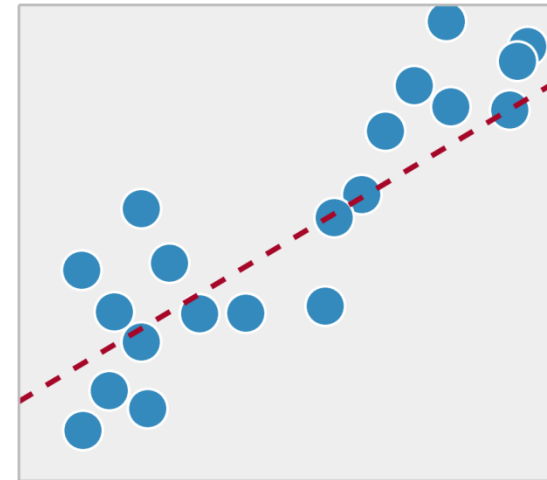


Regression vs. Classification

- Regression tries to predict a ***continuous*** variable
 - Examples?

Regression vs. Classification

- Regression tries to predict a ***continuous*** variable
 - *Income of a house hold* (depending on education, ...)
 - *Temperature* (depending on wind, humidity...)
 - *Housing price* (depending on e.g. size, age..)
 - Methods: e.g. *linear regression*
 - Mostly associated with statistics
 - Many classification methods can also be used for regression

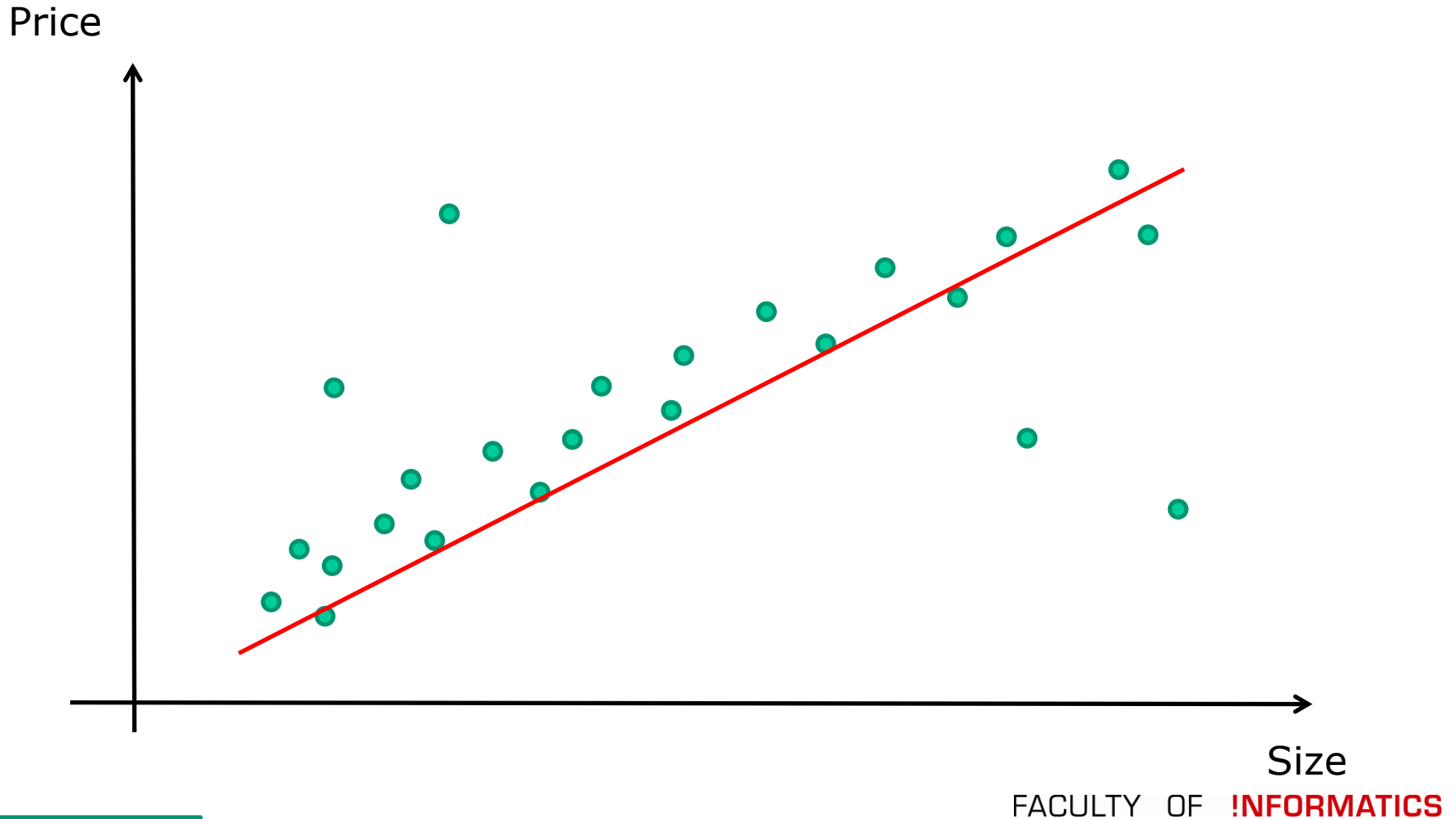


Regression example

- Data set about prices of houses (apartments)
- Input data e.g.
 - Size
 - Number of rooms
 - Size of garden/balcony/terrace
 - Year built / age of building
 - Location
- Goal: predict the expected price of the house in €
- Solution: multivariate regression
 - *If only one input: univariate regression*

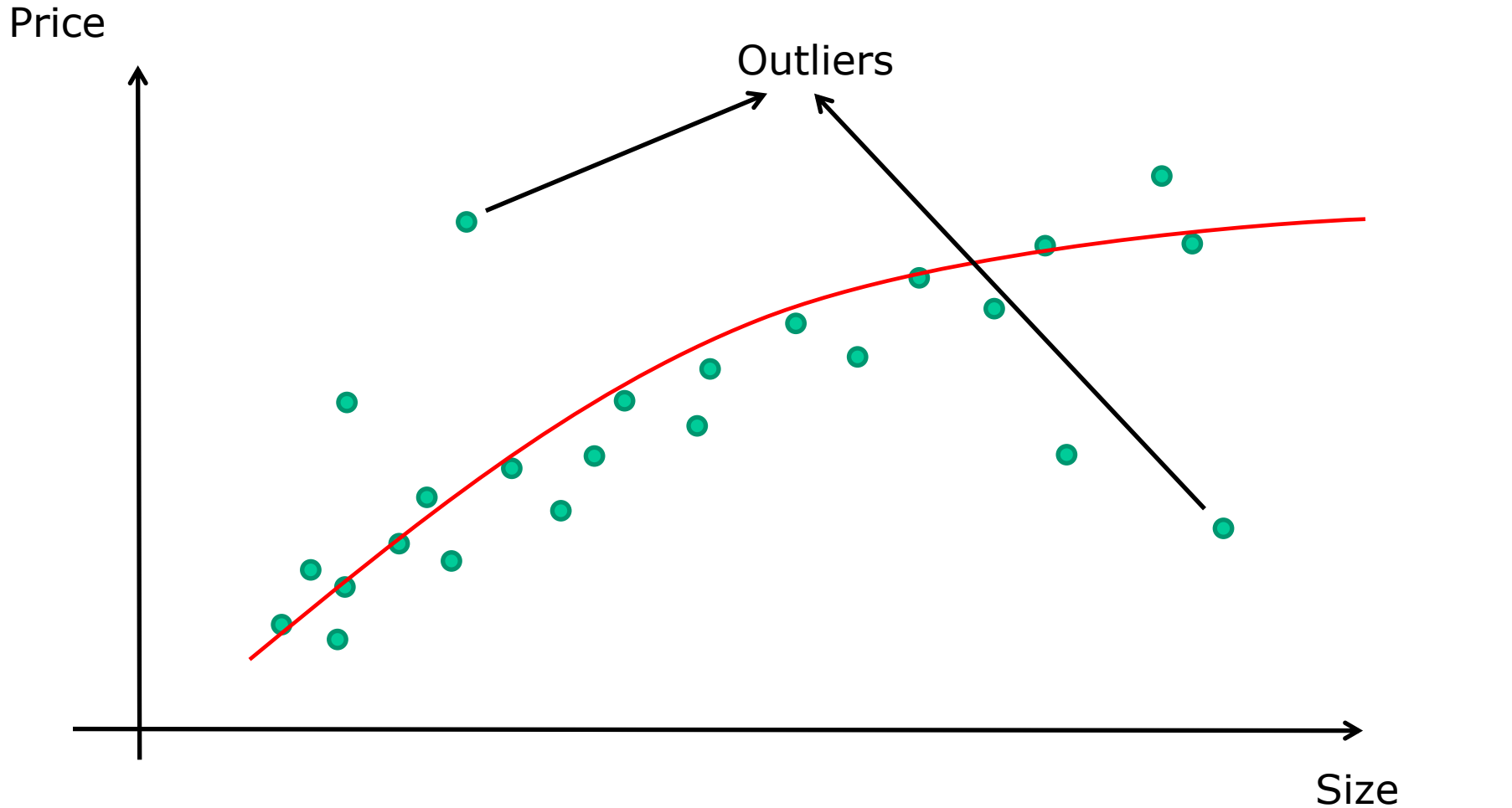
Regression example

Linear regression – find a “line” that describes the data



Regression example

Polynomial regression – find a “curve” that describes the data

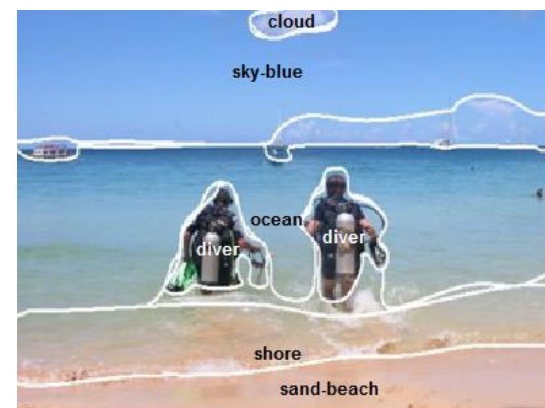
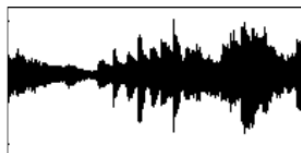


Regression vs. Classification

- If the output variable can take one of a predefined set of values: **classification** (*also: categorisation*)
 - Examples?

Regression vs. Classification

- If the output variable can take one of a predefined set of values: **classification** (also: categorisation)



Winter is here. Go to the store and buy some snow shovels.

Winter is here. Go to the store and buy some snow shovels.

Regression vs. Classification

- If the output variable can take one of a predefined set of values: **classification** (*also: categorisation*)

– **Information:** SPAM filtering



– **Image:** classification of hand-written letters for OCR; automatic labelling of images



– **Music:** classification of music into genres

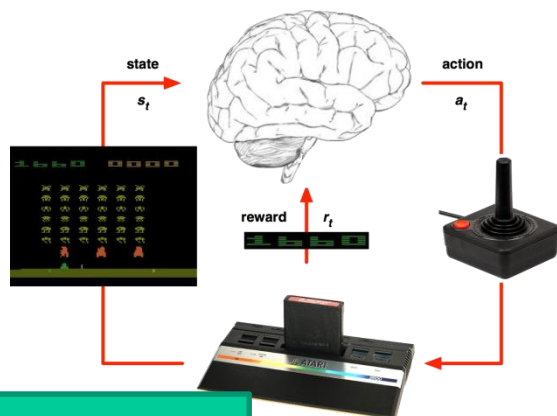
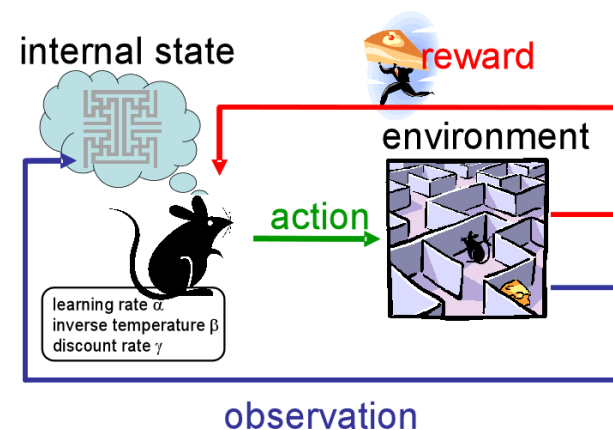


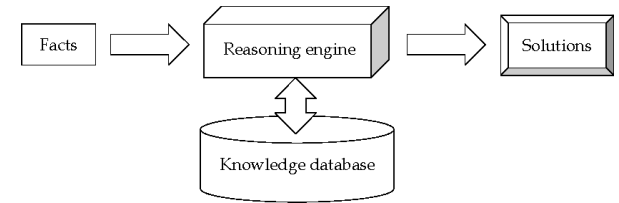
– **Medicine:** classification of whether a person has an illness, based e.g. on x-ray images, or other measurements



- Agent (e.g. computer program) takes actions in specified environment
 - Not explicitly presenting input/output pair
 - Reward (penalise) agent for actions
- ➔ Maximise cumulative reward

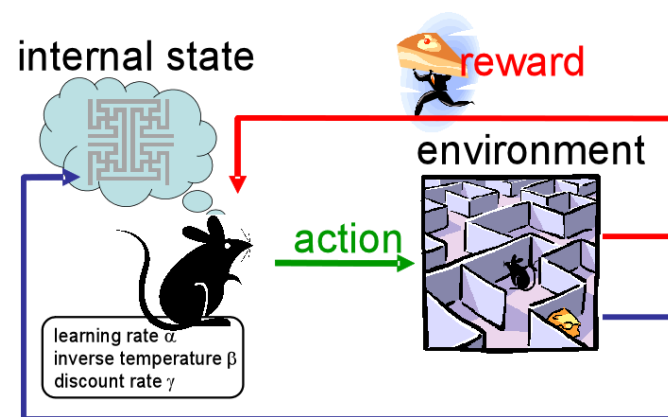
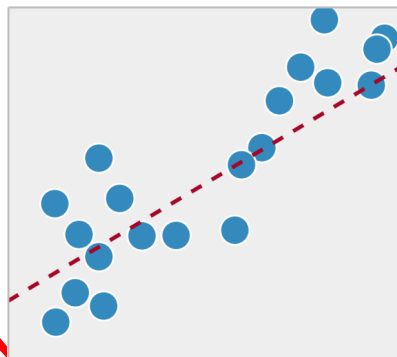
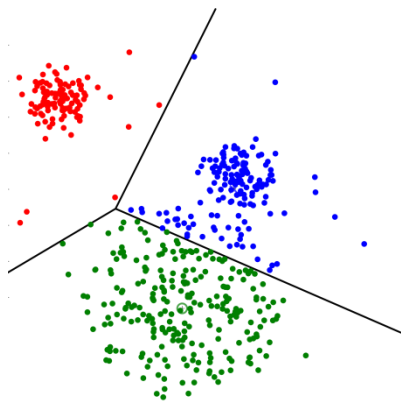
- Popular applications?





- *Rules often manually specified (by expert)*
- Alternative: learning rules
 - *Association rule mining – unsupervised*
 - Learning classifier systems
 - Can be supervised learning
 - Includes steps to discover new rules, selection of best rules, ...
- Related/similar approaches:
 - Ontologies (OWL: Web Ontology Language) & reasoning
 - *Bayesian Networks*

Scoping the lecture



- Example: data set describing characteristics of patients
 - Gender
 - Age
 - Smoker yes/no
 - Eye colour
- Want to predict whether a person has lung cancer
 - Available: some data with a label / annotation (cancer yes(no))



Classification: Setting

Attributes / Variables / Features

Output Variable

**T
r
a
i
n
i
n
g

S
a
m
p
l
e
s**

<i>gender</i>	<i>age</i>	<i>smoker</i>	<i>eye colour</i>
male	19	yes	green
female	44	yes	grey
male	49	yes	blue
male	12	no	brown
female	37	no	brown
female	60	no	brown
male	44	no	blue
female	27	yes	brown
female	51	yes	green
female	81	yes	grey
male	22	yes	brown
male	29	no	blue

<i>lung cancer</i>
No
yes
yes
no
no
yes
no
no
yes
no
no
no

Ground Truth /
Gold standard

→
→
→

male	77	yes	grey
male	19	yes	green
female	44	no	grey

?
?
?

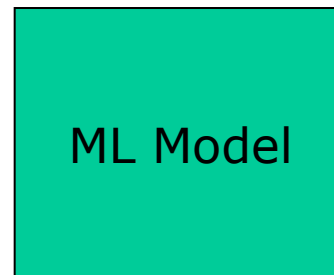
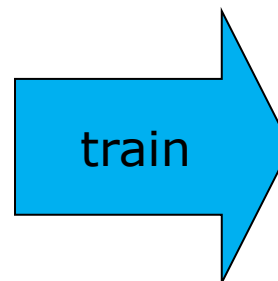
Prediction

Unlabelled samples

Classification: Setting

gender	age	smoker	eye colour
male	19	yes	green
female	44	yes	grey
male	49	yes	blue
male	12	no	brown
female	37	no	brown
female	60	no	brown
male	44	no	blue
female	27	yes	brown
female	51	yes	green
female	81	yes	grey
male	22	yes	brown
male	29	no	blue

lung cancer
No
yes
yes
no
no
yes
no
no
yes
no
no
no



male	77	yes	grey
male	19	yes	green
female	44	no	grey

?
?
?

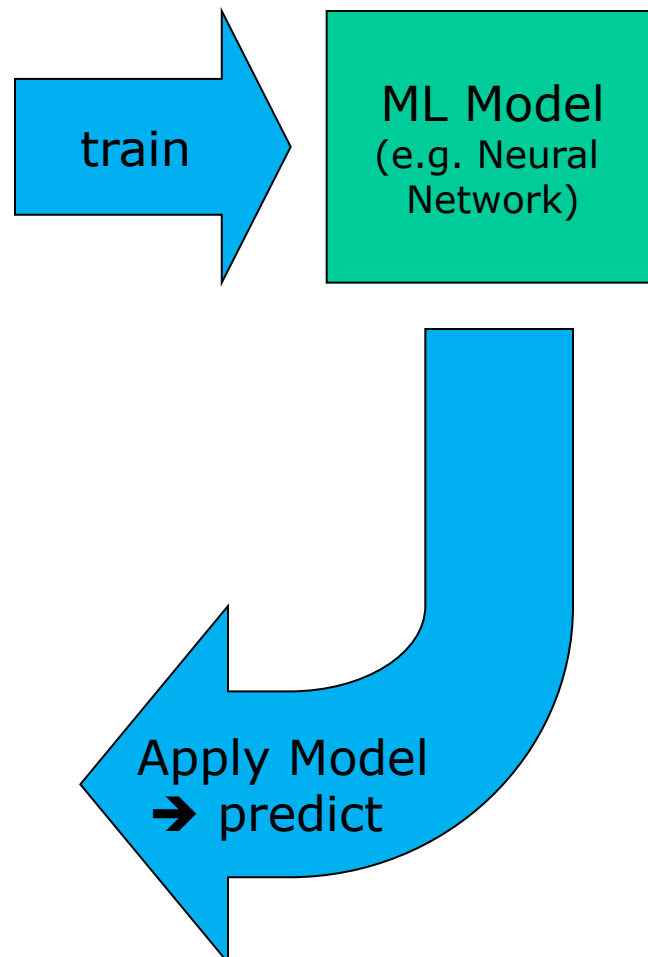
Classification: Setting

gender	age	smoker	eye colour
male	19	yes	green
female	44	yes	grey
male	49	yes	blue
male	12	no	brown
female	37	no	brown
female	60	no	brown
male	44	no	blue
female	27	yes	brown
female	51	yes	green
female	81	yes	grey
male	22	yes	brown
male	29	no	blue

lung cancer
No
yes
yes
no
no
yes
no
no
yes
no
no

male	77	yes	grey
male	19	yes	green
female	44	no	grey

yes
no
yes



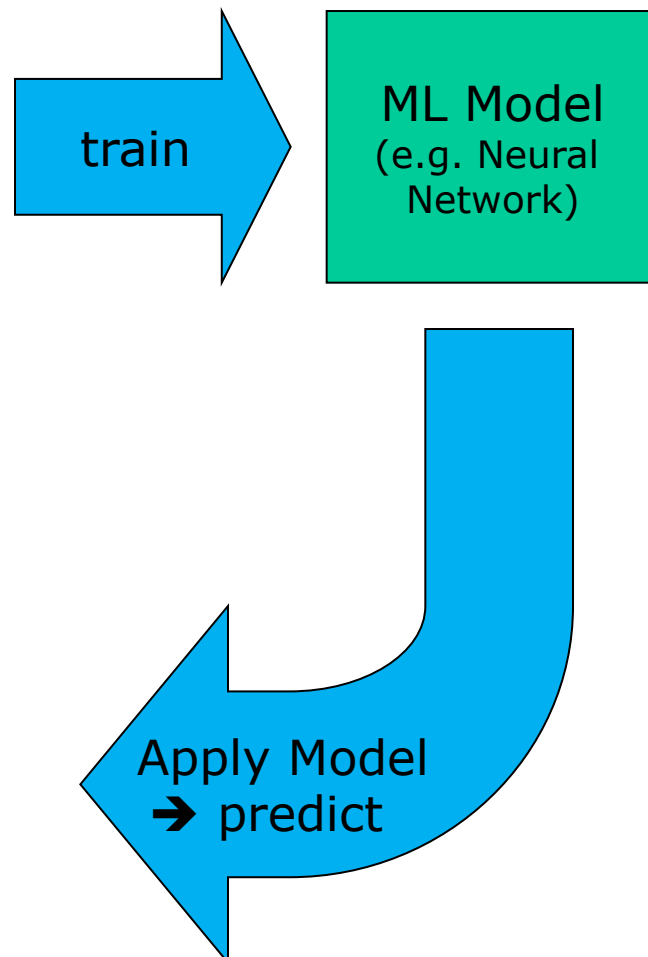
Regression: Setting

size	age	location	rooms
70	19	good	3
60	44	good	2
120	49	good	5
90	12	bad	3
80	37	bad	3
45	60	bad	2
52	44	bad	2
85	27	good	4
50	51	good	1
35	81	good	1
75	22	good	3
95	29	bad	4

price
370000
165000
575000
225000
268000
159000
185000
472000
164000
155000
391000
478000

35	77	good	4
25	19	good	1
85	44	bad	2

274000
122000
162000



gender	age	smoker	eye colour
male	19	yes	green
female	44	yes	grey
male	49	yes	blue
male	12	no	brown
female	37	no	brown
female	60	no	brown
male	44	no	blue
female	27	yes	brown
female	51	yes	green
female	81	yes	grey
male	22	yes	brown
male	29	no	blue

lung cancer
No
yes
yes
no
no
yes
no
no
yes
no
no
no

- Where do we get the data from?*

male	77	yes	grey
male	19	yes	green
female	44	no	grey

?
?
?

gender	age	smoker	eye colour
male	19	yes	green
female	44	yes	grey
male	49	yes	blue
male	12	no	brown
female	37	no	brown
female	60	no	brown
male	44	no	blue
female	27	yes	brown
female	51	yes	green
female	81	yes	grey
male	22	yes	brown
male	29	no	blue

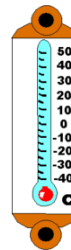
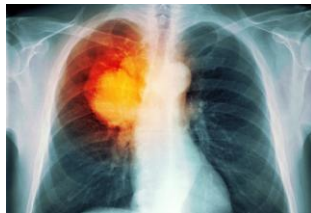
male	77	yes	grey
male	19	yes	green
female	44	no	grey

lung cancer
No
yes
yes
no
no
yes
no
no
yes
no
no
no

?
?
?

- Features:
 - Observation / measurement
 - *Extraction from media*
 - *(Transformation)*
- Labels (groundtruth):
 - Human experts
 - (Observation / measurement)

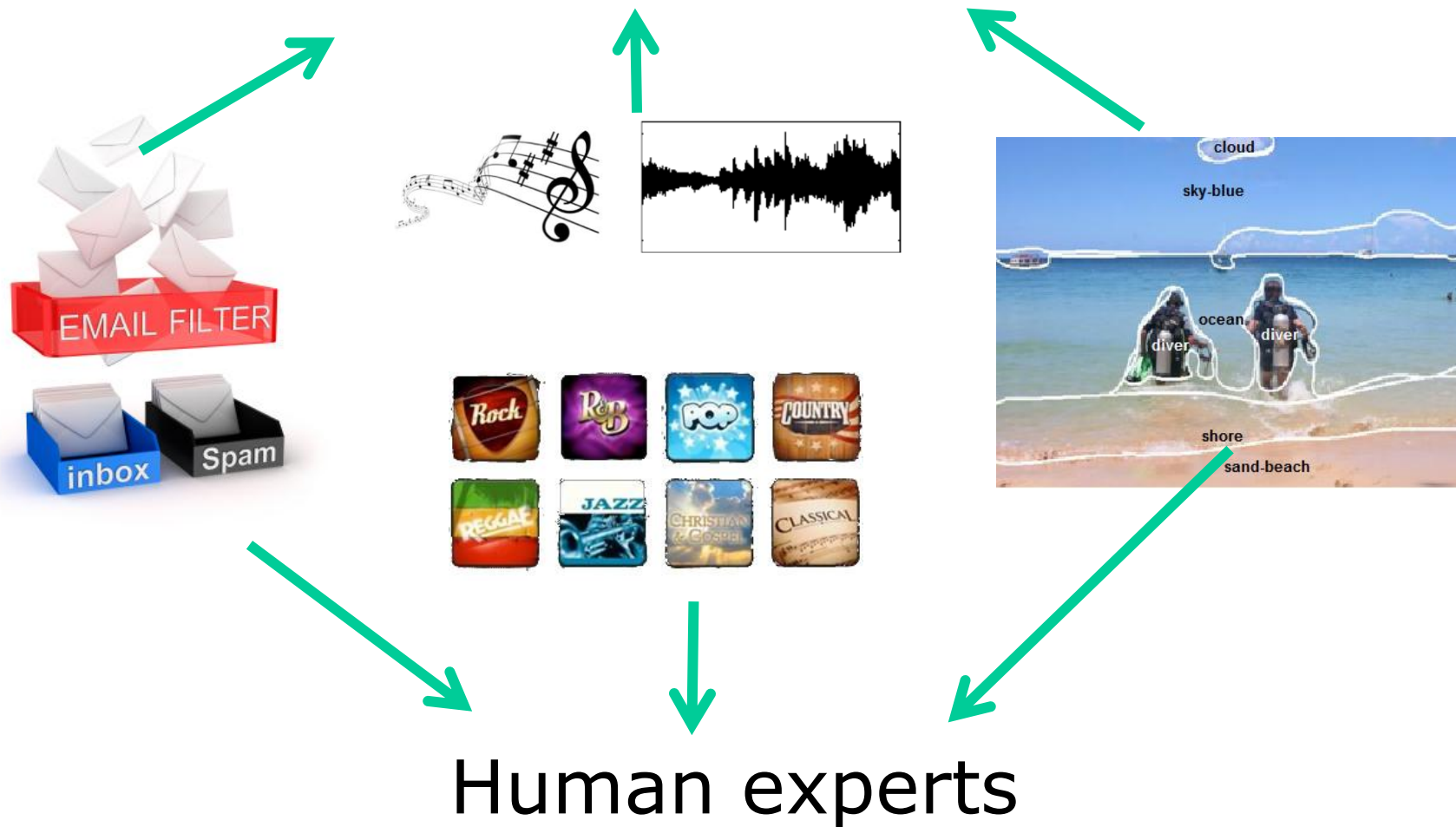
Observation / Measurements



Annotations
(human experts)

Measurement/observation

Feature Extraction



- Feature Extraction: description of complex content by derived values
 - Text: Bag of Words – counting term occurrences
 - Music: counting activity on various frequencies
 - Image: colour histograms, edge detection, “bag of visual words”,
- Important aspect of many data mining applications
 - Overview/intro in this course
 - *More details e.g. Information retrieval (188.412), ...*

.....

	Word 1	Word 2	Word 3	Word 4	Word 5	...	Word n	Σ
Doc 1	1	0	0	2	3		0	6
Doc 2	2	0	0	0	2		2	6
Doc 3	1	3	0	0	1		5	10
Doc 4	2	0	0	2	0		2	6
Doc 5	5	4	0	0	1		0	10
...							0	
Doc m	1	2	1	3	0		0	7

Text Feature
Extraction



Music Feature
Extraction

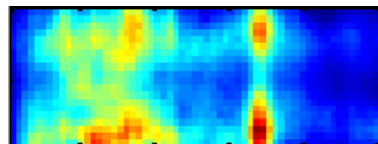
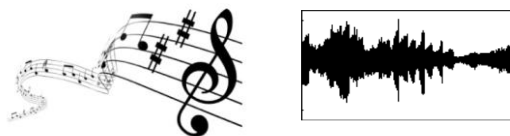
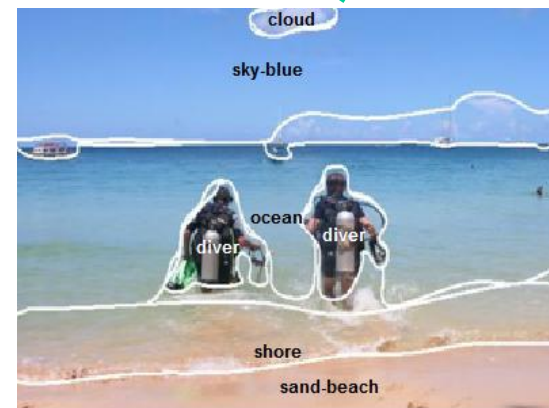
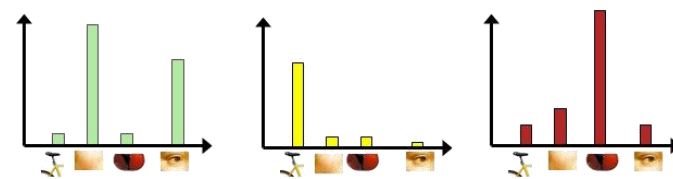
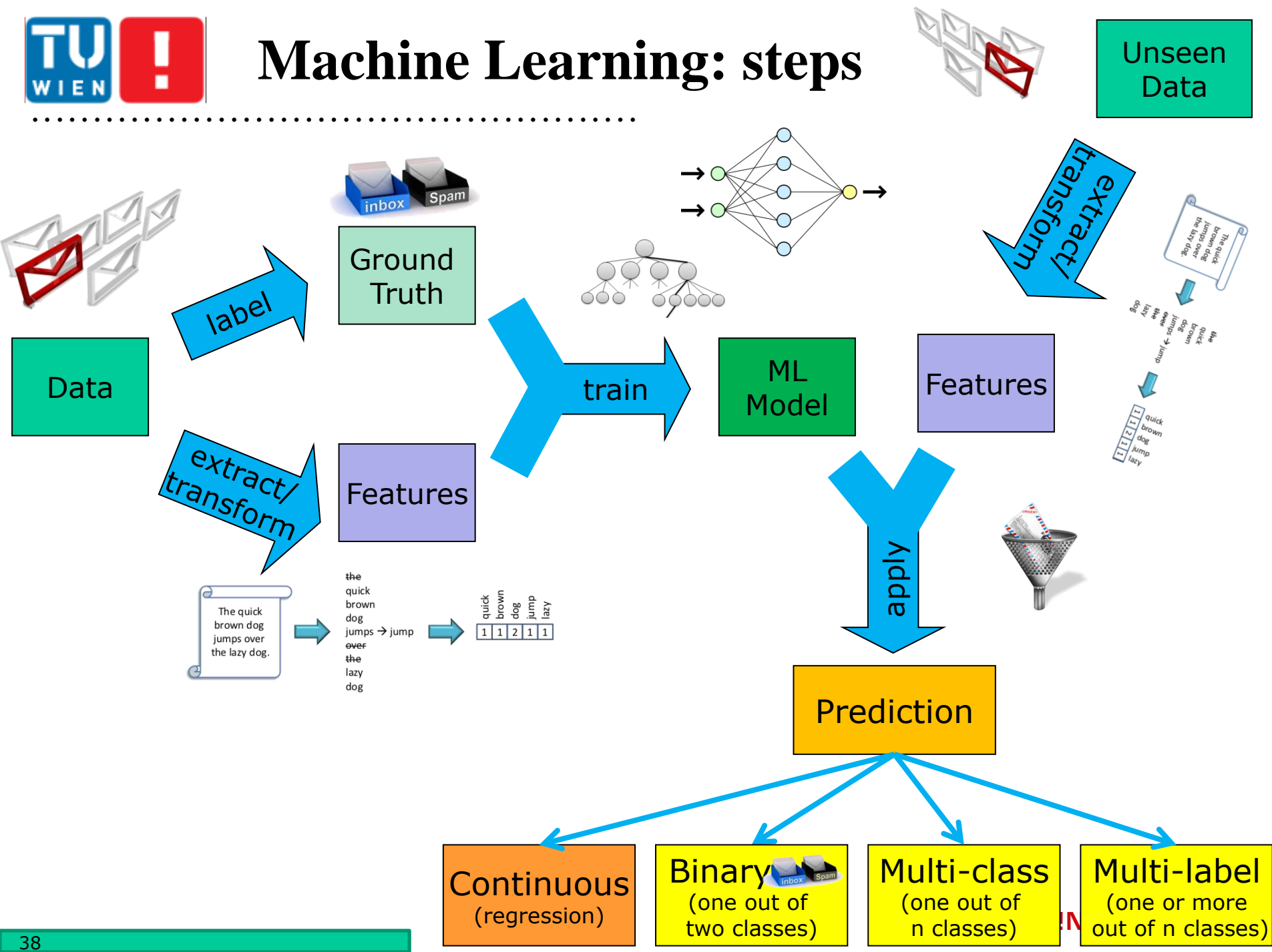


Image Feature
Extraction



Machine Learning: steps



- Feature Extraction: description of complex content by derived values
 - Text: Bag of Words – counting term occurrences
 - Music: counting activity on various frequencies
 - Image: colour histograms, edge detection,
- Important aspect of many ML applications
 - “Previews” in this lecture
 - More detailed, e.g. 188.412 “Information Retrieval”

- Algorithm can query the user to obtain labels for (unlabelled) specific samples (data points)
 - Often “difficult” samples, where the algorithm “is not sure” on the output variable
- **Motivation?**
 - Labelling data is expensive...
 - Goal: obtain a good model with less effort for labelling data
 - E.g. by selecting those samples that are difficult to decide for the algorithm (least certainty)
- Semi-Supervised learning

Independent of exact learning (data mining) method:

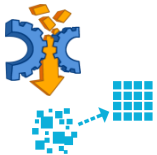
1. Identify problem / question



2. Identify & capture the available data



3. Prepare your data: clean & transform



4. Analyse your data → actual machine learning



5. Create report with results, visualisation, insights



6. Embed results into business / decision making



7. Plan for better data capturing in the future



1. Identify problem / question

- What is the problem
- Why does it need to be solved
- How can it be solved **(as machine learning problem!)**



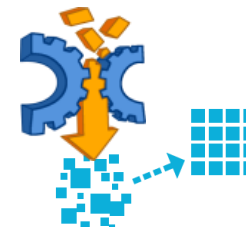
2. Identify & capture data

- Select your data
- Capture / extract your data



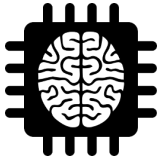
3. Clean & transform your data

- Deal with missing values
- Transform data: scale/normalise, attribute selection, ...



4. Analyse your data

- Select fitting algorithms
- Train and evaluate models, select best performing
- Improve your results – parameter tuning, ...



5. Document & analyse your results

- Problem definition, solutions, approach
- Present & compare your results in graphs, tables, ...
- Discoveries made during investigation
- Methods that did / did not work
- Limitations: when does the model not work? What questions can not be answered?





6. Embed results into business / decision making
 - Directly adjust your operation depending on the result
 - Use results for decision making
 - At various levels of the organisation

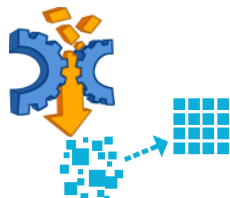
7. Plan for better data capturing in the future
 - Capture more data than currently available
 - Capture data in higher resolution (more detail, more frequent, ...)
 - Capture data in better quality



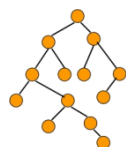
Scoping the lecture



Data types & feature extraction



Data preparation & transformation



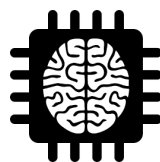
Decision Trees & Random Forests

Support Vector Machines

Neural Nets & Deep Learning

Ensemble learning

...



Evaluation

Model selection

Significance testing

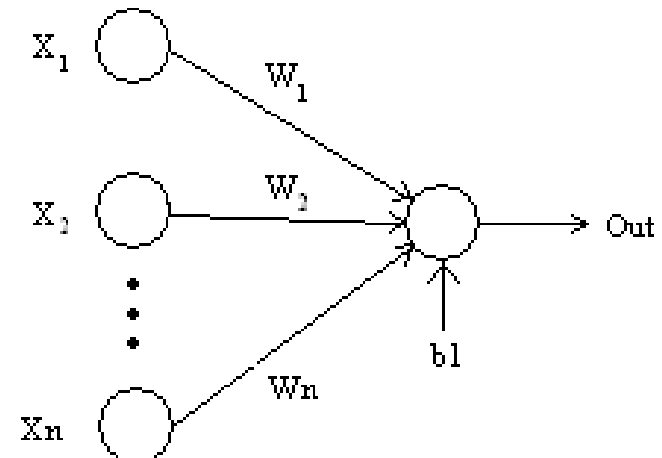


- Matching your expectations?
 - Anything not interesting at all?
 - Anything missing?



- Machine Learning: Intro
- Perceptron
- k-NN
- Performance evaluation (intro)
- Types of data & data preparation

- Proposed in 1957
- Artificial (neural) network
 - only one neuron; “single-layer perceptron”
- Input: continuous values $X (x_1, x_2, \dots, x_n)$
- Output: activation a
 - Boolean value 0 / 1



- Linear combination of inputs
 - Using weights W

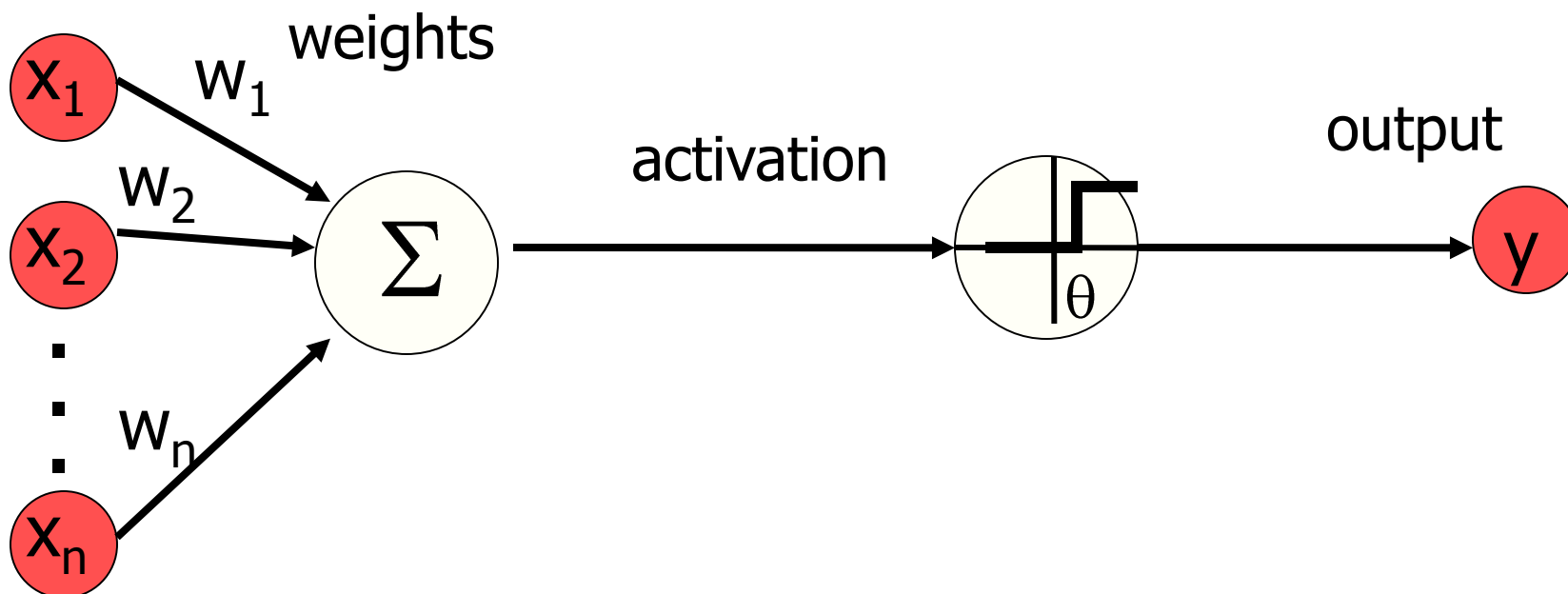
$$a = \sum_{i=1}^n w_i x_i$$

- Pass through threshold activation function with threshold θ

$$y = f(x) = \begin{cases} 1 & \text{if } a \geq \theta \\ 0 & \text{if } a < \theta \end{cases}$$

(Heaviside step function)

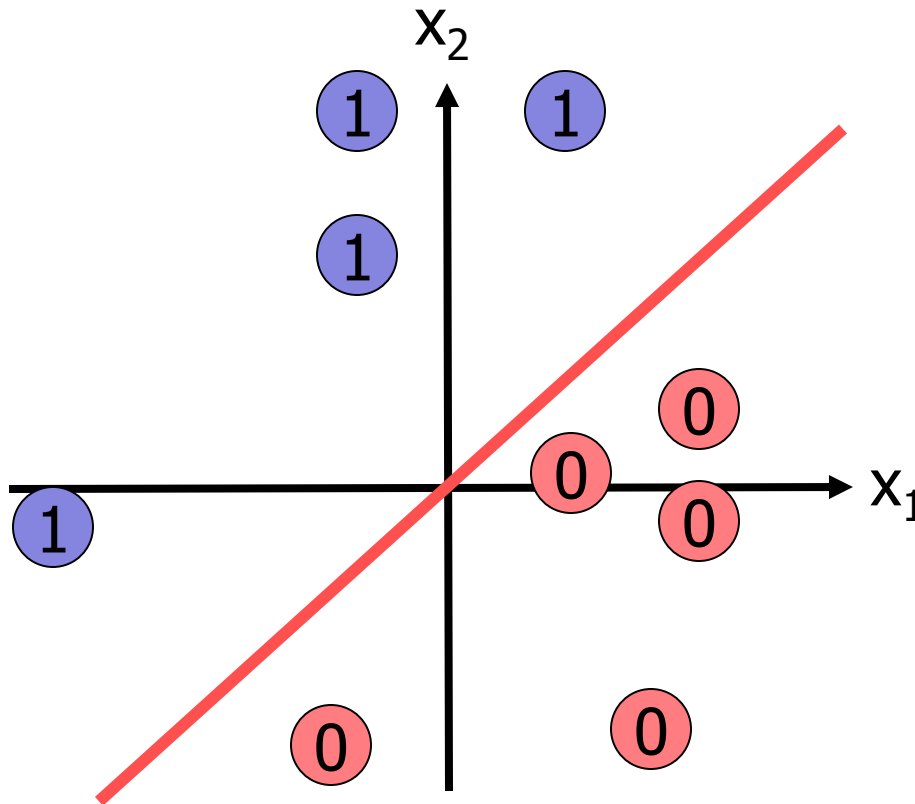
inputs



Perceptron: visual example

- Two-dimensional data set
 - Variables x_1 & x_2
 - E.g. position on a map
- Weight vector also has two components
 - w_1 & w_2
- Data set can be easily visualised in Cartesian coordinate system
- Two classes: **0** and **1**
 - E.g. **poor** and **rich** countries

Perceptron visualised



Decision line

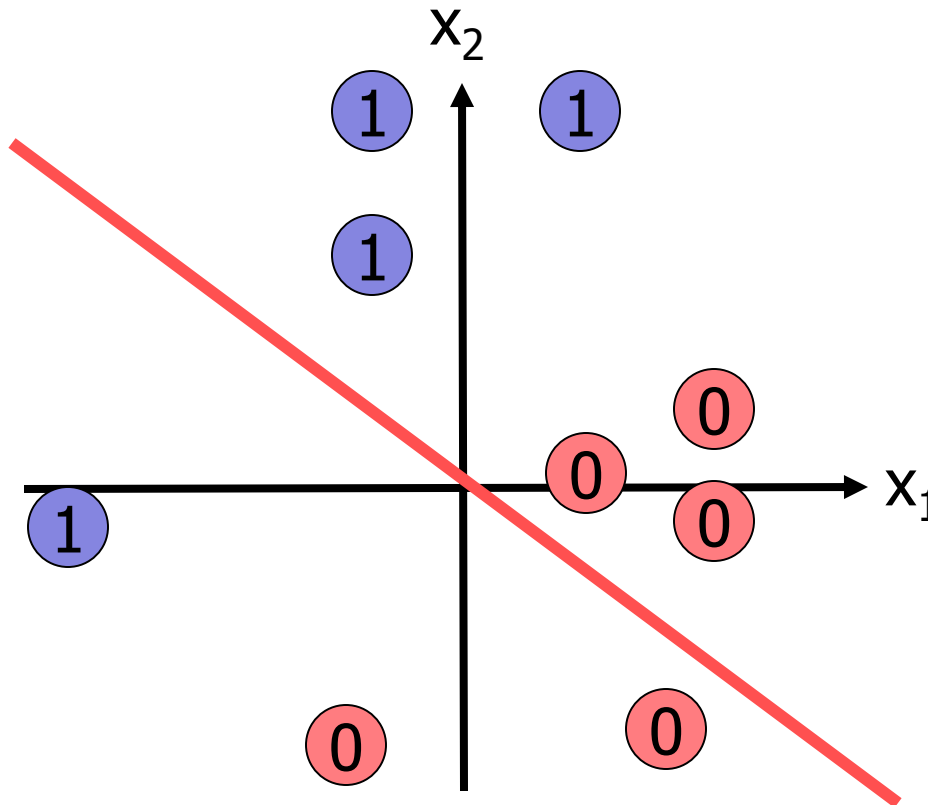
$$w_1 x_1 + w_2 x_2 = \theta$$

- Training the model: learning the weights from labelled samples (label: \mathbf{t})
 - Initialise weights
 - Repeat
 - Present training sample \mathbf{x}
 - Predict sample label: $\mathbf{y} = f(\mathbf{x})$
 - Prediction correct? Compare \mathbf{t} and \mathbf{y}
 - if $\mathbf{y} \neq \mathbf{t} \rightarrow$ Compute new weights \mathbf{w}' as $\mathbf{w}' = \mathbf{w} + \alpha (\mathbf{t} - \mathbf{y}) \mathbf{x}$
 - Until prediction correct ($\mathbf{y} = \mathbf{t}$) for all samples

- Parameter α : *learning rate*
 - determines the magnitude of weight updates
- If output is correct ($t=y$)
 - weights **not** changed
- If output is incorrect ($t \neq y$)
 - weights changed such that the output of the for the new weights \mathbf{w}' is **closer** to the input x_i

Perceptron training

- $t=0$: random weight



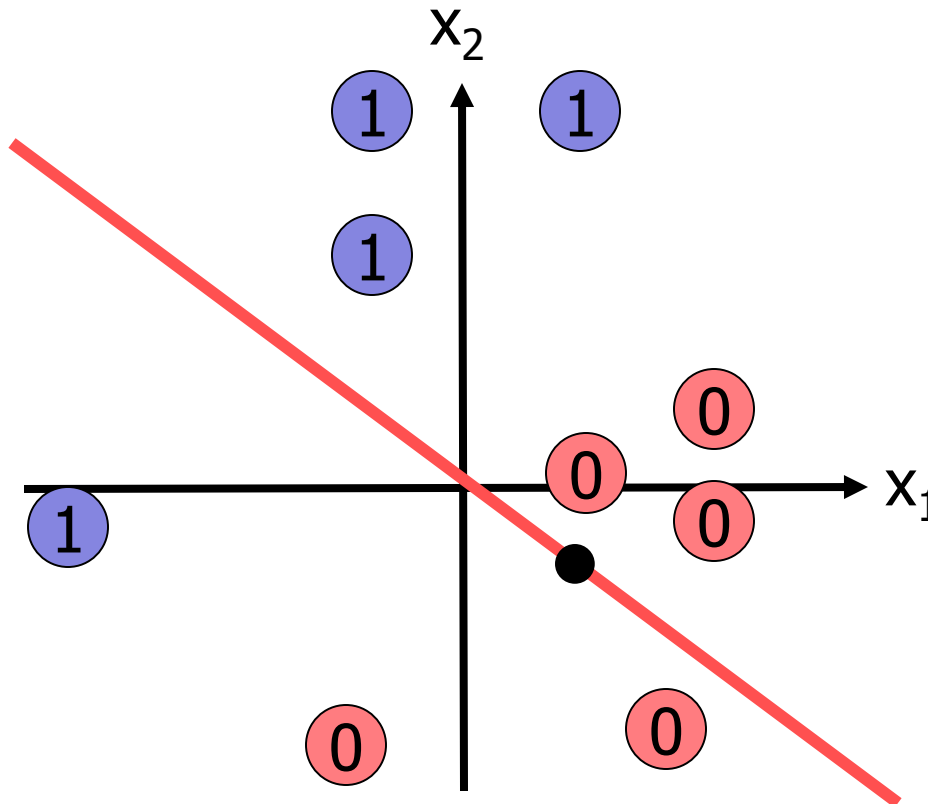
Decision line

$$w_1 x_1 + w_2 x_2 = \theta$$

(cf. $y = ax + b$)

Perceptron training

- $t=0$: random weight



Decision line

$$w_1 x_1 + w_2 x_2 = \theta$$

Initial weight vector:

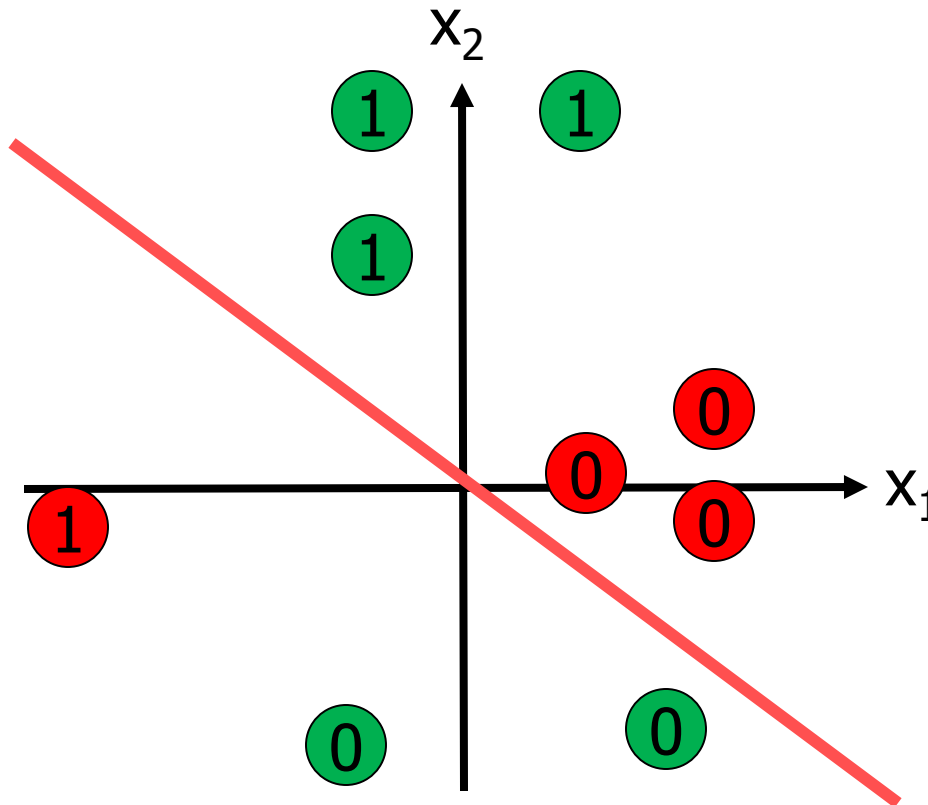
$$w_1=1, w_2=1$$

Point on decision line
(1, -1)

$$\begin{aligned} \rightarrow 1*1 + 1*(-1) &= \\ &= 1 - 1 = 0 = \theta \end{aligned}$$

Perceptron training

- $t=1$: compute $y = f(x)$, compare to t



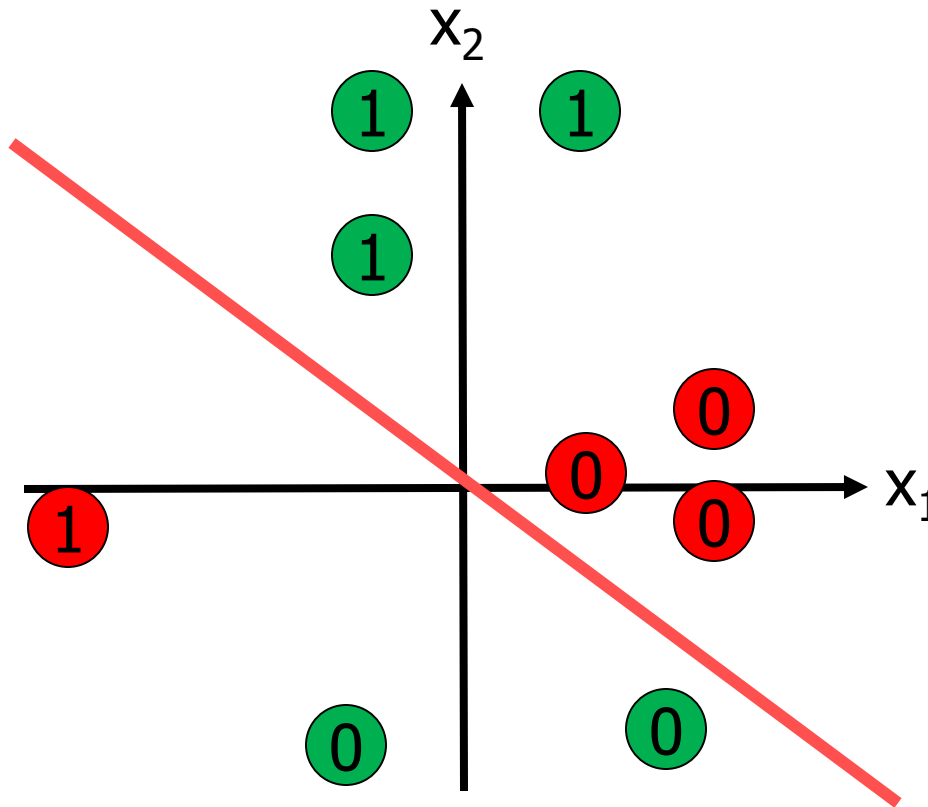
Decision line

$$w_1 x_1 + w_2 x_2 = \theta$$

(*green*: correctly , *red*: wrongly classified)

Perceptron training

- $t=1$: not all $f(x) = t \rightarrow$ adapt weights (decision line)



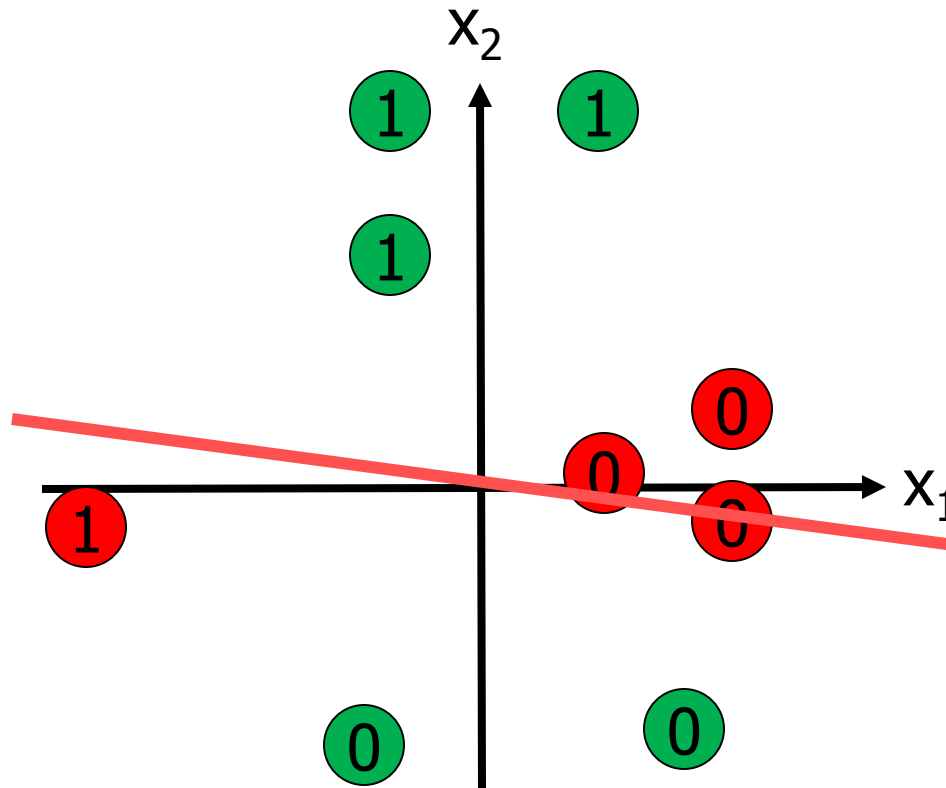
Decision line

$$w_1 x_1 + w_2 x_2 = \theta$$

(*green*: correctly , *red*: wrongly classified)

Perceptron training

- $t=1$: not all $f(x) = t \rightarrow$ adapt weights (decision line)



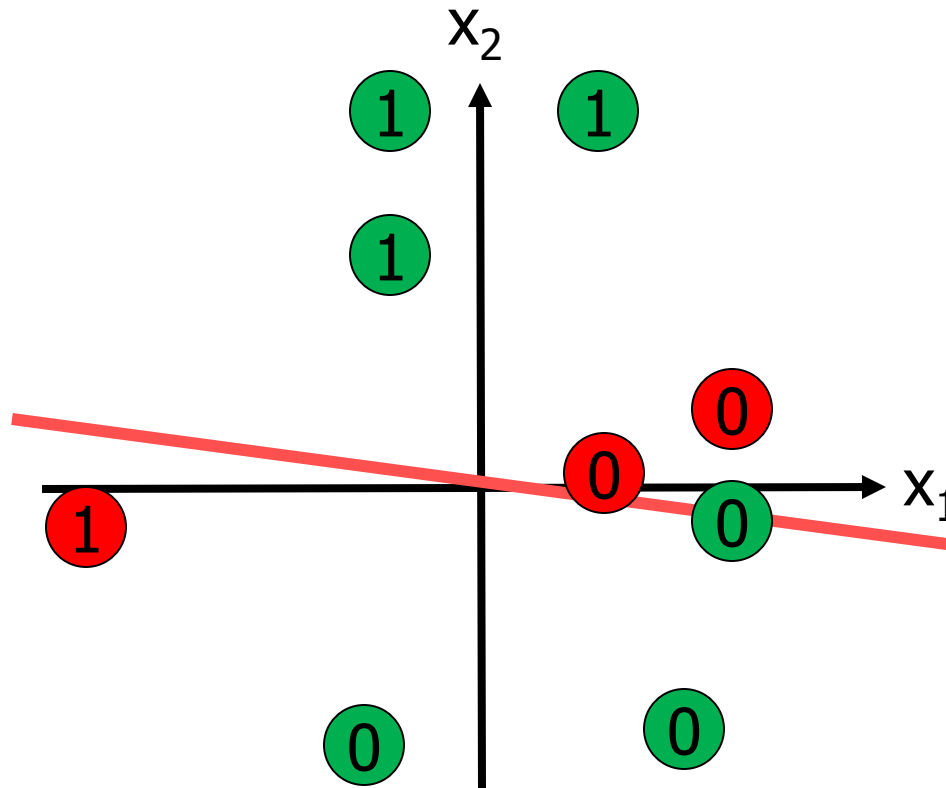
Decision line

$$w_1 x_1 + w_2 x_2 = \theta$$

(*green*: correctly , *red*: wrongly classified)

Perceptron training

- $t=2$: not all $f(x) = t \rightarrow$ adapt weights (decision line)



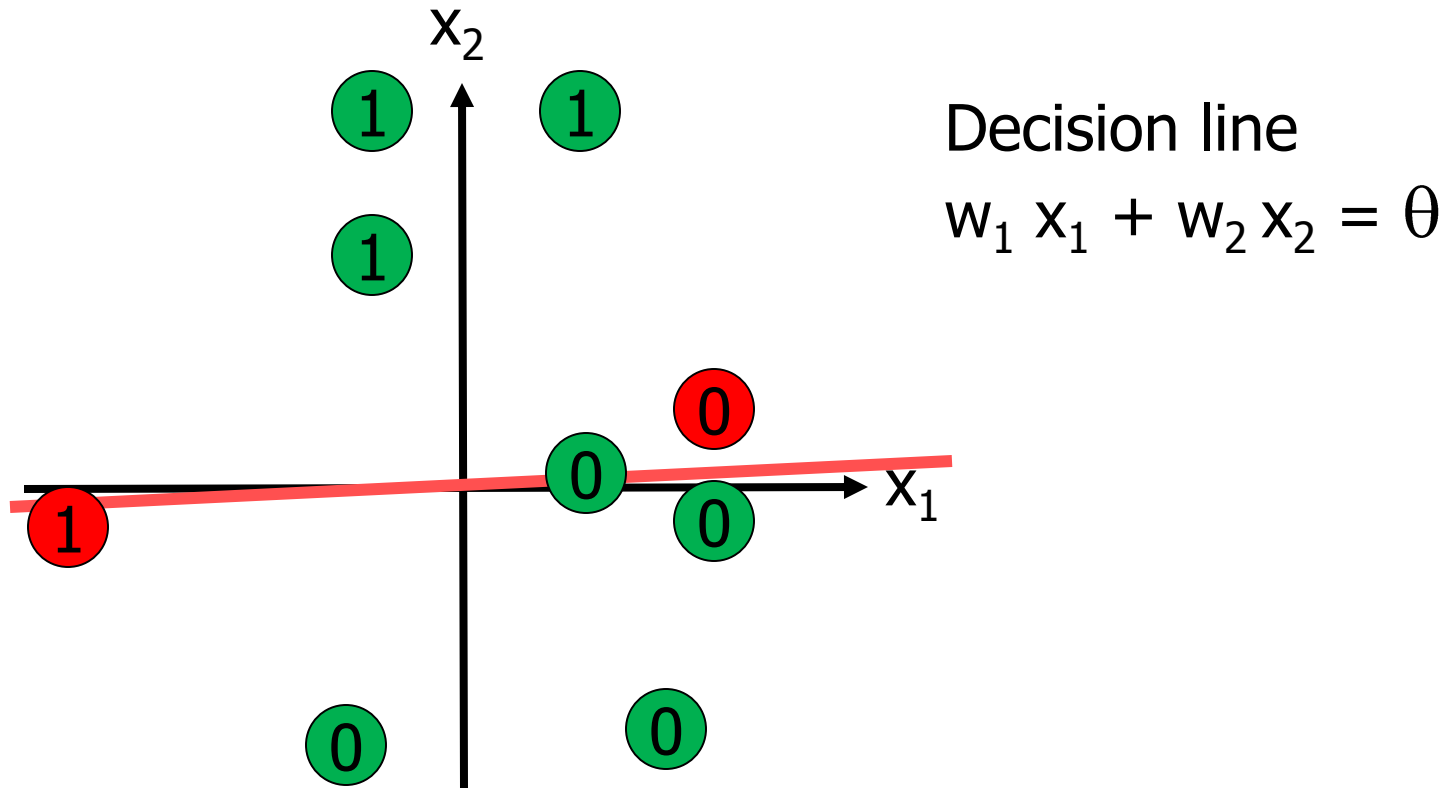
Decision line

$$w_1 x_1 + w_2 x_2 = \theta$$

(*green*: correctly , *red*: wrongly classified)

Perceptron training

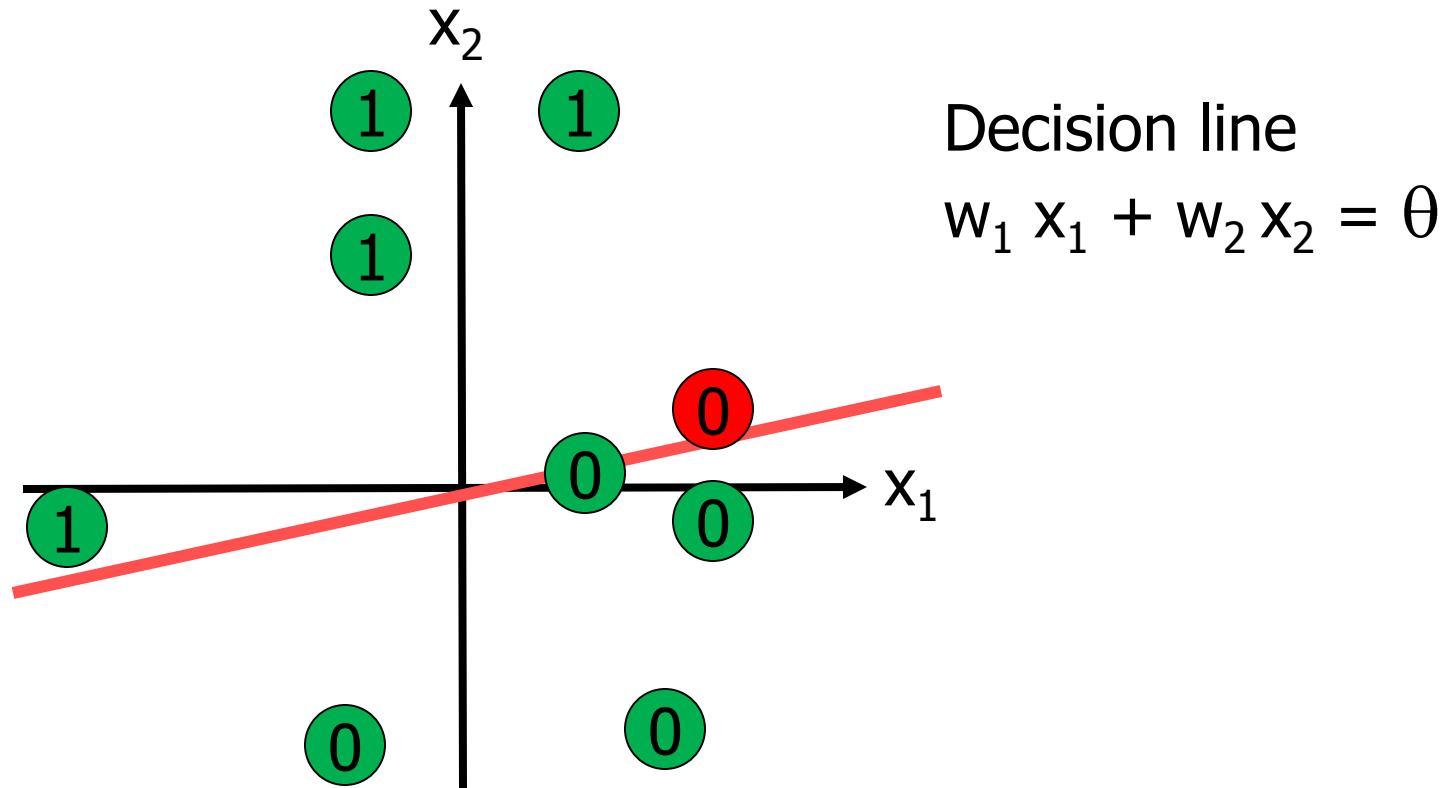
- $t=3$: not all $f(x) = t \rightarrow$ adapt weights (decision line)



(green: correctly , red: wrongly classified)

Perceptron training

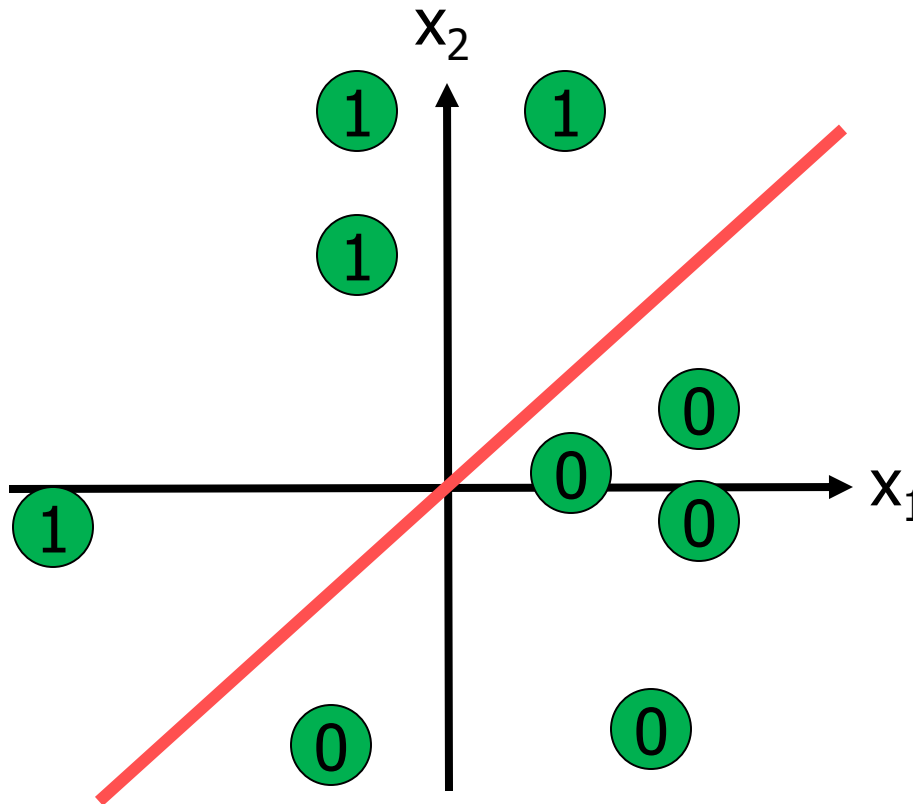
- $t=4$: not all $f(x) = t \rightarrow$ adapt weights (decision line)



(*green*: correctly , *red*: wrongly classified)

Perceptron training

- $t=5$: all $f(x) = t \rightarrow$ final state



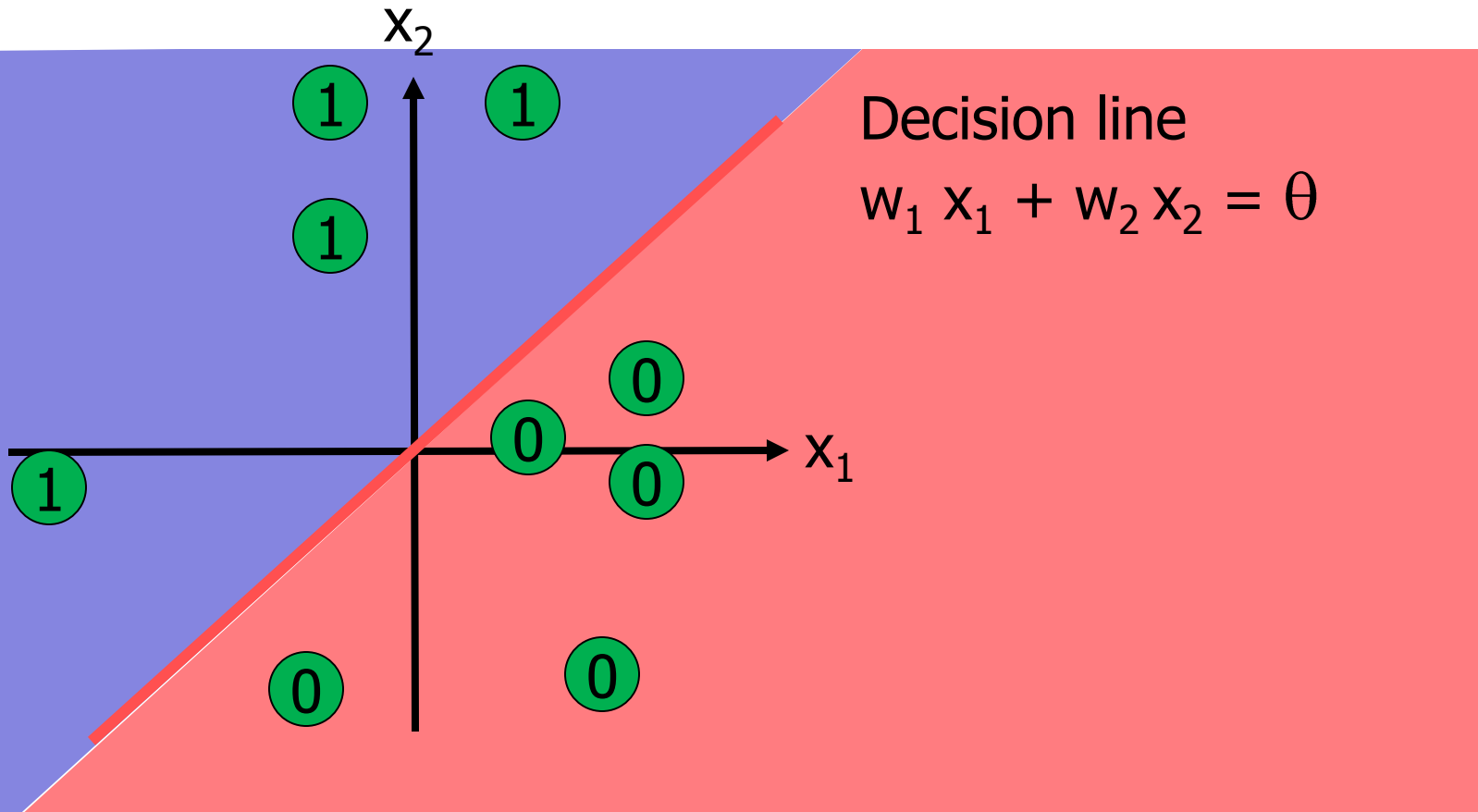
Decision line

$$w_1 x_1 + w_2 x_2 = \theta$$

(*green*: correctly , *red*: wrongly classified)

Perceptron training

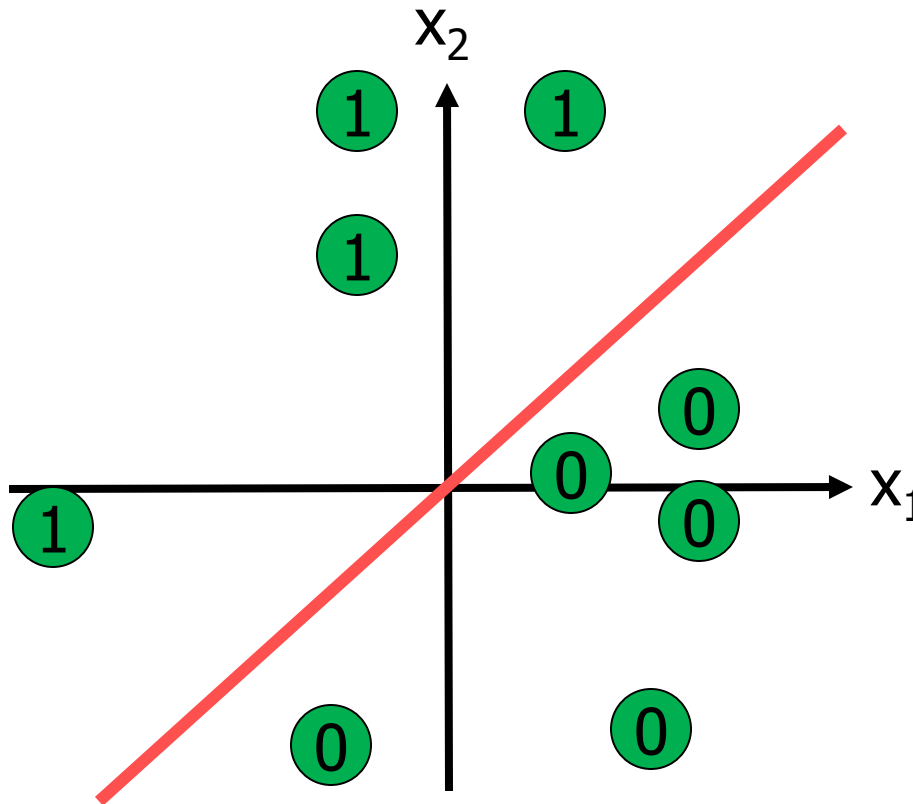
- $t=5$: all $f(x) = t \rightarrow$ final state; classification in **rich/poor**



(*green*: correctly , *red*: wrongly classified)

Perceptron training

- $t=5$: all $f(x) = t \rightarrow$ final state



Decision line

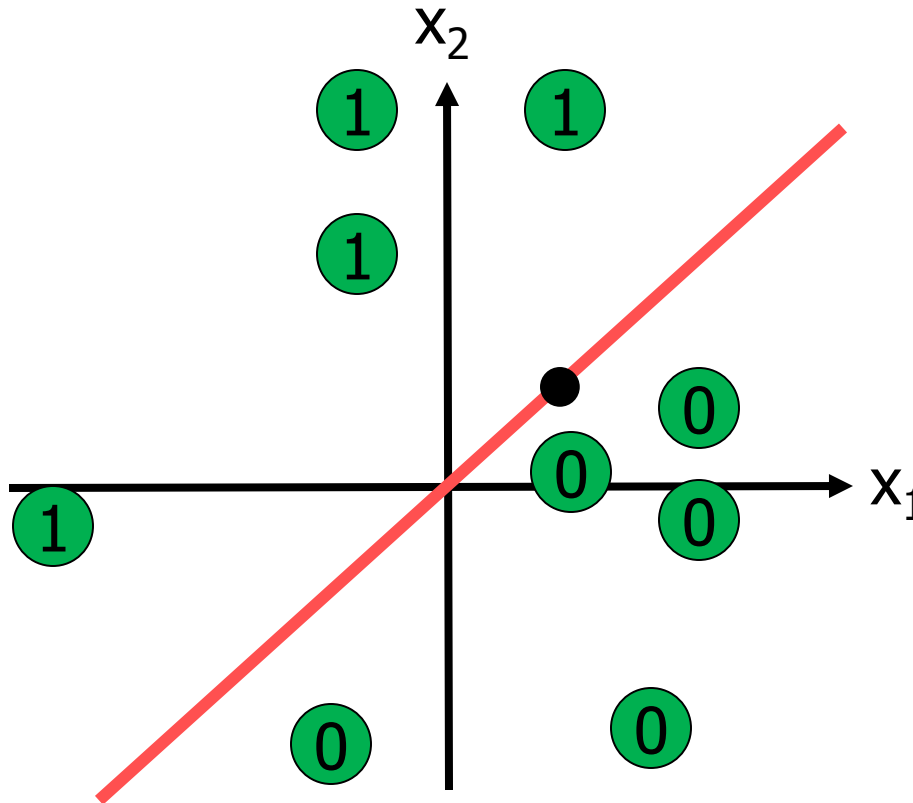
$$w_1 x_1 + w_2 x_2 = \theta$$

Weight vector ?

(*green*: correctly , *red*: wrongly classified)

Perceptron training

- $t=5$: all $f(x) = t \rightarrow$ final state



Decision line

$$w_1 x_1 + w_2 x_2 = \theta$$

Weight vector ?

$$w_1 = -1, w_2 = 1$$

Point on decision line

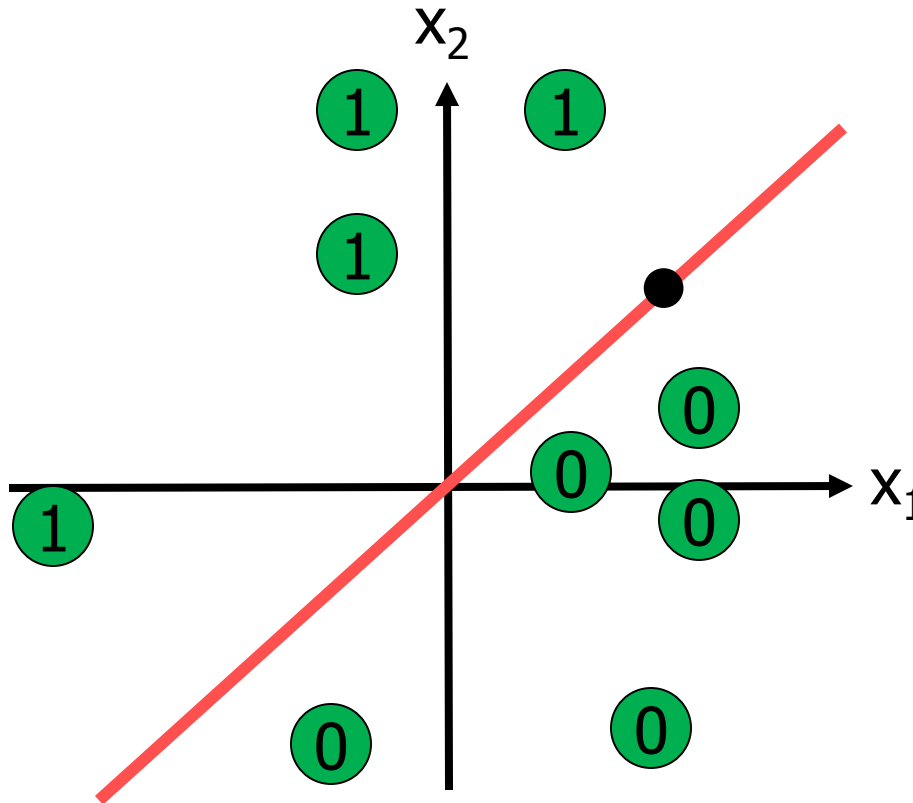
$$(1, 1)$$

$$\begin{aligned} \rightarrow & (-1) * 1 + 1 * 1 = \\ & = -1 + 1 = 0 = \theta \end{aligned}$$

(green: correctly , red: wrongly classified)

Perceptron training

- $t=5$: all $f(x) = t \rightarrow$ final state



Decision line

$$w_1 x_1 + w_2 x_2 = \theta$$

Weight vector ?

$$w_1 = -1, w_2 = 1$$

Point on decision line

(2,2)

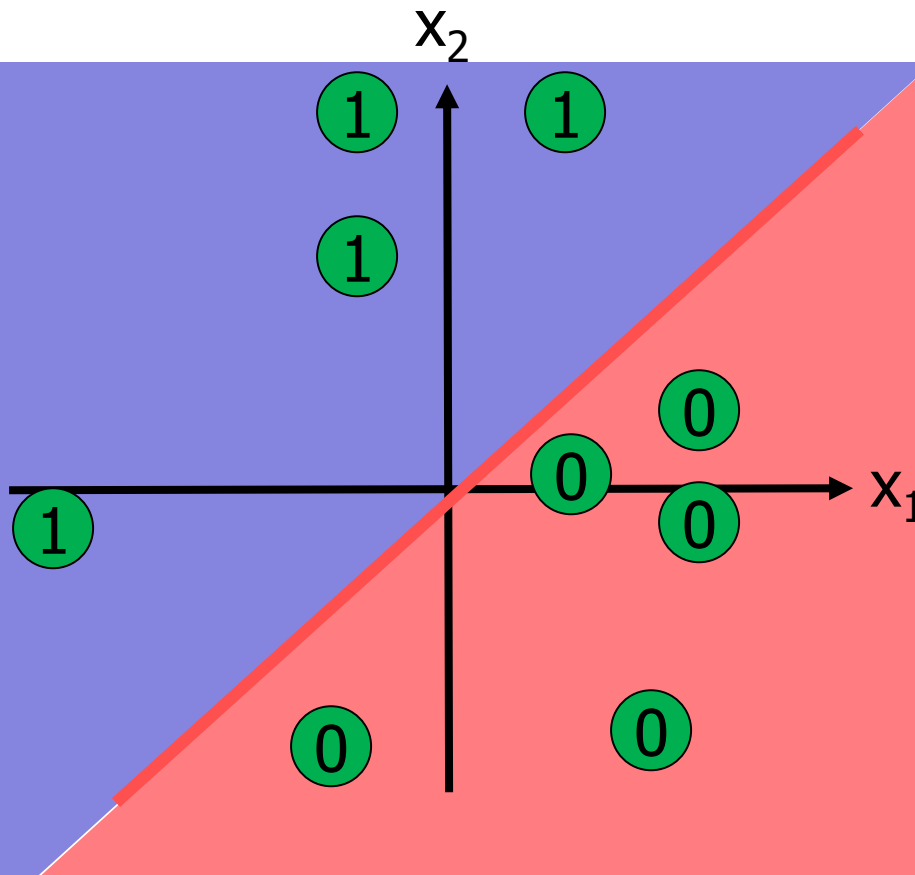
$$\begin{aligned} \rightarrow & (-1) \cdot 2 + 1 \cdot 2 = \\ & = -2 + 2 = 0 = \theta \end{aligned}$$

(green: correctly , red: wrongly classified)

- Separates linearly
 - Converges to a stable state when data is *linear separable*

Perceptron: properties

- Separates linearly
 - Converges to a stable state when data is *linear separable*



- Separates linearly
 - Converges to a stable state when data is *linear separable*
- Can predict binary decisions (true/false)
 - Can be extended for multi-class problems e.g. with different activation functions
 - *Training rule similar to other online learning algorithms, e.g. Self-Organising Maps*
- *More details next lecture*

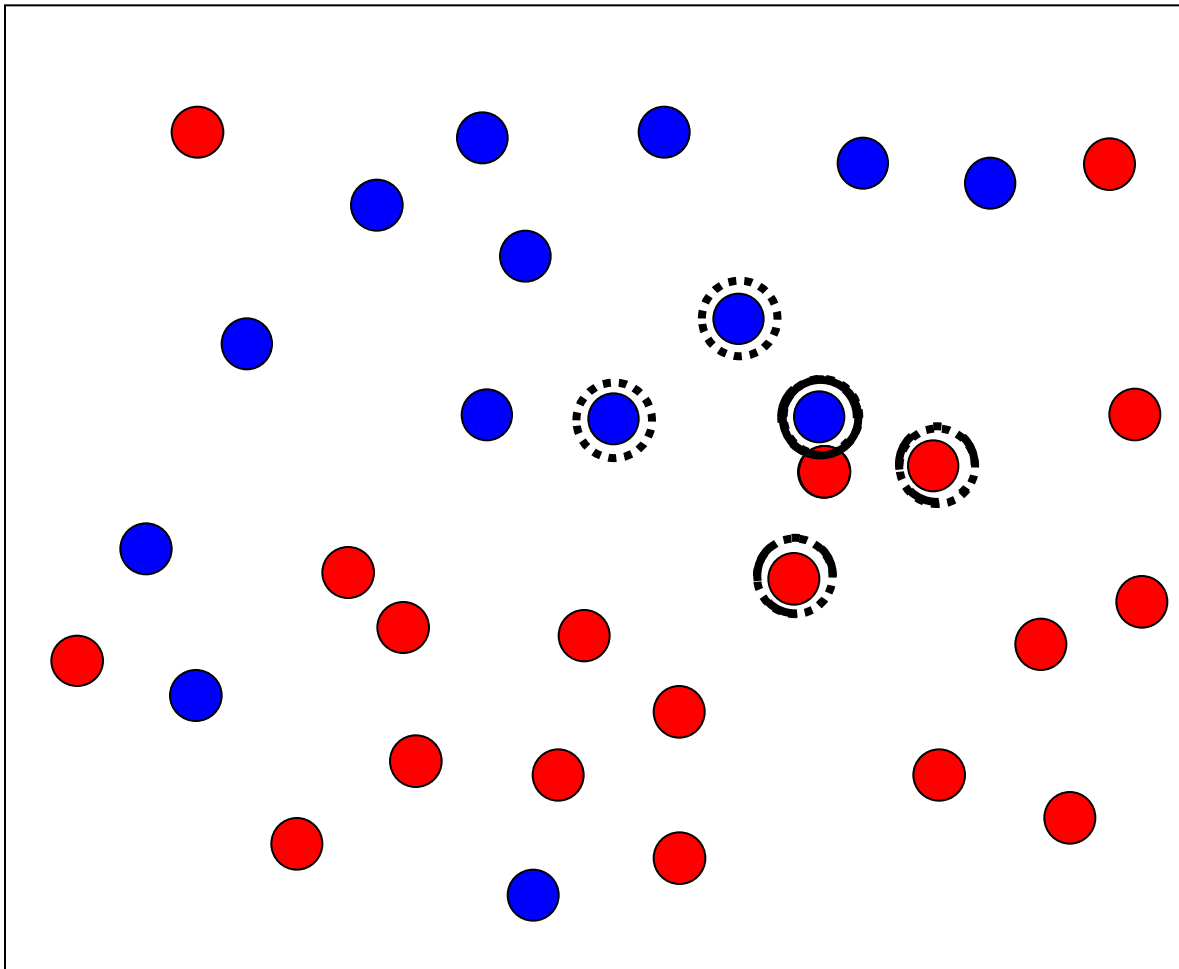
- <http://ditam.github.io/demos/perceptron/perceptronDemo.html>
- <https://www.cs.utexas.edu/~teammco/misc/perceptron/>

- Machine Learning: Intro
- Perceptron
- **k-NN**
- Performance evaluation (intro)
- Types of data & data preparation

- Simple algorithm, but well known
- Classify inputs based on k closest training examples
 - k to be chosen, high influence
 - Definition of “closest” – distance function

k-nn: Example

- 2-dimensional data, two classes (red & blue)



$k = 1$
 $k = 3$
 $k = 5$

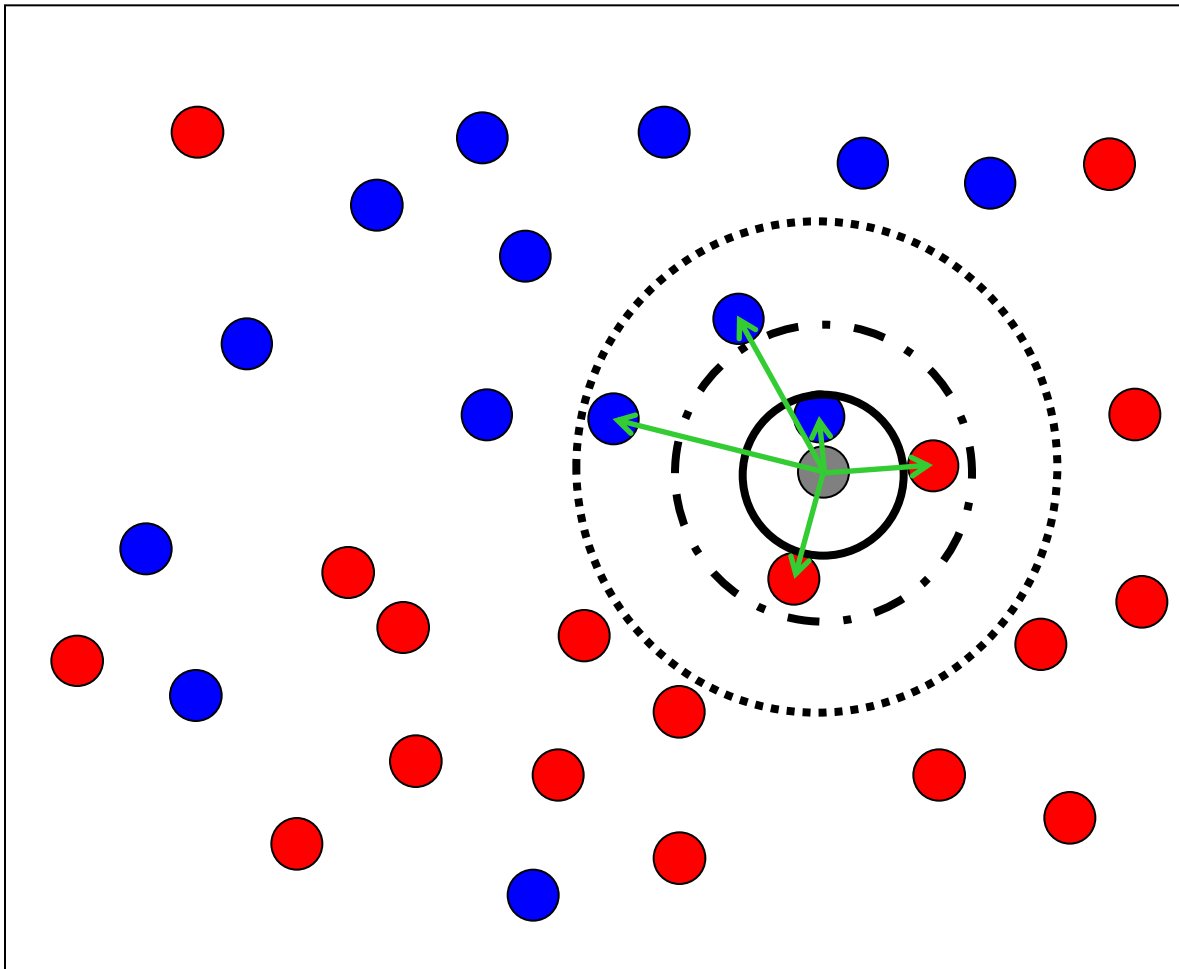
- Very sensitive to local noise
 - Little abstraction / generalisation
- Larger values of k reduce the effect of noise
 - makes boundaries between classes less distinct
- How to determine k ?
 - Good values for k vary (a lot!) with the data
- $k=1$ is called “nearest neighbour algorithm”

- *What is the maximum value for k ?*

- *What is the maximum value for k ?*
- *What happens if $k = n$ (number of samples)*

k-nn: distance function

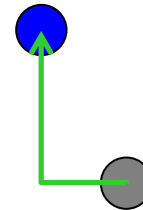
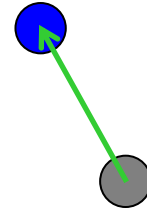
- In this example: Euclidean distance



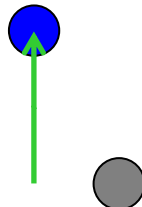
$k = 1$
 $k = 3$
 $k = 5$

k-nn: distance function

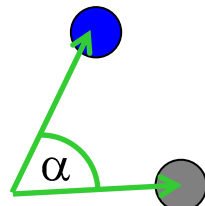
- Previous example: Euclidean distance
- Any other distance measure for numeric variables can be employed
 - Manhattan/City Block/L1
 - L_n



- Linfinity



- Cosine

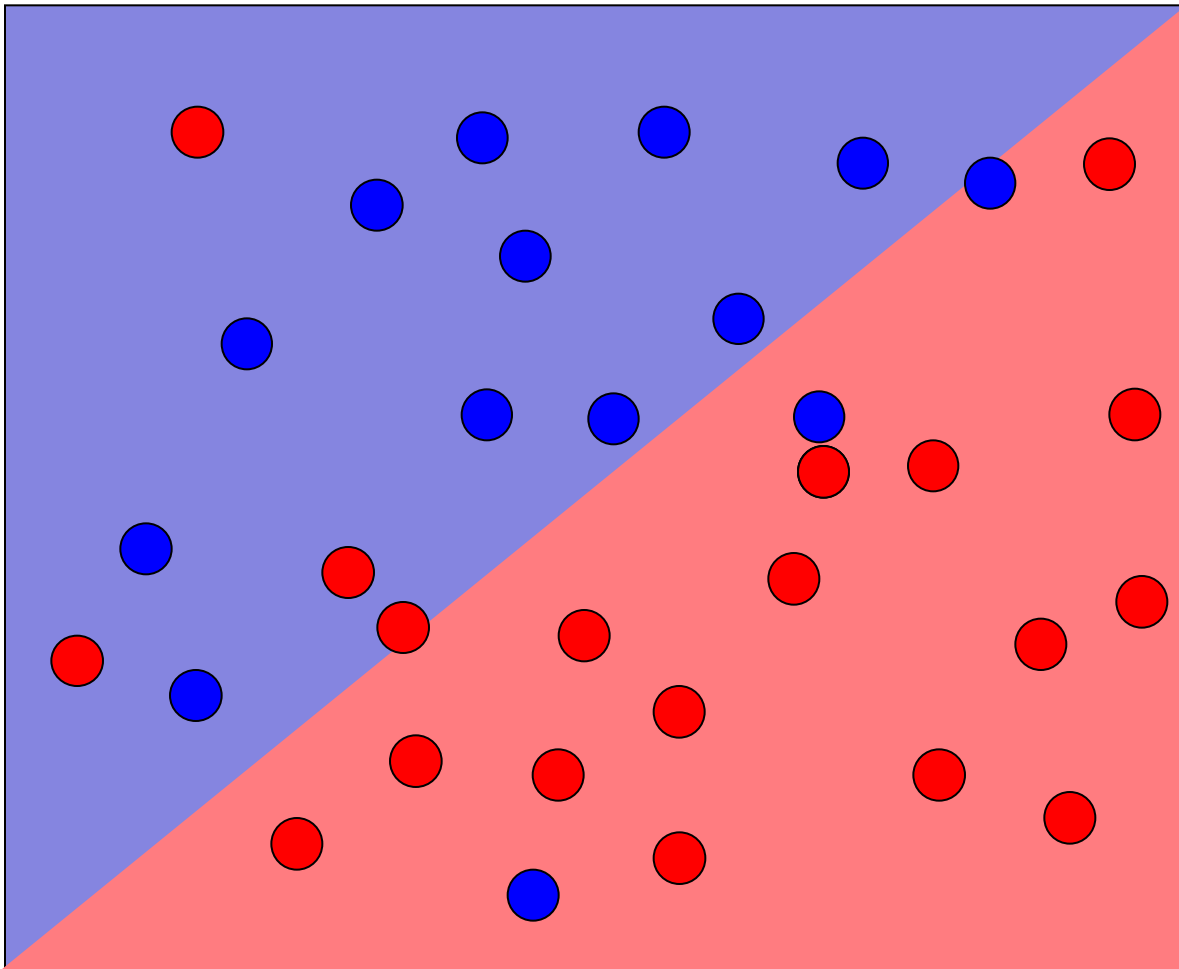


- Major differences between k-NN and perceptron algorithms?
- ***k-NN is a “Lazy learner”***
 - no model built beforehand
 - ➔ computation is done at classification step
 - Opposite is called “eager learning”

- Major differences between k-NN and perceptron algorithms?
- ***k-NN is a “Lazy learner”***
 - no model built beforehand
 - ➔ computation is done at classification step
 - Opposite is called “eager learning”
- k-NN can learn **complex** (almost arbitrary) **decision boundaries**

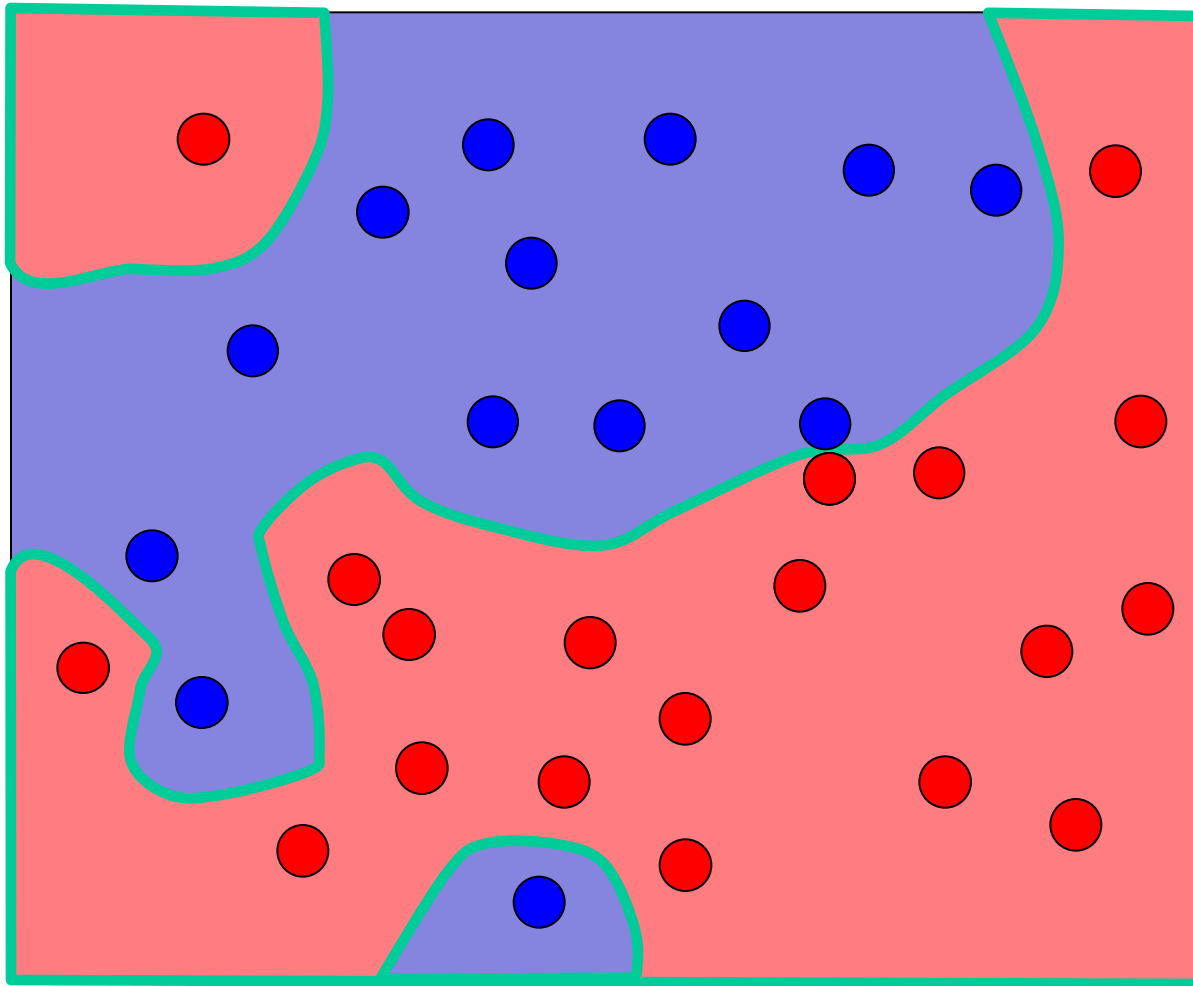
Perceptron: decision boundary

- *Only linear separation!*



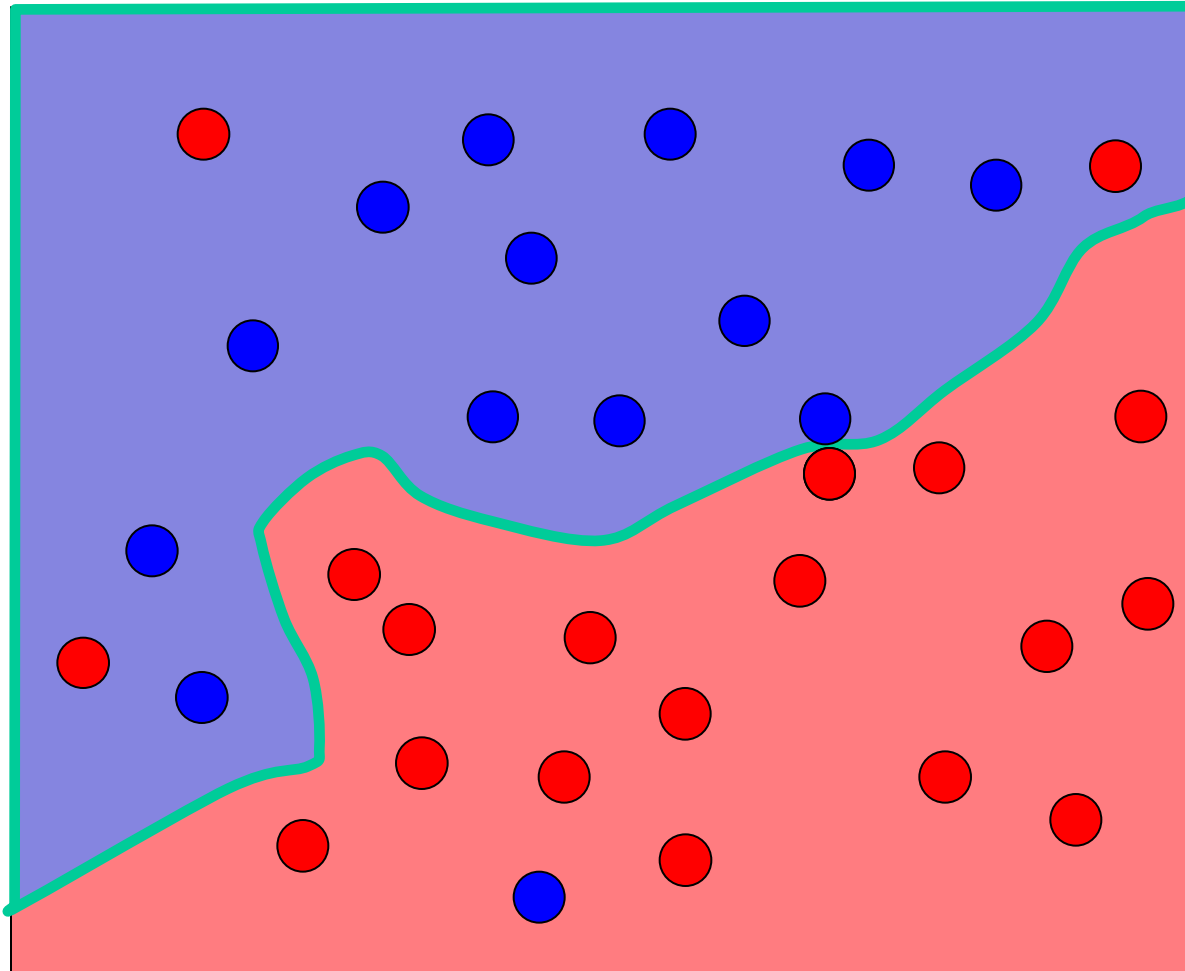
k-nn: decision boundary

- *(Almost) Arbitrary boundary*



$k = 1$

k-nn: decision boundary



$k = 3$

- Easy to understand and interpret
- Easy to implement
- No training time needed...

- Memory requirements depend on training data
- Becomes computationally expensive with many items to classify
 - Linear search: $O(Nd)$, N = # samples, d = dimension
 - Several optimisations proposed

- *More details next lecture*

- Machine Learning: Intro
- Perceptron
- k-NN
- Performance evaluation (intro)
- Types of data & data preparation

- Need to measure performance of an algorithm
 - To select the best algorithm
 - To estimate usefulness of models
 - ...
- Approach:
 - Train model on (labelled) training data
 - Test on (labelled) test data
- Measure performance
 - Several different measures ...

- Binary classification, classes true/false
- 4 possible outcomes prediction / groundtruth

		<i>Actual value (groundtruth)</i>	
		true	false
<i>Prediction (test outcome)</i>	true	True positive (TP)	False positive (FP, Type I error)
	false	False negative (FN, Type II error)	True negative (TN)

	true	false
true	True positive (TP)	False positive (FP, Type I error)
false	False negative (FN, Type II error)	True negative (TN)

- Accuracy: # correctly predicted samples

$$\frac{TP + TN}{TP + FP + TN + FN} = \frac{TP + TN}{\# \text{ samples}}$$

- Precision $\frac{TP}{TP + FP}$ ➔ How to optimise?

- Recall $\frac{TP}{TP + FN}$ ➔ How to optimise?

	true	false
true	True positive (TP)	False positive (FP, Type I error)
false	False negative (FN, Type II error)	True negative (TN)

- F-Measure: trade-off between precision and recall
- F1 Measure: equal weights, *harmonic mean*

$$\frac{2 * (precision * recall)}{(precision + recall)}$$

- Generally:
$$\frac{(1 + \beta^2) * (precision * recall)}{(\beta^2 * precision + recall)}$$

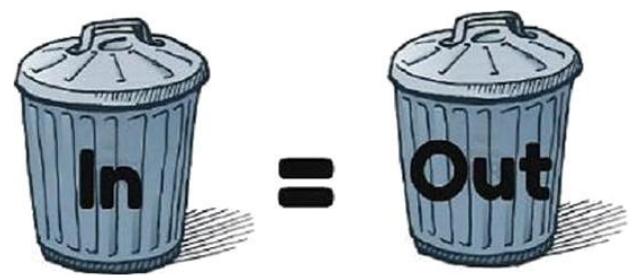
- Where to do evaluation on?
- The samples used for training?
 - *Why not?*
 - Would not tell us how good the model works for unknown data
 - Which is however why we train a model in first place ...
 - Will often achieve better (or perfect) results on training data
 - *Why?*

- Split data into training and test sets
 - E.g. ~80% training, 20% test, 66% - 33%
 - Linear, random, ...
- Results can vary a lot according to how the split is done
 - ➔ Cross validation (*discussed later*)

- Micro vs. macro averaging
 - Confusion matrix
 - Cost functions
- ROC curves
- Cross-validation & Bootstrapping
- Significance testing

- Evaluation measures for regression

- Machine Learning: Intro
- Perceptron
- k-NN
- Performance evaluation (intro)
- Types of data & data preparation

- Vital step for machine learning (supervised and unsupervised)
- ML algorithm will ***always*** give you a model
 - Quality of that model depends highly on the quality of the input data
- “Garbage in” → “Garbage out”
 
- One major goal of data preparation:
 - Eliminate “wrong influence” of variables

- Example data set from earlier (lung cancer)
- Potential issues ?
 - Missing values
 - Quantitative (continuous) data with different scales
 - Categorical data

<i>gender</i>	<i>age</i>	<i>height</i>	<i>smoker</i>	<i>eye colour</i>
male	19	170	yes	green
female	44	162	yes	grey
male	49	185	yes	blue
male	12	178	no	brown
female	37	165	no	brown
female	60	157	no	brown
male	44	19ß	no	blue
female	27	178	yes	brown
female	51	162	yes	green
female	81	168	yes	grey
male	22	184	yes	brown
male	29	176	no	blue

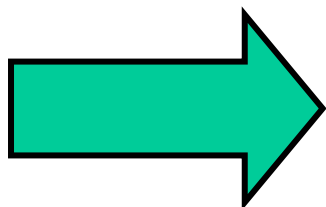
Categorical variables

.....

- Non-numeric variables with a finite number of levels
 - Also called ***nominal data***
- E.g. eye colour with values “green”, “grey”, “blue”, “brown”
- Some Machine Learning algorithms can ***only handle numeric variables !***
 - ***Which ones?***
- One solution: 1-to-N coding
 - Also called ***“one hot” (en)coding***

1-to-N Coding (one hot encoding)

<i>colour</i>
brown
blue
green
grey
brown
green
blue



<i>green</i>	<i>blue</i>	<i>brown</i>	<i>grey</i>
0	0	1	0
0	1	0	0
1	0	0	0
0	0	0	1
0	0	1	0
1	0	0	0
0	1	0	0

Categorical data: animals data set

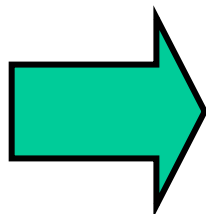
- Animal data set
 - Describes animal by some characteristics
 - Instances: cow, duck, bee, ...
- Variables
 - Size (tiny, small, medium, big)
 - Number of legs (2, 4, 6, 8)
 - Feathers (yes/no)
 - Eggs (yes/no)

	<i>size</i>	<i>legs</i>	<i>feathers</i>	<i>eggs</i>
duck	small	2	yes	yes
dog	medium	4	no	no
spider	tiny	8	no	yes
ladybird	tiny	6	no	yes
cow	large	4	no	no
bee	tiny	6	no	no
sparrow	small	2	yes	yes

1-to-N Coding animal data set

- Replace categorical attribute “size” with four binary attributes

	size	legs
duck	small	2
dog	medium	4
spider	tiny	8
ladybird	tiny	6
cow	large	4
bee	tiny	6
sparrow	small	2



tiny	small	med	large	legs
0	1	0	0	2
0	0	1	0	0
1	0	0	0	8
1	0	0	0	6
0	0	0	1	4
1	0	0	0	6
0	1	0	0	2

- N.b.: could also encode to numerical value*
 - E.g.: tiny=1, small=2, med=3, large=4 (if sizes are equally distant)*

Categorical data: another example

- *Any other pre-processing needed?*
- Variable “legs”
 - If considered categorical: defined order → ordinal data
 - Can compute similarity: 2 closer to 4 than to 6
 - Numerical value, can compute distance directly
 - *Does the number of legs denote similarity?*

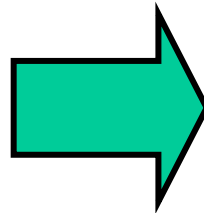
	<i>size</i>	<i>legs</i>
duck	small	2
dog	medium	4
spider	tiny	8
ladybird	tiny	6
cow	large	4
bee	tiny	6
sparrow	small	2

- *Does number of legs denote similarity?*
 - i.e., is an animal with 2 legs more similar to one with 4, or with 6?
 - And, is one with 4 equally similar to the one with 2 and 6?
 - Dog to monkey vs. dog to spider
 - One with 6 equally similar to one with 4 and 8?
 - Bee to spider vs. bee to cow

1-to-N Coding animal data set

.....

	<i>size</i>	<i>legs</i>
duck	small	2
dog	medium	4
spider	tiny	8
ladybird	tiny	6
cow	large	4
bee	tiny	6
sparrow	small	2



<i>2 legs</i>	<i>4 legs</i>	<i>6 legs</i>	<i>8 legs</i>
1	0	0	0
0	1	0	0
0	0	0	1
0	0	1	0
0	1	0	0
0	0	1	0
1	0	0	0

- *What to do with categorical data?*
 - a) 1-n coding – then apply any distance function mentioned earlier
 - b) Definition of custom distance functions that applies to categorical data

- Definition of custom distance functions

- E.g. adapt hamming distance

Hamming distance between two strings of equal length is the number of positions at which the corresponding symbols are different

→ count number of different nominal values

	<i>size</i>	<i>legs</i>	<i>feathers</i>	<i>eggs</i>	<i>distance</i>
horse	large	4	no	no	
duck	small	2	yes	yes	4
dog	medium	4	no	no	1
spider	tiny	8	no	yes	3
ladybird	tiny	6	no	yes	3
cow	large	4	no	no	
bee	tiny	6	no	no	2
sparrow	small	2	yes	yes	4

- Definition of custom distance functions

- E.g. adapt hamming distance

Hamming distance between two strings of equal length is the number of positions at which the corresponding symbols are different

→ count number of different nominal values

- Define custom distance for each attribute, aggregate e.g. via sum

- Different variables may exhibit significantly different value ranges
 - E.g. a length variable measured in cm, inch, or meters
 - Different types of measurements: length, speed, temperature, ...
 - Different types of measuring devices capturing different value ranges
 - ...
 - *Why is this a (potential) problem?*

- Some ML algorithms rely on *measuring the (numeric) distance* between samples
- ➔ There should be no impact by the value range
 - higher values in one attribute / variable would have unproportional effect on measure distance
 - ➔ would dominate distance metric
 - ➔ might thus dominate learning

- Remove effects of different value ranges in each attribute/variable
- Common method in statistics and machine learning
 - With many different notations
 - Scaling, Normalisation, Standardisation, ...
 - Often used interchangeably, different for each field ...

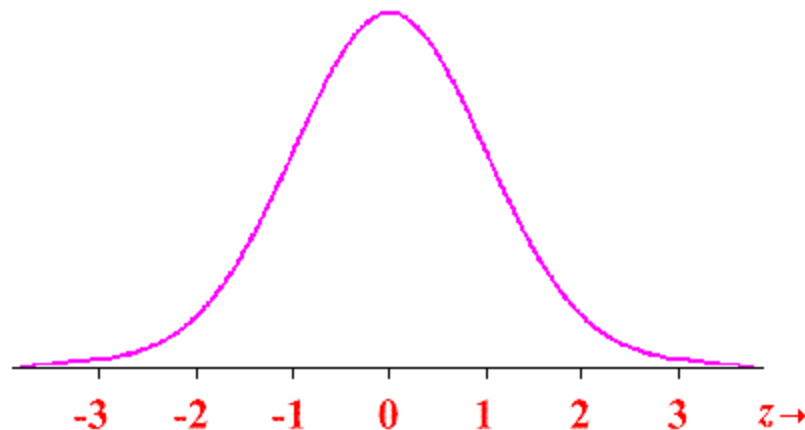
- Z-score standardisation / normalization
- Each variable x is converted to have a mean of 0 and a standard deviation of 1
 - subtracting the mean
 - dividing by standard deviation

$$z_i = \frac{x_i - \mu}{\sigma}$$

- Z-score standardisation / normalization

$$z_i = \frac{x_i - \mu}{\sigma}$$

- *What's the new value range after z-score?*



- Min-Max scaling
 - Scale all variables to the same (fixed) range
 - Often between 0 and 1
 - Subtract minimum value for each variable
 - Divide by value range of each variable

$$z_i = \frac{x_i - \min(X)}{\max(X) - \min(X)}$$

- *Multiply by new range (if different than 0..1)*

- *Is scaling an issue for*
 - *k-nn?*
 - *Perceptron?*

- Especially algorithms relying on distances
 - k-Nearest Neighbours
 - ...
- Scaling not needed for algorithms that don't use distances, e.g.
 - Naïve Bayes
 - Decision trees
 - ...

- Especially algorithms relying on distances
 - k-Nearest Neighbours
 - ...
- *Caveat*
 - Many implementations already do this pre-processing implicitly (e.g. WEKA)
 - Check default settings carefully

- For some samples, not all attribute values are known
- Some ML algorithms can handle missing values
- Solutions for other algorithms
 - Deletion of sample
 - bad when only few labelled samples
 - Imputation

- Substitution of a missing value
- Different methods
 - Mean value of the attribute (computed from other samples)
 - Random selection of value from another sample
 - Regression – using other attributes to predict
 - Clustering – values of cluster centroid
 - Nearest Neighbour – value of closest sample
- *Use in first exercise !*

- *When is imputation useful?*
- Most useful when the number of labelled samples (w/o missing values) is small
 - *Relatively easy to identify*
- When samples with missing values contain important information
 - *Difficult to identify*

Questions?

Upcoming topics:

- Decision trees
- Random Forests
- Support Vector Machines