

Course

Foundations of Artificial Intelligence

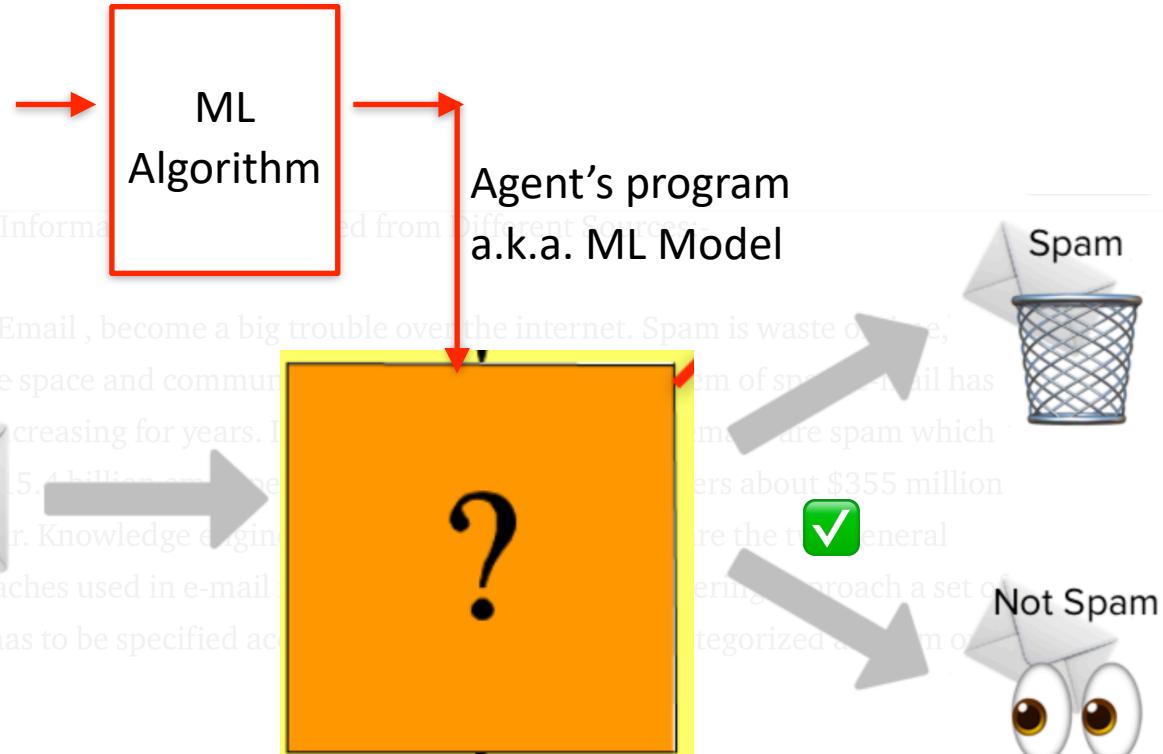
Lecture 09 – Supervised Machine Learning Algorithms

Dr. Mohsen Mesgar

Universität Duisburg-Essen

Recall ...

Examples of
problem/solution



Any other open questions?



In this lecture, you learn about ...

- Basic ML algorithms
 - Decision Trees
 - kNN
 - Naive Bayes classifier
 - Support Vector Machines (SVMs)

Decision Trees

Name gender

A network diagram illustrating name gender relationships. Nodes are represented by names, and connections between them indicate a relationship. The nodes are arranged in a roughly circular pattern:

- Top-left: Tom
- Moving clockwise:
 - Andrea
 - Dorchadas
 - Brigitte
 - Jamel
 - Lea
 - Wilhelm
 - Dafne
 - Simone
 - Diarmid
 - Lumen
 - Noor
 - Dietgard

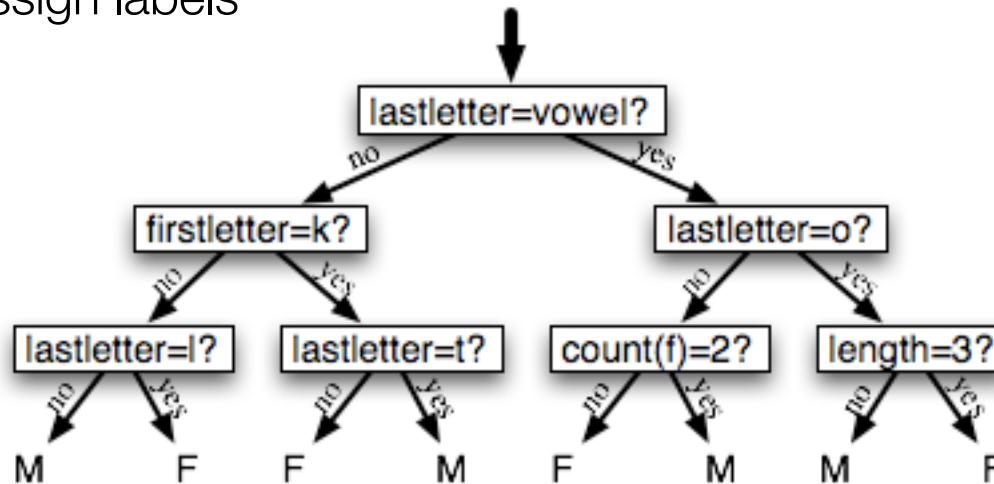
The connections are as follows:

- Tom is connected to Andrea.
- Andrea is connected to Dorchadas.
- Dorchadas is connected to Brigitte.
- Brigitte is connected to Jamel.
- Jamel is connected to Lea.
- Lea is connected to Wilhelm.
- Wilhelm is connected to Dafne.
- Dafne is connected to Simone.
- Simone is connected to Diarmid.
- Diarmid is connected to Lumen.
- Lumen is connected to Noor.
- Noor is connected to Dietgard.
- Dietgard is connected to Wilhelm.

Decision trees

Flowchart that selects labels for input values

- Decision nodes check feature values
- Leaf nodes assign labels



Start with the root node (decision stump)

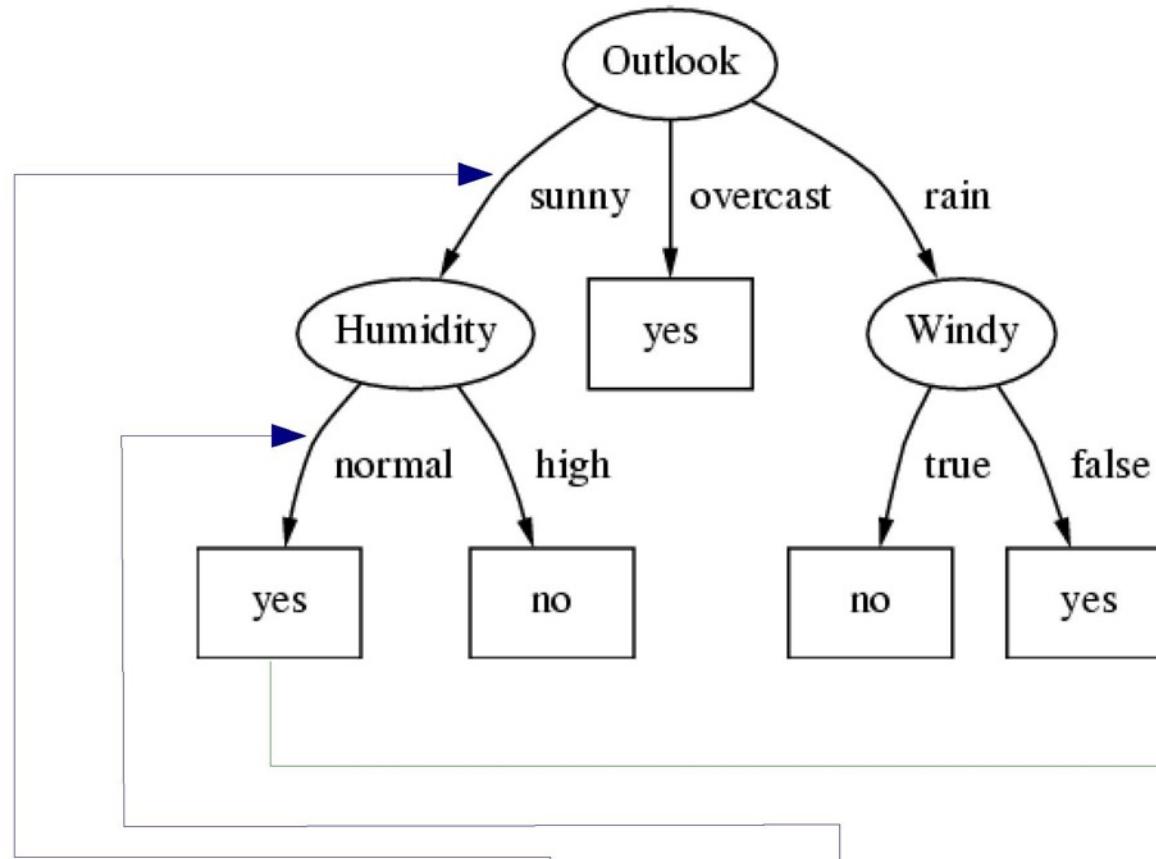
- Use feature value that classifies all input best
- Continue growing using most informative features

Play golf example

Day	Temperature	Outlook	Humidity	Windy	Play Golf?
07-05	hot	sunny	high	false	no
07-06	hot	sunny	high	true	no
07-07	hot	overcast	high	false	yes
07-09	cool	rain	normal	false	yes
07-10	cool	overcast	normal	true	yes
07-12	mild	sunny	high	false	no
07-14	cool	sunny	normal	false	yes
07-15	mild	rain	normal	false	yes
07-20	mild	sunny	normal	true	yes
07-21	mild	overcast	high	true	yes
07-22	hot	overcast	normal	false	yes
07-23	mild	rain	high	true	no
07-26	cool	rain	normal	true	no
07-30	mild	rain	high	false	yes

today	cool	sunny	normal	false	?
tomorrow	mild	sunny	normal	false	?

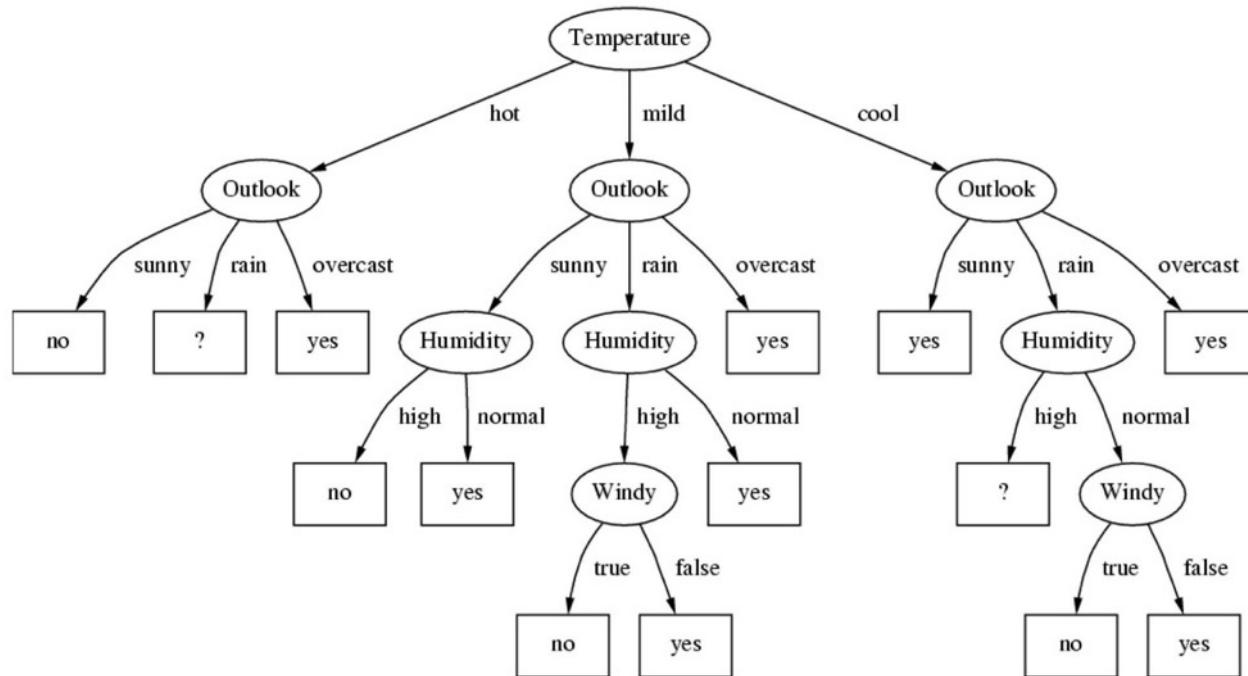
Final decision tree



tomorrow	mild	sunny	normal	false	?
----------	------	-------	--------	-------	---

How do we find the best decision tree?

- Search through all possible trees? → Too expensive.
- This tree also explains all the training data.
Will it generalize well to new data?
→ Determine good decision nodes based on information gain.



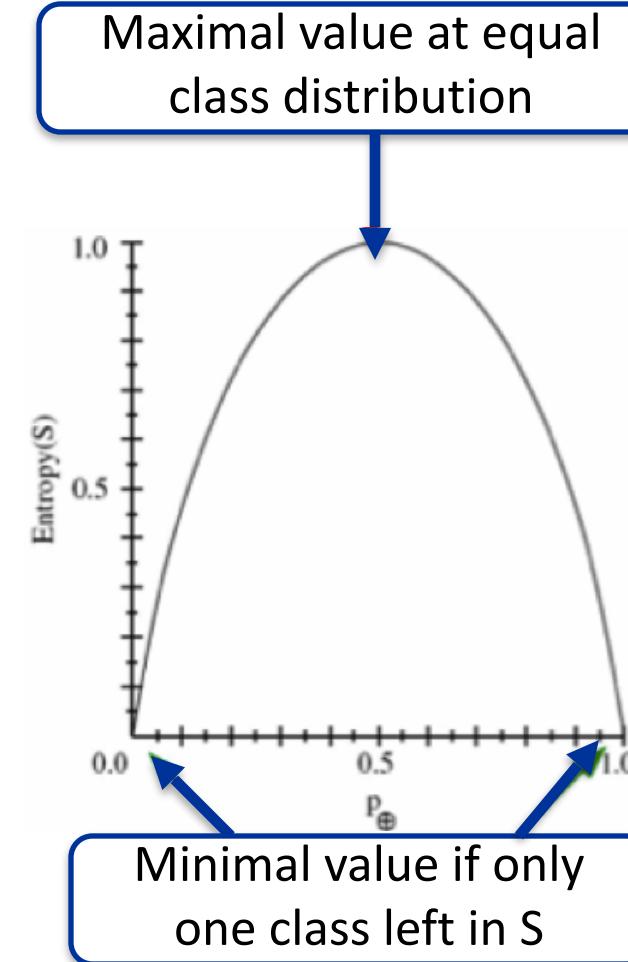
Entropy (for two classes)

- S is a set of examples
- P_i is the proportion of examples in class i
- $P_j = 1 - p_i$ is the proportion of examples in the other class j

$$E(S) = -p_i \log_2 p_i - p_j \log_2 p_j$$

Interpretation:

- Amount of unorderedness in the class distribution of S



Average entropy / information

Problem:

- Entropy only computes the quality of a single (sub-)set of examples
 - Corresponds to a single value
- How can we compute the quality of the entire split?
 - Corresponds to an entire attribute

Solution:

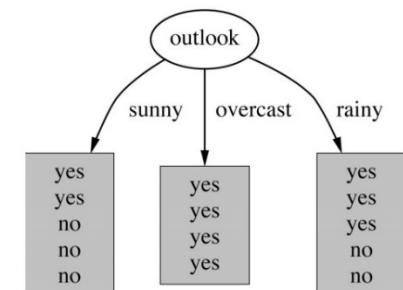
- Compute the weighted average over all sets resulting from the split
 - Weighted by their size

$$I(S, A) = \sum_i \frac{|S_i|}{|S|} \cdot E(S_i)$$

Example:

- Average entropy for attribute Outlook:

$$I(\text{Outlook}) = \frac{5}{14} \cdot 0.971 + \frac{4}{14} \cdot 0 + \frac{5}{14} \cdot 0.971 = 0.693$$



Decision tree example

outlook=sunny (3 yes/2 no)

$$\text{Entropy}(\text{sunny}) = -\frac{3}{5} * \log \frac{3}{5} - \frac{2}{5} * \log \frac{2}{5} = 0.97$$

log(0) is undefined

outlook=overcast (4 yes/0 no)

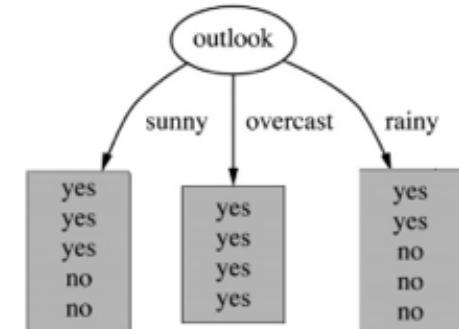
$$\text{Entropy}(\text{overcast}) = -1 * \log 1 - 0 * \log 0 = 0$$

outlook=rainy (2 yes/3 no)

$$\text{Entropy}(\text{rainy}) = -\frac{2}{5} * \log \frac{2}{5} - \frac{3}{5} * \log \frac{3}{5} = 0.97$$

Information(outlook) =

$$\frac{5}{14} * 0.97 + \frac{4}{14} * 0 + \frac{5}{14} * 0.97 = 0.69$$



Information of features

Feature humidity

- Entropy(high) = 0.99 (3 yes/4 no)
- Entropy(low) = 0.59 (6 yes/1 no)
- Information(humidity) = $\frac{7}{14} * 0.99 + \frac{7}{14} * 0.59 = 0.79$

Feature wind

- Entropy(weak) = 0.81 (6 yes/2 no)
- Entropy(strong) = 1.0 (3 yes/3 no)
- Information(wind) = $\frac{8}{14} * 0.81 + \frac{6}{14} * 1 = 0.89$

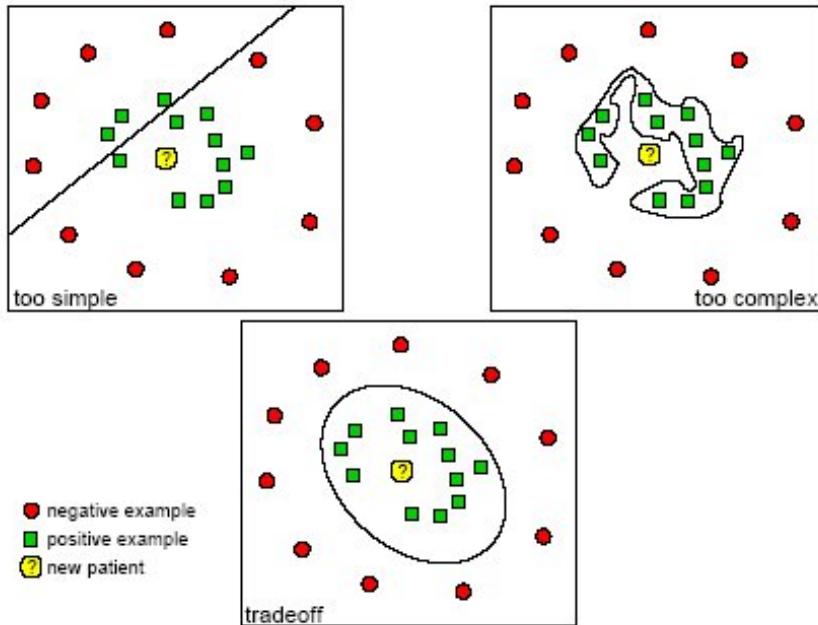
Feature temperature

- Entropy(hot) = 1.0 (2 yes/2 no)
- Entropy(mild) = 0.92 (4 yes/2 no)
- Entropy(cool) = 0.81 (3 yes/1 no)
- Information(temperature) = $\frac{4}{14} * 1 + \frac{6}{14} * 0.92 + \frac{4}{14} * 0.81 = 0.91$

Information Gain

- Entropy(before) = 0.94 (9 yes/5 no)
- InformationGain(outlook) = $0.94 - 0.69 = 0.25$
- InformationGain(humidity) = $0.94 - 0.79 = 0.15$
- InformationGain(wind) = $0.94 - 0.89 = 0.05$
- InformationGain(temperature) = $0.94 - 0.91 = 0.03$
- Select feature outlook as root node

Underfitting & overfitting



Underfitting

- the model is not capable of learning the (complex) patterns in the training set
→ use better features or classifier

Overfitting

- The model perfectly learns to classify the training data, but has bad generalization ability → high test error
- results in useless model for unseen data
- typical for small sample sizes combined with powerful models

Overfitting and pruning

The smaller the complexity of a concept, the less it overfits the data

- a polynomial of degree n can always fit $n+1$ points

Thus, learning algorithms try to keep the learned concepts simple

- note a “perfect” fit on the training data can always be found for a decision tree!
(except when data is contradictory)

Pre-pruning

- stop growing a branch when information becomes unreliable

Post-pruning

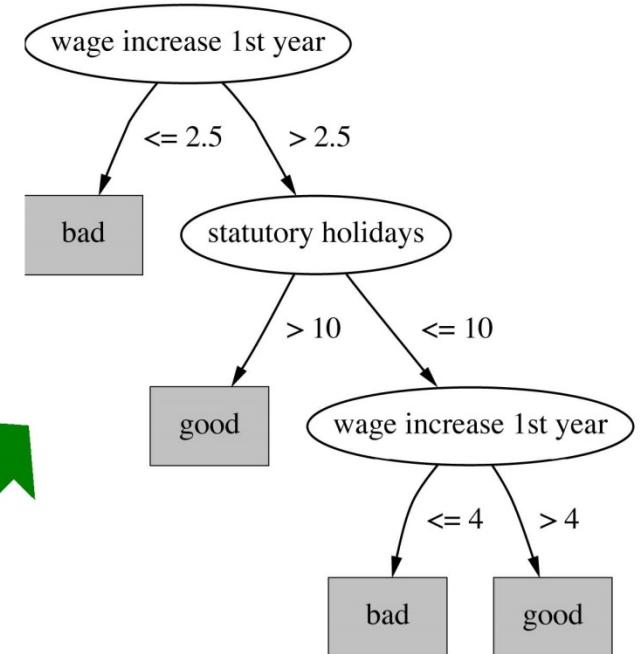
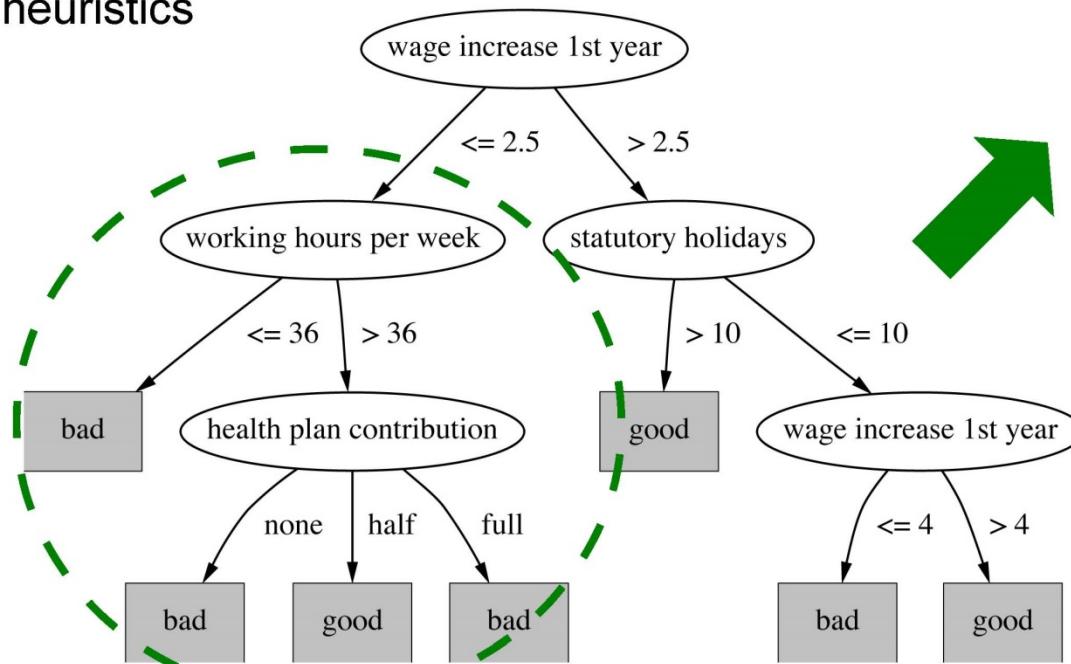
- grow a decision tree that correctly classifies all training data
- simplify it later by replacing some nodes with leaves

Post-pruning preferred in practice, because pre-pruning can “stop early”

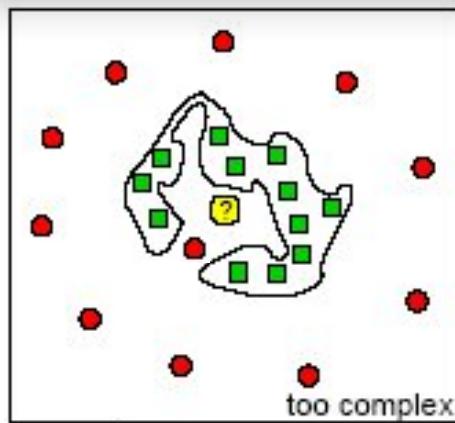
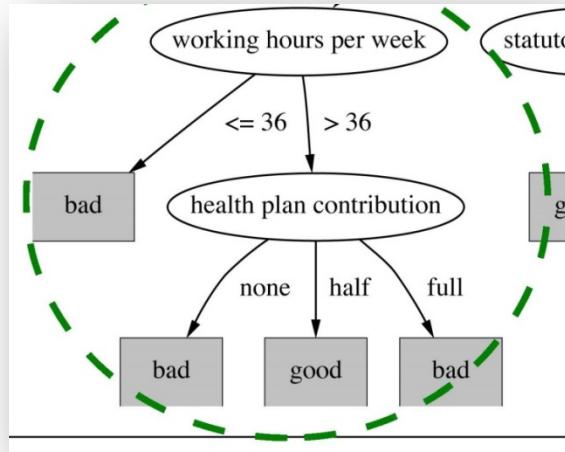
Sub-tree replacement

Proceeds Bottom-up:

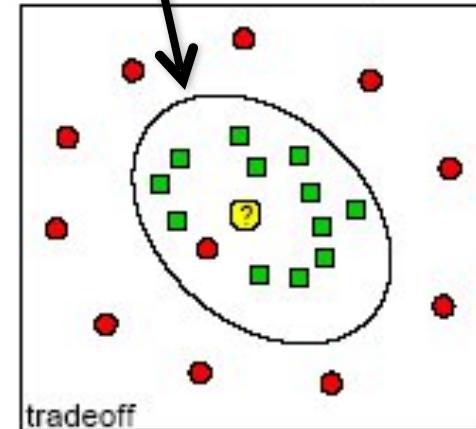
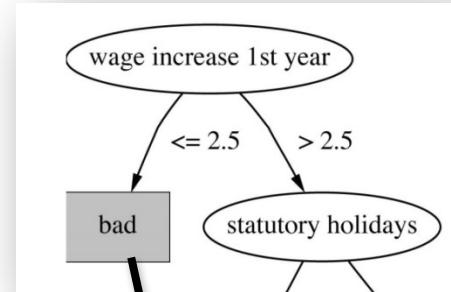
- consider replacing a tree only after considering all its subtrees
- may make a difference for complexity-based heuristics



Sub-tree replacement

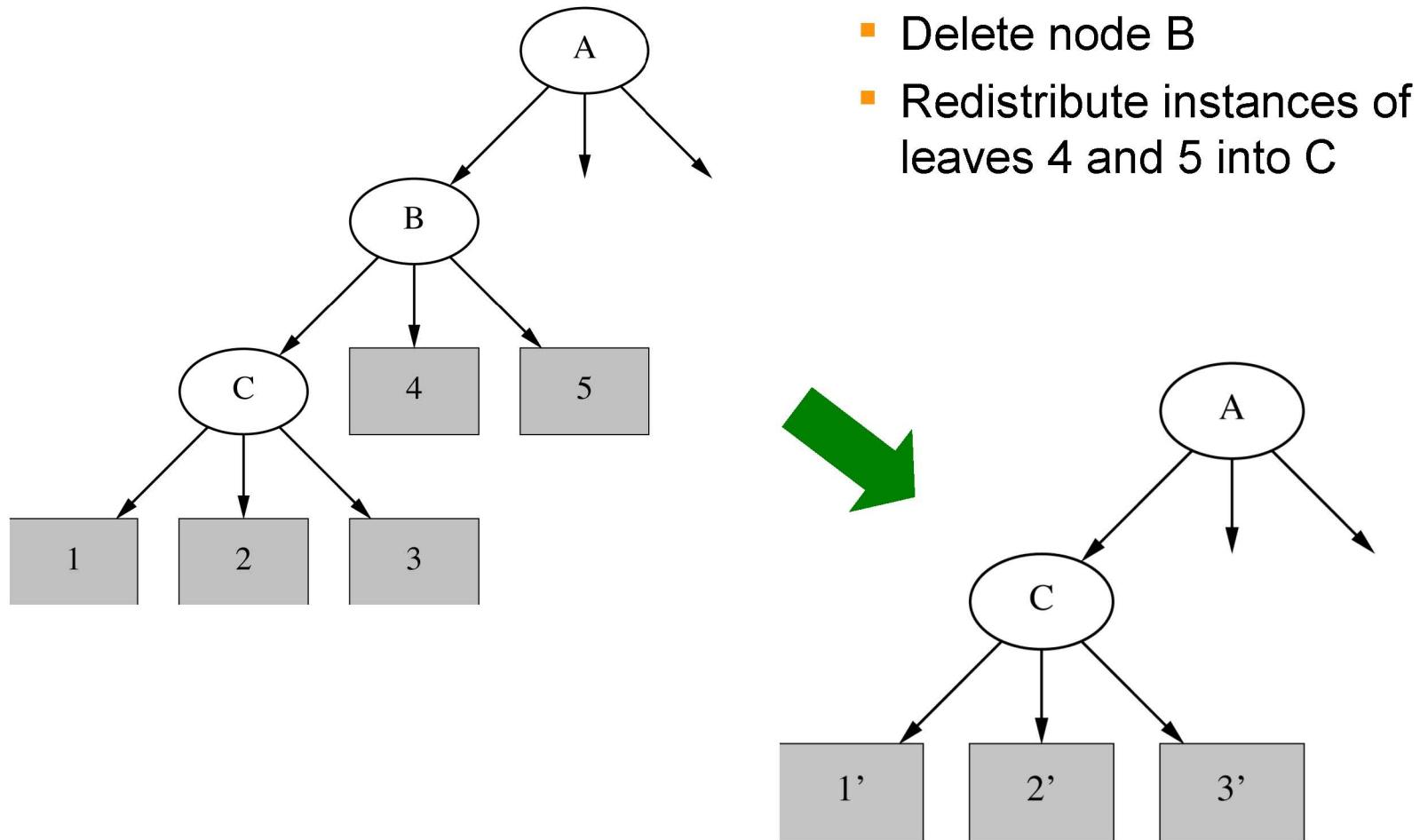


Overfitted



Pruned

Sub-tree raising



Advantages

- efficient to create (training is fast)
- simple to interpret → liked by users, explicit knowledge: white box model
- useful to learn models for data which can be hierarchically categorized

Problems

- overfitting is a serious problem (idiosyncrasies in the data set)
- simpler, seemingly less accurate trees are often preferable → pruning
- **amount of training data in lower nodes** is quite small
 - because it is split at every decision node
- **forces features to be checked in a certain order**, even if features are independent of each other

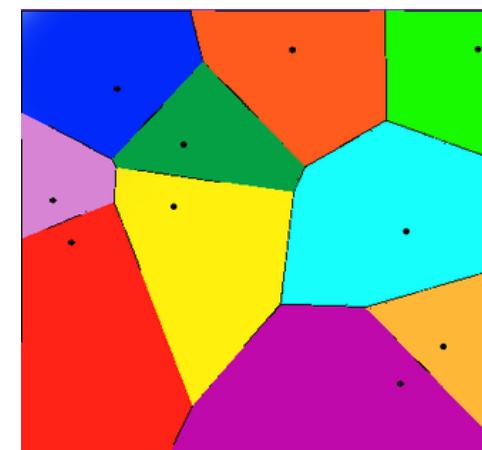
k-nearest neighbor (KNN)

Nearest neighbor prediction

Simply predict the class of the nearest neighbor of the unseen instance in the representation space.

Voronoi Diagram: nearest neighbor in 2-dimensional space

- example: 10 training instances
- colored area: neighbor region of each training instance
- determine color for test instance
- assign class of training stance
- How do we determine “nearness”?



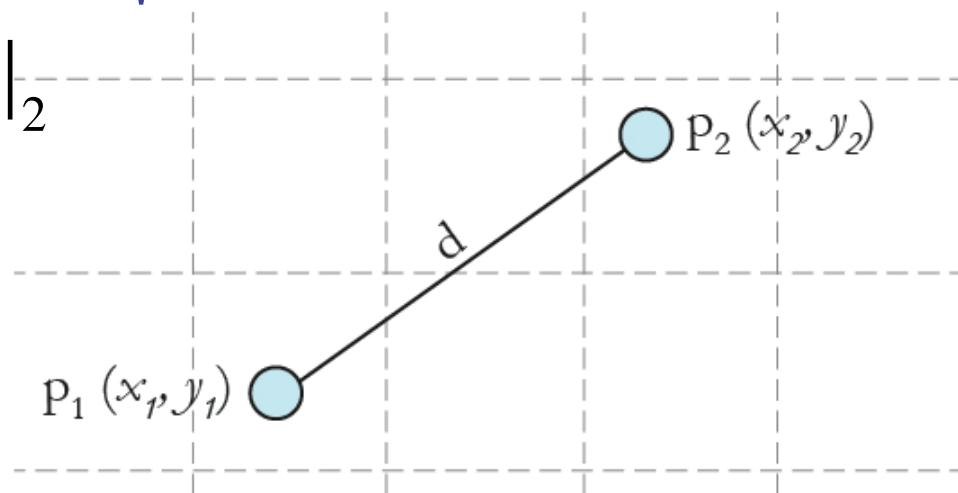
- Nearness in feature space = similarity between represented concepts
- **Assumption:** concepts that have similar values in many dimensions (= features) are similar.
- **How do we measure similarity?**
By measuring the **distance** between two feature vectors

Euclidean distance

- ◊ Straight line distance between two vectors
- ◊ In 2-dimensional space → Pythagoras

- ◊ In n-dimensional space: $d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$

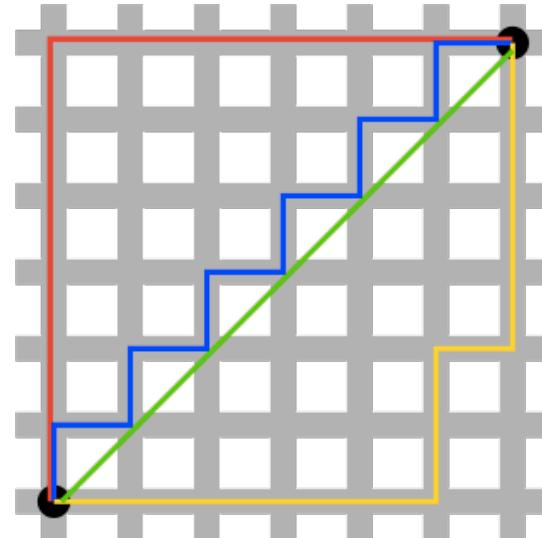
- ◊ also known as L₂-norm: $\|x - y\|_2$



$$\text{Euclidean distance } (d) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

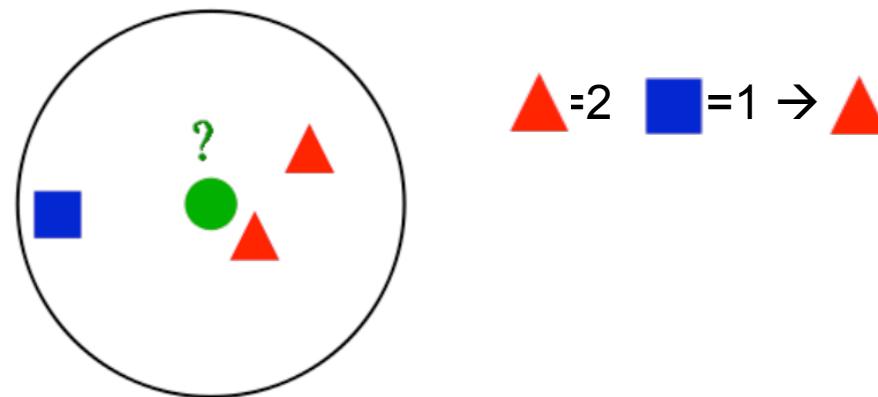
Manhattan distance

- Straight line is often not possible
- Step-wise distance between two points:
sum of the absolute differences of their Cartesian coordinates.
- $\mathbf{d}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n |x_i - y_i|$
- also known as L_1 -norm: $\|\mathbf{x} - \mathbf{y}\|_1$
- red, blue, yellow: Manhattan distance
- green: Euclidean distance



Assign majority class of the k-nearest neighbors.

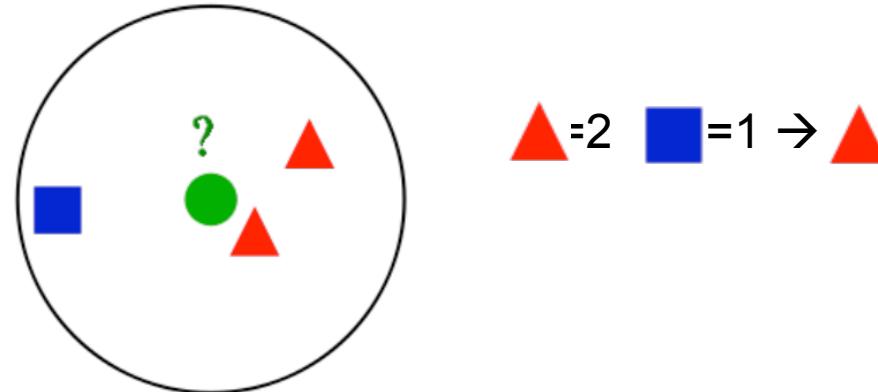
k=3



k-nearest neighbor classification

Assign majority class of the k-nearest neighbors.

$k=3$

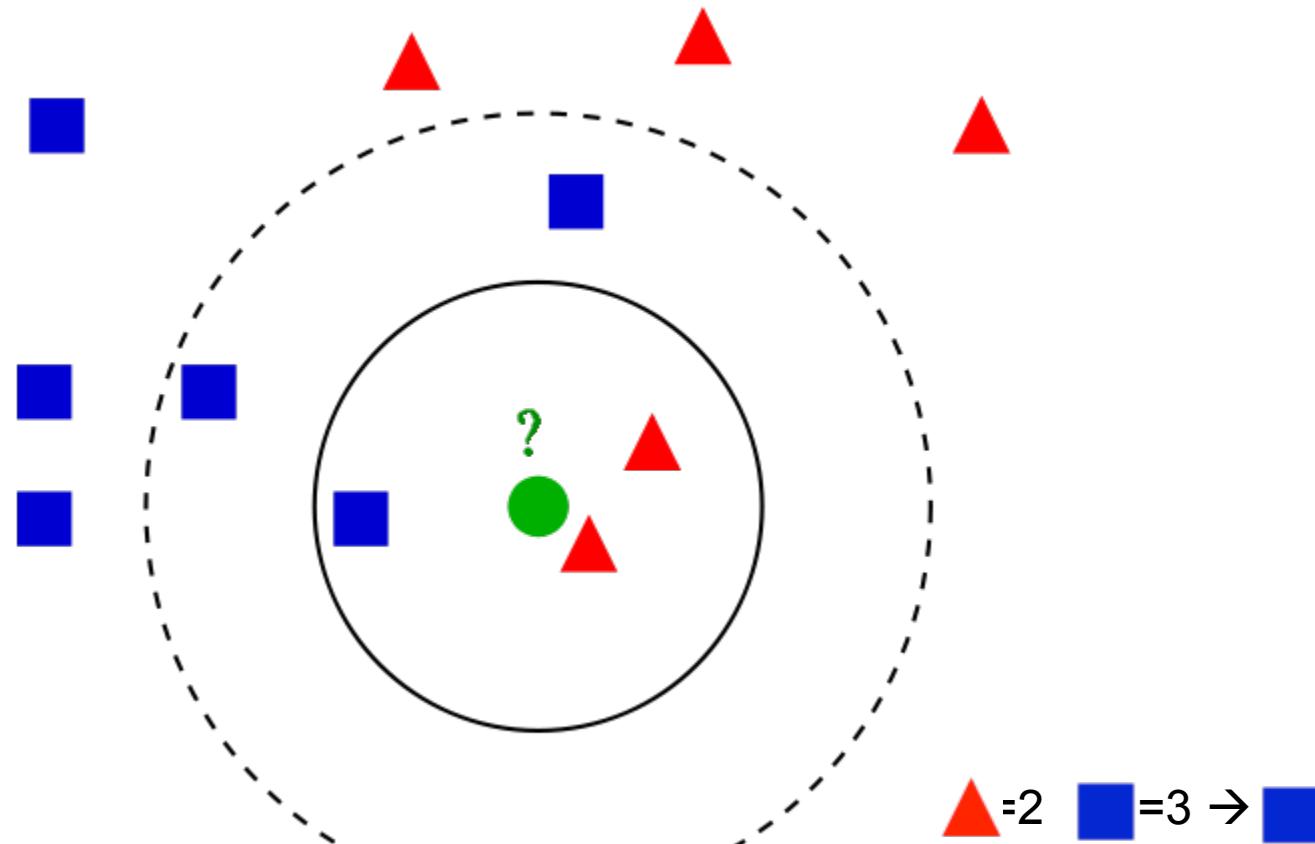


- No training!
- Just store all instances
 - might build special lookup data-structure to optimize comparison
- k is a hyperparameter

Also known as ***Collaborative Filtering*** or ***Instance-based Learning***.

Choice of k

3-nearest neighbors vs. 5-nearest neighbors prediction

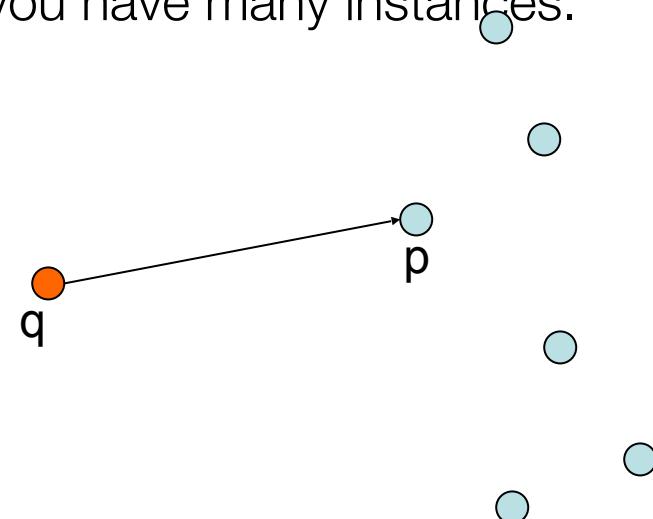


Nearest neighbor search

- **Given:** a set P of n points in \mathbb{R}^d , a test instance q
- **Goal:** find the nearest neighbor p of q in P
- **Naïve approach:** compute the distance from the query point to every other point in the database, keeping track of the “**best so far**”

The naïve approach is very expensive, if you have many instances.

- n comparisons for each test instance.
- Apply more efficient search algorithms.



Properties of kNN algorithm

Inductive bias

- similar classification of nearby instances
- similar to human heuristics

Curse of dimensionality

- similarity metric might be misled by irrelevant attributes

Naïve Bayes

Probabilistic classifier

- Calculate a probability distribution over all classes (not only a single class)
- More realistic setting for many tasks → model uncertainty

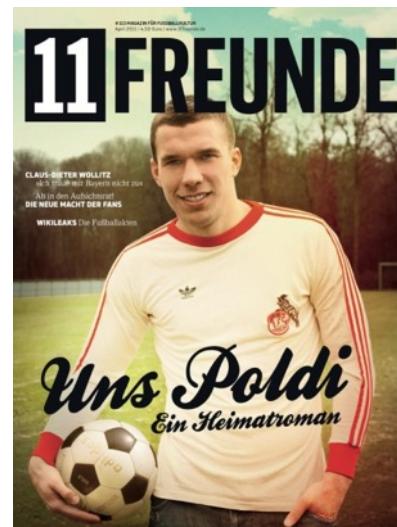
Example task: Document classification

We have a training dataset with documents from 3 different categories:

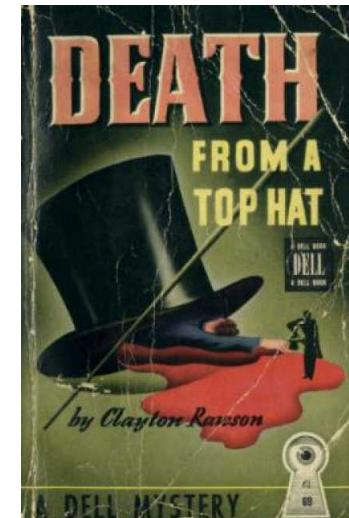
Automotive



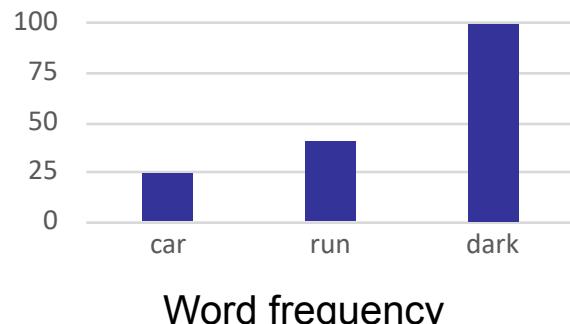
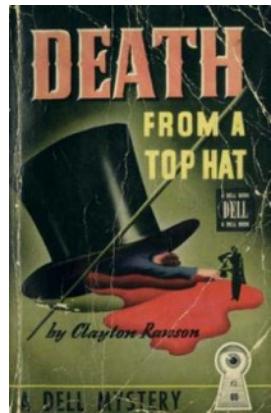
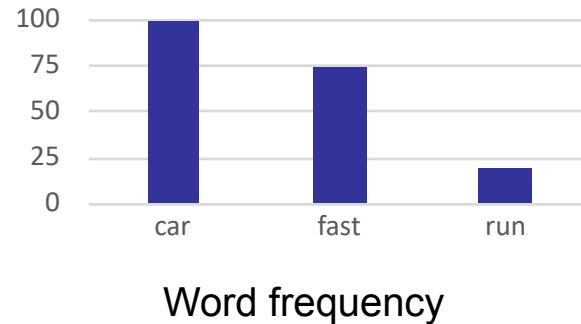
Sports



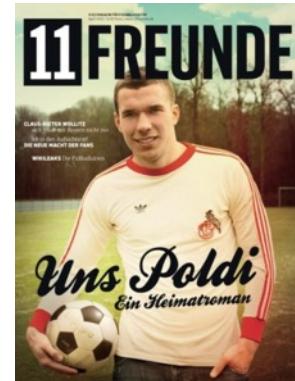
Murder Mystery



Feature representation: words in the document

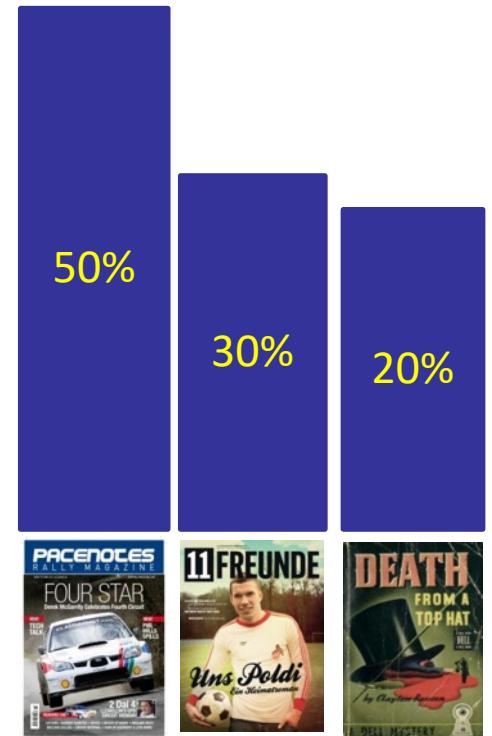


Some words are prototypical for a category. Each occurrence of a word in a document is an indicator for the category.



Prior probability

- How likely is a category?
 - Frequency of each category in the training data
- **If we assume that training data are representative,**
 - we expect that about half of the documents in the test data belong to the automotive category



Prior probability

Likelihood = Prior prob * feature contribution



* run * fast * ...



* run * fast * ...



* run * fast * ...



Naïve Bayes classifier

Goal: find the class with the maximum probability.

$$\max_{c \in \text{classes}} \frac{p([x_1, x_2, x_3, \dots x_n] | c) * p(c)}{p([x_1, x_2, x_3, \dots x_n])}$$

The denominator is a constant and can be dropped:

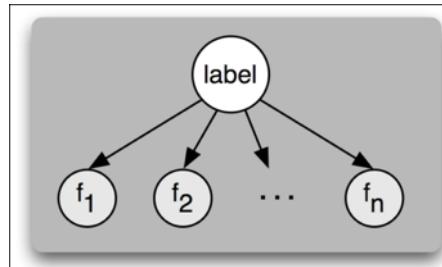
$$\max_{c \in \text{classes}} p([x_1, x_2, x_3, \dots x_n] | c) * p(c)$$

It is very unlikely that we find an instance with the exact same feature representation in the training data.

→ Assume independence of features!

Independence assumption

- Given a label, all features are independent of each other.
- This assumption makes it much easier to combine different features, as we do not need to worry about interactions between features.



The conditional probability that feature x_i occurred in class c in the training data.

$$p([x_1, x_2, x_3, \dots, x_n] | c) \dashrightarrow \text{independence assumption:}$$

$$\prod_{i=1}^n p(x_i | c)$$

Naïve Bayes Classifier:

$$\max_{c \in \text{classes}} \prod_{i=1}^n p(x_i | c) * p(c)$$

The naïvety of independence

It's unreasonable to assume that all features are independent.

- If *football* occurs, it's more likely that *goal* also occurs in the same text.

What happens if we ignore this dependence?

- Possibility of double-counting highly correlated features
- *football* and *goal* are highly correlated, but assumed independent, so each feature is counted separately in the training corpus
- This pushes the classifier closer to a given label than justified.

→ Correlated features are problematic for Naïve Bayes classifiers

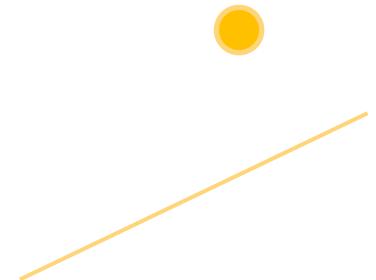
Naïve Bayes: Summary

- This is a probabilistic learning method which uses knowledge about the prior probabilities of classes and about the conditional probability of features occurring in a particular class.
- Training is very fast.
- Binary, numeric and nominal features can be mixed.
- If the independence assumption is violated too much, the classifier fails.
 - Identical or highly correlated features are problematic.
- Can be addressed with feature selection algorithms

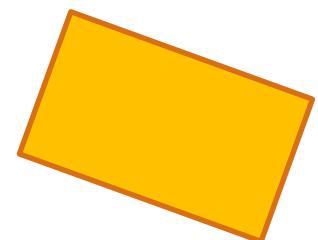
Support Vector Machines (SVM)

How to split a feature space into 2 classes?

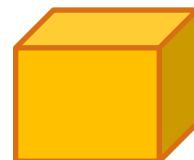
- 1-dimensional (= 1 feature)
→ By a point: if feature $x \geq 0.7$ then class 1 else class 2



- 2-dimensional (= 2 features)
→ By a line: $w x - b$
- 3-dimensional (= 3 features)
→ By a plane: $w_1 x_1 + w_2 x_2 - b$



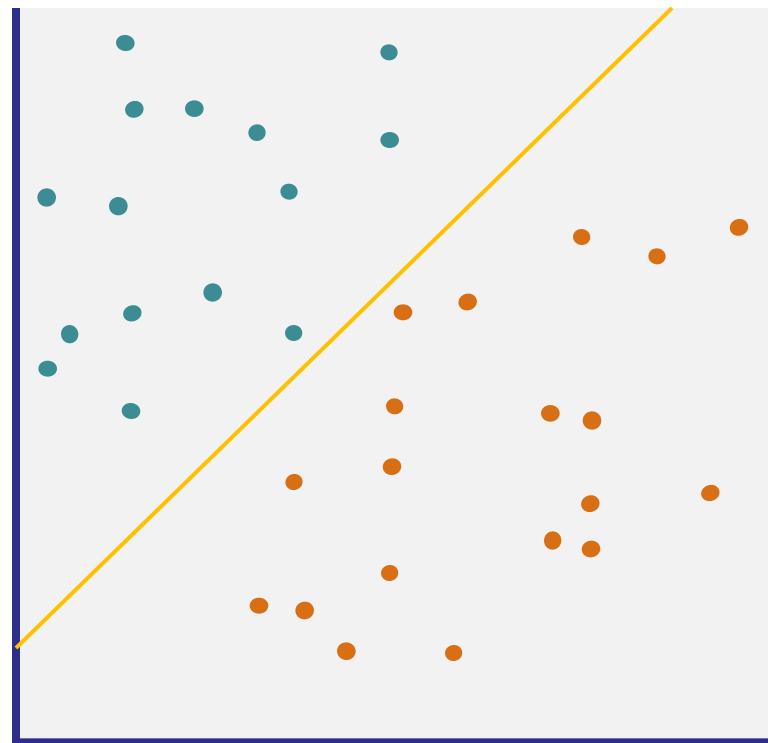
- n -dimensional (= n features)
→ By a hyperplane: $\vec{w} \cdot \vec{x} - b$
- A hyperplane in R^n is a $n-1$ dimensional subspace of R^n .



?

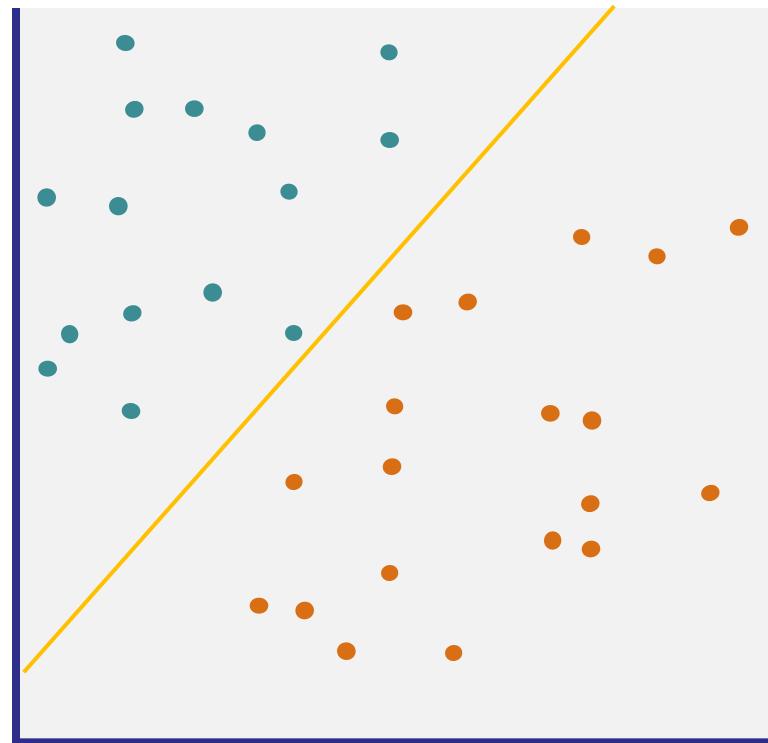
Linear separation

How would you classify these data points?



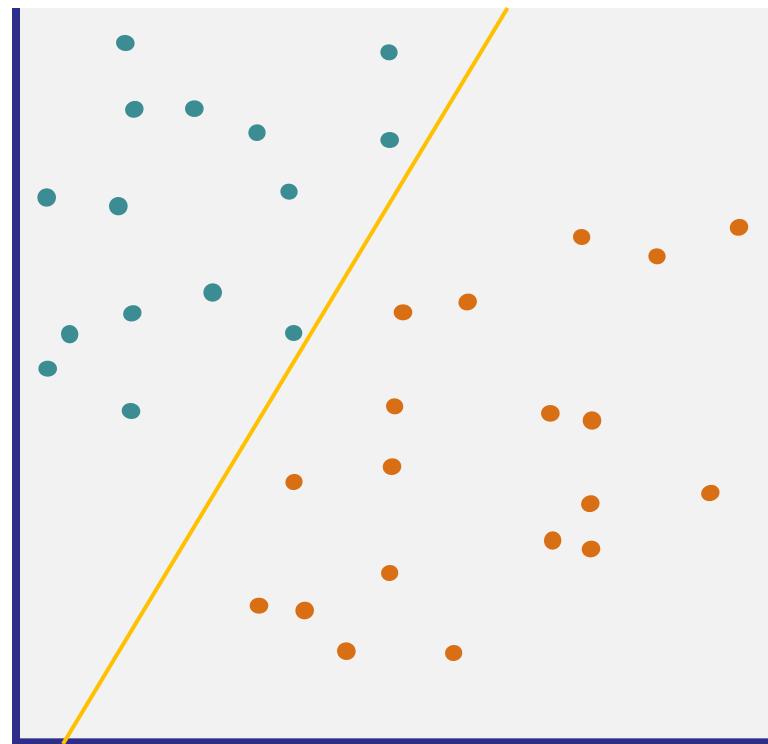
Linear separation

How would you classify these data points?



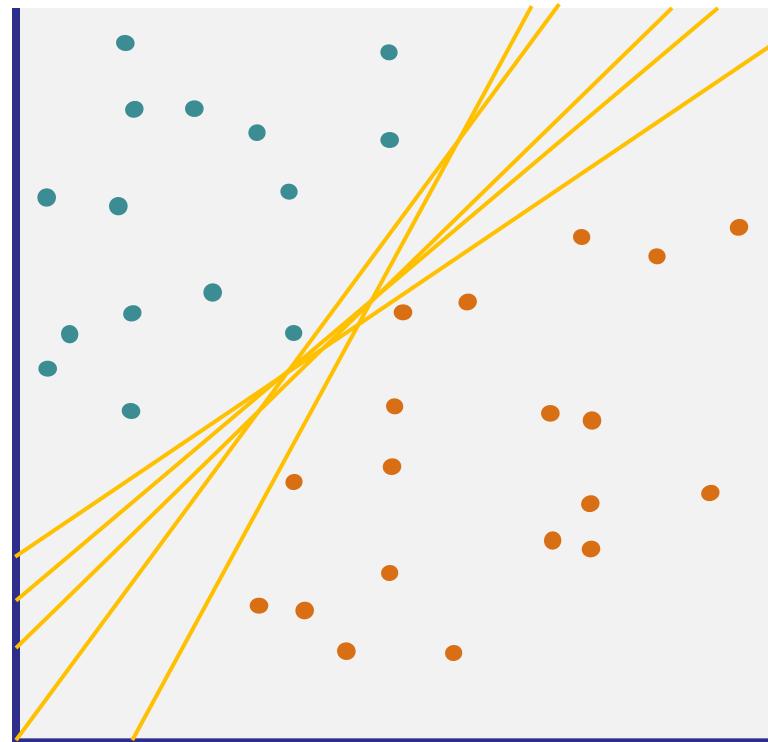
Linear separation

How would you classify these data points?



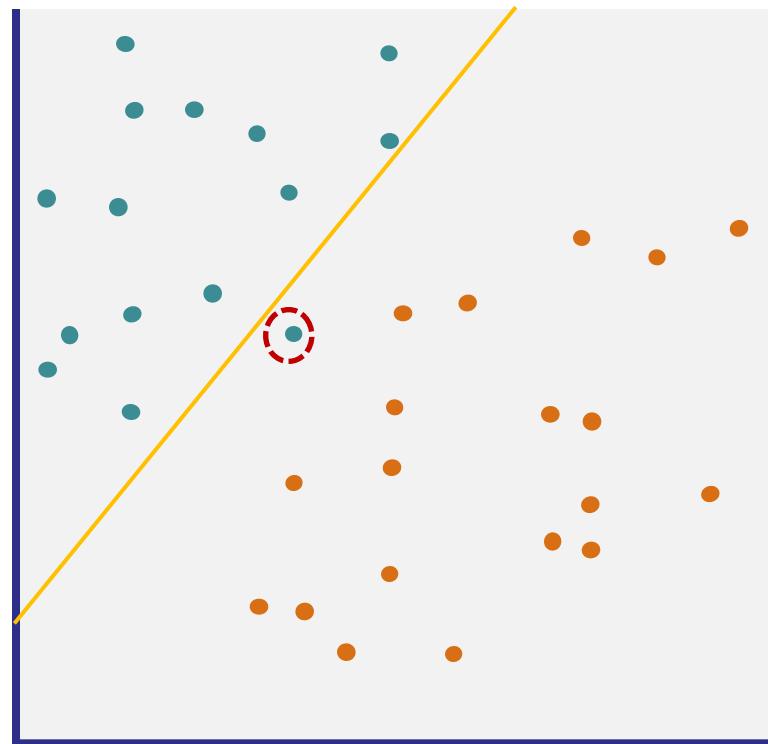
Linear separation

Any of these lines would be fine. How do we find the best?



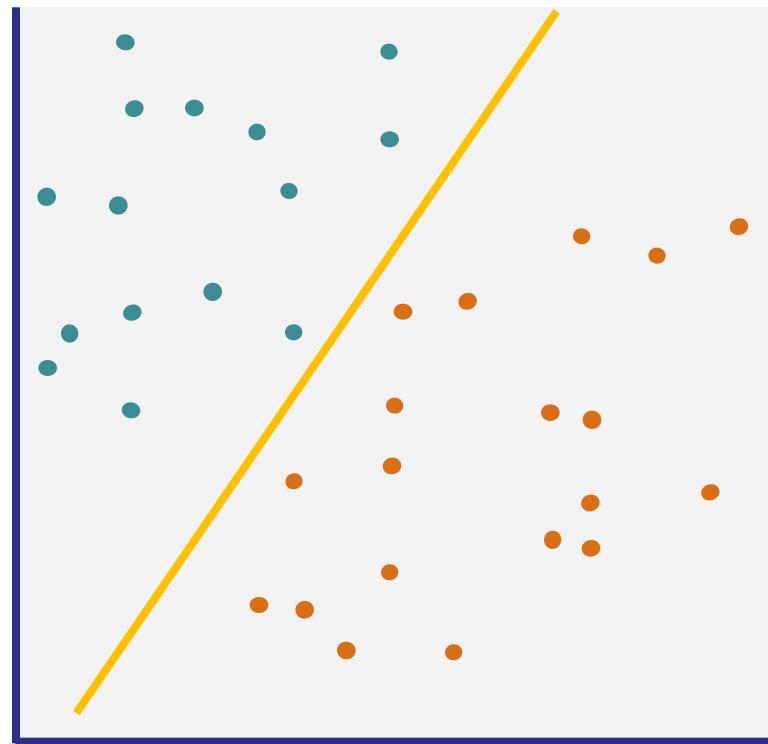
Linear separation

Avoid misclassification



Classifier margin

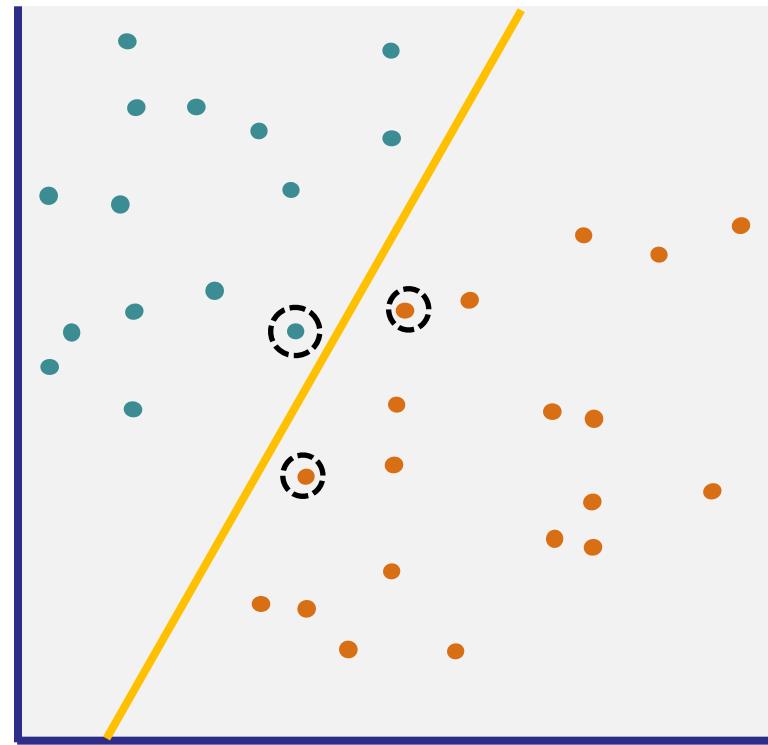
Margin of a linear classifier: the width by which the boundary could be increased before hitting a datapoint



Support vectors

Goal: Find the linear classifier with the maximum margin.

Support vectors: datapoints which are touched by the margin.



→ simplest form of support vector machines (SVM).

Formally

- $\vec{w} \cdot \vec{x} - b = 0$ defines a hyperplane, w is the normal vector of the hyperplane

- left boundary: $\vec{w} \cdot \vec{x} - b = 1$

- right boundary: $\vec{w} \cdot \vec{x} - b = -1$

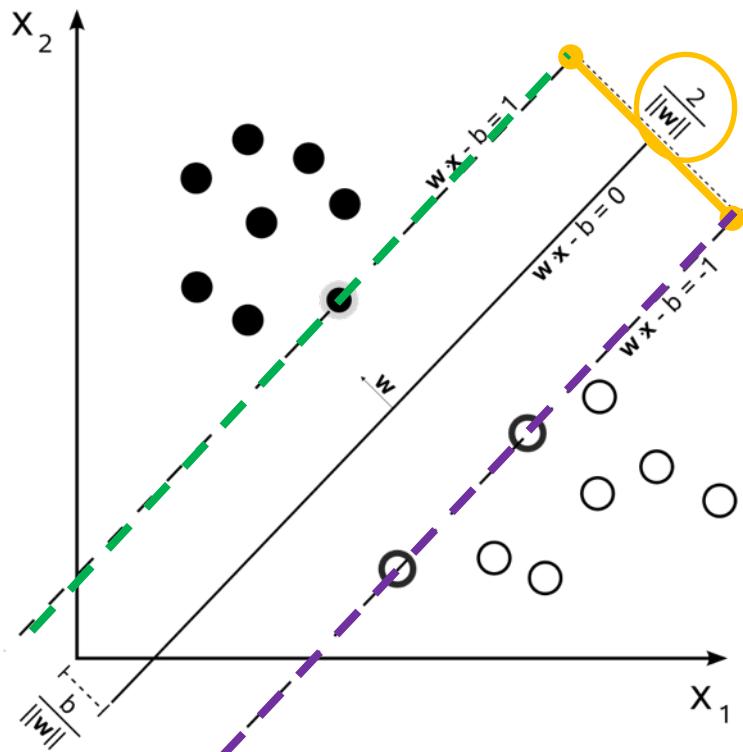
1. maximize the margin
= the distance between left and right boundary

$$= \frac{2}{\|\vec{w}\|}$$

→ minimize \vec{w}

2. determine the best \vec{w} and b , such that
for all points \vec{x}_i it holds that

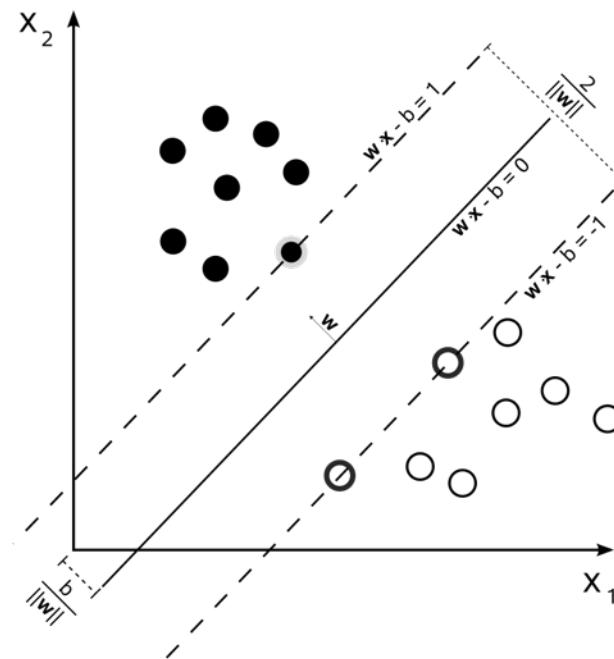
- $\vec{w} \cdot \vec{x}_i - b \geq 1, \text{ if } y_i = \text{class 1}$
- $\vec{w} \cdot \vec{x}_i - b \leq -1, \text{ if } y_i = \text{class 2}$



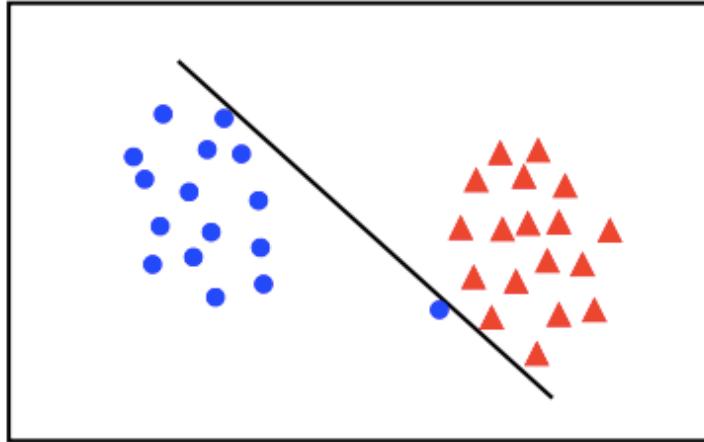
https://en.wikipedia.org/wiki/Support_vector_machine#/media/File:Svm_max_sep_hyperplane_with_margin.png

Maximum margin

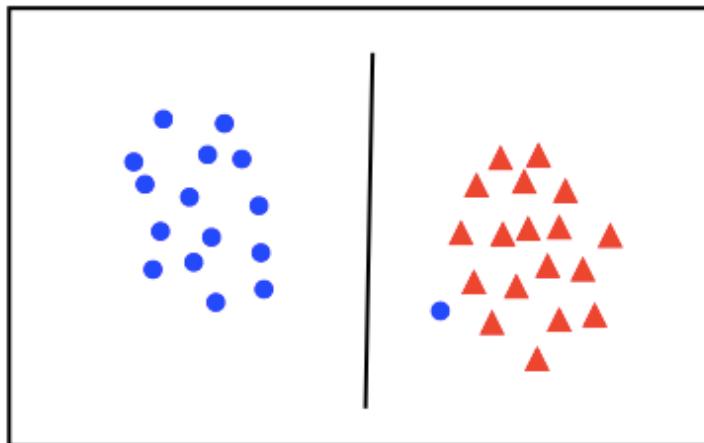
- Maximizing the margin is an intuitive approach.
- Consequence: only support vectors (the borderline cases) are important; other training examples can be skipped.
- Empirically it often works very well.



Handling Noise



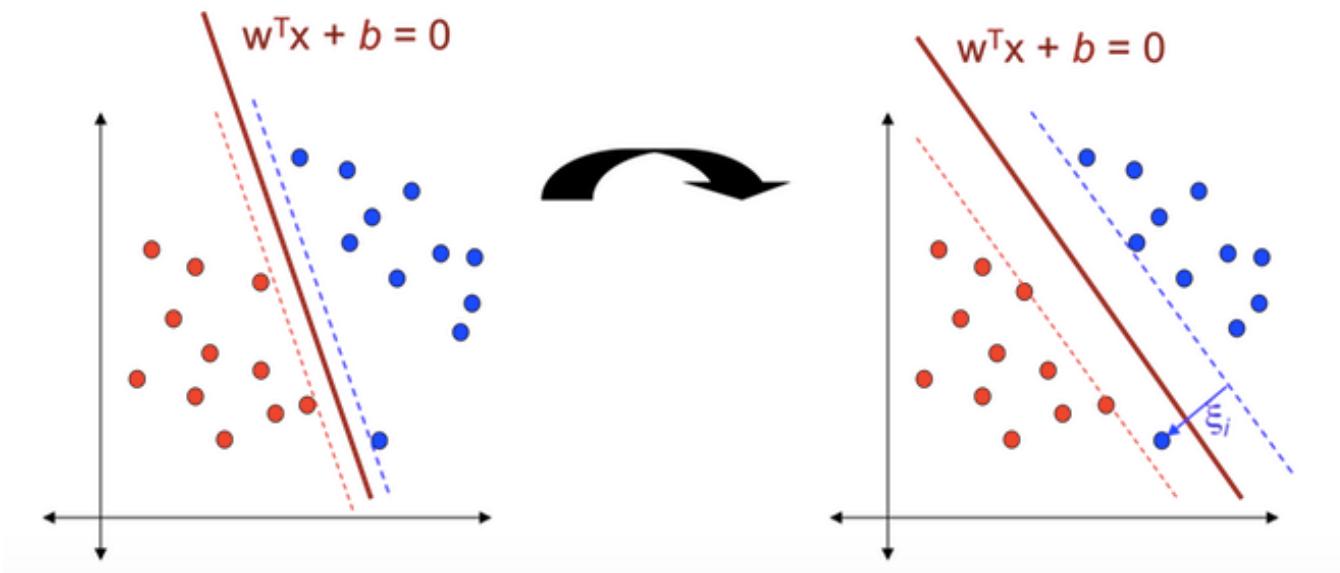
- Which solution is better?
- **Hard Margin:**
All data points should be classified correctly → No training error
- What if the training set is noisy?



Cut me some slack!

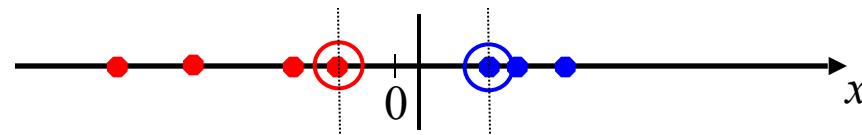
Soft margin classification:

- ◊ Slack variables ξ_i allow misclassification of difficult or noisy examples.
- A trade-off between the margin size and the number of mistakes on the training data.

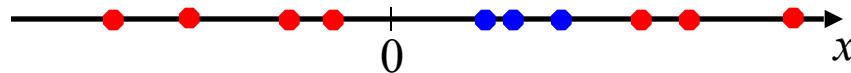


Non-linear SVMs

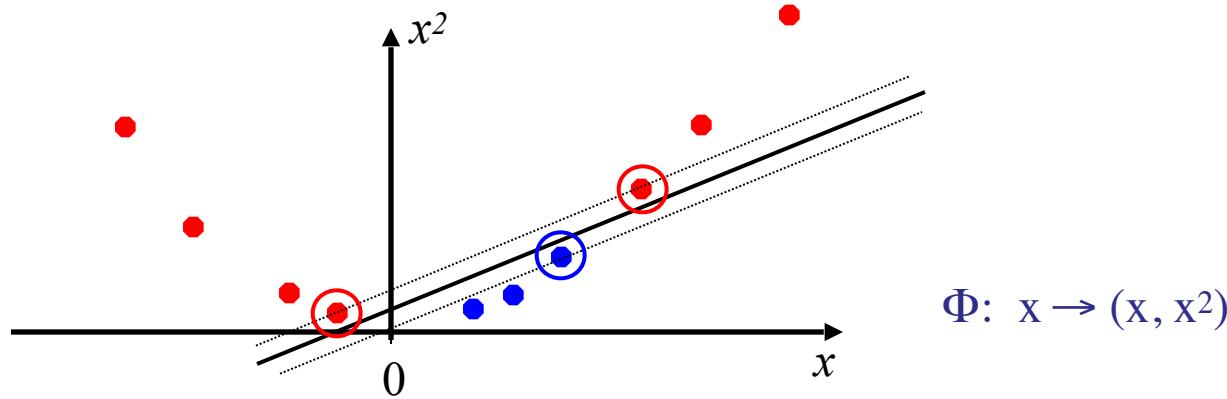
Datasets that are linearly separable with some noise work out great:



But what are we going to do if the dataset is just too hard?

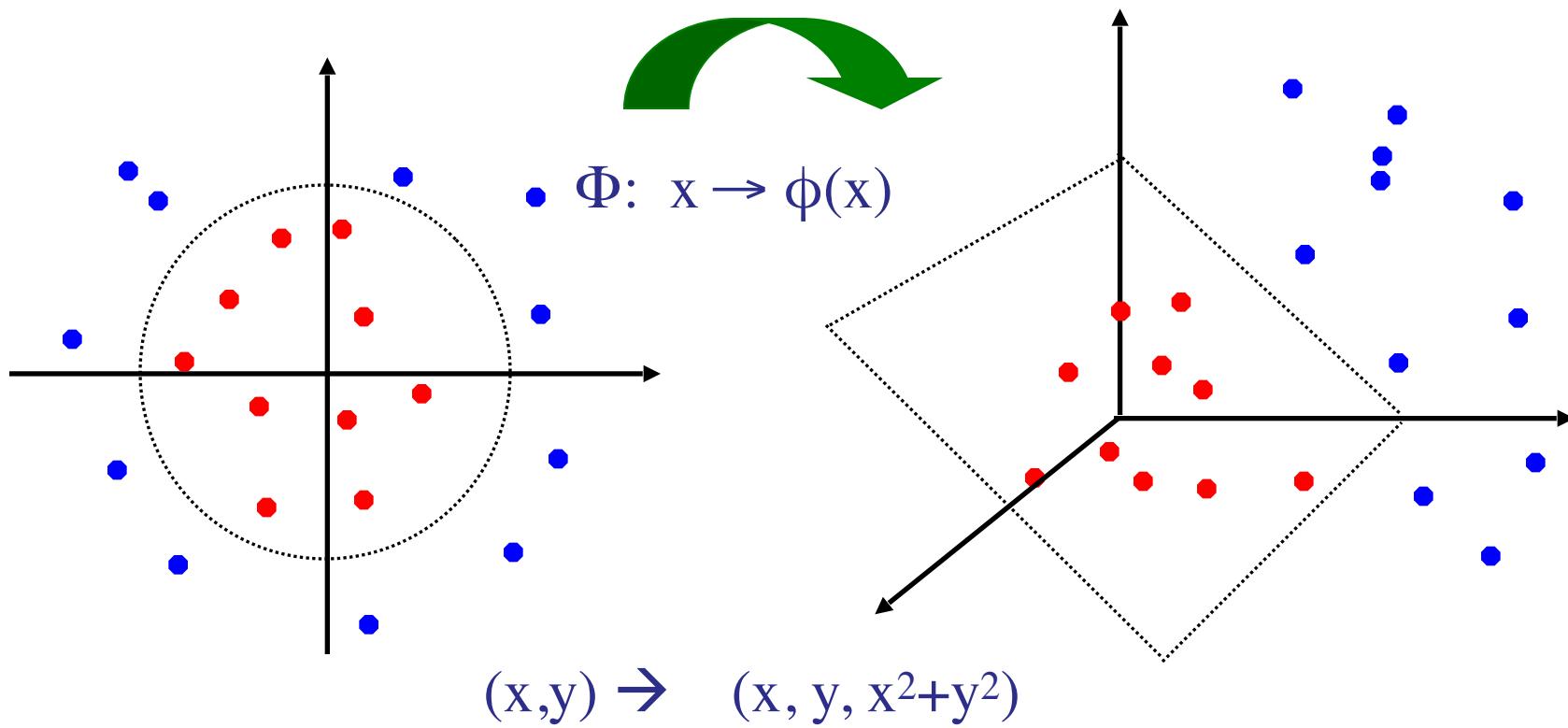


How about... mapping data to a higher-dimensional space:



Non-linear SVMs

Idea: the original input space can be mapped to some higher-dimensional feature space where the training set is separable.



Multi-class classification with SVMs

SVMs can only separate two classes. What if we have more than two?

One vs. all: with m classes, learn m SVMs

- ◊ SVM 1 learns “Output==1” vs “Output != 1”
- ◊ SVM 2 learns “Output==2” vs “Output != 2”
- ◊ ...
- ◊ SVM m learns “Output== m ” vs “Output != m ”

→ use prediction with largest distance to the decision plane

One vs. one: train $\frac{m * (m - 1)}{2}$ “one vs. one” classifiers

- ◊ SVM 1 learns “Output==1” vs “Output == 2”
 - ◊ SVM 2 learns “Output==1” vs “Output == 3”
 - ◊ ...
- choose the class that is predicted most often

Expensive, but often remarkably fast as one vs. one decision might be easier and uses less instances.

SVM advantages and limitations

Advantages

- deals well with high-dimensional, sparse vectors
- **very flexible:** different kernel functions, number of support vectors
- robust classifiers, noise-tolerant
- **overfitting** can be controlled by **soft margin approach**

Limitations

- only two classes: must use **1 vs. all** or other schemes for multiclass problems
- high **computational complexity**
 - but **efficient implementations available**
- **choice of kernel function and its parameters is a trial-and-error enterprise**

Summary

KNN, linear regression and decision trees

- are easy to understand and produce interpretable output.
- struggle with high-dimensional input.
- are prone to overfitting.

Naïve Bayes

- is a fast and often effective algorithm.
- can deal with mixed feature types.
- struggles with correlated features (independence assumption).

SVMs

- can deal with high-dimensional feature spaces.
- have many parameters to tune.

Reading

Mandatory

- <https://monkeylearn.com/blog/introduction-to-support-vector-machines-svm/>
→ advertised as a “7 min read” ☺

Optional

- any good machine learning introduction