

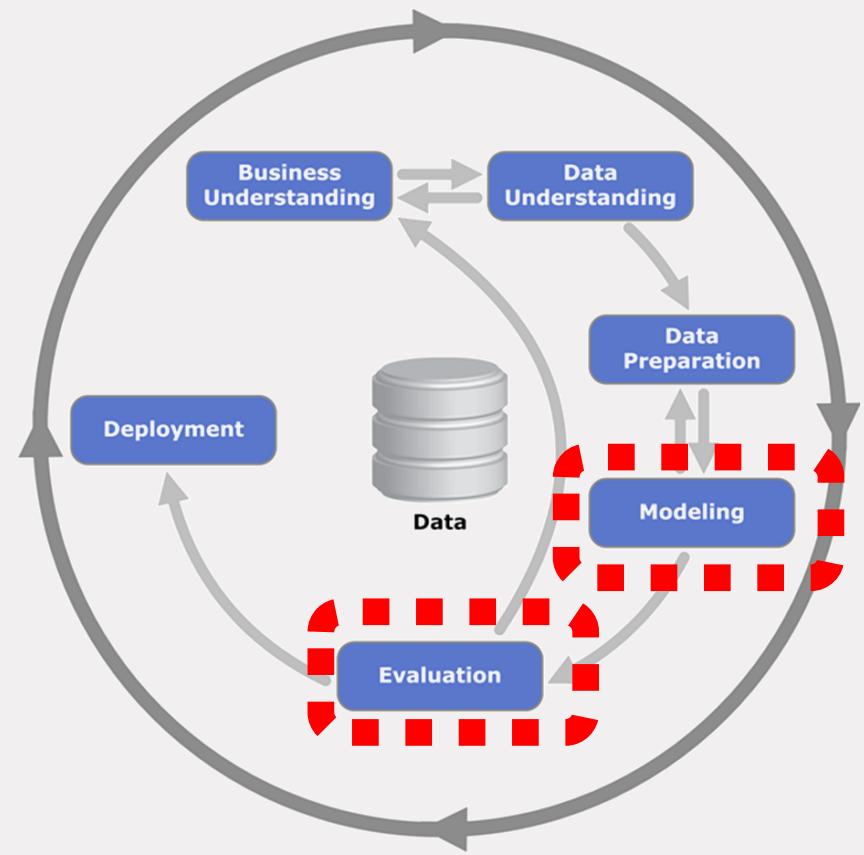
**1BM110: Data driven AI**

**Supervised learning 2: SVMs, Bias-Variance, & Ensemble methods**

**Unsupervised learning: clustering**

Dr. Yingqian Zhang

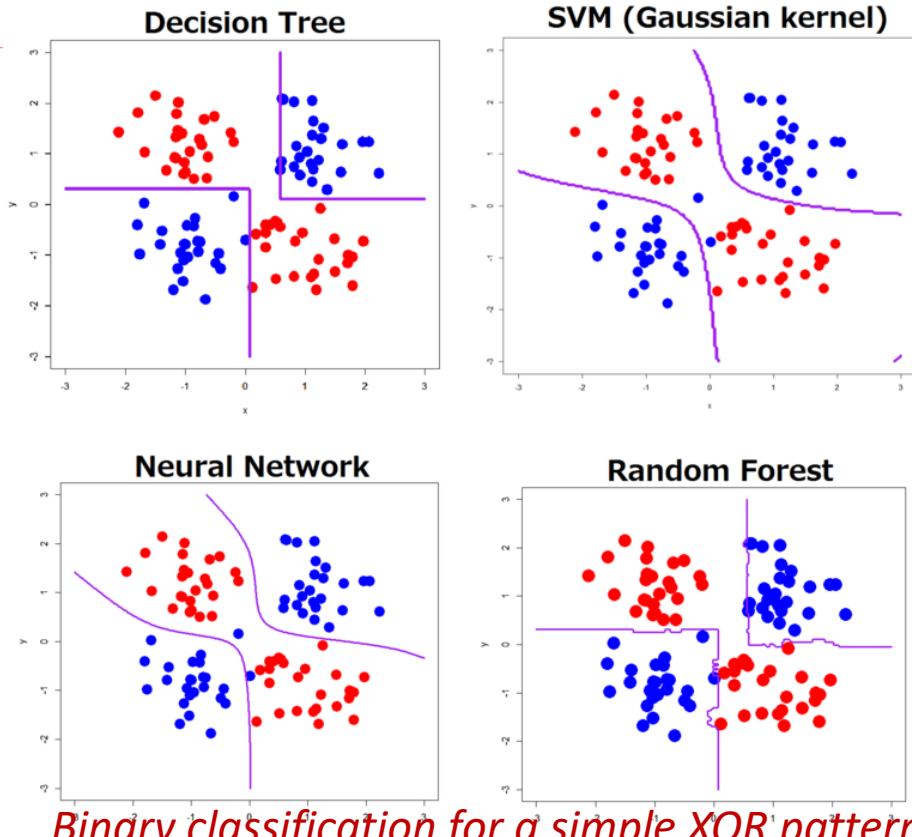
# Previous lecture



- Basic classification/regression
- Performance measurements
- Experimental setup (today)

# Classification Models

- Distance based models
  - K-nearest-neighbour classifier (KNN)
- Probabilistic models
  - Naive Bayes classifiers, Logistic regression
- Hyperplane models
  - Neural networks
  - Support vector machines (SVM)
- Tree-based models
  - Decision tree,
  - Random Forest

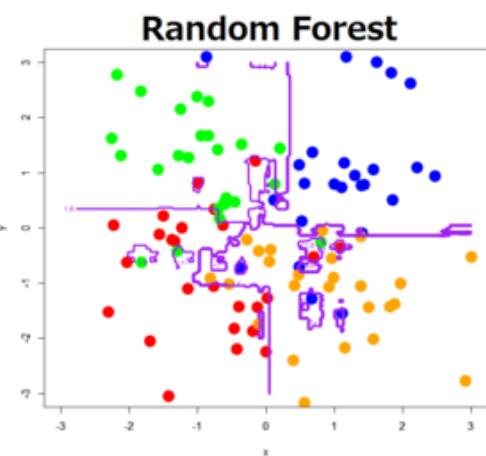
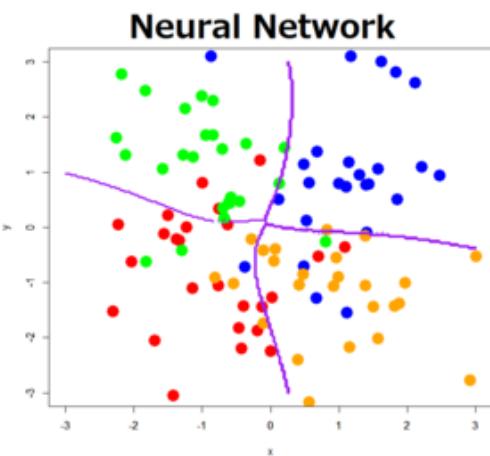
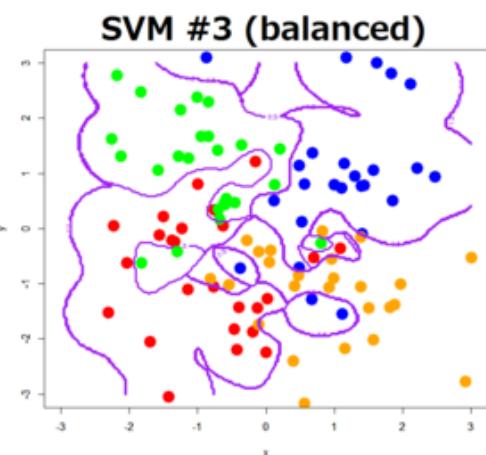
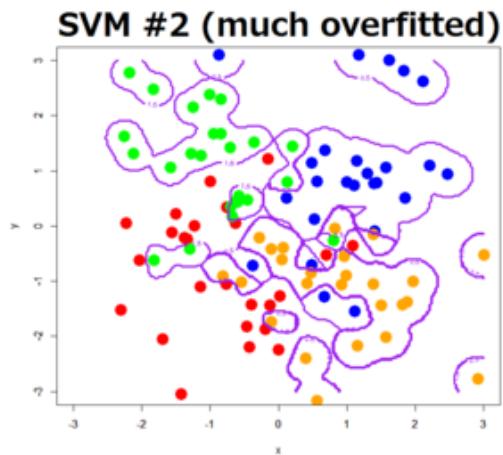
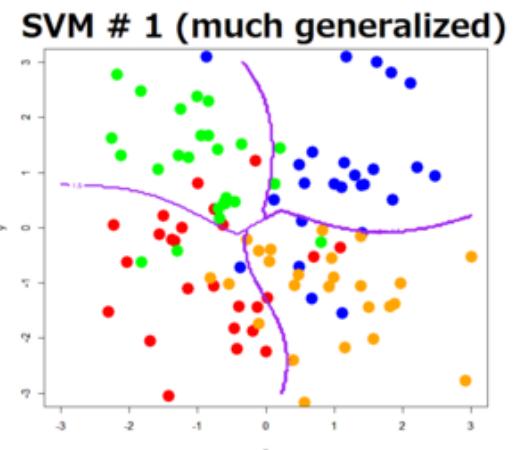
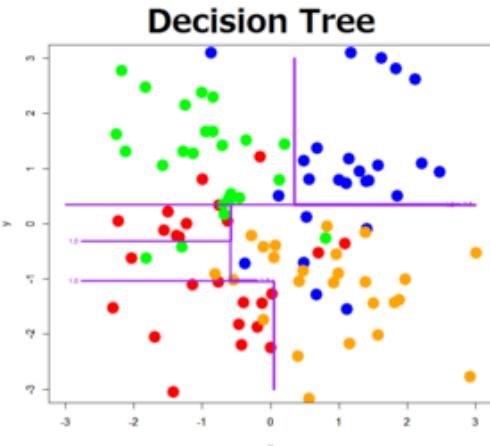


*Binary classification for a simple XOR pattern*

XOR( Exclusive OR)

x1	x2	y
0	0	0
0	1	1
1	0	1
1	1	0

## 4-classes classification for a complex pattern



(figures from the blog of TJO)

## Support Vector Machines (SVMs)

# Hyperplane Classifiers - SVM

Data set:

$D: (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_i, y_i) \dots$

$\mathbf{x}_i$ : training set with  
associated class labels

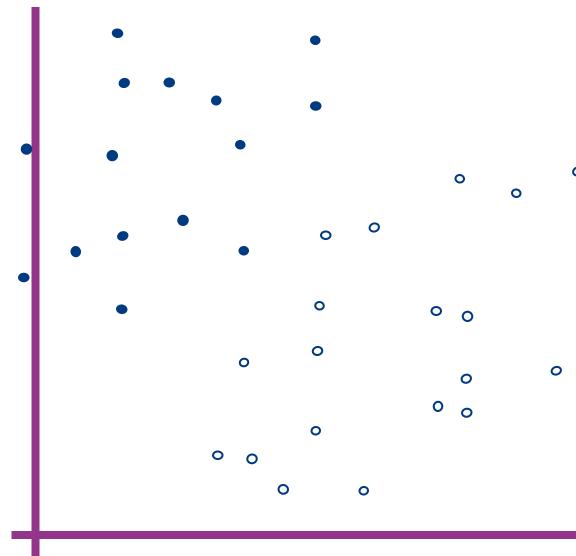
$y_i$

Each data point  $\mathbf{x}$  has  
2 input variables

Binary classification:  
 $y_i \in \{+1, -1\}$

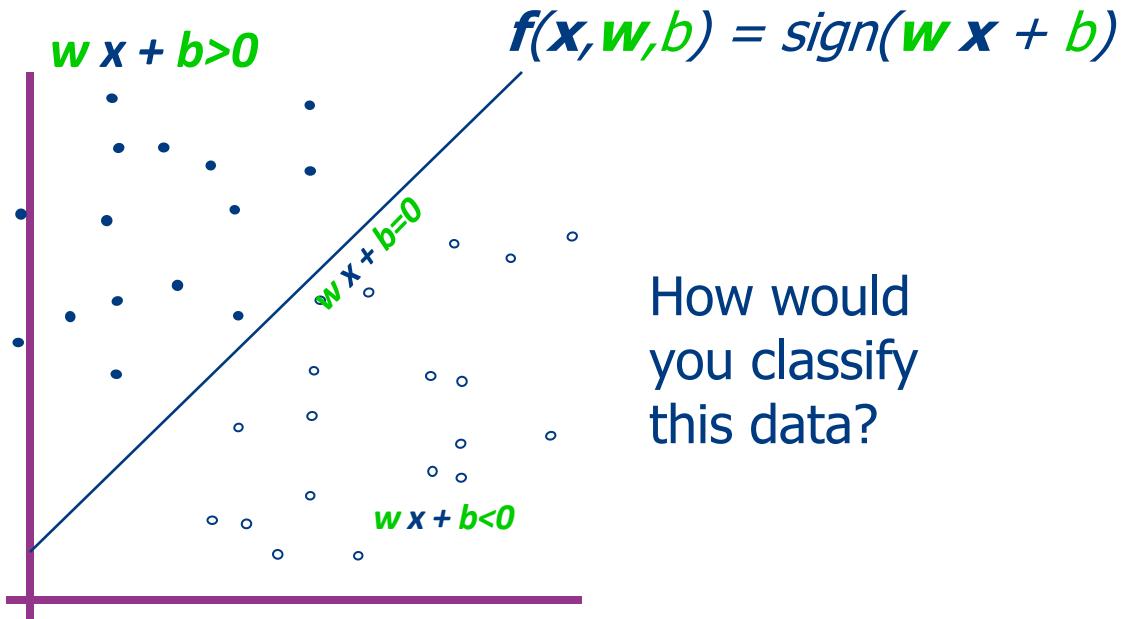
Solid circle: +1

Empty circle: -1



How would  
you classify  
this data?

# Hyperplane Classifiers - SVM



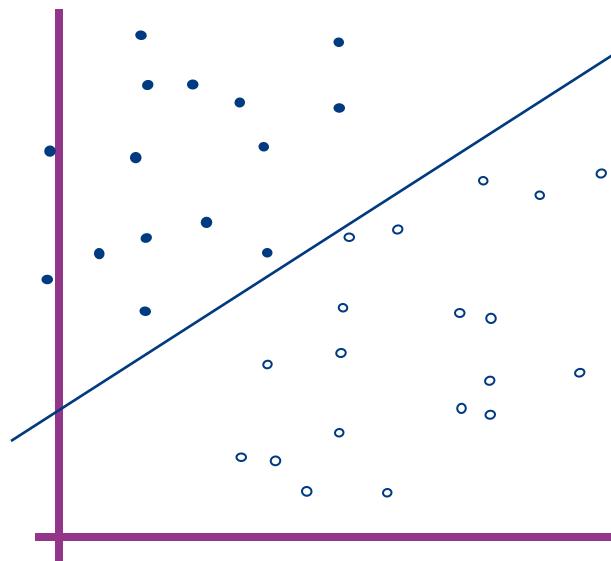
w: weight vector/slope

x: input vector

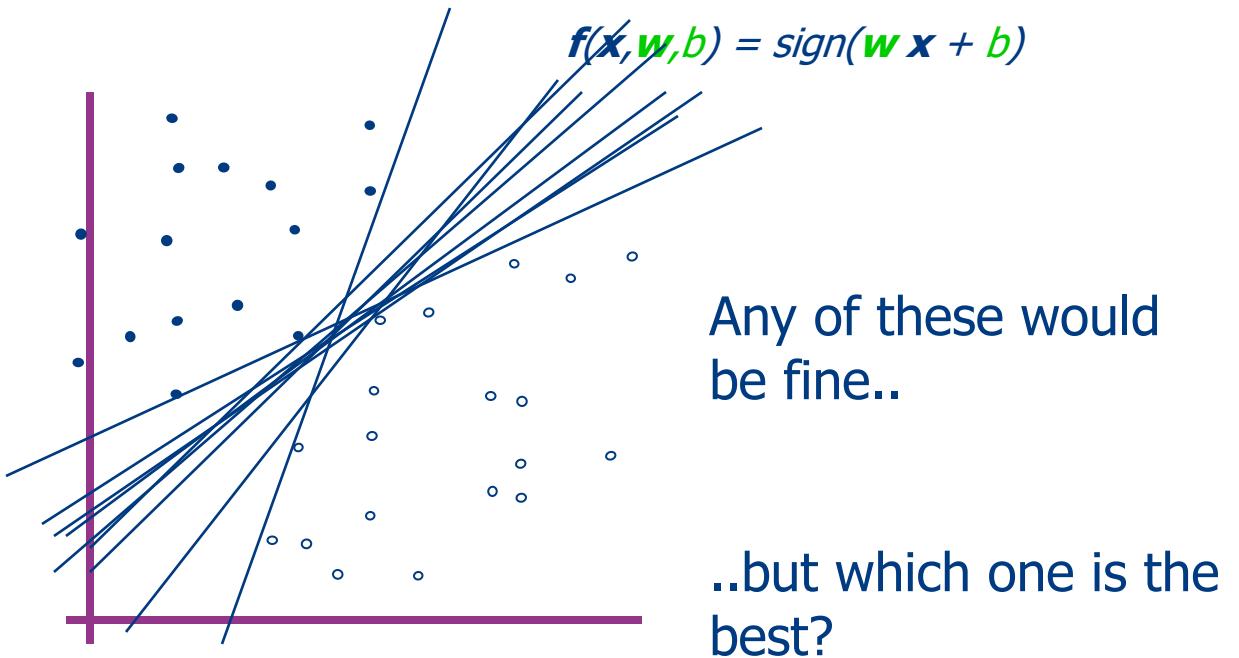
b: bias/intercept

# Hyperplane Classifiers - SVM

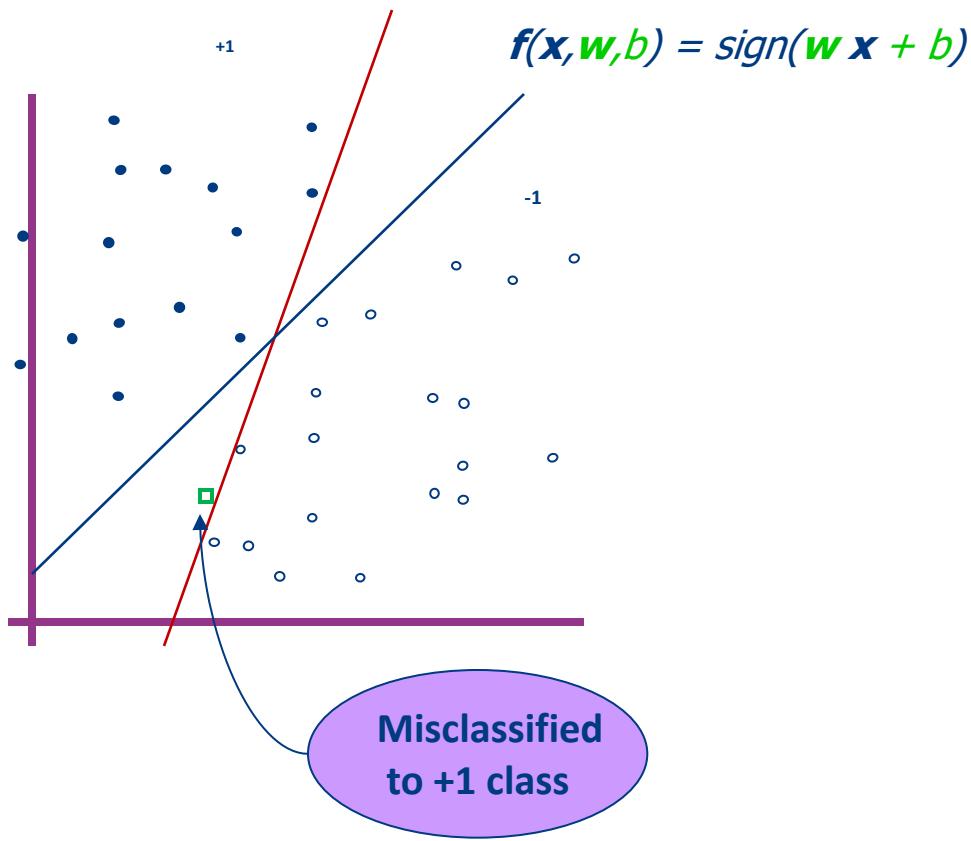
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$



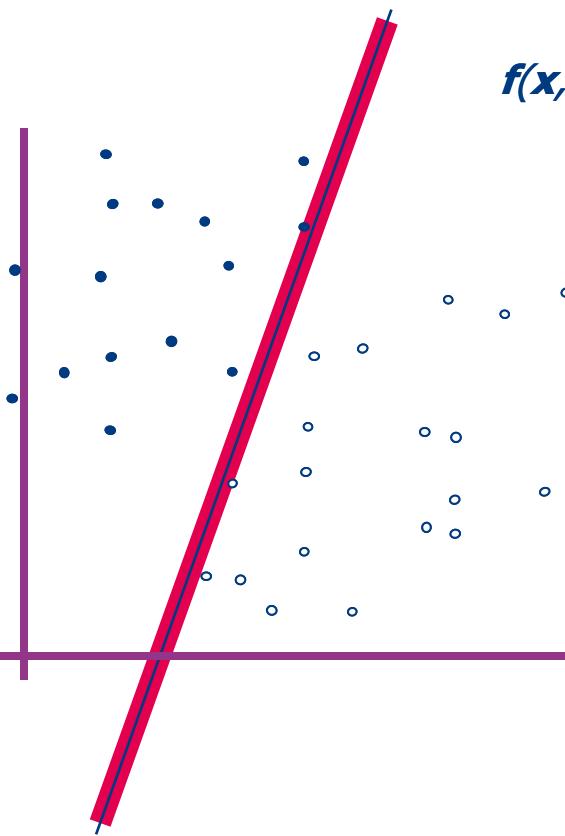
# Hyperplane Classifiers - SVM



# Hyperplane Classifiers - SVM



# Hyperplane Classifiers – Support Vector Machine



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

SVM searches for the maximum-margin line.

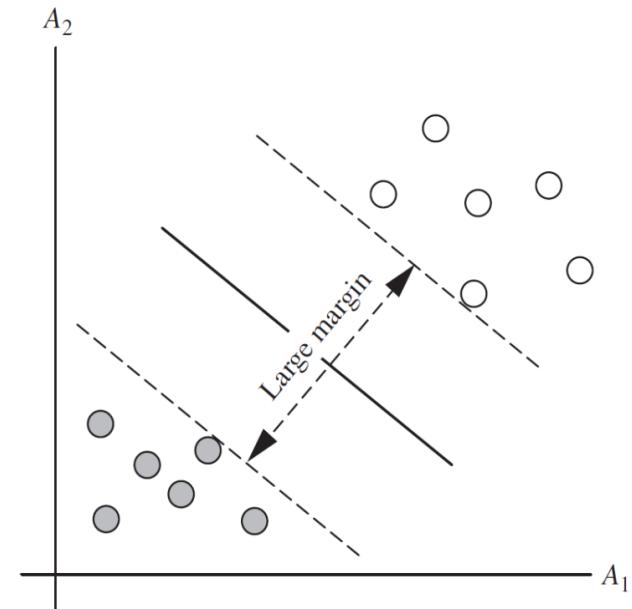
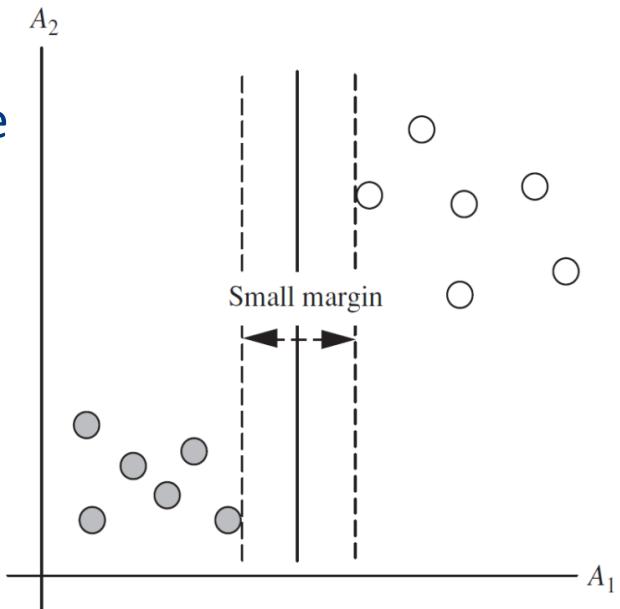
Define the margin of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

# Hyperplane Classifiers - SVM

SVM:

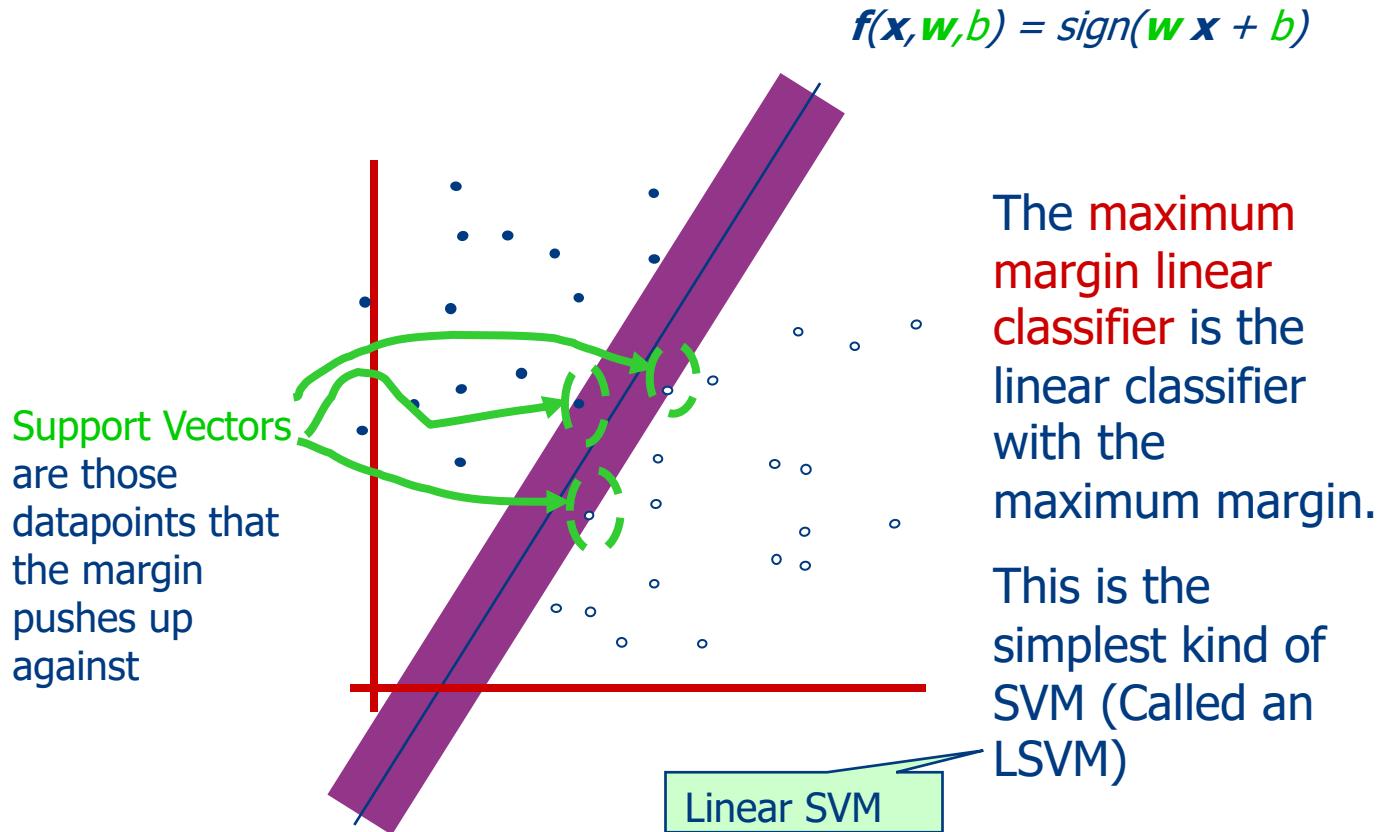
Finding maximum margin line/hyperplane

Larger margin  
-> larger separation



Source: Han et al.

# Hyperplane Classifiers - SVM



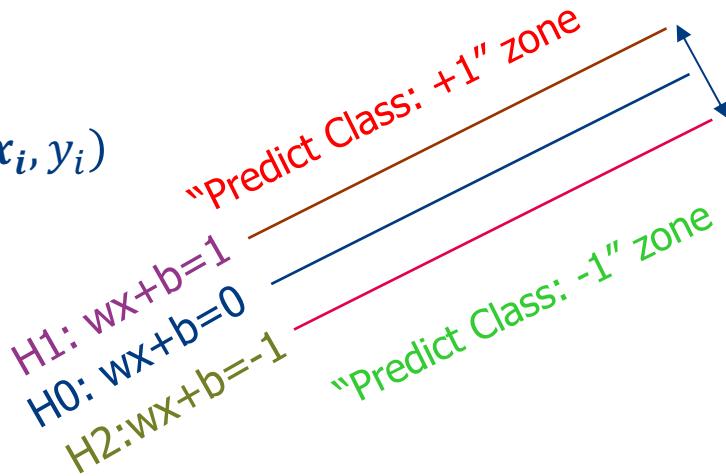
# Linear SVM: mathematically

Data set:

$$D: (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_i, y_i)$$

Binary classification:

$$y_i \in \{+1, -1\}$$



**M**=Margin Width

Separating hyperplane:

$$\mathbf{w} \cdot \mathbf{x} + b = 0 \quad \mathbf{w}: (w_1, w_2)$$

Recall the distance from point  $(x_1, x_2)$  to a line:  $\mathbf{w} \cdot \mathbf{x} + b = 0$  is:

$$|x_1w_1 + x_2w_2 + b| / \sqrt{w_1^2 + w_2^2}$$

So, the distance between H0 and H1 is:  $\frac{1}{\sqrt{w_1^2 + w_2^2}} = \frac{1}{\|\mathbf{w}\|}$

The distance between points on H1 and H2 is the margin width M:  $\frac{2}{\|\mathbf{w}\|}$

To maximize M, we need to minimize  $\|\mathbf{w}\|$ , with the condition no points between H1 and H2

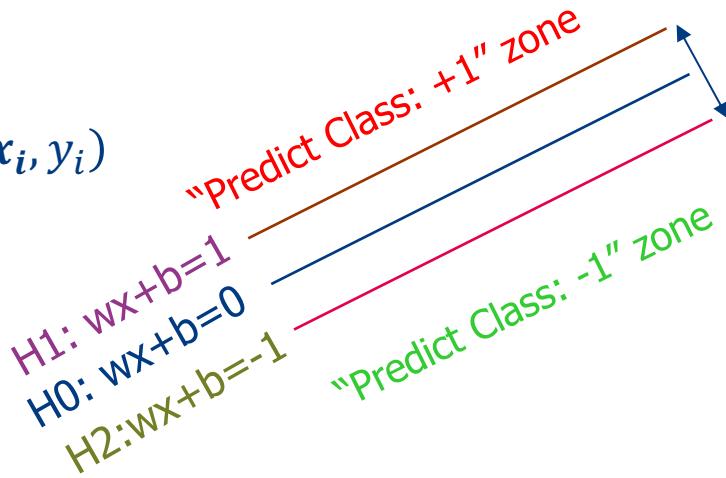
# Linear SVM: mathematically

Data set:

$$D: (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_i, y_i)$$

Binary classification:

$$y_i \in \{+1, -1\}$$



**M**=Margin Width

Separating hyperplane:

$$\mathbf{w} \mathbf{x} + b = 0 \quad \mathbf{w}: (w_1, w_2)$$

The distance between points on H1 and H2 is the margin width M:  $\frac{2}{\|\mathbf{w}\|}$

Optimization problem: Minimize  $f: \frac{1}{2} \|\mathbf{w}\|^2$ , subject to

$$\mathbf{w} \mathbf{x}_i + b \geq +1 \text{ when } y_i = +1$$

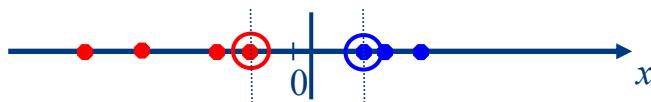
$$\mathbf{w} \mathbf{x}_i + b \leq -1 \text{ when } y_i = -1$$

$$\Rightarrow y_i(\mathbf{w} \mathbf{x}_i + b) \geq 1, \text{ for all } i$$

This is a quadratic (expensive!) optimization problem

# Non-linear SVMs

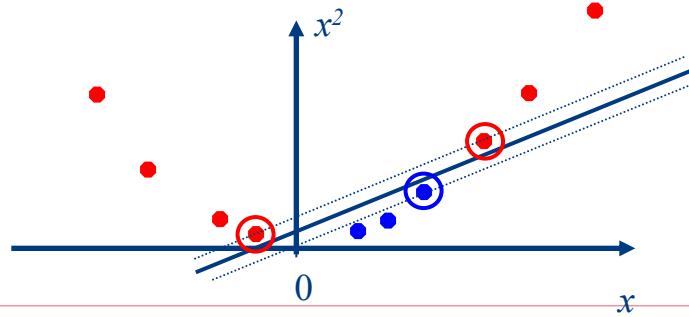
- Datasets that are linearly separable work out great:



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

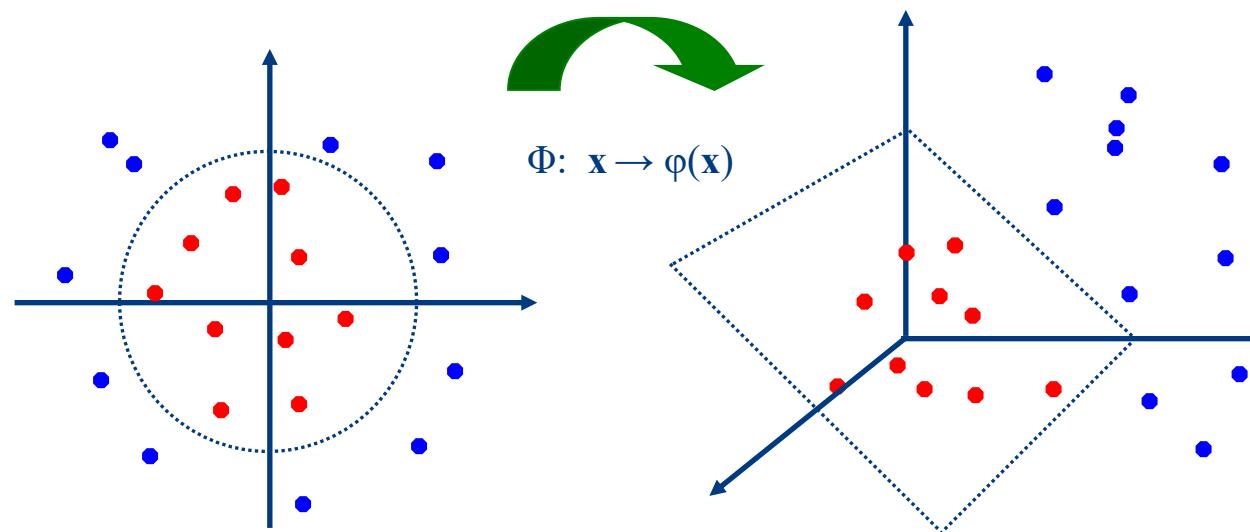


- How about... mapping data to a higher-dimensional space (1-D to 2-D):



# Non-linear SVMs: Feature spaces

- The approach for linear SVMs can be extended to create *nonlinear SVMs* for the classification of *nonlinear data*.
- The original input space can always be mapped to some higher-dimensional feature space where the training set is **linearly separable**.



# Non-linear SVMs

Two main steps :

1. Transform the original input data into a higher dimensional space using a nonlinear mapping.
2. Then, search for a linear separating hyperplane in the new space. We end up with a quadratic optimization problem that can be solved using the linear SVM formulation.

The maximal margin hyperplane found in the new space corresponds to a nonlinear separating hypersurface in the original space.

# The “Kernel Trick”

- Difficulties
  - How to choose this nonlinear mapping?
  - The computation becomes very costly!
- Notice: the training samples appear only in the form of dot products in solving the quadratic optimization problem:  
 $\varphi(x_i) \cdot \varphi(x_j)$ , where  $\varphi(x)$  is the nonlinear mapping
  - equivalent to apply a kernel function  $K(x_i, x_j)$  to the original input data:  
$$K(x_i, x_j) = \varphi(x_i) \cdot \varphi(x_j)$$
- Kernel trick: everywhere that  $\varphi(x_i) \cdot \varphi(x_j)$  appears in the training algorithm, we can replace it with  $K(x_i, x_j)$ 
  - All calculations are in the original (lower dimension) input space
  - Safely avoid mapping

# Examples of admissible kernel functions

Different kernel functions, each resulting in a different nonlinear classifier in the original input space.

- Polynomial kernel of degree  $h$ :

$$K(X_i, X_j) = (X_i \cdot X_j + 1)^h$$

- Gaussian (radial-basis function network):

$$K(X_i, X_j) = e^{-\|X_i - X_j\|^2 / 2\sigma^2}$$

- Sigmoid kernel:

$$K(X_i, X_j) = \tanh(\kappa X_i \cdot X_j - \delta)$$

- No golden selection rule; is a parameter we can tune in SVM

# SVMs

- Can be used for classification & regression problems.
- Highly accurate model (global optimization)
- Only support vectors are used to specify the separating hyperplane
  - ability to handle large feature spaces
  - less prone to overfitting than most of other methods
  - but, sensitive to noise
- Computationally expensive (slow)
  - need to solve an expensive quadratic optimization problem

## Bias-variance trade-off

# Bias of Learning Methods

Recall supervised learning methods try to find a function  $f: X \rightarrow Y$ , from training data.

- Bias: refers to the error of the estimation of  $f$ , usually introduced by modeling a real life problem by a much simpler problem.
  - For example, linear regression assumes that there is a linear relationship between output  $Y$  and input  $X$ . It is unlikely, in real life, the relationship is exactly linear so some bias will be present.
- The more complex a model is, the less bias it will generally have.

# Variance of Learning Methods

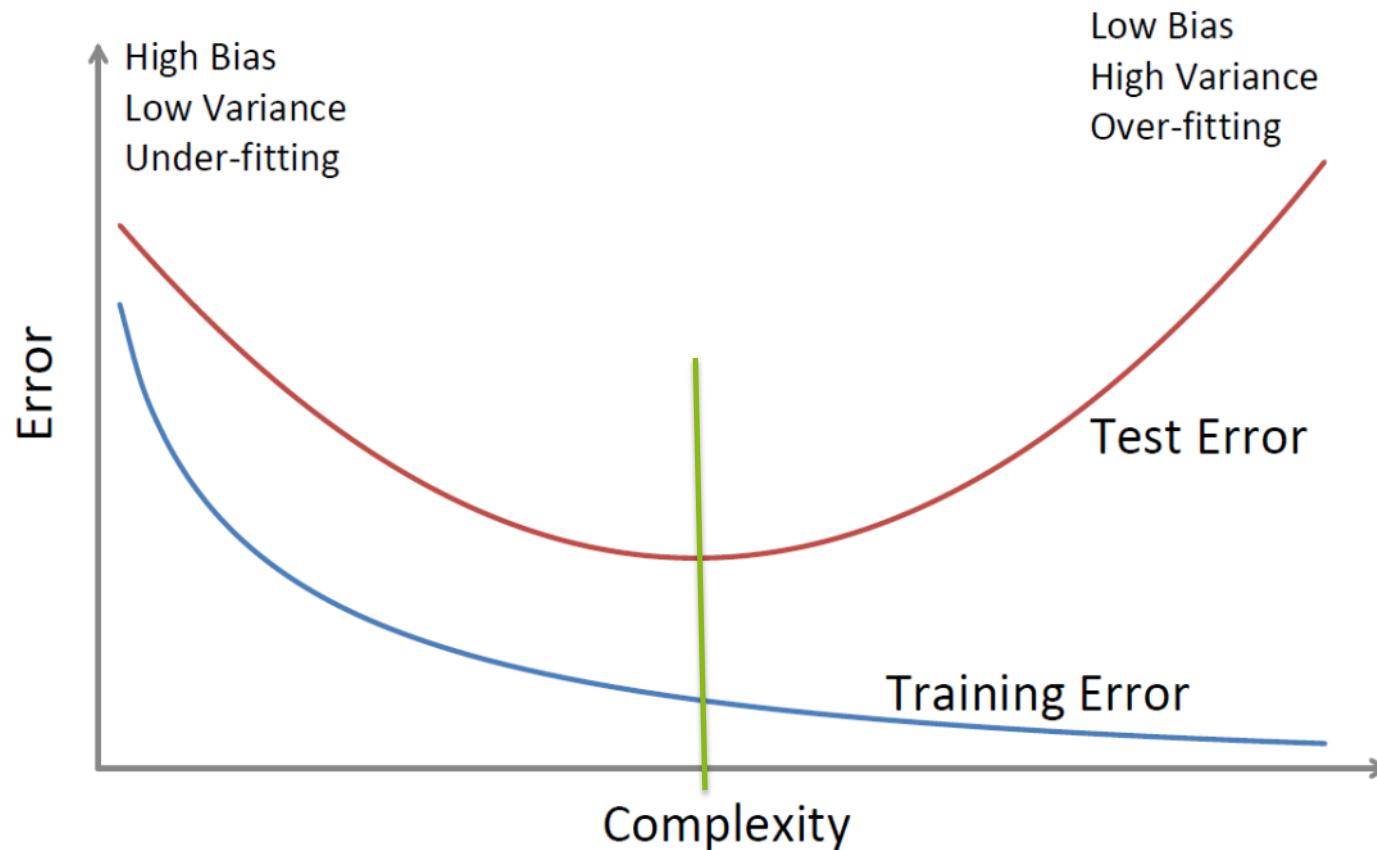
- Variance: refers to how much your estimate for function  $f$  would change if you had a different training data set.
- Generally, the better it can fit the train data, the more variance it has.

*Bias and variance are two competing forces that govern the choice of learning method/model*

# Bias/ Variance Tradeoff

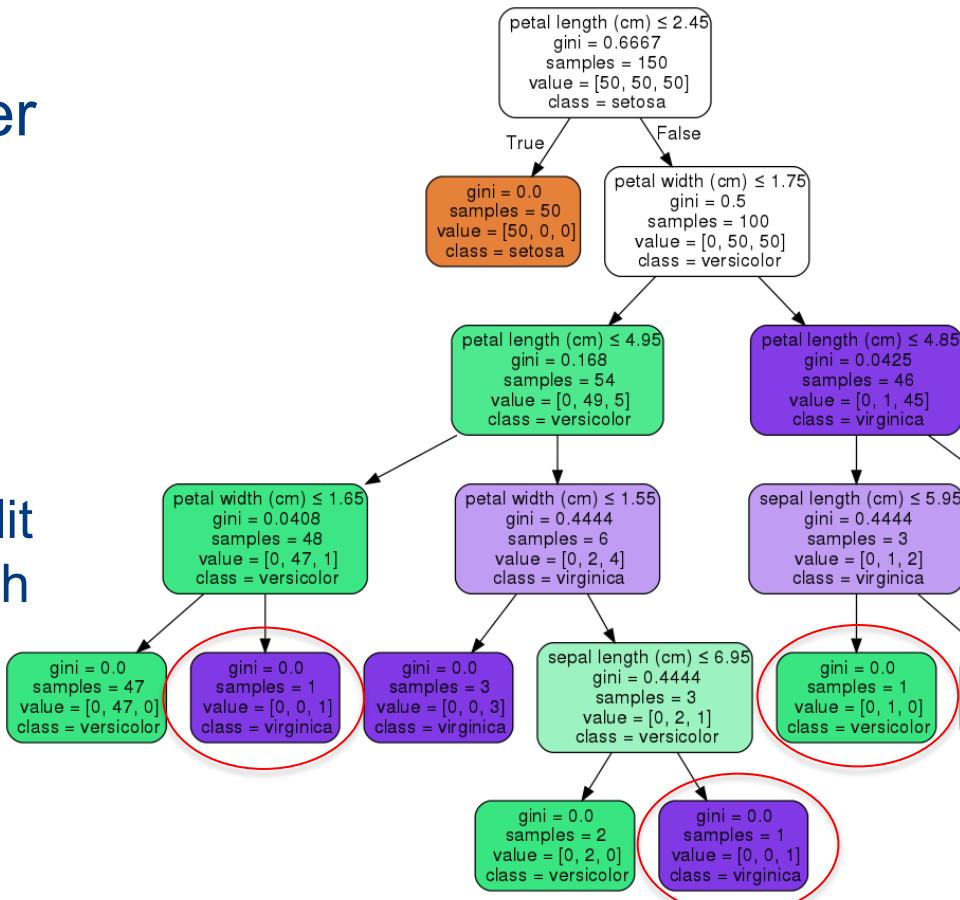
- How complex you want your model to be?
- The complexity of a model is measured by the **number of the free model parameters** which are estimated
- Thinking of a decision tree of depths 2, or 20
  - A model with **high complexity** will result in ***overfitting***: the **bias of the model is small** but the parameters are estimated with **high variance**
  - A simple model with **low complexity** will result in ***underfitting***: the parameters are estimated with **low variance** but the model has a **large bias**
  - In both cases the test error will be rather large as shown in the figure on the next slide

# Bias/ Variance Tradeoff



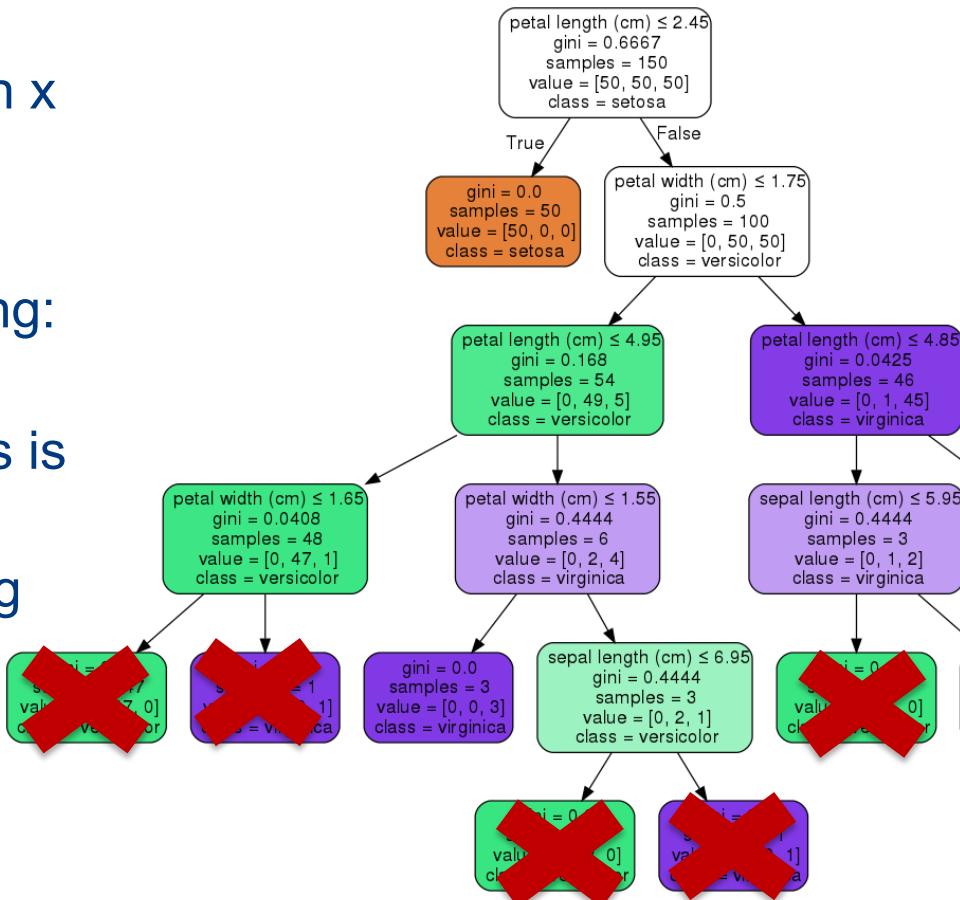
# Avoiding overfitting – decision trees

- Limiting the tree depth
- Not splitting nodes with fewer than x samples
- Post-pruning the tree, after learning:
  - Carefully evaluate whether a split resulting in two leaf nodes is worth the extra nodes
  - Trade-off between fit to training data and number of model parameters



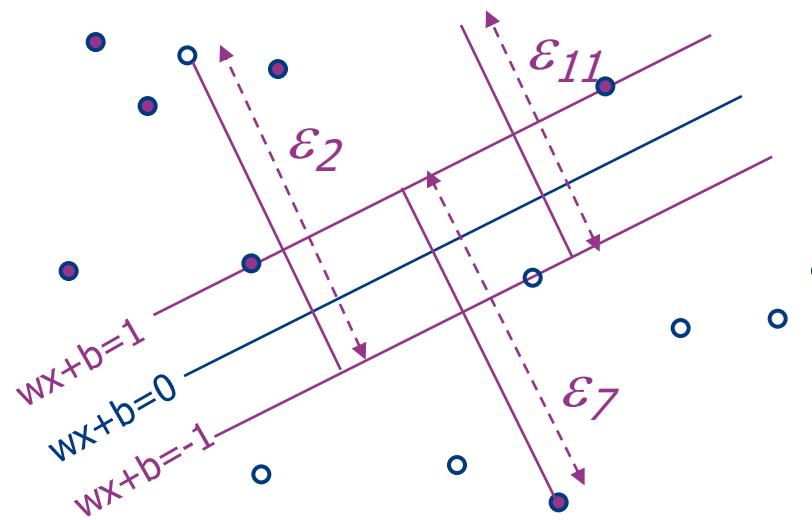
# Avoiding overfitting – decision trees

- Limiting the tree depth
- Not splitting nodes with fewer than  $x$  samples
- Post-pruning the tree, after learning:
  - Carefully evaluate whether a split resulting in two leaf nodes is worth the extra nodes
  - Trade-off between fit to training data and number of model parameters



# Avoiding overfitting - SVM

*Slack variables  $\varepsilon_i$  can be added to allow misclassification of difficult or noisy examples. Called *soft-margin*.*



## Ensemble methods

# Decision Trees: problem

- Decision trees suffer from high variance
  - If we randomly split the training data into 2 parts, and fit decision trees on both parts, the results could be quite different
- We would like to have models with low variance
- To solve this problem, we can use bagging (bootstrap aggregating).

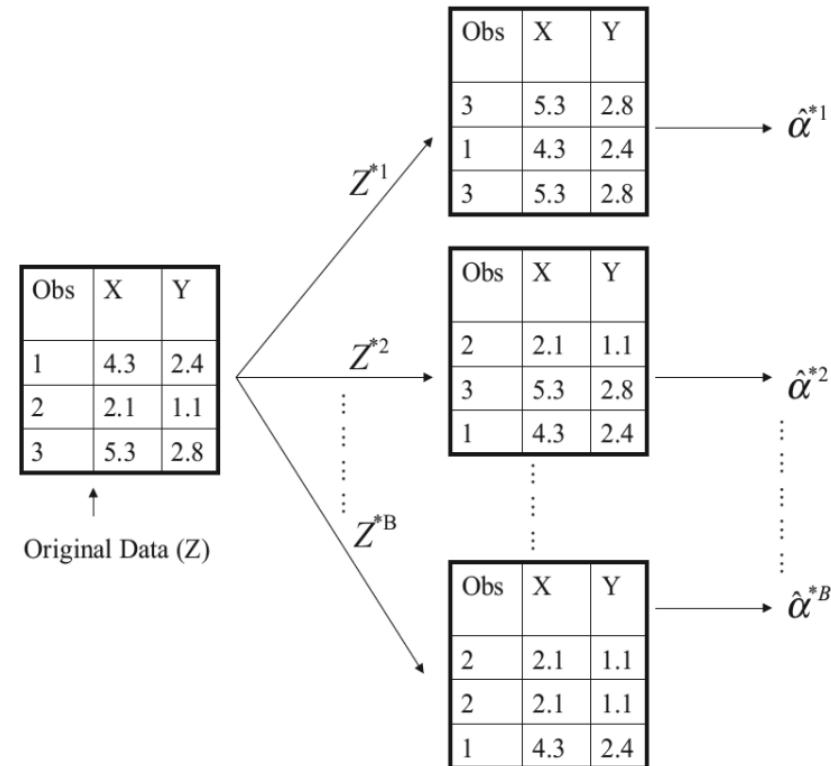
# How does bagging work?

## Bootstrapping: sampling method

- Given a training set  $D$  of size  $n$ , it generates  $B$  new training sets  $D_1, D_2, \dots, D_B$ , each of size  $m$ , by sampling from  $D$  uniformly and with replacement.

## Then, Bagging

- Generate  $B$  different bootstrapped training datasets
- Train the learning method on each of the  $B$  training datasets, and obtain the prediction
- For prediction:
  - Regression: average all predictions from all  $B$  trees
  - Classification: majority vote among all  $B$  trees, or average the probabilities of  $B$  trees



# How does bagging work?

- Bagging is an extremely powerful idea based on two things:
  - Bootstrapping: plenty of training datasets!
  - Averaging: reduces variance!
- Why does averaging reduces variance?
  - Averaging a set of observations reduces variance.  
Recall that given a set of  $n$  independent observations  $Z_1, \dots, Z_n$ , each with variance  $\sigma^2$ , the variance of the mean  $\bar{z}$  of the observations is given by  $\sigma^2/n$

# Bagging for Regression Trees

- Construct  $B$  regression trees using  $B$  bootstrapped training datasets
- Average the resulting predictions
- Note: These trees are not pruned, so each individual tree has **high variance but low bias**. Averaging these trees reduces variance, and thus we end up lowering both **variance and bias!**

# Bagging for Classification Trees

- Construct  $B$  classification trees using  $B$  bootstrapped training datasets
- For prediction, there are two approaches:
  1. Record the class that each bootstrapped data set predicts and provide an overall prediction to the most commonly occurring one (majority vote).
  2. If our classifier produces probability estimates we can just average the probabilities and then predict to the class with the highest probability.
- Both methods work well.

# Random Forests

- It is a very efficient learning method
- It builds on the idea of bagging, but it provides an improvement because it **de-correlates** the trees
- How does it work?
  - Build a number of decision trees on bootstrapped training datasets, but each time a split in a tree is considered, a random sample of  $m$  predictors is chosen as split candidates from the full set of  $p$  predictors (Usually  $m \approx \sqrt{p}$  )

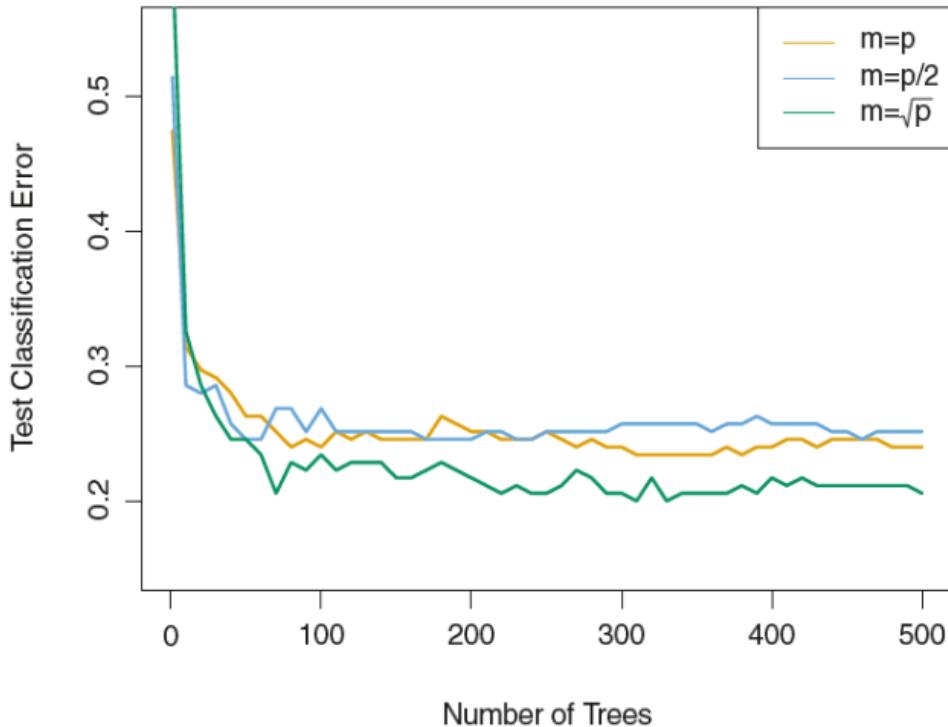
## Why are we considering a random sample of $m$ predictors instead of all $p$ predictors for splitting?

- Suppose that we have a very strong predictor in the data set along with a number of other moderately strong predictor, then in the collection of bagged trees, **most or all of them will use the very strong predictor for the first split!**
- Then, all bagged trees will look similar. Hence all the predictions from the bagged trees will be **highly correlated**
- Averaging many highly correlated quantities does not lead to a large variance reduction
- Random forests “de-correlates” the bagged trees leading to **more reduction in variance**

# Random Forest with different values of “m”

*Results from RF for the 15-class gene expression dataset with  $p = 500$  predictors.*

- Notice when random forests are built using  $m = p$ , then this amounts simply to bagging.
- *A single classification tree has an error rate of 45.7%.*

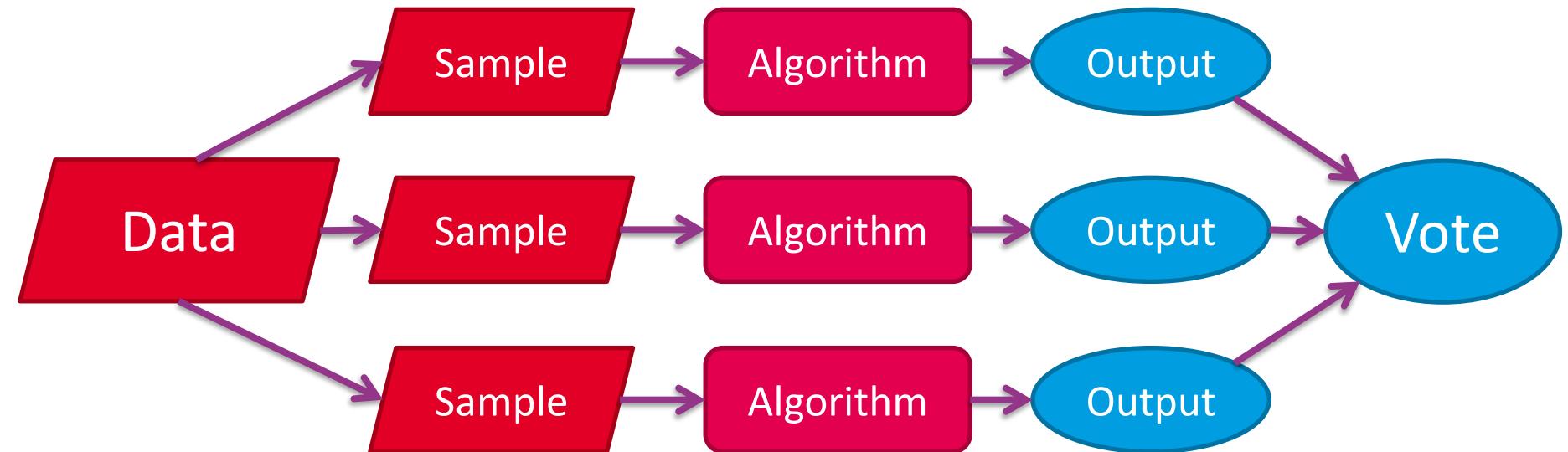


*m: number of features considered in node splitting decisions*

# More general: ensemble methods

- Key idea:
  - *Learn  $N$  models, and combine their predictions*
- Different strategies:
  - How to learn **different** models?
  - How to **combine** them?

# Bagging



- Learn one model by sampling with replacement and combining results by majority voting
- Example: random forest

# Basic random forest

## 1. Do X times:

1. Create data sample D
2. Learn tree T from D, but:
  - At every split, out of m features, consider a random feature subset of size  $\sqrt{m}$

## 2. For data row d:

1. Classify d using all X trees
2. Assign label using majority vote

# Extreme random forest

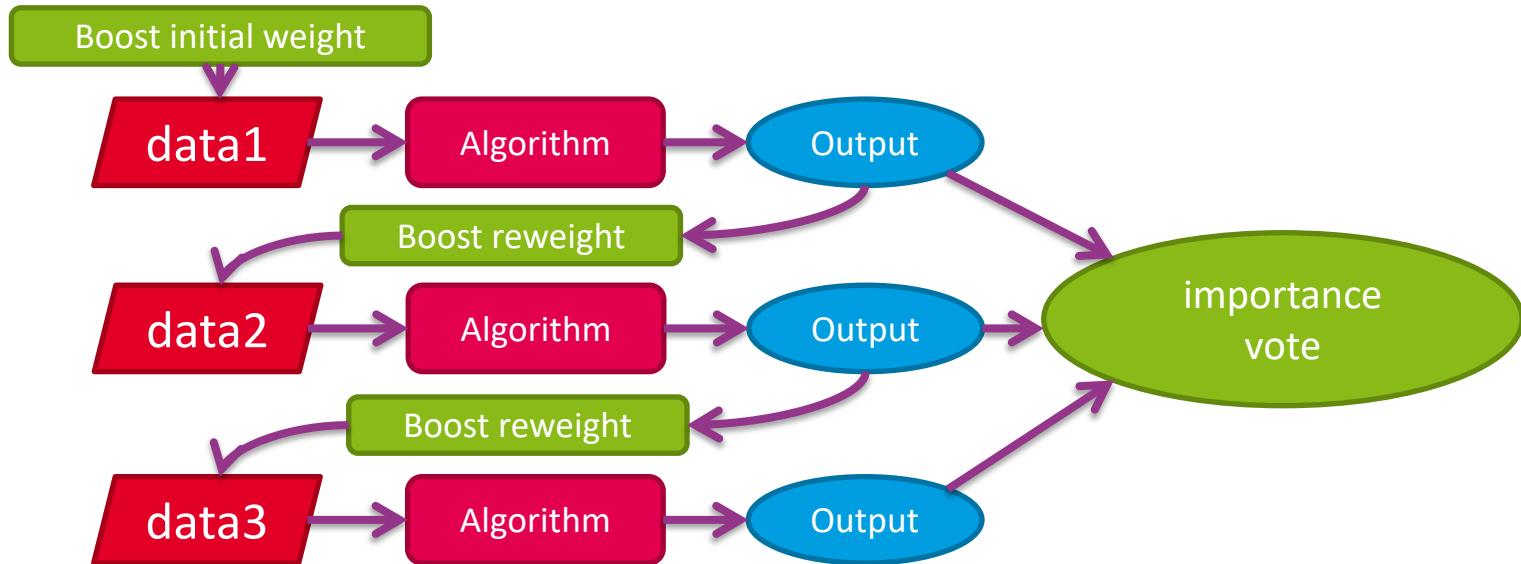
## 1. Do X times:

1. Create data sample D
2. Learn tree T from D, but:
  - At every split, out of m features, consider a random feature subset of size  $\sqrt{m}$
  - Randomize the heuristic (eg. select randomly one of Gini, information gain, ...)

## 2. For data row d:

1. Classify d using all X trees
2. Assign label using majority vote

# Boosting: AdaBoost



Combine multiple “weak classifiers” into a single strong classifier.

- Start with a weak classifier (slightly better than random guess for binary classification) using equal weights for observations

Repeat till terminate:

- Evaluate outputs, and then adjust the weights for observations: misclassified data get higher weights
- Train another classifier

Combine predictions from all classifiers

# AdaBoost for binary classification

- Initialize (equal) weights for n observations:

$$w = \frac{1}{n}$$

- For each boosting iteration t, calculate the actual influence (model importance) for this classifier in classifying the data points using

$$\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$$

$\epsilon_t$ : error rate (better than chance, hence between 0 and 0.5)

- Update weights for each data point i:

$$w_{t+1}(i) = \frac{w_t(i)}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } i \text{ was correctly classified} \\ e^{\alpha_t} & \text{otherwise} \end{cases}$$

$Z_t$  : normalization term, so the new weights sum to 1

# Boosting

- AdaBoost can be applied on top of any classifier; usually called the “best out-of-the-box classifier”.
  - Has small bias
  - Almost no parameters to tune, except iteration number
  - Provably effective
    - We simply need to find (weak) classifiers slightly better than random guessing
  - Sensitive to noisy data and outliers
- Other extremely popular boosting algorithms
  - XGBoost (eXtreme Gradient Boosting)
  - LightGBM (Light Gradient Boosting Machine)
  - Winners of many Kaggle competitions!

# Ensemble methods

- Boost a weak classifier to a strong classifier;  
mainly reduces bias.
- Bagging reduces variance.
- Beware of the black-box nature of ensembles
  - Can you still explain the outcome?

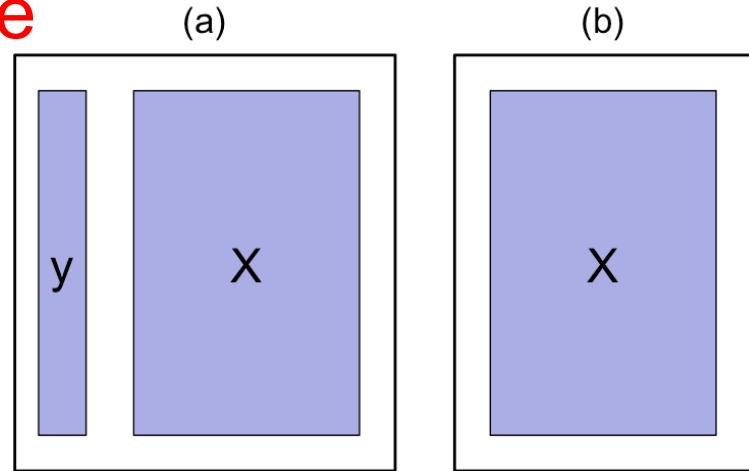
# Which methods to use?

- Is the learning problem easy or difficult?
- What is the most important criterion for your prediction model? Performance? Explanability?
- Do you data satisfy the assumptions of the chosen machine learning methods?
- Bias/ variance tradeoff
- ... ...

# Unsupervised Learning

# Unsupervised learning

- No output (target) variable



Supervised learning      Unsupervised learning

- The goal of unsupervised learning is to discover interesting patterns about data:
  - Can we discover subgroups among the observations?

# Unsupervised learning: clustering

- Goal of *clustering*: partitioning data into subsets (clusters), such that the observations in the clusters are rather homogeneous
  - A good clustering is one when the observations within a group are similar but between groups are very different

# Clustering

- Popular way to identify behaviors in business
- Example: online auction **bidding behavior** data
  - Each row represents a bidder
  - Each record has three *features*
  - *There is no output variable*
  - *Use clustering to find out what types of bidders were there → bidding behaviors!*

Total Number of Placed Bids	Time of First Bid	Time of Last Bid
3	10	20
5	200	235
20	20	200
40	5	220
2	24	50

Q1: Bidder 1 and Bidder 2: do they bid similarly?

Q2: which bidder is similar to Bidder 1 in terms of bidding behavior?

# Clustering: distances

How to measure “similarity” of variable values?

- For metric variables the most important distance is the Euclidean distance

$$d^2(x, z) = \sum_{j=1}^p (x_j - z_j)^2 = \|x - z\|^2$$

x and z are p-dimensional vectors

- For qualitative variables, Hamming distance is often used.
  - Translate variables to binary variables
  - Compute the number of different bits

Example: three observations with ‘temperature’ and ‘weather’

$x_1=(\text{high}, \text{sunny})$ ,  $x_2=(\text{low}, \text{cloudy})$ ,  $x_3=(\text{low}, \text{sunny})$

four dummy variables: temp\_low, temp\_high, w\_sunny, w\_cloudy

$X_1: (0,1,1,0)$ ,  $X_2: (1,0,0,1)$ ,  $X_3: (1,0,1,0)$

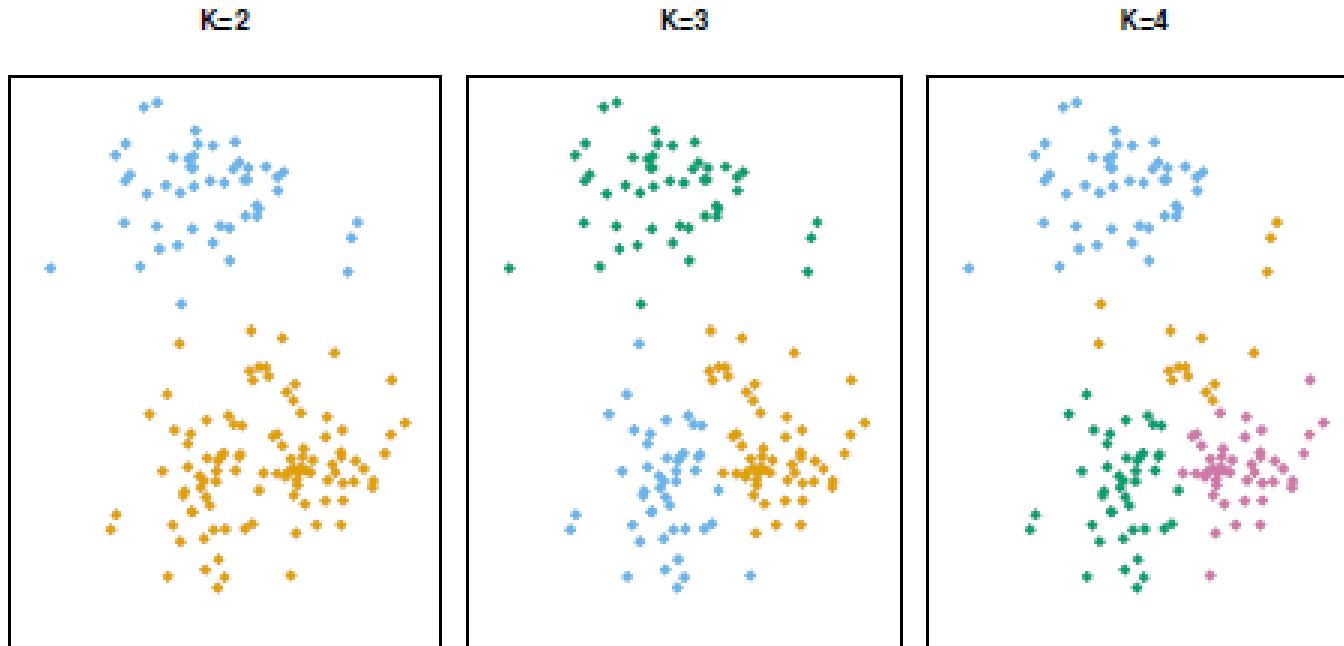
Distance between  $X_1$  (0110) and  $X_2$  (1001): 4

Between  $X_1$  (0110) and  $X_3$  (1010): 2

# Clustering algorithms: two famous ones

- K-means clustering
  - partition the observations into a pre-specified number of clusters.
- Hierarchical clustering
  - it ends up with a tree-like visual representation of the observations, called a **dendrogram**, that allows us to view at once the clusters obtained for each possible number of clusters, from 1 to n.

# K-means clustering



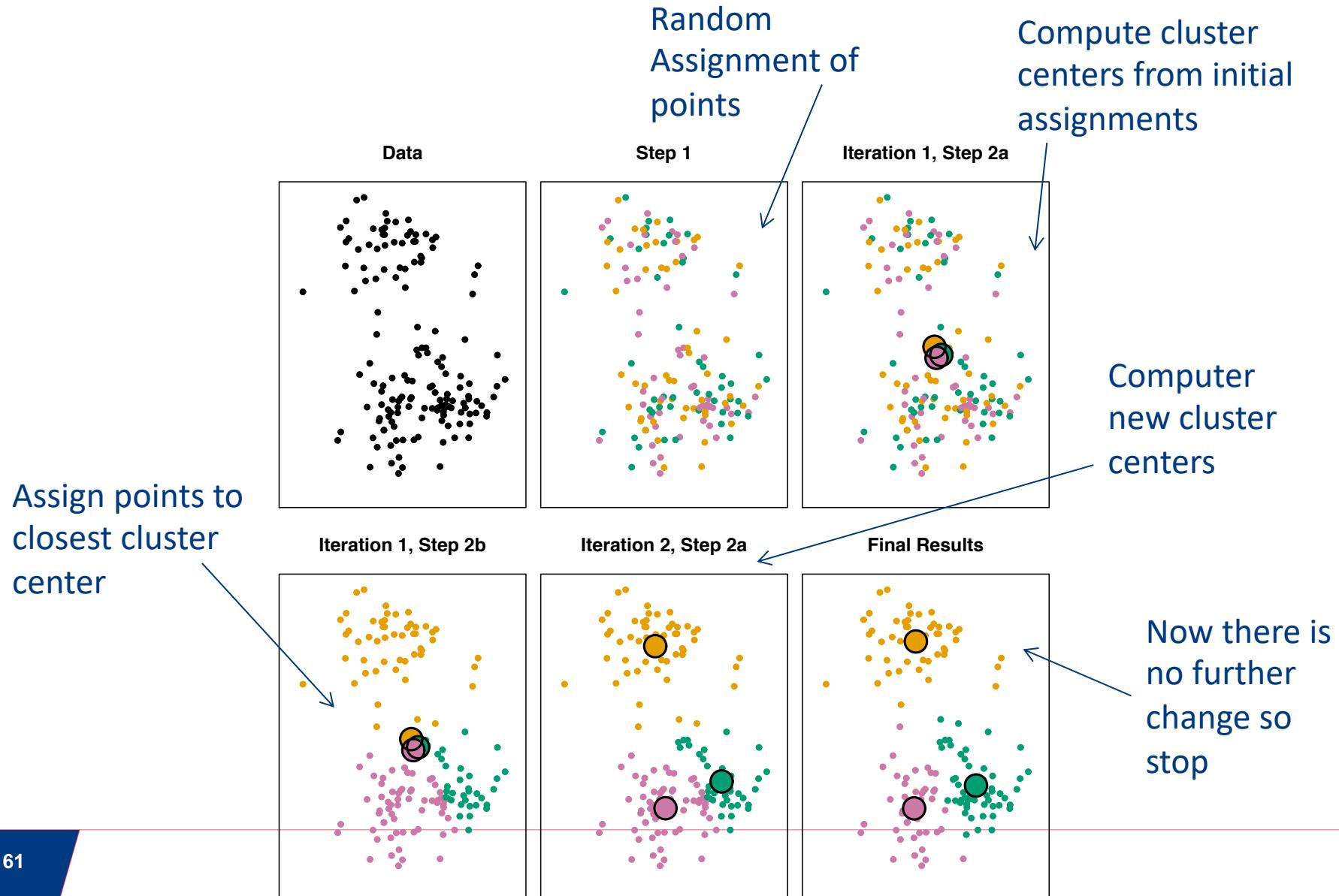
- A simulated data set with 150 observations in 2-dimensional space. Panels show the results of applying K-means clustering with different values of K, the number of clusters.

# K-means clustering: algorithm

Given a pre-defined  $K$

1. Randomly assign a number, from 1 to  $K$ , to each of the observations.
2. Iterate until the cluster assignments stop changing
  - a) For each of the  $K$  clusters, compute the cluster *centroid*. The  $k$ th cluster centroid is the vector of the  $p$  feature means for the observations in the  $k$ th cluster.
  - b) Assign each observation to the cluster with center *closest* to the observation (defined using Euclidean distance)

# Illustration of K-Means with 3 clusters



# K-means clustering: details

- The idea behind K-means clustering:
  - a good clustering is one for which the **within-cluster variation (WCV)** is as small as possible.
- The WCV for cluster k: measure the amount by which the observations within a cluster differ from each other.
- K-means clustering is to solve this optimization problem:

$$\underset{C_1, \dots, C_K}{\text{minimize}} \left\{ \sum_{k=1}^K \text{WCV}(C_k) \right\}.$$

# K-means clustering: details

- Typically use Sum of Squares to measure WCV

$$\text{WCV}(C_k) = \frac{1}{|C_k|} \sum_{i,i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2$$

$|C_k|$  : number of observations in the kth cluster

- The optimization problem of K-means clustering:

$$\underset{C_1, \dots, C_K}{\text{minimize}} \left\{ \sum_{k=1}^K \frac{1}{|C_k|} \sum_{i,i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2 \right\}$$

# K-means clustering: details

- We can rewrite the WCV for cluster k as:

$$\frac{1}{|C_k|} \sum_{i,i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2 = 2 \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - \bar{x}_{kj})^2$$

where  $\bar{x}_{kj} = \frac{1}{|C_k|} \sum_{i \in C_k} x_{ij}$  is the mean for feature j in kth cluster

- The algorithm is guaranteed to decrease the value of the objective at each iteration.
  - Why?

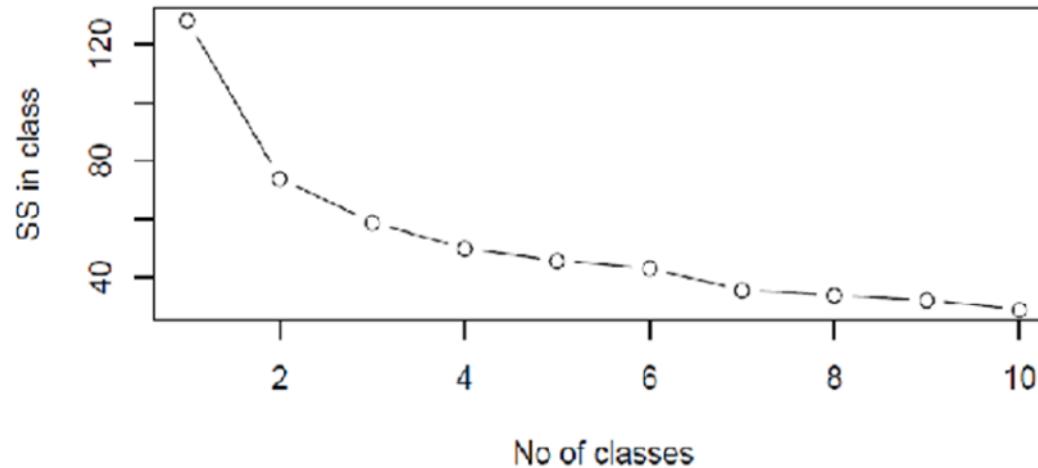
But it is not  
guaranteed to give  
the global minimum

Perform K-means  
clustering algorithm  
6 times,  
with different random  
starting assignments



# Determine the number of clusters

- Based on the sum of squares (SS) within the clusters (WCV) for a solution against the number of clusters
- The “optimal” number of clusters can be determined by the "elbow criterion".



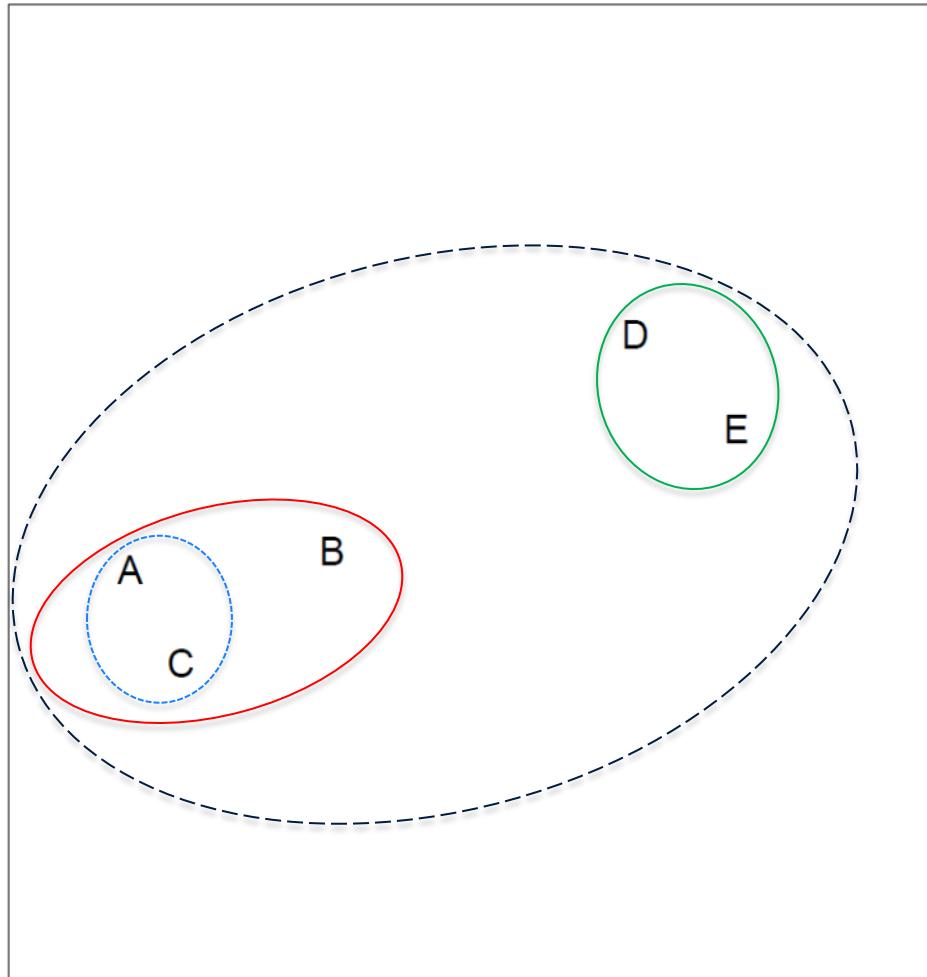
- This "elbow" cannot always be unambiguously identified

# Hierarchical Clustering

- We don't need to commit to a particular choice of K.
- Bottom-up or agglomerative clustering
  - A dendrogram is built starting from the leaves and combining clusters up to the trunk.

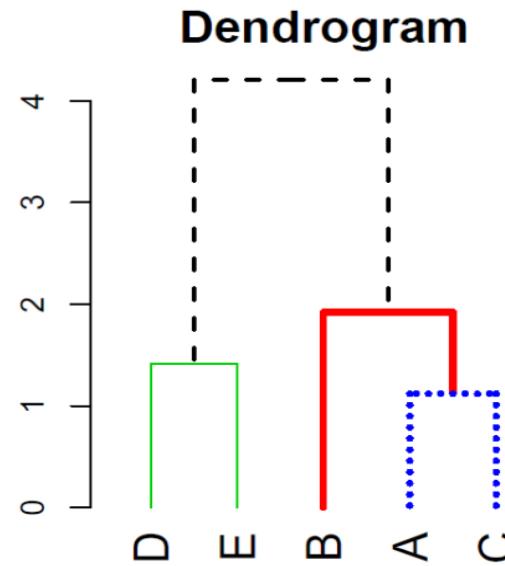
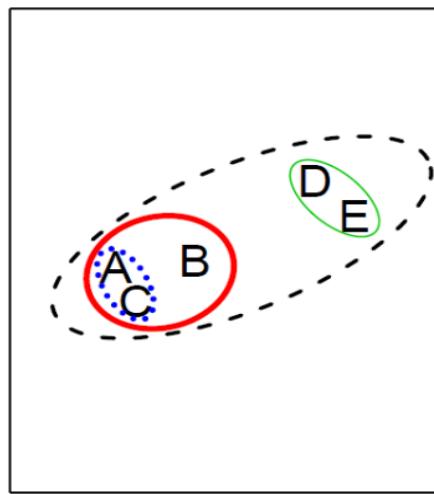
# Hierarchical Clustering: the idea

Builds a hierarchy in a “bottom-up” fashion...



# Hierarchical Clustering Algorithm

- Begin with  $n$  observations. Treat each observation as its own cluster.
- Identify the **closest** two clusters and merge them.
- Repeat.
- Ends when all points are in a single cluster.



# Hierarchical Clustering Algorithm

- Main modeling activities
  - Definition of the distance between the clusters
  - Determination of number of clusters as final clustering result

# Definition of the distance between the clusters

- The distance between the clusters is called *linkage*
- Different specifications for linkage (more details in the book!)
  - Single linkage: Distance between clusters is the distance of the closest points
  - Complete linkage: Distance between clusters is the distance of the farthest points
  - Average linkage: mean distance between all the points in the two clusters
  - Ward distance: difference between the total within cluster sum of squares for the two clusters separately, and the within cluster sum of squares resulting from merging the two clusters in one cluster

# Other clustering algorithms

Other frequently used general methods:

- Model based clustering which estimated a mixture distribution for the data using the EM-algorithm
- DBSCAN which groups points according to their density measured by the number of nearest neighbors
- Self organizing maps (SOM) which are a special type of neural nets and can be interpreted as k-means clustering defined on a distorted grid
- Two stage clustering combine the ideas of hierarchical clustering with the ideas of k-means by using a cluster-feature tree

# Applying clustering algorithms: special attention (1)

- Standardizing the variables matters!
  - Calculation of distances between observations characterized by variables with different scales needs special attention

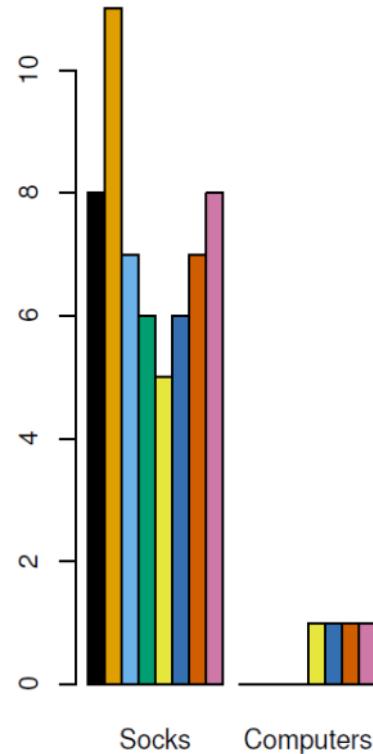
## Standardizing the variables matters: example

- *An online retailer sells two items: socks and computers.*  
She is interested in clustering shoppers based on their past shopping histories. She wants to identify subgroups of *similar* shoppers, so that shoppers within each subgroup can be shown items and advertisements that are particularly likely to interest them.

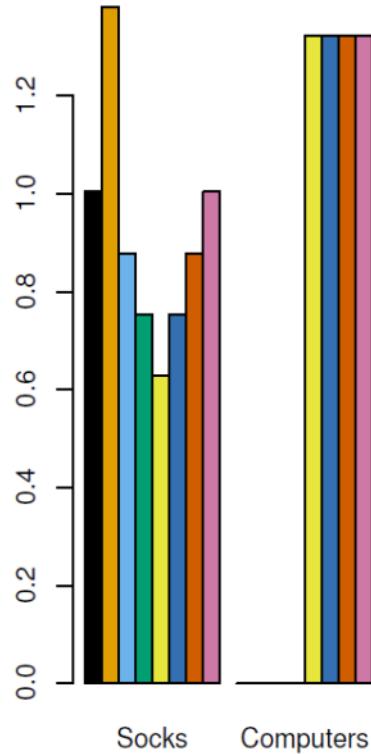
Suppose the data takes the form of a matrix, where the rows are the different shoppers and the columns are the *items* available for purchase; the elements indicate the *number of times* a given shopper has purchased a given item.

Some items may be purchased more frequently than others; for instance, a shopper might buy ten pairs of socks a year, but a computer very rarely.

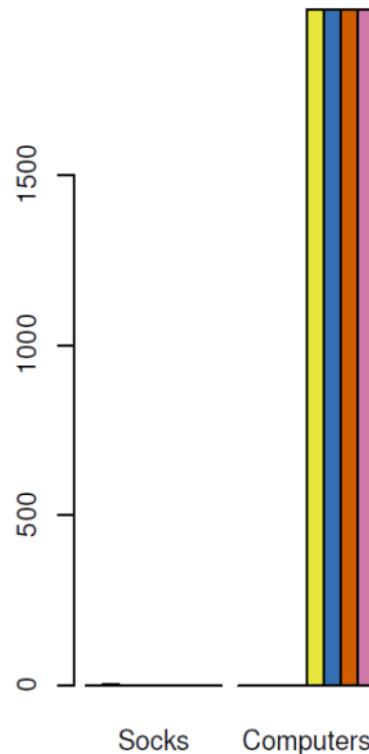
## Standardizing the variables matters: example



*Number of socks and computers purchased by eight shoppers. Socks have higher weight*



*Same data, after standardizing, socks and computers have equal weight*



*Same data, but y-axis represents the number of Euros. Computers have higher weight*

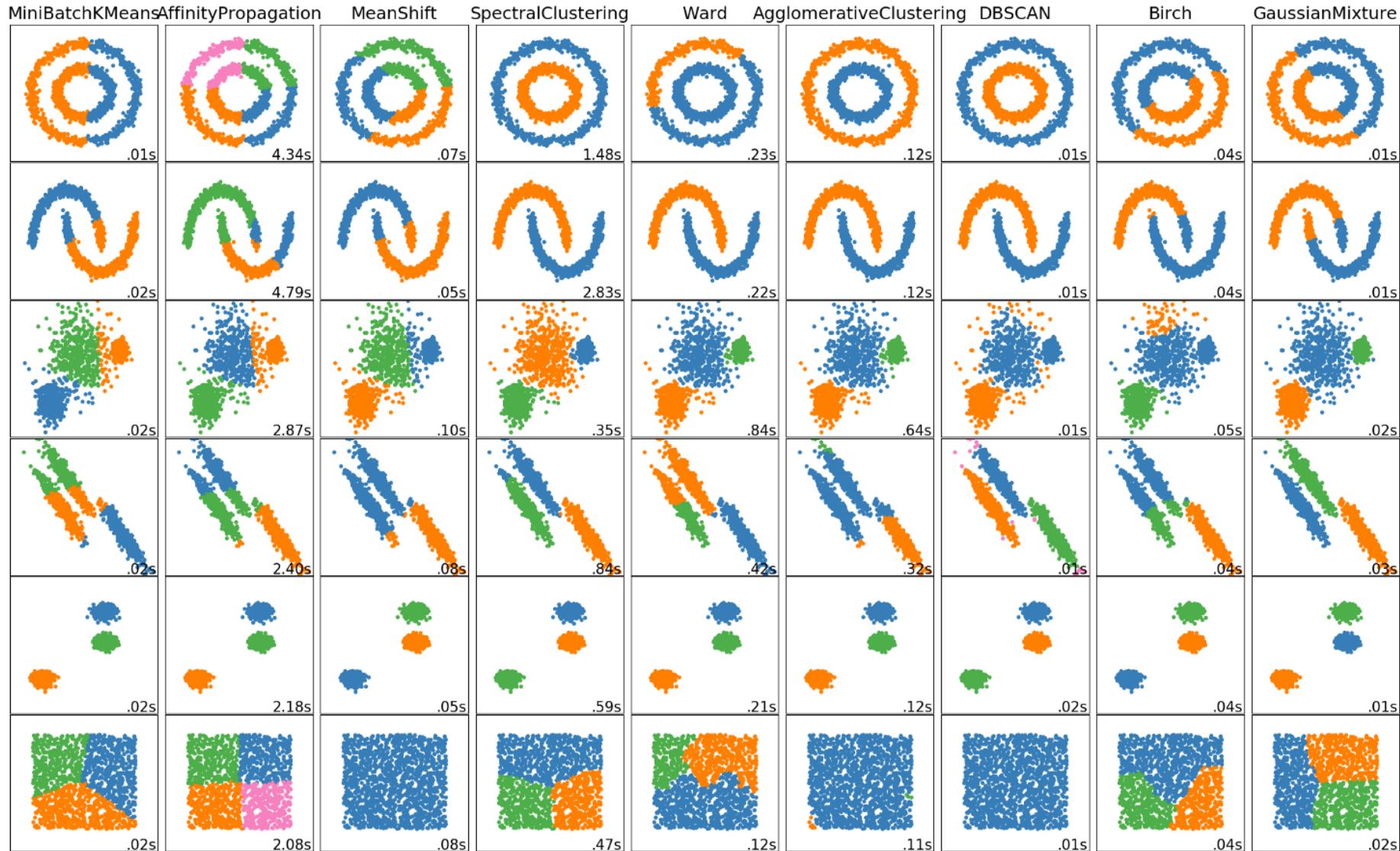
# Clustering: special attention 2: evaluation

## Evaluating clustering results:

### 1. Homogeneity measures for groups: internal evaluation

- Based on the data itself: Silhouette coefficient; Davies–Bouldin index; Dunn index (see book)
- Usually assign the highest score to the algorithm, which outputs clusters with high similarity within a cluster and low similarity between clusters
- **Potential problem:** evaluation is biased towards algorithms using the same cluster model!

## Special attention (2) -> evaluating clustering results



# Clustering: special attention (2)

Evaluating clustering results:

2. Validation with a test sample: idea from supervised learning
3. Validation wrt a subject matter explanation

# Special attention (2) -> evaluating clustering results

Evaluating clustering results:

2. Validation wrt a subject matter explanation
  - validating results with domain experts.

# Special attention (2) -> evaluating clustering results

- Example: online auction bidding behavior analysis
  - Each row represents a bidder; each with three *features*

Number of Bids	Time of First Bid	Time of Last Bid	Cluster
3	10	20	0
6	200	235	1
20	20	200	2
40	5	220	2
2	24	50	0

# Special attention (2) -> evaluating clustering results

- Interpreting results using **expert/domain knowledge** of bidder behavior

Number of Bids	Time of First Bid	Time of Last Bid	Cluster	<i>Definition from auction experts and literature</i>
3	10	20	0	Early evaluators
6	200	235	1	Opportunists
20	20	200	2	Participators
40	5	220	2	Participators
2	24	50	0	Early evaluators

# Special attention (2) -> evaluating clustering results

- Interpreting results using knowledge of bidder behavior

Number of Bids	Time of First	Time of Last	Cluster	<i>Definition from auction experts and literature</i>
3				Early evaluators
6				Opportunists
20				Participators
40	-	--	-	Participators
2	24	50	0	Early evaluators



# Summary of clustering: practical issues

- Should features first be standardized in some way?
- Most cluster methods use the definition of a distance between observations; which one to use?
- How many clusters to choose? (in both K-means or hierarchical clustering);
  - Difficult problem. No agreed-upon method.
  - The number of clusters can be determined by descriptive analysis of the cluster solution
  - For validation of a cluster solution interpretation from a business perspective is important
- How to choose different algorithms? Do evaluations properly!

# Further reading

- Han, Jiawei, Micheline Kamber, and Jian Pei, Data mining: concepts and techniques, Chapter 9, 10
- [https://www.sciencedirect.com/book/9780123814791/data-mining-  
concepts- and-techniques](https://www.sciencedirect.com/book/9780123814791/data-mining-concepts-and-techniques)  
(free via TU/e VPN)