

Learning from Data

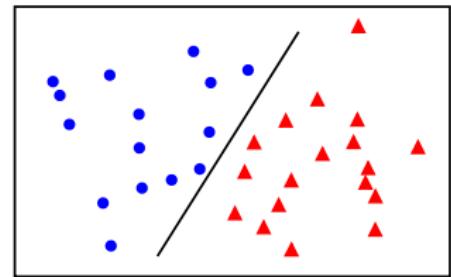
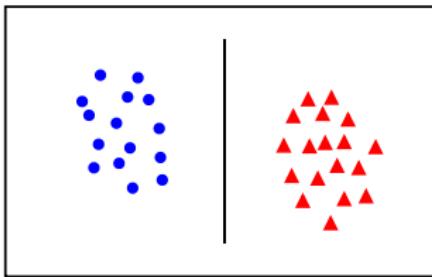
Lecture 5: Support Vector Machines, Regression, Unsupervised learning (Clustering)

Malvina Nissim
m.nissim@rug.nl
room 1311.421

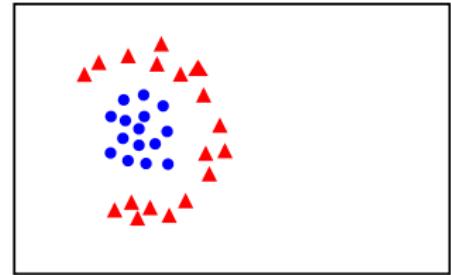
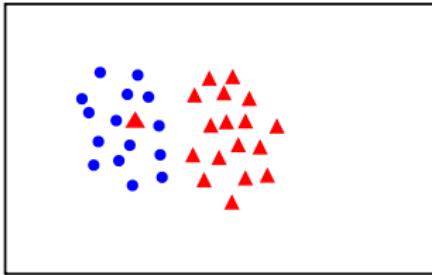
1 October 2018

Linear separability

linearly
separable



not
linearly
separable



Learning process

- the goal of the learning process is to come up with a “good” weight vector w
- the learning process will use examples to guide the search of a “good” w
- different notions of “goodness” exist, which yield different learning algorithms

Support Vector Machines

(also based on slides by Glenn Fung, O.L.Mangasarian, Rob Lothian, Kristin Bennett)

Support Vector Machines

vector space based machine-learning method aiming to find a decision boundary between two classes that is maximally far from any point in the training data (possibly discounting some points as outliers or noise)

Support Vector Machines

vector space based machine-learning method aiming to find a decision boundary between two classes that is maximally far from any point in the training data (possibly discounting some points as outliers or noise)

Support Vector Machines

vector space based machine-learning method aiming to find a **decision boundary between two classes** that is maximally far from any point in the training data (possibly discounting some points as outliers or noise)

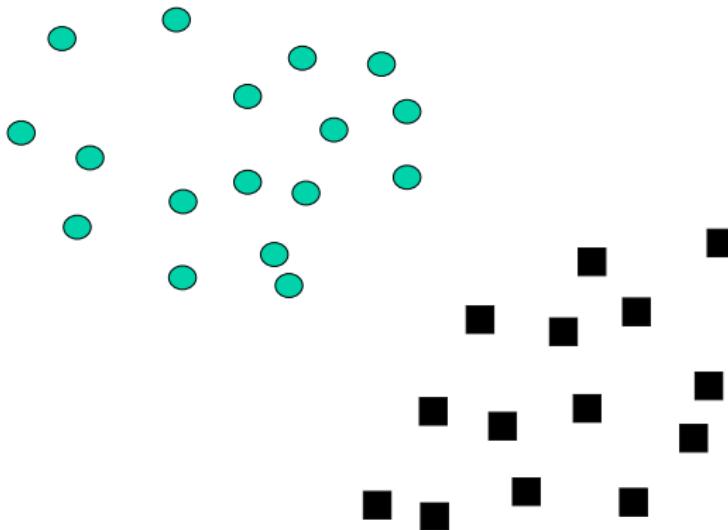
Support Vector Machines

vector space based machine-learning method aiming to find a decision boundary between two classes that is **maximally far from any point in the training data** (possibly discounting some points as outliers or noise)

Support Vector Machines

vector space based machine-learning method aiming to find a decision boundary between two classes that is maximally far from any point in the training data (**possibly discounting some points as outliers or noise**)

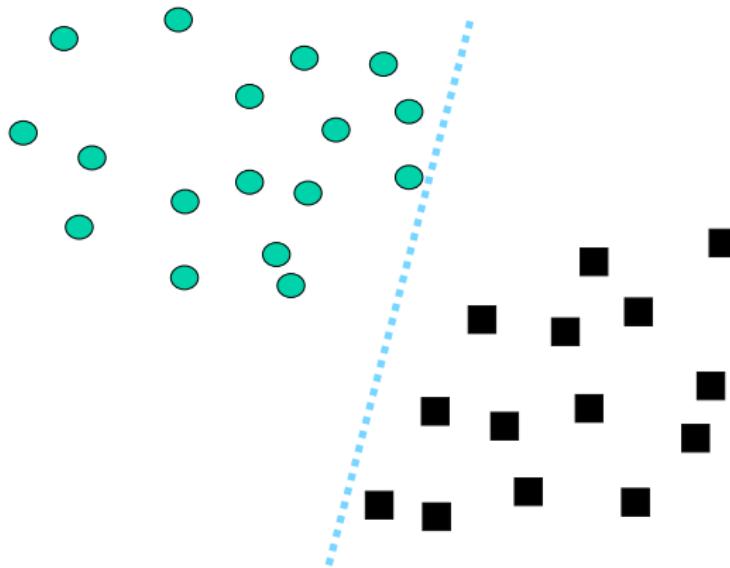
Which of the linear separators is optimal?



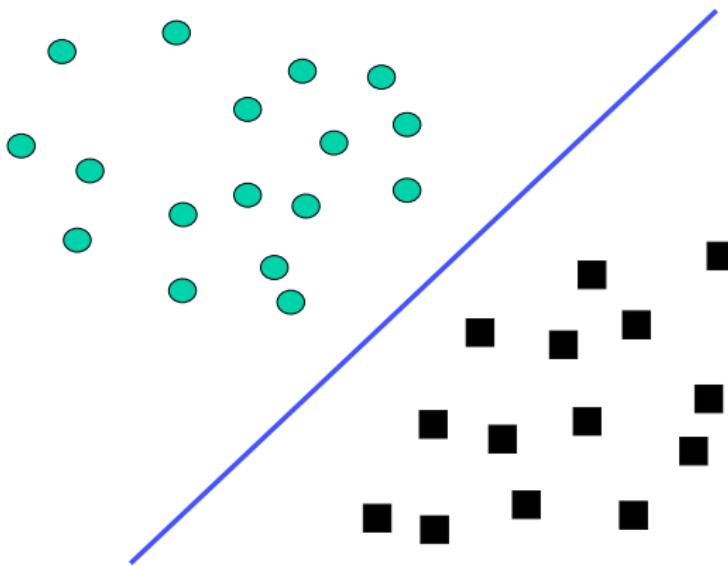
One intuition:

the fewer points near the decision surface, the fewer very uncertain classification decisions

Best Linear Separator?

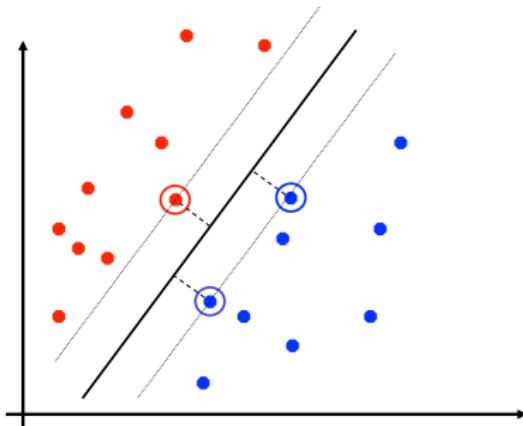


Best Linear Separator?



Maximum Margin Classification

- Maximizing the margin is good according to intuition and theory.
- Implies that only support vectors are important; other training examples are ignorable.



About the Name...

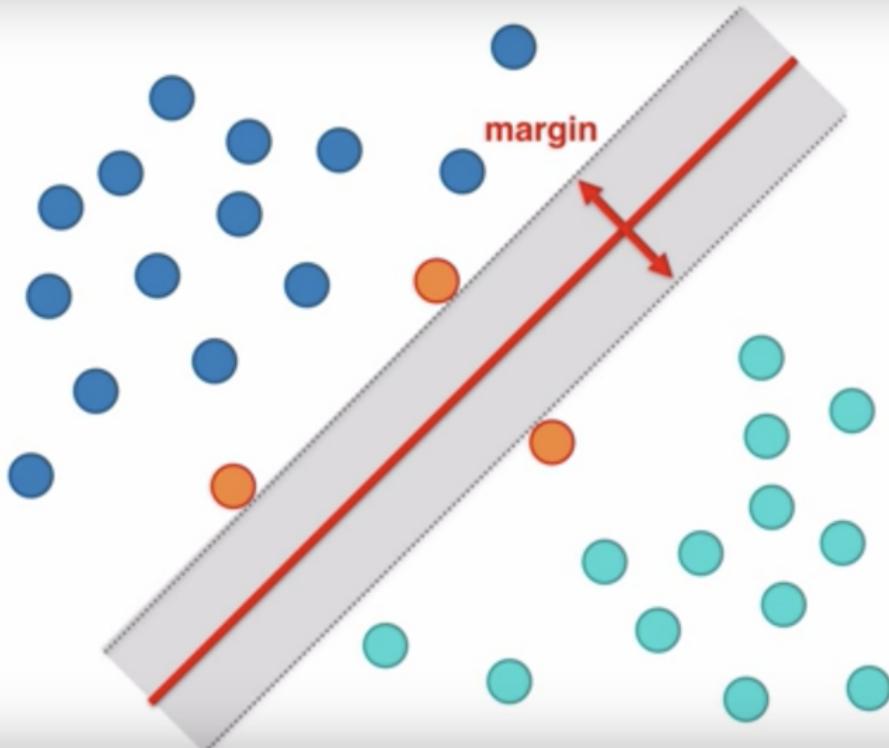
A Support Vector

A training sample used to define classification boundaries in SVMs

- located near class boundaries

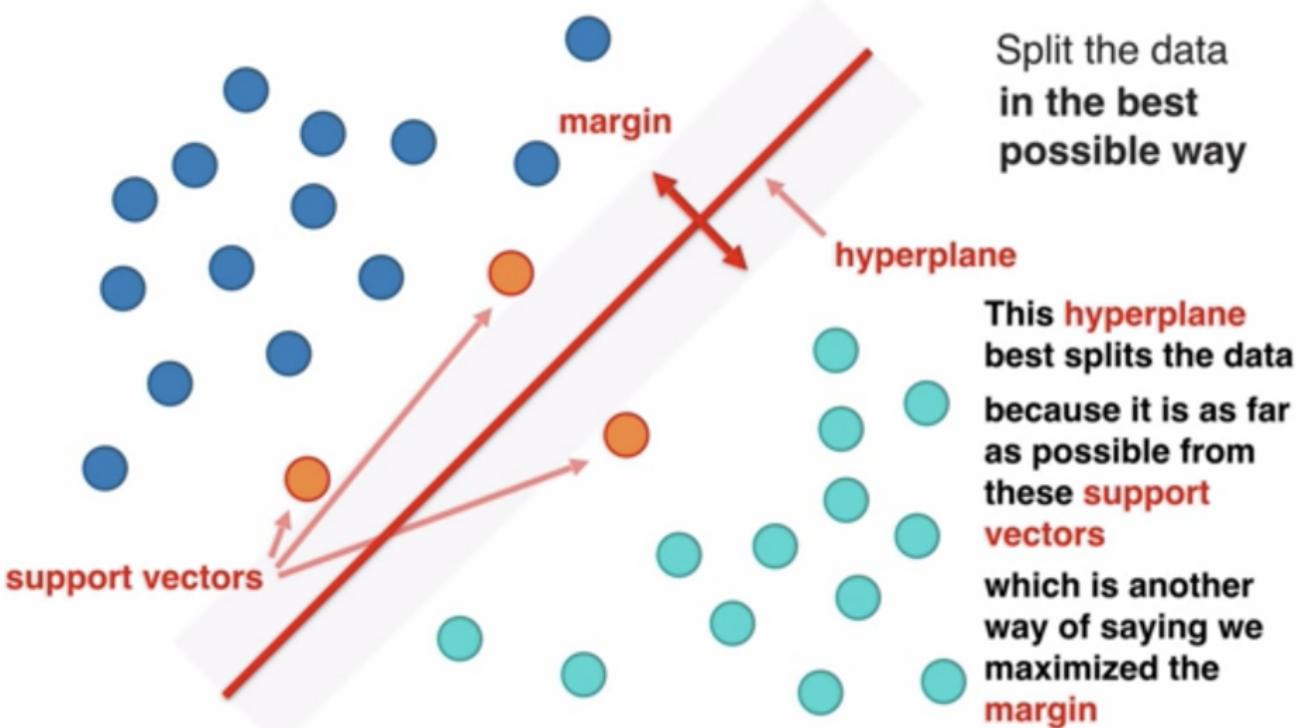
Support Vector Machines

Binary classifiers whose decision boundaries are defined by support vectors



Why is this the best split?

The distance between the support vectors and the hyperplane are as far as possible



More on ‘The Margin’

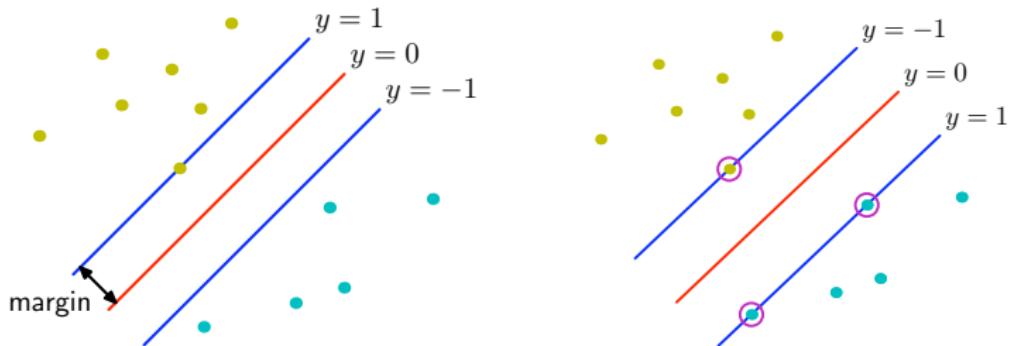
The Margin, b

Is the minimum distance of *any* sample to the decision boundary.

Training SVMs

= maximizing the margin, moving the decision boundary as far away from all training samples as possible.

Maximizing the Margin



- At left: a sub-optimal margin
- At right: optimal margin
- **y values:** for linear function defined by the SVM. For linearly separable data, all training instances correctly classified as -1 or 1 (locations in the margins have y values in (-1,1))

Support Vector Machines

- SVMs pick **best** separating hyperplane according to some criterion
 - e.g. maximum margin
- Training process is an **optimisation**
- Training set is effectively reduced to a relatively small number of **support vectors**

A bit deeper into SVMs

- C parameter:
- multiple classes:
- kernel trick:

A bit deeper into SVMs

- C parameter:
how strict should I be when drawing my hyperplane?
- multiple classes:
SVM is intrinsically a binary classifier, so what do I do with more classes?
- kernel trick:
what if my data is really *not* linearly separable?

Soft margin

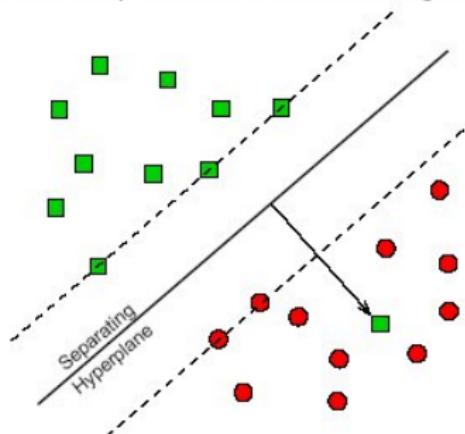
what if there isn't a margin that can split the data perfectly?
And how perfect should I be?

Soft margin

what if there isn't a margin that can split the data perfectly?
And how perfect should I be?

Non-separable training sets

Use linear separation, but admit training errors.



Penalty of error: distance to hyperplane multiplied by *error cost C*.

Soft margin

what if there isn't a margin that can split the data perfectly?
And how perfect should I be?

- introduce a parameter to allow for (a certain amount of) mislabelled examples
- the C parameter tells the SVM optimization how much you want to avoid misclassifying each training example.
The higher C, the largest the penalty to errors.
 - large values of C:
 - a very small value of C:
 - very very tiny values of C:

Soft margin

- the C parameter tells the SVM optimization how much you want to avoid misclassifying each training example.
The higher C, the largest the penalty to errors.
 - large values of C: large penalty to errors; the optimization will choose a smaller-margin hyperplane if that hyperplane does a better job of getting all the training points classified correctly
 - a very small value of C: the optimizer will look for a larger-margin separating hyperplane, even if that hyperplane misclassifies more points.
 - very very tiny values of C: you are likely to get misclassified examples, often even if your training data is linearly separable.

Soft margin

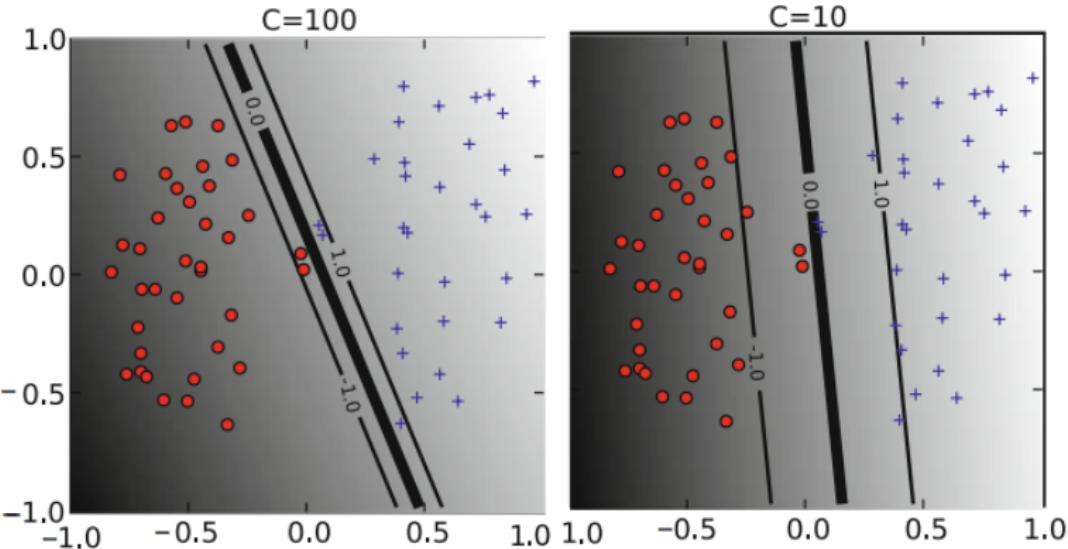


Fig. 13.3. The effect of the soft-margin constant, C , on the decision boundary. A smaller value of C (right) allows to ignore points close to the boundary and increases the margin. The decision boundary between negative examples (*circles*) and positive examples (*crosses*) is shown as a thick line. The lighter lines are on the margin (discriminant value equal to -1 or $+1$). The grayscale level represents the value of the discriminant function, dark for low values and a light shade for high values.

SVM is a binary classifier!

SVM is a binary classifier!

- OVR
- OVO

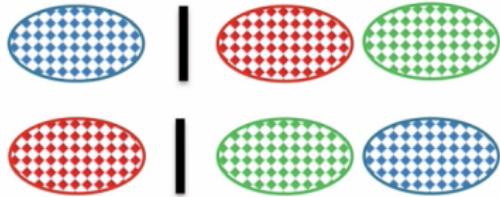
SVM is a binary classifier!

- OVR: one versus rest
- OVO: one versus one

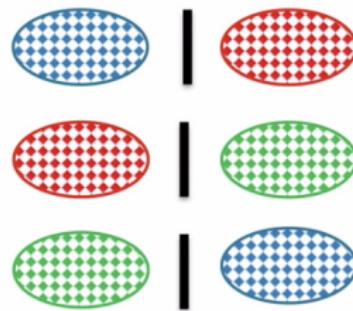
SVM is a binary classifier!

- OVR: one versus rest
- OVO: one versus one

OVR: One vs Rest

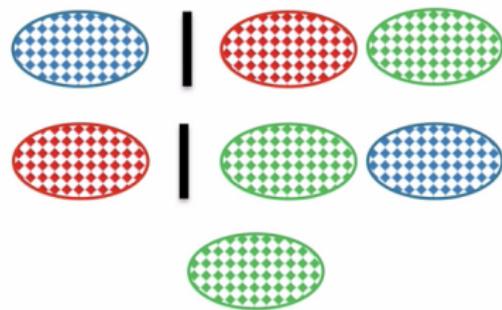


OVO: One vs One



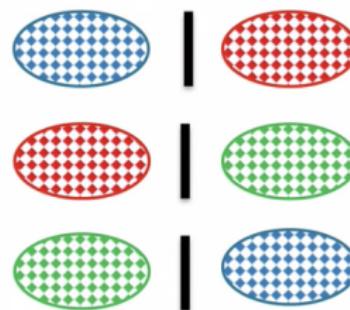
SVM is a binary classifier!

OVR: One vs Rest



Pros: Fewer classifications
Cons: Classes may be imbalanced

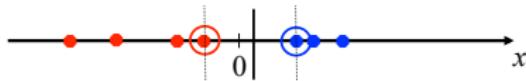
OVO: One vs One



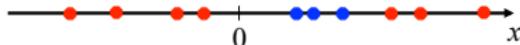
Pros: Less sensitive to imbalance
Cons: More classifications

Non-linear SVMs

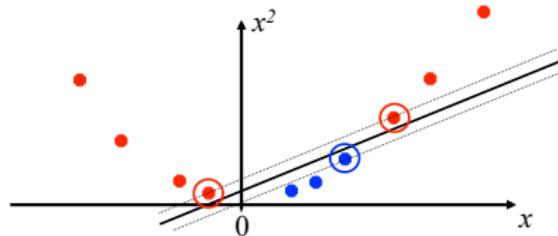
- Datasets that are linearly separable with some noise work out great:



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

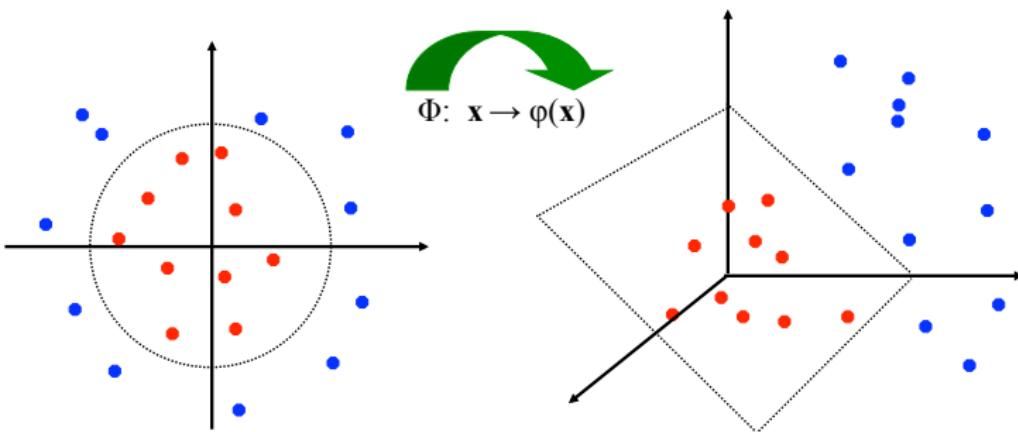


- How about... mapping data to a higher-dimensional space:



Non-linear SVMs: Feature spaces

- General idea: the original feature space can always be mapped to some higher-dimensional feature space where the training set is separable:



Kernels

- We may use Kernel functions to implicitly map to a new feature space
- Kernel fn:

$$K(\mathbf{x}_1, \mathbf{x}_2) \in \mathbf{R}$$

- Kernel must be equivalent to an inner product in some feature space

Example Kernels

Linear: $\langle \mathbf{x} \cdot \mathbf{z} \rangle$

Polynomial: $P(\langle \mathbf{x} \cdot \mathbf{z} \rangle)$

Gaussian: $\exp(-\|\mathbf{x} - \mathbf{z}\|^2 / \sigma^2)$

SVM in Scikit-learn

```
1 from sklearn import svm
2
3 C=1.0 # remember what this is?
4 cls = svm.SVC(kernel='linear', C=C)
5 classifier = Pipeline( [ ('vec', vec),
6                         ('cls', cls) ] )
7 classifier.fit(X, y)
```

Check: <http://scikit-learn.org/stable/modules/svm.html#tips-on-practical-use>

supervised classification algorithms

supervised classification algorithms

supervised **classification** algorithms

Classification vs Regression

create models of prediction from gathered data

- classification

the dependent variables are categorical

- input x : feature vector
- output: **discrete class label**

- regression

the dependent variables are numerical

- input x : feature vector
- output y : **continuous value**

regression

Classification vs Regression

create models of prediction from gathered data

- classification

the dependent variables are categorical

- input x : feature vector
- output: **discrete class label**

- regression

the dependent variables are numerical

- input x : feature vector
- output y : **continuous value**

Example: Predicting housing prices in the Netherlands

Data

training
data

size m ²	price
57	150,000
90	210,000
30	90,000

Data

training
data

size m ²	price
57	150,000
90	210,000
30	90,000

x (input) y (output)

Data

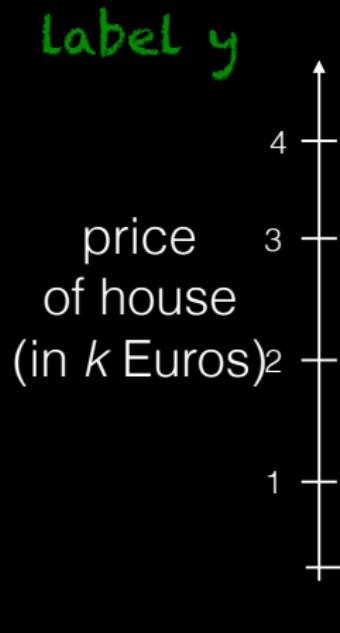
training
data

size m ²	price
57	150,000
90	210,000
30	90,000

x (input) y (output)

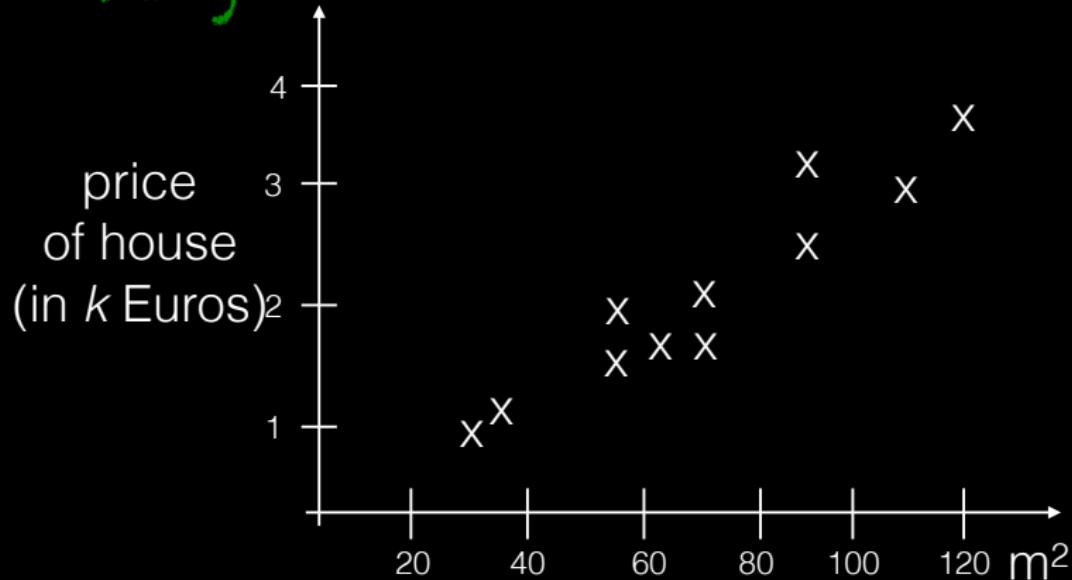
$\langle x, y \rangle$

Regression



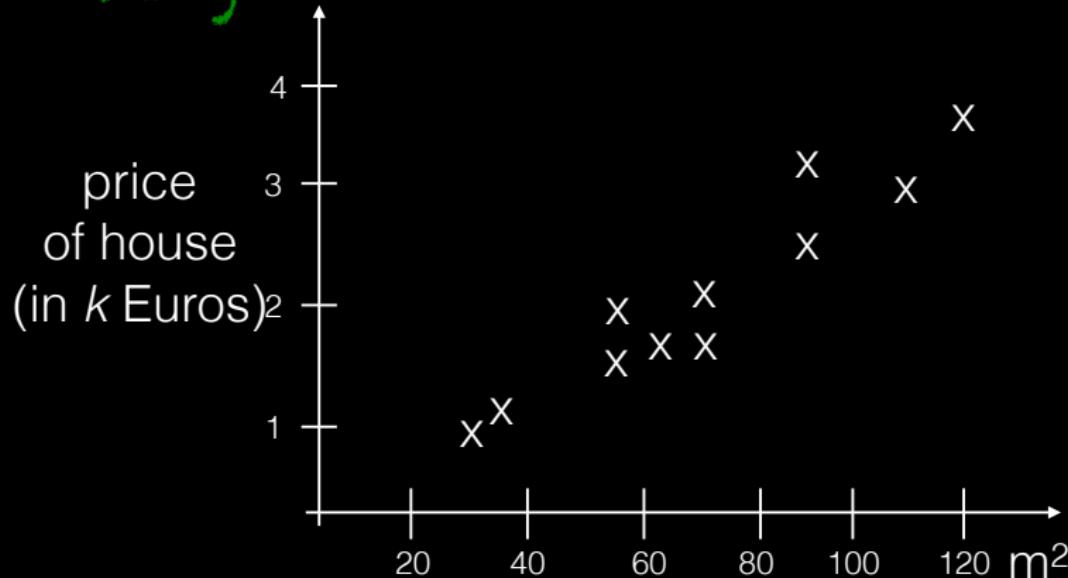
Regression

Label y



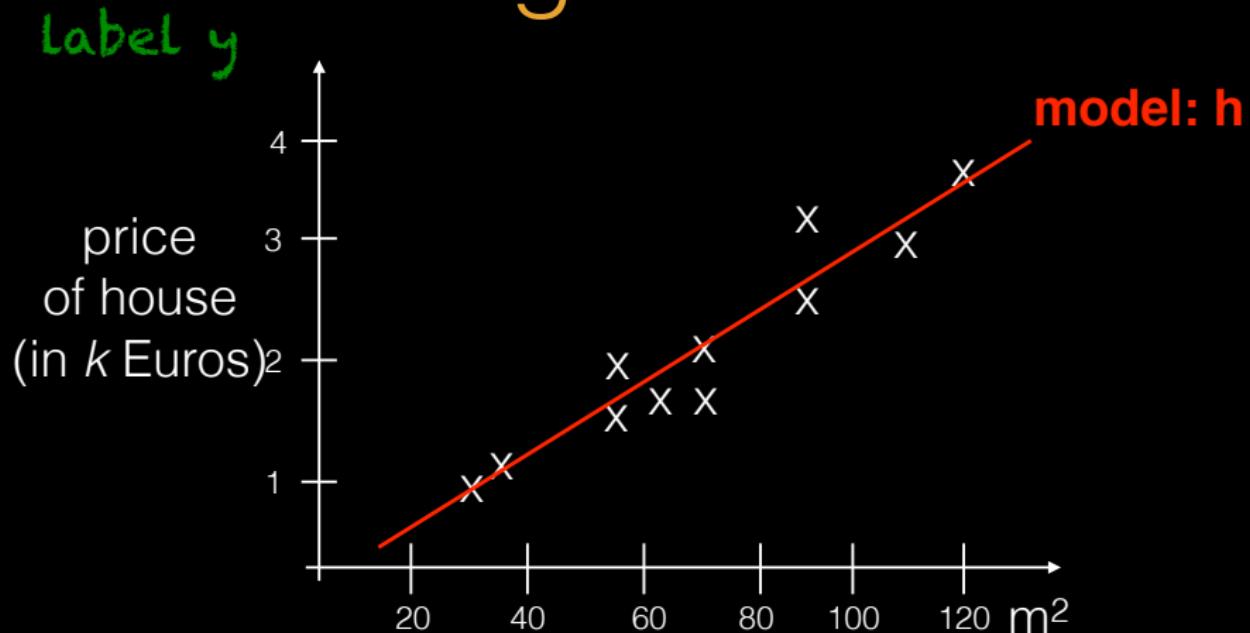
Regression

Label y



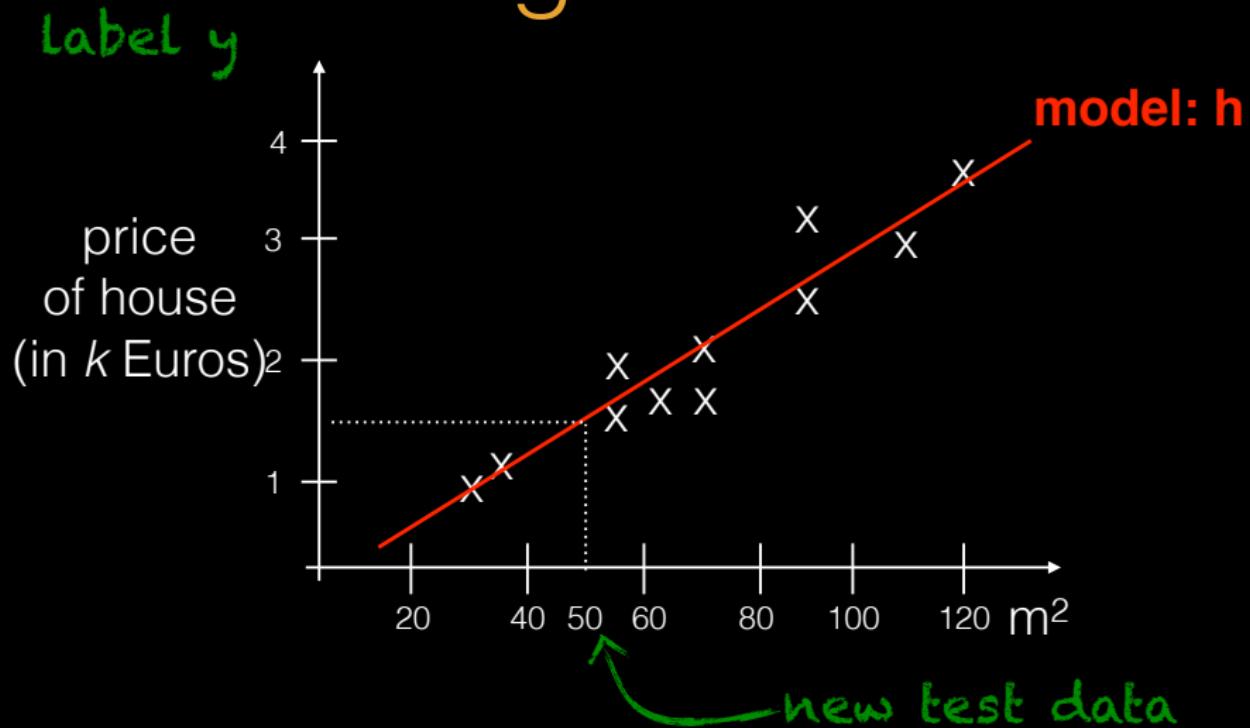
regression: predict numeric (continuous) output

Regression



regression: predict numeric (continuous) output

Regression



regression: predict numeric (continuous) output

Generalisation

experience
(training data)

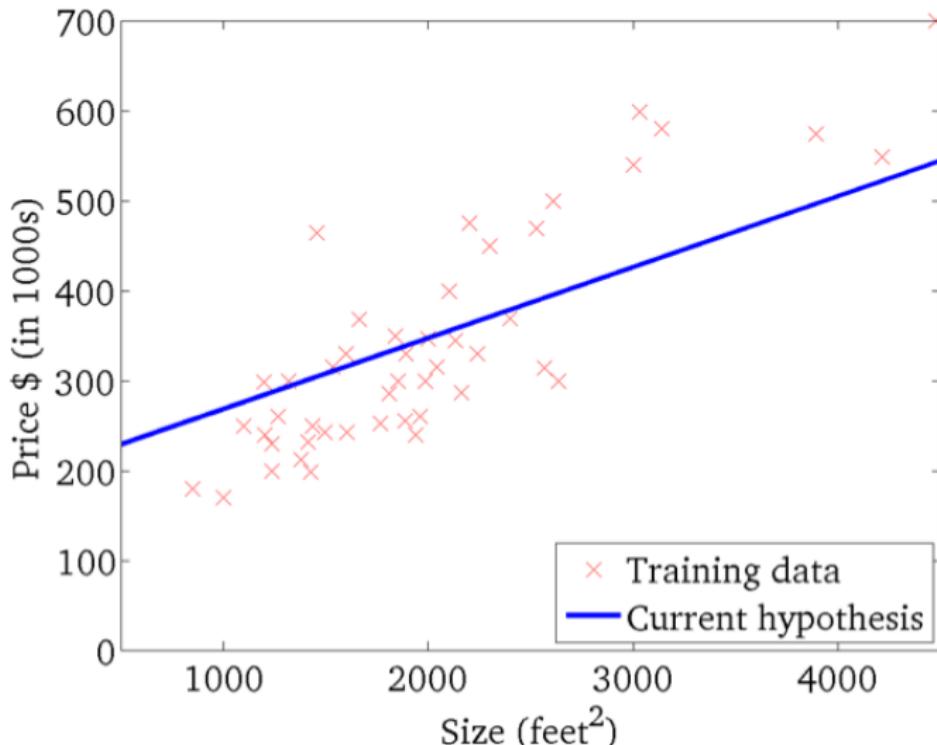


size m ²	price
57	150,000
90	210,000
30	90,000

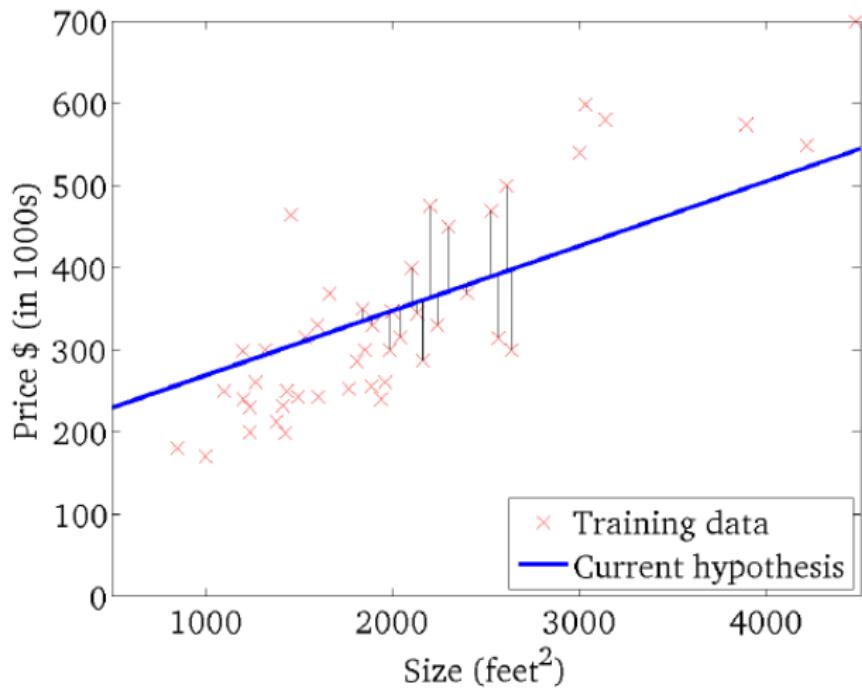
abstract
representation



Same one-feature example



Error as vertical lines



Mean Squared Error (MSE)

How close is a regression line to a set of points?

General idea:

- take the distances from the actual points to the regression line (distance = error)
- square them (necessary to remove any negative signs; also gives more weight to larger differences)
- take average

Mean Squared Error (MSE)

Steps to calculate the MSE from a set of X and Y values:

- find the regression line
- insert your X values into the linear regression equation to find the new Y values (Y)
- subtract the new Y value from the original to get the error
- square the errors
- add up the errors
- find the mean

interpretation: the smaller the MSE, the closer to the best fit

Regression in scikit

- training (fitting)
- testing
- evaluating

```
1 >>>from sklearn import linear_model
2 >>>from sklearn.metrics import mean_squared_error
3
4 >>>lr = linear_model.LinearRegression()
5 >>>lr.fit(Xtrain, Ytrain)
6
7 >>>Yguess = lr.predict(Xtest)
8 >>>mean_squared_error(Ytest, Yguess)
9
10 # worked out:
11 >>>np.mean((Yguess - Ytest) ** 2))
```

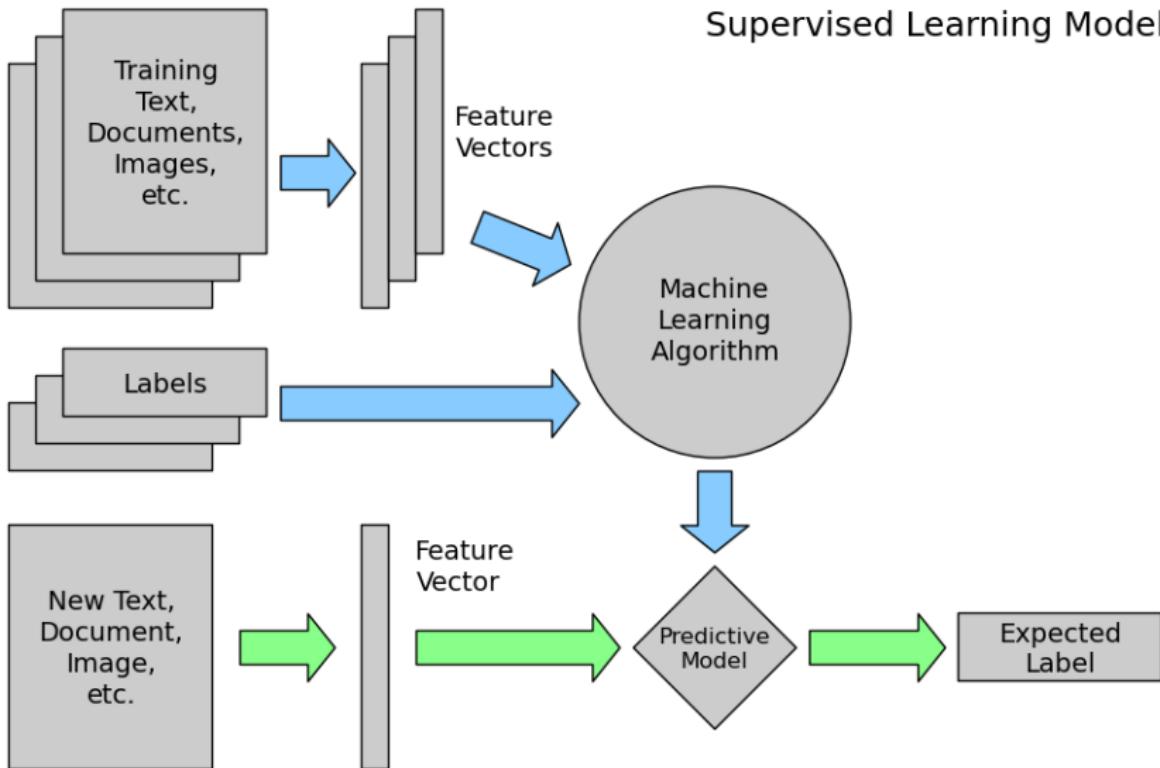
classification and regression are the most standard ways of doing
supervised learning

Supervised and Unsupervised Learning

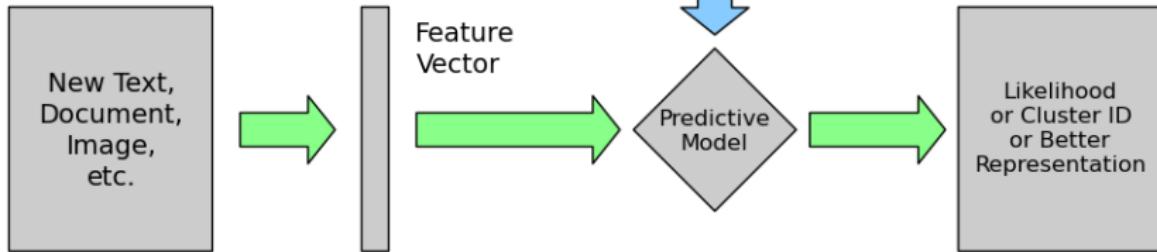
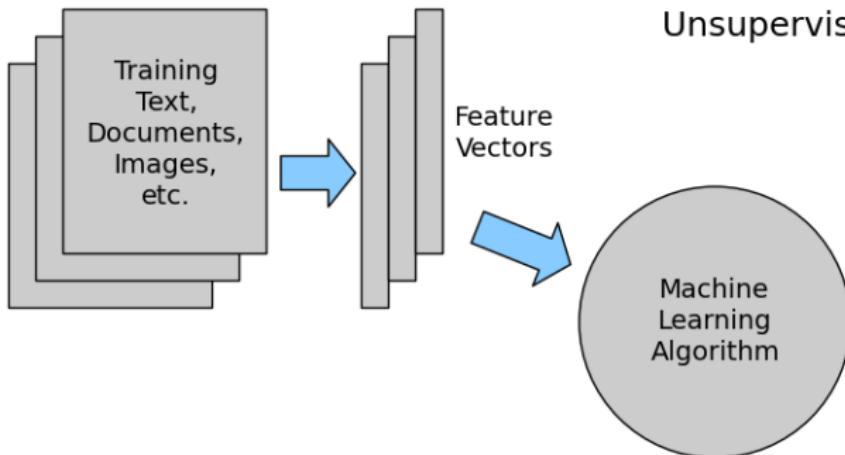
information about the **correct** distribution/label of the training examples

- in supervised learning it is **known**
 - fitting a model to labelled data which has the correct answer
- in unsupervised learning it is **not known**
 - finding structure in unlabelled data

Supervised Learning Model



Unsupervised Learning Model



Supervised learning

supervised learning – classification or regression

- in training, instances are associated with their class label
- based on features, the system must search for patterns and build a model
- the model must be able to *predict* the class of previously unseen instances

Unsupervised learning

unsupervised learning – clustering

- partitioning instances into subsets (clusters) that share similar characteristics
- subsets are not predefined
- a system can be told *how many* clusters it should form (K-means algorithm)

Classification vs. Clustering

- Classification: supervised learning
- Clustering: unsupervised learning
- Classification: Classes are **human-defined** and part of the input to the learning algorithm.
- Clustering: Clusters are **inferred from the data** without human input.
 - However, there are many ways of influencing the outcome of clustering: number of clusters, similarity measure, representation of documents, . . .

Soft vs Hard Clustering

- **hard**: partition the objects so that each object is in **exactly one** partition
- **soft**: assign a **degree** to which each object is in one cluster. It's a bit like fractional membership in several clusters: a document's assignment is a distribution over all clusters

Soft vs Hard Clustering

- **hard**: partition the objects so that each object is in **exactly one** partition
- **soft**: assign a **degree** to which each object is in one cluster. It's a bit like fractional membership in several clusters: a document's assignment is a distribution over all clusters (cf. probability scores)

Flat vs Hierarchical Clustering

- **flat**: a set of clusters without any explicit structure that would relate clusters to each other.
- **hierarchical**: clusters within clusters (a hierarchy of clusters)

Clustering: overview

take a collection of objects, like our documents

- assume that each object can be represented as a point in a geometric space (VSM!)
- in text documents, typically the points are constructed by using word frequencies (bag of words)

goal: create clusters that are *coherent internally, but clearly different from each other*

- documents within a cluster should be as similar as possible
- documents in one cluster should be as dissimilar as possible from documents in other clusters

K-means

K-means

- flat clustering
- number of clusters picked ahead of time
- iterative improvement
- notion of **centroid**
- typically uses Euclidean distance

K-means

- flat clustering
- number of clusters picked ahead of time
- iterative improvement
- notion of **centroid**
- typically uses Euclidean distance

objective: to minimize the average squared Euclidean distance of documents from their cluster centers, where a cluster center is defined as the mean or *centroid* of the documents in a cluster ω

K-means

Assume that there are k classes

- for every class, create a centroid: a point that is in the center of the class
- find centroids so that all the points in each class are as near as possible

K-means — iteration

- choose K (number of clusters)
- initialize K centroids

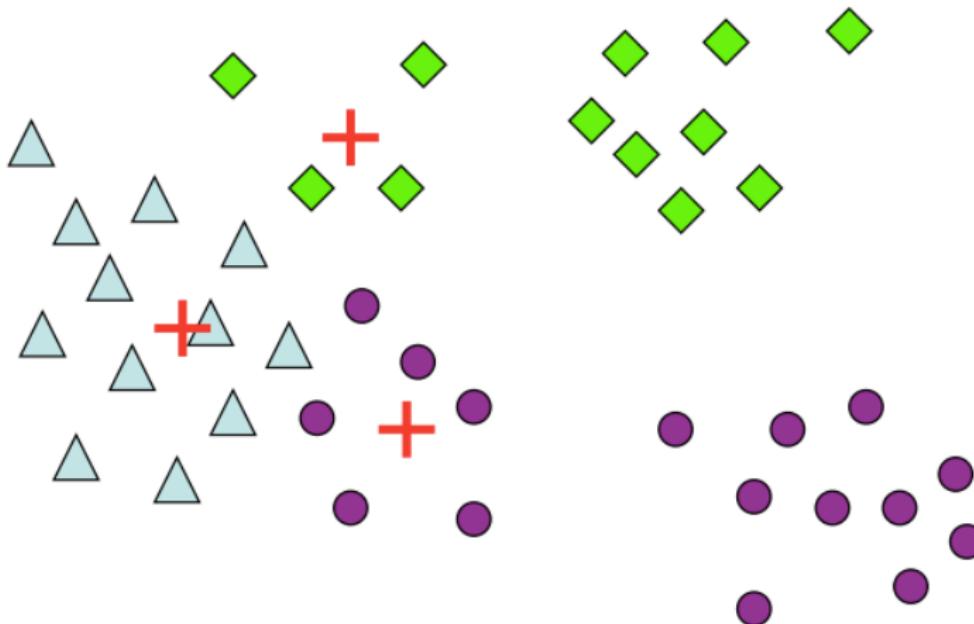
Then, repeatedly perform the following two steps:

- ➊ cluster assignment:
 - points are assigned the class of the closest centroid
- ➋ move centroid
 - Move centroid to the average of the points of the same cluster

After a number of iterations, centroid/assignments don't change anymore, then stop.

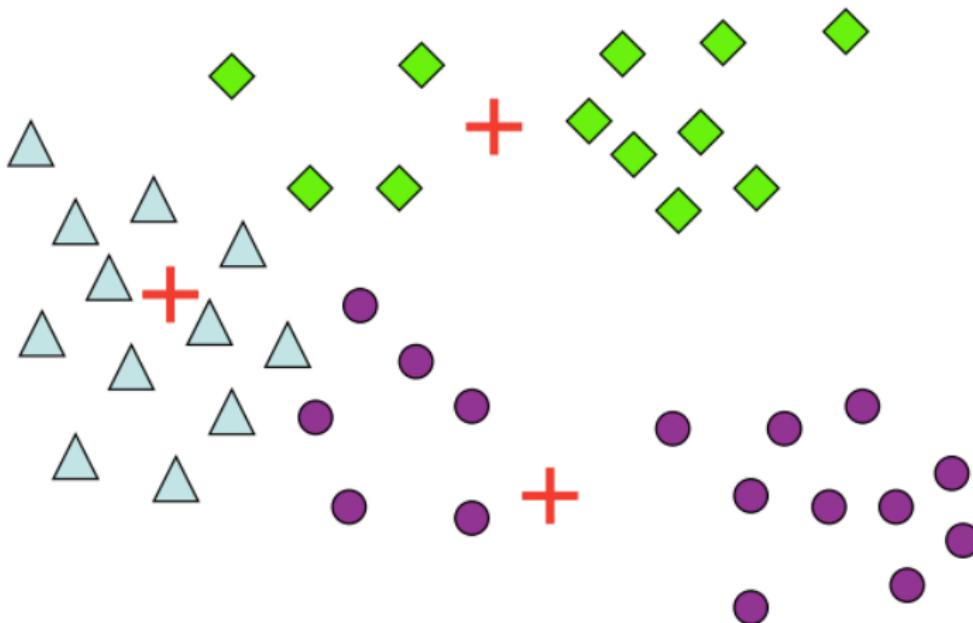
An Example

start: choose centroids and cluster



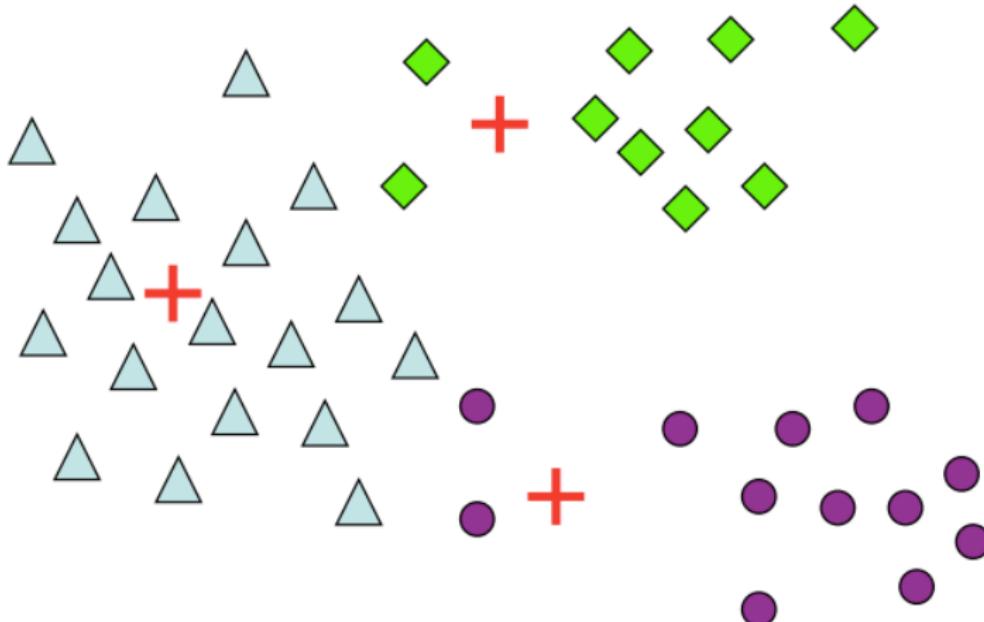
An Example

recompute centroids



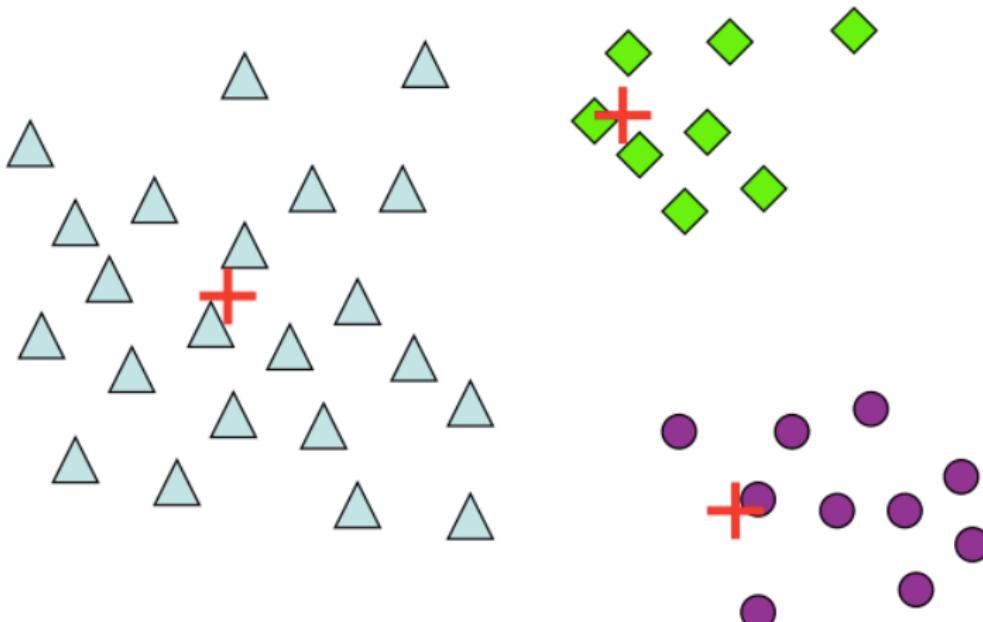
An Example

re-cluster around new centroids



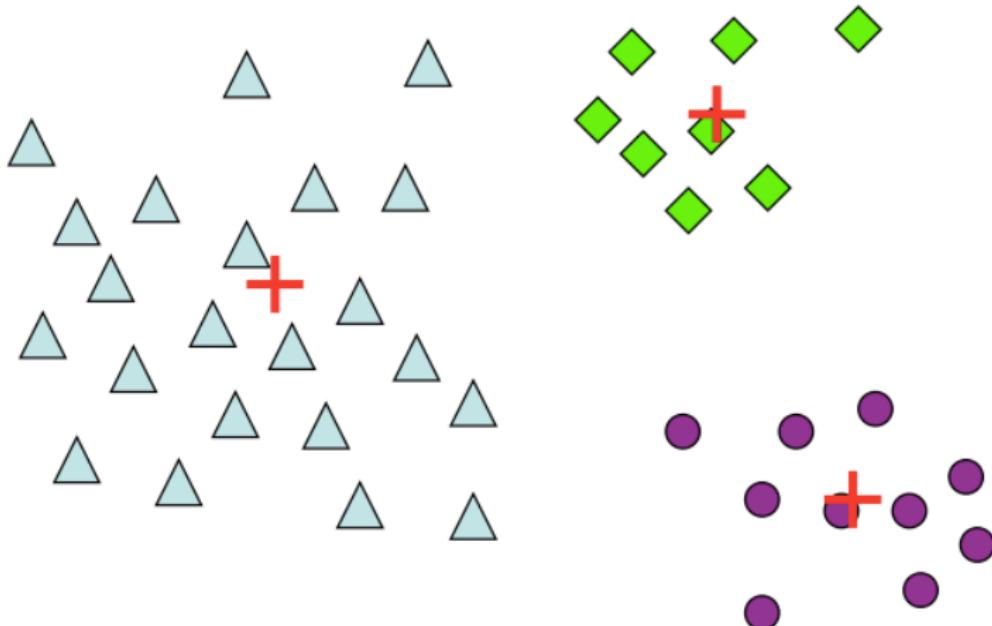
An Example

2nd recompute centroids and re-cluster



An Example

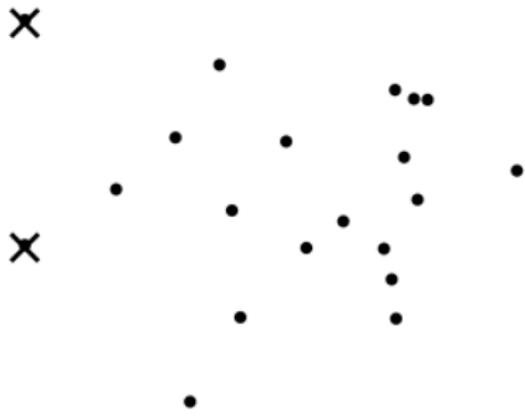
3rd (final) recompute and re-cluster



Worked Example: Set of to be clustered

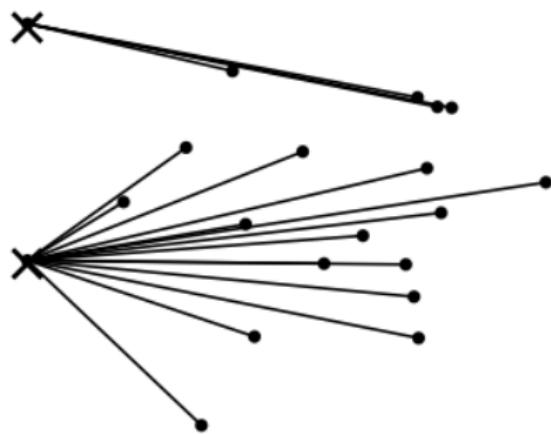


Worked Example: Random selection of initial centroids



Exercise: (i) Guess what the optimal clustering into two clusters is in this case; (ii) compute the centroids of the clusters

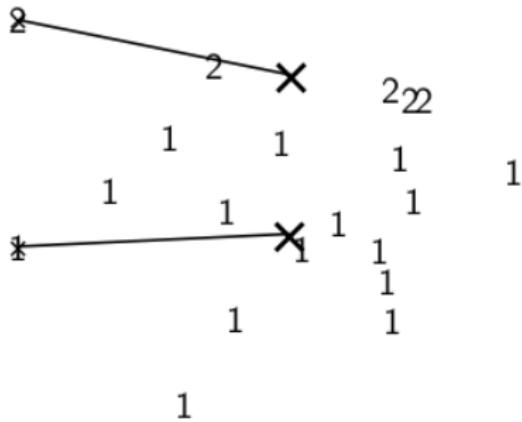
Worked Example: Assign points to closest center



Worked Example: Assignment

	*		
		2	22
	1	1	1
1	1	1	1
*		1	1
		1	1
			1

Worked Example: Recompute cluster centroids



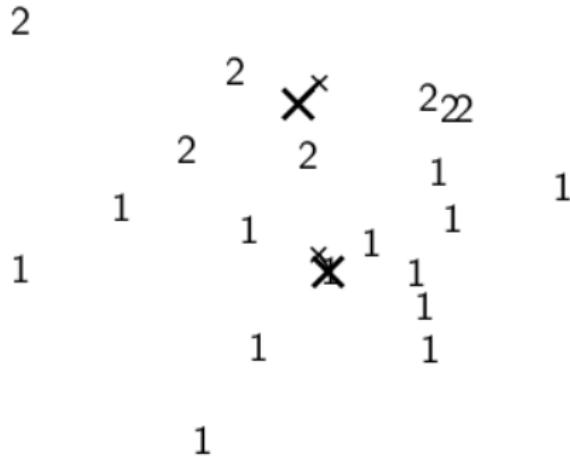
Worked Example: Assign points to closest centroid



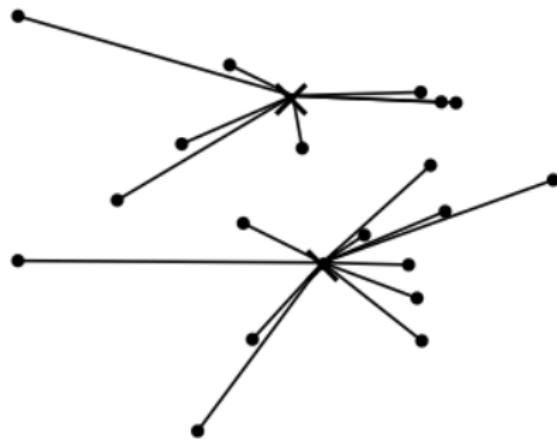
Worked Example: Assignment

$$\begin{array}{ccccccccc} & & 2 & & & & & & \\ & & \times & & & & & & \\ & 2 & & & & & 2 & & \\ & & 2 & & 2 & & & & \\ & 1 & & 1 & & 1 & & 1 & \\ & & \times & & 1 & & 1 & & \\ & & 1 & & & & 1 & & \\ & & & & & & 1 & & \\ & & & & & & & & \\ & & 1 & & & & & & \end{array}$$

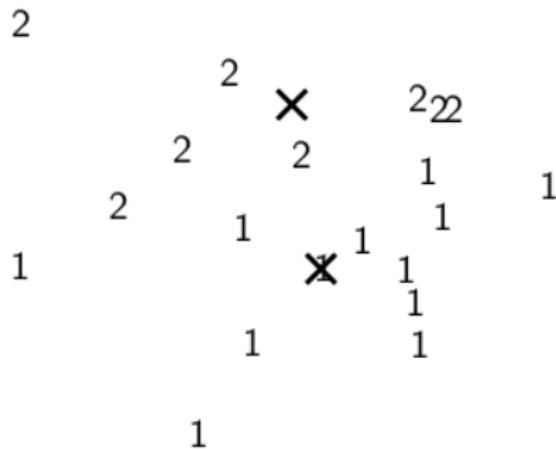
Worked Example: Recompute cluster centroids



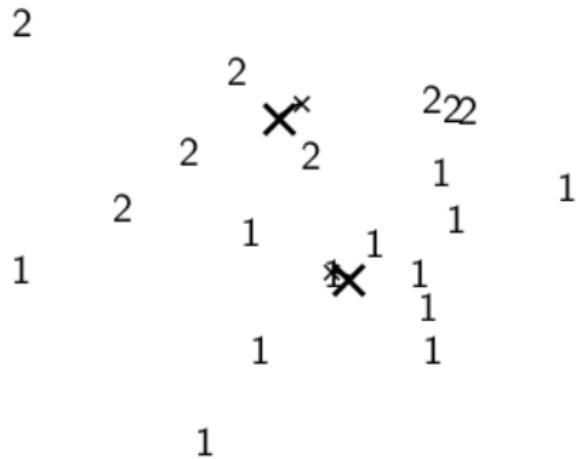
Worked Example: Assign points to closest centroid



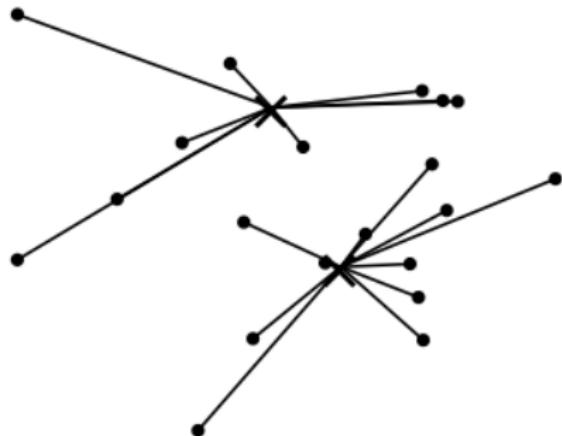
Worked Example: Assignment



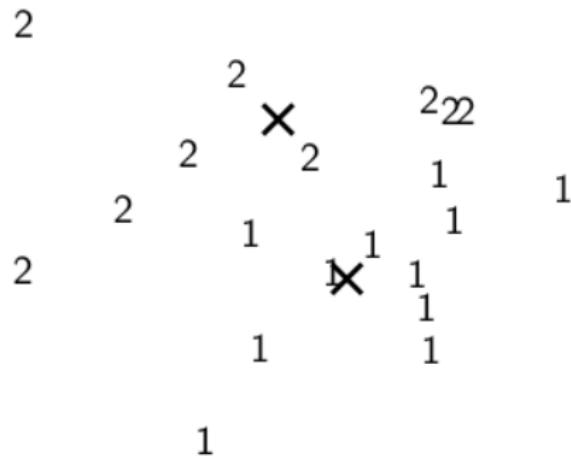
Worked Example: Recompute cluster centroids



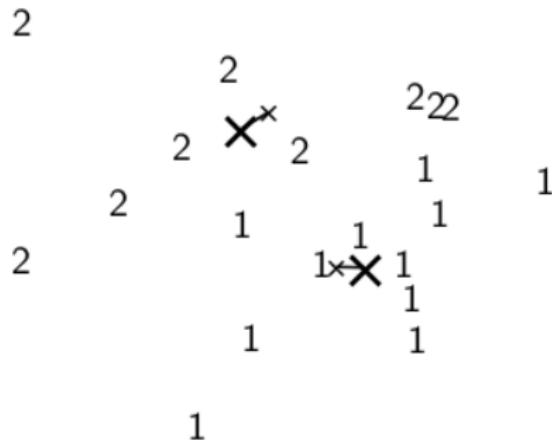
Worked Example: Assign points to closest centroid



Worked Example: Assignment



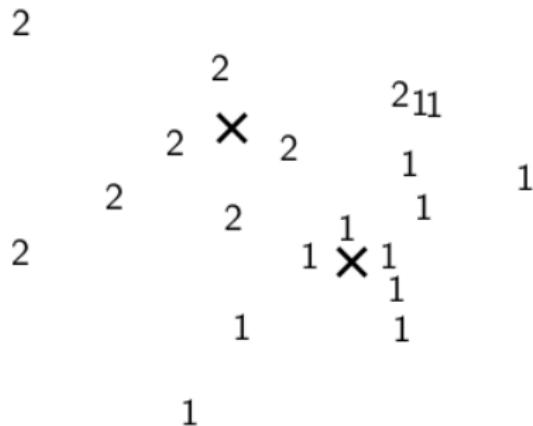
Worked Example: Recompute cluster centroids



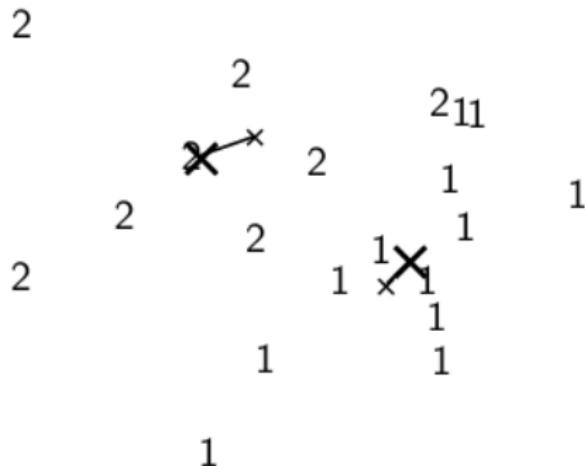
Worked Example: Assign points to closest centroid



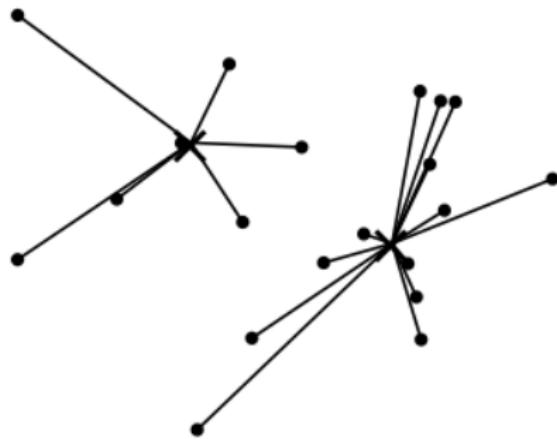
Worked Example: Assignment



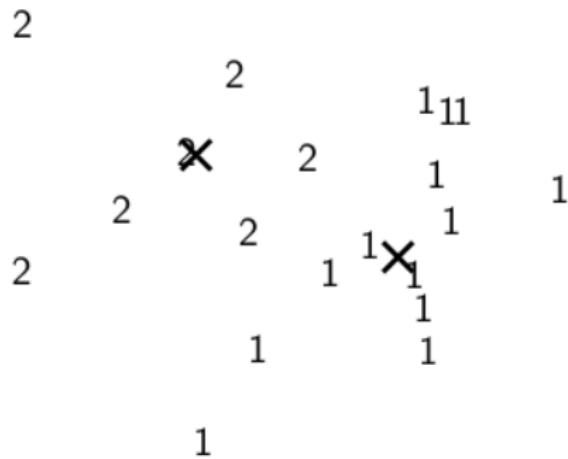
Worked Example: Recompute cluster centroids



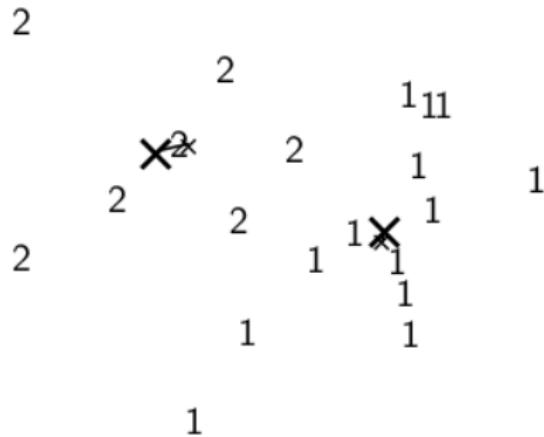
Worked Example: Assign points to closest centroid



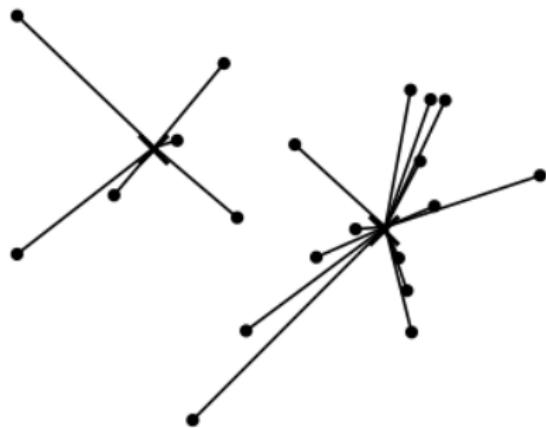
Worked Example: Assignment



Worked Example: Recompute cluster centroids



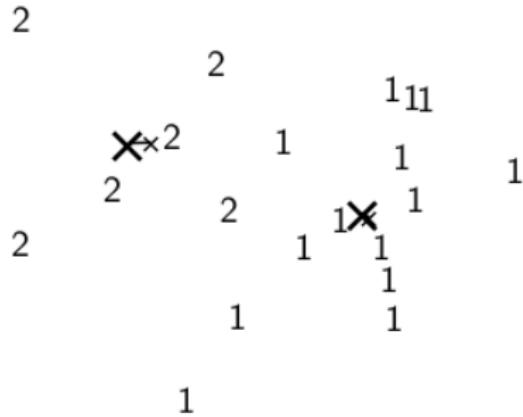
Worked Example: Assign points to closest centroid



Worked Example: Assignment

$$\begin{matrix} & & 2 \\ & & 2 \\ & & \cancel{\times}^2 & 1 & 1_{11} \\ & 2 & 2 & 1 & 1 & 1 \\ 2 & & 1 & 1 \cancel{\times}^1 & 1 & \\ & & & 1 & 1 & \\ & & & 1 & 1 & \\ & & & & 1 & \\ & & & & & 1 \\ & & & & & & 1 \\ & & & & & & & 1 \end{matrix}$$

Worked Example: Recompute cluster centroids



Demo

demo

Crucial K-means decisions

- $K = ?$
- initialisation (seeds)

Choosing K

- K can be given/known
 - experience for a given task
 - application/task requirements
- it might be not known, and finding the “right” number of clusters is part of the problem:
 - given docs, partition into an “appropriate” number of subsets
- tradeoff between having more clusters (more focus in each cluster) and having too many clusters
- (Elbow method)

Initialising centroids

- random
 - fast and easy, but often poor results
- run random multiple times, take best
 - slower, and still no guarantee of results
- pre-conditioning
 - remove outliers
- choose seeds algorithmically
 - run hierarchical clustering on sample points and use resulting centroids
 - works well on small samples and for few initial centroids

Evaluation

How do we evaluate clusters?

Evaluation

How do we evaluate clusters?

If we have labels, we can tell how well the clustering matches the gold standard classes (external evaluation).

Note: we only use *the partition* provided by the gold standard, *not the class labels*.

- **purity** is a simple and transparent evaluation measure.
- **rand index (RI)** penalizes both false positive and false negative decisions during clustering
- **F-measure (F-measure)** supports differential weighting of FPs and FNs
- **normalized mutual information (NMI)** can be information-theoretically interpreted.

Purity

- N = total number of documents
- each cluster is assigned to the class which is most frequent in the cluster
- accuracy of this assignment is measured by counting the number of correctly assigned documents and dividing by N .

Purity

- N = total number of documents
- each cluster is assigned to the class which is most frequent in the cluster
- accuracy of this assignment is measured by counting the number of correctly assigned documents and dividing by N .

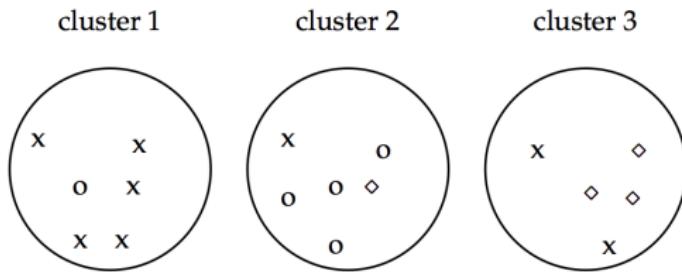
$$\text{purity}(\Omega, \mathbb{C}) = \frac{1}{N} \sum_K \max_j |\omega_k \cap c_j|$$

- $\Omega = \{\omega_1, \omega_2, \dots, \omega_K\}$ = set of clusters
- $\mathbb{C} = \{c_1, c_2, \dots, c_j\}$ = set of classes

Purity

$$purity(\Omega, \mathbb{C}) = \frac{1}{N} \sum_K \max_j |\omega_k \cap c_j|$$

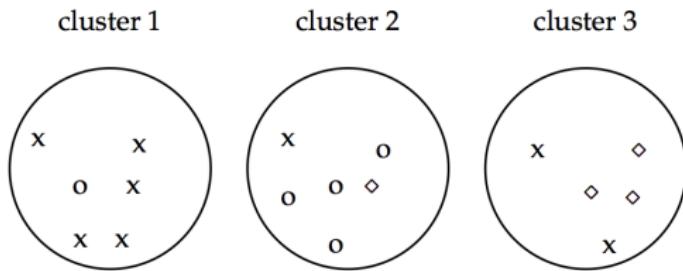
- $\Omega = \{\omega_1, \omega_2, \dots, \omega_K\}$ = set of clusters
- $\mathbb{C} = \{c_1, c_2, \dots, c_j\}$ = set of classes



Purity

$$purity(\Omega, \mathbb{C}) = \frac{1}{N} \sum_K \max_j |\omega_k \cap c_j|$$

- $\Omega = \{\omega_1, \omega_2, \dots, \omega_K\}$ = set of clusters
- $\mathbb{C} = \{c_1, c_2, \dots, c_j\}$ = set of classes

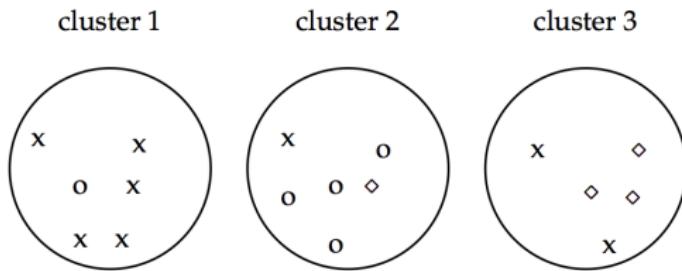


5 +

Purity

$$purity(\Omega, \mathbb{C}) = \frac{1}{N} \sum_K \max_j |\omega_k \cap c_j|$$

- $\Omega = \{\omega_1, \omega_2, \dots, \omega_K\}$ = set of clusters
- $\mathbb{C} = \{c_1, c_2, \dots, c_j\}$ = set of classes

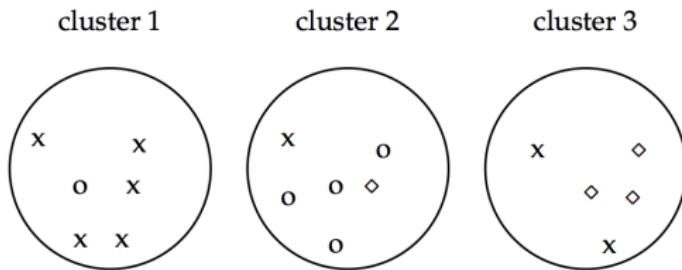


5 + 4 +

Purity

$$purity(\Omega, \mathbb{C}) = \frac{1}{N} \sum_K \max_j |\omega_k \cap c_j|$$

- $\Omega = \{\omega_1, \omega_2, \dots, \omega_K\}$ = set of clusters
- $\mathbb{C} = \{c_1, c_2, \dots, c_j\}$ = set of classes

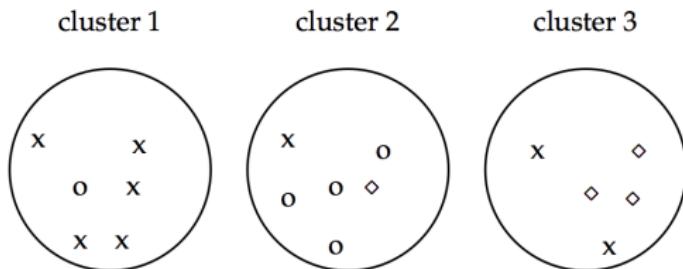


$5 + 4 + 3$

Purity

$$purity(\Omega, \mathbb{C}) = \frac{1}{N} \sum_K \max_j |\omega_k \cap c_j|$$

- $\Omega = \{\omega_1, \omega_2, \dots, \omega_K\}$ = set of clusters
- $\mathbb{C} = \{c_1, c_2, \dots, c_j\}$ = set of classes



$$\frac{(5 + 4 + 3)}{17} = 0.71$$

Purity

- N = total number of documents
- each cluster is assigned to the class which is most frequent in the cluster
- accuracy of this assignment is measured by counting the number of correctly assigned documents and dividing by N .

what is a sure case where *purity* = 1?

Purity

- N = total number of documents
- each cluster is assigned to the class which is most frequent in the cluster
- accuracy of this assignment is measured by counting the number of correctly assigned documents and dividing by N .

what is a sure case where *purity* = 1?

the larger the number of clusters...

Rand Index

RI measures the percentage of decisions that are correct

Rand Index

RI measures the percentage of decisions that are correct

Consider the errors:

- FP
- FN

Rand Index

RI measures the percentage of decisions that are correct

Consider the errors:

- FP – a false positive decision assigns two dissimilar documents to the same cluster.
- FN – a false negative decision assigns two similar documents to different clusters.

Rand Index

RI measures the percentage of decisions that are correct

Consider the errors:

- FP – a false positive decision assigns two dissimilar documents to the same cluster.
- FN – a false negative decision assigns two similar documents to different clusters.

$$\frac{TP + TN}{TP + FP + TN + FN}$$

Rand Index

RI measures the percentage of decisions that are correct

Consider the errors:

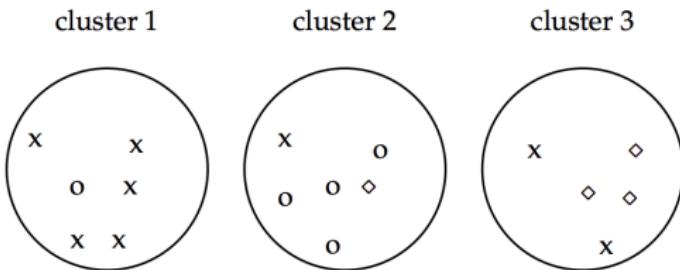
- FP – a false positive decision assigns two dissimilar documents to the same cluster.
- FN – a false negative decision assigns two similar documents to different clusters.

$$\frac{TP + TN}{TP + FP + TN + FN}$$

(accuracy)

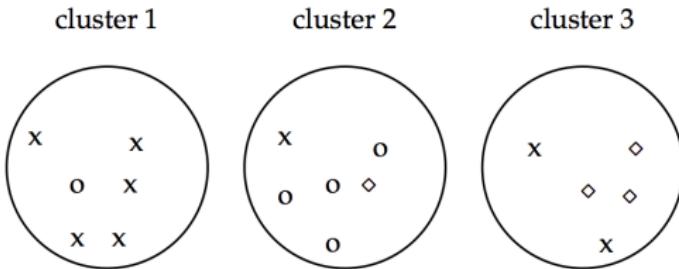
Rand Index

$$\frac{TP + TN}{TP + FP + TN + FN}$$



Rand Index

$$\frac{TP + TN}{TP + FP + TN + FN}$$



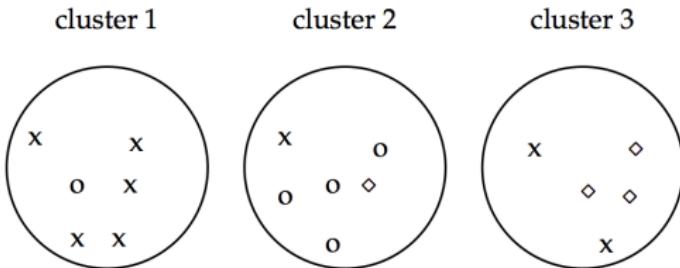
- positives = number of pairs that are in the same cluster

$$TP + FP = \binom{6}{2} + \binom{6}{2} + \binom{5}{2} = 40$$

$$TP = \binom{5}{2} + \binom{4}{2} + \binom{3}{2} + \binom{2}{2} = 20$$

Rand Index

$$\frac{TP + TN}{TP + FP + TN + FN}$$



	same cluster	different cluster
same class	TP=20	FN=24
different class	FP=20	TN=72

$$RI = \frac{20 + 72}{20 + 20 + 24 + 72} = 0.68$$

F-measure

$$RI = \frac{TP + TN}{TP + FP + TN + FN} = \text{equal weight to all errors (FPs and FNs)}$$

BUT: separating similar documents is sometimes worse than putting pairs of dissimilar documents in the same cluster
(especially true in IR: favour recall over precision)

F-measure

$$\frac{TP}{TP + FP}$$

$$\frac{TP}{TP + FN}$$

F-measure

$$\frac{TP}{TP + FP} \quad (\text{precision})$$

$$\frac{TP}{TP + FN} \quad (\text{recall})$$

$$F = \frac{2PR}{P + R} \quad (\text{F-measure, also called } F_1)$$

F-measure

$$F = \frac{2PR}{P + R} \quad (\text{F-measure, also called } F_1)$$

$$F = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

- if $\beta = 1$, then standard F-measure (F_1)
- if $\beta < 1$, then precision emphasised
- if $\beta > 1$, then recall emphasised

F-measure

$$F = \frac{(\beta^2 + 1)PR}{\beta^2P + R}$$

- if $\beta = 1$, then standard F-measure (F_1)
- if $\beta < 1$, then precision emphasised
- if $\beta > 1$, then recall emphasised

	same cluster	different cluster
same class	TP=20	FN=24
different class	FP=20	TN=72

$$\begin{aligned} F_1 &= 0.48 \\ F_5 &= 0.46 \end{aligned}$$

V-measure

- Homogeneity
Each cluster should include members of as few classes as possible
- Completeness
Each class should be represented in as few clusters as possible.

Normalised Mutual Information

- how much information does the clustering contain about the classification?
- singleton clusters (number of clusters = number of docs) have maximum MI
- therefore: normalise by entropy of clusters and classes

What do clusters mean?

if we cluster our reviews, will it be on topic or sentiment?
or possibly something else? How can we know they're not clustered on
say the author's gender? or Age?

Clustering with scikit-learn

<http://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>

Clustering with scikit-learn

<http://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>

`sklearn.metrics.homogeneity_completeness_v_measure`

`sklearn.metrics. homogeneity_completeness_v_measure (labels_true, labels_pred)`

[source]

Compute the homogeneity and completeness and V-Measure scores at once.

Those metrics are based on normalized conditional entropy measures of the clustering labeling to evaluate given the knowledge of a Ground Truth class labels of the same samples.

A clustering result satisfies homogeneity if all of its clusters contain only data points which are members of a single class.

A clustering result satisfies completeness if all the data points that are members of a given class are elements of the same cluster.

Both scores have positive values between 0.0 and 1.0, larger values being desirable.

Those 3 metrics are independent of the absolute values of the labels: a permutation of the class or cluster label values won't change the score values in any way.

V-Measure is furthermore symmetric: swapping `labels_true` and `label_pred` will give the same score. This does not hold for homogeneity and completeness.