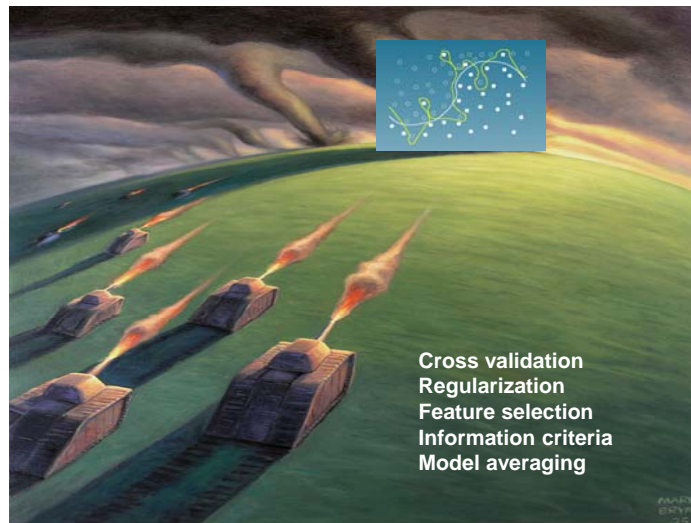


$$G_D < G_S + O\sqrt{\frac{1}{n}}$$

## The battle against overfitting



© Eric Xing @ CMU, 2006-2011

1

## 3. Feature Selection

$$\eta = (g^T x + \lambda) / \theta$$

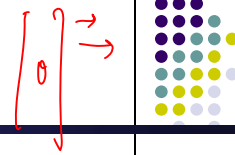
- Imagine that you have a supervised learning problem where the number of features  $d$  is very large (perhaps  $d \gg \text{\#samples}$ ), but you suspect that there is only a small number of features that are "**relevant**" to the learning task.
- VC-theory can tell you that this scenario is likely to lead to high generalization error – the learned model will potentially overfit unless the training set is fairly large.
- So lets get rid of useless parameters!

$$\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

© Eric Xing @ CMU, 2006-2011

2

# Feature selection schemes



- Given  $n$  features, there are  $2^n$  possible feature subsets (why?)
- Thus feature selection can be posed as a model selection problem over  $2^n$  possible models.
- For large values of  $n$ , it's usually too expensive to explicitly enumerate over and compare all  $2^n$  models. Some heuristic search procedure is used to find a good feature subset.
- Three general approaches:
  - Filter: i.e., direct feature ranking, but taking no consideration of the subsequent learning algorithm
    - add (from empty set) or remove (from the full set) features one by one based on  $S(i)$
    - Cheap, but is subject to local optimality and may be unrobust under different classifiers
  - Wrapper: determine the (inclusion or removal of) features based on performance under the learning algorithms to be used. See next slide
  - Simultaneous learning and feature selection.
    - E.x.  $L_1$  regularized LR, Bayesian feature selection (will not cover in this class), etc.

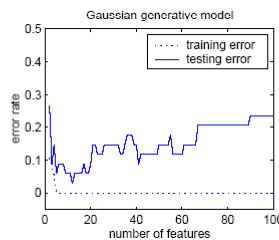
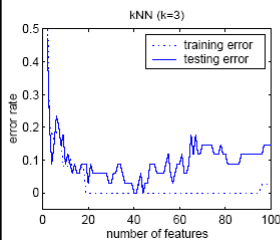
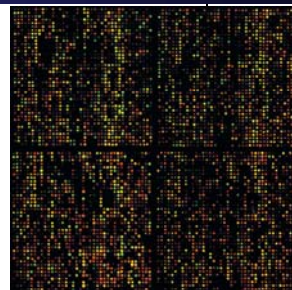
© Eric Xing @ CMU, 2006-2011

3

## Case study [Xing et al, 2001]



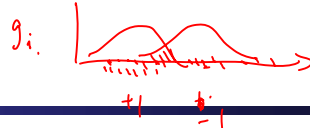
- The case:
  - 7130 genes from a microarray dataset
  - 72 samples
  - 47 type I Leukemias (called ALL) and 25 type II Leukemias (called AML)
- Three classifier:
  - kNN
  - Gaussian classifier
  - Logistic regression



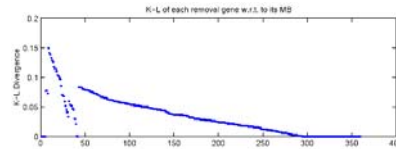
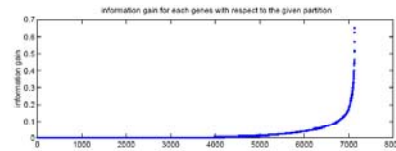
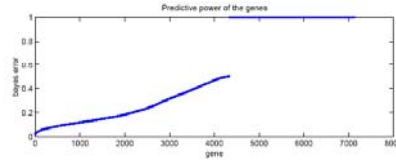
© Eric Xing @ CMU, 2006-2011

4

# Feature Ranking



- Bayes error of each gene
- information gain for each genes with respect to the given partition
- KL of each removal gene w.r.t. to its MB



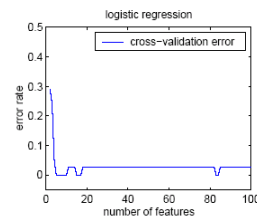
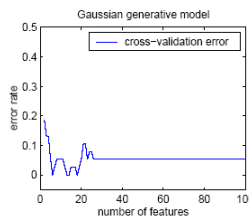
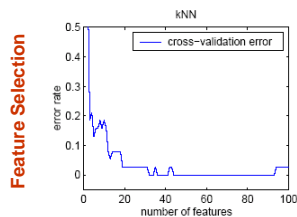
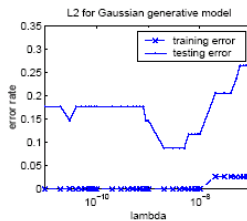
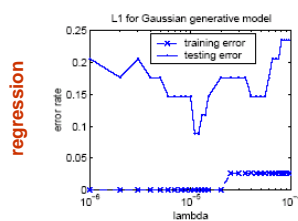
© Eric Xing @ CMU, 2006-2011

5

# Regularization vs. Feature Selection



- Explicit feature selection often outperform regularization



© Eric Xing @ CMU, 2006-2011

6

## 4. Information criterion



- Suppose we are trying select among several different models for a learning problem.
- The Problem:
  - Given model family  $\mathcal{F} = \{M_1, M_2, \dots, M_I\}$ , find  $M_i \in \mathcal{F}$  s.t.
 
$$M_i = \arg \max_{M \in \mathcal{F}} J(D, M)$$
- We can design  $J$  that not only reflect the predictive loss, but also the amount of information  $M_k$  can hold

$$E_y \left[ D(f \parallel g(x | \theta_{ML}(y))) \right]$$

$\uparrow$  true       $\uparrow$  est

© Eric Xing @ CMU, 2006-2011

7

## AIC and TIC



- AIC (An information criterion, not **Akaike** information criterion)

$$A = \log g(x | \hat{\theta}(y)) - k$$

where  $k$  is the number of parameters in the model

- TIC (Takeuchi information criterion)

$$A = \log g(x | \hat{\theta}(y)) - \text{tr}(I(\theta_0)\Sigma)$$

where

$$\theta_0 = \arg \min D(f \parallel g(\cdot | \theta)) \quad I(\theta_0) = -E_x \left[ \frac{\partial^2 \log g(x | \theta)}{\partial \theta \partial \theta^T} \right] \Big|_{\theta=\theta_0} \quad \Sigma = E_y (\hat{\theta}(y) - \theta_0)(\hat{\theta}(y) - \theta_0)^T$$

- We can approximate these terms in various ways (e.g., using the bootstrap)
- $\text{tr}(I(\theta_0)\Sigma) \approx k$

© Eric Xing @ CMU, 2006-2011

8

## 5. Bayesian Model Averaging



- Recall the Bayesian Theory: (e.g., for data  $D$  and model  $M$ )

$$P(M|D) = P(D|M)P(M)/P(D)$$

- the posterior equals to the likelihood times the prior, up to a constant.
- Assume that  $P(M)$  is uniform and notice that  $P(D)$  is constant, we have the following criteria:

$$P(D|M) = \int_{\theta} P(D|\theta, M)P(\theta|M)d\theta$$

- A few steps of approximations (you will see this in advanced ML class in later semesters) give you this:

$$P(D|M) \approx \log P(D|\hat{\theta}_{ML}) - \frac{k}{2} \log N$$

where  $N$  is the number of data points in  $D$ .

© Eric Xing @ CMU, 2006-2011

9

## Summary



- Structural risk minimization
- Bias-variance decomposition
- The battle against overfitting:
  - Cross validation
  - Regularization
  - Feature selection
  - Information criteria
  - Model averaging

© Eric Xing @ CMU, 2006-2011

10

# Machine Learning

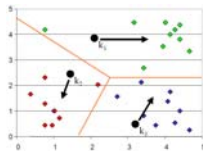
10-701/15-781, Fall 2011

## Clustering and Distance Metrics

Eric Xing

Lecture 8, October 5, 2011

Reading: Chap. 9, C.B book

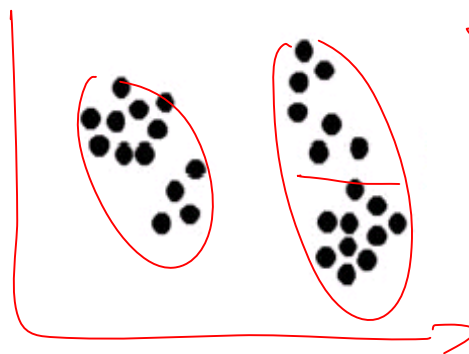


© Eric Xing @ CMU, 2006-2011

11



## What is clustering?



$$f: x \rightarrow y \quad (x^i, y^i)$$

$$f: x \rightarrow y \quad (x^i)$$

- Are there any “grouping” them ?
- What is each group ?
- How many ?
- How to identify them?

© Eric Xing @ CMU, 2006-2011

12



# What is clustering?



- Clustering: the process of grouping a set of objects into classes of similar objects
  - high intra-class similarity
  - low inter-class similarity
  - It is the commonest form of unsupervised learning
- Unsupervised learning = learning from raw (unlabeled, unannotated, etc) data, as opposed to supervised data where a classification of examples is given
- A common and important task that finds many applications in Science, Engineering, information Science, and other places
  - Group genes that perform the same function
  - Group individuals that has similar political view
  - Categorize documents of similar topics
  - Ideality similar objects from pictures

© Eric Xing @ CMU, 2006-2011

13

# Examples



- People



- Images



- Language

Piotr Pyotr Petros Pietro Pedro Pierre Piero Peter Peder Peka Peadar

- species



© Eric Xing @ CMU, 2006-2011

14

## Issues for clustering

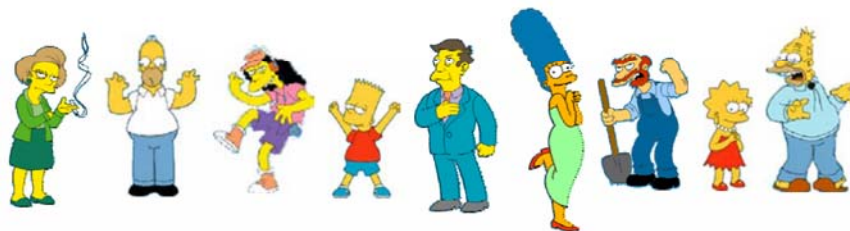


- What is a natural grouping among these objects?
  - Definition of "groupness"
- What makes objects "related"?
  - Definition of "similarity/distance"
- Representation for objects
  - Vector space? Normalization?
- How many clusters?
  - Fixed a priori?
  - Completely data driven?
    - Avoid "trivial" clusters - too large or small
- Clustering Algorithms
  - Partitional algorithms
  - Hierarchical algorithms
- Formal foundation and convergence

© Eric Xing @ CMU, 2006-2011

15

## What is a natural grouping among these objects?



Clustering is subjective



Simpson's Family



School Employees



Females



Males

© Eric Xing @ CMU, 2006-2011

16



## What is Similarity?



Hard to define!  
But we know it  
when we see it

- The real meaning of similarity is a philosophical question. We will take a more pragmatic approach
- Depends on representation and algorithm. For many rep./alg., easier to think in terms of a distance (rather than similarity) between vectors.

© Eric Xing @ CMU, 2006-2011

17

## What properties should a distance measure have?



- $D(A,B) = D(B,A)$  *Symmetry*
- $D(A,A) = 0$  *Constancy of Self-Similarity*
- $D(A,B) = 0$  iff  $A = B$  *Positivity Separation*
- $D(A,B) \leq D(A,C) + D(B,C)$  *Triangular Inequality*

© Eric Xing @ CMU, 2006-2011

18

## Intuitions behind desirable distance measure properties



- $D(A,B) = D(B,A)$  *Symmetry*
  - Otherwise you could claim "Alex looks like Bob, but Bob looks nothing like Alex"
- $D(A,A) = 0$  *Constancy of Self-Similarity*
  - Otherwise you could claim "Alex looks more like Bob, than Bob does"
- $D(A,B) = 0$  iff  $A = B$  *Positivity Separation*
  - Otherwise there are objects in your world that are different, but you cannot tell apart.
- $D(A,B) \leq D(A,C) + D(B,C)$  *Triangular Inequality*
  - Otherwise you could claim "Alex is very like Bob, and Alex is very like Carl, but Bob is very unlike Carl"

© Eric Xing @ CMU, 2006-2011

19

## Distance Measures: Minkowski Metric



- Suppose two object  $x$  and  $y$  both have  $p$  features

$$x = (x_1, x_2, \dots, x_p)$$

$$y = (y_1, y_2, \dots, y_p)$$

- The Minkowski metric is defined by

$$d(x, y) = \sqrt[r]{\sum_{i=1}^p |x_i - y_i|^r}$$

- Most Common Minkowski Metrics

1,  $r = 2$  (Euclidean distance)

$$d(x, y) = \sqrt{\sum_{i=1}^p |x_i - y_i|^2}$$

2,  $r = 1$  (Manhattan distance)

$$d(x, y) = \sum_{i=1}^p |x_i - y_i|$$

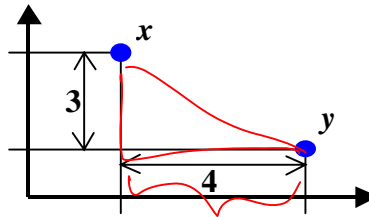
3,  $r = +\infty$  ("sup" distance)

$$d(x, y) = \max_{1 \leq i \leq p} |x_i - y_i|$$

© Eric Xing @ CMU, 2006-2011

20

## An Example



- 1: Euclidean distance:  $\sqrt{4^2 + 3^2} = 5$ .
- 2: Manhattan distance:  $4 + 3 = 7$ .
- 3: "sup" distance:  $\max\{4, 3\} = 4$ .

© Eric Xing @ CMU, 2006-2011

21

## Hamming distance

- Manhattan distance is called *Hamming distance* when all features are binary.
- [Gene Expression Levels Under 17 Conditions \(1-High,0-Low\)](#)

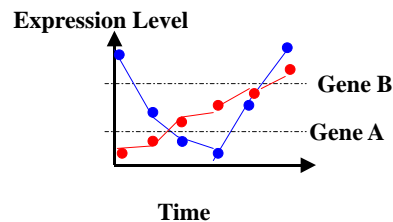
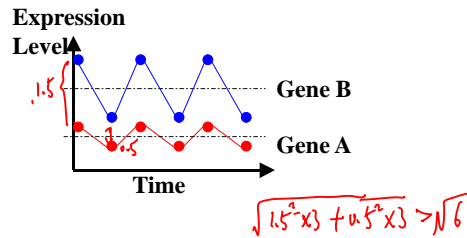
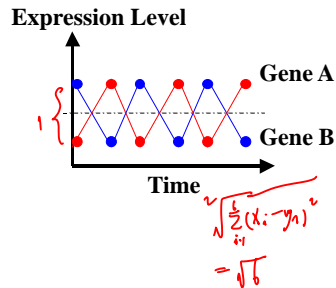
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
GeneA	0	1	1	0	0	1	0	0	1	0	0	1	1	1	0	0	1
GeneB	0	1	1	1	0	0	0	0	1	1	1	1	1	1	0	1	1

Hamming Distance:  $\#(01) + \#(10) = 4 + 1 = 5$ .

© Eric Xing @ CMU, 2006-2011

22

## Similarity Measures: Correlation Coefficient



© Eric Xing @ CMU, 2006-2011

23

## Similarity Measures: Correlation Coefficient



- Pearson correlation coefficient

$$s(x, y) = \frac{\sum_{i=1}^p (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^p (x_i - \bar{x})^2 \times \sum_{i=1}^p (y_i - \bar{y})^2}}$$

$$\text{where } \bar{x} = \frac{1}{p} \sum_{i=1}^p x_i \text{ and } \bar{y} = \frac{1}{p} \sum_{i=1}^p y_i.$$

$$|s(x, y)| \leq 1$$

- Special case: cosine distance

$$s(x, y) = \frac{\bar{x} \cdot \bar{y}}{|\bar{x}| \cdot |\bar{y}|}$$

$$\bar{x} = \frac{1}{N} \sum x_i$$



© Eric Xing @ CMU, 2006-2011

24

## Edit Distance:

### A generic technique for measuring similarity

- To measure the similarity between two objects, transform one of the objects into the other, and measure how much effort it took. The measure of effort becomes the distance measure.

#### The distance between Patty and Selma.

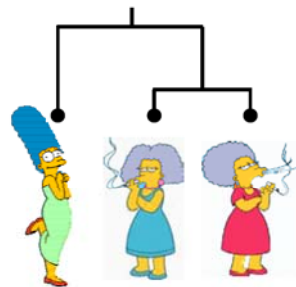
Change dress color, 1 point  
Change earring shape, 1 point  
Change hair part, 1 point

$$D(\text{Patty}, \text{Selma}) = 3$$

#### The distance between Marge and Selma.

Change dress color, 1 point  
Add earrings, 1 point  
Decrease height, 1 point  
Take up smoking, 1 point  
Lose weight, 1 point

$$D(\text{Marge}, \text{Selma}) = 5$$

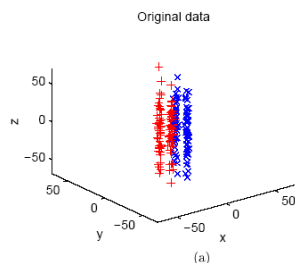


This is called the  
Edit distance  
or the  
Transformation distance

© Eric Xing @ CMU, 2006-2011

25

## Learning Distance Metric



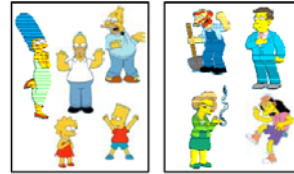
More later ...

© Eric Xing @ CMU, 2006-2011

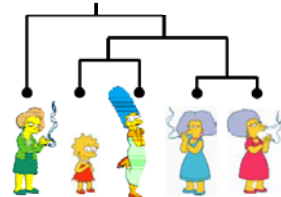
26

# Clustering Algorithms

- Partitional algorithms
  - Usually start with a random (partial) partitioning
  - Refine it iteratively
    - K means clustering
    - Mixture-Model based clustering



- Hierarchical algorithms
  - Bottom-up, agglomerative
  - Top-down, divisive

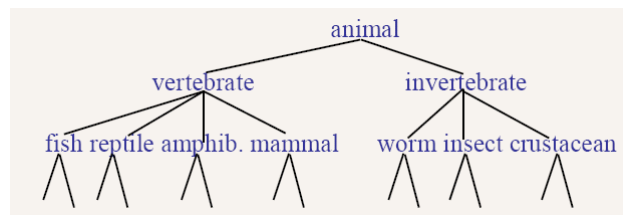


© Eric Xing @ CMU, 2006-2011

27

# Hierarchical Clustering

- Build a tree-based hierarchical taxonomy (dendrogram) from a set of documents.



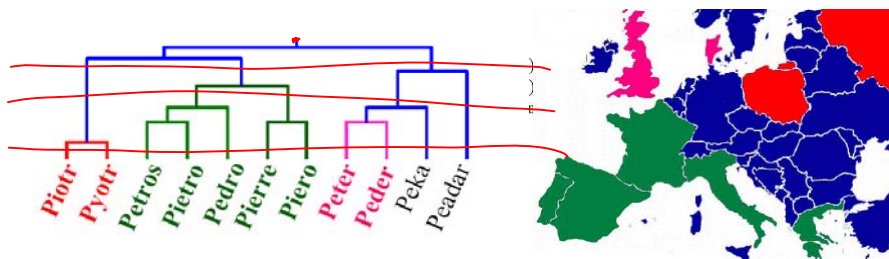
- Note that hierarchies are commonly used to organize information, for example in a web portal.
  - Yahoo! is hierarchy is manually created, we will focus on automatic creation of hierarchies in data mining.

© Eric Xing @ CMU, 2006-2011

28

# Dendrogram

- A Useful Tool for Summarizing Similarity Measurement
  - The similarity between two objects in a dendrogram is represented as the height of the lowest internal node they share.
- Clustering obtained by cutting the dendrogram at a desired level: each connected component forms a cluster.



© Eric Xing @ CMU, 2006-2011

29

# Hierarchical Clustering

- Bottom-Up Agglomerative Clustering
  - Starts with each obj in a separate cluster
  - then repeatedly joins the closest pair of clusters,
  - until there is only one cluster.

The history of merging forms a binary tree or hierarchy.
- Top-Down divisive
  - Starting with all the data in a single cluster,
  - Consider every possible way to divide the cluster into two. Choose the best division
  - And recursively operate on both sides.

© Eric Xing @ CMU, 2006-2011

30

# Closest pair of clusters

The distance between two clusters is defined as the distance between

- Single-Link
  - Nearest Neighbor: their closest members.
- Complete-Link
  - Furthest Neighbor: their furthest members.
- Centroid:
  - Clusters whose centroids (centers of gravity) are the most cosine-similar
- Average:
  - average of all cross-cluster pairs.

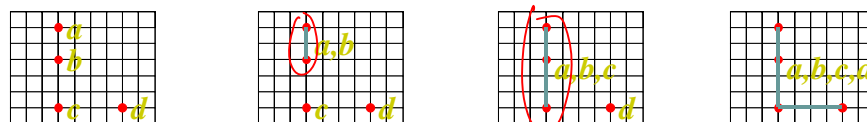


© Eric Xing @ CMU, 2006-2011

31

# Single-Link Method

## Euclidean Distance



(1)

(2)

(3)

	b	c	d
a	2	5	6
b		3	5
c			4

	b	c	d
a	2	5	6
b		3	5
c			4

	c	d
a, b	3	5
c		4

	d
a, b, c	4

Distance Matrix



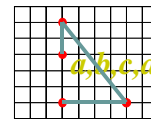
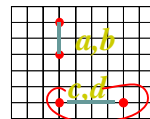
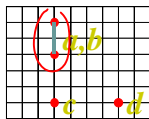
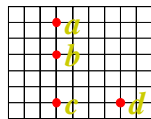
© Eric Xing @ CMU, 2006-2011

32



# Complete-Link Method

## Euclidean Distance



(1)

(2)

(3)

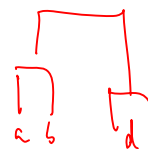
	<i>b</i>	<i>c</i>	<i>d</i>
<i>a</i>	2	5	6
<i>b</i>		3	5
<i>c</i>			4

	<i>b</i>	<i>c</i>	<i>d</i>
<i>a</i>	2	5	6
<i>b</i>		3	5
<i>c</i>			4

	<i>c</i>	<i>d</i>
<i>a, b</i>	5	6
<i>c</i>		4

	<i>c, d</i>
<i>a, b</i>	6

Distance Matrix

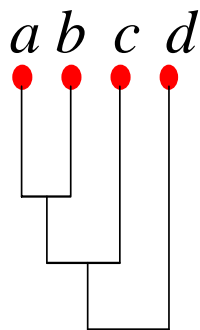


© Eric Xing @ CMU, 2006-2011

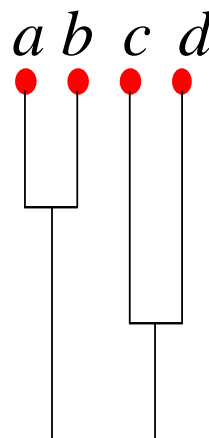
33

# Dendrograms

## Single-Link



## Complete-Link



© Eric Xing @ CMU, 2006-2011

34

## Computational Complexity

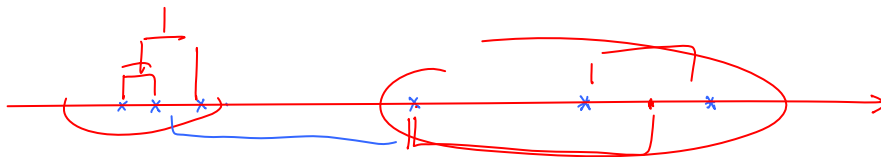


- In the first iteration, all HAC methods need to compute similarity of all pairs of  $n$  individual instances which is  $O(n^2)$ .
- In each of the subsequent  $n-2$  merging iterations, compute the distance between the most recently created cluster and all other existing clusters.
- In order to maintain an overall  $O(n^2)$  performance, computing similarity to each other cluster must be done in constant time.
- Else  $O(n^2 \log n)$  or  $O(n^3)$  if done naively

© Eric Xing @ CMU, 2006-2011

35

## Local-optimality of HAC



© Eric Xing @ CMU, 2006-2011

36

# Partitioning Algorithms



- Partitioning method: Construct a partition of  $n$  objects into a set of  $K$  clusters
- Given: a set of objects and the number  $K$
- Find: a partition of  $K$  clusters that optimizes the chosen partitioning criterion
  - Globally optimal: exhaustively enumerate all partitions
  - Effective heuristic methods: K-means and K-medoids algorithms

# K-Means



## Algorithm

1. Decide on a value for  $k$ .
2. Initialize the  $k$  cluster centers randomly if necessary.
3. Decide the class memberships of the  $N$  objects by assigning them to the nearest cluster centroids (aka the center of gravity or mean)

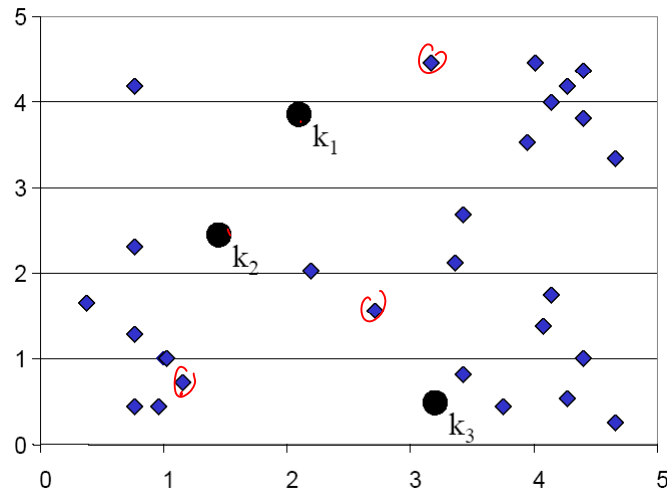
$$\vec{\mu}_k = \frac{1}{C_k} \sum_{i \in C_k} \vec{x}_i$$

*$d_{ik} =$*   *$C_i = \operatorname{argmax} d_{ik}$*

4. Re-estimate the  $k$  cluster centers, by assuming the memberships found above are correct.
5. If none of the  $N$  objects changed membership in the last iteration, exit. Otherwise go to 3.

## K-means Clustering: Step 1

$\mu_k$

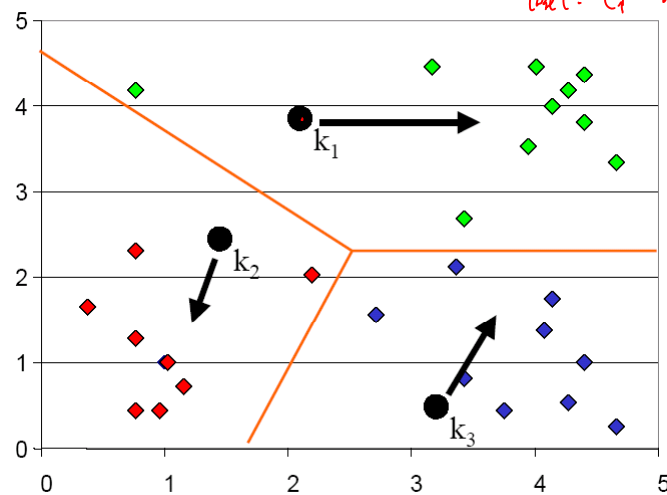


© Eric Xing @ CMU, 2006-2011

39

## K-means Clustering: Step 2

$v_i \text{ dist } k = \sqrt{(x_i - \mu_k)^2}$   
label:  $c_i = \text{argmin}_k \text{dist}$

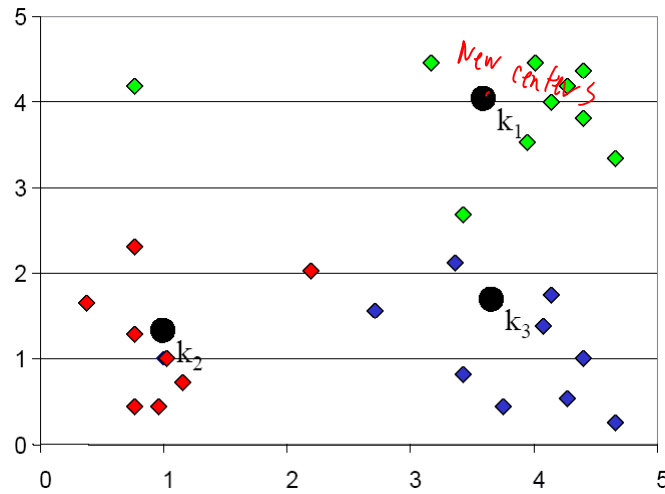


© Eric Xing @ CMU, 2006-2011

40

## K-means Clustering: Step 3

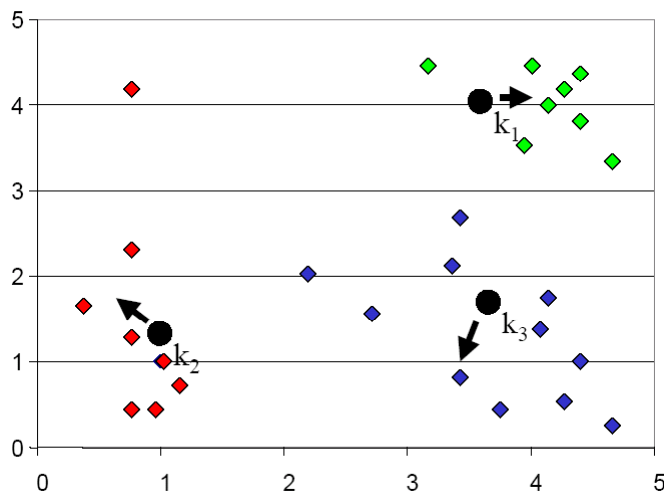
$$\mu_k = \frac{\sum x_i}{N_k}$$



© Eric Xing @ CMU, 2006-2011

41

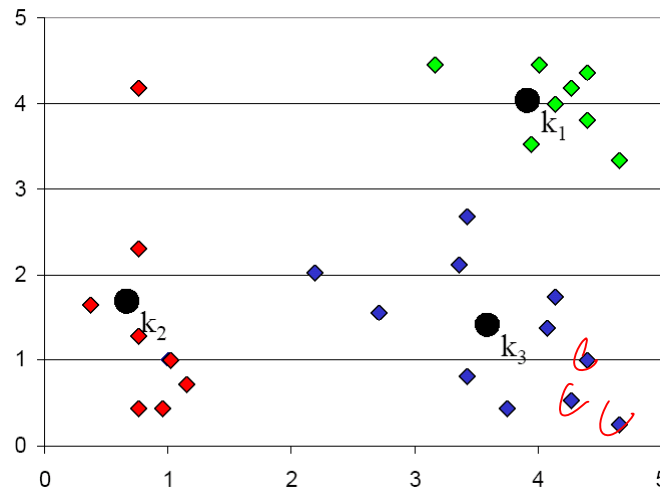
## K-means Clustering: Step 4



© Eric Xing @ CMU, 2006-2011

42

## K-means Clustering: Step 5



© Eric Xing @ CMU, 2006-2011

43

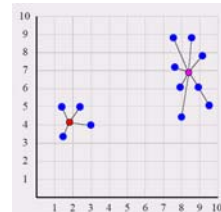
## Convergence

- Why should the K-means algorithm ever reach a fixed point?
  - -- A state in which clusters don't change.
- K-means is a special case of a general procedure known as the Expectation Maximization (EM) algorithm.
  - EM is known to converge.
  - Number of iterations could be large.

- Goodness measure

- sum of squared distances from cluster centroid:

$$SD_{K_i} = \sum_{j=1}^{m_k} \|x_{ij} - \mu_i\|^2 \quad SD_K = \sum_{i=1}^k SD_{K_i}$$



- Reassignment monotonically decreases SD since each vector is assigned to the closest centroid.

© Eric Xing @ CMU, 2006-2011

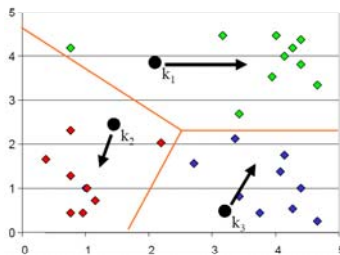
44

## Time Complexity

- Computing distance between two objs is  $O(m)$  where  $m$  is the dimensionality of the vectors.
- Reassigning clusters:  $O(Kn)$  distance computations, or  $O(Knm)$ .
- Computing centroids: Each doc gets added once to some centroid:  $O(nm)$ .  $O(kmn)$
- Assume these two steps are each done once for  $l$  iterations:  $O(lKnm)$ .

## Seed Choice

- Results can vary based on random seed selection.



- Some seeds can result in poor convergence rate, or convergence to sub-optimal clusterings.
  - Select good seeds using a heuristic (e.g., doc least similar to any existing mean)
  - Try out multiple starting points (very important!!!)
  - Initialize with the results of another method.

## How Many Clusters?



- Number of clusters  $K$  is given
  - Partition  $n$  docs into predetermined number of clusters
- Finding the “right” number of clusters is part of the problem
  - Given objs, partition into an “appropriate” number of subsets.
  - E.g., for query results - ideal value of  $K$  not known up front - though UI may impose limits.
- Solve an optimization problem: penalize having lots of clusters
  - application dependent, e.g., compressed summary of search results list.
  - Information theoretic approaches: model-based approach
- Tradeoff between having more clusters (better focus within each cluster) and having too many clusters
- Nonparametric Bayesian Inference

© Eric Xing @ CMU, 2006-2011

47

## What Is A Good Clustering?



- Internal criterion: A good clustering will produce high quality clusters in which:
  - the intra-class (that is, intra-cluster) similarity is high
  - the inter-class similarity is low
  - The measured quality of a clustering depends on both the obj representation and the similarity measure used
- External criteria for clustering quality
  - Quality measured by its ability to discover some or all of the hidden patterns or latent classes in gold standard data
  - Assesses a clustering with respect to ground truth
  - Example:
    - Purity
    - entropy of classes in clusters (or mutual information between classes and clusters)

© Eric Xing @ CMU, 2006-2011

48



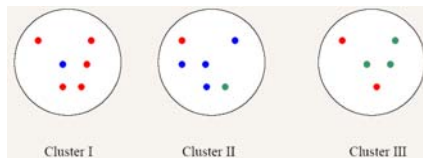
# External Evaluation of Cluster Quality



- Simple measure: **purity**, the ratio between the dominant class in the cluster and the size of cluster
- Assume documents with  $C$  gold standard classes, while our clustering algorithms produce  $K$  clusters,  $\omega_1, \omega_2, \dots, \omega_K$  with  $n_i$  members.

$$Purity(w_i) = \frac{1}{n_i} \max_j (n_{ij}) \quad j \in C$$

- Example



Cluster I: Purity =  $1/6 (\max(5, 1, 0)) = 5/6$

Cluster II: Purity =  $1/6 (\max(1, 4, 1)) = 4/6$

Cluster III: Purity =  $1/5 (\max(2, 0, 3)) = 3/5$

## Other measures



$$\frac{\sum_{i,j} ( \mathbb{I}(c(x_i) = c(x_j)) \cdot \mathbb{I}(\hat{c}(x_i) \neq \hat{c}(x_j)) )}{m(m-1)}$$

consistency.

$\hat{c}(x_i)$   $\hat{c}(x_j)$   
1 11  
2  
m total # of points

## Other partitioning Methods

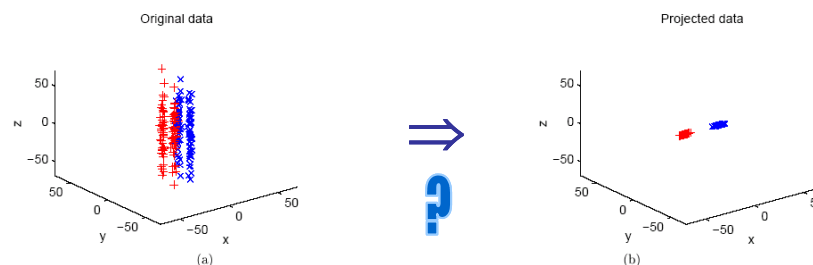


- Partitioning around mediods (PAM): instead of averages, use multidim medians as centroids (cluster “prototypes”). Dudoit and Freedland (2002).
- Self-organizing maps (SOM): add an underlying “topology” (neighboring structure on a lattice) that relates cluster centroids to one another. Kohonen (1997), Tamayo et al. (1999).
- Fuzzy k-means: allow for a “gradation” of points between clusters; soft partitions. Gash and Eisen (2002).
- Mixture-based clustering: implemented through an EM (Expectation-Maximization) algorithm. This provides soft partitioning, and allows for modeling of cluster centroids and shapes. Yeung et al. (2001), McLachlan et al. (2002)

© Eric Xing @ CMU, 2006-2011

51

## Semi-supervised Metric Learning



Xing et al, NIPS 2003

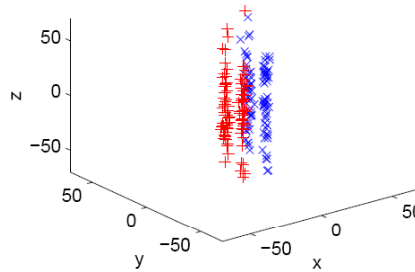
© Eric Xing @ CMU, 2006-2011

52

## What is a good metric?



- What is a good metric over the input space for learning and data-mining



- How to convey metrics sensible to a human user (e.g., dividing traffic along highway lanes rather than between overpasses, categorizing documents according to writing style rather than topic) to a computer data-miner using a systematic mechanism?

© Eric Xing @ CMU, 2006-2011

53

## Issues in learning a metric



- Data distribution is self-informing (E.g., lies in a sub-manifold)
  - Learning metric by finding an embedding of data in some space.
    - Con: does not reflect (changing) human subjectiveness.
- Explicitly labeled dataset offers clue for critical features
  - Supervised learning
    - Con: needs sizable homogeneous training sets.
- What about side information? (E.g., x and y look (or read) similar ...)
  - Providing small amount of qualitative and less structured side information is often much easier than stating explicitly a metric (what should be the metric for writing style?) or labeling a large set of training data.
- Can we learn a distance metric more informative than Euclidean distance using a small amount of side information?

© Eric Xing @ CMU, 2006-2011

54

## Distance Metric Learning



Side information:

Suppose for some set of points  $\{x_i\}_{i=1}^m \subseteq \mathbb{R}^n$ , we are given:

$\mathcal{S} : (x_i, x_j) \in \mathcal{S}$  if  $x_i$  and  $x_j$  are similar

$\mathcal{D} : (x_i, x_j) \in \mathcal{D}$  if  $x_i$  and  $x_j$  are dissimilar

Distance metric learning:

Learn a distance metric of the form

$$d(x, y) = d_A(x, y) = \|x - y\|_A = \sqrt{(x - y)^T A (x - y)},$$

such that pairs of points  $(x_i, x_j)$  in  $\mathcal{S}$  have small squared distance.

- In general,  $A$  parameterizes a family of Mahalanobis distances over  $\mathbb{R}^n$ .
- Learning  $A$  is equivalent to finding a rescaling of a data:  $x \rightarrow A^{1/2}x$ .

© Eric Xing @ CMU, 2006-2011

55

## Optimal Distance Metric



- Learning an optimal distance metric with respect to the side-information leads to the following optimization problem:

$$\min_A \sum_{(x_i, x_j) \in \mathcal{S}} \|x_i - x_j\|_A^2 \quad (1)$$

$$\text{s.t. } \sum_{(x_i, x_j) \in \mathcal{D}} \|x_i - x_j\|_A \geq 1, \quad (2)$$

$$A \geq 0. \quad (3)$$

- This optimization problem is **convex**. Local-minima-free algorithms exist.
- Xing et al 2003 provided an efficient **gradient descent + iterative constraint-projection** method

© Eric Xing @ CMU, 2006-2011

56

## Examples of learned distance metrics



- Distance metrics learned on three-cluster artificial data:

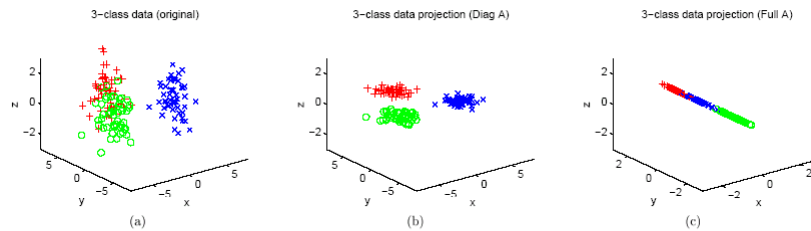


Figure 2: (a) Original data. (b) Rescaling corresponding to learned diagonal  $A$ . (c) Rescaling corresponding to full  $A$ .

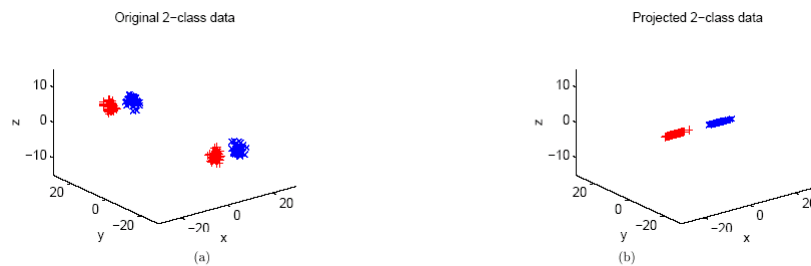
© Eric Xing @ CMU, 2006-2011

57

## Application to Clustering



- Artificial Data I: a difficult two-class dataset



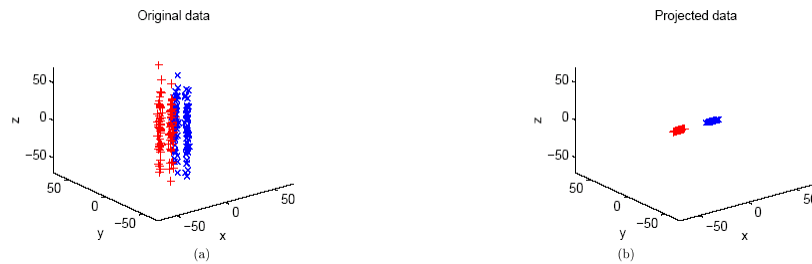
1. K-means: Accuracy = 0.4975
2. Constrained K-means: Accuracy = 0.5060
3. K-means + metric: Accuracy = 1
4. Constrained K-means + metric: Accuracy = 1

© Eric Xing @ CMU, 2006-2011

58

## Application to Clustering

- Artificial Data II: two-class data with strong irrelevant feature



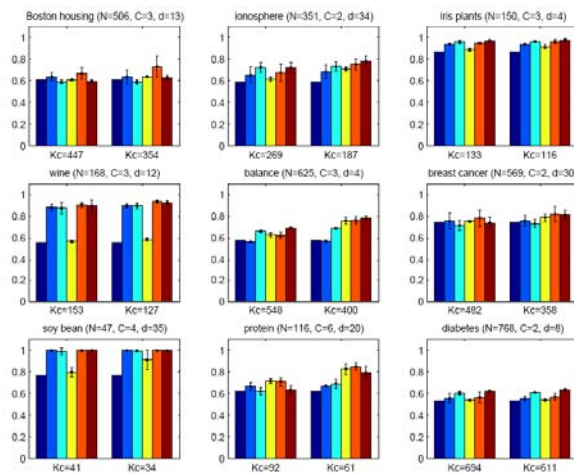
1. K-means: Accuracy = 0.4993
2. Constrained K-means: Accuracy = 0.5701
3. K-means + metric: Accuracy = 1
4. Constrained K-means + metric: Accuracy = 1

© Eric Xing @ CMU, 2006-2011

59

## Application to Clustering

- 9 datasets from the UC Irvine repository



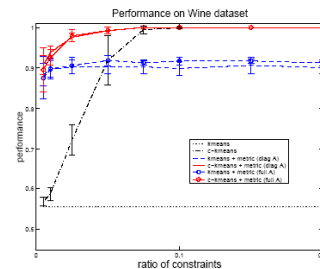
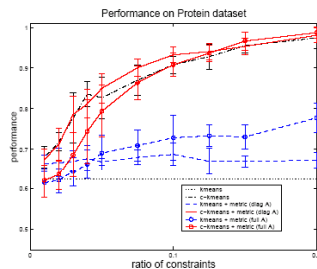
© Eric Xing @ CMU, 2006-2011

60

## Accuracy vs. amount of side-information



- Two typical examples of how the quality of the clusters found increases with the amount of side-information.



© Eric Xing @ CMU, 2006-2011

61

## Take home message



- Distance metric learning is an important problem in machine learning and data mining.
- A good distance metric can be learned from small amount of side-information in the form of similarity and dissimilarity constraints from data by solving a convex optimization problem.
- The learned distance metric can identify the most significant direction(s) in feature space that separates data well, effectively doing implicit Feature Selection.
- The learned distance metric can be used to improve clustering performance.

© Eric Xing @ CMU, 2006-2011

62