

CS 2750 Machine Learning

Lecture 22

Clustering

Milos Hauskrecht
milos@cs.pitt.edu
5329 Sennott Square

CS 2750 Machine Learning

Gaussian mixture model

Probability of occurrence of a data example x
is modeled as

$$p(\mathbf{x}) = \sum_{i=1}^k p(C = i) p(\mathbf{x} | C = i)$$

where

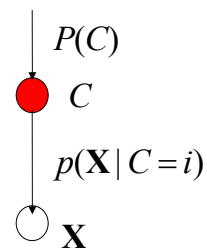
$$p(C = i)$$

= probability of a data point coming
from class $C=i$

$$p(\mathbf{x} | C = i) \approx N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

= class conditional density (modeled as a Gaussian)
for class i

Remember: C is hidden !!!!



CS 2750 Machine Learning

Generative classifier model

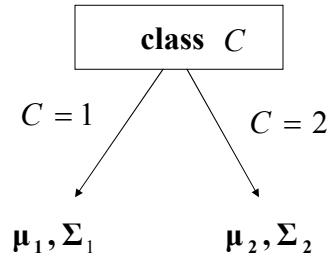
- Generative classifier model with Gaussian densities
- Assume the class labels are known. The ML estimate is

$$N_i = \sum_{j:C_j=i} 1$$

$$\tilde{\pi}_i = \frac{N_i}{N}$$

$$\tilde{\boldsymbol{\mu}}_i = \frac{1}{N_i} \sum_{j:C_j=i} \mathbf{x}_j$$

$$\tilde{\boldsymbol{\Sigma}}_i = \frac{1}{N_i} \sum_{j:C_j=i} (\mathbf{x}_j - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_i)^T$$



CS 2750 Machine Learning

Gaussian mixture model

- In the Gaussian mixture Gaussians are not labeled
- We can apply **EM algorithm**:
 - re-estimation based on the class posterior

$$h_{il} = p(C_l = i | \mathbf{x}_l, \Theta') = \frac{p(C_l = i | \Theta') p(\mathbf{x}_l | C_l = i, \Theta')}{\sum_{u=1}^m p(C_l = u | \Theta') p(\mathbf{x}_l | C_l = u, \Theta')}$$

$$N_i = \sum_l h_{il} \quad \leftarrow \text{Count replaced with the expected count}$$

$$\tilde{\pi}_i = \frac{N_i}{N}$$

$$\tilde{\boldsymbol{\mu}}_i = \frac{1}{N_i} \sum_l h_{il} \mathbf{x}_l$$

$$\tilde{\boldsymbol{\Sigma}}_i = \frac{1}{N_i} \sum_l h_{il} (\mathbf{x}_l - \boldsymbol{\mu}_i)(\mathbf{x}_l - \boldsymbol{\mu}_i)^T$$

CS 2750 Machine Learning

Gaussian mixture algorithm

- **Special case:** fixed covariance matrix for all hidden groups (classes) and uniform prior on classes

- **Algorithm:**

Initialize means μ_i for all classes i

Repeat two steps until no change in the means:

1. Compute the class posterior for each Gaussian and each point (a kind of responsibility for a Gaussian for a point)

Responsibility:
$$h_{il} = \frac{p(C_l = i | \Theta') p(x_l | C_l = i, \Theta')}{\sum_{u=1}^m p(C_l = u | \Theta') p(x_l | C_l = u, \Theta')}$$

2. Move the means of the Gaussians to the center of the data, weighted by the responsibilities

New mean:
$$\mu_i = \frac{\sum_{l=1}^N h_{il} \mathbf{x}_l}{\sum_{l=1}^N h_{il}}$$

K-means approximation to EM

Expectation-Maximization:

- posterior measures the responsibility of a Gaussian for every point

$$h_{il} = \frac{p(C_l = i | \Theta') p(x_l | C_l = i, \Theta')}{\sum_{u=1}^m p(C_l = u | \Theta') p(x_l | C_l = u, \Theta')}$$

K- Means

- Only the closest Gaussian is made responsible for a point

$$h_{il} = 1 \quad \text{If } i \text{ is the closest Gaussian}$$

$$h_{il} = 0 \quad \text{Otherwise}$$

Re-estimation of means

$$\mu_i = \frac{\sum_{l=1}^N h_{il} \mathbf{x}_l}{\sum_{l=1}^N h_{il}}$$

- Results in moving the means of Gaussians to the center of the data points it covered in the previous step

K-means algorithm

K-Means algorithm:

Initialize k values of means (centers)

Repeat two steps until no change in the means:

- Partition the data according to the current means (using the similarity measure)
- Move the means to the center of the data in the current partition

- **Used frequently for clustering data**

Clustering

Groups together “similar” instances in the data sample

Basic clustering problem:

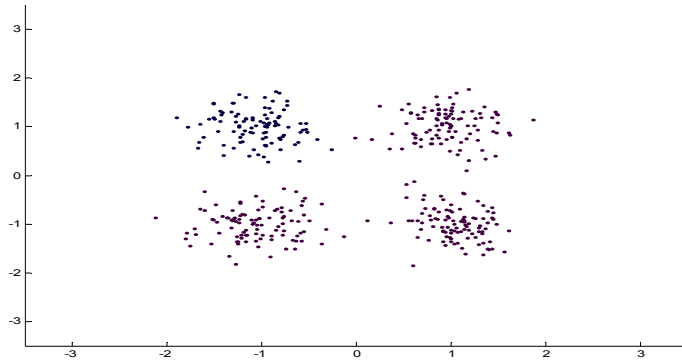
- distribute data into k different groups such that data points similar to each other are in the same group
- Similarity between data points is defined in terms of some distance metric (can be chosen)

Clustering is useful for:

- **Similarity/Dissimilarity analysis**
Analyze what data points in the sample are close to each other
- **Dimensionality reduction**
High dimensional data replaced with a group (cluster) label

Clustering example

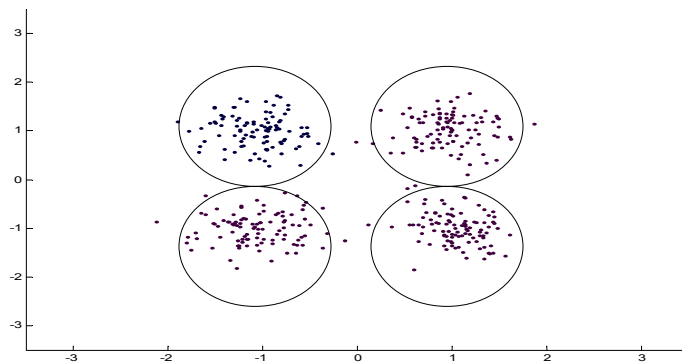
- We see data points and want to partition them into groups
- Which data points belong together?



CS 2750 Machine Learning

Clustering example

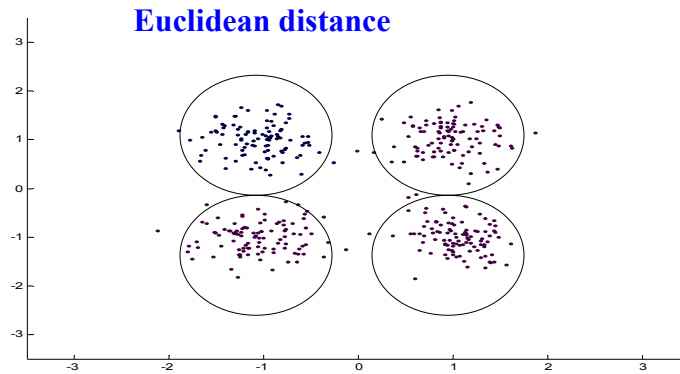
- We see data points and want to partition them into the groups
- Which data points belong together?



CS 2750 Machine Learning

Clustering example

- We see data points and want to partition them into the groups
- Requires a distance measure to tell us what points are close to each other and are in the same group



CS 2750 Machine Learning

Clustering example

- A set of patient cases
- We want to partition them into the groups based on similarities

Patient #	Age	Sex	Heart Rate	Blood pressure ...
Patient 1	55	M	85	125/80
Patient 2	62	M	87	130/85
Patient 3	67	F	80	126/86
Patient 4	65	F	90	130/90
Patient 5	70	M	84	135/85

CS 2750 Machine Learning

Clustering example

- A set of patient cases
- We want to partition them into the groups based on similarities

Patient #	Age	Sex	Heart Rate	Blood pressure ...
Patient 1	55	M	85	125/80
Patient 2	62	M	87	130/85
Patient 3	67	F	80	126/86
Patient 4	65	F	90	130/90
Patient 5	70	M	84	135/85

How to design the distance metric to quantify similarities?

CS 2750 Machine Learning

Clustering example. Distance measures.

In general, one can choose an arbitrary distance measure.

Properties of the distance measures:

Assume 2 data entries a, b

Positiveness: $d(a, b) \geq 0$

Symmetry: $d(a, b) = d(b, a)$

Identity: $d(a, a) = 0$

CS 2750 Machine Learning

Distance measures.

Assume pure real-valued data-points:

12	34.5	78.5	89.2	19.2
23.5	41.4	66.3	78.8	8.9
33.6	36.7	78.3	90.3	21.4
17.2	30.1	71.6	88.5	12.5

What distance metric to use?

Distance measures.

Assume pure real-valued data-points:

12	34.5	78.5	89.2	19.2
23.5	41.4	66.3	78.8	8.9
33.6	36.7	78.3	90.3	21.4
17.2	30.1	71.6	88.5	12.5

What distance metric to use?

Euclidian: works for an arbitrary k-dimensional space

$$d(a, b) = \sqrt{\sum_{i=1}^k (a_i - b_i)^2}$$

Distance measures.

Assume pure real-valued data-points:

12	34.5	78.5	89.2	19.2
23.5	41.4	66.3	78.8	8.9
33.6	36.7	78.3	90.3	21.4
17.2	30.1	71.6	88.5	12.5

What distance metric to use?

Squared Euclidian: works for an arbitrary k-dimensional space

$$d^2(a, b) = \sum_{i=1}^k (a_i - b_i)^2$$

Distance measures.

Assume pure real-valued data-points:

12	34.5	78.5	89.2	19.2
23.5	41.4	66.3	78.8	8.9
33.6	36.7	78.3	90.3	21.4
17.2	30.1	71.6	88.5	12.5

Assume that two variables are highly correlated k-dimensional space

$$d(a, b) = \sum_{i=1}^k |a_i - b_i|$$

Etc. ...

Distance measures.

Generalized distance metric:

$$d^2(\mathbf{a}, \mathbf{b}) = (\mathbf{a} - \mathbf{b})\Gamma^{-1}(\mathbf{a} - \mathbf{b})^T$$

Γ^{-1} is a matrix that weights attributes proportionally to their importance. Different weights lead to a different distance metric.

If $\Gamma = I$ we get squared Euclidean

$\Gamma = \Sigma$ Mahalanobis distance takes into account correlations among attributes

Distance measures.

Assume pure binary values data:

0 1 1 0 1

1 0 1 0 1

0 1 1 0 1

1 1 1 1 1

...

What distance metric to use?

Distance measures.

Assume pure binary values data:

```
0 1 1 0 1
1 0 1 0 1
0 1 1 0 1
1 1 1 1 1
...
```

What distance metric to use?

Edit distance: The number of attributes that need to be changed to make the entries the same

The same metric can be used for pure categorical values

Distance measures.

Combination of real-valued and categorical attributes

Patient #	Age	Sex	Heart Rate	Blood pressure ...
Patient 1	55	M	85	125/80
Patient 2	62	M	87	130/85
Patient 3	67	F	80	126/86
Patient 4	65	F	90	130/90
Patient 5	70	M	84	135/85

What distance metric to use?

Distance measures.

Combination of real-valued and categorical attributes

Patient #	Age	Sex	Heart Rate	Blood pressure ...
Patient 1	55	M	85	125/80
Patient 2	62	M	87	130/85
Patient 3	67	F	80	126/86
Patient 4	65	F	90	130/90
Patient 5	70	M	84	135/85

What distance metric to use?

A weighted sum approach: e.g. a mix of Euclidian and Edit distances for subsets of attributes

Clustering

Clustering is useful for:

- **Similarity/Dissimilarity analysis**
Analyze what data points in the sample are close to each other
- **Dimensionality reduction**
High dimensional data replaced with a group (cluster) label

Problems:

- Pick the correct similarity measure (problem specific)
- Choose the correct number of groups
 - Many clustering algorithms require us to provide the number of groups ahead of time

Clustering algorithms

Partitioning algorithms:

- **K-means algorithm**

- **suitable** only when data points have continuous values; groups are defined in terms of cluster centers (also called **means**).
- refinement of the method to categorical values: **K-medoids**

- **Probabilistic methods (with EM)**

- **Latent variable models**: class (cluster) is represented by a latent (hidden) variable value.
- **Examples**: mixture of Gaussians, Naïve Bayes with a hidden class

- **Hierarchical methods**

- **Agglomerative**
- **Divisive**

K-means

K-Means algorithm:

Initialize randomly k values of means (centers)

Repeat two steps until no change in the means:

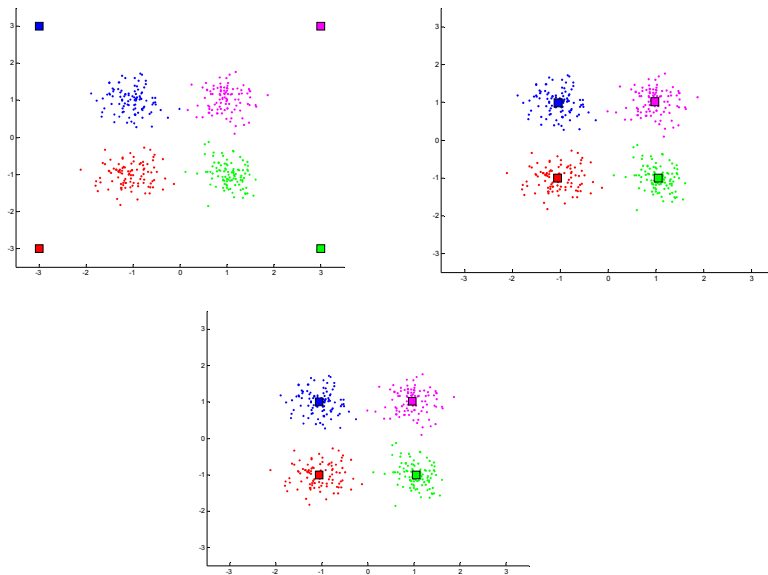
- Partition the data according to the current set of means (using the similarity measure)
- Move the means to the center of the data in the current partition

Stop when no change in the means

Properties:

- Minimizes the sum of **squared center-point distances** for all clusters
- The algorithm always converges (local optima).

K-Means example



CS 2750 Machine Learning

K-means algorithm

- **Properties:**
 - converges to centers minimizing the sum of squared center-point distances (still local optima)
 - The result is sensitive to the initial means' values
- **Advantages:**
 - Simplicity
 - Generality – can work for more than one distance measure
- **Drawbacks:**
 - Can perform poorly with overlapping regions
 - Lack of robustness to outliers
 - Good for attributes (features) with continuous values
 - Allows us to compute cluster means
 - k-medoid algorithm used for discrete data

CS 2750 Machine Learning

Probabilistic (EM-based) algorithms

- **Latent variable models**

**Examples: Naïve Bayes with hidden class
Mixture of Gaussians**

- **Partitioning:**

- the data point belongs to the class with the highest posterior

- **Advantages:**

- Good performance on overlapping regions
- Robustness to outliers
- Data attributes can have different types of values

- **Drawbacks:**

- EM is computationally expensive and can take time to converge
- Density model should be given in advance

Hierarchical clustering.

Uses an arbitrary similarity/dissimilarity measure.

Typical similarity measures $d(a,b)$:

Pure real-valued data-points:

- Euclidean, Manhattan, Minkowski distances

Pure binary values data:

- Number of matching values

Pure categorical data:

- Number of matching values

Combination of real-valued and categorical attributes

- Weighted approaches

Hierarchical clustering.

Approach:

- **Compute dissimilarity matrix for all pairs of points**
 - uses standard or other distance measures
- **Construct clusters greedily:**
 - **Agglomerative approach**
 - Merge pair of clusters in a bottom-up fashion, starting from singleton clusters
 - **Divisive approach:**
 - Splits clusters in top-down fashion, starting from one complete cluster
- **Stop the greedy construction** when some criterion is satisfied
 - E.g. fixed number of clusters

CS 2750 Machine Learning

Cluster merging

- **Construction of clusters through greedy agglomerative approach**
 - Merge pair of clusters in a bottom-up fashion, starting from singleton clusters
 - Merge clusters based on cluster (or linkage) distances.
Defined in terms of point distances. **Examples:**

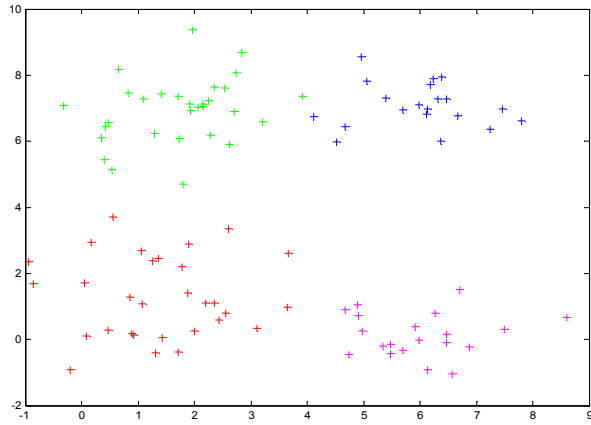
Min distance $d_{\min}(C_i, C_j) = \min_{p \in C_i, q \in C_j} |p - q|$

Max distance $d_{\max}(C_i, C_j) = \max_{p \in C_i, q \in C_j} |p - q|$

Mean distance $d_{\text{mean}}(C_i, C_j) = \left| \frac{1}{|C_i|} \sum_i p_i - \frac{1}{|C_j|} \sum_j q_j \right|$

CS 2750 Machine Learning

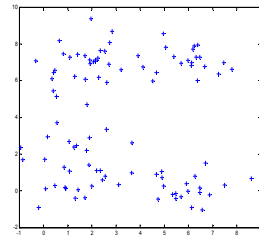
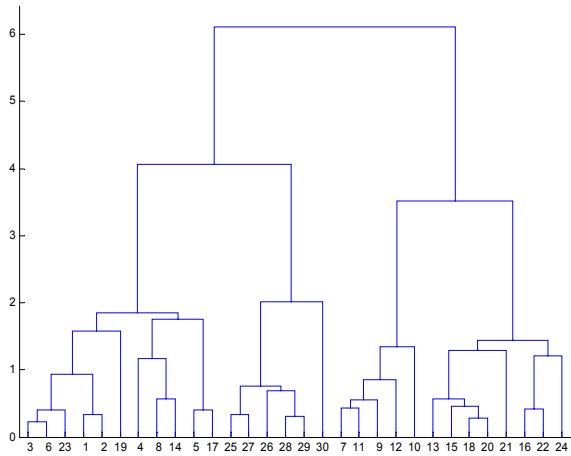
Hierarchical clustering example



CS 2750 Machine Learning

Hierarchical clustering example

- dendrogram



CS 2750 Machine Learning

Hierarchical clustering

- **Advantage:**

- Smaller computational cost; avoids scanning all possible clusterings

- **Disadvantage:**

- Greedy choice fixes the order in which clusters are merged; cannot be repaired

- **Partial solution:**

- combine hierarchical clustering with iterative algorithms like k-means

Other clustering methods

- **Spectral clustering**

- Uses similarity matrix

- **Multidimensional scaling**

- techniques often used in data visualization for exploring similarities or dissimilarities in data.