

Artificial Intelligence III:
Artificial Intelligence and Deep Learning

Lecture 6

# Unsupervised Learning

Dr. Patrick Chan
patrickchan@ieee.org
South China University of Technology, China



# Agenda

- Introduction
- Clustering
  - Similarity Measure
  - Criterion Function
  - Algorithm
- Feature Extraction
  - Principal Components Analysis

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture



## Supervised VS Unsupervised

- Supervised Learning
  - Label is given
  - Someone (a supervisor) provides the true answer



- No Label is given "Learning without a teacher"
- Much harder than supervised learning
- You never know the correct answer
  - How to evaluate the result?



## **Unsupervised Learning**

- No Supervision (No Label) How to evaluate the result?
  - External: Expert comments
    - Expert may be wrong
  - Internal: Objective functions
    - E.g. Distance between samples and centers
    - Very intuitive
- Different from Supervised Learning, the evaluation method is subjective



## Why No Label?

- Label is expensive
  - Especially for a huge dataset
  - E.g. Medical application
- Sometimes, the objective is not clear
  - Data Mining
- Gain some insight about the data structure before designing classifiers
  - E.g. Feature selection

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6



## **Unsupervised Learning Type**

- Parametric Approach
  - Assume distribution is known
  - Estimate parameters of distribution
    - E.g. Maximum-Likelihood Estimate
- Non-Parametric Approach
  - No assumption on the distribution
  - Group data into clusters
    - Samples in the same group share something in common

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture



## Clustering

How many clusters (groups) are there?



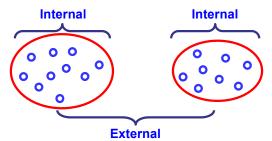


How can you know it?



## Clustering

How many clusters (groups) are there?

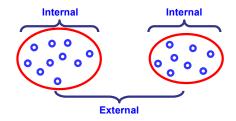


- Assumptions
  - Internal Characteristic (intra-cluster):
    - Distance within a cluster should be small
  - External Characteristic (inter-cluster):
    - Distance between clusters should be large



### Clustering Three Important Factors

- Distance (Similarity) Measure
  - How similar between two samples?
- Criterion Function
  - What kind of clustering result is expected?
- Clustering Algorithm
  - E.g. optimize the criterion function



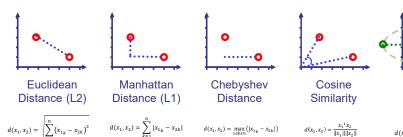
Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6



### Clustering **Similarity Measure**

- No best measure for all cases
- Application dependent



Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

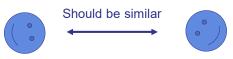
Mahalanobis

Distance



### Clustering Similarity Measure

- No best measure for all cases
- Application dependent
  - Examples:
    - Rotation Invariance in Face Recognition



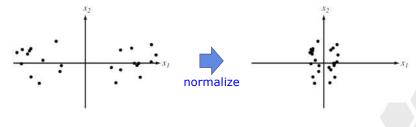
NO Rotation Invariance in Character Recognition

Should be different



### Clustering Similarity Measure

- Scale of features may be different
  - Different Ranges (Weight: 80 300, waist width: 28 45)
  - Different Units (Km VS mile, cm VS meter)
- May be solved by normalized, e.g. [0, 1]
  - Sometimes may not be suitable
  - normalization reduces cluster effect (right diagram)





## Naïve Clustering Algorithm

- A naïve clustering algorithm can be developed only based on similarity measure between samples
- Algorithm:
  - Calculate similarity for each sample pair
  - Group the samples in the same cluster if the measure between them is less than a threshold (d<sub>0</sub>)

Dr. Patrick Chan @ SCUT

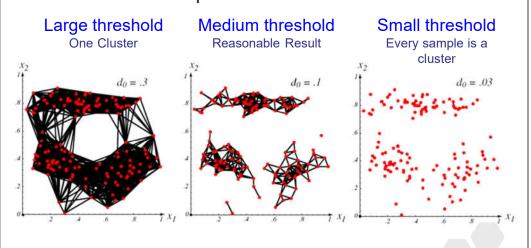
Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

### Clustering Naïve Clustering Algorithm Large threshold 1.0 2.0 3.2 3.0 $(d_0 = 4.0)$ 0 1.0 2.2 3.2 1 cluster 1.0 0 1.4 3.6 2.2 1.4 0 3.6 Medium threshold $(d_0 = 2.5)$ 2 clusters Small threshold $(d_0 = 0.5)$ 5 clusters Dr. Patrick Chan @ SCUT Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6



### Clustering Naïve Clustering Algorithm

Another Example



### Clustering Naïve Clustering Algorithm

- Advantage:
  - Easy to understand
  - Simple to implement
- Disadvantage:
  - Only local information is considered
  - Highly dependent on the threshold



### Clustering

## **Criterion Function**

- Commonly used criteria
  - Intra-cluster scatter (Variance of each cluster) Smaller is better

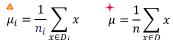
$$S_A = \sum_{i=1}^{c} \sum_{x \in D_i} (x - \mu_i)(x - \mu_i)^t$$

Inter-cluster scatter (Distance between clusters) Larger is better

$$S_E = \sum_{i=1}^{c} n_i (\mu_i - \mu) (\mu_i - \mu)^t$$

Combination

$$S = \alpha S_A - (1 - \alpha) S_E$$



c: the number of cluster

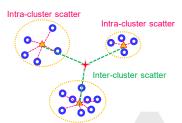
the number of samples

 $n_i$ : the number of samples in cluster i

D: the set of all samples

 $D_i$ : the set of samples in cluster i

 $\alpha$ : the tradeoff



### Dr. Patrick Chan @ SCUT

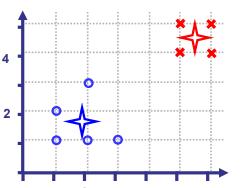
Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

### Clustering

## **Criterion Function: Example**

Intra-cluster scatter

$$S_A = \sum_{i=1}^{c} \sum_{x \in D_i} (x - \mu_i)(x - \mu_i)^t$$



Dr. Patrick Chan @ SCUT

$$\mu_1 = \frac{1}{4} \begin{pmatrix} [5 & 4] + [5 & 5] \\ +[6 & 4] + [6 & 5] \end{pmatrix} = [5.5 & 4.5]$$

$$\mu_2 = \frac{1}{5} \begin{pmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 2 \end{bmatrix} \\ + \begin{bmatrix} 2 & 1 \end{bmatrix} + \begin{bmatrix} 2 & 3 \end{bmatrix} \end{pmatrix} = \begin{bmatrix} 1.8 & 1.6 \end{bmatrix}$$

$$= \begin{pmatrix} \sqrt{(5-5.5)^2 + (4-4.5)^2} + \sqrt{(5-5.5)^2 + (5-4.5)^2} \\ + \sqrt{(6-5.5)^2 + (4-4.5)^2} + \sqrt{(6-5.5)^2 + (5-4.5)^2} \end{pmatrix} \\ + \begin{pmatrix} \sqrt{(1-1.8)^2 + (1-1.6)^2} + \sqrt{(1-1.8)^2 + (2-1.6)^2} \\ + \sqrt{(2-1.8)^2 + (1-1.6)^2} + \sqrt{(2-1.8)^2 + (3-1.6)^2} \end{pmatrix} \\ + \sqrt{(3-1.8)^2 + (1-1.6)^2} \end{pmatrix}$$

$$= 2.83 + 5.28 = 8.11$$

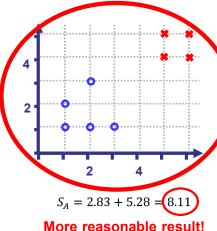
Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

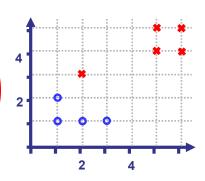


### Clustering

## **Criterion Function: Example**

 $\bullet$  Smaller  $S_A$  is preferred





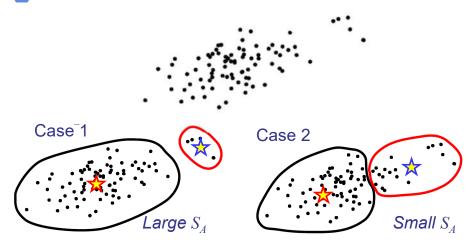
$$S_4 = 6.81 + 3.48 = 10.29$$

### Clustering **Criterion Function**

- Is  $S_A$  (Intra-cluster scatter) a good criterion for all situations?
- How to separate the following samples into two clusters?



# Criterion Function



Case 1 is more reasonable

However, it has a larger value of  $S_A$  due to the large cluster

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

# Criterion Function

- $\bullet$   $S_A$  (Intra-cluster scatter) is
- Appropriate:
  - The clusters form compact groups
  - Equally sized clusters
- Not Appropriate
  - When natural groupings have very different sizes

2 Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture



# Clustering Algorithm

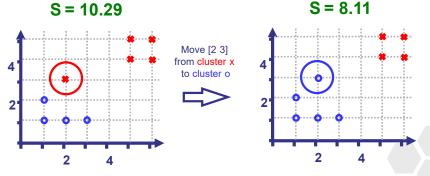
- Find the optimal clustering result
- Exhaustive search is impossible
  - ~C<sup>n</sup> possible partitions
  - C: class number, n: sample number
- Methods:
  - Iterative Optimization Algorithm
    - K-means
  - Hierarchical Clustering
    - Bottom Up Approach
    - Top Down Approach



# Iterative Optimization Algorithm

Euclidean Distance is used

- 1. Find a reasonable initial cluster result
- 2. Move sample(s) from one cluster to another such that the objective function is improved the most
- 3. Goto 2 until stable





### **Clustering: Iterative Optimization Algorithm**

### K-means

- ◆ A well-known technique: K-means
  - Assume there are k clusters
  - Minimize Criterion Function (Intra-class scatter):

$$S = \sum_{i=1}^{k} \sum_{x \in D_i} \|x - \mu_i\|^2 \qquad \mu_i = \frac{1}{n_i} \sum_{x \in D_i} x$$

$$k: \quad \text{the number of cluster}$$

$$\mu_i = \frac{1}{n_i} \sum_{x \in D_i} x$$

 $n_i$ : the number of samples in cluster i $D_i$ : the set of samples in cluster i

 In each iteration, assign a sample to its closest cluster

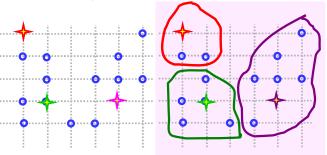
Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

### **Clustering: Iterative Optimization Algorithm**

### K-means

◆ Example: *k*=3



### 1. Initialization

Randomly assign the center of each cluster

### 2. Assign Samples

Assign samples to closest center

3. Re-calculate mean

Compute the new means using new samples

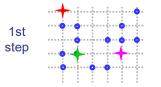
Repeat until stable (no sample moves again)

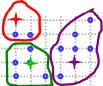
Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

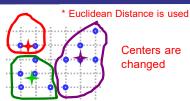


### Clustering: Iterative Optimization Algorithm

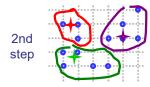
### K-means





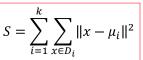


Centers are changed

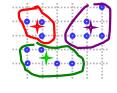


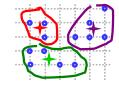


Centers are









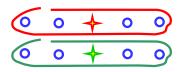
Centers are not changed K-means Stops

### **Clustering: Iterative Optimization Algorithm**

### K-means

Dr. Patrick Chan @ SCUT

- Pros:
  - Optimize the objective function efficiently
  - Algorithm converges
- Cons:
  - May be trapped at local minimum (similar to gradient descent)







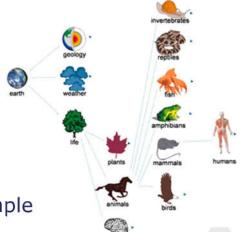
Trapped at Local Minimum

Global Minimum



### Clustering **Hierarchical Clustering**

- Sometimes, clusters have subclusters, and so on
  - A cluster can further be broken down into smaller clusters
- Hierarchical cluster
- Taxonomy is an example



Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6



### Clustering **Hierarchical Clustering**

- Two types:
  - Top Down Approach
    - Start with 1 cluster
      - One cluster contains all samples
    - Form hierarchy by splitting the most dissimilar clusters

### Bottom Up Approach

- Start with n clusters
  - Each cluster contains one sample
- Form hierarchy by merging the most similar clusters
- Not efficient if a large number of samples but a number of clusters is needed

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6



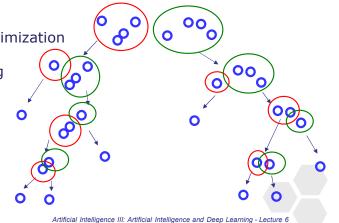
### Clustering: Hierarchical Clustering

## Top Down Approach

Start from one cluster

 Break down a cluster with more than one sample into two

 Any Iterative Optimization Algorithm can be applied by setting c = 2





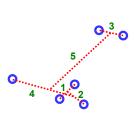
## **Clustering: Hierarchical Clustering**

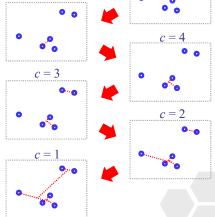
**Bottom Up Approach** 

\* Fuclidean Distance is used

Initially each sample forms a cluster

 Merge the nearest two clusters until one cluster left







### **Clustering: Hierarchical Clustering**

## **Bottom Up Approach**

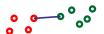
How to calculate distance between clusters?





Mean Distance

$$d_{mean}(D_i, D_j) = ||m_i - m_i||$$



Minimum Distance

$$d_{min}(D_i, D_j) = \min_{x \in D_i, z \in D_j} ||x - z||$$



Maximum Distance

$$d_{max}(D_i, D_j) = \max_{x \in D_i, z \in D_j} ||x - z||$$



Average Distance

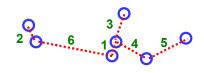
$$d_{avg}(D_i, D_j) = \frac{1}{n_i n_j} \sum_{x \in D_i} \sum_{x \in D_j} ||x - z||$$

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

**Clustering: Hierarchical Clustering Bottom Up Approach** 

- Single Linkage (Nearest-Neighbor)
  - Minimum Distance is used
  - Encourage growth of elongated clusters
- Complete Linkage (Farthest Neighbor)
  - Maximum Distance is used
  - Encourages compact clusters



Single Linkage

Min distance between points of each cluster

Complete Linkage

Max distance between points of each cluster

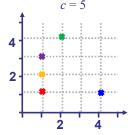
Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6



Clustering: Hierarchical Clustering: Bottom Up Approach

## Single Linkage: Example 1/4



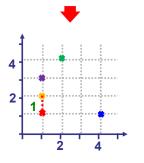


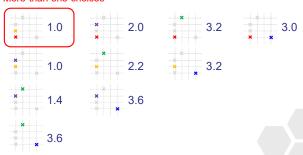


$$d((1,1),(2,4)) = \sqrt{(1-2)^2 + (1-4)^2}$$
  
= 3.2

Euclidean Distance is used

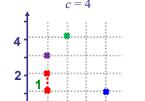
$$min(3.2) = 3.2$$



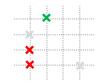


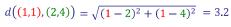
Clustering: Hierarchical Clustering: Bottom Up Approach

## Single Linkage: Example 2/4





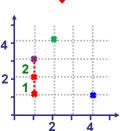




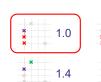
$$d((1,2),(2,4)) = \sqrt{(1-2)^2 + (2-4)^2} = 2.2$$



Euclidean Distance is used



Dr. Patrick Chan @ SCUT







Clustering: Hierarchical Clustering: Bottom Up Approach

## Single Linkage: Example 3/4



Example:

Euclidean Distance is used



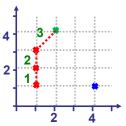
$$d((1,1),(2,4)) = \sqrt{(1-2)^2 + (1-4)^2} = 3.2$$

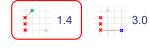
$$d((1,2),(2,4)) = \sqrt{(1-2)^2 + (2-4)^2} = 2.2$$

$$d((1,3),(2,4)) = \sqrt{(1-2)^2 + (3-4)^2} = 1.4$$

$$\min(3.2, 2.2, 1.4) = 1.4$$







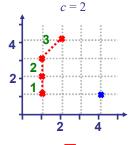


Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

Clustering: Hierarchical Clustering: Bottom Up Approach

## Single Linkage: Example 4/4



Example:

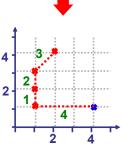
$$d((1,1),(4,1)) = \sqrt{(1-4)^2 + (1-1)^2} = 3.0$$

$$d((1,2),(4,1)) = \sqrt{(1-4)^2 + (2-1)^2} = 3.2$$

$$d((1,3),(4,1)) = \sqrt{(1-4)^2 + (3-1)^2} = 3.6$$

$$d((2,4),(4,1)) = \sqrt{(2-4)^2 + (4-1)^2} = 3.6$$

min(3.2, 2.2, 1.4) = 3.0



3.0

\* This step is unnecessary as only one candidate

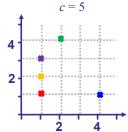
Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6



Clustering: Hierarchical Clustering: Bottom Up Approach

## Complete Linkage: Example 1/4



Example:



More than one choices

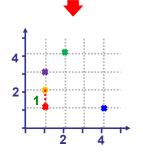
$$d((1,1),(2,4)) = \sqrt{(1-2)^2 + (1-4)^2}$$
  
= 3.2

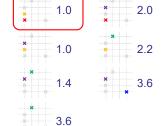
$$max(3.2) = 3.2$$

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

\* This step is the same as Single Linkage since distance measure of clusters with one point

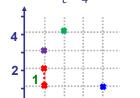
\* Euclidean Distance is used





Clustering: Hierarchical Clustering: Bottom Up Approach

## Complete Linkage: Example 2/4



Example:

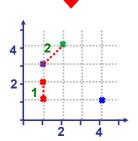


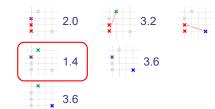
 $d((1,1),(2,4)) = \sqrt{(1-2)^2 + (1-4)^2} = 3.2$ 

$$d((1,2),(2,4)) = \sqrt{(1-2)^2 + (2-4)^2} = 2.2$$

max(3.2, 2.2) = 3.2

\* Euclidean Distance is used





Dr. Patrick Chan @ SCUT

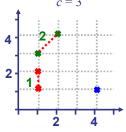




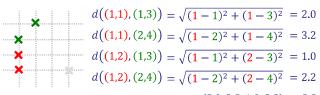
Clustering: Hierarchical Clustering: Bottom Up Approach

## Complete Linkage: Example 3/4

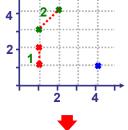




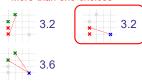
### Example:











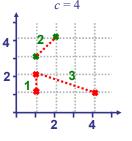
Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6



### Clustering: Hierarchical Clustering: Bottom Up Approach

## Complete Linkage: Example 4/4







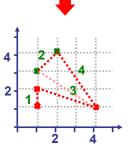
$$d((1,1),(1,3)) = \sqrt{(1-1)^2 + (1-3)^2} = 2.0$$
  
$$d((1,1),(2,4)) = \sqrt{(1-2)^2 + (1-4)^2} = 3.2$$

$$d((1,2),(1,3)) = \sqrt{(1-1)^2 + (2-3)^2} = 1.0$$

$$d((1,2),(2,4)) = \sqrt{(1-2)^2 + (2-4)^2} = 2.2$$

$$d((4,1),(1,3)) = \sqrt{(4-1)^2 + (1-3)^2} = 3.6$$
$$d((4,1),(2,4)) = \sqrt{(4-2)^2 + (1-4)^2} = 3.6$$

$$\max(2.0, 3.2, 1.0, 2.2, 3.6, 3.6) = 3.6$$



3.6

\* This step is unnecessary as only one candidate

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

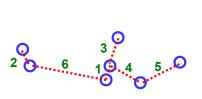


### Clustering: Hierarchical Clustering

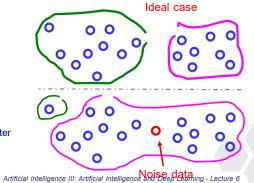
## **Bottom Up Approach**

### Single Linkage (Nearest-Neighbor)

- Minimum Distance is used
- Encourage growth of elongated clusters
- Disadvantage: Sensitive to noise



Min distance between points of each cluster



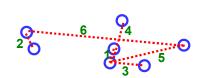
## Clustering: Hierarchical Clustering

## **Bottom Up Approach**

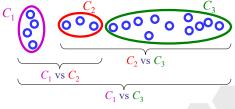
### Complete Linkage (Farthest Neighbor)

- Maximum Distance is used
- Encourages compact clusters
- Disadvantage: Does not work well if elongated clusters present

Ideally, C2 and C3 should be merged



Max distance between points of each cluster



However, C<sub>1</sub> and C<sub>2</sub> will be merged



### **Clustering: Hierarchical Clustering**

## **Bottom Up Approach**

- Minimum and maximum distance are noise sensitive (especially, minimum)
- More robust result to outlier when average or mean are used

$$d_{avg}(D_i, D_j) = \frac{1}{n_i n_j} \sum_{x \in D_i} \sum_{x \in D_j} ||x - z||$$

$$d_{mean}(D_i, D_i) = ||m_i - m_i||$$

 Mean is less time consumed than Average distance

Dr. Patrick Chan @ SCUT

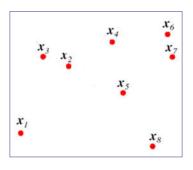
Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

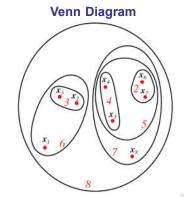


### **Clustering: Hierarchical Clustering** Venn

- Venn diagram can show hierarchical clustering
- No quantitative information is provided

### Sample points





Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

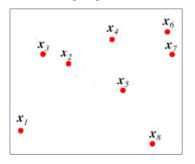


### **Clustering: Hierarchical Clustering**

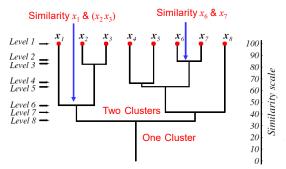
## Dendogram

- Dendogram is another way to represent a hierarchical clustering
- Able to indicate the similarity value

### Sample points



### Dendogram

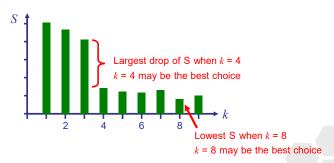


### Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture 6

### Clustering

## **Number of Clusters**

- How to decide the number of clusters?
- Possible solution:
  - Try a range of k and see which one has the lowest or largest drop criterion value (S)
  - Example:



Dr. Patrick Chan @ SCUT



## **Curse of Dimensionality**

- Real data usually have plenty of features
  - E.g., documents, images...
- Huge number of features causes problems
  - Sparsity
  - Complexity (storage and process)



Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture



## **Curse of Dimensionality**

Can the data be described with fewer dimensions, without losing much of its original meaning?

### Dimensionality Reduction

- Not just reduce the amount of data
- Often brings out the useful part of the data

D

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture



## Feature Reduction

For unsupervised learning, which feature is more useful to represent a dataset?

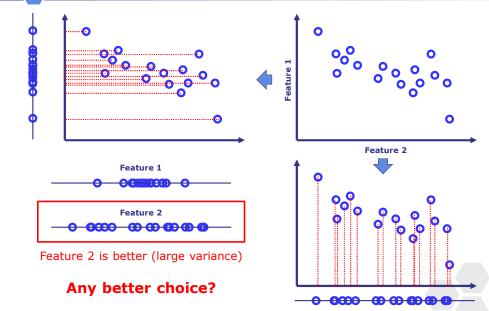
Feature A Values are simil

Feature B O O O O O O Values are very different

- A feature with different values provides more information
  - Variance is one of measures

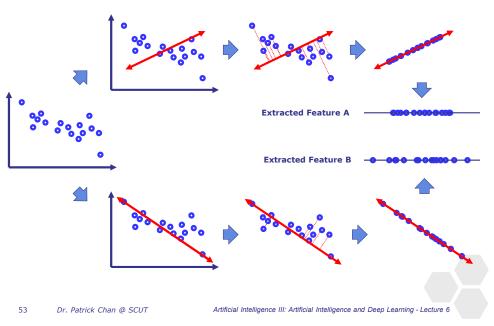


## **Feature Reduction**



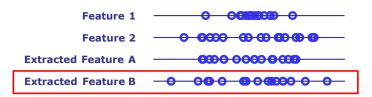


## **Feature Reduction**

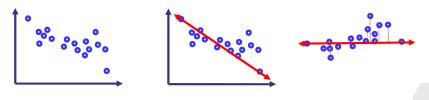




## **Feature Reduction**



 Extracted Feature B is the best for representing the data



Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture



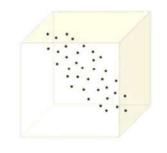
## **Principal Components Analysis**

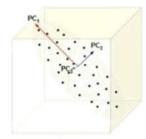
- PCA reduces data by geometrically projecting them onto lower dimensions called principal components (PCs)
- Project a dataset into a new set of features such that:
  - The features have zero covariance to each other (they are orthogonal)
  - Each feature captures the most remaining variance in the data, while orthogonal to the existing feature

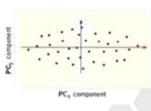


## **Principal Components Analysis**

- First principal component (PC) is the direction of greatest variability (variance) in the data
- Second PC is the next orthogonal (uncorrelated) direction of greatest variability
- And so on...



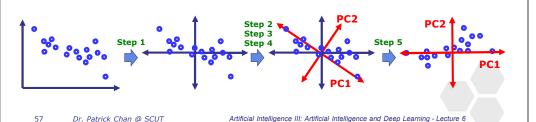






## **Principal Components Analysis**

- **Standardize** the dataset
- Calculate the **covariance matrix** for the features in the dataset
- 3. Calculate the **eigenvalues** and **eigenvectors** for the covariance matrix.
- 4. Pick top k eigenvalues and form their eigenvectors
- Transform the original matrix



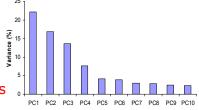


### **Principal Components Analysis**

## **PC** number Determination

 A feature with small eigenvalue contains small information

- n dimensions in original data
- Choose only the first p eigenvectors, based on their eigenvalues (p < n)
- Final data has only p dimensions



How to determine k?

Dr. Patrick Chan @ SCUT

Artificial Intelligence III: Artificial Intelligence and Deep Learning - Lecture



### **Principal Components Analysis**

## PC number Determination

Determine by projection error

$$\frac{\frac{1}{m}\sum_{i=1}^{n} \left(x^{(i)} - z_k(x^{(i)})\right)^2}{\frac{1}{m}\sum_{i=1}^{n} (x^{(i)})^2} \le \varepsilon$$

Determine by variation ratio

$$\frac{\sum_{j=1}^{k} \sigma_j^2}{\sum_{j=1}^{n} \sigma_j^2} \approx r$$

 $\sigma_i$ : the j<sup>th</sup> variance in descending order

r: expected ratio (e.g. 85%)



### **Principal Components Analysis**

## Limitation

Orthogonal Feature





Non-linear projection





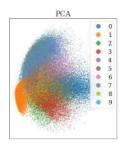


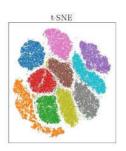




## **Principal Components Analysis**

Other feature extraction methods











## References

http://users.umiacs.umd.edu/~jbg/teachin g/INST\_414/

