

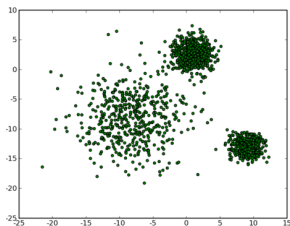
K -means 聚类算法

王石平 博士

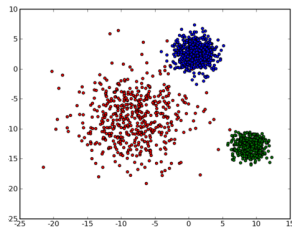
Outline

- Motivation
- K -means algorithm
- Face image clustering

An example

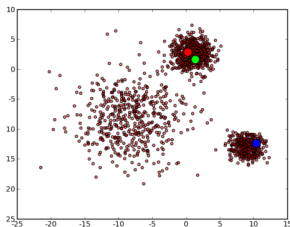


(a) Original data

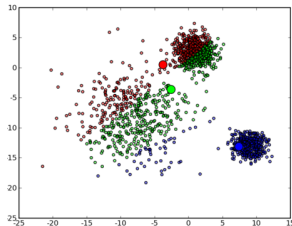


(b) Semantic clustering

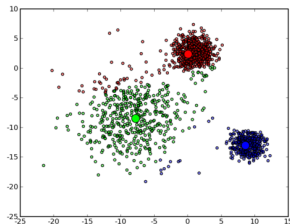
Iteration process



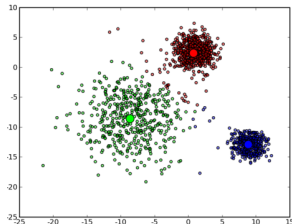
(c) 1th-iteration



(d) 2th-iteration



(e) 3th-iteration



(f) 4th-iteration

Algorithm 1 *K*-means algorithm

Input: data samples $X = \{x_1, \dots, x_n\}$ and the number of clusters k .

Output: Clustering results of all samples.

- 1: Randomly generate k cluster centroids μ_1, \dots, μ_k ;
 - 2: **repeat**
 - 3: Assign each sample to the cluster having smallest distance to the corresponding cluster centroid;
 - 4: Update cluster centroids with the means of clusters;
 - 5: **until** convergence
 - 6: **return** Clustering result.
-

Algorithm 2 K -means algorithm

Input: data samples $X = \{x_1, \dots, x_n\}$ where $x_i \in \mathbb{R}^d$ and the number of clusters k .

Output: Clustering results of all samples.

1: Initialize k cluster centroids μ_1, \dots, μ_k where $\mu_i \in \mathbb{R}^d$;

2: Initialize clustering indicator array $cluster \in \mathbb{R}^n$;

3: **repeat**

4: Set $cluster_i = \arg \min_j ||x_i - \mu_j||$;

5: Update $\mu_j = \frac{\sum_{i=1}^n 1(cluster_i=j)x_i}{\sum_{i=1}^n 1(cluster_i=j)}$;

6: **until** convergence

7: **return** $cluster$.

Question

- Objective function?
- Initialization sensitivity?
- Convergence?
- Performance evaluation?

Objective function of K -means clustering

$$J = \sum_{i=1}^n \sum_{j=1}^k r_{ij} \|x_i - \mu_j\|$$

where $r_{ij} = 1$ if x_i is assigned to cluster j ; otherwise, $r_{ij} = 0$. 计算可得

$$\mu_j = \frac{\sum_{i=1}^n r_{ij} x_i}{\sum_{i=1}^n r_{ij}}.$$

Clustering accuracy

Denote p_i and q_i be the obtained clustering label and the label provided by the dataset, respectively. The ACC is defined as follows

$$ACC = \frac{\sum_{i=1}^n \delta(p_i, \text{map}(q_i))}{n} \quad (1)$$

where $\delta(a, b) = 1$ if $a = b$; otherwise $\delta(a, b) = 0$. $\text{map}(\bullet)$ is the best permutation mapping function that matches the obtained clustering label to the equivalent label of the dataset.

Normalized mutual information

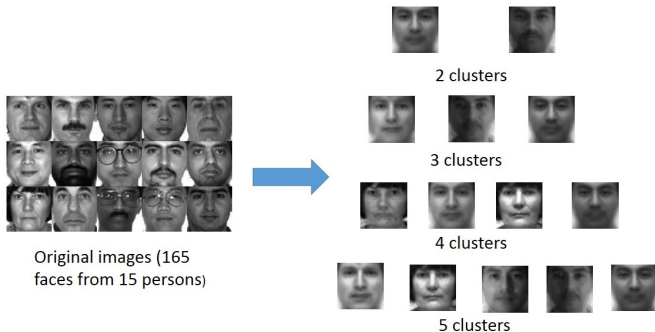
Given two random variables P and Q , NMI of P and Q is defined as

$$NMI(P, Q) = \frac{I(P; Q)}{\sqrt{H(P)H(Q)}} \quad (2)$$

where $I(P; Q)$ is the mutual information of P and Q , and $H(P)$ and $H(Q)$ are the entropies of P and Q , respectively. Here, the clustering results $\tilde{C} = \{\tilde{C}_i\}_{i=1}^{\tilde{c}}$ of and the ground truth labels $C = \{C_j\}_{j=1}^c$ of all samples are viewed as two discrete random variables. NMI is specified as

$$NMI(C, \tilde{C}) = \frac{\sum_{i=1}^{\tilde{c}} \sum_{j=1}^c |\tilde{C}_i \cap C_j| \log \frac{n |\tilde{C}_i \cap C_j|}{|\tilde{C}_i| |C_j|}}{\sqrt{(\sum_{i=1}^{\tilde{c}} |\tilde{C}_i| \log \frac{|\tilde{C}_i|}{n}) (\sum_{j=1}^c |C_j| \log \frac{|C_j|}{n})}}. \quad (3)$$

K-means clustering result from face images



Thanks!