

Supervised Learning

Review

- Given training data and label.
- Discriminative: estimate $\hat{\mathbb{P}}\{Y = y|X = x\}$ to classify.
- Generative: estimate $\hat{\mathbb{P}}\{X = x|Y = y\}$ and Bayes rule to classify.

Naive Bayes

Review

- Naive Bayes: $X_j \leftarrow Y$.

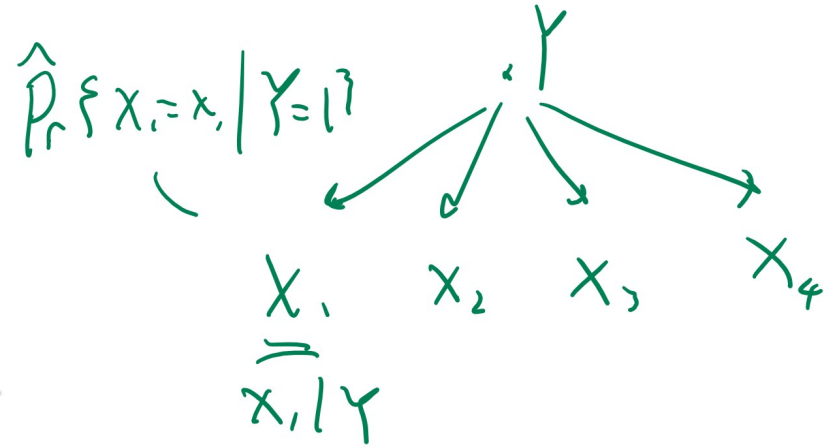
$P_r\{Y=0\}$ —————

$$\mathbb{P}\{Y = 1 | X_1 = x_1, \dots, X_m = x_m\}$$

$$\hat{\mathbb{P}}\{Y = 1\} \prod_{j=1}^m \hat{\mathbb{P}}\{X_j = x_j | Y = 1\}$$

$$= \frac{\hat{\mathbb{P}}\{Y = 1\} \prod_{j=1}^m \hat{\mathbb{P}}\{X_j = x_j | Y = 1\}}{\mathbb{P}\{X_1 = x_1, \dots, X_m = x_m\}} = \frac{1}{1}$$

$$= \frac{1 + \exp\left(\underbrace{-\log\left(\frac{\mathbb{P}\{Y = 1\}}{\mathbb{P}\{Y = 0\}}\right)}_b - \sum_{j=1}^m \underbrace{\log\left(\frac{\mathbb{P}\{X_j = x_j | Y = 1\}}{\mathbb{P}\{X_j = x_j | Y = 0\}}\right)}_{w^T x}\right)}{1}$$



$P_r\{Y=1\} P_r\{X=x | Y=1\} + P_r\{Y=0\} P_r\{X=x | Y=0\}$

Logistic Regression

Review

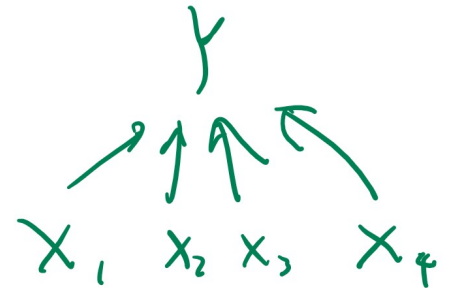
$$\frac{1}{1 + \exp \left(- \log \left(\frac{\hat{\mathbb{P}}\{Y = 1\}}{\hat{\mathbb{P}}\{Y = 0\}} \right) - \sum_{j=1}^m \log \left(\frac{\hat{\mathbb{P}}\{X_j = x_j | Y = 1\}}{\hat{\mathbb{P}}\{X_j = x_j | Y = 0\}} \right) \right)}$$

b
 $w^T x$

- Logistic Regression: $X_j \rightarrow Y$.

$$\tilde{\mathbb{P}}\{Y = 1 | X_1 = x_1, \dots, X_m = x_m\} = \frac{1}{1 + \exp \left(- \left(b + \sum_{j=1}^m w_j x_j \right) \right)}$$

$a \in (0, 1]$



Naive Bayes v Logistic Regression Derivation

Review

Generative Adversarial Network

Review

- Generative Adversarial Network (GAN): two competitive neural networks.
- ① Generative network input random noise and output fake images.
- ② Discriminative network input real and fake images and output label real or fake.

Generative Adversarial Network Diagram

Review

Midterm

Admin

$W1 \rightarrow W5$

Supervised.

- Materials: END HERE
- Calculator: pay 2 points *out of 40.*
- Formula sheet: will post
- Additional formula sheet: 2 points each
- NO examples, quiz questions, homework questions: 2 points each

request room

Unsupervised Learning

Motivation

- Supervised learning: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$.
- Unsupervised learning: x_1, x_2, \dots, x_n .
- There are a few common tasks without labels.
- ① Clustering: separate instances into groups.
- ② Novelty (outlier) detection: find instances that are different.
- ③ Dimensionality reduction: represent each instance with a lower dimensional feature vector while maintaining key characteristics.

Unsupervised Learning Applications

Motivation

- 1 Google News
- 2 Google Photo
- 3 Image Segmentation
- 4 Text Processing

group similar words into one type.

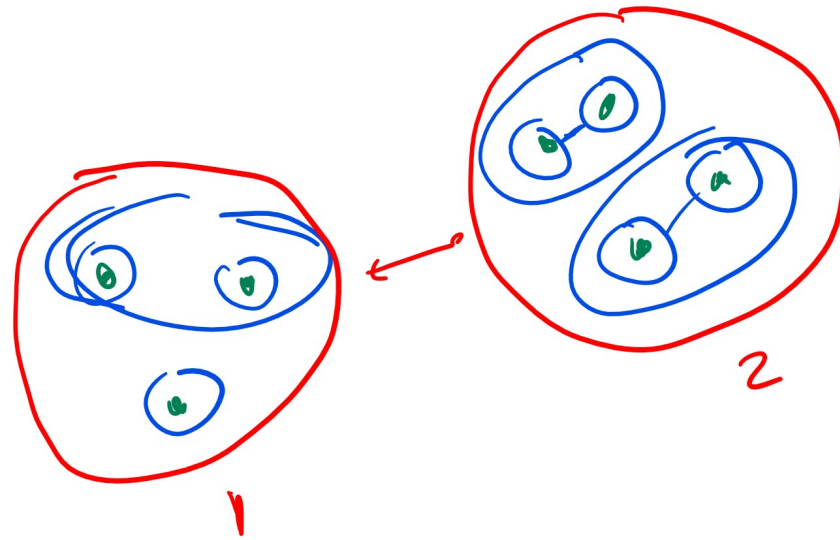
Hierarchical Clustering

Description

- Start with each instance as a cluster.
- Merge clusters that are closest to each other.
- Result in a binary tree with close clusters as children.

Hierarchical Clustering Diagram

Description



7 clusters
6
2

Clusters

Definition

- A cluster is a set of instances.

$$C_k \subseteq \{x_i\}_{i=1}^n$$

- A clustering is a partition of the set of instances into clusters.

$$C = C_1, C_2, \dots, C_K$$

$$C_k \cap C_{k'} = \emptyset \quad \forall k' \neq k, \quad \bigcup_{k=1}^K C_k = \{x_i\}_{i=1}^n$$

Distance between Points

Definition

- Usually, the distance between two instances is measured by the Euclidean distance or L_2 distance.

$$\rho(x_i, x_{i'}) = \|x_i - x_{i'}\|_2 = \sqrt{\sum_{j=1}^m (x_{ij} - x_{i'j})^2}$$

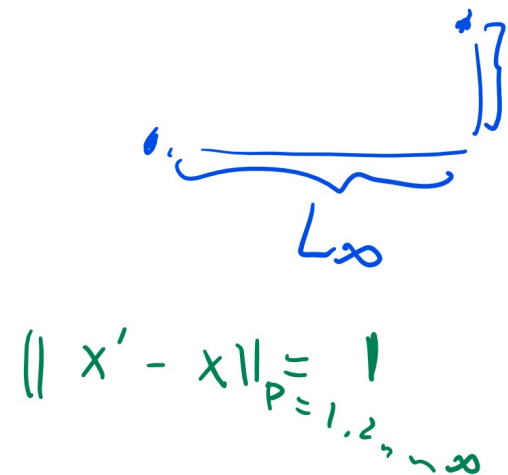
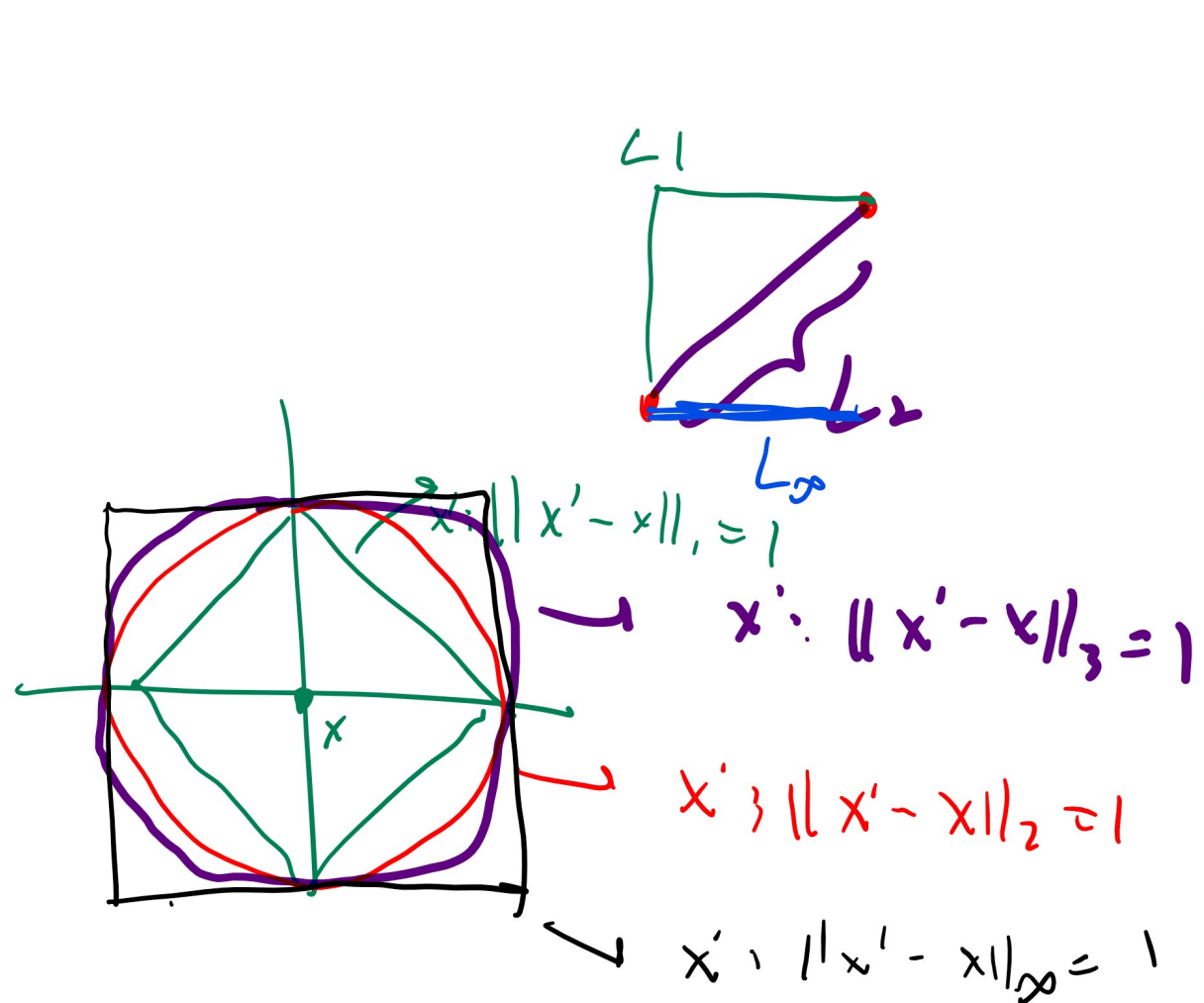
- Other examples include: L_1 distance and L_∞ distance.

$$\rho_1(x_i, x_{i'}) = \|x_i - x_{i'}\|_1 = \sum_{j=1}^m |x_{ij} - x_{i'j}|$$

$$\rho_\infty(x_i, x_{i'}) = \|x_i - x_{i'}\|_\infty = \max_{j=1,2,\dots,m} \{|x_{ij} - x_{i'j}|\}$$

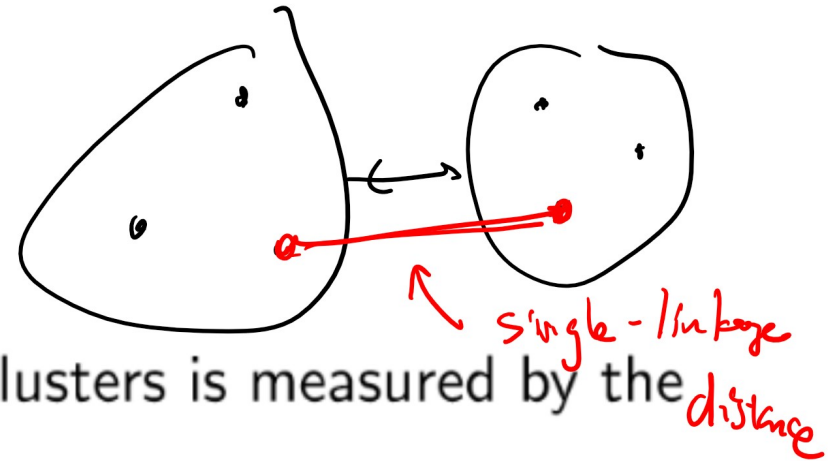
L_p Distance Diagram

Definition



Single Linkage Distance

Definition



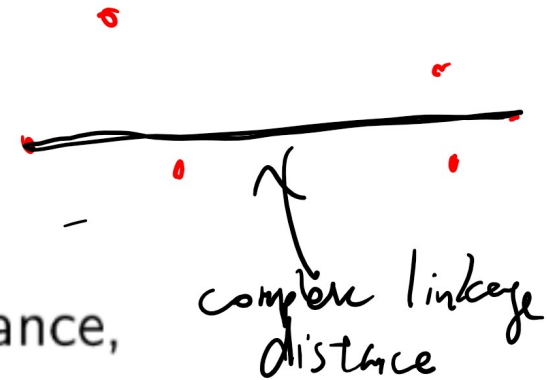
- Usually, the distance between two clusters is measured by the single-linkage distance.

$$\rho(C_k, C_{k'}) = \min \{ \rho(x_i, x_{i'}) : x_i \in C_k, x_{i'} \in C_{k'} \}$$

- It is the shortest distance from any instance in one cluster to any instance in the other cluster.

Complete Linkage Distance

Definition



- Another measure is complete-linkage distance,

$$\rho(C_k, C_{k'}) = \max \{ \rho(x_i, x_{i'}) : x_i \in C_k, x_{i'} \in C_{k'} \}$$

- It is the longest distance from any instance in one cluster to any instance in the other cluster.

Average Linkage Distance Diagram

Definition

- Another measure is average-linkage distance.

$$\rho(C_k, C_{k'}) = \frac{1}{|C_k| |C_{k'}|} \sum_{x_i \in C_k, x_{i'} \in C_{k'}} \rho(x_i, x_{i'})$$

- It is the average distance from any instance in one cluster to any instance in the other cluster.

Hierarchical Clustering Example 1 Part I

Quiz (Graded)

$$\sqrt{(x_1 - x_2)^2} = |x_1 - x_2|$$

Q2

- Spring 2018 Midterm Q5
- Given three clusters $A = \{0, \underline{4}, 6\}$, $B = \{\underline{3}, \underline{9}\}$, $C = \{\underline{11}\}$.
What is the next iteration of hierarchical clustering with Euclidean distance and single linkage?



• A: Merge A and B.

- B: Merge A and C.
- C: Merge B and C.
- D: No change, E: Do not choose.

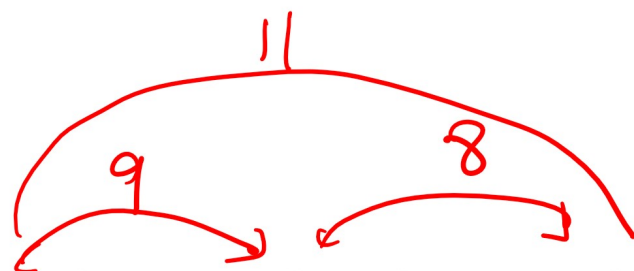
CS540S1
min distance between points in cluster

Hierarchical Clustering Example 1 Part II

Quiz (Graded)

Q4

- Spring 2018 Midterm Q5
- Given three clusters $A = \{0, 2, 6\}$, $B = \{3, 9\}$, $C = \{11\}$. What is the next iteration of hierarchical clustering with Euclidean distance and complete linkage?
- A: Merge A and B.
- B: Merge A and C.
- C: Merge B and C.
- D: No change, E: Do not choose.



complete linkage
 max to find dist.
 merge with dist.

Hierarchical Clustering Example 2

Quiz (Participation)

- Spring 2017 Midterm Q4
- Given the distance between the clusters so far. Which pair (choose 2) of clusters will be merged using single linkage.

—	A	B	C	D	E
A	0	1075	2013	2054	996
B	1075	0	3272	2687	2037
C	2013	3272	0	808	1307
D	2054	2687	808	0	1059

merge CD

	A	B	CD	E
A	0	1075	2013	996
B	1075	0	2687	2037
CD	2013	2687	0	1059

Hierarchical Clustering

Algorithm

- Input: instances: $\{x_i\}_{i=1}^n$, the number of clusters K , and a distance function ρ .
- Output: a list of clusters $C = C_1, C_2, \dots, C_K$.
- Initialize for $t = 0$.

$$C^{(0)} = C_1^{(0)}, \dots, C_n^{(0)}, \text{ where } C_k^{(0)} = \{x_k\}, k = 1, 2, \dots, n$$

- Loop for $t = 1, 2, \dots, n - k + 1$.

$$(k_1^*, k_2^*) = \arg \min_{k_1, k_2} \rho \left(C_{k_1}^{(t-1)}, C_{k_2}^{(t-1)} \right)$$

$$C^{(t)} = \left(C_{k_1^*}^{(t-1)} \cup C_{k_2^*}^{(t-1)} \right), C_1^{(t-1)}, \dots, \text{no } k_1^*, k_2^*, \dots, C_n^{(t-1)}$$

—

,

]

.

K Means Clustering

Description

- This is not K Nearest Neighbor.
- Start with random cluster centers.
- Assign each point to its closest center.
- Update all cluster centers as the center of its points.

K Means Clustering Diagram

Description

Center

Definition

- The center is the average of the instances in the cluster,

$$c_k = \frac{1}{|C_k|} \sum_{x \in C_k} x$$

Distortion

Distortion

- Distortion for a point is the distance from the point to its cluster center.
- Total distortion is the sum of distortion for all points.

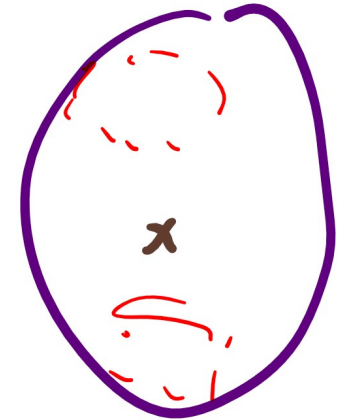
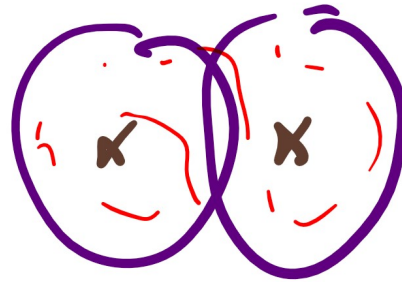
$$D_K = \sum_{i=1}^n \rho(x_i, c_{k^*(x_i)})$$

$$k^*(x) = \arg \min_{k=1,2,\dots,K} \rho(x, c_k)$$

Handwritten annotations:
 - An arrow points from "all training data" to the summation index $i=1$ to n .
 - A red box encloses $c_{k^*(x_i)}$ with an arrow pointing to "cluster center".
 - A red box encloses $k^*(x_i)$ with an arrow pointing to "cluster x_i is assigned to".
 - A red arrow points from x_i to $c_{k^*(x_i)}$.

Objective Function

Definition



← means converge L
local min

- This algorithm stop in finite steps.
- This algorithm is trying to minimize the total distortion but fails.

Gradient Descent

Definition

- When ρ is the Euclidean distance. K Means algorithm is the gradient descent when distortion is the objective (cost) function.

$$\frac{\partial}{\partial c_k} \sum_{k=1}^K \sum_{x \in C_k} \|x - c_k\|_2^2 = 0$$

$$\Rightarrow -2 \sum_{x \in C_k} (x - c_k) = 0$$

$$\Rightarrow c_k = \frac{1}{|C_k|} \sum_{x \in C_k} x$$

gradient descent for total distortion as cost.

same as k means

Gradient Descent Derivation

Derivation

K Means Clustering Example Part I

Quiz (Graded)

Q5

• Spring 2018 Midterm Q5

• Given data $\{5, 7, 10, 12\}$ and initial cluster centers $c_1 = 3, c_2 = 13$, what is the initial clusters?

• A: $\{5, 7\}$ and $\{10, 12\}$

• B: $\{5\}$ and $\{7, 10, 12\}$

• C: $\{5, 7, 10\}$ and $\{12\}$

• D: none of the above, E: do not choose.

$$\begin{array}{l} \{5, 7\} \\ \hline \uparrow \\ c_1 = 6 \end{array} \quad , \quad \begin{array}{l} \{10, 12\} \\ \hline \uparrow \\ c_2 = 11 \end{array}$$

K Means Clustering Example Part II

Quiz (Graded)

- Spring 2018 Midterm Q5
- Given data $\{5, 7, 10, 12\}$ and initial cluster centers $c_1 = 3, c_2 = 13$, what are the cluster in the next iteration?
- A: $\{5, 7\}$ and $\{10, 12\}$ 6. 11
- B: $\{5\}$ and $\{7, 10, 12\}$
- C: $\{5, 7, 10\}$ and $\{12\}$
- D: none of the above, E: do not choose.

K Means Clustering

Algorithm

- Input: instances: $\{x_i\}_{i=1}^n$, the number of clusters K , and a distance function ρ .
- Output: a list of clusters $C = C_1, C_2, \dots, C_K$.
- Initialize $t = 0$.

$$c_k^{(0)} = K \text{ random points}$$

- Loop until $c^{(t)} = c^{(t-1)}$.

$$C_k^{(t-1)} = \left\{ x : k = \arg \min_{k' \in \{1, 2, \dots, K\}} \rho(x, c_{k'}^{(t-1)}) \right\} \rightarrow \text{label}$$

$$c_k^{(t)} = \frac{1}{|C_k^{(t-1)}|} \sum_{x \in C_k^{(t-1)}} x \rightarrow \text{recenter}$$

Number of Clusters

Discussion

- There are a few ways to pick the number of clusters K .

① K can be chosen using prior knowledge about X .

② ~~K can be the one that minimizes distortion? No, when $K = n$, distortion = 0.~~

③ K can be the one that minimizes distortion + regularizer.

$$K^* = \arg \min_k (D_k + \lambda \cdot m \cdot k \cdot \log n)$$

dimension of each x

of x instances

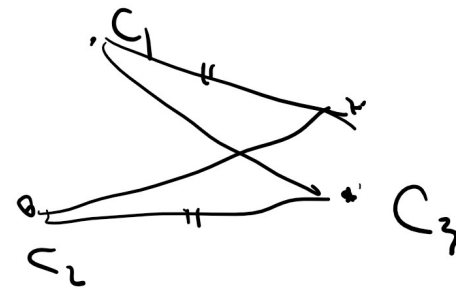
large cost for
large # of cluster k

- λ is a fixed constant chosen arbitrarily.

Initial Clusters

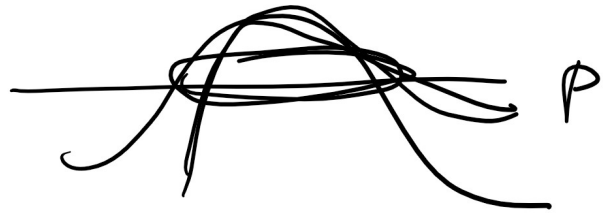
Discussion

- There are a few ways to initialize the clusters.
- ① K uniform random points in $\{x_i\}_{i=1}^n$.
- ② 1 uniform random point in $\{x_i\}_{i=1}^n$ as $c_1^{(0)}$, then find the farthest point in $\{x_i\}_{i=1}^n$ from $c_1^{(0)}$ as $c_2^{(0)}$, and find the farthest point in $\{x_i\}_{i=1}^n$ from the closer of $c_1^{(0)}$ and $c_2^{(0)}$ as $c_3^{(0)}$, and repeat this K times.

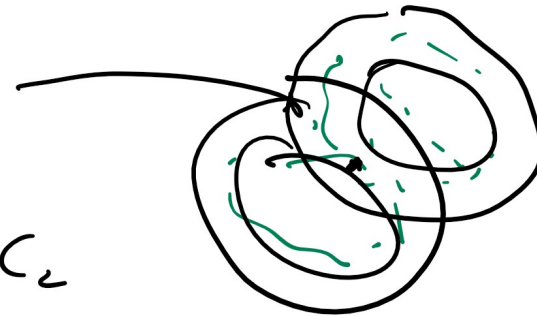


Gaussian Mixture Model

Discussion



prob belong
to C_1 and C_2



- In K means, each instance belong to one cluster with certainty.
- One continuous version is called the Gaussian mixture model: each instance belongs to one of the clusters with a positive probability.
- The model can be trained using Expectation Maximization Algorithm (EM Algorithm).

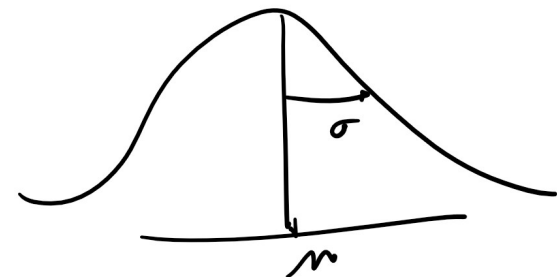
EM Algorithm, Part I

Discussion

- The means μ_k and variances σ_k^2 for each cluster need to be trained. The mixing probability π_k also needs to be trained.

$$\underbrace{(\mu_1, \sigma_1^2, \pi_1), (\mu_2, \sigma_2^2, \pi_2), \dots, (\mu_K, \sigma_K^2, \pi_K)}$$

- Initialize by random guesses of clusters means and variances.



EM Algorithm, Part II

Discussion

- Expectation Step. Compute responsibilities for $i = 1, 2, \dots, n$ and $k = 1, 2, \dots, K$.

prob of point i assigned to cluster k ,

$$\hat{\gamma}_{i,k} = \frac{\hat{\pi}_k \phi_k(x_i)}{\sum_{k'=1,2,\dots,K} \hat{\pi}_{k'} \phi_{k'}(x_i)}$$

$$\phi_k(x) = \frac{1}{\sqrt{2\pi\hat{\sigma}_k}} \exp\left(-\frac{(x - \hat{\mu}_k)^2}{2\hat{\sigma}_k^2}\right)$$

EM Algorithm, Part III

Discussion

- Maximization Step. Compute means and variances for each $k = 1, 2, \dots, K$.

$$\hat{\mu}_k = \frac{\sum_{i=1}^n \hat{\gamma}_{i,k} x_i}{\sum_{i=1}^n \hat{\gamma}_i}, \text{ and } \hat{\sigma}_k^2 = \frac{\sum_{i=1}^n \hat{\gamma}_{i,k} (x_i - \hat{\mu}_k)^2}{\sum_{i=1}^n \hat{\gamma}_i}$$

$$\hat{\pi}_k = \frac{1}{n} \sum_{i=1}^n \hat{\gamma}_{i,k}$$

Handwritten notes and calculations:

- A circle containing "a.c.f."
- A circle containing "0.7", "0.1", "0.2", and "x_i".
- Calculation: $\mu_1 = \frac{0.9 \cdot x_1 + 0.9 x_2 + 0.1 x_4}{0.1 + 0.9 + 0.9 + 0.1}$

- Repeat until convergent.

$$C_1 = \frac{x_1 + x_2 + x_3}{3}$$

Gaussian Mixture Model Diagram

Discussion