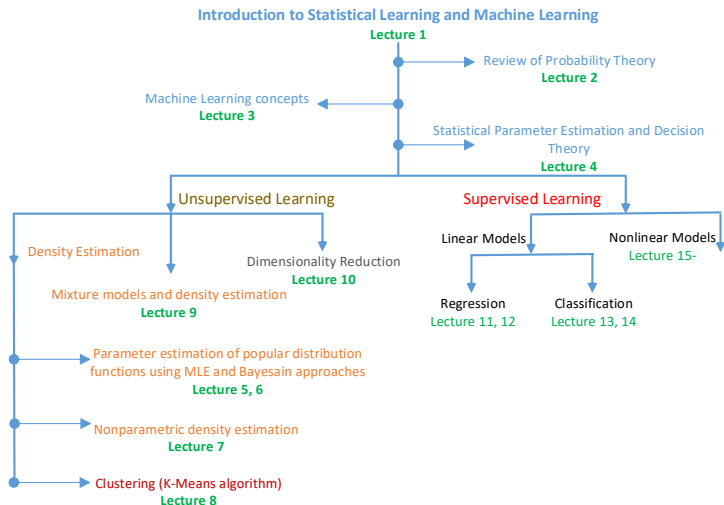


Statistical Learning and Machine Learning

Lecture 9 - Expectation Maximization

September 25, 2021

Course overview and where do we stand



Objectives of the lecture

- Review of mixture of Gaussians model
- Parameter estimation for mixture of Gaussians using Maximum Likelihood Estimation (and EM algorithm)
- Nearest Centroid vs K-Nearest Neighbour Classification

Mixture of Gaussians

Reminder: When using a superposition of K Gaussians, the resulting distribution has the form:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (9)$$

and is called **mixture of Gaussian**.

The parameters π_k are called *mixing coefficients* and satisfy:

$$0 \leq \pi_k \leq 1, \quad \sum_{k=1}^K \pi_k = 1.$$

Thus, the mixing coefficients have the form of probabilities ($\pi_k = p(k)$).

Mixture of Gaussians: Latent variables

Let \mathbf{z} be a binary random variable having a 1-of- K representation satisfying:

$$z_k \in \{0, 1\}, \quad \sum_{k=1}^K z_k = 1.$$

\mathbf{z} is associated with the marginal distribution $p(\mathbf{x})$, which is specified as:

$$p(z_k = 1) = \pi_k$$

Reminder: Since \mathbf{z} is a 1-of- K vector:

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

Mixture of Gaussians: Latent variables

The conditional distribution of \mathbf{x} for the k -th Gaussian is:

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \Sigma_k)$$

and with respect to \mathbf{z} is:

$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \Sigma_k)^{z_k}$$

We can now write $p(\mathbf{x})$:

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \Sigma_k) \quad (10)$$

Compare (10) with (9)!

Mixture of Gaussians: Latent variables

The conditional probability of \mathbf{z} given \mathbf{x} is (using Bayes' theorem):

$$\begin{aligned}\gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}\end{aligned}\tag{11}$$

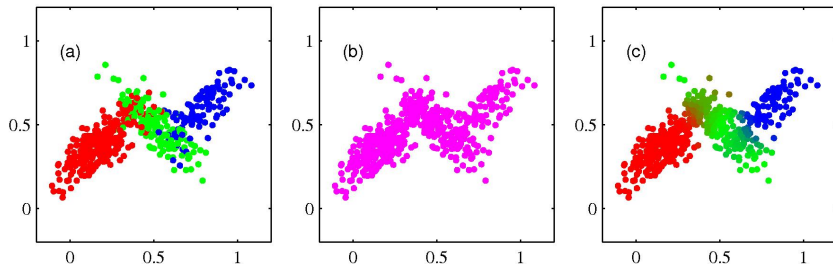
Thus, we consider π_k as the prior probability of $z_k = 1$ and $\gamma(z_k)$ as the corresponding posterior probability after observing \mathbf{x} (also called *responsibility* of component k for 'explaining' \mathbf{x}).

Mixture of Gaussians: Latent variables

Technique to generate random samples described by the Gaussian mixture model:

- 1 Generate a random vector $\hat{\mathbf{z}}$ using $p(\mathbf{z})$ (multinomial distribution)
- 2 Generate \mathbf{x} using $p(\mathbf{x}|\hat{\mathbf{z}}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \Sigma_k)^{\hat{z}_k}$

Use $p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$, $p(\mathbf{x})$ and $\gamma(\mathbf{z})$ for visualization.



Mixture of Gaussians: Parameter estimation using MLE

Given a set of i.i.d. data points $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, we want to estimate the parameters of the mixture of Gaussians:

$$\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k\}, \quad k = 1, \dots, K$$

The log-likelihood function is:

$$\begin{aligned} \ln p(\mathcal{D} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \ln \prod_{n=1}^N p(\mathbf{x}_n) \\ &= \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\} \end{aligned}$$

Maximization of the log-likelihood function is more complex than in the case of a single Gaussian. Why?

Mixture of Gaussians: Parameter estimation using MLE

Problem of singularity

- Consider the case where the covariance matrices of Gaussian mixture model are diagonal: $\Sigma_k = \sigma_k^2 \mathbf{I}$
- Further consider that $\mu_j = \mathbf{x}_n$
- The data point will contribute the following term to the likelihood function:

$$\mathcal{N}(\mathbf{x}_n | \mathbf{x}_n, \sigma_j^2 \mathbf{I}) = \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma_j}$$

For $\sigma_j \rightarrow 0$, the above term will go to infinity and thus maximization of the likelihood function will not make sense.

- **Remedy:** Detecting the singularity and resetting the mean μ_j to random value and resetting the covariance to some large value.

Mixture of Gaussians: Parameter estimation using MLE

We set the derivative of the log-likelihood function w.r.t. each parameter equal to zero.

$$\begin{aligned} 0 &= \frac{\partial \ln p(\mathcal{D}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}_k} \\ 0 &= - \sum_{n=1}^N \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^N \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \boldsymbol{\Sigma}_k (\mathbf{x}_n - \boldsymbol{\mu}_k) \end{aligned} \quad (12)$$

By multiplying with $\boldsymbol{\Sigma}_k^{-1}$ (assuming $\boldsymbol{\Sigma}_k$ is nonsingular matrix):

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

with $N_k = \sum_{n=1}^N \gamma(z_{nk})$ being the effective number of data points assigned to cluster k .

Mixture of Gaussians: Parameter estimation using MLE

We set the derivative of the log-likelihood function w.r.t. Σ_k equal to zero:

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk})(\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T$$

Thus, Σ_k is the covariance matrix of a single Gaussian fitted to all data, with each data point weighted by the corresponding responsibility value and divided with the effective number of points.

Mixture of Gaussians: Parameter estimation using MLE

To optimize w.r.t. π_k , we need to also consider the constraint that $\sum_{k=1}^K \pi_k = 1$. To do so, we use a Lagrange multiplier and maximize for:

$$\ln p(\mathcal{D}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

We set the derivative of the above Lagrangian function equal to zero and:

$$\sum_{n=1}^N \frac{\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} + \lambda = 0$$

leading to $\pi_k = N_k/N$ (π_k is the average responsibility of the component to explain the data).

Mixture of Gaussians: Expectation-Maximization

EM for Gaussian Mixtures

1. Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k , and evaluate the initial value of the log likelihood.
2. **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}.$$

3. **M step.** Re-estimate the parameters using the current responsibilities

$$\mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k^{\text{new}}) (\mathbf{x}_n - \mu_k^{\text{new}})^T$$

$$\pi_k^{\text{new}} = \frac{N_k}{N}$$

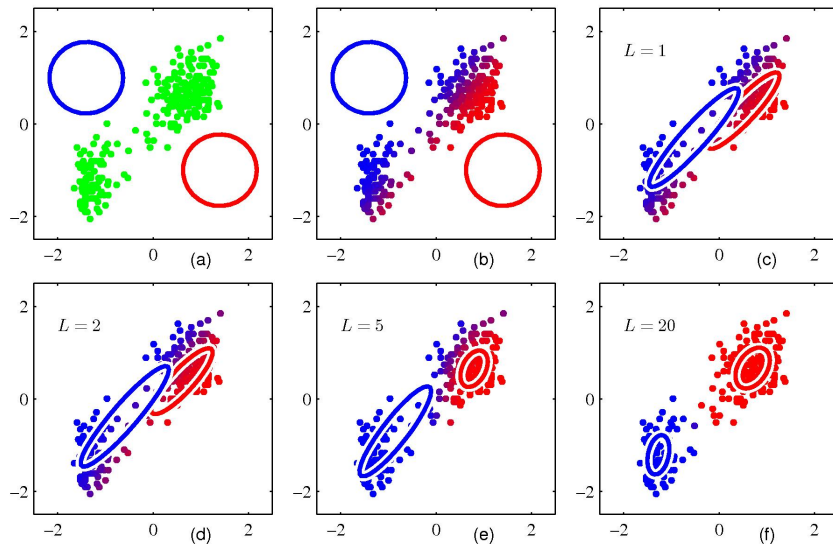
where
$$N_k = \sum_{n=1}^N \gamma(z_{nk}).$$

4. Evaluate the log likelihood

$$\ln p(\mathbf{X} | \mu, \Sigma, \pi) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

and check for convergence of either the parameters or the log likelihood. If the convergence criterion is not satisfied return to step 2.

Mixture of Gaussians: Expectation-Maximization



Classification: Formal Definition

- Given a **training set**
 - ① Observations in the form of data **features**
 - ② **class labels** corresponding to each set of observation(s)
- Generate a **function/classifier**
 - ① input: observation/features
 - ② output: class label
- The classifier used for classifying previously unseen records

Nearest Centroid-based classification

- **Training:** compute center or centroid for each of K classes
 - ① center of all points of that class
 - ② can be viewed as centroid for a cluster as in K-means
- **Classification:**
 - ① assign each data point to the class corresponding to the nearest centroid
- How is this different from k-nearest neighbour classifier?

Nearest Centroid-based classification



Figure: 2D data points generated from Gaussian distribution with identity covariance matrix and means $(-1, 0)$ and $(1, 0)$.

Source: STOR 390: Introduction to Data Science by Iain Carmichael.

Nearest Centroid-based classification

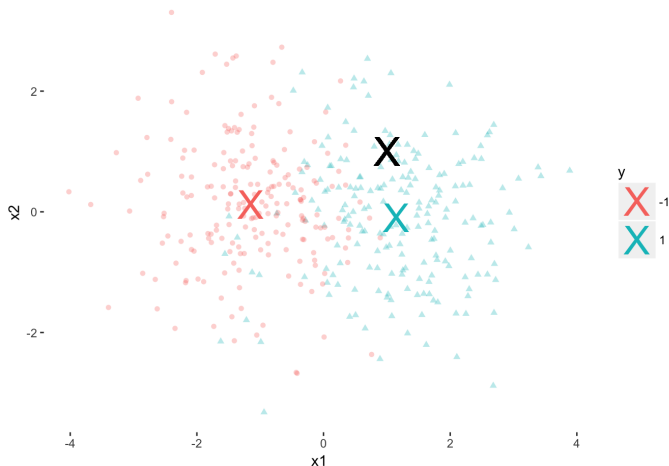


Figure: Generation of centroids of two classes by taking means of the training data points corresponding to the classes.

Source: STOR 390: Introduction to Data Science by Iain Carmichael.

Nearest Centroid-based classification

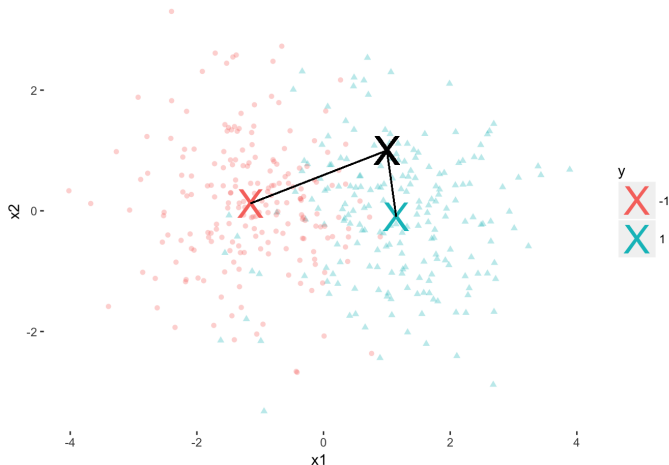


Figure: Find the distances of the test data point with the two means. Assign the test data point (black cross) to the class 1 since the distance of the test data point to the mean of class 1 is smaller.

Source: STOR 390: Introduction to Data Science by Iain Carmichael.

Nearest Centroid-based classification

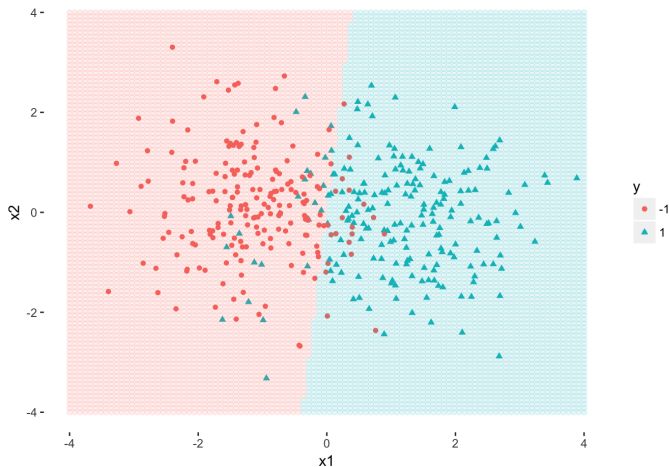


Figure: Nearest Centroid-based classification is an example of Linear classification since the the 2 classes are separated by a line.

Source: STOR 390: Introduction to Data Science by Iain Carmichael.

K-Nearest Neighbour Classification

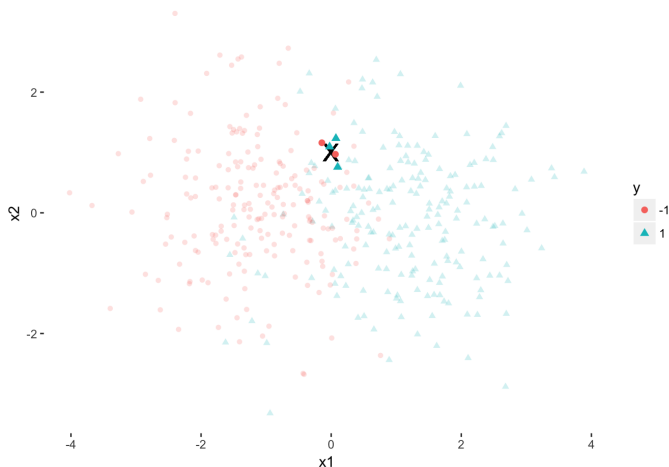


Figure: K-Nearest Neighbour Classification on the previous data set.

Source: [STOR 390: Introduction to Data Science](#) by Iain Carmichael.

K-Nearest Neighbour Classification

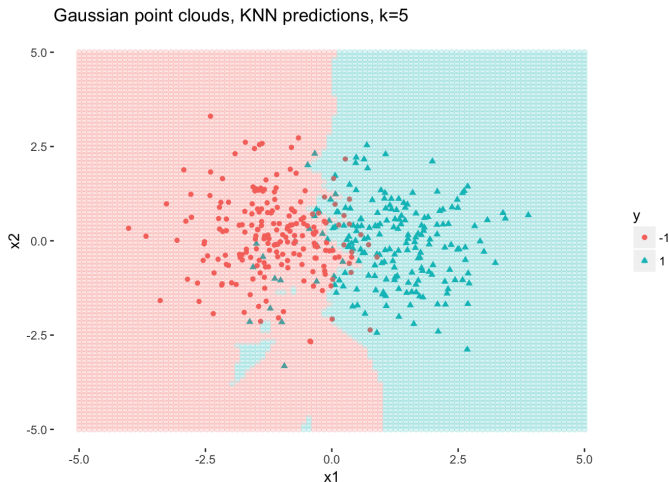


Figure: K-Nearest Neighbour Classification on the previous data set ($K=5$).

Source: [STOR 390: Introduction to Data Science](#) by Iain Carmichael.

Nearest Centroid vs K-Nearest Neighbour

Which of the two classifiers do you think will perform better in the following cases?

- Nonlinear classification
- Label noise (some training data points are mislabelled)
- Presence of outliers in data

Nearest Centroid vs K-Nearest Neighbour

- Nearest Centroid

- ① requires a lot of memory (all training data must be stored at all times)
- ② slow to compute
- ③ robust to outliers

- K-Nearest Neighbour

- ① requires little memory
- ② fast
- ③ robust to class label noise