

Foundations of Machine Learning
African Masters in Machine Intelligence



AIMS | African Institute for
Mathematical Sciences
RWANDA


**Imperial College
London**

Density Estimation with Gaussian Mixture Models

Marc Deisenroth

Quantum Leap Africa
African Institute for Mathematical
Sciences, Rwanda

Department of Computing
Imperial College London

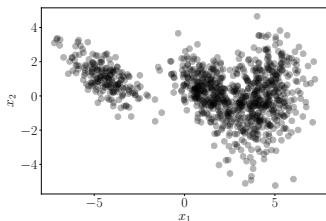
 @mpd37
mdeisenroth@aimsammi.org

October 30, 2018

Reading Material

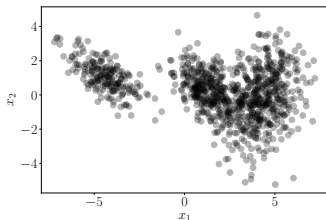
- ▶ Mathematics for Machine Learning (Chapter 11): mml-book.com
- ▶ Pattern Recognition and Machine Learning (Chapter 9)

Problem Statement



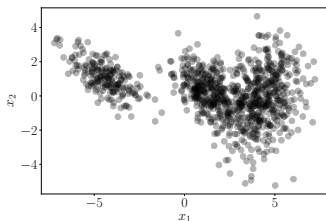
- **Density estimation:** Given a dataset (unlabeled), find a probability density function from which the data could have plausibly been generated

Problem Statement



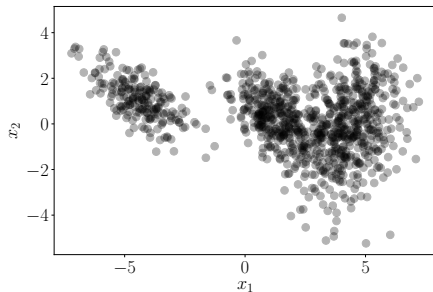
- ▶ **Density estimation:** Given a dataset (unlabeled), find a probability density function from which the data could have plausibly been generated
- ▶ Typically: Fix the class/model of densities and find optimal parameters given this class

Problem Statement



- ▶ **Density estimation:** Given a dataset (unlabeled), find a probability density function from which the data could have plausibly been generated
- ▶ Typically: Fix the class/model of densities and find optimal parameters given this class
- ▶ Example. Class: Gaussian; Find mean and variance
 - ▶▶ MLE/MAP estimation

Problem Statement (2)



- ▶ Gaussians (or similarly all other distributions we encountered so far) have very limited modeling capabilities: Too simple
 - ▶ **Mixture models** are more flexible

Overview

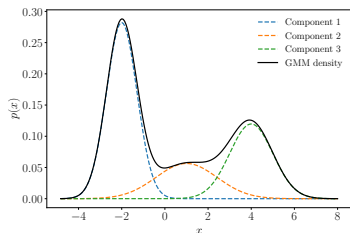
Gaussian Mixture Models

Parameter Learning

Implementation

Probabilistic Perspective

Gaussian Mixture Model



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

$$0 \leq \pi_k \leq 1$$

$$\sum_{k=1}^K \pi_k = 1$$

- ▶ Individual components are Gaussian distributions
- ▶ Each component is weighted by π_k (mixture weights)

Overview

Gaussian Mixture Models

Parameter Learning

Implementation

Probabilistic Perspective

Parameter Learning for GMMs

- ▶ Objective: Maximum likelihood estimate of model parameters θ given a dataset \mathcal{X}
- ▶ $\theta := \{\pi_k, \mu_k, \Sigma_k, k = 1, \dots, K\}$
- ▶ Maximum likelihood estimate:

Parameter Learning for GMMs

- ▶ Objective: Maximum likelihood estimate of model parameters θ given a dataset \mathcal{X}
- ▶ $\theta := \{\pi_k, \mu_k, \Sigma_k, k = 1, \dots, K\}$
- ▶ Maximum likelihood estimate:

$$\theta^* = \arg \max_{\theta} p(\mathcal{X} | \theta) = \arg \max_{\theta} \sum_{k=1}^K \pi_k \mathcal{N}(\mathcal{X} | \mu_k, \Sigma_k)$$

$$\stackrel{\log}{=} \arg \max_{\theta} \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathcal{X} | \mu_k, \Sigma_k)$$

Parameter Learning for GMMs

- ▶ Objective: Maximum likelihood estimate of model parameters θ given a dataset \mathcal{X}
- ▶ $\theta := \{\pi_k, \mu_k, \Sigma_k, k = 1, \dots, K\}$
- ▶ Maximum likelihood estimate:

$$\begin{aligned}\theta^* &= \arg \max_{\theta} p(\mathcal{X} | \theta) = \arg \max_{\theta} \sum_{k=1}^K \pi_k \mathcal{N}(\mathcal{X} | \mu_k, \Sigma_k) \\ &\stackrel{\log}{=} \arg \max_{\theta} \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathcal{X} | \mu_k, \Sigma_k)\end{aligned}$$

- ▶ **Problem:** We cannot move the log into the sum
- ▶ **Difficult optimization problem**

Parameter Learning for GMMs

- ▶ Objective: Maximum likelihood estimate of model parameters θ given a dataset \mathcal{X}
- ▶ $\theta := \{\pi_k, \mu_k, \Sigma_k, k = 1, \dots, K\}$
- ▶ Maximum likelihood estimate:

$$\theta^* = \arg \max_{\theta} p(\mathcal{X} | \theta) = \arg \max_{\theta} \sum_{k=1}^K \pi_k \mathcal{N}(\mathcal{X} | \mu_k, \Sigma_k)$$

$$\stackrel{\log}{=} \arg \max_{\theta} \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathcal{X} | \mu_k, \Sigma_k)$$

- ▶ **Problem:** We cannot move the log into the sum
 - ▶ **Difficult optimization problem**
 - ▶ Iterative scheme (**EM Algorithm**) for learning parameters

GMM Likelihood

Assume an i.i.d. data set $\mathcal{X} = \mathbf{x}_1, \dots, \mathbf{x}_N$ is given, and we want to determine the optimal parameters $\boldsymbol{\theta}^*$ of the GMM via Maximum Likelihood

1. Likelihood:

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

GMM Likelihood

Assume an i.i.d. data set $\mathcal{X} = \mathbf{x}_1, \dots, \mathbf{x}_N$ is given, and we want to determine the optimal parameters $\boldsymbol{\theta}^*$ of the GMM via Maximum Likelihood

1. Likelihood:

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

2. Log-likelihood:

$$\log p(\mathcal{X}|\boldsymbol{\theta}) = \sum_{i=1}^N \log p(\mathbf{x}_i|\boldsymbol{\theta}) = \underbrace{\sum_{i=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{=:L}$$

Necessary Optimality Conditions

Learning Objective

Find parameters θ^* that maximize the log-likelihood

$$\log p(\mathcal{X}|\theta) = \sum_{i=1}^N \log p(x_i|\theta) = \underbrace{\sum_{i=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k)}_{=:L}$$

Necessary Optimality Conditions

Learning Objective

Find parameters θ^* that maximize the log-likelihood

$$\log p(\mathcal{X}|\theta) = \sum_{i=1}^N \log p(x_i|\theta) = \underbrace{\sum_{i=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k)}_{=:L}$$

$$\frac{\partial L}{\partial \mu_k} = \mathbf{0}^\top \iff \sum_{i=1}^N \frac{\partial \log p(x_i|\theta)}{\partial \mu_k} = \mathbf{0}^\top$$

$$\frac{\partial L}{\partial \Sigma_k} = \mathbf{0} \iff \sum_{i=1}^N \frac{\partial \log p(x_i|\theta)}{\partial \Sigma_k} = \mathbf{0}$$

$$\frac{\partial L}{\partial \pi_k} = 0 \iff \sum_{i=1}^N \frac{\partial \log p(x_i|\theta)}{\partial \pi_k} = 0$$

High-Level Gradients

We need to compute gradients of the form

$$\frac{\partial \log p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

High-Level Gradients

We need to compute gradients of the form

$$\frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

In general,

$$\frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{1}{p(\mathbf{x}_i|\boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

High-Level Gradients

We need to compute gradients of the form

$$\frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

In general,

$$\frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{1}{p(\mathbf{x}_i|\boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

With

$$p(\mathbf{x}_i|\boldsymbol{\theta}) =$$

High-Level Gradients

We need to compute gradients of the form

$$\frac{\partial \log p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

In general,

$$\frac{\partial \log p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{1}{p(\mathbf{x}_i | \boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

With

$$p(\mathbf{x}_i | \boldsymbol{\theta}) = \sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

High-Level Gradients

We need to compute gradients of the form

$$\frac{\partial \log p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

In general,

$$\frac{\partial \log p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{1}{p(\mathbf{x}_i | \boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

With

$$p(\mathbf{x}_i | \boldsymbol{\theta}) = \sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

we get

$$\frac{\partial p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} =$$

High-Level Gradients

We need to compute gradients of the form

$$\frac{\partial \log p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

In general,

$$\frac{\partial \log p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{1}{p(\mathbf{x}_i | \boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

With

$$p(\mathbf{x}_i | \boldsymbol{\theta}) = \sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

we get

$$\frac{\partial p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k}$$

In More Detail

$$\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} =$$

In More Detail

$$\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k}$$

In More Detail

$$\begin{aligned}\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} &= \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \\ &= \pi_k (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\end{aligned}$$

In More Detail

$$\begin{aligned}\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} &= \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \\ &= \pi_k (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\end{aligned}$$

Overall:

$$\frac{\partial L}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \frac{1}{p(\mathbf{x}_i|\boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k}$$

In More Detail

$$\begin{aligned}\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} &= \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \\ &= \pi_k (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\end{aligned}$$

Overall:

$$\begin{aligned}\frac{\partial L}{\partial \boldsymbol{\mu}_k} &= \sum_{i=1}^N \frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \frac{1}{p(\mathbf{x}_i|\boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} \\ &= \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{:=r_{ik}}\end{aligned}$$

In More Detail

$$\begin{aligned}\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} &= \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \\ &= \pi_k (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\end{aligned}$$

Overall:

$$\begin{aligned}\frac{\partial L}{\partial \boldsymbol{\mu}_k} &= \sum_{i=1}^N \frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \frac{1}{p(\mathbf{x}_i|\boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} \\ &= \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{:=r_{ik}} = \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} = \mathbf{0}^\top\end{aligned}$$

In More Detail

$$\begin{aligned}\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} &= \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \\ &= \pi_k (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\end{aligned}$$

Overall:

$$\begin{aligned}\frac{\partial L}{\partial \boldsymbol{\mu}_k} &= \sum_{i=1}^N \frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \frac{1}{p(\mathbf{x}_i|\boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} \\ &= \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{:=r_{ik}} = \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} = \mathbf{0}^\top \\ &\iff \sum_{i=1}^N r_{ik} \mathbf{x}_i = \sum_{i=1}^N r_{ik} \boldsymbol{\mu}_k\end{aligned}$$

In More Detail

$$\begin{aligned}\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} &= \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \\ &= \pi_k (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\end{aligned}$$

Overall:

$$\begin{aligned}\frac{\partial L}{\partial \boldsymbol{\mu}_k} &= \sum_{i=1}^N \frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \frac{1}{p(\mathbf{x}_i|\boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} \\ &= \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{:=r_{ik}} = \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} = \mathbf{0}^\top \\ \iff \sum_{i=1}^N r_{ik} \mathbf{x}_i &= \sum_{i=1}^N r_{ik} \boldsymbol{\mu}_k \iff \boldsymbol{\mu}_k^* = \frac{\sum_{i=1}^N r_{ik} \mathbf{x}_i}{\sum_{i=1}^N r_{ik}} = \frac{1}{N_k} \sum_{i=1}^N r_{ik} \mathbf{x}_i\end{aligned}$$

Similarly...

$$\frac{\partial L}{\partial \mathbf{\Sigma}_k} = \mathbf{0} \iff \mathbf{\Sigma}_k^* = \frac{1}{N_k} \sum_{i=1}^N r_{ik} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top$$

$$\frac{\partial L}{\partial \pi_k} = \mathbf{0}^\top \iff \pi_k^* = \frac{N_k}{N} \quad \blacktriangleright \text{Requires Lagrange multipliers}$$

\blacktriangleright See Chapter 11 of “Mathematics for Machine Learning” for details

Similarly...

$$\frac{\partial L}{\partial \mathbf{\Sigma}_k} = \mathbf{0} \iff \mathbf{\Sigma}_k^* = \frac{1}{N_k} \sum_{i=1}^N r_{ik} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top$$

$$\frac{\partial L}{\partial \pi_k} = \mathbf{0}^\top \iff \pi_k^* = \frac{N_k}{N} \quad \blacktriangleright \text{Requires Lagrange multipliers}$$

- ▶ See Chapter 11 of “Mathematics for Machine Learning” for details
- ▶ **Bad news:** These results do not constitute a closed-form solution of the parameters $\boldsymbol{\mu}_k, \mathbf{\Sigma}_k, \pi_k$ of the mixture model because the responsibilities r_{ik} depend on those parameters in a complex way.

Similarly...

$$\frac{\partial L}{\partial \Sigma_k} = \mathbf{0} \iff \Sigma_k^* = \frac{1}{N_k} \sum_{i=1}^N r_{ik} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top$$

$$\frac{\partial L}{\partial \pi_k} = \mathbf{0}^\top \iff \pi_k^* = \frac{N_k}{N} \quad \blacktriangleright \text{Requires Lagrange multipliers}$$

- ▶ See Chapter 11 of “Mathematics for Machine Learning” for details
 - ▶ **Bad news:** These results do not constitute a closed-form solution of the parameters $\boldsymbol{\mu}_k, \Sigma_k, \pi_k$ of the mixture model because the responsibilities r_{ik} depend on those parameters in a complex way.
 - ▶ **Good news:** Results suggest a simple **iterative** scheme for finding a solution to the MLE problem: Compute responsibilities and then update one parameter at a time while keeping the other ones fixed ▶ **Expectation Maximization** algorithm

Overview

Gaussian Mixture Models

Parameter Learning

Implementation

Probabilistic Perspective

Expectation Maximization (EM) Algorithm

- ▶ Iterative scheme for learning parameters in mixture models and latent-variable models
 1. Choose initial values for μ_k, Σ_k, π_k
 2. Until convergence, alternate between
 - ▶ **E-step:** Evaluate the responsibilities r_{ik} (posterior probability of data point i belonging to mixture component k)
 - ▶ **M-step:** Use the updated responsibilities to re-estimate the parameters μ_k, Σ_k, π_k
- ▶ Every step in the EM algorithm increases the likelihood function
- ▶ Convergence: Check log-likelihood or the parameters

Implementation

1. Initialize μ_k, Σ_k, π_k

Implementation

1. Initialize $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$
2. **E-step:** Evaluate responsibilities for every data point \mathbf{x}_i using current parameters $\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$:

$$r_{ik} = \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

Implementation

1. Initialize $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$
2. **E-step:** Evaluate responsibilities for every data point \mathbf{x}_i using current parameters $\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$:

$$r_{ik} = \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

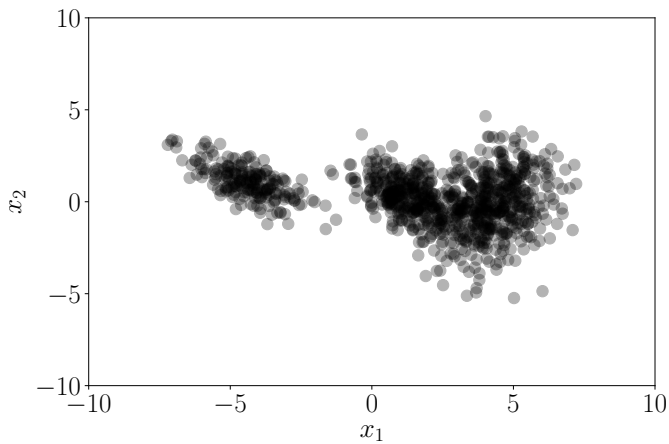
3. **M-step:** Re-estimate parameters $\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ using the current responsibilities r_{ik} (from E-step):

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} \mathbf{x}_i$$

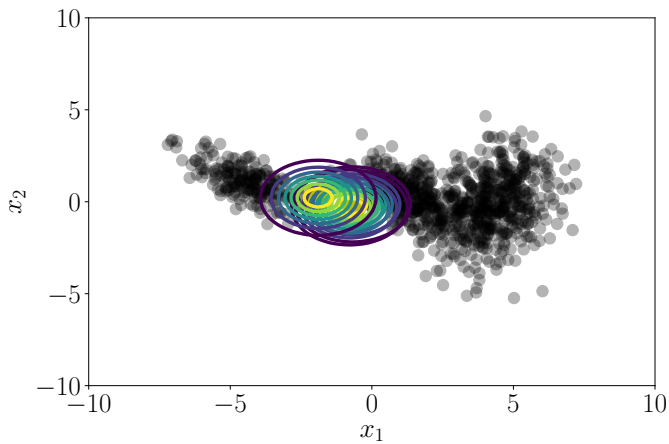
$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top$$

$$\pi_k = \frac{N_k}{N}$$

Example

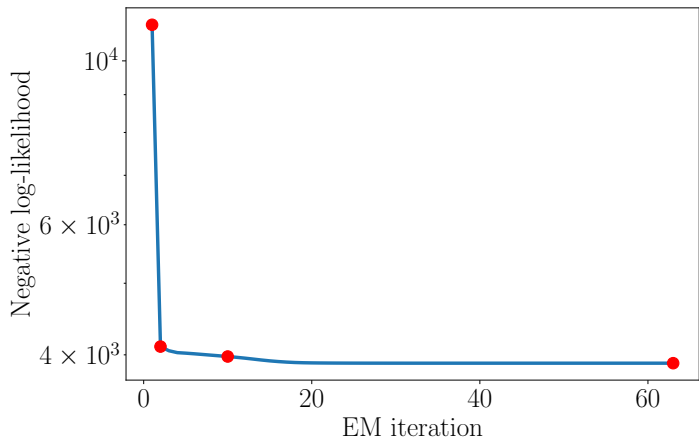


Example



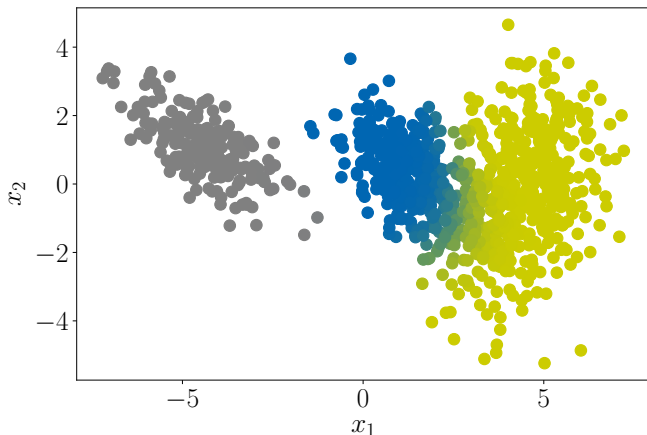
Example

Example (2)



- Negative log likelihood never increases

Visualizing the Responsibilities



- Soft assignments of data points between obvious clusters

Overview

Gaussian Mixture Models

Parameter Learning

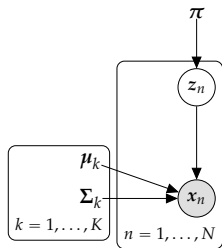
Implementation

Probabilistic Perspective

Probabilistic Perspective

- ▶ We will be very explicit about **how data is generated** for a given set of model parameters
- ▶ **Justification** of some ad-hoc choices we made earlier (e.g., definition of responsibilities)
- ▶ **Interpretation** of some model parameters as prior/posterior probabilities
- ▶ Can be used for a **principled derivation of the EM algorithm** (which generally allows for maximum likelihood estimation in latent variable models)
 - ▶▶ Not covered here. See Bishop (2006) for more details

Probabilistic Perspective on one Slide



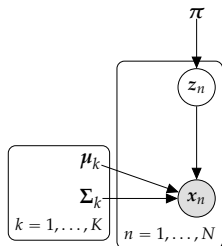
$$p(z_{nk} = 1) = \pi_k, \quad 0 \leq \pi_k \leq 1, \quad \sum_{k=1}^K \pi_k = 1$$

$$p(\mathbf{x}_n | z_{nk} = 1) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$p(\mathbf{x}_n) = \sum_{k=1}^K p(z_{nk} = 1) p(\mathbf{x}_n | z_{nk} = 1)$$

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

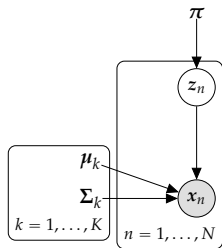
Probabilistic Perspective on one Slide



$$p(z_{nk} = 1) = \pi_k, \quad 0 \leq \pi_k \leq 1, \quad \sum_{k=1}^K \pi_k = 1$$
$$p(\mathbf{x}_n | z_{nk} = 1) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$
$$p(\mathbf{x}_n) = \sum_{k=1}^K p(z_{nk} = 1) p(\mathbf{x}_n | z_{nk} = 1)$$
$$= \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- ▶ $\mathbf{z}_n = (z_{n1}, \dots, z_{nK})$ is a discrete latent variable. Exactly one entry of \mathbf{z}_n is 1, all others are 0 **▶▶ 1-of-K code/One-hot encoding**
- ▶ For every data point \mathbf{x}_n there is a corresponding latent variable \mathbf{z}_n that **indicates which mixture component generated \mathbf{x}_n**

Probabilistic Perspective on one Slide



$$p(z_{nk} = 1) = \pi_k, \quad 0 \leq \pi_k \leq 1, \quad \sum_{k=1}^K \pi_k = 1$$
$$p(\mathbf{x}_n | z_{nk} = 1) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$
$$p(\mathbf{x}_n) = \sum_{k=1}^K p(z_{nk} = 1) p(\mathbf{x}_n | z_{nk} = 1)$$
$$= \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- ▶ $\mathbf{z}_n = (z_{n1}, \dots, z_{nK})$ is a discrete latent variable. Exactly one entry of \mathbf{z}_n is 1, all others are 0 **▶▶ 1-of-K code/One-hot encoding**
- ▶ For every data point \mathbf{x}_n there is a corresponding latent variable \mathbf{z}_n that **indicates which mixture component generated \mathbf{x}_n**
- ▶ Posterior $p(z_k = 1 | \mathbf{x}_i) = r_{ik}$ corresponds to the “responsibility” (see earlier) that mixture component k generated data point i .

Prior $p(\mathbf{z})$

- ▶ $\pi_k = p(z_k = 1)$ is the (prior) probability that the k th mixture component generates a data point \mathbf{x}
- ▶ $\boldsymbol{\pi} := [\pi_1, \dots, \pi_K]^\top$, $\sum_k \pi_k = 1$
- ▶ $p(\mathbf{z}) = \boldsymbol{\pi}$ is a probability vector of length K

Generative Process

Ancestral sampling from a GMM (generative process) is simple:

$z^{(i)} \sim p(z)$ Select mixture component

$x_i \sim p(x|z^{(i)} = 1)$ Draw sample from this component

Discard sampled $z^{(i)}$ and end up with valid data samples x_i from the GMM

Likelihood $p(\mathbf{x}|\boldsymbol{\theta})$

- ▶ $\boldsymbol{\theta} := \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k : k = 1, \dots, K\}$ contains all model parameters
- ▶ Likelihood $p(\mathbf{x}|\boldsymbol{\theta})$ does not depend on latent variables
 - ▶▶ Marginalize out latent variable \mathbf{z} :

$$\begin{aligned} p(\mathbf{x}|\boldsymbol{\theta}) &= \sum_{\mathbf{z}} p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{z})p(\mathbf{z}|\boldsymbol{\theta}) \\ &= \sum_{k=1}^K p(\mathbf{x}|\boldsymbol{\theta}, z_k = 1)p(z_k = 1|\boldsymbol{\theta}) \\ &= \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \end{aligned}$$

Likelihood $p(\mathbf{x}|\boldsymbol{\theta})$

- ▶ $\boldsymbol{\theta} := \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k : k = 1, \dots, K\}$ contains all model parameters
- ▶ Likelihood $p(\mathbf{x}|\boldsymbol{\theta})$ does not depend on latent variables
 - ▶▶ Marginalize out latent variable \mathbf{z} :

$$\begin{aligned} p(\mathbf{x}|\boldsymbol{\theta}) &= \sum_{\mathbf{z}} p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{z})p(\mathbf{z}|\boldsymbol{\theta}) \\ &= \sum_{k=1}^K p(\mathbf{x}|\boldsymbol{\theta}, z_k = 1)p(z_k = 1|\boldsymbol{\theta}) \\ &= \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \end{aligned}$$

- ▶▶ Classical GMM likelihood function

Posterior $p(\mathbf{z}|\mathbf{x})$

- ▶ Bayes' theorem:

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}{p(\mathbf{x})}$$

Posterior $p(\mathbf{z}|\mathbf{x})$

- ▶ Bayes' theorem:

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}{p(\mathbf{x})}$$

- ▶ Posterior distribution

Posterior $p(\mathbf{z}|\mathbf{x})$

- ▶ Bayes' theorem:

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}{p(\mathbf{x})}$$

- ▶ Posterior distribution

Posterior $p(\mathbf{z}|\mathbf{x})$

- ▶ Bayes' theorem:

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}{p(\mathbf{x})}$$

- ▶ Posterior distribution

$$p(z_k = 1 | \mathbf{x}) = \frac{p(\mathbf{x} | z_k = 1)p(z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x} | z_j = 1)} = \frac{\pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

Posterior $p(\mathbf{z}|\mathbf{x})$

- ▶ Bayes' theorem:

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}{p(\mathbf{x})}$$

- ▶ Posterior distribution

$$p(z_k = 1 | \mathbf{x}) = \frac{p(\mathbf{x} | z_k = 1)p(z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x} | z_j = 1)} = \frac{\pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

- ▶▶ “Responsibility” of the k th mixture component for data point \mathbf{x}
- ▶▶ Posterior probability that the k th mixture component generated data point \mathbf{x}

Some Comments

- ▶ If we make “hard assignments” of data points to clusters (responsibilities r_{ik} or posteriors are either 0 or 1), we get *K*-means clustering

Some Comments

- ▶ If we make “hard assignments” of data points to clusters (responsibilities r_{ik} or posteriors are either 0 or 1), we get **K-means clustering**
- ▶ **Overfitting** as with other maximum likelihood approaches (what does overfitting in GMMs look like?)

Some Comments

- ▶ If we make “hard assignments” of data points to clusters (responsibilities r_{ik} or posteriors are either 0 or 1), we get **K-means clustering**
- ▶ **Overfitting** as with other maximum likelihood approaches (what does overfitting in GMMs look like?)
- ▶ **Bayesian extensions** exist with priors on the parameters

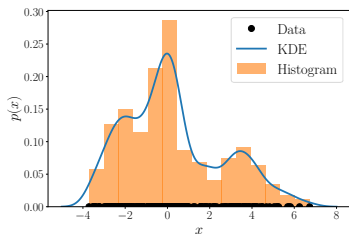
Some Comments

- ▶ If we make “hard assignments” of data points to clusters (responsibilities r_{ik} or posteriors are either 0 or 1), we get **K-means clustering**
- ▶ **Overfitting** as with other maximum likelihood approaches (what does overfitting in GMMs look like?)
- ▶ **Bayesian extensions** exist with priors on the parameters
- ▶ Choosing the number of components can be done via model selection or a Bayesian prior (Dirichlet process, Görür (2007))

Some Comments

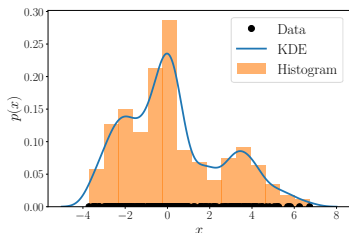
- ▶ If we make “hard assignments” of data points to clusters (responsibilities r_{ik} or posteriors are either 0 or 1), we get **K-means clustering**
- ▶ **Overfitting** as with other maximum likelihood approaches (what does overfitting in GMMs look like?)
- ▶ **Bayesian extensions** exist with priors on the parameters
- ▶ Choosing the number of components can be done via model selection or a Bayesian prior (Dirichlet process, Görür (2007))
- ▶ EM can generally be used for **parameter learning in latent variable models** (e.g., Ghahramani & Roweis (1999)) or reinforcement learning (e.g., Barber (2012))
 - ▶▶ Latent variable perspective is useful to derive EM in a principled way

Other Density Estimation Methods



- ▶ Histograms (Pearson1895)

Other Density Estimation Methods

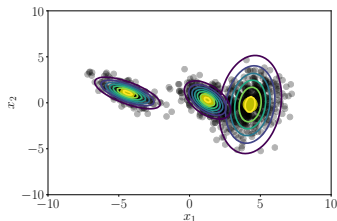
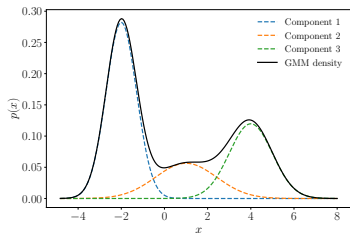


- ▶ Histograms (Pearson 1895)
- ▶ Kernel density estimation (Rosenblatt 1956)

$$p(\mathbf{x}) = \frac{1}{Nh} \sum_{n=1}^N k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right), \quad h > 0,$$

$$k(\cdot) \geq 0, \quad \int k(\mathbf{x}) d\mathbf{x} = 1$$

Summary



- ▶ Density estimation with Gaussian mixture models
- ▶ No closed-form solution to maximum likelihood estimation
- ▶ EM algorithm for an iterative solution
- ▶ Latent variable perspective

References I

- [1] D. Barber. *Bayesian Reasoning and Machine Learning*. Cambridge University Press, 2012.
- [2] C. M. Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer-Verlag, 2006.
- [3] Z. Ghahramani and S. T. Roweis. Learning Nonlinear Dynamical Systems using an EM Algorithm. In M. S. Kearns, S. A. Solla, and D. A. Cohn, editors, *Advances in Neural Information Processing Systems*, volume 11. The MIT Press, 1999.
- [4] D. Görür. *Nonparametric Bayesian Discrete Latent Variable Models for Unsupervised Learning*. PhD thesis, Technische Universität Berlin, 2007.
- [5] K. P. Murphy. *Machine Learning: A Probabilistic Perspective*. MIT Press, Cambridge, MA, USA, 2012.
- [6] K. Pearson. Contributions to the Mathematical Theory of Evolution. II. Skew Variation in Homogeneous Material. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 186:343–414, 1895.
- [7] M. Rosenblatt. Remarks on Some Nonparametric Estimates of a Density Function. *The Annals of Mathematical Statistics*, 27(3):832–837, 1956.