Machine Learning CS 4641



Gaussian Mixture Model

Nakul Gopalan Georgia Tech

Some of the slides are based on slides from Jiawei Han Chao Zhang, Mahdi Roozbahani and Barnabás Póczos.

Outline

- Overview -
- Gaussian Mixture Model
- The Expectation-Maximization Algorithm

Recap

Conditional probabilities:

$$p(A,B) = p(A|B)p(B) = p(B|A)p(A)$$

Bayes rule:

$$p(A|B) = \frac{p(A,B)}{p(B)} = \frac{p(B|A)p(A)}{p(B)}$$

 $p(A = 1) = \sum_{i=1}^{K} p(A = 1, B_i) = \sum_{i=1}^{K} p(A|B_i) p(B_i)$

| | Tomorrow=Rainy | Tomorrow=Cold | P(Today) |
|-------------|-------------------|-------------------|-------------------|
| Today=Rainy | 4/9 | 2/9 | [4/9 + 2/9] = 2/3 |
| Today=Cold | 2/9 | 1/9 | [2/9 + 1/9] = 1/3 |
| P(Tomorrow) | [4/9 + 2/9] = 2/3 | [2/9 + 1/9] = 1/3 | |

P(Tomorrow = Rainy) =

Hard Clustering Can Be Difficult

• Hard Clustering: K-Means, Hierarchical Clustering, DBSCAN



Towards Soft Clustering

• K-means

-hard assignment: each object belongs to only one cluster

$$\theta_i \in \{\theta_1, \ldots, \theta_K\}$$

Mixture modeling

-soft assignment: probability that an object belongs to a cluster



Outline

- Overview
- Gaussian Mixture Model
- The Expectation-Maximization Algorithm

Gaussian Distribution



What is a Gaussian?

For *d* dimensions, the Gaussian distribution of a vector $x = (x^1, x^2, ..., x^d)^T$ is defined by:

$$N(x \mid \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)$$

where μ is the mean and Σ is the covariance matrix of the Gaussian.



Mixture Models

• Formally a Mixture Model is the weighted sum of a number of pdfs where the weights are determined by a distribution, π

$$p(x) = \pi_0 f_0(x) + \pi_1 f_1(x) + \pi_2 f_2(x) + \ldots + \pi_k f_k(x)$$

where
$$\sum_{i=0}^k \pi_i = 1$$
$$p(x) = \sum_{i=0}^k \pi_i f_i(x)$$
$$\downarrow$$
$$What is f in GMM?$$

$$p(x) = \pi_0 f_0(x) + \pi_1 f_1(x) + \pi_2 f_2(x)$$



Why p(x) is a pdf?

$$p(x) = \pi_0 f_0(x) + \pi_1 f_1(x) + \pi_2 f_2(x) + \ldots + \pi_k f_k(x)$$

where
$$\sum_{i=0}^k \pi_i = 1$$

Why GMM?

It creates a new pdf for us to generate random variables. It is a generative model.

It clusters different components using a Gaussian distribution.

So it provides us the inferring opportunity. Soft assignment!!

Some notes:

Is summation of a bunch of Gaussians a Gaussian itself?

p(x) is a Probability density function or it is also called a marginal distribution function.

p(x) = the density of selecting a data point from the pdf which is created from a mixture model. Also, we know that the area under a density function is equal to 1.

Mixture Models are Generative

• Generative simply means dealing with joint probability p(x,z)

$$p(x) = \pi_0 f_0(x) + \pi_1 f_1(x) + \dots + \pi_k f_k(x)$$

Let's say f(.) is a Gaussian distribution

 $p(\mathbf{x}) = \pi_0 N(X|\mu_0, \sigma_0) + \pi_1 N(X|\mu_1, \sigma_1) + \dots + \pi_k N(X|\mu_k, \sigma_k)$

$$p(x) = \sum_{k} N(x|\mu_k, \sigma_k) \pi_k$$

$$p(x) = \sum_{k} p(x|z_k) p(z_k)$$

 z_k is component k

$$p(x) = \sum_{k} p(x, z_k)$$

 $P(x) = \sum_{i=1}^{n} \frac{F_{i}(x)}{F_{i}(x)}$ Test \sim K T T; ~1 5 unkonour 3 P. (22) - sis onland $P(n) \rightarrow \mathcal{N}(\mathcal{M}; \overline{z};)$ UNENOWN

GMM with graphical model concept



What is soft assignment?



What is the probability of a datapoint x in each component?

How many components we have here? 3

How many probability distributions? 3

What is the sum value of the 3 probabilities for each datapoint?

How to calculate the probability of datapoints in the first component (inferring)?

$$p(\mathbf{x}) = \pi_0 N(X|\mu_0, \sigma_0) + \pi_1 N(X|\mu_1, \sigma_1) + \pi_2 N(X|\mu_2, \sigma_2)$$

Let's calculate the responsibility of the first component among the rest for one point x

Let's call that au_0

Given a datapoint X, what is probability of that datapoint in component 0 If I have 100 datapoints and 3 components, what is the size of τ ? 1002

Inferring Cluster Membership

- We have representations of the joint $p(x, z_{nk}|\theta)$ and the marginal, $p(x|\theta)$
- The conditional of $p(z_{nk}|x,\theta)$ can be derived using Bayes rule.
 - The responsibility that a mixture component takes for explaining an observation x.

$$\tau(z_k) = p(z_k = 1|x) = \frac{p(z_k = 1)p(x|z_k = 1)}{\sum_{j=1}^{K} p(z_j = 1)p(x|z_j = 1)}$$

Some K^T
$$= \frac{\pi_k N(x|\mu_k, \Sigma_k)}{\sum_{j=1}^{K} \pi_j N(x|\mu_j, \Sigma_j)}$$

Mixtures of Gaussians



What are GMM parameters?

Mean μ_k Variance σ_k Size π_k Marginal probability distribution $p(x|\theta) = \sum_{k} p(x, z_{k}|\theta) = \sum_{k} p(x|z_{k}, \theta) p(z_{k}|\theta) = \sum_{k} N(x|\mu_{k}, \sigma_{k})\pi_{k}$ $f_{k}(x)$ π_{k} $f_{k}(x)$ $f_{k}(x)$ $p(z_k|\theta) = \pi_k$ Select a mixture component with probability π $p(x|z_k,\theta) = N(x|\mu_k,\sigma_k)$ Sample from that component's Gaussian

 π_1

 π_0

 χ

 π_2

How about GMM for multimodal distribution?

- What if we know the data consists of a few Gaussians
- What if we want to fit parametric models



Gaussian Mixture Model

 A density model p(X) may be multi-modal: model it as a mixture of uni-modal distributions (e.g. Gaussians)



Why having "Latent variable"

- A variable can be unobserved (latent) because:
 - it is an imaginary quantity meant to provide some simplified and abstractive view of the data generation process.
 - e.g., speech recognition models, mixture models (soft clustering)...
 - it is a real-world object and/or phenomena, but difficult or impossible to measure
 - e.g., the temperature of a star, causes of a disease, evolutionary ancestors ...
 - it is a real-world object and/or phenomena, but sometimes wasn't measured, because of faulty sensors, etc.
 - Discrete latent variables can be used to partition/cluster data into sub-groups.
 - Continuous latent variables (factors) can be used for dimensionality reduction (factor analysis, etc).

Latent variable representation

$$p(\mathbf{x}|\theta) = \sum_{k} p(x, z_{nk}|\theta) = \sum_{k} p(z_{nk}|\theta) p(x|z_{nk}, \theta) = \sum_{k=0}^{K} \pi_{k} N(x|\mu_{k}, \Sigma_{k})$$

$$\underset{k=0}{\overset{\text{Max gindize'}}{\overset{\text{Max gindize'}}{\overset{Max gindize'}}{\overset{Max gindize'}}{\overset{Max gindize'}}{\overset{Max gindize'}}{\overset{Max gindize'}}{\overset{Max gindize'}}{\overset{Max gindize'}}{\overset{Max$$

The distribution that we can model using a mixture of Gaussian components is much more expressive than what we could have modeled using a single component.

2010

Well, we don't know π_k, μ_k, Σ_k What should we do?

We use a method called "Maximum Likelihood Estimation" (MLE) to solve the problem.

$$p(\mathbf{x}) = p(\mathbf{x}|\theta) = \sum_{k} p(x, z_k|\theta) = \sum_{k} p(z_k|\theta) p(x|z_k, \theta) = \sum_{k=0}^{K} \pi_k N(x|\mu_k, \Sigma_k)$$

Let's identify a likelihood function, why?

Because we use likelihood function to optimize the probabilistic model parameters!

$$\arg\max p(x|\theta) = p(x|\pi,\mu,\Sigma) = \prod_{n=1}^{N} p(x_n|\theta) = \prod_{n=1}^{N} \sum_{k=0}^{K} \pi_k N(x_n|\mu_k,\Sigma_k)$$

• And set partials to zero...



Maximum Likelihood of a GMM

Optimization of covariance

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

$$\Sigma_k = \frac{1}{\sum_{n=1}^N \tau(z_{nk})} \sum_{n=1}^N \tau(z_{nk}) (x_n - \mu_k) (x_n - \mu_k)^T$$

Maximum Likelihood of a GMM



MLE of a GMM



$$\sum_{k} = \frac{1}{N_k} \sum_{n=1}^{N} \tau(z_{nk}) (x_n - \mu_k) (x_n - \mu_k)^T$$

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

$$N_k = \sum_{n=1}^N \tau(z_{nk})$$

Not a closed form solution!! τ is not known exactly What next?

Outline

- Overview
- Gaussian Mixture Model
- The Expectation-Maximization Algorithm

EM for GMMs

• E-step: Evaluate the Responsibilities

$$\tau(z_{nk}) = \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \Sigma_j)}$$

EM for GMMs

M-Step: Re-estimate Parameters

$$\mu_k^{new} = \frac{\sum_{n=1}^N \tau(z_{nk}) x_n}{N_k}$$
$$\Sigma_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \tau(z_{nk}) (x_n - \mu_k^{new}) (x_n - \mu_k^{new})^T$$

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

Expectation Maximization

- Expectation Maximization (EM) is a general algorithm to deal with hidden variables.
- Two steps:
 - 。 E-Step: Fill-in hidden values using inference
 - ^o M-Step: Apply standard MLE method to estimate parameters
- EM always converges to a local minimum of the likelihood.



EM for Gaussian Mixture Model:



After 1st iteration



After 2nd iteration



After 3rd iteration



After 4th iteration



After 5th iteration



After 6th iteration



After 20th iteration



Demo

 Demo link: https://lukapopijac.github.io/gaussian-mixturemodel/

EM Algorithm for GMM (matrix form)

Given a Gaussian mixture model, the goal is to maximize the likelihood function with respect to the parameters comprising the means and covariances of the components and the mixing coefficients).

- 1. Initialize the means μ_{j} , covariances \sum_{j} and mixing coefficients π_{j} , and evaluate the initial value of the log likelihood.
- 2. E step. Evaluate the responsibilities using the current parameter values

$$\underline{\tau(z_{nk})} = \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \Sigma_j)}$$

Book : C.M. Bishop, Pattern Recognition and Machine Learning, Springer, 2006

EM for GMMs

M-Step: Re-estimate Parameters

$$\mu_k^{new} = \frac{\sum_{n=1}^N \tau(z_{nk}) x_n}{N_k}$$
$$\Sigma_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \tau(z_{nk}) (x_n - \mu_k^{new}) (x_n - \mu_k^{new})^T$$
$$\pi_k^{new} = \frac{N_k}{N}$$

EM Algorithm for GMM (matrix form)

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

$$\mu_k^{new} = \frac{\sum_{n=1}^N \tau(z_{nk}) x_n}{N_k} \left(\sum_{k=1}^{new} = \frac{1}{N_k} \sum_{n=1}^N \tau(z_{nk}) (x_n - \mu_k^{new}) (x_n - \mu_k^{new})^T \right) \left(\pi_k^{new} = \frac{N_k}{N} \right)$$

4. Evaluate log likelihood

$$\ln \mathbf{p}(\mathbf{X} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \mathbf{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

If there is no convergence, return to step 2.

Relationship to K-means

- K-means makes hard decisions.
 - 。 Each data point gets assigned to a single cluster.
- GMM/EM makes soft decisions.
 - $_{\circ}$ Each data point can yield a posterior p(z|x)
- K-means is a special case of EM.

General form of EM

- Given a joint distribution over observed and latent variables: $p(X, Z|\theta)$ • Want to maximize: $p(X|\theta)$ 1. Initialize parameters: θ^{old} 2. E Step: Evaluate: $p(Z|X, \theta^{old})$
- 3. M-Step: Re-estimate parameters (based on expectation of completedata log likelihood)

$$\theta^{new} = \operatorname{argmax}_{\theta} \sum_{Z} p(Z|X, \theta^{old}) \ln p(X, Z|\theta) = \underset{P(Z|X, \theta^{2k_d})}{\operatorname{argmax}_{\theta}} \mathbb{E}[\ln(p(x, Z|\theta))]$$

4. Check for convergence of params or likelihood

EM improves loglikelihood in both steps



$$\theta^{new} = \operatorname{argmax}_{\theta} \sum_{Z} p(Z|X, \theta^{old}) \ln p(X, Z|\theta)$$

Only true if F(IE(n)) Jensen's inequality $\boldsymbol{\ell}(\boldsymbol{\theta};\boldsymbol{X}) = \log \boldsymbol{p}(\boldsymbol{X} \mid \boldsymbol{\theta})$ $= \log \sum p(\boldsymbol{x}, \boldsymbol{z} \mid \boldsymbol{\theta})$ $\| F(F(x)) \|$ $= \log \sum_{z} q(z \mid x) \frac{p(x, z \mid \theta)}{q(z \mid x)}$ Will lead to maximizing this $\geq \sum_{z} q(z \mid x) \log \frac{p(x, z \mid \theta)}{q(z \mid x)}$ Maximizing this

$$F(q,\theta) = \sum_{z} q(z \mid x) \log \frac{p(x,z \mid \theta)}{q(z \mid x)}$$
$$= \sum_{z} q(z \mid x) \log p(x,z \mid \theta) - \sum_{z} q(z \mid x) \log q(z \mid x)$$
$$= \left\langle \ell_{c}(\theta;x,z) \right\rangle_{q} + H_{q} \right\rangle$$

The first term is the expected complete log likelihood and the second term, which does not depend on θ , is the entropy.

Thus, in the M-step, maximizing with respect to θ for fixed q we only need to consider the first term:

$$\theta^{t+1} = \arg \max_{\theta} \left\langle \ell_{c}(\theta; \boldsymbol{x}, \boldsymbol{z}) \right\rangle_{q^{t+1}} = \arg \max_{\theta} \sum_{\theta} q(\boldsymbol{z} \mid \boldsymbol{x}) \log \boldsymbol{p}(\boldsymbol{x}, \boldsymbol{z} \mid \theta)$$

covariance_type="diag" or "spherical" or "full"



Source: Python Data Science Handbook by Jake VanderPlas



Silhouette Coefficient

Define the silhoutte coefficient of a point \mathbf{x}_i as

$$\boldsymbol{s}_{i} = \frac{\mu_{out}^{\min}(\mathbf{x}_{i}) - \mu_{in}(\mathbf{x}_{i})}{\max\left\{\mu_{out}^{\min}(\mathbf{x}_{i}), \mu_{in}(\mathbf{x}_{i})\right\}}$$

where $\mu_{in}(\mathbf{x}_i)$ is the mean distance from \mathbf{x}_i to points in its own cluster \hat{y}_i :

$$\mu_{in}(\mathbf{x}_i) = \frac{\sum_{\mathbf{x}_j \in C_{\hat{y}_i}, j \neq i} \delta(\mathbf{x}_i, \mathbf{x}_j)}{n_{\hat{y}_i} - 1}$$

and $\mu_{out}^{\min}(\mathbf{x}_i)$ is the mean of the distances from \mathbf{x}_i to points in the closest cluster:

$$\mu_{out}^{\min}(\mathbf{x}_i) = \min_{j \neq \hat{y}_i} \left\{ \frac{\sum_{\mathbf{y} \in C_j} \delta(\mathbf{x}_i, \mathbf{y})}{n_j} \right\}$$

The Silhouette Coefficient for clustering C: $SC = \frac{1}{n} \sum_{i=1}^{n} s_i$.

SC close to 1 implies a good clustering (Points are close to their own clusters but far from other clusters)

Take-Home Messages

- The generative process of Gaussian Mixture Model
- Inferring cluster membership based on a learned GMM
- The general idea of Expectation-Maximization
- Expectation-Maximization for GMM
- Silhouette Coefficient