

# Introduction to Machine Learning

Mixture Models

Varun Chandola

April 24, 2019

## Outline

## Contents

<b>1 Latent Variable Models</b>	<b>1</b>
1.1 Latent Variable Models - Introduction . . . . .	2
<b>2 Mixture Models</b>	<b>2</b>
2.1 Using Mixture Models . . . . .	3
<b>3 Parameter Estimation</b>	<b>3</b>
3.1 Issues with Direct Optimization of the Likelihood or Posterior	5
<b>4 Expectation Maximization</b>	<b>5</b>
4.1 EM Operation . . . . .	8
4.2 EM for Mixture Models . . . . .	9
4.3 K-Means as EM . . . . .	10

## 1 Latent Variable Models

- Consider a probability distribution parameterized by  $\theta$
- Generates samples ( $\mathbf{x}$ ) with probability  $p(\mathbf{x}|\theta)$

### 2-step generative process

1. Distribution generates the hidden variable
2. Distribution generates the observation, given the hidden variable

### Magazine Example - Sampling an Article

- Assume that the editor has access to  $p(\mathbf{x})$
- $\mathbf{x}$  - a random variable that denotes an article

### Direct Model

- Sample from  $p(\mathbf{x})$  for an article

### Latent Variable Model

1. First sample a topic  $z$  from a topic distribution  $p(z)$
2. Pick an article from the topic-wise distribution  $p(\mathbf{x}|z)$

### 1.1 Latent Variable Models - Introduction

- The observed random variable  $\mathbf{x}$  depends on a hidden random variable  $\mathbf{z}$
- $\mathbf{z}$  is generated using a *prior* distribution -  $p(\mathbf{z})$
- $\mathbf{x}$  is generated using  $p(\mathbf{x}|\mathbf{z})$
- Different combinations of  $p(\mathbf{z})$  and  $p(\mathbf{x}|\mathbf{z})$  give different latent variable models
  1. **Mixture Models**
  2. Factor analysis
  3. Probabilistic Principal Component Analysis (PCA)
  4. Latent Dirichlet Allocation (LDA)

## 2 Mixture Models

- A latent discrete state

$$z \in \{1, 2, \dots, K\}$$

- $p(z) \sim \text{Multinomial}(\boldsymbol{\pi})$
- For every state  $k$ , we have a probability distribution for  $\mathbf{x}$

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

- Overall, probability for  $\mathbf{x}$

$$p(\mathbf{x}|\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}|\boldsymbol{\theta})$$

- A **convex combination** of  $p_k$ 's
- $\pi_k$  is the probability of  $k^{\text{th}}$  mixture component to be true
  - Or, contribution of the  $k^{\text{th}}$  component
  - Or, the mixing weight

### 2.1 Using Mixture Models

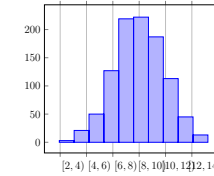
#### 1. Black-box Density Model

- Use  $p(\mathbf{x}|\boldsymbol{\theta})$  for many things
- Example: *class conditional density*

#### 2. Clustering

- *Soft clustering*
  1. First learn the parameters of the mixture model
    - Each mixture component corresponds to a cluster  $k$
  2. Compute  $p(z = k|\mathbf{x}, \boldsymbol{\theta})$  for every input point  $\mathbf{x}$  (*Bayes Rule*)

$$p(z = k|\mathbf{x}, \boldsymbol{\theta}) = \frac{p(z = k|\boldsymbol{\theta})p(\mathbf{x}|z = k, \boldsymbol{\theta})}{\sum_{k'=1}^K p(z = k'|\boldsymbol{\theta})p(\mathbf{x}|z = k', \boldsymbol{\theta})}$$



## 3 Parameter Estimation

### Simple Parameter Estimation

- **Given:** A set of scalar observations

$$x_1, x_2, \dots, x_n$$

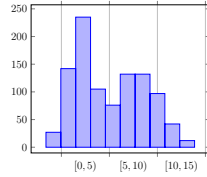
- **Task:** Find the generative model (form and parameters)

1. Observe empirical distribution of  $x$
2. Make choice of the *form* of the probability distribution (**Gaussian**)
3. Estimate parameters from the data using MLE or MAP ( **$\mu$  and  $\sigma$** )

In the above example we choose the random variable  $x$  to be distributed as a Gaussian random variable.

### When Data has Multiple Modes

- Single mode is not sufficient
- In reality data is generated from two Gaussians
- How to estimate  $\mu_1, \sigma_1, \mu_2, \sigma_2$ ?
- What if we knew  $z_i \in \{1, 2\}$ ?
  - $z_i = 1$  means that  $x_i$  comes from first mixture component
  - $z_i = 2$  means that  $x_i$  comes from second mixture component



- **Issue:**  $z_i$ 's are not known beforehand
- Need to explore  $2^N$  possibilities

Obviously, if  $z_i$ 's were known, you can create two data sets corresponding to the two values that  $z_i$  can take, and then estimate parameters for each set.

### 3.1 Issues with Direct Optimization of the Likelihood or Posterior

- For direct optimization, we find parameters that maximize (log-)likelihood (or (log-)posterior)
- Easy to optimize if  $z_i$  were all known
- What happens when  $z_i$ 's are not known
  - Likelihood and posterior will have multiple modes
  - Non-convex function - harder to optimize

## 4 Expectation Maximization

- Recall the we want to maximize the log-likelihood of a data set with respect to  $\theta$ :

$$\hat{\theta} = \underset{\theta}{\text{maximize}} \ell(\theta)$$

- Log-likelihood for a mixture model can be written as:

$$\begin{aligned} \ell(\theta) &= \sum_{i=1}^N \log p(\mathbf{x}_i | \theta) \\ &= \sum_{i=1}^N \log \left[ \sum_{k=1}^K p(z_k) p_k(\mathbf{x}_i | \theta) \right] \end{aligned}$$

- Hard to optimize (a summation inside the log term)

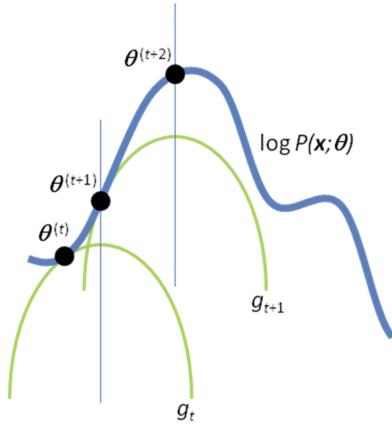
Note that the above equation for log-likelihood is for mixture models only. In general, it maybe written as:

$$\begin{aligned} \ell(\theta) &= \sum_{i=1}^N \log p(\mathbf{x}_i | \theta) \\ &= \sum_{i=1}^N \log \left[ \sum_{\forall \mathbf{z}_k} p(\mathbf{x}_i, \mathbf{z}_k | \theta) \right] \\ &= \sum_{i=1}^N \log \left[ \sum_{\forall \mathbf{z}_k} p(\mathbf{z}_k) p(\mathbf{x}_i | \mathbf{z}_k, \theta) \right] \end{aligned}$$

- Repeat until converged:
  1. Start with some guess for  $\theta$  and compute the most likely value for  $z_i, \forall i$
  2. Given  $z_i, \forall i$ , update  $\theta$
- Does not explicitly maximize the log-likelihood of mixture model
- Can we come up with a better algorithm?
  - Repeat until converged:
    1. Start with some guess for  $\theta$  and compute the probability of  $z_i = k, \forall i, k$
    2. Combine probabilities to update  $\theta$

## Expectation Maximization Algorithm

- A principled approach to maximize a function with latent variables
- At iteration  $t$ , for a given value of  $\theta^{(t)}$ , let  $Q$  be a convex function that is a lower bound of  $l(\theta)$



**Supplementary Figure 1** Convergence of the EM algorithm. Starting from initial parameters  $\theta^{(t)}$ , the E-step of the EM algorithm constructs a function  $g_t$  that lower-bounds the objective function  $\log P(x; \theta)$ . In the M-step,  $\theta^{(t+1)}$  is computed as the maximum of  $g_t$ . In the next E-step, a new lower-bound  $g_{t+1}$  is constructed; maximization of  $g_{t+1}$  in the next M-step gives  $\theta^{(t+2)}$ , etc.

## Steps in EM

- EM is an iterative procedure
- Start with some value for  $\theta$
- At every iteration  $t$ , update  $\theta$  such that the log-likelihood of the data goes up

- Move from  $\theta^{t-1}$  to  $\theta$  such that:

$$\ell(\theta) - \ell(\theta^{t-1})$$

is maximized

- **Complete log-likelihood** for any LVM

$$\ell(\theta) = \sum_{i=1}^N \log p(\mathbf{x}_i, \mathbf{z}_i | \theta)$$

- Cannot be computed as we do not know  $\mathbf{z}_i$

## Expected complete log-likelihood

$$Q(\theta, \theta^{t-1}) = \mathbb{E}[\ell(\theta | D, \theta^{t-1})]$$

- Expected value of  $\ell(\theta | D, \theta^{t-1})$  for all possibilities of  $\mathbf{z}_i$

Recall that expected value of a function  $f(x)$  for a random variable  $x$  is given by:

$$\mathbb{E}[f(x)] = \sum_{x'} f(x') p(x')$$

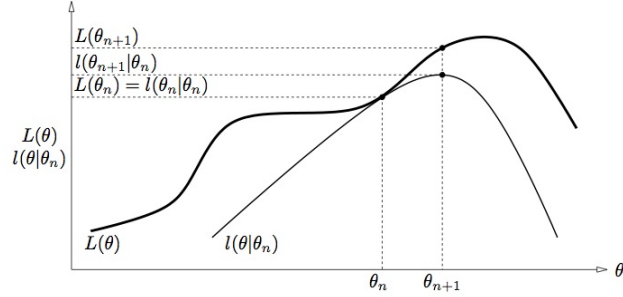
If  $x$  is continuous, the sum is replaced by an integral

$$\mathbb{E}[f(x)] = \int f(x) p(x) dx$$

In the case of EM, we are interested in computing the expected value over all possibilities of  $\mathbf{z}$ . The probability of each possibility is computed using the current estimate  $\theta^{t-1}$ .

## 4.1 EM Operation

1. Initialize  $\theta$
2. At iteration  $t$ , compute  $Q(\theta, \theta^{t-1})$
3. Maximize  $Q(\cdot)$  with respect to  $\theta$  to get  $\theta^t$
4. Goto step 2



## 4.2 EM for Mixture Models

- EM formulation is generic
- Calculating (E) and maximizing (M)  $Q()$  needs to be done for specific instances

$Q$  for MM

$$\begin{aligned}
 Q(\theta, \theta^{t-1}) &= \mathbb{E} \left[ \sum_{i=1}^N \log p(\mathbf{x}_i, z_i | \theta) \right] \\
 &= \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log \pi_k + \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log p(\mathbf{x}_i | \theta_k) \\
 r_{ik} &\triangleq p(z_i = k | \mathbf{x}_i, \theta^{t-1})
 \end{aligned}$$

The quantity  $r_{ik}$  can be thought of as the responsibility that cluster  $k$  takes for data point  $\mathbf{x}_i$  in iteration  $t$ .

**E-Step**

- Compute  $r_{ik}, \forall i, k$

$$\begin{aligned}
 r_{ik} &= p(z_i = k | \mathbf{x}_i, \theta^{t-1}) \\
 &= \frac{\pi_k p(\mathbf{x}_i | \theta_k^{t-1})}{\sum_{k'} \pi_{k'} p(\mathbf{x}_i | \theta_{k'}^{t-1})}
 \end{aligned}$$

- Compute  $Q()$

**M-Step**

- Maximize  $Q()$  w.r.t.  $\theta$
- $\theta$  consists of  $\pi = \{\pi_1, \pi_2, \dots, \pi_K\}$  and  $\theta = \{\theta_1, \theta_2, \dots, \theta_K\}$
- For Gaussian Mixture Model (GMM) ( $\theta_k \equiv (\mu_k, \Sigma_k)$ ):

$$\pi_k = \frac{1}{N} \sum_i r_{ik} \quad (1)$$

$$\mu_k = \frac{\sum_i r_{ik} \mathbf{x}_i}{\sum_i r_{ik}} \quad (2)$$

$$\Sigma_k = \frac{\sum_i r_{ik} \mathbf{x}_i \mathbf{x}_i^\top}{\sum_i r_{ik}} - \mu_k \mu_k^\top \quad (3)$$

## 4.3 K-Means as EM

- Similar to GMM
  1.  $\Sigma = \sigma^2 \mathbf{I}_D$
  2.  $\pi_k = \frac{1}{K}$
  3. The most probable cluster for  $\mathbf{x}_i$  is computed as the prototype closest to it (hard clustering)

## References