

Disclaimer

- The material provided in this document is not my original work and is a summary of some one else's work(s).
- A simple Google search of the title of the document will direct you to the original source of the material.
- I do not guarantee the accuracy, completeness, timeliness, validity, non-omission, merchantability or fitness of the contents of this document for any particular purpose.
- Downloaded from najeebkhan.github.io

THE EM ALGORITHM

Presented by

Najeeb

July 14, 2014

Introduction

- Given the true densities for all classes, one may determine the optimum Bayes classifier

$$\arg \max_k P(w_k | x) = \arg \max_k P(x | w_k) P(w_k)$$

- Given parameters trained to maximize the likelihood for the data in each class, we can approximate the Bayes classifier
- However, in many problems, there is no analytical solution for the parameters of the statistical classifier given the training data

Introduction

- An example of such a problem is finding the means, variances, and mixture weights for a sum of Gaussian densities, given a sample of data points
- We will postulate a model in which each individual Gaussian generates some of the observed data points
- Given that there is no way to directly compute estimates for these quantities, we need an optimization technique to find the parameters for each Gaussian that will maximize the likelihood of all of the observed data

Introduction

- The dominant approach to such problems is called the Expectation Maximization (EM) algorithm
- The parameter estimation problem is structured to incorporate variables representing information that is not directly observed, but that is assumed to be part of the model that generated the data
- Such a variable is often called hidden or missing
- In the Gaussian mixture case, a hidden variable could be the index of the Gaussian that generated a data point

Introduction

- The key idea of EM is to
 - Estimate the densities by taking an expectation of the logarithm of the joint density between the known and unknown components
 - Maximize this function by updating the parameters that are used in the probability estimation

Expectation

- The log likelihood can be expressed in terms of hidden variables as

$$\ell(X | \Theta) = \log P(X | \Theta)$$

$$P(k, X | \Theta) = \frac{P(k, X, \Theta)}{P(\Theta)} = \frac{P(k | X, \Theta)P(X, \Theta)}{P(\Theta)}$$

$$P(k, X | \Theta) = P(k | X, \Theta)P(X | \Theta)$$

$$P(X | \Theta) = \frac{P(k, X | \Theta)}{P(k | X, \Theta)}$$

$$\log P(X | \Theta) = \log P(k, X | \Theta) - \log P(k | X, \Theta)$$

Expectation

$$\log P(X | \Theta) = \log P(k, X | \Theta) - \log P(k | X, \Theta)$$

$$\sum_{k=1}^K P(k | X, \Theta^i) \log P(X | \Theta) =$$

$$\sum_{k=1}^K P(k | X, \Theta^i) \log P(k, X | \Theta) - P(k | X, \Theta^i) \log P(k | X, \Theta)$$

$$\log P(X | \Theta) \sum_{k=1}^K P(k | X, \Theta^i) = Q(\Theta, \Theta^i) - H(\Theta, \Theta^i)$$

$$Q(\Theta, \Theta^i) = \sum_{k=1}^K P(k | X, \Theta^i) \log P(k, X | \Theta)$$

$$H(\Theta, \Theta^i) = \sum_{k=1}^K P(k | X, \Theta^i) \log P(k | X, \Theta)$$

Q function

- If the Q function is maximized the likelihood is also maximized

$$\begin{aligned} Q(\Theta, \Theta^i) &= \sum_{k=1}^K P(k | X, \Theta^i) \log P(k, X | \Theta) \\ &= \sum_{k=1}^K P(k | X, \Theta^i) \log [P(k | X, \Theta) P(X | \Theta)] \\ &= \sum_{k=1}^K P(k | X, \Theta^i) \log P(k | X, \Theta) + \log P(X | \Theta) \sum_{k=1}^K P(k | X, \Theta^i) \\ &= \sum_{k=1}^K P(k | X, \Theta^i) \log P(k | X, \Theta) + \log P(X | \Theta) \end{aligned}$$

Q function

$$Q(\Theta, \Theta^i) = \sum_{k=1}^K P(k | X, \Theta^i) \log P(k | X, \Theta) + \log P(X | \Theta) \quad (1)$$

$$Q(\Theta^i, \Theta^i) = \sum_{k=1}^K P(k | X, \Theta^i) \log P(k | X, \Theta^i) + \log P(X | \Theta^i) \quad (2)$$

$$\log P(X | \Theta) - \log P(X | \Theta^i) = Q(\Theta, \Theta^i) - Q(\Theta^i, \Theta^i) + NNC$$

Q function

$$Q(\Theta, \Theta^i) = \sum_{k=1}^K P(k | X, \Theta^i) \log P(k | X, \Theta) + \log P(X | \Theta) \quad (1)$$

$$Q(\Theta^i, \Theta^i) = \sum_{k=1}^K P(k | X, \Theta^i) \log P(k | X, \Theta^i) + \log P(X | \Theta^i) \quad (2)$$

$$\log P(X | \Theta) - \log P(X | \Theta^i) = Q(\Theta, \Theta^i) - Q(\Theta^i, \Theta^i) + NNC$$

- If a change to Θ increases Q , $\log P(x | \Theta)$ increases

Q function

- In principle, we could maximize the expectation for each value of Θ and then re-estimate Θ
- In many cases, it is possible to structure the problem so that we can analytically determine the choice for the parameters that will maximize the expectation in each iteration
- Then a new expression for the expectation can be determined, followed by a new parameter estimation, and so on

EM For Mixture Gaussian Density Estimation

- We can always decompose an unknown probability density $P(X | \Theta)$ as

$$P(X | \Theta) = \sum_{k=1}^K P(X, k | \Theta) = \sum_{k=1}^K P(k | \Theta) P(X | k, \Theta)$$

- We can write an expression for the log joint density for observed and hidden variables

$$Z = \log P(x, k | \Theta) = \log[P(k | \Theta) P(x | k, \Theta)]$$

- With an expected value over N samples and K mixture components of

$$Q = \sum_{k=1}^K \sum_{n=1}^N P(k | x_n, \Theta^i) \log[P(k | \Theta) P(x_n | k, \Theta)]$$

- Θ^i is the old parameter

EM For Mixture Gaussian Density Estimation

- Using properties of the log

$$Q = \sum_{k=1}^K \sum_{n=1}^N P(k | x_n, \Theta^i) \log P(k | \Theta) + \sum_{k=1}^K \sum_{n=1}^N P(k | x_n, \Theta^i) \log P(x_n | k, \Theta)$$

- We will often assume models for which the Θ parameters of the two terms are disjointed and can be optimized separately

EM For Mixture Gaussian Density Estimation

- Assuming that each component $P(x|k)$ is Gaussian

$$Q = \sum_{k=1}^K \sum_{n=1}^N P(k | x_n, \Theta^i) \log P(k | \Theta) + \sum_{k=1}^K \sum_{n=1}^N P(k | x_n, \Theta^i) \left[-\log \sigma_k - \frac{(x_n - \mu_k)^2}{2\sigma_k^2} + C \right]$$

- Given this choice of parametric form for the mixture density, we can use standard optimization methods to find the best value for the parameters

EM For Mixture Gaussian Density Estimation

- Solving for the means

$$\frac{\partial Q}{\partial \mu_j} = 0$$

$$\sum_{n=1}^N P(j | x_n, \Theta^i) \left(\frac{x_n}{\sigma_j^2} - \frac{\mu_j}{\sigma_j^2} \right) = 0$$

$$\sum_{n=1}^N P(j | x_n, \Theta^i) x_n = \sum_{n=1}^N P(j | x_n, \Theta^i) \mu_j$$

$$\mu_j = \frac{\sum_{n=1}^N P(j | x_n, \Theta^i) x_n}{\sum_{n=1}^N P(j | x_n, \Theta^i)}$$

EM For Mixture Gaussian Density Estimation

- Solving for the means

$$\frac{\partial Q}{\partial \mu_j} = 0$$

$$\sum_{n=1}^N P(j | x_n, \Theta^i) \left(\frac{x_n}{\sigma_j^2} - \frac{\mu_j}{\sigma_j^2} \right) = 0$$

$$\sum_{n=1}^N P(j | x_n, \Theta^i) x_n = \sum_{n=1}^N P(j | x_n, \Theta^i) \mu_j$$

$$\mu_j = \frac{\sum_{n=1}^N P(j | x_n, \Theta^i) x_n}{\sum_{n=1}^N P(j | x_n, \Theta^i)}$$

EM For Mixture Gaussian Density Estimation

- Solving for the variances

$$\frac{\partial Q}{\partial \sigma_j} = 0$$

$$\sum_{n=1}^N P(j | x_n, \Theta^i) \left(-\frac{1}{\sigma_j} + \frac{(x_n - \mu_j)^2}{\sigma_j^3} \right) = 0$$

$$\frac{1}{\sigma_j} \sum_{n=1}^N P(j | x_n, \Theta^i) = \frac{1}{\sigma_j^3} \sum_{n=1}^N P(j | x_n, \Theta^i) (x_n - \mu_j)^2$$

$$\sigma_j^2 = \frac{\sum_{n=1}^N P(j | x_n, \Theta^i) (x_n - \mu_j)^2}{\sum_{n=1}^N P(j | x_n, \Theta^i)}$$

EM For Mixture Gaussian Density Estimation

- Solving for the Mixture Weights

$$Q^* = \sum_{k=1}^K \sum_{n=1}^N P(k | x_n, \Theta^i) \log P(k | \Theta) +$$
$$\sum_{k=1}^K \sum_{n=1}^N P(k | x_n, \Theta^i) \left[-\log \sigma_k - \frac{(x_n - \mu_k)^2}{2\sigma_k^2} + C \right] +$$
$$\lambda \left[\sum_{k=1}^K P(k | \Theta) - 1 \right]$$

EM For Mixture Gaussian Density Estimation

- Solving for the Mixture Weights

$$\frac{\partial Q^*}{\partial P(j | \Theta)} = 0$$

$$\frac{1}{P(j | \Theta)} \sum_{n=1}^N P(j | x_n, \Theta^i) + \lambda = 0$$

$$P(j | \Theta) = \frac{1}{\lambda} \sum_{n=1}^N P(j | x_n, \Theta^i)$$

$$P(j | \Theta) = \frac{1}{N} \sum_{n=1}^N P(j | x_n, \Theta^i)$$

EM For Mixture Gaussian Density Estimation

- Solving for the Mixture Weights

$$P(j | x_n, \Theta^i) = \frac{P(x_n | j, \Theta^i) P(j | \Theta^i)}{P(x_n | \Theta^i)}$$

$$P(j | x_n, \Theta^i) = \frac{P(x_n | j, \Theta^i) P(j | \Theta^i)}{\sum_{k=1}^K P(x_n | k, \Theta^i) P(k | \Theta^i)}$$

HMM Training

$$Q = \sum_{k=1}^K \sum_{n=1}^N P(k | x_n, \Theta^i) \log[P(k | \Theta)P(x_n | k, \Theta)]$$

- For the case of HMM we use the state sequence as the hidden variable

$$Q = \sum_Q P(Q | X_1^N, \Theta_{old}, M) \log \left[P(X_1^N | Q, \Theta, M) P(Q | \Theta, M) \right]$$

$$\begin{aligned} Q = & \sum_{n=1}^N \sum_{k=1}^{L(M)} P(q_k^n | X_1^N, \Theta_{old}, M) \log P(x_n | q_k^n, \Theta, M) \\ & + \sum_{k=1}^{L(M)} P(q_k^1 | X_1^N, \Theta_{old}, M) \log P(x_n | q_k^1, \Theta, M) \\ & + \sum_{n=2}^N \sum_{k=1}^{L(M)} \sum_{l=1}^{L(M)} P(q_l^n, q_k^{n-1} | X_1^N, \Theta_{old}, M) \log P(q_l^n | q_k^{n-1}, \Theta, M) \end{aligned}$$

HMM Training

- Assuming single univariate Gaussian is associated with each hmm state

- $$\mu_j = \frac{\sum_{n=1}^N P(q_j^n | X_1^N, \Theta_{old}, M) x_n}{\sum_{n=1}^N P(q_j^n | X_1^N, \Theta^i)}$$

$$\sigma_j^2 = \frac{\sum_{n=1}^N P(q_j^n | X_1^N, \Theta_{old}, M) (x_n - \mu_j)^2}{\sum_{n=1}^N P(q_j^n | X_1^N, \Theta^i)}$$