

Expectation Maximisation

The EM algorithm is a widely used approach for learning in the presence of unobserved variables, provided the general probability distribution governing these variables is known.

Estimating Means of k Gaussians

The problem basically is that there is a set of data D that originates from k distinct Gaussian distributions. Each instance is generated by first selecting a Gaussian at random and then a single instance x is selected from that distribution. The variance of the Gaussians is known, and the task is to output a hypothesis that describes the means of each of the k distributions.

For a single distribution, this mean is given by

$$mean_{ML} = \underset{mean}{argmin} \sum_{i=1}^m (x_i - \mu)^2$$

But since the data comes from a mixture of Gaussians this is not valid.

We define each instance as a triplet $\langle x, z_1, z_2 \dots z_N \rangle$ where x is the observed value of an instance and z_1 indicates the normal distribution that was used to generate the value x .

If z_i has the value 1, the value x was generated from the i th Gaussian.

The EM algorithm searches for a maximum likelihood hypothesis by repeatedly re-estimating the values of the hidden variables z given its current hypothesis, then recalculating the maximum likelihood hypothesis using these expected values for the hidden variables.

The algorithm is given as

1. Calculate the expected value $E[z_i]$ of each hidden variable z , assuming the current hypothesis $\langle \text{mean}_1, \text{mean}_2 \dots \rangle$ holds.
2. Calculate a new maximum likelihood hypothesis h' assuming the value taken by each variable z is the expected value.
3. Replace h by h' and re-iterate.

The expected value $E[z_i]$ is calculated as

$$E[z_{ij}] = \frac{p(x = x_i | \mu = \mu_j)}{\sum_{n=1}^k p(x = x_i | \mu = \mu_n)}$$

$$E[z_{ij}] = \frac{e^{-\frac{1}{2\sigma^2}(x - \mu_i)^2}}{\sum_{i=1}^k e^{-\frac{1}{2\sigma^2}(x_i - \mu_n)^2}}$$

The maximum likelihood hypothesis is a collection of all the means calculated using this expected value, where each mean is given by

$$\mu_j = \frac{\sum_{i=1}^m E[z_{ij}]x_i}{\sum_{i=1}^m E[z_{ij}]}$$

The General Statement of the EM algorithm

1. Estimation step E

$$Q(h' | h) = E[\ln P(Y | h') | h, X]$$

2. Maximization step M

$$h = \underset{h'}{\operatorname{argmax}} Q(h' | h)$$

Where

θ = Set of parameters = <mean1, mean2...>

θ' = Set of revised parameters

X = observed data in a set of m independently drawn instances

Z = unobserved data in these same instances

Y = random variable defined in terms of Z

h = Initial hypothesis

h' = Revised hypothesis

When the function Q is continuous, the EM algorithm converges to a stationary point of the likelihood function P . When this likelihood function has a single maximum, EM converges to this global maximum likelihood for h' . Otherwise, it guarantees only a local maximum.