

Document clustering:

Derivation of Expectation and Maximization:

The training algorithm of EM is to find the maximum likelihood.

$$\ln P(X|\theta) = \ln \sum P(X, Z|\theta)$$

For computing likelihood, we will need Z , now we will calculate posterior probability, $P(Z|X, \theta)$.

The Expectation is:

$$Q(\theta, \theta^{\text{old}}) = \sum_z P(z|X, \theta^{\text{old}}) \cdot \ln P(X, z|\theta)$$

For Generative Model:

- toss K -face(dice) with ϕ to choose the face K that n^{th} document belongs to.
- For each placeholder:
 - generate the word by tossing dice with parameter μ_K , corresponding to K .

The parameters of model are:

- cluster proportion ϕ where $\sum_{k=1}^K \phi_k = 1$
- word proportion μ_k for cluster k , where $\sum_{k=1}^K \mu_k = 1$

The probability of generating a pair k & cluster (k, d) :

$$\begin{aligned} p(k, d) &= p(k) p(d|k) \\ &= \phi_k \cdot \prod_{w \in d} \mu_{k,w}^{c(w,d)} \end{aligned}$$

where $c(w, d)$ is no. of occurrences.

For incomplete data, the document cluster is not given, the latent variable Z_n .

\therefore The probability of document is,

$$\begin{aligned} p(d_1, d_2 \dots d_N) &= \sum_{n=1}^N \ln \sum_{k=1}^K p(Z_n, k=1, d_n) \\ &= \sum_{n=1}^N \ln \sum_{k=1}^K (\phi_k \prod_{w \in d_n} \mu_{k,w}^{c(w,d_n)}) \end{aligned}$$

The expectation function, $Q(\theta, \theta^{\text{old}})$:

$$\begin{aligned} Q(\theta, \theta^{\text{old}}) &= \sum_{n=1}^N \sum_{k=1}^K p(Z_n, k=1 | d_n, \theta^{\text{old}}) \ln p(Z_n, k=1, d_n) \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma(Z_n, k) (\ln \phi_k + \sum_{w \in A} c(w, d_n) \ln \mu_{k,w}) \end{aligned}$$

where $\gamma(Z_n, k) = p(Z_n, k=1 | d_n, \theta^{\text{old}})$ represents the responsibility factor, assuming all words came from a dictionary A .

The probability that a doc. belongs to a cluster, gamma,

$$\gamma(z_n, k) = \frac{\prod_{n=1}^N \phi_k \prod_{w \in A} \mu_{k,w}^{c(w,d)}}{\prod_{n=1}^N \sum_{k=1}^K (\phi_k \prod_{w \in A} \mu_{k,w}^{c(w,d)})}$$

To maximize the Q function,

→ Mixing components : $\phi_k = N_k / N$ where $N_k = \sum_{n=1}^N \gamma(z_n, k)$

→ The word proportion parameter:

$$\mu_{k,w} = \frac{\sum_{n=1}^N \gamma(z_n, k) c(w, d_n)}{\sum_{w' \in A} \sum_{n=1}^N \gamma(z_n, k) c(w', d_n)}$$

Hence, soft EM algo for document clustering is:

- choose an initial setting : $\theta^{\text{old}} = (\phi^{\text{old}}, \mu_1^{\text{old}}, \dots, \mu_K^{\text{old}})$
- while convergence is not met:
 - E step : set $\forall n, \forall k : \gamma(z_n, k)$ based on θ^{old}
 - M step : set θ^{new} based on above equations.
 - $\theta^{\text{old}} \leftarrow \theta^{\text{new}}$

Hard EM algo for document clustering is:

- choose an initial setting for θ^{old}
- while convergence not met:
 - E step : set $z_n \leftarrow \arg\max_z P(z|x, \theta^{\text{old}})$
 - M step : set $\theta^{\text{new}} \leftarrow \arg\max_{\theta} P(x, z^* | \theta)$
 - $\theta^{\text{old}} \leftarrow \theta^{\text{new}}$