

Lecture 8: Graphical models for Text

4F13: Machine Learning

Joaquin Quiñonero-Candela and Carl Edward Rasmussen

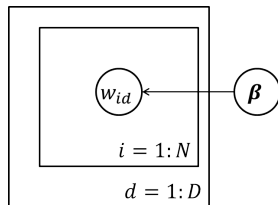
Department of Engineering
University of Cambridge

<http://mlg.eng.cam.ac.uk/teaching/4f13/>

A really simple document model

Consider a collection of D documents with dictionary of M unique words.

- N_d : number of (non-unique) words in document d .
- w_{id} : i -th word in document d ($w_{id} \in \{1 : M\}$).
- $\beta = [\beta_1, \dots, \beta_M]^T$: parameters of a Multinomial distribution over the dictionary of M unique words.



We can fit β by maximising the likelihood:

$$\begin{aligned}\hat{\beta} &= \operatorname{argmax} \prod_{d=1}^D \operatorname{Mult}(c_{1d}, \dots, c_{Md} | \beta, N_d) \\ &= \operatorname{argmax} \operatorname{Mult}(c_1, \dots, c_M | \beta, N)\end{aligned}$$

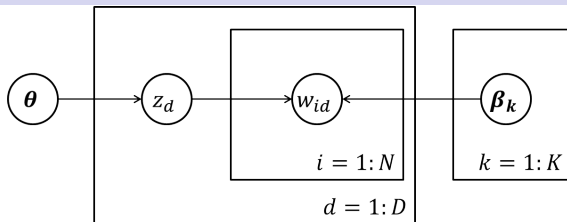
$$\hat{\beta}_j = \frac{c_j}{\sum_{l=1}^M c_l}$$

- $N = \sum_{d=1}^D N_d$: total number of (non-unique) words in the collection.
- c_{jd} : count of occurrences of unique word j in document d .
- $c_j = \sum_{d=1}^D c_{jd}$: count of total occurrences of unique word j in the collection.

Limitations of the really simple document model

- Document d is the result of sampling N_d words from the Multinomial β .
- β estimated by maximum likelihood reflects the aggregation of all documents.
- All documents are therefore modelled by the global word frequency distribution.
- The generative model wastes mass, because it cannot specialize.
- All unique words do not necessarily *co-occur* in a given document.
- It possible that documents might be about different *topics*.

A mixture of Multinomials model



We want to allow for a mixture of K Multinomials parametrised by β_1, \dots, β_K . Each of those Multinomials corresponds to a *document category*.

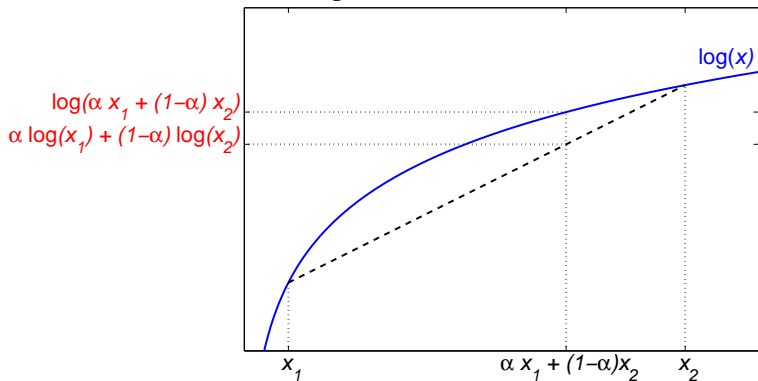
- $z_d \in \{1 : K\}$ assigns all words in document d to one of the K categories.
- $\theta_j = p(z_d = j)$ is the probability any document d is assigned to category j .
- $\theta = [\theta_1, \dots, \theta_K]$ is also the parameter of a Multinomial over the K categories.

We have introduced a new set of *hidden* variables z_d .

- How do we fit those variables? What do we do with them?
- We are actually not interested in them: We are only interested in θ and β .

Jensen's Inequality

For any concave function, such as $\log(x)$



For $\alpha_i \geq 0$, $\sum \alpha_i = 1$ and any $\{x_i > 0\}$

$$\log\left(\sum_i \alpha_i x_i\right) \geq \sum_i \alpha_i \log(x_i)$$

Equality if and only if $\alpha_i = 1$ for some i (and therefore all others are 0).

The Expectation Maximization (EM) algorithm

Given a set of observed (visible) variables V , a set of unobserved (hidden / latent / missing) variables H , and model parameters θ , optimize the log likelihood:

$$\mathcal{L}(\theta) = \log p(V|\theta) = \log \int p(H, V|\theta) dH, \quad (1)$$

where we have written the marginal for the visibles in terms of an integral over the joint distribution for hidden and visible variables.

Using *Jensen's inequality* for **any** distribution of hidden states $q(H)$ we have:

$$\mathcal{L} = \log \int q(H) \frac{p(H, V|\theta)}{q(H)} dH \geq \int q(H) \log \frac{p(H, V|\theta)}{q(H)} dH = \mathcal{F}(q, \theta), \quad (2)$$

defining the $\mathcal{F}(q, \theta)$ functional, which is a **lower bound** on the log likelihood.

In the EM algorithm, we alternately optimize $\mathcal{F}(q, \theta)$ wrt q and θ , and we can prove that this will never decrease \mathcal{L} .

The E and M steps of EM

The lower bound on the log likelihood:

$$\mathcal{F}(\mathbf{q}, \theta) = \int \mathbf{q}(\mathbf{H}) \log \frac{p(\mathbf{H}, \mathbf{V}|\theta)}{\mathbf{q}(\mathbf{H})} d\mathbf{H} = \int \mathbf{q}(\mathbf{H}) \log p(\mathbf{H}, \mathbf{V}|\theta) d\mathbf{H} + \mathcal{H}(\mathbf{q}), \quad (3)$$

where $\mathcal{H}(\mathbf{q}) = - \int \mathbf{q}(\mathbf{H}) \log \mathbf{q}(\mathbf{H}) d\mathbf{H}$ is the **entropy** of \mathbf{q} . We iteratively alternate:

E step: optimize $\mathcal{F}(\mathbf{q}, \theta)$ wrt the distribution over hidden variables given the parameters:

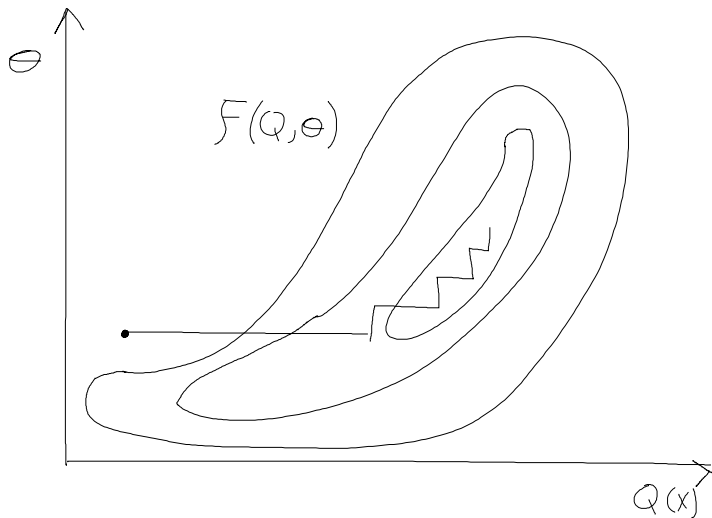
$$\mathbf{q}^{(k)}(\mathbf{H}) := \operatorname{argmax}_{\mathbf{q}(\mathbf{H})} \mathcal{F}(\mathbf{q}(\mathbf{H}), \theta^{(k-1)}). \quad (4)$$

M step: maximize $\mathcal{F}(\mathbf{q}, \theta)$ wrt the parameters given the hidden distribution:

$$\theta^{(k)} := \operatorname{argmax}_{\theta} \mathcal{F}(\mathbf{q}^{(k)}(\mathbf{H}), \theta) = \operatorname{argmax}_{\theta} \int \mathbf{q}^{(k)}(\mathbf{H}) \log p(\mathbf{H}, \mathbf{V}|\theta) d\mathbf{H}, \quad (5)$$

which is equivalent to optimizing the expected complete-data likelihood $p(\mathbf{H}, \mathbf{V}|\theta)$, since the **entropy of $\mathbf{q}(\mathbf{H})$** does not depend on θ .

EM as Coordinate Ascent in \mathcal{F}



The EM algorithm never decreases the log likelihood

The difference between the cost functions:

$$\begin{aligned}\mathcal{L}(\theta) - \mathcal{F}(q, \theta) &= \log p(V|\theta) - \int q(H) \log \frac{p(H, V|\theta)}{q(H)} dH \\ &= \log p(V|\theta) - \int q(H) \log \frac{p(H|V, \theta)p(V|\theta)}{q(H)} dH \\ &= - \int q(H) \log \frac{p(H|V, \theta)}{q(H)} dH = \mathcal{KL}(q(H), p(H|V, \theta)),\end{aligned}$$

is called the Kullback-Liebler divergence; it is non-negative and only zero if and only if $q(H) = p(H|V, \theta)$ (thus this is the E step). Although we are working with the **wrong cost function**, the likelihood is still increased in every iteration:

$$\mathcal{L}(\theta^{(k-1)}) \underset{\text{E step}}{=} \mathcal{F}(q^{(k)}, \theta^{(k-1)}) \underset{\text{M step}}{\leq} \mathcal{F}(q^{(k)}, \theta^{(k)}) \underset{\text{Jensen}}{\leq} \mathcal{L}(\theta^{(k)}),$$

where the first equality holds because of the E step, and the first inequality comes from the M step and the final inequality from Jensen. Usually EM converges to a local optimum of \mathcal{L} (although there are exceptions).

EM and Mixtures of Multinomials

The the Mixture model for text, the latent variables are

$$z_d \in \{1, \dots, K\}, \text{ where } d = 1, \dots, D$$

which for each document encodes which mixture component generated it.

E-step: for each document d , set q to the posterior

$$q_d(z_d) \propto p(z_d = k|\theta) \prod_{i=1}^{N_d} p(w_i|\beta_{w_i k}) = \theta_k \text{Mult}(c_{1d}, \dots, c_{Md}|\beta_k, N_d) = r_{kd}$$

M-step: Maximize

$$\begin{aligned} \sum_{k=1}^K q_d(z_d = k) \log p(\{w_{id}\}, z_d) &= \sum_k r_{kd} \log \prod_{d=1}^D \prod_{i=1}^{N_d} p(w_i|\beta_{w_i k}) p(z_d = k) \\ &= \sum_k r_{kd} \left(\sum_{d=1}^D \log \prod_{j=1}^M \beta_{jk}^{c_{jd}} + \log \theta_k \right) \\ &= \sum_{k,d} r_{kd} \left(\sum_{j=1}^M c_{jd} \log \beta_{jk} + \log \theta_k \right) = F(\theta, \beta) \end{aligned}$$

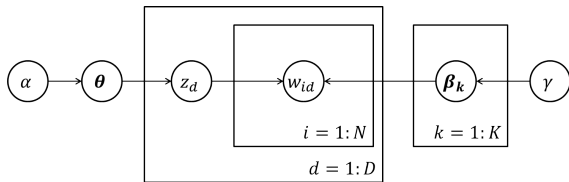
EM: M step for mixture model

Need Lagrange multipliers to constrain the maximization and ensure proper distributions.

$$\begin{aligned}\theta_k &= \operatorname{argmax} F(\theta, \beta) + \lambda(1 - \sum_{k=1}^K \theta_k) \\ &= \frac{\sum_{d=1}^D r_{kd}}{\sum_{k'=1}^K \sum_{d=1}^D r_{k'd}}\end{aligned}$$

$$\begin{aligned}\beta_{jk} &= \operatorname{argmax} F(\theta, \beta) + \sum_{j=k}^K \lambda_k(1 - \sum_{j=1}^M \beta_{jk}) \\ &= \frac{\sum_{d=1}^D r_{kd} c_{jd}}{\sum_{j'=1}^M \sum_{d=1}^D r_{kd} c_{j'd}}\end{aligned}$$

A Bayesian mixture of Multinomials model



With the EM algorithm we have essentially estimated α and β by maximum likelihood. An alternative, Bayesian treatment is to introduce *hyperpriors*.

- $\theta \sim \text{Dir}(\alpha)$ is a symmetric Dirichlet over category probabilities.
- $\beta_k \sim \text{Dir}(\gamma)$ is a symmetric Dirichlet over unique word probabilities.

What is different?

- We no longer want to compute a point estimate of θ or β .
- We are now interested in computing the *posterior* distributions.