

Lecture 12: Graphical models for Text

Machine Learning 4F13, Spring 2014

Carl Edward Rasmussen and Zoubin Ghahramani

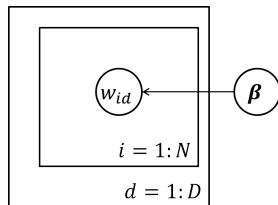
CUED

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

A really simple document model

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

- N_d : number of (non-unique) words in document d .
- w_{id} : i -th word in document d ($w_{id} \in \{1 \dots M\}$).
- $w_{id} \sim \text{Cat}(\boldsymbol{\beta})$: each word is drawn from a discrete categorical distribution with parameters $\boldsymbol{\beta}$
- $\boldsymbol{\beta} = [\beta_1, \dots, \beta_M]^\top$: parameters of a categorical / multinomial distribution¹ over the dictionary of M unique words.

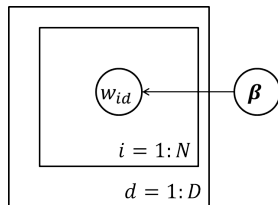


¹It's a categorical distribution if we observe the sequence of words in the document, it's a multinomial if we only observe the counts.

A really simple document model

Modelling D documents from a dictionary of M unique words.

- N_d : number of (non-unique) words in document d .
- w_{id} : i -th word in document d ($w_{id} \in \{1 \dots M\}$).
- $w_{id} \sim \text{Cat}(\boldsymbol{\beta})$: each word is drawn from a discrete categorical distribution with parameters $\boldsymbol{\beta}$



We can fit $\boldsymbol{\beta}$ by maximising the likelihood:

$$\hat{\boldsymbol{\beta}} = \operatorname{argmax}_{\boldsymbol{\beta}} \prod_{d=1}^D \text{Mult}(c_{1d}, \dots, c_{Md} | \boldsymbol{\beta}, N_d)$$

$$= \operatorname{argmax}_{\boldsymbol{\beta}} \text{Mult}(c_1, \dots, c_M | \boldsymbol{\beta}, N)$$

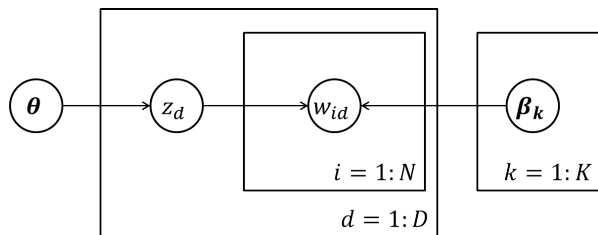
$$\hat{\beta}_j = \frac{c_j}{N} = \frac{c_j}{\sum_{l=1}^M c_l}$$

- $N = \sum_{d=1}^D N_d$: total number of (non-unique) words in the collection.
- $c_{jd} = \sum_{i=1}^{N_d} \mathbb{I}(w_{id} = j)$: count 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.
- This generative model does not specialise.
- It possible that different documents might be about different *topics*.

A mixture of multinomials model



$$z_d \sim \text{Cat}(\theta)$$
$$w_{id}|z_d \sim \text{Cat}(\beta_{z_d})$$

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*.

- $z_d \in \{1, \dots, K\}$ assigns document d to one of the K categories.
- $\theta_k = p(z_d = k)$ is the probability any document d is assigned to category k .
- so $\theta = [\theta_1, \dots, \theta_K]$ is the parameter of a multinomial over 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?
- Are these variables interesting? Or are we only interested in θ and β ?

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:

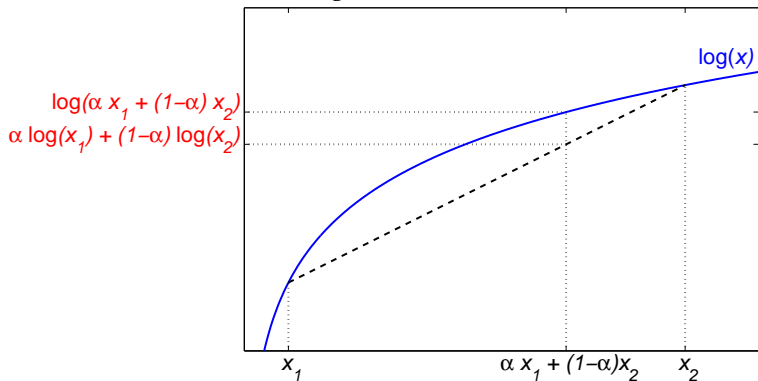
$$\mathcal{L}(\theta) = \log \int q(H) \frac{p(H, V|\theta)}{q(H)} dH \geq \int q(H) \log \frac{p(H, V|\theta)}{q(H)} dH = \mathcal{F}(q, \theta), \quad (2)$$

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 $\mathcal{L}(\theta)$.

Jensen's Inequality

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



For $\alpha_i \geq 0$, $\sum_i \alpha_i = 1$ 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 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: maximize $\mathcal{F}(\mathbf{q}, \theta)$ wrt the distribution over hidden variables given the parameters:

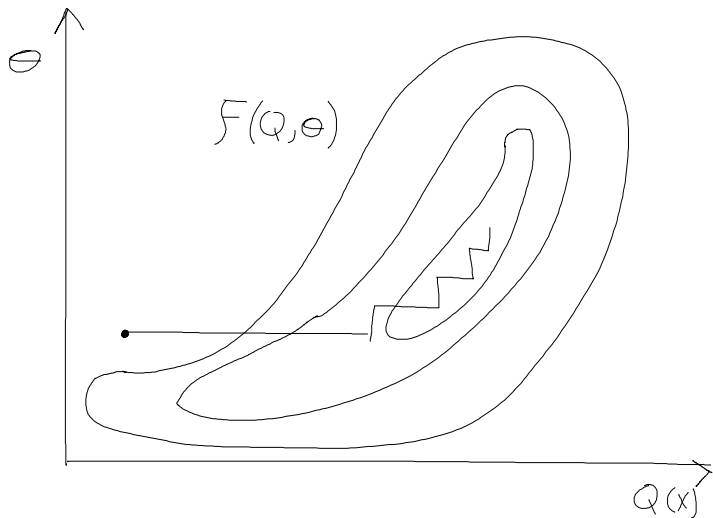
$$\mathbf{q}^{(k)}(\mathbf{H}) := \operatorname{argmax}_{\mathbf{q}(\mathbf{H})} \mathcal{F}(\mathbf{q}(\mathbf{H}), \theta^{(k-1)}). \quad (4)$$

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 objective 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 zero if and only if $q(H) = p(H|V, \theta)$ (thus this is the E step). Although we are optimising a **lower bound**, \mathcal{F} , the likelihood \mathcal{L} 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

In 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

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

M-step: Maximize

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

EM: M step for mixture model

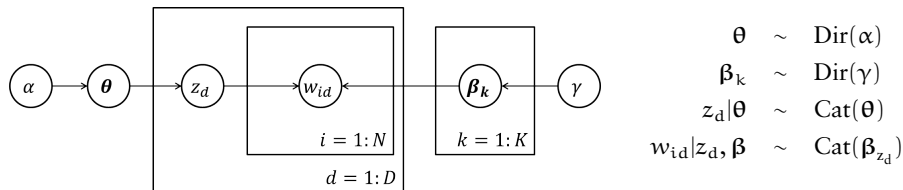
$$F(\mathbf{R}, \theta, \beta) = \sum_{k,d} r_{kd} \left(\sum_{j=1}^M c_{jd} \log \beta_{jk} + \log \theta_k \right)$$

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

$$\begin{aligned} \hat{\theta}_k &\leftarrow \operatorname{argmax}_{\theta_k} F(\mathbf{R}, \theta, \beta) + \lambda \left(1 - \sum_{k'=1}^K \theta_{k'} \right) \\ &= \frac{\sum_{d=1}^D r_{kd}}{\sum_{k'=1}^K \sum_{d=1}^D r_{k'd}} = \frac{\sum_{d=1}^D r_{kd}}{D} \end{aligned}$$

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

A Bayesian mixture of Multinomials model



With the EM algorithm we have essentially estimated θ and β by maximum likelihood. An alternative, Bayesian treatment infers the parameters starting from priors:

- $\theta \sim \text{Dir}(\alpha)$ 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.

Variational Bayesian Learning

Lower Bounding the Marginal Likelihood

Let the hidden latent variables be \mathbf{x} , data \mathbf{y} and the parameters θ .

Lower bound the **marginal likelihood (Bayesian model evidence)** using Jensen's inequality:

$$\begin{aligned}\log P(\mathbf{y}) &= \log \int d\mathbf{x} d\theta P(\mathbf{y}, \mathbf{x}, \theta) && |m \\ &= \log \int d\mathbf{x} d\theta Q(\mathbf{x}, \theta) \frac{P(\mathbf{y}, \mathbf{x}, \theta)}{Q(\mathbf{x}, \theta)} \\ &\geq \int d\mathbf{x} d\theta Q(\mathbf{x}, \theta) \log \frac{P(\mathbf{y}, \mathbf{x}, \theta)}{Q(\mathbf{x}, \theta)}.\end{aligned}$$

Use a simpler, factorised approximation to $Q(\mathbf{x}, \theta)$:

$$\begin{aligned}\log P(\mathbf{y}) &\geq \int d\mathbf{x} d\theta Q_{\mathbf{x}}(\mathbf{x}) Q_{\theta}(\theta) \log \frac{P(\mathbf{y}, \mathbf{x}, \theta)}{Q_{\mathbf{x}}(\mathbf{x}) Q_{\theta}(\theta)} \\ &= \mathcal{F}(Q_{\mathbf{x}}(\mathbf{x}), Q_{\theta}(\theta), \mathbf{y}).\end{aligned}$$

Maximize this lower bound.

Variational Bayesian Learning ...

Maximizing this **lower bound**, \mathcal{F} , leads to **EM-like** updates:

$$Q_x^*(\mathbf{x}) \propto \exp \langle \log P(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) \rangle_{Q_\theta(\boldsymbol{\theta})} \quad E\text{-like step}$$

$$Q_\theta^*(\boldsymbol{\theta}) \propto P(\boldsymbol{\theta}) \exp \langle \log P(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) \rangle_{Q_x(\mathbf{x})} \quad M\text{-like step}$$

Maximizing \mathcal{F} is equivalent to minimizing KL-divergence between the *approximate posterior*, $Q(\boldsymbol{\theta})Q(\mathbf{x})$ and the *true posterior*, $P(\boldsymbol{\theta}, \mathbf{x} | \mathbf{y})$.

$$\begin{aligned} \log P(\mathbf{y}) - \mathcal{F}(Q_x(\mathbf{x}), Q_\theta(\boldsymbol{\theta}), \mathbf{y}) &= \\ \log P(\mathbf{y}) - \int d\mathbf{x} d\boldsymbol{\theta} Q_x(\mathbf{x}) Q_\theta(\boldsymbol{\theta}) \log \frac{P(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta})}{Q_x(\mathbf{x}) Q_\theta(\boldsymbol{\theta})} &= \\ \int d\mathbf{x} d\boldsymbol{\theta} Q_x(\mathbf{x}) Q_\theta(\boldsymbol{\theta}) \log \frac{Q_x(\mathbf{x}) Q_\theta(\boldsymbol{\theta})}{P(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})} &= \text{KL}(Q \| P) \end{aligned}$$