

# Lecture 12: Models for documents

Machine Learning 4F13, Michaelmas 2015

Zoubin Ghahramani

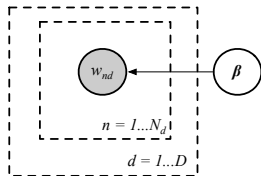
Department of Engineering  
University of Cambridge

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

# A really simple document model

Consider a collection of  $D$  documents from a vocabulary of  $M$  words.

- $N_d$ : number of words in document  $d$ .
- $w_{nd}$ :  $n$ -th word in document  $d$  ( $w_{nd} \in \{1 \dots M\}$ ).
- $w_{nd} \sim \text{Cat}(\beta)$ : each word is drawn from a discrete categorical distribution with parameters  $\beta$
- $\beta = [\beta_1, \dots, \beta_M]^\top$ : parameters of a categorical / multinomial distribution<sup>1</sup> over the  $M$  vocabulary words.



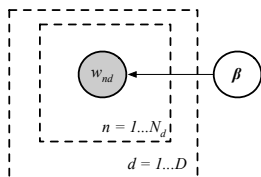
---

<sup>1</sup>It'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 vocabulary of  $M$  unique words.

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



We can fit  $\beta$  by maximising the likelihood:

$$\begin{aligned}\hat{\beta} &= \operatorname{argmax}_{\beta} \prod_{d=1}^D \prod_{n=1}^{N_d} \text{Cat}(w_{nd} | \beta) \\ &= \operatorname{argmax}_{\beta} \text{Mult}(c_1, \dots, c_M | \beta, N)\end{aligned}$$

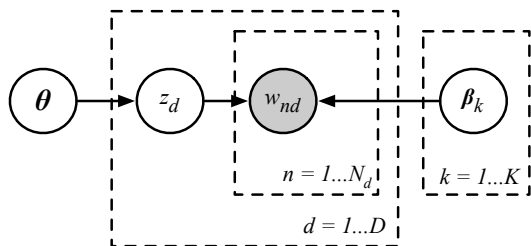
$$\hat{\beta}_m = \frac{c_m}{N} = \frac{c_m}{\sum_{\ell=1}^M c_{\ell}}$$

- $N = \sum_{d=1}^D N_d$ : total number of words in the collection.
- $c_m = \sum_{d=1}^D \sum_{n=1}^{N_d} \mathbb{I}(w_{nd} = m)$ : total count of vocabulary word  $m$ .

# Limitations of the really simple document model

- Document  $d$  is the result of sampling  $N_d$  words from the categorical distribution with parameters  $\beta$ .
- $\beta$  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.
- We would like a model where different documents might be about different *topics*.

# A mixture of categoricals model



$$z_d \sim \text{Cat}(\boldsymbol{\theta})$$
$$w_{nd} | z_d \sim \text{Cat}(\boldsymbol{\beta}_{z_d})$$

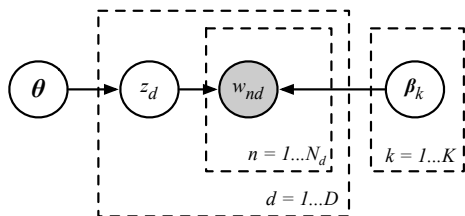
We want to allow for a mixture of  $K$  categoricals parametrised by  $\beta_1, \dots, \beta_K$ . Each of those categorical distributions corresponds to a *document category*.

- $z_d \in \{1, \dots, K\}$  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  $\boldsymbol{\theta} = [\theta_1, \dots, \theta_K]$  is the parameter of a categorical distribution 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  $\boldsymbol{\theta}$  and  $\boldsymbol{\beta}$ ?

# A mixture of categoricals model: the likelihood



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

$$\begin{aligned} p(\mathbf{w}|\theta, \beta) &= \prod_{d=1}^D p(\mathbf{w}_d|\theta, \beta) \\ &= \prod_{d=1}^D \sum_{k=1}^K p(\mathbf{w}_d, z_d = k|\theta, \beta) \\ &= \prod_{d=1}^D \sum_{k=1}^K p(z_d = k|\theta) p(\mathbf{w}_d|z_d = k, \beta_k) \\ &= \prod_{d=1}^D \sum_{k=1}^K p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|z_d = k, \beta_k) \end{aligned}$$

# 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  $\theta$ , 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:

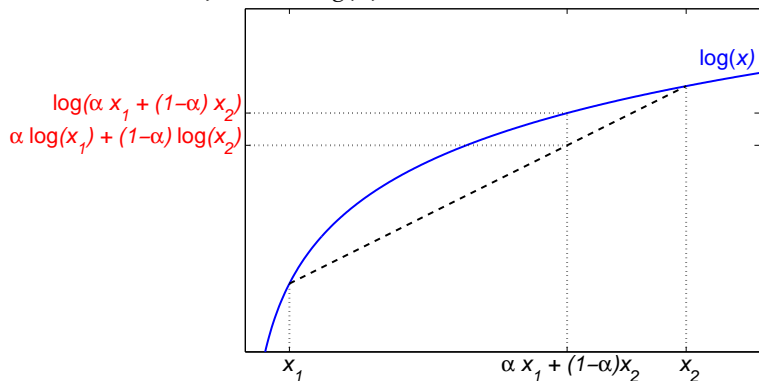
$$\mathcal{L}(\theta) = \log \int q(H) \frac{p(H, V|\theta)}{q(H)} dH \geq \int q(H) \log \frac{p(H, V|\theta)}{q(H)} dH = \mathcal{F}(q, \theta), \quad (2)$$

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  $\theta$ , 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:

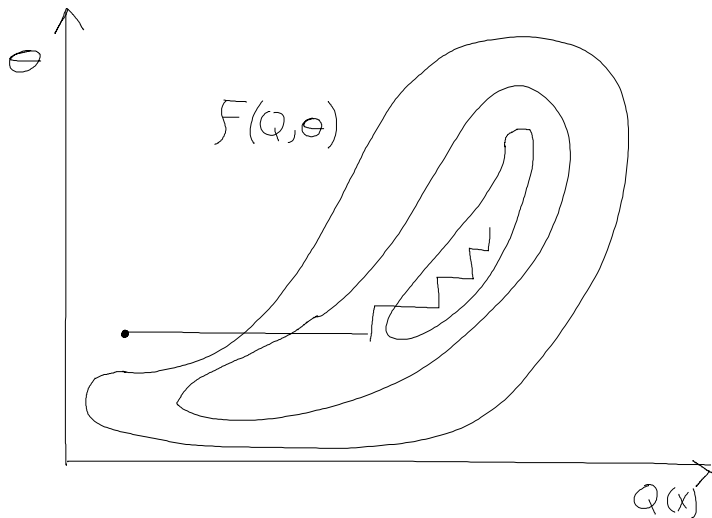
$$\mathbf{q}^{(k)}(\mathbf{H}) := \operatorname{argmax}_{\mathbf{q}(\mathbf{H})} \mathcal{F}(\mathbf{q}(\mathbf{H}), \theta^{(k-1)}). \quad (4)$$

**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  $\theta$ .

# 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 Categoricals

In the mixture model, the likelihood is:

$$p(\mathbf{w}|\boldsymbol{\theta}, \boldsymbol{\beta}) = \prod_{d=1}^D \sum_{k=1}^K p(z_d = k|\boldsymbol{\theta}) \prod_{n=1}^{N_d} p(w_{nd}|z_d = k, \boldsymbol{\beta}_k)$$

**E-step:** for each  $d$ , set  $q$  to the posterior (where  $c_{md} = \sum_{n=1}^{N_d} \mathbb{I}(w_{nd} = m)$ ):

$$q(z_d = k) \propto p(z_d = k|\boldsymbol{\theta}) \prod_{n=1}^{N_d} p(w_{nd}|\boldsymbol{\beta}_{k,w_n}) = \theta_k \text{Mult}(c_{1d}, \dots, c_{Md}|\boldsymbol{\beta}_k, N_d) \stackrel{\text{def}}{=} r_{kd}$$

**M-step:** Maximize

$$\begin{aligned} \sum_{d=1}^D \sum_{k=1}^K q(z_d = k) \log p(\mathbf{w}, z_d) &= \sum_{k,d} r_{kd} \log \left[ p(z_d = k|\boldsymbol{\theta}) \prod_{n=1}^{N_d} p(w_{nd}|\boldsymbol{\beta}_{k,w_{nd}}) \right] \\ &= \sum_{k,d} r_{kd} \left( \log \prod_{m=1}^M \beta_{km}^{c_{md}} + \log \theta_k \right) \\ &= \sum_{k,d} r_{kd} \left( \sum_{m=1}^M c_{md} \log \beta_{km} + \log \theta_k \right) \stackrel{\text{def}}{=} F(\mathbf{R}, \boldsymbol{\theta}, \boldsymbol{\beta}) \end{aligned}$$

# EM: M step for mixture model

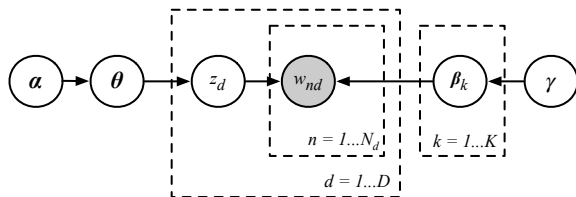
$$F(\mathbf{R}, \theta, \beta) = \sum_{k,d} r_{kd} \left( \sum_{m=1}^M c_{m,d} \log \beta_{km} + \log \theta_k \right)$$

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

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

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

# A Bayesian mixture of categoricals model



$$\begin{aligned}\theta &\sim \text{Dir}(\alpha) \\ \beta_k &\sim \text{Dir}(\gamma) \\ z_d | \theta &\sim \text{Cat}(\theta) \\ w_{nd} | z_d, \beta &\sim \text{Cat}(\beta_{z_d})\end{aligned}$$

With the EM algorithm we have essentially estimated  $\theta$  and  $\beta$  by maximum likelihood. An alternative, Bayesian treatment infers these parameters starting from priors, e.g.:

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

What is different?

- We no longer want to compute a point estimate of  $\theta$  or  $\beta$ .
- 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  $\theta$ .

**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:

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

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

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

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