

PROBABILISTIC MACHINE LEARNING
LECTURE 24
VARIATIONAL INFERENCE

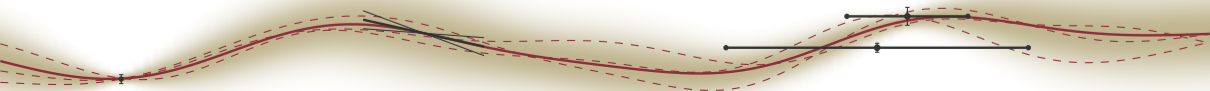
Philipp Hennig

13 July 2021

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



FACULTY OF SCIENCE
DEPARTMENT OF COMPUTER SCIENCE
CHAIR FOR THE METHODS OF MACHINE LEARNING





$$\log p(x | \theta) = \mathcal{L}(q, \theta) + D_{\text{KL}}(q \| p(z | x, \theta))$$

$$\mathcal{L}(q, \theta) = \sum_z q(z) \log \left(\frac{p(x, z | \theta)}{q(z)} \right) \quad D_{\text{KL}}(q \| p(z | x, \theta)) = - \sum_z q(z) \log \left(\frac{p(z | x, \theta)}{q(z)} \right)$$

- ▶ For EM, we minimized KL-divergence to find $q = p(z | x, \theta)$ (E), then maximized $\mathcal{L}(q, \theta)$ in θ .
- ▶ What if we treated the parameters θ as a *probabilistic* variable for full Bayesian inference?

$$z \leftarrow z \cup \theta$$

- ▶ Then we could just maximize $\mathcal{L}(q(z))$ wrt. q (not z !) to implicitly minimize $D_{\text{KL}}(q \| p(z | x))$, because $\log p(x)$ is constant. This is an **optimization in the space of distributions** q , not (necessarily) in parameters of such distributions, and thus a very powerful notion.
- ▶ In general, this will be intractable, because the optimal choice for q is exactly $p(z | x)$. But maybe we can help out a bit with approximations. Amazingly, we often don't need to impose strong approximations. Sometimes we can get away with just imposing restrictions on the **factorization** of q , not its analytic form.

$$\log p(x) = \mathcal{L}(q) + D_{\text{KL}}(q \| p(z | x))$$

$$\mathcal{L}(q) = \int q(z) \log \left(\frac{p(x, z)}{q(z)} \right) dz \quad D_{\text{KL}}(q \| p(z | x)) = - \int q(z) \log \left(\frac{p(z | x)}{q(z)} \right) dz$$

- ▶ For EM, we minimized KL-divergence to find $q = p(z | x, \theta)$ (E), then maximized $\mathcal{L}(q, \theta)$ in θ .
- ▶ What if we treated the parameters θ as a *probabilistic* variable for full Bayesian inference?

$$z \leftarrow z \cup \theta$$

- ▶ Then we could just maximize $\mathcal{L}(q(z))$ wrt. q (not z !) to implicitly minimize $D_{\text{KL}}(q \| p(z | x))$, because $\log p(x)$ is constant. This is an **optimization in the space of distributions** q , not (necessarily) in parameters of such distributions, and thus a very powerful notion.
- ▶ In general, this will be intractable, because the optimal choice for q is exactly $p(z | x)$. But maybe we can help out a bit with approximations. Amazingly, we often don't need to impose strong approximations. Sometimes we can get away with just imposing restrictions on the **factorization** of q , not its analytic form.

Consider a joint distribution $p(x, z)$ with $z \in \mathbb{R}^n$

- ▶ to find a “good” but tractable approximation $q(z)$, assume that it factorizes $q(z) = \prod_i q_i(z_i)$.
- ▶ Initialize all q_i to some initial *distribution*
- ▶ Iteratively compute

$$\begin{aligned}\mathcal{L}(q) &= \int q_j \log \tilde{p}(x, z_j) dz_j - \int q_j \log q_j dz_j + \text{const.} \\ &= -D_{\text{KL}}(q_j(z) \parallel \tilde{p}(x, z_j)) + \text{const.}\end{aligned}$$

and maximize wrt. q_j . Doing so *minimizes* $D_{\text{KL}}(q(z_j) \parallel \tilde{p}(x, z_j))$, thus the minimum is at q_j^* with

$$\log q_j^*(z_j) = \log \tilde{p}(x, z_j) = \mathbb{E}_{q, i \neq j}(\log p(x, z)) + \text{const.} \quad (\star)$$

- ▶ note that this expression identifies a **function** q_j , not some parametric form.
- ▶ the optimization converges, because $-\mathcal{L}(q)$ can be shown to be *convex* wrt. q .

In physics, this trick is known as **mean field theory** (because an n -body problem is separated into n separate problems of individual particles who are affected by the “mean field” \tilde{p} summarizing the expected effect of all other particles).

Variational Inference

- ▶ is a general framework to construct approximating **probability distributions** $q(z)$ to non-analytic posterior distributions $p(z | x)$ by minimizing the **functional**

$$q^* = \arg \min_{q \in \mathcal{Q}} D_{\text{KL}}(q(z) \| p(z | x)) = \arg \max_{q \in \mathcal{Q}} \mathcal{L}(q)$$

- ▶ the beauty is that we get to *choose* q , so one can nearly always find a tractable approximation.
- ▶ If we impose the *mean field approximation* $q(z) = \prod_i q(z_i)$, get

$$\log q_j^*(z_j) = \mathbb{E}_{q, i \neq j}(\log p(x, z)) + \text{const.}$$

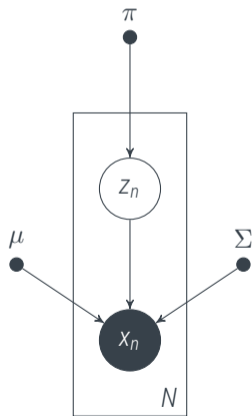
- ▶ for Exponential Family p things are particularly simple: we only need the expectation under q of the sufficient statistics.

Variational Inference is an extremely flexible and powerful approximation method. Its downside is that constructing the bound and update equations can be tedious. For a quick test, variational inference is often not a good idea. But for a deployed product, it can be the most powerful tool in the box.



- Remember EM for Gaussian mixtures $\theta := (\pi, \mu, \Sigma)$

$$\begin{aligned} p(x, z \mid \mu, \Sigma, \pi) &= \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \cdot \mathcal{N}(x_n; \mu_k, \Sigma_k)^{z_{nk}} \\ &= \prod_{n=1}^N p(z_{n:} \mid \pi) \cdot p(x_n \mid z_{n:}, \mu, \Sigma) \end{aligned}$$





- Remember EM for Gaussian mixtures $\theta := (\pi, \mu, \Sigma)$

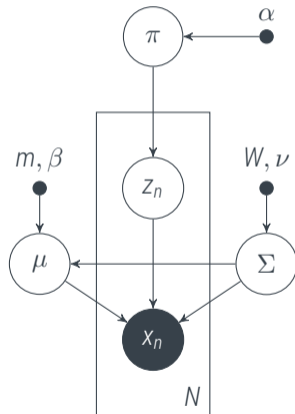
$$p(x, z \mid \mu, \Sigma, \pi) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \cdot \mathcal{N}(x_n; \mu_k, \Sigma_k)^{z_{nk}}$$

- For Bayesian inference, turn parameters into variables

$$p(x, z, \pi, \mu, \Sigma) = p(x, z \mid \pi, \mu, \Sigma) \cdot p(\pi) \cdot p(\mu \mid \Sigma) \cdot p(\Sigma)$$

$$p(\pi) = \mathcal{D}(\pi \mid \alpha) = \frac{\Gamma\left(\sum_{k=1}^K \alpha_k\right)}{\prod_k \Gamma(\alpha_k)} \prod_k \pi_k^{\alpha_k - 1}$$

$$p(\mu \mid \Sigma) \cdot p(\Sigma) = \prod_{k=1}^K \mathcal{N}(\mu_k; m, \Sigma_k / \beta) \cdot \mathcal{W}(\Sigma_k^{-1}; W, \nu)$$



- ▶ We know that the full posterior $p(z, \pi, \mu, \Sigma | x)$ is intractable (check the graph!)
- ▶ But let's consider an approximation $q(z, \pi, \mu, \Sigma)$ with the factorization

$$q(z, \pi, \mu, \Sigma) = q(z) \cdot q(\pi, \mu, \Sigma)$$

- ▶ from (\star) , we have

$$\begin{aligned} \log q^*(z) &= \mathbb{E}_{q(\pi, \mu, \Sigma)} (\log p(x, z, \pi, \mu, \Sigma)) + \text{const.} \\ &= \mathbb{E}_{q(\pi)} (\log p(z | \pi)) + \mathbb{E}_{q(\mu, \Sigma)} (\log p(x | z, \mu, \Sigma)) + \text{const.} \\ &= \sum_n^N \sum_k^K z_{nk} \underbrace{\left(\mathbb{E}_{q(\pi)} (\log \pi_k) + \frac{1}{2} \mathbb{E}_{q(\mu, \Sigma)} (\log |\Sigma^{-1}| - (x_n - \mu_k)^\top \Sigma_k^{-1} (x - \mu_k)) \right)}_{=:\log \rho_{nk}} + \text{const.} \end{aligned}$$

$$q^*(z) \propto \prod_n \prod_k \rho_{nk}^{z_{nk}} \quad \text{define } r_{nk} = \frac{\rho_{nk}}{\sum_{j=1}^K \rho_{nj}}, \text{ then } q^*(z) = \prod_n \prod_k r_{nk}^{z_{nk}} \text{ with } \mathbb{E}_{q(z)} [z] = r_{nk}$$

using $q(z, \pi, \mu, \Sigma) = q(z) \cdot q(\pi, \mu, \Sigma)$

[Exposition from Bishop, PRML 2006, Chapter 10.2]

- ▶ Define some convenient notation:

$$N_k := \sum_{n=1}^N r_{nk} \quad \bar{x}_k := \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n \quad S_k := \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \bar{x}_k)(x_n - \bar{x}_k)^\top$$

- ▶ from (\star) , we have

$$\log q^*(\pi, \mu, \Sigma) = \mathbb{E}_{q(z)} (\log p(x, z, \pi, \mu, \Sigma)) + \text{const.}$$

$$= \mathbb{E}_{q(z)} \left(\log p(\pi) + \sum_k \log p(\mu_k, \Sigma_k) + \log p(z | \pi) + \sum_n \log p(x_n | z, \mu, \Sigma) \right)$$

$$= \log p(\pi) + \sum_{k=1}^K \log p(\mu_k, \Sigma_k) + \mathbb{E}_{q(z)} (\log p(z | \pi))$$

$$+ \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}(z_{nk}) \log \mathcal{N}(x_n; \mu_k, \Sigma_k) + \text{const.}$$

$$\begin{aligned} \log q^*(\pi, \mu, \Sigma) &= \log p(\pi) + \sum_{k=1}^K \log p(\mu_k, \Sigma_k) + \mathbb{E}_{q(z)}(\log p(z | \pi)) \\ &\quad + \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}(z_{nk}) \log \mathcal{N}(x_n; \mu_k, \Sigma_k) + \text{const.} \end{aligned}$$

- The bound exposes an **induced factorization** into $q(\pi, \mu, \Sigma) = q(\pi) \cdot \prod_{k=1}^K q(\mu_k, \Sigma_k)$

$$\text{where } \log q(\pi) = \log p(\pi) + \mathbb{E}_{q(z)}(\log p(z | \pi)) + \text{const.}$$

$$= (\alpha - 1) \sum_k \log \pi_k + \sum_k \sum_n r_{nk} \log \pi_k + \text{const.}$$

$$q(\pi) = \mathcal{D}(\pi, \alpha_k := \alpha + N_k) \quad \text{with } N_k = \sum_n r_{nk}$$

$$\begin{aligned} \log q^*(\pi, \mu, \Sigma) &= \log p(\pi) + \sum_{k=1}^K \log p(\mu_k, \Sigma_k) + \mathbb{E}_{q(z)}(\log p(z | \pi)) \\ &+ \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}(z_{nk}) \log \mathcal{N}(x_n; \mu_k, \Sigma_k) + \text{const.} \end{aligned}$$

- The bound exposes an **induced factorization** into $q(\pi, \mu, \Sigma) = q(\pi) \cdot \prod_{k=1}^K q(\mu_k, \Sigma_k)$

where (leaving out some tedious algebra) $q^*(\mu_k, \Sigma_k) = \mathcal{N}(\mu_k; m_k, \Sigma_k / \beta_k) \mathcal{W}(\Sigma_k^{-1}; W_k, \nu_k)$

$$\text{with } \beta_k := \beta + N_k \quad m_k := \frac{1}{\beta_k} (\beta m + N_k \bar{x}_k) \quad \nu_k := \nu + N_k$$

$$W_k^{-1} := W^{-1} + N_k S_k + \frac{\beta N_k}{\beta + N_k} (\bar{x}_k - m)(\bar{x}_k - m)^\top$$

Some tabulated identities, required for the concrete algorithm

$$p(x \mid \alpha) = \mathcal{D}(x; \alpha) = \frac{\Gamma(\hat{\alpha})}{\prod_d \Gamma(\alpha_d)} \prod_d x^{\alpha_d - 1} = \frac{1}{B(\alpha)} \prod_d x^{\alpha_d - 1} \quad \hat{\alpha} := \sum_d \alpha_d$$

- ▶ $\mathbb{E}_p(x_d) = \frac{\alpha_d}{\hat{\alpha}}$
- ▶ $\text{var}_p(x_d) = \frac{\alpha_d(\hat{\alpha} - \alpha_d)}{\hat{\alpha}^2(\hat{\alpha} + 1)}$
- ▶ $\text{cov}(x_d, x_i) = -\frac{\alpha_d \alpha_i}{\hat{\alpha}^2(\hat{\alpha} + 1)}$
- ▶ $\text{mode}(x_d) = \frac{\alpha_d - 1}{\hat{\alpha} - D}$
- ▶ $\mathbb{E}_p(\log x_d) = F(\alpha_d) - F(\hat{\alpha})$
- ▶ $\mathbb{H}(p) = -\int p(x) \log p(x) dx = -\sum_d (\alpha_d - 1)(F(\alpha_d) - F(\hat{\alpha})) + \log B(\alpha)$

Where $F(z) = \frac{d}{dz} \log \Gamma(z)$ (the “digamma-function”).

`scipy.special.digamma(z)`

<https://dlmf.nist.gov/5>

- Recall from above:

$$\log q^*(z) = \sum_n^N \sum_k^K z_{nk} \underbrace{\left(\mathbb{E}_{q(\pi)}(\log \pi_k) + \frac{1}{2} \mathbb{E}_{q(\mu, \Sigma)}(\log |\Sigma^{-1}| - (x_n - \mu_k)^\top \Sigma_k^{-1} (x_n - \mu_k)) \right)}_{=:\log \rho_{nk}} + \text{const.}$$

- now we can evaluate ρ_{nk} , using tabulated identities

$$\log \tilde{\pi}_k := \mathbb{E}_{\mathcal{D}(\pi; \alpha_k)}(\log \pi_k) = F(\alpha_k) - F\left(\sum_k \alpha_k\right)$$

and for the Wishart:

$$\log |\tilde{\Sigma}^{-1}|_k := \mathbb{E}_{\mathcal{W}(\Sigma_k^{-1}; W_k, \nu_k)}(\log |\Sigma_k^{-1}|) = \sum_{d=1}^D F\left(\nu_k + \frac{1-d}{2}\right) + D \log 2 + \log |W_k|$$

$$\mathbb{E}_{\mathcal{N}(\mu_k; m_k, \Sigma_k / \beta_k) \mathcal{W}(\Sigma_k^{-1}; W_k, \nu_k)}((x_n - \mu_k)^\top \Sigma^{-1} (x_n - \mu_k)) = D \beta_k^{-1} + \nu_k (x_n - m_k)^\top W_k (x_n - m_k)$$

- ▶ this yields the update equation

$$\mathbb{E}_q(z_{nk}) = r_{nk} \propto \tilde{\pi}_k |\Sigma^{-1}|^{1/2} \exp\left(-\frac{D}{2\beta_k} - \frac{\nu_k}{2} (x_n - m_k)^\top W_k (x_n - m_k)\right)$$

compare this with the EM-update

$$r_{nk} \propto \pi_k |\Sigma^{-1}|^{1/2} \exp\left(-\frac{1}{2} (x_n - \mu_k)^\top \Sigma_k^{-1} (x_n - \mu_k)\right)$$

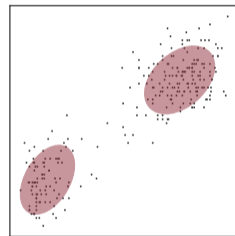
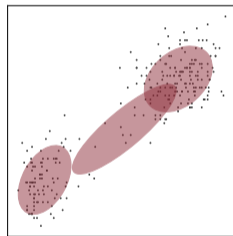
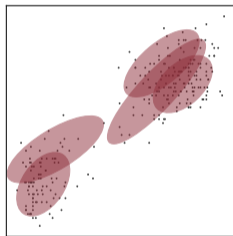
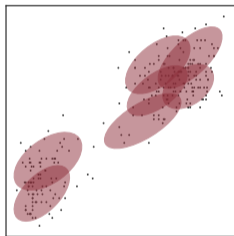
- ▶ Here, variational Inference is the Bayesian version of EM: Instead of maximizing the likelihood for $\theta = (\mu, \Sigma, \pi)$, we maximize a variational bound.
- ▶ One advantage of this is that the posterior can actually “decide” to ignore components:

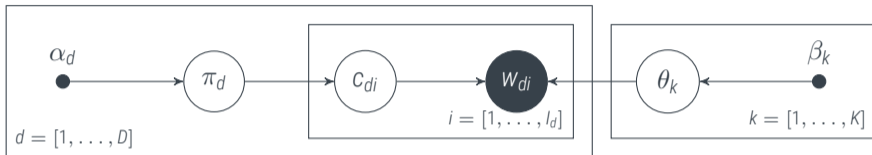
Example

The Old Faithful Dataset, using $\alpha = 10^{-3}$



[from Bishop, PRML 2006, Fig. 10.6 / Ann-Kathrin Schalkamp]





To draw I_d words $w_{di} \in [1, \dots, V]$ of document $d \in [1, \dots, D]$:

- ▶ Draw K topic distributions θ_k over V words from
- ▶ Draw D document distributions over K topics from
- ▶ Draw topic assignments c_{ik} of word w_{di} from
- ▶ Draw word w_{di} from

$$p(\Theta | \beta) = \prod_{k=1}^K \mathcal{D}(\theta_k; \beta_k)$$

$$p(\Pi | \alpha) = \prod_{d=1}^D \mathcal{D}(\pi_d; \alpha_d)$$

$$p(C | \Pi) = \prod_{i,d,k} \pi_{dk}^{c_{dik}}$$

$$p(w_{di} = v | c_{di}, \Theta) = \prod_k \theta_{kv}^{c_{dik}}$$

Useful notation: $n_{dkv} = \#\{i : w_{di} = v, c_{ijk} = 1\}$. Write $n_{dk} := [n_{dk1}, \dots, n_{dkV}]$ and $n_{dk\cdot} = \sum_v n_{dkv}$, etc.

$$p(C, \Pi, \Theta, W) = \left(\prod_{d=1}^D \mathcal{D}(\boldsymbol{\pi}_d; \boldsymbol{\alpha}_d) \right) \cdot \left(\prod_{d=1}^D \prod_{i=1}^{I_d} \left(\prod_{k=1}^K \pi_{dk}^{c_{dik}} \right) \right) \cdot \left(\prod_{d=1}^D \prod_{i=1}^{I_d} \left(\prod_{k=1}^K \theta_{kw_{di}}^{c_{dik}} \right) \right) \cdot \left(\prod_{k=1}^K \mathcal{D}(\boldsymbol{\theta}_k; \boldsymbol{\beta}_k) \right)$$

- ▶ The posterior $p(\Pi, \Theta, C | W)$ is intractable. We want an approximation q that *factorises*

$$q(\Pi, \Theta, C) = q(C) \cdot q(\Pi, \Theta)$$

- ▶ To find the *best* such approximation – the one that *minimizes* $D_{\text{KL}}(q || p(\Pi, \Theta, C | W))$, we *maximize* the **ELBO** (minimize variational free energy)

$$\mathcal{L}(q) = \int q(C, \Theta, \Pi) \log \left(\frac{p(C, \Pi, \Theta, W)}{q(C, \Theta, \Pi)} \right) dC d\Theta d\Pi$$

$$p(C, \Pi, \Theta, W) = \left(\prod_{d=1}^D \mathcal{D}(\boldsymbol{\pi}_d; \boldsymbol{\alpha}_d) \right) \cdot \left(\prod_{d=1}^D \prod_{i=1}^{I_d} \left(\prod_{k=1}^K \pi_{dk}^{c_{dik}} \right) \right) \cdot \left(\prod_{d=1}^D \prod_{i=1}^{I_d} \left(\prod_{k=1}^K \theta_{kw_{di}}^{c_{dik}} \right) \right) \cdot \left(\prod_{k=1}^K \mathcal{D}(\boldsymbol{\theta}_k; \boldsymbol{\beta}_k) \right)$$

Recall from above: To maximize the ELBO of a factorized approximation, compute the **mean field**

$$\log q^*(z_i) = \mathbb{E}_{z_j, j \neq i} (\log p(x, z)) + \text{const.}$$

$$\log q^*(C) = \mathbb{E}_{q(\Pi, \Theta)} \left(\sum_{d,i,k} c_{dik} \log(\pi_{dk} \theta_{kw_{di}}) \right) + \text{const.} = \sum_{d,i} \sum_{k=1}^K c_{dik} \underbrace{\left(\mathbb{E}_{q(\Pi, \Theta)} (\log \pi_{dk} \theta_{dw_{di}}) \right)}_{=:\log \gamma_{dik}} + \text{const.}$$

Thus, $q(C) = \prod_{d,i} q(\mathbf{c}_{di})$ with $q(\mathbf{c}_{di}) = \prod_k \tilde{\gamma}_{dik}^{c_{dik}}$, where $\tilde{\gamma}_{dik} = \gamma_{dik} / \sum_k \gamma_{dik}$

(Note: Thus, $\mathbb{E}_q(c_{dik}) = \tilde{\gamma}_{dik}$)

$$p(C, \Pi, \Theta, W) = \left(\prod_{d=1}^D \mathcal{D}(\boldsymbol{\pi}_d; \boldsymbol{\alpha}_d) \right) \cdot \left(\prod_{d=1}^D \prod_{i=1}^{I_d} \left(\prod_{k=1}^K \pi_{dk}^{c_{dik}} \right) \right) \cdot \left(\prod_{d=1}^D \prod_{i=1}^{I_d} \left(\prod_{k=1}^K \theta_{kw_{di}}^{c_{dik}} \right) \right) \cdot \left(\prod_{k=1}^K \mathcal{D}(\boldsymbol{\theta}_k; \boldsymbol{\beta}_k) \right)$$

Recall from above: To maximize the ELBO of a factorized approximation, compute the **mean field**

$$\log q^*(z_i) = \mathbb{E}_{z_j, j \neq i} (\log p(x, z)) + \text{const.}$$

$$\begin{aligned} \log q^*(\Pi, \Theta) &= \mathbb{E}_{\Pi_{d,j}, q(c_{di})} \left(\sum_{d,k} (\alpha_{dk} - 1 + n_{dk\cdot}) \log \pi_{dk} + \sum_{k,v} (\beta_{kv} - 1 + n_{\cdot kv}) \log \theta_{kv} \right) + \text{const.} \\ &= \sum_{d=1}^D \sum_{k=1}^K (\alpha_{dk} - 1 + \mathbb{E}_{q(C)}(n_{dk\cdot})) \log \pi_{dk} + \sum_{k=1}^K \sum_{v=1}^V (\beta_{kv} - 1 + \mathbb{E}_{q(C)}(n_{\cdot kv})) \log \theta_{kv} + \text{const.} \\ q^*(\Pi, \Theta) &= \prod_{d=1}^D \mathcal{D}(\boldsymbol{\pi}_d; \tilde{\boldsymbol{\alpha}}_d := [\alpha_d + \tilde{\gamma}_d]) \cdot \prod_{k=1}^K \mathcal{D}(\boldsymbol{\theta}_k; \tilde{\boldsymbol{\beta}}_k := [\beta_k + \sum_d \sum_{i=1}^{I_d} \tilde{\gamma}_{di} \mathbb{I}(w_{di} = v)]_{v=1, \dots, V})). \end{aligned}$$

$$q(\boldsymbol{\pi}_d) = \mathcal{D} \left(\boldsymbol{\pi}_d; \tilde{\alpha}_{dk} := \left[\alpha_{dk} + \sum_{i=1}^{l_d} \tilde{\gamma}_{dik} \right]_{k=1, \dots, K} \right) \quad \forall d = 1, \dots, D$$

$$q(\boldsymbol{\theta}_k) = \mathcal{D} \left(\boldsymbol{\theta}_k; \tilde{\beta}_{kv} := \left[\beta_{kv} + \sum_d^D \sum_{i=1}^{l_d} \tilde{\gamma}_{dik} \mathbb{I}(w_{di} = v) \right]_{v=1, \dots, V} \right) \quad \forall k = 1, \dots, K$$

$$q(\mathbf{c}_{di}) = \prod_k \tilde{\gamma}_{dik}^{c_{dik}}, \quad \forall d \quad i = 1, \dots, l_d$$

where $\tilde{\gamma}_{dik} = \gamma_{dik} / \sum_k \gamma_{dik}$ and (note that $\sum_k \tilde{\alpha}_{dk} = \text{const.}$)

$$\begin{aligned} \gamma_{dik} &= \exp \left(\mathbb{E}_{q(\boldsymbol{\pi}_{dk})} (\log \pi_{dk}) + \mathbb{E}_{q(\boldsymbol{\theta}_{di})} (\log \theta_{kw_{di}}) \right) \\ &= \exp \left(F(\tilde{\alpha}_{jk}) + F(\tilde{\beta}_{kw_{di}}) - F \left(\sum_v \tilde{\beta}_{kv} \right) \right) \end{aligned}$$

$$p(C, \Pi, \Theta, W) = \left(\prod_{d=1}^D \frac{\Gamma(\sum_k \alpha_{dk})}{\prod_k \Gamma(\alpha_{dk})} \prod_{k=1}^K \pi_{dk}^{\alpha_{dk}-1+n_{dk}} \right) \cdot \left(\prod_{k=1}^K \frac{\Gamma(\sum_v \beta_{kv})}{\prod_v \Gamma(\beta_{kv})} \prod_{v=1}^V \theta_{kv}^{\beta_{kv}-1+n_{kv}} \right)$$

We need

$$\begin{aligned} \mathcal{L}(q, W) &= \mathbb{E}_q(\log p(W, C, \Theta, \Pi)) + \mathbb{H}(q) \\ &= \int q(C, \Theta, \Pi) \log p(W, C, \Theta, \Pi) dC d\Theta d\Pi - \int q(C, \Theta, \Pi) \log q(C, \Theta, \Pi) dC d\Theta d\Pi \\ &= \int q(C, \Theta, \Pi) \log p(W, C, \Theta, \Pi) dC d\Theta d\Pi + \sum_k \mathbb{H}(\mathcal{D}(\theta_k \tilde{\beta}_k)) + \sum_d \mathbb{H}(\mathcal{D}(\pi_d \tilde{\alpha}_d)) + \sum_{di} \mathbb{H}(\tilde{\gamma}_{di}) \end{aligned}$$

The entropies can be computed from the tabulated values. For the expectation, we use $\mathbb{E}_{q(C)}(n_{dkv}) = \sum_i \gamma_{dik} \mathbb{I}(W_{di} = v)$ and use $\mathbb{E}_{\mathcal{D}(\pi_d; \tilde{\alpha})}(\log \pi_d) = F(\tilde{\alpha}_d) - F(\hat{\alpha})$ from above.

Dirty secret: In practice, the ELBO itself isn't strictly necessary.

```
1 procedure LDA( $W, \alpha, \beta$ )
2    $\tilde{\gamma}_{dik} \leftarrow \text{DIRICHLET\_RAND}(\alpha)$  // initialize
3    $\mathcal{L} \leftarrow -\infty$ 
4   while  $\mathcal{L}$  not converged do
5     for  $d = 1, \dots, D; k = 1, \dots, K$  do
6        $\tilde{\alpha}_{dk} \leftarrow \alpha_{dk} + \sum_i \tilde{\gamma}_{dik}$  // update document-topics distributions
7     end for
8     for  $k = 1, \dots, K; v = 1, \dots, V$  do
9        $\tilde{\beta}_{kv} \leftarrow \beta_{kv} + \sum_{d,i} \tilde{\gamma}_{dik} \mathbb{I}(w_{di} = v)$  // update topic-word distributions
10    end for
11    for  $d = 1, \dots, D; k = 1, \dots, K; i = 1, \dots, l_d$  do
12       $\tilde{\gamma}_{dik} \leftarrow \exp(F(\tilde{\alpha}_{dk}) + F(\tilde{\beta}_{kw_{di}}) - F(\sum_v \tilde{\beta}_{kv}))$  // update word-topic assignments
13       $\tilde{\gamma}_{dik} \leftarrow \tilde{\gamma}_{dik} / \tilde{\gamma}_{di}$ 
14    end for
15     $\mathcal{L} \leftarrow \text{BOUND}(\tilde{\gamma}, w, \tilde{\alpha}, \tilde{\beta})$  // update bound
16  end while
17 end procedure
```

$$\mathcal{L}(q) = \int q(\mathcal{C}, \Theta, \Pi) \log \left(\frac{p(\mathcal{C}, \Pi, \Theta, W)}{q(\mathcal{C}, \Theta, \Pi)} \right) d\mathcal{C} d\Theta d\Pi$$

Variational Inference is a powerful mathematical tool to construct efficient approximations to intractable *probability distributions* (not just point estimates, but entire distributions). Often, just imposing factorization is enough to make things tractable. The downside of variational inference is that constructing the bound can take significant ELBOw grease. However, the resulting algorithms are often highly efficient compared to tools that require less derivation work, like Monte Carlo.

“Derive your variational bound in the time it takes for your Monte Carlo sampler to converge.”

The Toolbox

Framework:

$$\int p(x_1, x_2) dx_2 = p(x_1)$$

$$p(x_1, x_2) = p(x_1 | x_2)p(x_2)$$

$$p(x | y) = \frac{p(y | x)p(x)}{p(y)}$$

Modelling:

- ▶ graphical models
- ▶ Gaussian distributions
- ▶ (deep) learnt representations
- ▶ Kernels
- ▶ Markov Chains
- ▶ Exponential Families / Conjugate Priors
- ▶ Factor Graphs & Message Passing

Computation:

- ▶ Monte Carlo
- ▶ Linear algebra / Gaussian inference
- ▶ maximum likelihood / MAP
- ▶ Laplace approximations
- ▶ EM (iterative maximum likelihood)
- ▶ variational inference / mean field