

PROBABILISTIC INFERENCE AND LEARNING

LECTURE 24

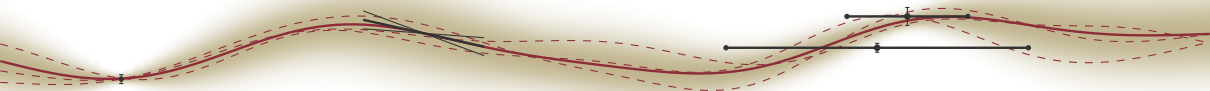
AN EXTENSIVE EXAMPLE II

Philipp Hennig
23 January 2019

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

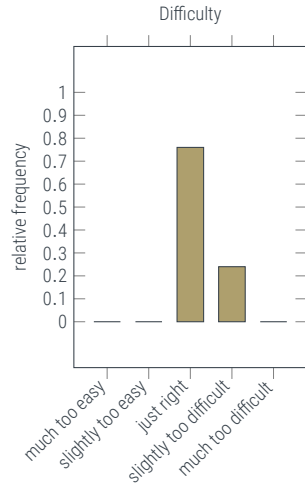
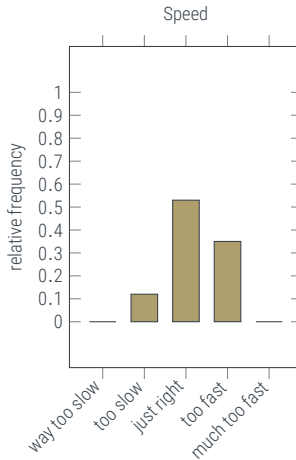
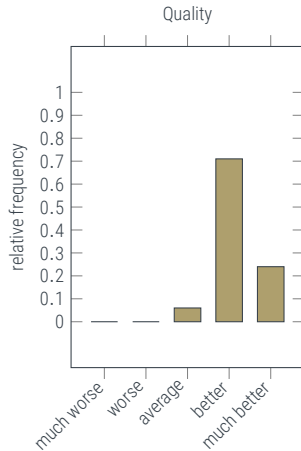


FACULTY OF SCIENCE
DEPARTMENT OF COMPUTER SCIENCE
CHAIR FOR THE METHODS OF MACHINE LEARNING



Last Lecture: Debrief

Feedback dashboard





Things you did not like:

- ✦ not caring about word order
- ✦ the indices / indexing errors
- ✦ using both Π and π
- ✦ putting the equations together
- ✦ Exercise 11 was too hard, more exam-like questions, please.

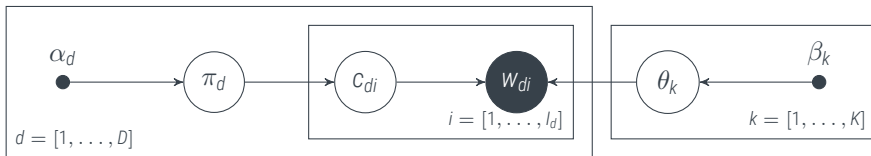
Things you did not understand:

- ✦ what are π_{dk} and θ_{kv} , intuitively?
- ✦ how can we measure the performance of the approaches?

Things you enjoyed:

- ✦ **the Example**
- ✦ the “thumb of uncertainty”
- ✦ the directed graph
- ✦ trying out naïve solutions first

0. Introduction to Reasoning under Uncertainty
 1. Probabilistic Reasoning
 2. Probabilities over Continuous Variables
 3. Gaussian Probability Distributions
 4. Gaussian Parametric Regression
 5. More on Parametric Regression
 6. Gaussian Processes
 7. More on Kernels & GPs
 8. A practical GP example
 9. Markov Chains, Time Series, Filtering
 10. Classification
 11. Empirical Example of Classification
 12. Bayesianism and Frequentism
 13. Stochastic Differential Equations
 14. Exponential Families
 15. Graphical Models
 16. Factor Graphs
 17. The Sum-Product Algorithm
 18. Mixture Models
 19. The EM Algorithm
 20. Variational Inference
 21. Monte Carlo
 22. Markov Chain Monte Carlo
 23. Advanced Modelling Example I
 24. **Advanced Modelling Example II**
 25. Advanced Modelling Example III
 26. Advanced Modelling Example IV
 27. Some Wild Stuff
 28. Revision



To draw l_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_{id} from
- ★ Draw word w_{id} from

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

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

$$p(\mathcal{C} \mid \Pi) = \prod_{i,d,k} \pi_{id}^{c_{dik}}$$

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

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

Recall definitions

$$\mathcal{D}(\boldsymbol{\pi}; \boldsymbol{\alpha}) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k - 1} \quad \text{and} \quad n_{dkv} = \#\{i : w_{di} = v, c_{ijk} = 1\}, n_{dk\cdot} = \sum_v n_{dkv}$$

Thus

$$\begin{aligned} p(\mathcal{C}, \Pi, \Theta, W) &= \left(\prod_{d=1}^D p(\boldsymbol{\pi}_d \mid \boldsymbol{\alpha}_d) \right) \cdot \left(\prod_{d=1}^D \prod_{i=1}^{l_d} p(c_{di} \mid \pi_d) \right) \cdot \left(\prod_{d=1}^D \prod_{i=1}^{l_d} p(w_{di} \mid c_{di}, \Theta) \right) \cdot \left(\prod_{k=1}^K p(\boldsymbol{\theta}_k \mid \boldsymbol{\beta}_k) \right) \\ &= \left(\prod_{d=1}^D \mathcal{D}(\boldsymbol{\pi}_d; \boldsymbol{\alpha}_d) \right) \cdot \left(\prod_{d=1}^D \prod_{i=1}^{l_d} \left(\prod_{k=1}^K \pi_{dk}^{c_{dik}} \right) \right) \cdot \left(\prod_{d=1}^D \prod_{i=1}^{l_d} \left(\prod_{k=1}^K \theta_{kw_{di}}^{c_{dik}} \right) \right) \cdot \left(\prod_{k=1}^K \mathcal{D}(\boldsymbol{\theta}_k; \boldsymbol{\beta}_k) \right) \\ &= \left(\prod_{d=1}^D \frac{\Gamma(\sum_k \alpha_{dk})}{\prod_k \Gamma(\alpha_{dk})} \prod_{k=1}^K \pi_{dk}^{\alpha_{dk} - 1 + n_{dk\cdot}} \right) \cdot \left(\prod_{k=1}^K \frac{\Gamma(\sum_v \beta_{kv})}{\prod_v \Gamma(\beta_{kv})} \prod_{v=1}^V \theta_{kv}^{\beta_{kv} - 1 + n_{\cdot kv}} \right) \end{aligned}$$

Designing a Probabilistic Machine Learning Model

1. Take a close look at the **Data**
2. Think about modelling goals, decide on data types
3. Design the **Model**
4. Design the **Algorithm**
 - ✦ For conditionally independent parts, conjugate priors may give analytic inference
 - ✦ for linearly related variables, consider Gaussians
 - ✦ consider revisiting the model to simplify as much as possible
 - ✦ if there is still no analytic answer, use **The Toolbox!**

The Toolbox

Five principal methods for dealing with computational complexity in probabilistic inference

1. **Maximum Likelihood (ML) / Maximum A-Posteriori (MAP)** estimation:

To estimate θ in $p(D \mid \theta)$ or $p(\theta \mid D)$, set $\hat{\theta} = \arg \max_{\theta} p$.

2. **Laplace** Approximation: $p(\theta \mid D) \approx \mathcal{N} \left(\theta; \hat{\theta}, -(\nabla \nabla^{\top} \log p(\theta \mid D))^{-1} \right)$

3. **Variational Inference**:

To approximate $p(\theta \mid D)$, impose structure on $q(\theta)$, then minimize $D_{\text{KL}}(q \parallel p)$

4. **Monte Carlo**: $\int f(x)p(x) dx \approx \sum_i f(x_i)$ where $x_i \sim p$

5. **Numerical Quadrature**: $\int p(f \mid \theta)p(\theta) d\theta \approx \sum_i w_i \cdot p(f \mid \theta_i)$

Disclaimer: The listed items are neither mutually exclusive nor collectively exhaustive. Some of the methods are intricately interrelated.

Find the *best simple* way to describe a *complex* thing

reminder: $n_{dkv} = \#\{i : w_{di} = v, c_{ijk} = 1\}$. $n_{dk\cdot} = \sum_v n_{dkv}$, etc.

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

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

$$q(\Pi, \Theta, C) = \left(\prod_{d,i} q(c_{di}) \right) \left(\prod_d q(\pi_d) \right) \left(\prod_k q(\theta_k) \right)$$

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

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

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

Recall from Lecture 20: To maximize the ELBO of a factorized approximation, compute the **mean field**

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

$$\begin{aligned} \log q^*(c_{di}) &= \mathbb{E}_{\prod_{d,j \neq i} q(c_{dj}) q(\pi_d) \prod_k q(\theta_k)} \left(\sum_{d,k} (\alpha_{dk} - 1 + n_{dk\cdot}) \log \pi_{dk} + \sum_{k,v} (\beta_{kv} - 1 + n_{\cdot kv}) \log \theta_{kv} \right) + \text{const.} \\ &= \sum_{k=1}^K c_{dik} \underbrace{\left(\mathbb{E}_{q(\pi_{dk})} (\log \pi_{dk}) + \mathbb{E}_{q(\theta_{di})} (\log \theta_{kw_{di}}) \right)}_{=: \log \gamma_{dik}} + \text{const.} \end{aligned}$$

Thus, $q(c_{di}) = \prod_k \tilde{\gamma}_{dik}^{c_{dik}}$, where $\tilde{\gamma}_{dik} = \gamma_{dik} / \sum_k \gamma_{dik}$ (Note: Thus, $\mathbb{E}_q(c_{dik}) = \tilde{\gamma}_{dik}$)

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

Recall from Lecture 20: To maximize the ELBO of a factorized approximation, compute the **mean field**

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

$$\begin{aligned} \log q^*(\boldsymbol{\pi}_d) &= \mathbb{E}_{\prod_{e \neq d, i} q(c_{ei}) q(\boldsymbol{\pi}_e) \prod_k q(\boldsymbol{\theta}_k)} \left(\sum_{d,k} (\alpha_{dk} - 1 + n_{dk\cdot}) \log \pi_{dk} + \sum_{k,v} (\beta_{kv} - 1 + n_{\cdot kv}) \log \theta_{kv} \right) + \text{const.} \\ &= \sum_{k=1}^K (\alpha_{dk} - 1 + \mathbb{E}_{q(C)}(n_{dk\cdot})) \log \pi_{dk} + \text{const.} \end{aligned}$$

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

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

Recall from Lecture 20: To maximize the ELBO of a factorized approximation, compute the **mean field**

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

$$\begin{aligned} \log q^*(\boldsymbol{\theta}_k) &= \mathbb{E}_{\prod_{d,i} q(c_{di}) q(\boldsymbol{\pi}_d) \prod_{\ell \neq k} q(\boldsymbol{\theta}_\ell)} \left(\sum_{d,k} (\alpha_{dk} - 1 + n_{dk\cdot}) \log \pi_{dk} + \sum_{k,v} (\beta_{kv} - 1 + n_{\cdot kv}) \log \theta_{kv} \right) + \text{const.} \\ &= \sum_{v=1}^V (\beta_{kv} - 1 + \mathbb{E}_{q(C)}(n_{\cdot kv})) \log \theta_{kv} + \text{const.} \end{aligned}$$

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

Properties of the Dirichlet

(let $\hat{\alpha} := \sum_d \alpha_d$)

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

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

Where $F(z) = \frac{d}{dz} \log \Gamma(z)$ (the “digamma-function”). `scipy.special.digamma(z)` <https://dlmf.nist.gov/5>

$$\begin{aligned} q(\boldsymbol{\pi}_d) &= \mathcal{D} \left(\boldsymbol{\pi}_d; \tilde{\alpha}_{dk} := \left[\alpha_{dk} + \sum_{i=1}^{l_d} \tilde{\gamma}_{dik} \right]_{k=1, \dots, K} \right) \\ q(\boldsymbol{\theta}_k) &= \mathcal{D} \left(\boldsymbol{\theta}_k; \tilde{\beta}_{kv} := \left[\beta_{kv} + \sum_d^D \sum_{i=1}^{l_d} \tilde{\gamma}_{dik} \mathbb{I}(w_{di} = v) \right]_{v=1, \dots, V} \right) \\ q(\mathbf{c}_{di}) &= \prod_k \tilde{\gamma}_{dik}^{c_{dik}}, \quad \text{where} \quad \tilde{\gamma}_{dik} = \gamma_{dik} / \sum_k \gamma_{dik} \quad \text{and (note that } \sum_k \tilde{\alpha}_{dk} = \text{const.)} \\ \gamma_{dik} &= \exp \left(\mathbb{E}_{q(\boldsymbol{\pi}_{dk})} (\log \pi_{dk}) + \mathbb{E}_{q(\boldsymbol{\theta}_{di})} (\log \theta_{kw_{di}}) \right) \\ &= \exp \left(F(\tilde{\alpha}_{jk}) + F(\tilde{\beta}_{kw_{di}}) - F \left(\sum_v \tilde{\beta}_{kv} \right) \right) \end{aligned}$$

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

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.

Lectures 18–20.

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

Designing a Probabilistic Machine Learning Model

1. Take a close look at the **Data**
2. Think about modelling goals, decide on data types
3. Design the **Model**
4. Design the **Algorithm**
 - ✦ For conditionally independent parts, conjugate priors may give analytic inference
 - ✦ for linearly related variables, consider Gaussians
 - ✦ consider revisiting the model to simplify as much as possible
 - ✦ if there is still no analytic answer, use **The Toolbox!**