

# PROBABILISTIC MACHINE LEARNING

## LECTURE 22

### SUMMARY AND CLEANUP

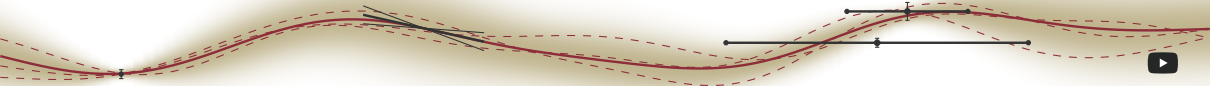
Philipp Hennig

17 July 2023

EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN



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





# The course so far

Intermediate summary, more to come in the final lectures



# Probabilities: the language of reasoning under uncertainty

Lectures 1–3

- ▶ We can describe all inference tasks by assigning **probabilities**, (or, for continuous variables probability density functions) jointly to all variables in the problem.
- ▶ Probabilities and pdfs satisfy “the rules of probability”:

$$\int_{\mathbb{R}^d} p(x) dx = 1$$

$$p_{X_1}(x_1) = \int_{\mathbb{R}} p_X(x_1, x_2) dx_2$$

sum rule

$$p(x_1 | x_2) = \frac{p(x_1, x_2)}{p(x_2)}$$

product rule

$$p(x_1 | x_2) = \frac{p(x_1) \cdot p(x_2 | x_1)}{\int p(x_1) \cdot p(x_2 | x_1) dx_1}$$

Bayes' Theorem.



## Definition (Exponential Family, simplified form)

Consider a random variable  $X$  taking values  $x \in \mathbb{X} \subset \mathbb{R}^n$ . A probability distribution for  $X$  with pdf of the functional form

$$p_w(x) = h(x) \exp [\phi(x)^\top w - \log Z(w)] = \frac{h(x)}{Z(w)} e^{\phi(x)^\top w} = p(x \mid w)$$

is called an **exponential family** of probability measures. The function  $\phi : \mathbb{X} \rightarrow \mathbb{R}^d$  is called the **sufficient statistics**. The parameters  $w \in \mathbb{R}^d$  are the **natural parameters** of  $p_w$ . The normalization constant  $Z(w) : \mathbb{R}^d \rightarrow \mathbb{R}$  is the **partition function**. The function  $h(x) : \mathbb{X} \rightarrow \mathbb{R}_+$  is the **base measure**. For notational convenience, it can be useful to re-parametrize the natural parameters  $w$  as  $w := \eta(\theta)$  in terms of *canonical parameters*  $\theta$ .

# Exponential Families: typed reasoning

Lectures 4–6

Exponential Families have Conjugate Priors

- Consider the exponential family  $p_w(x | w) = h(x) \exp [\phi(x)^\top w - \log Z(w)]$
- its conjugate prior is the exponential family  $F(\alpha, \nu) = \int \exp(\alpha^\top w - \nu \log Z(w)) dw$

$$p_\alpha(w | \alpha, \nu) = \exp \left[ \left( \begin{matrix} w \\ -\log Z(w) \end{matrix} \right)^\top \begin{pmatrix} \alpha \\ \nu \end{pmatrix} - \log F(\alpha, \nu) \right]$$

$$\text{because } p_\alpha(w | \alpha, \nu) \prod_{i=1}^n p_w(x_i | w) \propto p_\alpha \left( w \middle| \alpha + \sum_i \phi(x_i), \nu + n \right)$$

- and the predictive is

$$\begin{aligned}
 p(x) &= \int p_w(x | w) p_\alpha(w | \alpha, \nu) dw = h(x) \int e^{(\phi(x) + \alpha)^\top w + (\nu + 1) \log Z(w) - \log F(\alpha, \nu)} dw \\
 &= h(x) \frac{F(\phi(x) + \alpha, \nu + 1)}{F(\alpha, \nu)}
 \end{aligned}$$

Computing  $F(\alpha, \nu)$  can be tricky. In general, this is **the** challenge when constructing an EF.

# Gaussians: Inference as Linear Algebra

Lectures 5–6

- ▶ products of Gaussians are Gaussians

$$\begin{aligned}\mathcal{N}(x; a, A) \mathcal{N}(x; b, B) \\ = \mathcal{N}(x; c, C) \mathcal{N}(a; b, A + B)\end{aligned}$$

$$C := (A^{-1} + B^{-1})^{-1} \quad c := C(A^{-1}a + B^{-1}b)$$

- ▶ linear projections of Gaussians are Gaussians

$$\begin{aligned}p(z) &= \mathcal{N}(z; \mu, \Sigma) \\ \Rightarrow p(Az) &= \mathcal{N}(Az, A\mu, A\Sigma A^T)\end{aligned}$$

- ▶ marginals of Gaussians are Gaussians

$$\int \mathcal{N} \left[ \begin{pmatrix} x \\ y \end{pmatrix}; \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix} \right] dy = \mathcal{N}(x; \mu_x, \Sigma_{xx})$$

- ▶ (linear) conditionals of Gaussians are Gaussians

$$\begin{aligned}p(x | y) &= \frac{p(x, y)}{p(y)} \\ &= \mathcal{N} \left( x; \mu_x + \Sigma_{xy} \Sigma_{yy}^{-1} (y - \mu_y), \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx} \right)\end{aligned}$$

## Bayesian inference becomes linear algebra

If  $p(x) = \mathcal{N}(x; \mu, \Sigma)$  and  $p(y | x) = \mathcal{N}(y; A^T x + b, \Lambda)$ , then

$$p(B^T x + c | y) = \mathcal{N}[B^T x + c; B^T \mu + c + B^T \Sigma A (A^T \Sigma A + \Lambda)^{-1} (y - A^T \mu - b), B^T \Sigma B - B^T \Sigma A (A^T \Sigma A + \Lambda)^{-1} A^T \Sigma B]$$

$$\text{prior} \quad p(w) = \mathcal{N}(w; \mu, \Sigma) \quad \Rightarrow \quad p(f) = \mathcal{N}(f_x; \phi_x^\top \mu, \phi_x \Sigma \phi_x)$$

$$\text{likelihood} \quad p(y \mid w, \phi_x) = \mathcal{N}(y; \phi_x^\top w, \sigma^2 l) = \mathcal{N}(y; f_x, \sigma^2 l)$$

$$\begin{aligned} \text{posterior on } w \quad p(w \mid y, \phi_x) &= \mathcal{N}(w; \mu + \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 l)^{-1} (y - \phi_x^\top \mu), \\ &\quad \Sigma - \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 l)^{-1} \phi_x^\top \Sigma) \\ &= \mathcal{N}\left(w; (\Sigma^{-1} + \sigma^{-2} \phi_x \phi_x^\top)^{-1} \left( \Sigma^{-1} \mu + \sigma^{-2} \phi_x y \right), \right. \\ &\quad \left. (\Sigma^{-1} + \sigma^{-2} \phi_x \phi_x^\top)^{-1} \right) \end{aligned}$$

$$\begin{aligned} \text{posterior on } f \quad p(f_x \mid y, \phi_x) &= \mathcal{N}(f_x; \phi_x^\top \mu + \phi_x^\top \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 l)^{-1} (y - \phi_x^\top \mu), \\ &\quad \phi_x^\top \Sigma \phi_x - \phi_x^\top \Sigma \phi_x (\phi_x^\top \Sigma \phi_x + \sigma^2 l)^{-1} \phi_x^\top \Sigma \phi_x) \\ &= \mathcal{N}\left(f_x; \phi_x (\Sigma^{-1} + \sigma^{-2} \phi_x \phi_x^\top)^{-1} \left( \Sigma^{-1} \mu + \sigma^{-2} \phi_x y \right), \right. \\ &\quad \left. \phi_x (\Sigma^{-1} + \sigma^{-2} \phi_x \phi_x^\top)^{-1} \phi_x^\top \right) \end{aligned}$$



$$\begin{aligned} p(f(\bullet) \mid \mathbf{w}) &= \mathcal{N}(f(\bullet); \phi(\bullet)^\top \mathbf{w}, \sigma I) \\ p(f(\bullet)) &= \int p(f(\bullet) \mid \mathbf{w}) p(\mathbf{w}) d\mathbf{w} \\ &= \mathcal{N}(f(\bullet); \phi(\bullet)^\top \boldsymbol{\mu}, \phi(\bullet)^\top \boldsymbol{\Sigma} \phi(\bullet) + \sigma I) \end{aligned}$$

using the abstraction / encapsulation

$$m(\bullet) := \phi(\bullet)^\top \boldsymbol{\mu}$$

$$m : \mathbb{X} \rightarrow \mathbb{R}$$

mean function

$$k(\bullet, \circ) := \phi(\bullet)^\top \boldsymbol{\Sigma} \phi(\circ)$$

$$k : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$$

covariance function, aka. **kernel**



# Algorithm 1 Linear Algebra for/as Gaussian (process) inference – efficient data-loading and book-keeping

**Input:** sufficient statistics  $K = k_{XX} + \sigma^2 I$ ,  $\bar{y} = y - \mu_X$ , initial guesses  $\alpha_0, C_0$

**Output:** defragmented statistics  $S, C_i, \alpha_i$

```

1  procedure TRAIN( $K, y, C_0 = 0, \alpha_0 = 0$ )
2    for  $i \in \{1, \dots, n\}$  do
3       $s_i \leftarrow \text{POLICY}(S = [s_{j < i}], Z = [z_{j < i}])$  // Action – load,  $s_i \in \mathbb{R}^{N \times k_i}$ 
4       $z_i \leftarrow K s_i$  // Observation – compute,  $z_i \in \mathbb{R}^{N \times k_i}$ 
5       $d_i \leftarrow (I - C_{i-1} K) s_i = s_i - C_{i-1} z_i$  // low-rank update,  $d_i \in \mathbb{R}^{N \times k_i}$ 
6       $H_i \leftarrow s_i^\top K d_i = z_i^\top d_i$  // Schur complement,  $H_i \in \mathbb{R}^{k_i \times k_i}$ 
7       $C_i \leftarrow C_{i-1} + d_i H_i^{-1} d_i^\top$  // Inverse estimate
8       $\alpha_i \leftarrow C_i y = \alpha_{i-1} + d_i H_i^{-1} d_i^\top \bar{y}$  // Solution Estimate
9    end for
10   return  $S = [s_j]_{j \leq i}, \alpha_i, C_i$ 
11 end procedure
12 procedure PREDICT( $x, S, \alpha, C$ )
13    $k_{xS} \leftarrow k[x, S]$  // Covariance to Observations
14    $\mu_x \leftarrow \mu_x + k_{xS} \alpha$  // Point estimate
15    $v_{xx} \leftarrow k_{xx} - k_{xS} C k_{Sx}$  // Uncertainty
16 end procedure

```

$$p(f) = \mathcal{GP}(f; m, k) \quad p(y | f_x) = \sigma(yf_x) = \begin{cases} \sigma(f) & \text{if } y = 1 \\ 1 - \sigma(f) & \text{if } y = -1 \end{cases} \quad \text{using } \sigma(x) = 1 - \sigma(-x).$$

- Find maximum posterior probability for **latent  $f$**  at **training points**

$$\hat{f} = \arg \max \log p(\mathbf{f}_X | y)$$

- Assign approximate Gaussian posterior at training points

$$q(\mathbf{f}_X) = \mathcal{N}(\mathbf{f}_X; \hat{f}, -(\nabla \nabla^\top \log p(\mathbf{f}_X | y)|_{\mathbf{f}_X = \hat{f}})^{-1}) =: \mathcal{N}(\mathbf{f}_X; \hat{f}, \hat{\Sigma})$$

- approximate posterior **predictions** at  $f_x$  for **latent function**

$$\begin{aligned} q(f_x | y) &= \int p(f_x | \mathbf{f}_X) q(\mathbf{f}_X) d\mathbf{f}_X = \int \mathcal{N}(f_x; m_x + k_{xX} K_{XX}^{-1} (\mathbf{f}_X - m_X), k_{xx} - k_{xX} K_{XX}^{-1} k_{Xx}) q(\mathbf{f}_X) d\mathbf{f}_X \\ &= \mathcal{N}(f_x; m_x + k_{xX} K_{XX}^{-1} (\hat{f} - m_X), k_{xx} - k_{xX} K_{XX}^{-1} k_{Xx} + k_{xX} K_{XX}^{-1} \hat{\Sigma} K_{XX}^{-1} k_{Xx}) \end{aligned}$$

- compute predictions for **label probabilities**:

$$\mathbb{E}_{p(f|y)}[\pi_x] \approx \mathbb{E}_q[\pi_x] = \int \sigma(f_x) q(f_x | y) df_x$$

# Deep Learning: Any deep network as a GP

Laplace approximations as a general tool

1. Realise that the loss is a **negative log-posterior**

$$\mathcal{L}(\boldsymbol{\theta}) = \left( \underbrace{\frac{1}{N} \sum_{i=1}^N \ell(y_i; \overbrace{f(x_i, \boldsymbol{\theta})}^{\text{deep net}})}_{\text{empirical risk}} + \underbrace{r(\boldsymbol{\theta})}_{\text{regularizer}} \right) = - \sum_{i=1}^N \log p(y | \boldsymbol{\theta}) - \log p(\boldsymbol{\theta}) = - \log p(\boldsymbol{\theta} | y) + \text{const.}$$

2. Train the deep net *as usual* to find  $\boldsymbol{\theta}_* = \arg \max_{\boldsymbol{\theta} \in \mathbb{R}^D} p(\boldsymbol{\theta} | y)$
3. At  $\boldsymbol{\theta}_*$ , compute a **Laplace approximation** of the log-posterior, with  $\Psi := \nabla \nabla^\top \log p(\boldsymbol{\theta}_* | y)$

$$\log p(\boldsymbol{\theta} | y) + \text{const.} = \mathcal{L}(\boldsymbol{\theta}) \approx \mathcal{L}(\boldsymbol{\theta}_*) + 1/2 (\boldsymbol{\theta} - \boldsymbol{\theta}_*)^\top \Psi (\boldsymbol{\theta} - \boldsymbol{\theta}_*) = \log \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\theta}_*, -\Psi^{-1})$$

4. Linearize  $f(x, \boldsymbol{\theta})$  around  $\boldsymbol{\theta}_*$ , with  $[J(x)]_{ij} = \frac{\partial f_i(x, \boldsymbol{\theta}_*)}{\partial \theta_j}$  as  $f(x, \boldsymbol{\theta}) \approx f(x, \boldsymbol{\theta}_*) + J(x, \boldsymbol{\theta}_*)(\boldsymbol{\theta} - \boldsymbol{\theta}_*)$

$$\text{thus } p(f(\bullet) | \mathcal{D}) = \int p(f | w) dp(w) \approx \mathcal{GP}(f(\bullet); f(\bullet, \boldsymbol{\theta}_*), -J(\bullet)\Psi^{-1}J(\circ)) \quad \text{with}$$

$$\mathbb{E}(f(\bullet)) = f(\bullet, \boldsymbol{\theta}_*) \quad \text{the trained net as the mean function}$$

$$\text{cov}(f(\bullet), f(\circ)) = -J(\bullet)\Psi^{-1}J(\circ)^\top \quad \text{the Laplace tangent kernel as the covariance function}$$



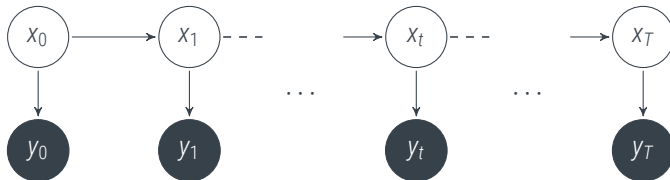
# Markov Chains: $\mathcal{O}(T)$ inference in time series

Lectures 20/21

Assume:

$$p(x_t | X_{0:t-1}) = p(x_t | x_{t-1})$$

and  $p(y_t | X) = p(y_t | x_t)$



$$p(x_t | Y_{1:t-1}) = \int p(x_t | x_{t-1}) p(x_{t-1} | Y_{1:t-1}) dx_{t-1}$$

$$p(x_t | Y_{0:t}) = \frac{p(y_t | x_t) p(x_t | Y_{0:t-1})}{\int p(y_t | x_t) p(x_t | Y_{0:t-1}) dx_t}$$

$$p(x_t | Y) = p(x_t | Y_{0:t}) \int p(x_{t+1} | x_t) \frac{p(x_{t+1} | Y)}{p(x_{t+1} | Y_{1:t})} dx_{t+1}$$

# 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:

- ▶ Directed Graphical Models
- ▶ Exponential Families (also as likelihoods)
- ▶ Gaussian Distributions
- ▶ Kernels
- ▶ Markov Chains
- ▶ Deep Networks

## Computation:

- ▶ autodiff
- ▶ MAP with Laplace approximations
- ▶ Linear algebra as a computational primitive
- ▶
- ▶



# Bayesian Hierarchical Learning

beyond analytical inference



# Bayesian Hierarchical Learning

the formalism for training hyperparameters / fit architectures & models

$$p(f | y, x, \theta) = \frac{p(y | f, x, \theta)p(f |, \theta)}{\int p(y | f, x, \theta)p(f |, \theta) df} = \frac{p(y | f, x, \theta)p(f |, \theta)}{p(y | x, \theta)}$$

► Model parameters like  $\theta$  are also known as **hyper-parameters**.

► This is largely a computational, practical distinction:

**data** are observed

→ condition

**variables** are the things we care about

→ full probabilistic treatment

**parameters** are the things we have to deal with to get the model right

→ integrate out

**hyper-parameters** are the top-level, too expensive to properly infer

→ fit

The **model evidence** in Bayes' Theorem is the (marginal) **likelihood** for the model. So we would like

$$p(\theta | y) = \frac{p(y | \theta)p(\theta)}{\int p(y | \theta')p(\theta') d\theta'}$$



# The Gaussian Base Case

Bayesian Hierarchical Inference for GP regression, from Lecture 10

$$p(f | \theta) = \mathcal{GP}(f; m_\theta, k_\theta) \quad \text{e.g.} \quad m_\theta(\bullet) = \phi(\bullet)^\top \theta, \text{ or } k_\theta(\bullet, \circ) = \theta_1 \exp\left(-\frac{(\bullet - \circ)^2}{2\theta_2^2}\right).$$

- The **evidence** in Bayes' theorem is the **marginal likelihood for the model**

$$p(f | y, x, \theta) = \frac{p(y | f, x, \theta) p(f | \theta)}{\int p(y | f, x, \theta) p(f | \theta) df} = \frac{p(y | f, x, \theta) p(f | \theta)}{p(y | x, \theta)}$$

- For Gaussians *and Gaussian processes*, the evidence has **analytic form**:

$$\underbrace{\mathcal{N}(y; \phi_X^{\theta^\top} w + b, \Lambda)}_{p(y|f,x,\theta)} \cdot \underbrace{\mathcal{N}(w, \mu, \Sigma)}_{p(f)} = \underbrace{\mathcal{N}(w; m_{\text{post}}^\theta, V_{\text{post}}^\theta)}_{p(f|y,x,\theta)} \cdot \underbrace{\mathcal{N}(y; \phi_X^{\theta^\top} \mu + b, \phi_X^{\theta^\top} \Sigma \phi_X^\theta + \Lambda)}_{p(y|\theta,x)}$$

$$\mathcal{N}(y; f^\theta(X), \Lambda^\theta) \cdot \mathcal{GP}(f, \mu^\theta, k^\theta) = \mathcal{GP}(f; m_{\text{post}}^\theta, V_{\text{post}}^\theta) \cdot \mathcal{N}(y; \mu^\theta(X), \Lambda^\theta + k^\theta(X, X))$$



$$\begin{aligned}
 \hat{\theta} &= \arg \max_{\theta} p(\mathbf{y} \mid \mathbf{x}, \theta) = \arg \max_{\theta} \int \mathcal{N}(\mathbf{y}; f(\mathbf{X}), \Lambda_{\theta}) \cdot \mathcal{N}(f_{\mathbf{X}}, \mu_{\mathbf{X}}, k_{\mathbf{X}\mathbf{X}}) df_{\mathbf{X}} \\
 &= \arg \max_{\theta} \mathcal{N}(\mathbf{y}; \mu_{\mathbf{X}}; k_{\mathbf{X}\mathbf{X}} + \Lambda_{\theta}) \int \mathcal{N}(f_{\mathbf{X}}; \mu_{\mathbf{y}, \mathbf{X}}, v_{\mathbf{y}, \mathbf{X}\mathbf{X}}) df_{\mathbf{X}} \\
 &= \arg \max_{\theta} \mathcal{N}(\mathbf{y}; \mu_{\mathbf{X}}^{\theta}, k_{\mathbf{X}\mathbf{X}}^{\theta} + \Lambda^{\theta}) \\
 &= \arg \max_{\theta} \log \mathcal{N}(\mathbf{y}; \mu_{\mathbf{X}}^{\theta}, k_{\mathbf{X}\mathbf{X}}^{\theta} + \Lambda^{\theta}) \\
 &= \arg \min_{\theta} -\log \mathcal{N}(\mathbf{y}; \mu_{\mathbf{X}}^{\theta}, k_{\mathbf{X}\mathbf{X}}^{\theta} + \Lambda^{\theta}) \\
 &= \arg \min_{\theta} \frac{1}{2} \left( \underbrace{(\mathbf{y} - \mu_{\mathbf{X}}^{\theta})^{\top} (k_{\mathbf{X}\mathbf{X}}^{\theta} + \Lambda^{\theta})^{-1} (\mathbf{y} - \phi_{\mathbf{X}}^{\theta \top} \mu)}_{\text{square error}} + \underbrace{\log |k_{\mathbf{X}\mathbf{X}}^{\theta} + \Lambda^{\theta}|}_{\text{model complexity / Occam factor}} \right) + \frac{N}{2} \log 2\pi
 \end{aligned}$$

*Numquam ponenda est pluralitas sine necessitate.  
Plurality must never be posited without necessity.*

William of Occam  
(1285 (Occam, Surrey)–1349 (Munich, Bavaria))  
stained-glass window by Lawrence Lee

# The Evidence framework beyond GP regression

Laplace approximations yield approximate evidences.

- For general likelihoods, the evidence  $p(\mathbf{y} \mid \boldsymbol{\theta}) = \int p(\mathbf{y} \mid \mathbf{f}, \boldsymbol{\theta}) p(\mathbf{f} \mid \boldsymbol{\theta}) d\mathbf{f}$  will be intractable
- If we approximate it with a Gaussian by a Laplace approximation, we need to be careful with the constant term. Say we have found  $\mathbf{f}^* = \arg \max \log p(\mathbf{f} \mid \mathbf{y})$  and  $\Psi = -\nabla \nabla \log p(\mathbf{f} \mid \mathbf{y})$ . Then

$$\begin{aligned} p(\mathbf{y} \mid \boldsymbol{\theta}) &= \int p(\mathbf{y} \mid \mathbf{f}) p(\mathbf{f}) d\mathbf{f} \\ &\approx \int \exp \left( \log p(\mathbf{y} \mid \mathbf{f}^*) + \log p(\mathbf{f}^*) - \frac{1}{2} (\mathbf{f} - \mathbf{f}^*) \Psi (\mathbf{f} - \mathbf{f}^*) \right) d\mathbf{f} \\ &= p(\mathbf{y} \mid \mathbf{f}^*) p(\mathbf{f}^*) \cdot (2\pi)^{D/2} |\Psi|^{1/2} \end{aligned}$$

- If the prior happens to be Gaussian  $p(\mathbf{f}) = \mathcal{N}(\mathbf{f}; \mathbf{m}, K)$ , we get

$$\log p(\mathbf{y} \mid \boldsymbol{\theta}) \approx \log p(\mathbf{y} \mid \mathbf{f}^*) - \frac{1}{2} (\mathbf{f}^* - \mathbf{m}) K^{-1} (\mathbf{f}^* - \mathbf{m}) - \frac{1}{2} \log |B|$$

$$\text{with } |B| = |K| \cdot |K^{-1} + \underbrace{\nabla \nabla^\top \log p(\mathbf{y} \mid \mathbf{f}^*)}_{=: W}|^{-1} = |I + K \cdot W|.$$

# A special option: The EM algorithm

Maximizing model Evidence

The general recipe for hyperparameter inference:

Consider a model with **parameters**  $\theta$ , **observed data**  $y$  and **latent variables**  $z$

- Ideally, we would like to maximize the **marginal (log-) likelihood (evidence)**

$$\log p(y | \theta) = \log \left( \int p(y, z | \theta) dz \right) \quad (\star)$$

- if we can not do this integral, we can try **Laplace** (as above). This is nearly always *possible* (if  $\log p(y | z)$  is twice differentiable), but it is fundamentally an approximation.
- however, *in some cases*, we may be able to compute the **Expectation** of the “complete data” log **likelihood** (for a fixed value  $\theta_*$ )

$$q(\theta, \theta_*) = \int p(z | y, \theta_*) \log p(y, z | \theta) dz$$

and then **Maximize**  $q(\theta, \theta_*)$  **with respect to**  $\theta$ . This can be easier than  $(\star)$  because the log “simplifies things” (e.g. turns products into sums, thus factors into components).

Definition: The Expectation Maximization (EM) algorithm:

Consider a model with parameters  $\theta$ , observed data  $y$  and latent variables  $z$ .

while not converged, do:

E compute the Expected complete data log-likelihood

$$q(\theta, \theta_t) = \int p(z \mid y, \theta_t) \log p(y, z \mid \theta) dz$$

M Set  $\theta_{t+1}$  to Maximize  $\theta_{t+1} = \arg \max_{\theta} q(\theta, \theta_{t+1})$ .

We will see on Thursday why this is a meaningful thing to do.

# Example: EM for Gauss-Markov Models

completed in this week's exercise

Consider the Gauss-Markov Model with  $\mathbf{z} := [\mathbf{x}_t]_{t=0,\dots,T}$  and  $\boldsymbol{\theta} := (A, Q, H, R)$ .

$$p(\mathbf{z}, \mathbf{y} \mid \boldsymbol{\theta}) = p(\mathbf{x}_0) \prod_{t=1}^T p(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \boldsymbol{\theta}) p(\mathbf{y}_t \mid \mathbf{x}_t, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_0; m_0, P_0) \prod_{t=1}^T \mathcal{N}(\mathbf{x}_t; A\mathbf{x}_{t-1}, Q) \mathcal{N}(\mathbf{y}_t; H\mathbf{x}_t, R)$$

Here it is actually possible to compute  $p(\mathbf{y} \mid \boldsymbol{\theta})$  (see last lectures), but that term can only be maximized numerically. Instead, notice that the complete-data log-likelihood neatly separates into local terms

$$\log p(\mathbf{z}, \mathbf{y} \mid \boldsymbol{\theta}) = \log \mathcal{N}(\mathbf{x}_0) + \sum_{t=1}^T \log \mathcal{N}(\mathbf{x}_t; A\mathbf{x}_{t-1}, Q) + \sum_{t=1}^T \log \mathcal{N}(\mathbf{y}_t; H\mathbf{x}_t, R)$$

And the expectation of a log-Gaussian is easy to compute, because

$$\int (\mathbf{W}\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{W}\mathbf{x} - \boldsymbol{\mu}) \cdot \mathcal{N}(\mathbf{x}; \mathbf{m}, V) d\mathbf{x} = (\boldsymbol{\mu} - \mathbf{W}\mathbf{m})^\top \Sigma^{-1} (\boldsymbol{\mu} - \mathbf{W}\mathbf{m}) + \text{tr}(\mathbf{W}^\top \Sigma^{-1} \mathbf{W} V)$$

Rest: Homework.

## Summary:

- The *Evidence* in Bayes' theorem is the *marginal likelihood* of the model

$$p(\mathbf{y} \mid \theta) = \int p(\mathbf{y}, \mathbf{z} \mid \theta) d\mathbf{z}$$

- In principle it provides the “next level” for Bayesian inference on  $\theta$ . But it is often intractable.
- Laplace approximations provide a *general, approximate* way to approximate the Evidence
- in *some* models, the **EM** algorithm offers a tractable iterative solution.

Please cite this course, as

```
@techreport{Tuebingen_ProbML23,
  title =
    {Probabilistic Machine Learning},
  author = {Hennig, Philipp},
  series = {Lecture Notes
    in Machine Learning},
  year = {2023},
  institution = {Tübingen AI Center}}
```

