

Inverse Problems

Sommersemester 2022

Vesa Kaarnioja
vesa.kaarnioja@fu-berlin.de

FU Berlin, FB Mathematik und Informatik

Tenth lecture, July 4, 2022

Today's lecture

- Optimization perspective (chapter 3 in Sanz-Alonso, Stuart, Taeb 2018)
- Gaussian approximation (chapter 4 in Sanz-Alonso, Stuart, Taeb 2018)

The setting

We work in the inverse problem setting of finding $x \in \mathbb{R}^d$ from $y \in \mathbb{R}^k$ given by

$$y = F(x) + \eta$$

with noise $\eta \sim \nu$ and prior $x \sim \pi$ such that $\eta \perp x$. The posterior π^y on $x|y$ is given by Bayes' theorem

$$\pi^y(x) = \frac{1}{Z} \nu(y - F(x)) \pi(x).$$

We have the negative log-likelihood:

$$L(x) = -\log \nu(y - F(x)),$$

and a regularizer

$$R(x) = -\log \pi(x).$$

When added together these two functions of x comprise an objective function of the form

$$J(x) = L(x) + R(x).$$

Furthermore

$$\pi^y(x) = \frac{1}{Z} \nu(y - F(x)) \pi(x) \propto e^{-J(x)}.$$

We see that minimizing the objective function J is equivalent to maximizing the posterior π^y . Therefore, the MAP estimator can be rewritten in terms of J as follows:

$$\hat{x}_{\text{MAP}} = \arg \max_{x \in \mathbb{R}^d} \pi^y(x) = \arg \min_{x \in \mathbb{R}^d} J(x).$$

Let us consider conditions under which the MAP estimator is attained, and characterize the MAP estimator in terms of small ball probabilities – this interpretation generalizes the definition of MAP estimators to measures that do not possess a Lebesgue density.

For any optimization problem for an objective function with a finite infimum, it is of interest to determine whether the infimum is attained.

Theorem (Attainable MAP estimator)

Assume that J is non-negative, continuous and that $J(x) \rightarrow \infty$ as $|x| \rightarrow \infty$. Then J attains its infimum. Therefore, the MAP estimator of x based on the posterior $\pi^y(x) \propto \exp(-J(x))$ is attained.

Proof.

By the assumed growth and non-negativity of J , there is R such that $\inf_{x \in \mathbb{R}^d} J(x) = \inf_{x \in \bar{B}(0, R)} J(x)$ where $\bar{B}(0, R)$ denotes the closed ball of radius R around the origin. Since J is assumed to be continuous, its infimum over $\bar{B}(0, R)$ is attained and the proof is complete. \square

Remark. The assumption that $J(x) \rightarrow \infty$ is not restrictive: this condition needs to hold in order to be able to normalize $\pi^y(x) \propto \exp(-J(x))$ into a PDF, which is implicitly assumed in the second part of the theorem statement.

Example. Suppose that

- 1 $F: \mathbb{R}^d \rightarrow \mathbb{R}^k$ is continuous and $\eta \sim \mathcal{N}(0, \Gamma)$;
- 2 the objective function $J(x) = L(x) + R(x)$ has Γ -weighted L^2 loss

$$L(x) = \frac{1}{2} \|y - F(x)\|_{\Gamma^{-1}}^2$$

and L^p regularizer

$$R(x) = \frac{\lambda}{p} \|u\|_p^p, \quad p \in (0, \infty).$$

Then the assumptions on J in the previous theorem are satisfied, and the infimum of J is attained at the MAP estimator of the corresponding Bayesian problem with posterior PDF proportional to $\exp(-J(u))$.

Intuitively the MAP estimator maximizes posterior probability. We make this precise in the following theorem which links the objective function J to small ball probabilities.

Theorem (Objective function and posterior probability)

Assume that J is non-negative, continuous and that $J(x) \rightarrow \infty$ as $|x| \rightarrow \infty$. Let

$$\alpha(x, \delta) := \int_{B(x, \delta)} \pi^y(v) dv = \mathbb{P}^{\pi^y}(B(x, \delta)),$$

be the posterior probability of a ball with radius δ centered at x . Then, for all $x_1, x_2 \in \mathbb{R}^d$, we have

$$\lim_{\delta \rightarrow 0} \frac{\alpha(x_1, \delta)}{\alpha(x_2, \delta)} = e^{J(x_2) - J(x_1)}.$$

Remark: For fixed x_2 , the right-hand side is maximized at point x_1 that minimizes J . Independently of the choice of any fixed x_2 , the above result shows that the probability of a small ball of radius δ centered at x_1 is, approximately, maximized by choosing the centre at a minimizer of J .

This result essentially characterizes the MAP estimate and, since it makes no reference to Lebesgue density, it can be generalized to infinite dimensions.

Proof. Let $x_1, x_2 \in \mathbb{R}^d$, $\varepsilon > 0$. By continuity of J , for all sufficiently small δ :

$$x \in \bar{B}(x_j, \delta) \Rightarrow |J(x) - J(x_j)| \leq \varepsilon, \quad j \in \{1, 2\},$$

and therefore

$$e^{-J(x_1)-\varepsilon} \leq e^{-J(v)} \leq e^{-J(x_1)+\varepsilon} \quad \text{for all } v \in B(x_1, \delta),$$

$$e^{-J(x_2)-\varepsilon} \leq e^{-J(v)} \leq e^{-J(x_2)+\varepsilon} \quad \text{for all } v \in B(x_2, \delta).$$

It follows, for all δ sufficiently small, that

$$B_\delta e^{-J(x_1)-\varepsilon} \leq \int_{B(x_1, \delta)} e^{-J(v)} \, dv \leq B_\delta e^{-J(x_1)+\varepsilon},$$

$$B_\delta e^{-J(x_2)-\varepsilon} \leq \int_{B(x_2, \delta)} e^{-J(v)} \, dv \leq B_\delta e^{-J(x_2)+\varepsilon},$$

where B_δ is the Lebesgue measure of a ball with radius δ . Taking the ratio of α 's and using the above bounds we obtain that, for all δ sufficiently small,

$$e^{J(x_2)-J(x_1)-2\varepsilon} \leq \frac{\alpha(x_1, \delta)}{\alpha(x_2, \delta)} \leq e^{J(x_2)-J(x_1)+2\varepsilon}.$$

Since ε was arbitrary, the desired result follows. □

Unimodal distributions

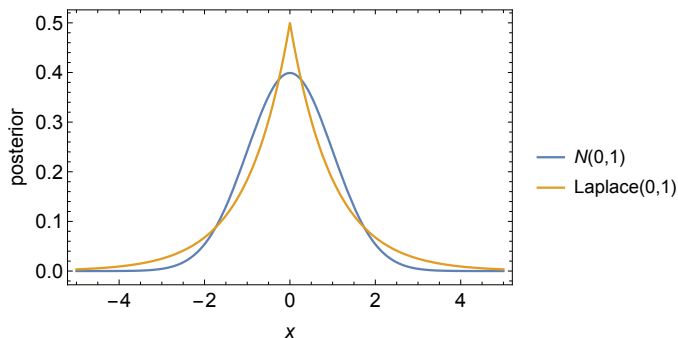


Figure: If the posterior is single-peaked, the MAP estimator reasonably summarizes the most likely value of the unknown parameter.

Problems with uneven distributions

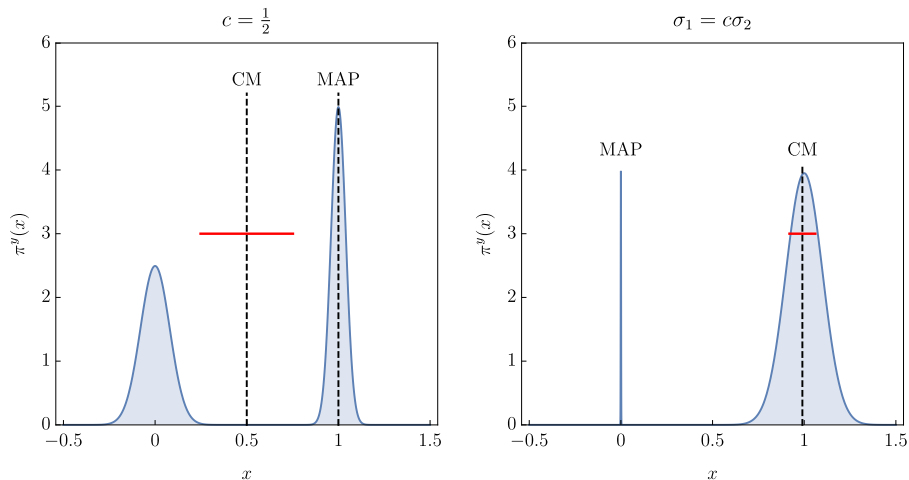


Figure: If the posterior is unevenly distributed, then it is less clear that the MAP or CM estimators usefully summarize the posterior.

Problems with rough distributions

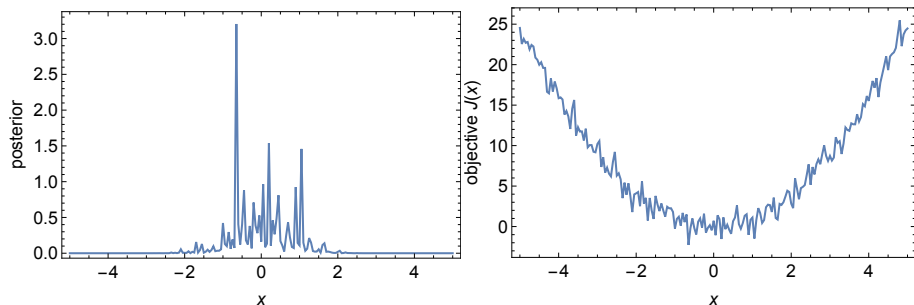


Figure: If the objective function $J(x)$ is very rough (here it is a quadratic function contaminated with white noise), then the resulting posterior density is very rough.

The objective function has small-scale roughness, but it has a larger pattern. The MAP estimator cannot capture this larger pattern as it is found by minimizing the objective function. Arguably, $x = 0$ might be a better point estimate.

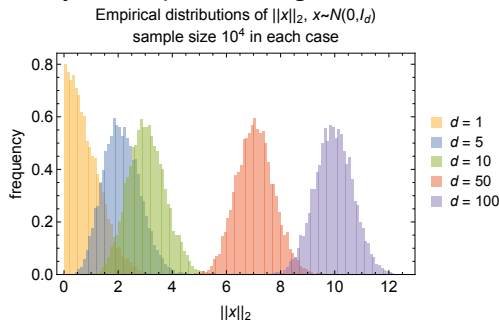
Problems with high dimension

Gaussian Annulus Theorem: Nearly all the probability of a d -dimensional spherical Gaussian distribution with unit variance is concentrated in a thin annulus of width $\mathcal{O}(1)$ at radius \sqrt{d} .

For example, if $x \sim \mathcal{N}(0, I_d)$, then

d	$\mathbb{P}(\ x\ < 5)$
10	0.99465
50	0.00119
100	1.135e-15

A point estimator may not capture enough information about the density.



Gaussian approximation

We explored the idea of obtaining a point estimator using an optimization perspective arising from maximizing the posterior PDF. This idea reduces the complexity of Bayesian inference from determination of an entire distribution to determination of a single point, however, the approach has some limitations, in particular for noisy, multi-peaked or high-dimensional posterior distributions.

Instead of seeking a point estimator, we can try seeking a Gaussian distribution $p = \mathcal{N}(\mu, \Sigma)$ that minimizes the Kullback–Leibler divergence from the posterior $\pi^y(x)$. Since the Kullback–Leibler divergence is not symmetric this leads to two distinct problems, which we will consider separately.

The Kullback–Leibler divergence

Definition

Let $\pi, \pi' > 0$ be two probability distributions on \mathbb{R}^d . The *Kullback–Leibler (KL) divergence*, or *relative entropy*, of π with respect to π' is defined by

$$\begin{aligned}d_{\text{KL}}(\pi \|\pi') &:= \int_{\mathbb{R}^d} \log \left(\frac{\pi(x)}{\pi'(x)} \right) \pi(x) \, dx \\ &= \mathbb{E}^{\pi} \left[\log \left(\frac{\pi}{\pi'} \right) \right] \\ &= \mathbb{E}^{\pi'} \left[\log \left(\frac{\pi}{\pi'} \right) \frac{\pi}{\pi'} \right].\end{aligned}$$

Kullback–Leibler is a divergence in that $d_{\text{KL}}(\pi \|\pi') \geq 0$, with equality if and only if $\pi = \pi'$ a.e. However, unlike Hellinger and total variation, it is not a distance. In particular, the KL divergence is not symmetric: in general

$$d_{\text{KL}}(\pi \|\pi') \neq d_{\text{KL}}(\pi' \|\pi).$$

The KL divergence is useful for at least the following reasons:

- it provides an upper bound for many distances;
- its logarithmic structure allows explicit computations that are difficult using actual distances;
- it satisfies many convenient analytical properties such as being convex in both arguments and lower-semicontinuous in the topology of weak convergence;
- it has an information-theoretic and physical interpretation.

Lemma

The KL divergence provides the following upper bounds for Hellinger and total variation distance:

$$d_H(\pi, \pi')^2 \leq \frac{1}{2} d_{\text{KL}}(\pi \parallel \pi'), \quad d_{\text{TV}}(\pi, \pi')^2 \leq d_{\text{KL}}(\pi \parallel \pi').$$

Proof. Recall from Week 6 that $\frac{1}{\sqrt{2}} d_{\text{TV}}(\pi, \pi') \leq d_H(\pi, \pi')$
 $\Leftrightarrow d_{\text{TV}}(\pi, \pi')^2 \leq 2d_H(\pi, \pi')^2$. Thus the second inequality follows from the first one. We prove only the first inequality.

Consider the function $\phi: \mathbb{R}_+ \rightarrow \mathbb{R}$ defined by

$$\phi(x) = x - 1 - \log x.$$

Note that

$$\phi'(x) = 1 - \frac{1}{x},$$

$$\phi''(x) = \frac{1}{x^2},$$

$$\lim_{x \rightarrow +\infty} \phi(x) = \phi(0) = \infty.$$

Thus the function is convex on its domain. As the minimum of ϕ is attained at $x = 1$, and as $\phi(1) = 0$, we deduce that $\phi(x) \geq 0$ for all $x \in (0, \infty)$. Hence,

$$x - 1 \geq \log x \quad \text{for all } x > 0,$$

$$\sqrt{x} - 1 \geq \frac{1}{2} \log x \quad \text{for all } x > 0.$$

We can use this last inequality to bound the Hellinger distance:

$$\begin{aligned}d_{\text{H}}(\pi, \pi')^2 &= \frac{1}{2} \int_{\mathbb{R}^d} \left(1 - \sqrt{\frac{\pi'}{\pi}}\right)^2 \pi \, dx \\&= \frac{1}{2} \int_{\mathbb{R}^d} \left(1 + \frac{\pi'}{\pi} - 2\sqrt{\frac{\pi'}{\pi}}\right) \pi \, dx \\&= 1 - \int_{\mathbb{R}^d} \sqrt{\frac{\pi'}{\pi}} \pi \, dx \\&= \int_{\mathbb{R}^d} \left(1 - \sqrt{\frac{\pi'}{\pi}}\right) \pi \, dx \\&\leq -\frac{1}{2} \int_{\mathbb{R}^d} \log\left(\frac{\pi'}{\pi}\right) \pi \, dx = \frac{1}{2} \int_{\mathbb{R}^d} \log\left(\frac{\pi}{\pi'}\right) \pi \, dx = \frac{1}{2} d_{\text{KL}}(\pi \parallel \pi').\end{aligned}$$

□

Best Gaussian approximation

Let π be the target distribution, e.g., the posterior. We consider two different minimization problems, both leading to a “best Gaussian”:

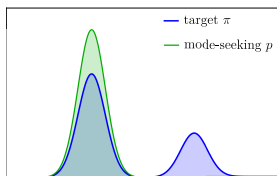
$$\inf_{p \in \mathcal{A}} d_{\text{KL}}(p \parallel \pi) \quad (\text{“Mode-seeking Gaussian approximation”})$$

and

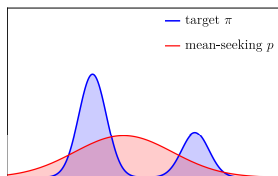
$$\inf_{p \in \mathcal{A}} d_{\text{KL}}(\pi \parallel p), \quad (\text{“Mean-seeking Gaussian approximation”})$$

where the minimization is performed over the set of Gaussian distributions on \mathbb{R}^d with positive definite covariance, i.e.,

$$\mathcal{A} := \{\mathcal{N}(\mu, \Sigma) \mid \mu \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d} \text{ positive definite}\}.$$



(a) Minimizing $d_{\text{KL}}(p||\pi)$



(b) Minimizing $d_{\text{KL}}(\pi||p)$

- Fig. (a): Minimizing $d_{\text{KL}}(p||\pi)$ may miss out components of π – we want $\log \frac{p}{\pi}$ to be small, which can happen when $p \approx \pi$ or $p \ll \pi$. Minimizing $d_{\text{KL}}(p||\pi)$ over Gaussians p can only give a single mode approximation which is achieved by matching one of the modes; we may think of this as “mode-seeking”.
- Fig. (b): Minimizing $d_{\text{KL}}(\pi||p)$ over Gaussians p we want $\log \frac{\pi}{p}$ to be small where p appears as the denominator. Wherever π has some mass we must let p also have some mass there in order to keep $\frac{\pi}{p}$ as close as possible to one. The mass of p is allocated in a way such that on average the divergence between p and π attains its minimum; hence, it can be thought of as “mean-seeking”.

Different applications will favor different choices between the mean and mode seeking approaches to Gaussian approximation.

Best Gaussian fit by minimizing $d_{\text{KL}}(p||\pi)$ (“mode-seeking”)

Theorem (Best Gaussian approximation / “mode-seeking”)

Suppose that the loss function $L(x) := -\log \nu(y - F(x))$ is non-negative and bounded above and that the prior $\pi \sim \mathcal{N}(0, \lambda^{-1}I)$. Then there exists at least one probability distribution $p \in \mathcal{A}$ at which the infimum

$$\inf_{p \in \mathcal{A}} d_{\text{KL}}(p||\pi^y)$$

is attained.

Proof. Let $p(x) = \frac{1}{(2\pi)^{d/2}(\det \Sigma)^{1/2}} e^{-\frac{1}{2}|x-\mu|_{\Sigma}^2}$, $\pi^y(x) = \frac{1}{Z} e^{-L(x) - \frac{\lambda}{2}|x|^2}$.

Then

$$\begin{aligned} d_{\text{KL}}(p||\pi^y) &= \mathbb{E}^p \left[\log \left(\frac{1}{(2\pi)^{d/2}(\det \Sigma)^{1/2}} e^{-\frac{1}{2}|x-\mu|_{\Sigma}^2} \right) - \log \left(\frac{1}{Z} e^{-L(x) - \frac{\lambda}{2}|x|^2} \right) \right] \\ &= -\frac{d}{2} \log(2\pi) - \frac{1}{2} \log \det \Sigma + \log Z + \mathbb{E}^p \left[-\frac{1}{2}|x-\mu|_{\Sigma}^2 + L(x) + \frac{\lambda}{2}|x|^2 \right]. \end{aligned}$$

$$d_{\text{KL}}(\rho \parallel \pi^y) = -\frac{d}{2} \log(2\pi) - \frac{1}{2} \log \det \Sigma + \log Z + \mathbb{E}^{\rho} \left[-\frac{1}{2} |x - \mu|_{\Sigma^{-1}}^2 + \mathsf{L}(x) + \frac{\lambda}{2} |x|^2 \right].$$

Note that Z is the normalization constant for π and is independent of ρ and hence of μ and Σ . We can represent a given random variable $x \sim \rho$ by writing $x = \mu + \Sigma^{1/2}\xi$, where $\xi \sim \mathcal{N}(0, I)$, and hence

$$|x - \mu|_{\Sigma^{-1}}^2 = |\Sigma^{1/2}\xi|_{\Sigma^{-1}}^2 = |\xi|^2 \quad \Rightarrow \quad \mathbb{E}^{\rho} \left[-\frac{1}{2} |x - \mu|_{\Sigma^{-1}}^2 \right] = -\frac{d}{2}.$$

Moreover,

$$\begin{aligned} \mathbb{E}^{\rho}[|x|^2] &= \int_{\mathbb{R}^d} |x - \mu + \mu|^2 \rho(x) \, dx \\ &= \int_{\mathbb{R}^d} |x - \mu|^2 \rho(x) \, dx + 2\langle \mu, \int_{\mathbb{R}^d} x \rho(x) \, dx \rangle - 2\langle \mu, \int_{\mathbb{R}^d} \mu \rho(x) \, dx \rangle + \int_{\mathbb{R}^d} |\mu|^2 \rho(x) \, dx \\ &= \text{tr}(\Sigma) + 2\langle \mu, \mu \rangle - 2\langle \mu, \mu \rangle + |\mu|^2 = \text{tr}(\Sigma) + |\mu|^2. \end{aligned}$$

We obtain

$$d_{\text{KL}}(p||\pi^y) = -\frac{d}{2} - \frac{d}{2} \log(2\pi) - \frac{1}{2} \log \det \Sigma + \mathbb{E}^p \mathbf{L}(x) + \frac{\lambda}{2} |\mu|^2 + \frac{\lambda}{2} \text{tr}(\Sigma) + \log Z.$$

Define $\mathcal{I}(\mu, \Sigma) = \mathbb{E}^p \mathbf{L}(x) + \frac{\lambda}{2} |\mu|^2 + \frac{\lambda}{2} \text{tr}(\Sigma) - \frac{1}{2} \log \det \Sigma$. Note that there is a correspondence between minimizing $d_{\text{KL}}(p||\pi^y)$ over $p \in \mathcal{A}$ and minimizing $\mathcal{I}(\mu, \Sigma)$ over $\mu \in \mathbb{R}^d$ and positive definite Σ . Moreover:

- $\mathcal{I}(0, I) < \infty$.
- For any Σ , $\mathcal{I}(\mu, \Sigma) \rightarrow \infty$ as $|\mu| \rightarrow \infty$.
- For any μ , $\mathcal{I}(\mu, \Sigma) \rightarrow \infty$ as $\text{tr}(\Sigma) \rightarrow 0$ or $\text{tr}(\Sigma) \rightarrow \infty$.

Therefore, there are $M, r, R > 0$ such that the infimum of $\mathcal{I}(\mu, \Sigma)$ over $\mu \in \mathbb{R}^d$ and positive definite Σ is equal to the infimum of $\mathcal{I}(\mu, \Sigma)$ over

$$\tilde{\mathcal{A}} := \{(\mu, \Sigma) : \mu \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d} \text{ positive-definite symmetric, } |\mu| \leq M, r \leq \text{tr}(\Sigma) \leq R\}.$$

Since \mathcal{I} is continuous in $\tilde{\mathcal{A}}$ it achieves its infimum and the proof is complete. □

We remark that the theorem establishes the existence of a best Gaussian approximation. However, minimizers need not be unique.

Best Gaussian fit by minimizing $d_{\text{KL}}(\pi \| p)$ (“mean-seeking”)

The best Gaussian approximation in Kullback–Leibler with respect to its second argument is unique and given by moment matching.

Theorem (Best Gaussian by moment matching / “mean-seeking”)

Assume that $\bar{\mu} := \mathbb{E}^{\pi}[x]$ is finite and that $\bar{\Sigma} := \mathbb{E}^{\pi}[(x - \bar{\mu})(x - \bar{\mu})^{\text{T}}]$ is positive definite. (Here, π denotes the target distribution, e.g., the posterior.) Then the infimum

$$\inf_{p \in \mathcal{A}} d_{\text{KL}}(\pi \| p)$$

is attained by $p = \mathcal{N}(\bar{\mu}, \bar{\Sigma})$.

Proof. Note that $d_{\text{KL}}(\pi \| p) = -\mathbb{E}^{\pi}[\log p] + \overbrace{\mathbb{E}^{\pi}[\log \pi]}^{\text{independent of } p}$. Since we want a Gaussian minimizer, write $p(x) = ((2\pi)^d |\det \Sigma|)^{-1/2} \exp(-\frac{1}{2}|x - \mu|_{\Sigma^{-1}}^2)$
 $\Rightarrow -\mathbb{E}^{\pi}[\log p] = -\mathbb{E}^{\pi}[\log((2\pi)^{-d/2}(\det \Sigma)^{-1/2} e^{-\frac{1}{2}|x - \mu|_{\Sigma^{-1}}^2})]$
 $= \frac{1}{2}\mathbb{E}^{\pi}[|x - \mu|_{\Sigma^{-1}}^2] + \frac{1}{2} \log \det \Sigma + \frac{d}{2} \log(2\pi).$

Note that the final term is irrelevant for the optimization problem.

Let $\Lambda := \Sigma^{-1}$. Our task is equivalent to finding the minimizer of

$$I(\mu, \Lambda) := \frac{1}{2} \mathbb{E}^\pi [(x - \mu) \Lambda (x - \mu)^\top] - \frac{1}{2} \log \det \Lambda.$$

Let $\Lambda = (\Lambda_{ij})_{i,j=1}^d$. We can view the above functional as the $d + d^2$ variate function $I(\mu_1, \dots, \mu_d, \Lambda_{11}, \Lambda_{12}, \dots, \Lambda_{dd})$. Thus, we only need to show that

$$\nabla I(\bar{\mu}, \bar{\Sigma}^{-1}) = 0 \quad \text{and} \quad \nabla^2 I(\mu, \Sigma^{-1}) > 0 \quad \text{for all } \mu, \Sigma.$$

(($\bar{\mu}, \bar{\Sigma}^{-1}$) is the critical point and the objective function is convex.)

By defining the notations $\partial_\mu f := \left(\frac{\partial f}{\partial \mu_i}\right)_{i=1}^d$ (gradient w.r.t. vector μ) and $\partial_\Lambda f := \left(\frac{\partial f}{\partial \Lambda_{ji}}\right)_{i,j=1}^d$ (gradient w.r.t. vector $(\Lambda_{11}, \Lambda_{12}, \dots, \Lambda_{dd})$, reshaped into a $d \times d$ matrix), we easily see that $\nabla I = 0$ can be expressed as the pair

$$\begin{cases} 0 = \partial_\mu I &= -\mathbb{E}^\pi [\Lambda (x - \mu)] = 0 \\ 0 = \partial_\Lambda I &= \frac{1}{2} \partial_\Lambda (\mathbb{E}^\pi [(x - \mu) \Lambda (x - \mu)^\top]) - \frac{1}{2 \det \Lambda} \partial_\Lambda \det \Lambda \\ &= \frac{1}{2} \mathbb{E}^\pi [(x - \mu)(x - \mu)^\top] - \frac{1}{2} \Lambda^{-1}, \end{cases}$$

where we used a special case of **Jacobi's formula** $\partial_\Lambda \det \Lambda = \det \Lambda \cdot \Lambda^{-1}$.

Clearly, $(x, \Lambda) = (\bar{\mu}, \bar{\Sigma}^{-1})$ is the critical point satisfying the above condition.

Finally, we need to show that $\nabla^2 l(\mu, \Sigma^{-1})$ is positive definite. To this end, we note that

$$\begin{aligned} p(x) &= \sqrt{\frac{\det \Lambda}{(2\pi)^d}} e^{-\frac{1}{2}(x-\mu)^T \Lambda (x-\mu)} = \sqrt{\frac{\det \Lambda}{(2\pi)^d}} e^{-\frac{1}{2}x^T \Lambda x + \mu^T \Lambda x - \frac{1}{2}\mu^T \Lambda \mu} \\ &= \sqrt{\frac{\det \Lambda}{(2\pi)^d}} e^{-\frac{1}{2}\mu^T \Lambda \mu} e^{-\frac{1}{2}x^T \Lambda x + \mu^T \Lambda x} = \frac{e^{-\frac{1}{2}x^T \Lambda x + \mu^T \Lambda x}}{\int_{\mathbb{R}^d} e^{-\frac{1}{2}x^T \Lambda x + \mu^T \Lambda x} dx}. \end{aligned}$$

Noting that $x^T \Lambda x = \sum_{i,j=1}^d \Lambda_{ij} x_i x_j = \sum_{i,j=1}^d \Lambda_{ij} (xx^T)_{ij}$, we can write $x^T \Lambda x = \text{vec}(\Lambda) \cdot \text{vec}(xx^T)$, where we define

$$\text{vec}(M) := (M_{11}, M_{12}, \dots, M_{dd})^T \quad \text{for } M \in \mathbb{R}^{d \times d}.$$

In particular,

$$-\frac{1}{2}x^T \Lambda x + \mu^T \Lambda x = \underbrace{\begin{bmatrix} \Lambda \mu \\ -\frac{1}{2} \text{vec}(\Lambda) \end{bmatrix}}_{=: \theta}^T \underbrace{\begin{bmatrix} x \\ \text{vec}(xx^T) \end{bmatrix}}_{=: T(x)}$$

and we can write $p_\theta(x) := p(x) = \frac{1}{Z(\theta)} e^{\theta^T T(x)}$, $Z(\theta) := \int_{\mathbb{R}^d} e^{\theta^T T(x)} dx$.

The importance of the characterization

$$p_{\theta}(x) = \frac{1}{Z(\theta)} e^{\theta^T T(x)}, \quad Z(\theta) := \int_{\mathbb{R}^d} e^{\theta^T T(x)} dx,$$

lies in the fact that *every possible Gaussian PDF* can be parameterized by the vector $\theta = (\theta_1, \dots, \theta_{d+d^2})^T$. Thus, the KL divergence $d_{\text{KL}}(\pi \| p_{\theta})$ that we are interested in can be recast as

$$\begin{aligned} H(\theta) &:= d_{\text{KL}}(\pi \| p_{\theta}) = -\mathbb{E}^{\pi}[\log p_{\theta}] + \mathbb{E}^{\pi}[\log \pi] \\ &= -\theta^T \mathbb{E}^{\pi}[T(x)] + \log Z(\theta) + \mathbb{E}^{\pi}[\log \pi]. \end{aligned}$$

Noting that $\nabla_{\theta}^2(\theta^T \mathbb{E}^{\pi}[T(x)]) = 0$ and $\frac{\partial \log Z(\theta)}{\partial \theta_i} = \frac{1}{Z(\theta)} \int_{\mathbb{R}^d} \frac{\partial}{\partial \theta_i} e^{\theta^T T(x)} dx = \frac{1}{Z(\theta)} \int_{\mathbb{R}^d} T_i(x) e^{\theta^T T(x)} dx$, we compute

$$\begin{aligned} [\nabla_{\theta}^2 H(\theta)]_{ij} &= \frac{\partial^2 \log Z(\theta)}{\partial \theta_i \partial \theta_j} = \frac{\partial}{\partial \theta_j} \left(\frac{1}{Z(\theta)} \int_{\mathbb{R}^d} T_i(x) e^{\theta^T T(x)} dx \right) \\ &= -\frac{1}{Z(\theta)^2} \left(\int_{\mathbb{R}^d} T_i(x) e^{\theta^T T(x)} dx \right) \left(\int_{\mathbb{R}^d} T_j(x) e^{\theta^T T(x)} dx \right) + \frac{1}{Z(\theta)} \int_{\mathbb{R}^d} T_i(x) T_j(x) e^{\theta^T T(x)} dx \\ &= \mathbb{E}^{p_{\theta}}[T_i T_j] - \mathbb{E}^{p_{\theta}}[T_i] \mathbb{E}^{p_{\theta}}[T_j] = [\text{Cov}^{p_{\theta}}(T)]_{ij}, \end{aligned}$$

which is positive definite. □

Remark. Notice that the preceding proof of convexity holds for any distribution p that can be parameterized by the following more general expression:

$$p_{\theta}(x) = h(x)\exp\left(\theta^T T(x) - A(\theta)\right) \quad (1)$$

$$\text{with } A(\theta) = \log \left[\int_{\mathbb{R}^d} h(x)\exp\left(\theta^T T(x)\right) dx \right].$$

Since $h(x)$ is independent of θ , the conclusion of the previous theorem carries over to distributions with the form of (1). Such distributions belong to the *exponential family* in the statistics literature. Here, θ is called the natural parameter, $T(x)$ the sufficient statistic, $h(x)$ the base measure, and $A(\theta)$ the log-partition.

The Gaussian distribution is a special case in which $h(x)$ is constant with respect to x .

Variational formulation of Bayes' theorem

We have been concerned with finding the best Gaussian approximations to a measure with respect to KL divergences. Bayes' theorem itself can be formulated through a closely related minimization principle. Consider a posterior $\pi^y(x)$ in the following form:

$$\pi^y(x) = \frac{1}{Z} \exp(-L(x)) \pi(x),$$

where $\pi(x)$ is the prior, $L(x)$ is the negative log-likelihood, and Z the normalization constant. We assume here for exposition that all densities are positive. Let p be an *arbitrary* PDF. Then we can express $d_{\text{KL}}(p \parallel \pi^y)$ as

$$\begin{aligned} d_{\text{KL}}(p \parallel \pi^y) &= \int_{\mathbb{R}^d} \log\left(\frac{p}{\pi^y}\right) p \, dx = \int_{\mathbb{R}^d} \log\left(\frac{p}{\pi} \frac{\pi}{\pi^y}\right) p \, dx \\ &= \int_{\mathbb{R}^d} \log\left(\frac{p}{\pi} \exp(L(x)) Z\right) p \, dx \\ &= d_{\text{KL}}(p \parallel \pi) + \mathbb{E}^p[L(x)] + \log Z. \end{aligned}$$

If we define

$$\mathcal{J}(p) = d_{\text{KL}}(p \parallel \pi) + \mathbb{E}^p[\mathcal{L}(x)]$$

then we have the following:

Theorem (Bayes' theorem as an optimization principle)

The posterior distribution π^y is given by the following minimization principle:

$$\pi^y = \arg \min_{p \in \mathcal{P}} \mathcal{J}(p),$$

where \mathcal{P} contains all probability densities on \mathbb{R}^d .

Proof.

Since Z is the normalization constant for π^y and is independent of p , the minimizer of $d_{\text{KL}}(p \parallel \pi^y)$ will also be the minimizer of $\mathcal{J}(p)$. Since the global minimizer of $d_{\text{KL}}(p \parallel \pi^y)$ is attained at $p = \pi^y$ the result follows. \square

Why is it useful to view the posterior as the minimizer of an energy?

- The variational formulation provides a natural way to approximate the posterior by restricting the minimization problem to distributions satisfying some computationally desirable property.
 - For instance, variational Bayes methods often restrict the minimization to densities with product structure and in this chapter we have studied restriction to the class of Gaussian distributions.
- Variational formulations allow to show convergence of posterior distributions indexed by some parameter of interest by studying the Γ -convergence of the associated objective functionals.
- Variational formulations provide natural paths, defined by a gradient flow, towards the posterior. Understanding these flows and their rates of convergence is helpful in the choice of sampling algorithms.