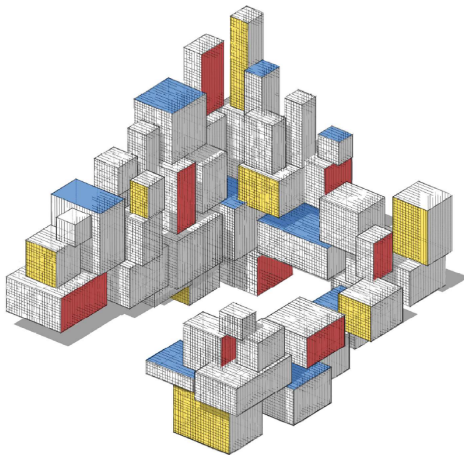


# Bayesian Learning



# Purpose

In this lecture we discuss:

- Bayesian unsupervised learning
- Prior, likelihood, posterior pdfs
- Bayesian normal model
- MAP estimates and posterior mean

# Bayesian Unsupervised Learning

In Bayesian unsupervised learning, we seek to approximate the unknown joint density  $f(\mathbf{x}_1, \dots, \mathbf{x}_n)$  of the training data  $\mathcal{T}_n = \{\mathbf{X}_1, \dots, \mathbf{X}_n\}$  via a joint pdf of the form

$$\int \left( \prod_{i=1}^n g(\mathbf{x}_i | \boldsymbol{\theta}) \right) w(\boldsymbol{\theta}) \, \mathrm{d}\boldsymbol{\theta}, \quad (1)$$

where

- $g(\cdot | \boldsymbol{\theta}) \in \mathcal{G}_p := \{g(\cdot | \boldsymbol{\theta}), \boldsymbol{\theta} \in \boldsymbol{\Theta} \subset \mathbb{R}^p\},$
- $w(\boldsymbol{\theta}) \in \mathcal{W}_p.$

Note that this includes the case of exchangeable training data.

# Notation

Following standard practice in a Bayesian context, instead of writing  $f_X(x)$  and  $f_{X|Y}(x|y)$  for the pdf of  $X$  and the conditional pdf of  $X$  given  $Y$ , one simply writes  $f(x)$  and  $f(x|y)$ . If  $Y$  is a different random variable, its pdf (at  $y$ ) is thus denoted by  $f(y)$ .

Thus, we will use the same symbol

- $g$  for (conditional) **approximating** probability densities, and
- $f$  for (conditional) **true** and unknown probability densities.

Using Bayesian notation, we can write  $g(\tau | \theta) = \prod_{i=1}^n g(\mathbf{x}_i | \theta)$  and thus the approximating joint pdf (1) can then be written as

$\int g(\tau | \theta) w(\theta) d\theta$  and the true unknown joint pdf as  
 $f(\tau) = f(\mathbf{x}_1, \dots, \mathbf{x}_n)$ .

# Kullback–Leibler Risk

We can use the Kullback–Leibler risk to measure the discrepancy between the proposed approximation (1) and the true  $f(\tau)$ :

$$\ell(g) = \mathbb{E} \ln \frac{f(\mathcal{T})}{\int g(\mathcal{T} | \boldsymbol{\theta}) w(\boldsymbol{\theta}) d\boldsymbol{\theta}} = \int f(\tau) \ln \frac{f(\tau)}{\int g(\tau | \boldsymbol{\theta}) w(\boldsymbol{\theta}) d\boldsymbol{\theta}} d\tau.$$

Because the training data is not necessarily iid, the expectation must be with respect to the *joint* density of  $\mathcal{T}$ , not with respect to the *marginal*  $f(\mathbf{x})$  (as in the iid case).

# Training Loss

As the true KL risk is not known, we minimize the corresponding training loss, which is equivalent to maximizing the likelihood of the training data  $\tau$ ; that is, solving the optimization problem

$$\max_{w \in \mathcal{W}_p} \int g(\tau | \boldsymbol{\theta}) w(\boldsymbol{\theta}) d\boldsymbol{\theta},$$

where the maximization is over the class of densities  $\mathcal{W}_p$ .

Suppose that we have a rough guess, denoted  $w_0(\boldsymbol{\theta})$ , for the best  $w \in \mathcal{W}_p$  that minimizes the Kullback–Leibler risk.

We can always increase the likelihood  $L_0 := \int g(\tau | \boldsymbol{\theta}) w_0(\boldsymbol{\theta}) d\boldsymbol{\theta}$  by instead using the density  $w_1(\boldsymbol{\theta}) := w_0(\boldsymbol{\theta}) g(\tau | \boldsymbol{\theta}) / L_0$ , giving a likelihood  $L_1 := \int g(\tau | \boldsymbol{\theta}) w_1(\boldsymbol{\theta}) d\boldsymbol{\theta}$ .

## Increasing the Likelihood

To see this, write  $L_0$  and  $L_1$  as expectations with respect to  $w_0$ . In particular, we can write

$$L_0 = \mathbb{E}_{w_0} g(\tau | \theta) \quad \text{and} \quad L_1 = \mathbb{E}_{w_1} g(\tau | \theta) = \mathbb{E}_{w_0} g^2(\tau | \theta) / L_0.$$

It follows that

$$L_1 - L_0 = \frac{1}{L_0} \mathbb{E}_{w_0} [g^2(\tau | \theta) - L_0^2] = \frac{1}{L_0} \text{Var}_{w_0} [g(\tau | \theta)] \geq 0. \quad (2)$$

We may thus expect to obtain better predictions using  $w_1$  instead of  $w_0$ , because  $w_1$  has taken into account the observed data  $\tau$  and increased the likelihood of the model.

If we iterate this process, we obtain a sequence of densities  $w_1, w_2, \dots$  that degenerates to the point mass at the **maximum likelihood** estimator  $\hat{\theta}$ .

# Bayesian Learner

In many situations, the maximum likelihood estimate  $g(\tau | \hat{\theta})$  is either not an appropriate approximation to  $f(\tau)$  or simply fails to exist.

In such cases, given an initial non-degenerate guess  $w_0(\theta) = g(\theta)$ , one can obtain a more appropriate and non-degenerate approximation to  $f(\tau)$  by taking  $w(\theta) = w_1(\theta) = g(\tau | \theta) g(\theta) / g(\tau) = g(\theta | \tau)$  in (1).

This gives the following Bayesian learner of  $f(\mathbf{x})$ :

$$g_{\tau}(\mathbf{x}) := \int g(\mathbf{x} | \theta) \frac{g(\tau | \theta) g(\theta)}{g(\tau)} d\theta. \quad (3)$$



# Prior, Likelihood, Posterior

With this notation, we have the following definitions.

## Definition: Prior, Likelihood, and Posterior

Let  $\tau$  and  $\mathcal{G}_p := \{g(\cdot | \theta), \theta \in \Theta\}$  be the training set and family of approximating functions.

- A pdf  $g(\theta)$  that reflects our *a priori* beliefs about  $\theta$  is called the **prior** pdf.
- The conditional pdf  $g(\tau | \theta)$  is called the **likelihood**.
- Inference about  $\theta$  is given by the **posterior** pdf  $g(\theta | \tau)$ , which is proportional to the product of the prior and the likelihood:

$$g(\theta | \tau) \propto g(\tau | \theta) g(\theta).$$

# Prior and Posterior Predictive Density

The initial guess  $g(\boldsymbol{\theta})$  conveys the *a priori* (prior to training the Bayesian learner) information about the optimal density in  $\mathcal{W}_p$  that minimizes the KL risk. Using this prior  $g(\boldsymbol{\theta})$ , the Bayesian approximation to  $f(\mathbf{x})$  is the **prior predictive density**:

$$g(\mathbf{x}) = \int g(\mathbf{x} | \boldsymbol{\theta}) g(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

The posterior pdf conveys improved knowledge about this optimal density in  $\mathcal{W}_p$  after training with  $\tau$ . Using the posterior  $g(\boldsymbol{\theta} | \tau)$ , the Bayesian learner of  $f(\mathbf{x})$  is the **posterior predictive density**:

$$g_\tau(\mathbf{x}) = g(\mathbf{x} | \tau) = \int g(\mathbf{x} | \boldsymbol{\theta}) g(\boldsymbol{\theta} | \tau) d\boldsymbol{\theta},$$

where we have assumed that  $g(\mathbf{x} | \boldsymbol{\theta}, \tau) = g(\mathbf{x} | \boldsymbol{\theta})$ ; that is, the likelihood depends on  $\tau$  only through the parameter  $\boldsymbol{\theta}$ .

# Choice of the Prior

The choice of the prior is typically governed by two considerations:

1. The prior should be simple enough to facilitate the computation or simulation of the posterior pdf.
2. The prior should be general enough to model ignorance of the parameter of interest.

Priors that do not convey much knowledge of the parameter are said to be **uninformative**. The uniform or **flat** prior is frequently used.

For the purpose of analytical and numerical computations, we can view  $\theta$  as a random vector with prior density  $g(\theta)$ , which after training is updated to the posterior density  $g(\theta | \tau)$ .

The above thinking allows us to write

$g(\mathbf{x} | \tau) \propto \int g(\mathbf{x} | \theta) g(\tau | \theta) g(\theta) d\theta$ , for example, thus ignoring any constants that do not depend on the argument of the densities.

## Example: Normal Model

Suppose that the training data  $\mathcal{T} = \{X_1, \dots, X_n\}$  is modeled using the likelihood  $g(x | \theta)$  that is the pdf of

$$X | \theta \sim \mathcal{N}(\mu, \sigma^2),$$

where  $\theta := [\mu, \sigma^2]^\top$ . For the prior distribution of  $\theta$  we can specify prior distributions for  $\mu$  and  $\sigma^2$  separately and then assume independence. A possible prior distribution for  $\mu$  is

$$\mu \sim \mathcal{N}(\nu, \phi^2).$$

A common prior distribution for  $1/\sigma^2$  is

$$\frac{1}{\sigma^2} \sim \text{Gamma}(\alpha, \beta).$$

The smaller  $\alpha$  and  $\beta$  are, the less informative is the prior.

# Posterior Pdf

Under this prior,  $\sigma^2$  is said to have an **inverse gamma** distribution, with pdf proportional to  $\exp(-\beta/z) / z^{\alpha+1}$ .

The Bayesian posterior is then:

$$\begin{aligned} g(\mu, \sigma^2 | \tau) &\propto g(\mu) \times g(\sigma^2) \times g(\tau | \mu, \sigma^2) \\ &\propto \exp \left\{ -\frac{(\mu - \nu)^2}{2\phi^2} \right\} \times \frac{\exp \{-\beta/\sigma^2\}}{(\sigma^2)^{\alpha+1}} \times \frac{\exp \{-\sum_i (x_i - \mu)^2 / (2\sigma^2)\}}{(\sigma^2)^{n/2}} \\ &\propto (\sigma^2)^{-n/2-\alpha-1} \exp \left\{ -\frac{(\mu - \nu)^2}{2\phi^2} - \frac{\beta}{\sigma^2} - \frac{(\mu - \bar{x}_n)^2 + S_n^2}{2\sigma^2/n} \right\}, \end{aligned}$$

where  $S_n^2 := \frac{1}{n} \sum_i x_i^2 - \bar{x}_n^2 = \frac{1}{n} \sum_i (x_i - \bar{x}_n)^2$  is the (scaled) sample variance.

## Posterior Mean

The conditional pdf of  $\mu$  given  $\sigma^2$  and  $\tau$  is

$$g(\mu | \sigma^2, \tau) \propto \exp \left\{ -\frac{(\mu - \nu)^2}{2\phi^2} - \frac{(\mu - \bar{x}_n)^2}{2\sigma^2/n} \right\},$$

which can be recognized as the pdf of

$$(\mu | \sigma^2, \tau) \sim \mathcal{N} \left( \gamma_n \bar{x}_n + (1 - \gamma_n) \nu, \gamma_n \sigma^2 / n \right), \quad (4)$$

where  $\gamma_n := \frac{n}{\sigma^2} \bigg/ \left( \frac{1}{\phi^2} + \frac{n}{\sigma^2} \right)$ .

The posterior mean  $\mathbb{E}[\mu | \sigma^2, \tau] = \gamma_n \bar{x}_n + (1 - \gamma_n) \nu$  is a weighted linear combination of the prior mean  $\nu$  and the sample average  $\bar{x}_n$ .

As  $n \rightarrow \infty$ , the weight  $\gamma_n \rightarrow 1$  and thus the posterior mean approaches the maximum likelihood estimate  $\bar{x}_n$ .

## Normal Model (cont.)

An example of an **improper prior** is obtained when we take prior  $g(\mu) \propto 1$ .

Nevertheless, the posterior is a proper density, and in particular the conditional posterior of  $(\mu \mid \sigma^2, \tau)$  simplifies to

$$(\mu \mid \sigma^2, \tau) \sim \mathcal{N}(\bar{x}_n, \sigma^2/n).$$

In addition,

$$g(\sigma^2 \mid \tau) = \int g(\mu, \sigma^2 \mid \tau) d\mu \propto (\sigma^2)^{-(n-1)/2-\alpha-1} \exp\left\{-\frac{\beta + nS_n^2/2}{\sigma^2}\right\},$$

which we recognize as the density corresponding to

$$\frac{1}{\sigma^2} \mid \tau \sim \text{Gamma}\left(\alpha + \frac{n-1}{2}, \beta + \frac{n}{2}S_n^2\right).$$

## Normal Model (cont.)

In addition to  $g(\mu) \propto 1$ , we can also use an improper prior  $g(\sigma^2) \propto 1/\sigma^2$  for  $\sigma^2$ .

In this case, the posterior marginal density for  $\sigma^2$  implies that:

$$\frac{nS_n^2}{\sigma^2} \Big| \tau \sim \chi_{n-1}^2$$

and the posterior marginal density for  $\mu$  implies that:

$$\frac{\mu - \bar{x}_n}{S_n/\sqrt{n-1}} \Big| \tau \sim t_{n-1}. \quad (5)$$

In general, deriving a simple formula for the posterior density of  $\theta$  is either impossible or too tedious. Instead, **Monte Carlo methods** are used to simulate (approximately) from the posterior for the purposes of inference and prediction.



# Credible Interval

One way in which a distributional result such as (5) can be useful is in the construction of a 95% **credible interval**  $\mathcal{I}$  for the parameter  $\mu$ ; that is, an interval  $\mathcal{I}$  such that the probability  $\mathbb{P}[\mu \in \mathcal{I} \mid \tau]$  is equal to 0.95.

For example, the symmetric 95% credible interval is

$$\mathcal{I} = \left[ \bar{x}_n - \frac{S_n}{\sqrt{n-1}}\gamma, \bar{x}_n + \frac{S_n}{\sqrt{n-1}}\gamma \right],$$

with  $\gamma$  the 0.975-quantile of the  $t_{n-1}$  distribution (conditional on  $\tau$ ).

Note that the parameter  $\mu$  is interpreted as a random variable with a distribution.

This is unlike the case of classical confidence intervals, where the parameter is nonrandom, but the interval is a random object.

## Example

Consider analysing the number of deaths during birth in a maternity ward. Suppose that the data is given by  $\tau = \{x_1, \dots, x_n\}$ , where  $x_i = 1$  if the  $i$ -th baby has died during birth and  $x_i = 0$  otherwise, for  $i = 1, \dots, n$ .

A possible Bayesian model for the data is:

$$\begin{cases} \theta \sim \mathcal{U}(0, 1) & \text{(uniform prior),} \\ (X_1, \dots, X_n \mid \theta) \stackrel{\text{iid}}{\sim} \text{Ber}(\theta). \end{cases}$$

The likelihood is therefore

$$g(\tau \mid \theta) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} = \theta^s (1 - \theta)^{n-s},$$

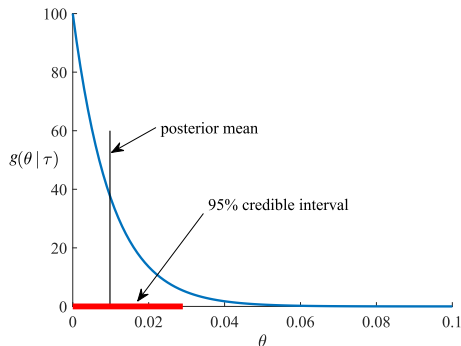
where  $s = x_1 + \dots + x_n$  is the total number of deaths.

# Posterior Pdf

Since  $g(\theta) = 1$ , the posterior pdf is

$$g(\theta | \tau) \propto \theta^s (1 - \theta)^{n-s}, \quad \theta \in [0, 1],$$

which is the pdf of the  $\text{Beta}(s + 1, n - s + 1)$  distribution.



**Figure:** Posterior pdf for  $\theta$ , with  $n = 100$  and  $s = 0$ .

# MAP Estimate and Posterior Mean

The **maximum a posteriori** (MAP) estimate of  $\theta$  (the mode or maximizer of the posterior density) is

$$\operatorname{argmax}_{\theta} g(\theta \mid \tau) = \frac{s}{n},$$

which agrees with the maximum likelihood estimate.

When  $(s, n) = (0, 100)$ , the maximum likelihood estimate  $\hat{\theta} = 0$  infers (wrongly) that deaths at birth are not possible (same for the MAP).

In contrast the **posterior mean**  $\mathbb{E}[\theta \mid \tau] = (s + 1)/(n + 2)$  is 0.0098 for this case and provides the more reasonable point estimate.