

# Bayesian Regression

Assaf Dvora

June 2019

## 1 The Linear Regression Problem

Assume we have a training set  $\mathcal{D}$  of  $n$  observations,  $\mathcal{D} = \{(\mathbf{x}_i, y_i) | i = 1, \dots, n\}$ , where  $\mathbf{x} \in \mathbb{R}^d$  denotes an input vector (covariates) of dimension  $\mathcal{D}$  and  $y$  denotes a scalar output or target (dependent variable). The columns vector inputs for all  $n$  cases are aggregated in the  $D \times n$  design matrix  $X \in \mathbb{R}^{D \times n}$  and the targets are collected in the vector  $\mathbf{y}$ . We can write  $\mathcal{D} = (X, \mathbf{y})$ . We are interested in making inferences about the relationship between inputs and targets, i.e. the conditional distribution of the targets given the inputs (but we are not interested in modeling the input distribution itself).

## 2 Bayesian Linear Regression

Linear regression assumes that there is a linear relationship between inputs and outputs, up to some noise term  $\epsilon$ :

$$y = \mathbf{x}^T \mathbf{w} + \epsilon, \quad (1)$$

where  $\mathbf{w}$  is a vector of weights (parameters) of the linear model. Often bias term is added, but as this can be implemented by augmenting the input vector  $\mathbf{x}$  with additional element whose value is always one, we do not explicitly include it in our notation. We further assume that this noise term follows an IID Gaussian distribution with zero mean and variance  $\sigma_\epsilon^2$

$$\epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2). \quad (2)$$

The likelihood results in

$$\begin{aligned} p(\mathbf{y} | X, \mathbf{w}) &= \prod_{i=1}^n p(y_i | \mathbf{x}_i, \mathbf{w}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_\epsilon^2}} \exp \left[ -\frac{(y_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma_\epsilon^2} \right] = \\ &= \frac{1}{(2\pi\sigma_\epsilon^2)^{n/2}} \exp \left[ -\frac{1}{2\sigma_\epsilon^2} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \mathbf{w})^2 \right] = \mathcal{N}(X^T \mathbf{w}, \sigma_\epsilon^2 I). \end{aligned} \quad (3)$$

In the Bayesian formalism we need to specify a *prior* over the parameters. We put a zero mean Gaussian prior with covariance matrix  $\Sigma_p$  on the weights

$$\mathbf{w} \sim \mathcal{N}(0, \Sigma_p). \quad (4)$$

The posterior distribution over the weights is given by,

$$p(\mathbf{w}|\mathbf{y}, X) = \frac{p(\mathbf{y}|X, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|X)} = \frac{p(\mathbf{y}|X, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|X, \mathbf{w})p(\mathbf{w})d\mathbf{w}}. \quad (5)$$

Since the prior and the likelihood are conjugate Gaussian distributions, the posterior is Gaussian:

$$p(\mathbf{w}|\mathbf{y}, X) = \mathcal{N}(\boldsymbol{\mu}_w, \Sigma_w), \quad (6)$$

and “completing the square” we obtain,

$$\boldsymbol{\mu}_w = \frac{1}{\sigma_\epsilon^2} A^{-1} X \mathbf{y}, \quad \Sigma_w = A^{-1}, \quad (7)$$

where  $A = \sigma_\epsilon^{-2} X X^T + \Sigma_p^{-1}$ . We note that for Gaussian posterior the mean vector is also the MAP estimate of  $\mathbf{w}$ .

To make predictions for a test case we average over all possible parameter values, weighted by their posterior probability. Thus, the predictive probability is given by,

$$p(y_*|\mathbf{x}_*, X, \mathbf{y}) = \int p(y_*|\mathbf{x}_*, \mathbf{w})p(\mathbf{w}|\mathbf{y}, X)d\mathbf{w} \quad (8)$$

in the case of linear model with Gaussian posterior, the predictive probability is also Gaussian,

$$p(y_*|\mathbf{x}_*, X, \mathbf{y}) = \mathcal{N}\left(\frac{1}{\sigma_\epsilon^2} \mathbf{x}_*^T A^{-1} X \mathbf{y}, \mathbf{x}_*^T A^{-1} \mathbf{x}_*\right). \quad (9)$$

### 3 Bayesian Linear Regression Using Gaussian Process.

Lets assume the following additive model for the response variable

$$y = f(\mathbf{x}) + \varepsilon, \quad (10)$$

where  $f(\mathbf{x})$  is unobserved latent **stochastic process** indexed by the input vector  $\mathbf{x}$ , and  $\varepsilon$  is i.i.d Gaussian noise with variance  $\sigma_\epsilon^2$ . Hence, the response variable satisfies

$$p(y|f(\mathbf{x}), \sigma_\epsilon^2) = \mathcal{N}(f(\mathbf{x}), \sigma_\epsilon^2). \quad (11)$$

We further assume that  $\sigma_\epsilon^2$  is known and hence we use  $p(y|f(\mathbf{x}))$  instead.

Denote with  $\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)]^T \in R^n$  the vector of realizations of the random process. The likelihood of the model is given by

$$p(\mathcal{D}|\mathbf{f}) = p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^n p(y_i|f(\mathbf{x}_i), \sigma_\epsilon^2) = \mathcal{N}(\mathbf{f}, \sigma_\epsilon^2 I). \quad (12)$$

To define a prior over  $\mathbf{f}$  we use the Bayesian Linear Regression model,  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ , with prior  $\mathbf{w} \sim \mathcal{N}(0, \Sigma_w)$ . We show in Appendix A that this results in a Gaussian process for  $f(\mathbf{x})$  with the mean and covariance satisfying

$$\begin{aligned} E[f(\mathbf{x})] &= 0, \\ E[f(\mathbf{x})f(\mathbf{x}')]\triangleq \mathcal{K}(\mathbf{x}, \mathbf{x}') &= \mathbf{x}^T \Sigma_w \mathbf{x}'. \end{aligned} \quad (13)$$

Thus, the prior probability of the latent vector  $\mathbf{f}$  is given by

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, K), \quad (14)$$

where  $K = X^T \Sigma_w X$ . Based on Bayes theorem, the posterior probability can be written as

$$p(\mathbf{f}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{f})p(\mathbf{f})}{p(\mathcal{D})} \quad (15)$$

since the prior probability  $p(\mathbf{f})$  and the likelihood function  $p(\mathcal{D}|\mathbf{f})$  are conjugate Gaussian distributions, the posterior is Gaussian and satisfies

$$p(\mathbf{f}|\mathcal{D}) = \mathcal{N}(\boldsymbol{\mu}_p, \Sigma_p). \quad (16)$$

We show in Appendix B that the mean vector,  $\boldsymbol{\mu}_p$ , and the covariance matrix,  $\Sigma_p$  are given by

$$\boldsymbol{\mu}_p = K(\sigma_\varepsilon^2 I + K)^{-1} \mathbf{y} \quad (17)$$

$$\Sigma_p = \sigma_\varepsilon^2(\sigma_\varepsilon^2 I + K)^{-1} K \quad (18)$$

For a test case,  $\mathbf{x}_*$ , we would like to predict the latent variable  $f_* \triangleq f(\mathbf{x}_*)$  given the knowledge we have about our data. The posterior predictive distribution  $p(f_*|\mathbf{x}_*\mathcal{D})$  can be written as

$$p(f_*|\mathcal{D}) = \int p(f_*|\mathbf{f})p(\mathbf{f}|\mathcal{D})d\mathbf{f} \quad (19)$$

where  $p(f_*|\mathbf{f})$  is the (prior) conditional distribution of  $f_*$  given  $\mathbf{f}$  and  $p(\mathbf{f}|\mathcal{D})$  is the posterior distribution of  $\mathbf{f}$ . We note that  $\mathbf{f}$  and  $f_*$  are jointly Gaussian

$$\begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix} \sim \mathcal{N}\left[\begin{pmatrix} \mathbf{0} \\ 0 \end{pmatrix}, \begin{pmatrix} K & \mathbf{k}_* \\ \mathbf{k}_*^T & \mathcal{K}(\mathbf{x}_*, \mathbf{x}_*) \end{pmatrix}\right] \quad (20)$$

where  $\mathbf{k}_* = [\mathcal{K}(\mathbf{x}_1, \mathbf{x}_*), \mathcal{K}(\mathbf{x}_2, \mathbf{x}_*), \dots, \mathcal{K}(\mathbf{x}_n, \mathbf{x}_*)]$ . The conditional distribution  $p(f_*|\mathbf{f})$  is also Gaussian and is given by (Rasmussen et al.)

$$p(f_*|\mathbf{f}) = \mathcal{N}(\mathbf{k}_*^T K^{-1} \mathbf{f}, \mathcal{K}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T K^{-1} \mathbf{k}_*). \quad (21)$$

In Appendix C we show that the posterior predictive distribution,  $p(f_*|\mathcal{D})$  is Gaussian and satisfies

$$\mathbb{E}[f_*|\mathcal{D}] = \mathbf{k}_*^T K^{-1} \boldsymbol{\mu}_p \quad (22)$$

$$\text{var}[f_*|\mathcal{D}] = \mathcal{K}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T (K^{-1} - K^{-1} \Sigma_p K^{-1}) \mathbf{k}_* \quad (23)$$

where  $\boldsymbol{\mu}_p$  and  $\Sigma_p$  are the mean vector and covariance matrix of the posterior  $p(\mathbf{f}|\mathcal{D})$  (respectively). Substituting (17) to (22) we obtain the mean value of the posterior predictive distribution

$$\mathbb{E}[f_*|\mathcal{D}] = \mathbf{k}_*^T (K + \sigma_\varepsilon^2 I)^{-1} \mathbf{y}. \quad (24)$$

To obtain an expression for the posterior predictive variance we first note that the posterior covariance matrix,  $\Sigma_p$  can be rearranged to

$$\begin{aligned} \Sigma_p &= \sigma_\varepsilon^2 (\sigma_\varepsilon^2 I + K)^{-1} K = (I + \sigma_\varepsilon^{-2} K)^{-1} K \\ &= [K^{-1} (I + \sigma_\varepsilon^{-2} K)]^{-1} = (K^{-1} + \sigma_\varepsilon^{-2} I)^{-1}. \end{aligned} \quad (25)$$

Substituting (25) to (23) we obtain,

$$\text{var}[f_*|D] = \mathcal{K}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T (K^{-1} - K^{-1} (K^{-1} + \sigma_\varepsilon^{-2} I)^{-1} K^{-1}) \mathbf{k}_* \quad (26)$$

Using the matrix inversion lemma (Rasmussen et al., Appendix A), we obtain

$$\text{var}[f_*|D] = \mathcal{K}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T (K + \sigma_\varepsilon^2 I)^{-1} \mathbf{k}_*. \quad (27)$$

Thus, the posterior predictive distribution satisfies

$$p(\mathbf{f}|\mathcal{D}) = \mathcal{N}(\mathbf{k}_*^T (K + \sigma_\varepsilon^2 I)^{-1} \mathbf{y}, \mathcal{K}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T (K + \sigma_\varepsilon^2 I)^{-1} \mathbf{k}_*) \quad (28)$$

Next we show that the posterior predictive distribution (28) is equal to the predictive probability (9) obtained for the Bayesian Linear Regression problem (Section 2). Using the linear covariance defined in (13) we obtain

$$K = X^T \Sigma_w X \quad (29)$$

$$\mathbf{k}_* = X^T \Sigma_w \mathbf{x}_* \quad (30)$$

$$\mathcal{K}(\mathbf{x}_*, \mathbf{x}_*) = \mathbf{x}_*^T \Sigma_w \mathbf{x}_* \quad (31)$$

Substituting into the variance expression (27) and using the matrix inverse lemma (see Appendix D) we obtain,

$$\begin{aligned} \text{var}[f_*|D] &= \mathbf{x}_*^T \Sigma_w \mathbf{x}_* - \mathbf{x}_*^T \Sigma_w X (X^T \Sigma_w X + \sigma_\varepsilon^2 I)^{-1} X^T \Sigma_w \mathbf{x}_* = \\ &= \mathbf{x}_*^T (\Sigma_w - X (X^T \Sigma_w X + \sigma_\varepsilon^2 I)^{-1} X^T) \mathbf{x}_* = \\ &= \mathbf{x}_*^T (\Sigma_w^{-1} + \sigma_\varepsilon^{-2} X X^T)^{-1} \mathbf{x}_* \end{aligned}$$

which is equal to the variance of the predictive probability in (9). Similarly, substituting (29)-(31) into the mean (24) we obtain,

$$\begin{aligned} \mathbb{E}[f_*|\mathcal{D}] &= \mathbf{k}_*^T (K + \sigma_\varepsilon^2 I)^{-1} \mathbf{y} = \\ &= \mathbf{x}_*^T \Sigma_w X (X^T \Sigma_w X + \sigma_\varepsilon^2 I)^{-1} \mathbf{y} = \end{aligned}$$

Since  $\Sigma_w$  is symmetric and positive-definite we can use a variant of the matrix inverse lemma (Appendix D) to obtain

$$\mathbb{E}[f_*|\mathcal{D}] = \frac{1}{\sigma_\varepsilon^2} \mathbf{x}_*^T (\sigma_\varepsilon^{-2} X X^T - \Sigma_w^{-1})^{-1} X \mathbf{y}$$

which is equal to the mean of the predictive probability in (9).

## Appendix A

Using the Bayesian Linear Regression model, the process  $f(\mathbf{x})$  is defined as

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}, \quad \mathbf{w} \sim \mathcal{N}(0, \Sigma_w).$$

We would like to show that  $f(\mathbf{x})$  is Gaussian process.

Definition A.1: A Gaussian process is a collection of random variables, any finite number of which have a joint Gaussian distribution.

Lemma A.1: A collection of Gaussian random variables are jointly Gaussian and form a joint Gaussian distribution if any linear combination of the variables is a Gaussian random variable.

Denote  $f(x_i), i = 1, \dots, n$  the set of random variables corresponding to  $n$  input point  $x_i$ . Further denote with  $a_i \in \mathbb{R}, i = 1, \dots, n$  a set of deterministic values. Thus,

$$\sum_{i=1}^n a_i f(\mathbf{x}_i) = \sum_{i=1}^n a_i \mathbf{w}^T \mathbf{x}_i = \sum_{i=1}^n \mathbf{w}^T a_i \mathbf{x}_i = \mathbf{w}^T \underbrace{\sum_{i=1}^n a_i \mathbf{x}_i}_{\mathbf{Y}} = \mathbf{w}^T \mathbf{Y}.$$

Since  $\mathbf{w}$  is Gaussian random vector,  $\mathbf{w}^T \mathbf{Y}$  is Gaussian variable (linear combination of Gaussian variables) and according to Lemma A.1  $f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)$  are jointly Gaussian. Hence, according to Definition A.1  $f(\mathbf{x})$  is Gaussian process with mean and covariance function satisfying

$$\begin{aligned} E[f(\mathbf{x})] &= 0 \\ E[f(\mathbf{x})f(\mathbf{x}')] &= \mathbf{x}^T \Sigma_w \mathbf{x}' \end{aligned}$$

## Appendix B

From Eq. (15) we obtain

$$\begin{aligned} \log p(\mathbf{f}|\mathcal{D}) &= C - \frac{1}{2\sigma_\epsilon^2} (\mathbf{y} - \mathbf{f})^T (\mathbf{y} - \mathbf{f}) - \frac{1}{2} \mathbf{f}^T K^{-1} \mathbf{f} = \\ &= C - \frac{1}{2\sigma_\epsilon^2} (\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{f} + \mathbf{f}^T \mathbf{f}) - \frac{1}{2} \mathbf{f}^T K^{-1} \mathbf{f} = \\ &= C + \frac{1}{\sigma_\epsilon^2} \mathbf{y}^T \mathbf{f} - \frac{1}{2} \mathbf{f}^T [\sigma_\epsilon^{-2} I + K^{-1}] \mathbf{f} = \\ &= C - \frac{1}{2} \{ \mathbf{f}^T [\sigma_\epsilon^{-2} I + K^{-1}] \mathbf{f} - 2 \frac{1}{\sigma_\epsilon^2} \mathbf{y}^T \mathbf{f} \} \end{aligned}$$

Denotes  $A = [\sigma_\epsilon^{-2} I + K^{-1}]$  and  $\mathbf{b} = \frac{1}{\sigma_\epsilon^2} \mathbf{y}$ . Completing the square yields

$$\begin{aligned} \log p(\mathbf{f}|\mathcal{D}) &= C - \frac{1}{2} (\mathbf{f}^T A \mathbf{f} - 2\mathbf{b}^T \mathbf{f}) \\ &= C - \frac{1}{2} [(\mathbf{f} - A^{-1} \mathbf{b})^T A (\mathbf{f} - A^{-1} \mathbf{b}) + \mathbf{b}^T A^{-1} \mathbf{b}] \end{aligned}$$

and hence,

$$\begin{aligned}\mathbb{E}[\mathbf{f}|\mathcal{D}] &= K(\sigma_\varepsilon^2 I + K)^{-1} \mathbf{y} \\ \text{cov}[\mathbf{f}|\mathcal{D}] &= \sigma_\varepsilon^2 (\sigma_\varepsilon^2 I + K)^{-1} K\end{aligned}$$

## Appendix C

We would like to compute the posterior predictive distribution defined in Eq. (19) (reproduced here for clarity)

$$p(f_*|\mathcal{D}) = \int p(f_*|\mathbf{f})p(\mathbf{f}|\mathcal{D})d\mathbf{f}$$

where the distributions  $p(f_*|\mathbf{f})$  and  $p(\mathbf{f}|\mathcal{D})$  are Gaussian with

$$\begin{aligned}p(\mathbf{f}|\mathcal{D}) &= \mathcal{N}(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p) \\ p(f_*|\mathbf{f}) &= \mathcal{N}(\mathbf{k}_*^T K^{-1} \mathbf{f}, \mathcal{K}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T K^{-1} \mathbf{k}_*)\end{aligned}$$

We use the following property (Bishop et al., Section 2.3.3) of Gaussian distribution. Property B.1: Consider a marginal Gaussian  $p(\mathbf{x})$  and a Gaussian conditional distribution  $p(\mathbf{y}|\mathbf{x})$  in which  $p(\mathbf{y}|\mathbf{x})$  has a mean that is a linear function of  $\mathbf{x}$ , and a covariance which is independent of  $\mathbf{x}$

$$\begin{aligned}p(\mathbf{x}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Lambda^{-1}), \\ p(\mathbf{y}|\mathbf{x}) &= \mathcal{N}(\mathbf{y}|A\mathbf{x} + \mathbf{b}, L^{-1}).\end{aligned}$$

Here  $\boldsymbol{\mu}, A, \mathbf{b}$  are parameters governing the means, and  $\Lambda, L$  are the precision matrices (inverse covariance matrices). If  $\mathbf{x}$  has dimensionality  $M$  and  $\mathbf{y}$  has dimensionality  $D$  then the matrix  $A$  has size  $D \times M$ . The marginal distribution  $p(\mathbf{y})$  is Gaussian with mean and covariance given by

$$\begin{aligned}\mathbb{E}[\mathbf{y}] &= A\boldsymbol{\mu} + \mathbf{b}, \\ \text{cov}[\mathbf{y}] &= L^{-1} + A\Lambda^{-1}A^T\end{aligned}$$

using property B.1 we obtain

$$\begin{aligned}\mathbb{E}[f_*] &= \mathbf{k}_*^T K^{-1} \boldsymbol{\mu}_p \\ \text{var}[f_*] &= \mathcal{K}(\mathbf{x}_*, \mathbf{x}_*) + \mathbf{k}_*^T (K^{-1} \boldsymbol{\Sigma}_p K^{-1} - K^{-1}) \mathbf{k}_*\end{aligned}$$

## Appendix D

The *matrix inverse lemma*, also known as the Woodbury, Sherman & Morrison formula states that

$$(Z + UWV^T)^{-1} = Z^{-1} - Z^{-1}U(W^{-1} + V^T Z^{-1}U)^{-1}V^T Z^{-1}, \quad (32)$$

assuming the relevant inverses all exist. Here  $Z$  is  $n \times n$ ,  $W$  is  $m \times m$  and  $U$  and  $V$  are both of size  $n \times m$ .

For positive-definite matrices  $P$  and  $R$  we obtain the following variant of the matrix inverse lemma

$$(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = P B^T (B P B^T + R)^{-1}.$$

We note that if  $P$  and  $R$  are also symmetric,  $(P^{-1} + B^T R^{-1} B)^{-1}$  is symmetric function (symmetry is closed under summation and inverse) and hence

$$R^{-1} B (P^{-1} + B^T R^{-1} B)^{-1} = (B P B^T + R)^{-1} B P^T$$