# Bayesian Linear Regression

Why there is a need for a Bayesian formulation?

- MLE (maximum likelihood estimation) is prone over-fitting (especially when the observed data is scarce)

- Oftentimes it desirable to produce uncertainty estimates

<u>Idea</u> : Place a "prior" over the unknown model parameters $\theta$,

$$\hookrightarrow p(\theta)$$

and the use Bayes rule to estimate the optimal paramete via the principle of maximum a-posteriori estimation.

$$
\underbrace{p(\theta|D)}_{\substack{\text{posterior} \\ \text{pdf over the} \\ \text{model params}}} = \frac{\overbrace{p(D|\theta)}^{\text{likelihood}}\ \overbrace{p(\theta)}^{\text{prior}}}{\underbrace{p(D)}_{\substack{\text{model} \\ \text{evidence}}}} = \frac{p(D|\theta)\,p(\theta)}{\underbrace{\int p(D|\theta)\,p(\theta)\,d\theta}_{\substack{\text{marginal} \\ \text{likelihood}}}}
$$

- Recall linear regression :

<u>Setup</u> : Given $D := \{(x_1,y_1), \ldots, (x_m,y_m)\}$, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$

<u>Model</u> : $y_i = w^T x_i + \varepsilon$, if we assume that $\varepsilon \sim N(0, \sigma^2)$

$\qquad\qquad \underbrace{\phantom{xxx}}_{} w^T \phi(x_i)$ in case we use basis fun

$$\implies_{\text{likelihood}} y_i \overset{\text{i.i.d.}}{\sim} N(y_i | w^T x_i, \alpha^{-1}), \text{ where } \alpha = \frac{1}{\sigma^2} \text{ is the precision.}$$

- Unkown params : $\theta := \{w_1, \ldots, w_d, \cancel{\alpha}\}$ (for now let us assume

$\qquad\qquad\qquad\qquad\qquad\qquad \searrow$ that the noise precisi is known)

In a Bayesion formulation we assume

a prior : $w \overset{\text{i.i.d.}}{\sim} N(0, b^{-1} I_d)$, $I_d := \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix}$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\ \ d\times d$

i.e.   a   multi-variate   Gaussian   prior   on   $w \in \mathbb{R}^d$.

<u>Likelihood</u> :  $p(\mathcal{D}|w) \propto \exp\left[-\frac{\alpha}{2}(y-Xw)^T(y-Xw)\right]$  (see lectu

where   $\overset{n \times 1}{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$,  $\overset{n \times d}{X} = \underset{\substack{\text{design} \\ \text{matrix}}}{\begin{bmatrix} x_{11} & - & \cdots & x_{1d} \\ \vdots & & & \\ x_{n1} & - & \cdots & x_{nd} \end{bmatrix}}$,  $\underset{d \times 1}{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$

<u>Posterior</u> :

$$\underbrace{p(w|\mathcal{D})}_{\text{posterior}} \overset{\text{Bayes}}{\propto} \underbrace{p(\mathcal{D}|w)}_{\text{likelihood}} \underbrace{p(w)}_{\text{prior}}$$  $\left(\begin{array}{l}\text{i.e. we ommited} \\ \text{the denominator} \\ \text{in Bayes rule}\end{array}\right.$

$$\implies p(w|\mathcal{D}) \propto \exp\left[-\frac{a}{2}(y-Xw)^T(y-Xw) - \frac{b}{2}w^Tw\right] \text{(verify} \textcircled{1}\text{)}$$

$\mathcal{N}$otice that the exponent is quadratic in $w$. This hints that the posterior $p(w|\mathcal{D})$ is Gaussian. To see this, let us derive the result by " caupleting the square".

First, let us re-writte :

$$\alpha(y-Xw)^T(y-Xw) + bw^Tw \overset{\text{(verify!)}}{=} ay^Ty - 2\boxed{aw^TX^Ty} + w^T(\alpha X^TX + bI)$$

Recall the form of the <u>exponent</u> of a multi-variate Gaussian

$X \sim \mathcal{N}(\mu, \Lambda^{-1})$,  then  $\overset{\curvearrowright}{(X-\mu)^T\Lambda(X-\mu)} =$

$$= X^T\Lambda X - 2X^T\Lambda\mu + \underbrace{\mu^T\Lambda\mu}_{\text{constant (i.e. does not involve } X)}$$

To "match" the terms, let :

$$\boxed{\Lambda := \alpha X^TX + bI \quad \text{(precision)}}$$

we also want :  $\alpha w^TX^Ty = w^T\Lambda\mu$  $\implies$ $\boxed{\mu := \alpha\Lambda^{-1}X^Ty \text{ (mea}}$

Based on these newly defined variables, we can re-write

① $\Rightarrow$ $p(w|D) \propto \exp[(w-\mu)^T \Lambda (w-\mu)]$ ②

$$\boxed{p(w|D) \propto \mathcal{N}(w|\underset{=}{\mu}, \Lambda^{-1}), \quad \begin{cases} \Lambda = a X^T X + bI \\ \mu = \alpha \Lambda^{-1} X^T y \end{cases}}$$

This derivation was possible because:

(i) We assumed a linear model (i.e. model depends linearly or

(ii) We assumed a Gaussian likelihood (i.e. we assumed a Gaussian model for the observation noise)

(iii) We assumed a Gaussian prior over w.

• Maximum a-posteriori estimation for w (MAP):

$$W_{MAP} = \underset{w}{\arg\max} \underbrace{p(w|D)}_{posterior} \quad , \quad W_{MLE} = \underset{w}{\arg\max} \underbrace{p(D|w)}_{likelihood}$$

Recall, :

$\Rightarrow \quad W_{MAP} = \mu = \alpha (\alpha X^T X + bI)^{-1} X^T y$

$$\Rightarrow \boxed{W_{MAP} = (X^T X + \frac{b}{\alpha} I)^{-1} X^T y} \longrightarrow$$

The Bayesian approach natur introduces regularization.

Compare this to

$$\boxed{W_{MLE} = (X^T X)^{-1} X^T y}$$

✱ Equivalently one can see the distinction between MLE vs MA by noticing the following:

$$W_{MLE} = \underset{w}{\arg\min} \|y - Xw\|_2^2$$

$$W_{MAP} = \underset{w}{\arg\min} \|y - Xw\|_2^2 + \underbrace{\lambda \|w\|_2^2}_{regularization} \quad , \quad \lambda = \frac{b}{\alpha}$$

At the end, all we really care about is making predictions (ideally with quantified uncertainty), i.e.

$$p(y^* | x^*, D)$$

$$\underbrace{\hat{y} = f(x^*)}_{\substack{\text{point} \\ \text{prediction}}} \quad , \quad \underbrace{\qquad\qquad}_{\substack{\text{statistical} / \text{probabilistic} \\ \text{prediction.}}}$$

Specifically to the Bayesian linear regression model defined abov

$$P(y^* \mid x^*, \mathcal{D}) = \int \underbrace{P(y^* \mid x^*, \mathcal{D}, w)}_{\substack{\text{likelihood} \\ (\text{Gaussian})}} \underbrace{P(w \mid \mathcal{D})}_{\substack{\text{posterior} \\ (\text{Gaussian})}} \, dw$$

$$\propto \int \exp\left[-\frac{\alpha}{2} (y^* - X w)^T (y^* - X^* w)\right] \exp\left[-\frac{1}{2} (w-\mu)^T \Lambda (w-\mu)\right]$$

$$\cdots \implies \boxed{P(y^* \mid x^*, \mathcal{D}) = \mathcal{N}\left(y^* \mid u, 1/\lambda\right)}, \quad \text{where}$$

$$\underset{\substack{\text{predictive posterior} \\ \text{distribution}}}{}$$

$$\begin{cases} u = \mu^T x^* \\ \frac{1}{\lambda} = \frac{1}{\alpha} + x^{*T} \Lambda^{-1} x^* \end{cases}$$

$$\text{where} \quad \begin{cases} \mu = W_{MAP} = \left(X^T X + \frac{b}{\alpha} I\right)^{-1} X^T y \\ \Lambda = \alpha X^T X + b I \end{cases}$$