

Ref for Bayesian estimation: Bishop 2.3.6

From Last Time (with corrections): Bayesian Regression relies on integration

$$p(w|y) = N(w|m_n, S_n)$$

$$\text{if } p(w) = N(w, 0, \alpha^{-1}I)$$

then

$$m_n = \beta S_n \Phi^T y$$

$$S_n^{-1} = \alpha I + \beta \Phi^T \Phi$$

substituting

$$m_n = \beta (\alpha I + \beta \Phi^T \Phi)^{-1} \Phi^T y$$

$$= \left(\frac{1}{\beta}\right)^{-1} (\alpha I + \beta \Phi^T \Phi)^{-1} \Phi^T y$$

$$= \left(\frac{\alpha}{\beta} I + \Phi^T \Phi\right)^{-1} \Phi^T y$$

$$\text{if not scalar } (AB)^{-1} = B^{-1} A^{-1}$$

We can also directly maximize the posterior:

Maximum a posteriori (MAP) estimate: \rightarrow based on optimizing something

$$\max_w p(w|y)$$

$$\min_w -\log p(w|y)$$

$$= \min_w \frac{\beta}{2} \sum (y_i - w^T \phi(x_i))^2 + \frac{\alpha}{2} w^T w$$

Covariance Matrix Review

What is covariance?

Def.

$$\begin{aligned} \text{Cov}(x, y) &= \mathbb{E}[(x - \mathbb{E}x)(y - \mathbb{E}y)] \\ &= \mathbb{E}[xy] - \mathbb{E}x \mathbb{E}y \end{aligned}$$

ex: If we have a population

$$w_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix}, \quad i = 1 \dots n$$

$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\begin{aligned} \text{if } x_i > \bar{x} \text{ and } y_i > \bar{y} &\text{ then } \text{Cov}(x, y) > 0 \\ x_i < \bar{x} \text{ and } y_i < \bar{y} &\text{ " } > 0 \\ x_i > \bar{x} \text{ and } y_i < \bar{y} &\text{ " } < 0 \end{aligned}$$

With some data with $\bar{x} = 10$, $\bar{y} = 20$

$$w_1 = \begin{bmatrix} 0 \\ 20 \end{bmatrix}, w_2 = \begin{bmatrix} 11 \\ 21 \end{bmatrix}, w_3 = \begin{bmatrix} 9 \\ 19 \end{bmatrix}$$



$$\text{cov}(x, y) = \frac{1}{3} \left((0)(0) + (1)(1) + (-1)(-1) \right) = \frac{2}{3}$$

Now

$$w_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, w_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, w_3 = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \quad \text{cov}(x, y) = \dots = \frac{2}{3}$$

We can also calculate the covariance matrix

$$\frac{1}{n} \sum (w_i - \bar{w})(w_i - \bar{w})^T = 2 \times 2 \text{ matrix}$$

first case

$$\frac{1}{3} \left[\begin{bmatrix} 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 & 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} + \begin{bmatrix} -1 \\ -1 \end{bmatrix} \begin{bmatrix} -1 & -1 \end{bmatrix} \right] = \frac{1}{3} \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}$$

Symmetric

second case

$$\begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} = \frac{1}{3} \begin{bmatrix} 2 & -2 \\ -2 & 2 \end{bmatrix} \quad \text{where non-diagonal entries are covariances}$$

Independent \rightarrow uncorrelated $\text{cov}(x, y) = 0$

uncorrelated \nrightarrow independent

... but if $w \sim N(\mu, \Sigma)$, Σ diagonal (all off diag entries are zero)

then uncorrelated \rightarrow independent

reason:

$$\text{Since } N(w | \mu, \Sigma) = \prod_{i=1}^m N(w_i | \mu_i, \sigma_i^2)$$

Recall $p(x, y) = p(x)p(y)$

$$\mu = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_m \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_m^2 \end{bmatrix}$$

uncorrelated: $\mathbb{E}(x, y) = \mathbb{E}x \mathbb{E}y$

Back to Bayesian Regression

We no longer have a point prediction for \vec{w}

We instead have a distribution \Rightarrow this will allow us to put levels of confidence on our estimates

Predictive distribution

defn.

$$p(\hat{y} | y) = \int \underbrace{p(\hat{y} | w)}_{\text{posterior}} p(w | y) dw \quad N(w | m_n, S_n)$$

$$N(\hat{y} | w^T \phi(x), \beta^2)$$

\leftarrow remember this is a

result of the noise

$$p(\hat{y} | y) = N(\hat{y} | m_n^T \phi(x), \sigma_n^2(x))$$

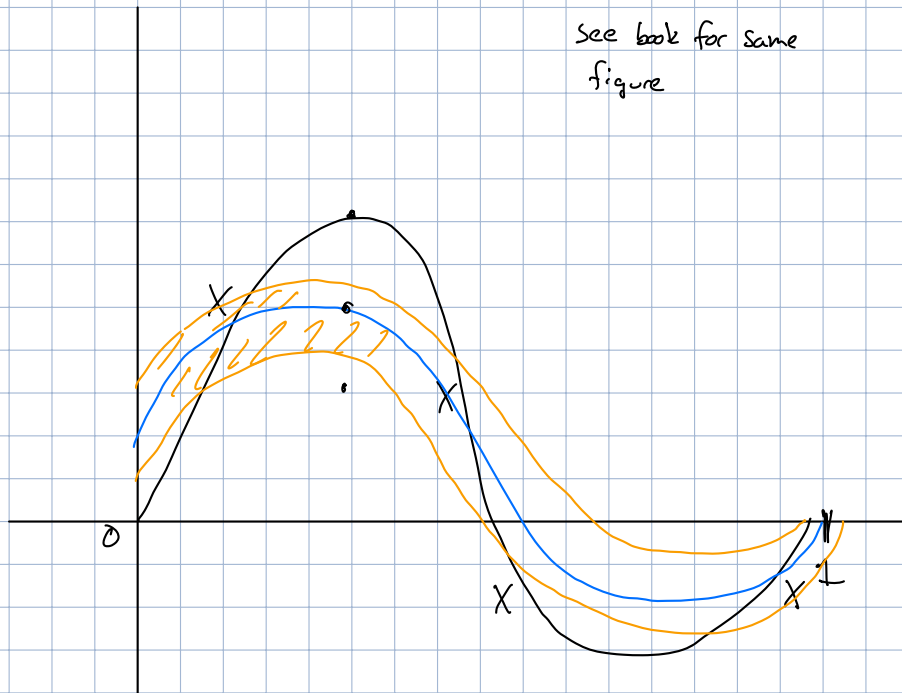
these are all scalars

where $\sigma_n^2 = \frac{1}{\beta} + \phi(x)^T S_n \phi(x)$

↑
noise

↑
arises from uncertainty
in w

See book for same
figure



in the demo we see that we have larger uncertainty in areas where we don't have data

For different values of x we have predictions with different variances $\sigma_n^2(x)$

Instead of a single curve that fits the data we have a family of curves

The mean of the family is $m^T \phi(x)$

another example of the family is $w^T \phi(x)$

$$w \sim N(m, S_n)$$

Unpacking the predictive distribution some more

Take

$$\hat{y} = w_0 + w_1 x$$

$$= \int (w_0 p(w_0 | y) + w_1 p(w_1 | y) x) dw_0 dw_1$$

$$= \int w^T \phi(x) p(w | y) dw$$

\Rightarrow the integral is prediction averaged over all w
like formula for expectation