

Data Science 423

Homework 5

Dennis Lee

November 26, 2019

1 For Bayesian Linear Regression, we make use of Bayes' Theorem:

$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Normalization}}$$

* I'm using the books notation.

Where our data is X & $X^T \vec{\varphi} \in \Theta = [\vec{\varphi}, \sigma^2]$

$$\Pr(\vec{\varphi} | X, \vec{w}) = \frac{\Pr(\vec{w} | X, \vec{\varphi}) P(\vec{\varphi})}{\Pr(\vec{w} | X)}$$

We model the prior as a multivariate normal distribution:

$$\text{Prior: } \Pr(\vec{\varphi}) = N_{\vec{\varphi}}[\vec{0}, \sigma_p^2 I]$$

For the likelihood with we write the likelihood as:

$$\text{Likelihood: } \Pr(\vec{w} | X, \theta) = N_{\vec{w}}[X^T \vec{\varphi}, \sigma^2 I]$$

Focusing on the numerator for now:

$$\Pr(\vec{\varphi} | X, \vec{w}) \propto N_{\vec{\varphi}}[X^T \vec{\varphi}, \sigma^2 I] N_{\vec{\varphi}}[\vec{0}, \sigma_p^2 I]$$

we then make use of several relations.

$$\Pr(\varphi | X, \vec{\omega}) \propto N_{\vec{\omega}} \underbrace{[X^T \vec{\varphi}, \sigma^2 I]}_{N_{\vec{\varphi}}[\emptyset, \sigma_p^2 I]}$$

$$\text{Using: } N_x[Ay + b, \Sigma] = \kappa \cdot N_y[A'x + b', \Sigma']$$

we find that in our case:

$$\Sigma = \sigma^2 I \rightarrow \Sigma' = (X \sigma^2 X^T)^{-1}$$

$$A = X^T \rightarrow A' = (X \sigma^2 X^T)^{-1} X \sigma^{-2}$$

$$b = \emptyset \rightarrow b' = \emptyset$$

Thus:

$$\Pr(\varphi | X, \vec{\omega}) \propto N_{\vec{\varphi}} \left[(\sigma^2 X X^T)^{-1} X \sigma^{-2} \vec{\omega}, (\sigma^2 X X^T)^{-1} \right]$$

$$\times N_{\vec{\varphi}} \left[\emptyset, \sigma_p^2 I \right]$$

Combining using:

$$N_x[a, A] N_x[b, B] = \kappa N_x \left[(A^{-1} + B^{-1})^{-1} (A^{-1} a + B^{-1} b), (A^{-1} + B^{-1})^{-1} \right]$$

$$\text{For us: } a = (\sigma^2 X X^T)^{-1} X \sigma^{-2} \vec{\omega} \quad b = \emptyset$$

$$A = (\sigma^2 X X^T)^{-1} \quad B = \sigma_p^2 I$$

$$\Pr(\varphi | X, \vec{\omega}) \propto \kappa N_{\vec{\varphi}} \left[\left[(\sigma^2 X X^T) + (\sigma_p^2 I) \right]^{-1} \left((\sigma^2 X X^T)^{-1} X \sigma^{-2} \vec{\omega} \right) \right.$$

$$\left. , \left[(\sigma^2 X X^T) + (\sigma_p^2 I) \right]^{-1} \right]$$

$$\Pr(\varphi | X, \omega) \propto \kappa N_\varphi \left[\left[(\sigma^2 XX^T) + (\sigma_p^2 I)^T \right]^{-1} \left((\sigma^2 XX^T) (\sigma^2 XX^T)^T X \sigma_w^2 \right) \right.$$

$$\left. , \left[(\sigma^2 XX^T) + (\sigma_p^2 I)^T \right]^{-1} \right]$$

$$\propto \kappa N_\varphi \left[\underbrace{\left[(\sigma^2 XX^T) + \sigma_p^2 I \right]^{-1}}_{\text{If we define } A = \sigma^2 XX^T + \sigma_p^2 I} (X \sigma^2 \omega), \underbrace{\left[(\sigma^2 XX^T) + \sigma_p^2 I \right]^{-1}}_{\text{where } A = \sigma^2 XX^T + \sigma_p^2 I} \right]$$

If we define $A = \sigma^2 XX^T + \sigma_p^2 I$

$$\Pr(\varphi | X, \omega) \propto N_\varphi \left[\sigma^2 A^{-1} X \omega, A^{-1} \right]$$

$$\Pr(\varphi | X, \omega) \propto N_\varphi \left[\sigma^2 A^{-1} X \omega, A^{-1} \right] \text{ where } A = \sigma^2 XX^T + \sigma_p^2 I$$

Since I used the book notation, to use the notation of the homework, we get:

$$\Pr(\tilde{\omega} | \tilde{X}, \bar{y}) = N_{\tilde{\omega}} \left[\sigma^2 \tilde{A}^{-1} \tilde{X} \bar{y}, \tilde{A}^{-1} \right]$$

$$\tilde{A} = \sigma^2 \tilde{X} \tilde{X}^T + \sigma_p^2 I$$

2

To show that:

$$\Pr(y^* | \bar{x}^*, \tilde{X}, \bar{y}) = N_{y^*} \left[\sigma^2 \bar{x}^{*T} A^{-1} \tilde{X} \bar{y}, \bar{x}^{*T} A^{-1} \bar{x}^* + \sigma^2 \right]$$

we start with the integral:

$$\Pr(y^* | \bar{x}^*, \tilde{X}, \bar{y}) = \int N_{y^*} \left[\tilde{w}^T \bar{x}^*, \sigma^2 \right] N_{\tilde{w}} \left[\frac{1}{\sigma^2} A^{-1} \tilde{X} \bar{y}, A^{-1} \right] d\tilde{w}$$

predictions
of each \tilde{w}
posterior
distribution

Since each of the predictions are Gaussian, we can argue that the predictive distribution will also likewise be gaussian. Alternatively, we can see that if we go through all the steps, as in problem 1, we will receive a constant defined as a Gaussian when computing the product of two normals. Either way, knowing this, we simply compute the mean (μ_{pred}) & variance (σ^2_{pred}) to define the Gaussian.

Looking at the integral, we see that the expectation value/mean of y^* is \bar{x}^* multiplied by the expectation value of the weight \tilde{w} . Since that's just the mean:

$$\begin{aligned} \mu_{\text{pred}} &= E[y^*] = E[\tilde{w}^T \bar{x}^*] = \bar{x}^{*T} E[\tilde{w}] \\ &= \bar{x}^{*T} \left(\frac{1}{\sigma^2} A^{-1} \tilde{X} \bar{y} \right) \end{aligned}$$

From the Computer Vision textbook (Ch4)

we can write:

$$\sigma_{\text{pred}}^2 = \sigma^2 + \bar{x}^* T A^{-1} \bar{x}^*$$

Combining all of this, we have:

$$P(y^* | \bar{x}^*, \hat{x}, \bar{y}) = N y^* \left[\sigma^{-2} \bar{x}^* T A^{-1} \bar{x}^*, \bar{x}^* T A^{-1} \bar{x}^* + \sigma^2 \right]$$

Given the logistic regression model:

$$L = \sum_{i=1}^P y_i \log[\sigma(a_i)] + \sum_{i=1}^P (1-y_i) \log[1-\sigma(a_i)]$$

where $a_i = \bar{x}_i^T \tilde{w}$, the gradient with respect to \tilde{w} is given by:

$$\begin{aligned}\nabla_{\tilde{w}} L &= \sum_{i=1}^P y_i \frac{1}{\sigma(a_i)} \nabla_w \sigma(a_i) + \sum_{i=1}^P (1-y_i) \frac{1}{1-\sigma(a_i)} [-\nabla_w \sigma(a_i)] \\ &= \sum_{i=1}^P y_i (1-\sigma(a_i)) \nabla_w a_i + \sum_{i=1}^P (1-y_i) (-1) \sigma(a_i) \nabla_w a_i \\ &= \sum_{i=1}^P y_i (1-\sigma(a_i)) \bar{x}_i + \sum_{i=1}^P (1-y_i) (-1) \sigma(a_i) \bar{x}_i \\ &= \sum_{i=1}^P [-1(y_i + \sigma(a_i)y_i) + (-1)(\sigma(a_i) - \sigma(a_i)y_i)] \bar{x}_i \\ &= - \sum_{i=1}^P (-y_i + \sigma(a_i)y_i + \sigma(a_i) - \sigma(a_i)y_i) \bar{x}_i \\ &= - \sum_{i=1}^P (-y_i + \sigma(a_i)) \bar{x}_i\end{aligned}$$

$$\boxed{\nabla_{\tilde{w}} L = - \sum_{i=1}^P (\sigma(a_i) - y_i) \bar{x}_i}$$

4

Starting with:

$$\nabla_{\tilde{\omega}} L = - \sum_{i=1}^P (\sigma(a_i) - y_i) \bar{x}_i$$

we find the Hessian:

$$\begin{aligned}\nabla^2_{\tilde{\omega}} L &= - \sum_{i=1}^P \nabla_{\tilde{\omega}} \sigma(a_i) \bar{x}_i \\ &= - \sum_{i=1}^P \sigma(a_i)(1-\sigma(a_i)) \bar{x}_i \nabla_{\tilde{\omega}} a_i \\ &= - \sum_{i=1}^P \sigma(a_i)(1-\sigma(a_i)) \bar{x}_i \bar{x}_i^T\end{aligned}$$

$$\boxed{\nabla^2_{\tilde{\omega}} L = - \sum_{i=1}^P \sigma(a_i)(1-\sigma(a_i)) \bar{x}_i \bar{x}_i^T}$$