

Lecture 9: Bayesian Inference, Continued

Lecturer: Matthew Wicker

1 Learning Objectives

In our last lecture we introduced Bayes theorem and how we can use it to express uncertainty particularly in density estimation. In this lecture, we are going to specifically derive how Bayes theorem prescribes inference in a conditional density setup (i.e., supervised learning) using our linear regression model. We will end these notes by thinking about how having a posterior distribution over our parameters yields predictive uncertainty and in lecture we discussed how uncertainty allows us to perform model selection.

2 Background

It has been a few lectures since we last reviewed our linear regression model. So here, we will first write out some brief review notes, and then will recall Bayes theorem in the context of conditional density estimation.

2.1 Linear Regression Model

Recall that our basis expanded linear regression model takes the form:

$$\hat{\mathbf{y}} = \phi(\mathbf{X})\theta$$

Where $\mathbf{X} \in \mathbb{R}^{n \times d}$ is our feature matrix and $\phi : \mathbb{R}^n \mapsto \mathbb{R}^k$ is a basis expansion function and finally $\theta \in \mathbb{R}^{k \times m}$ is the parameter that we would like to fit. We will throughout make the simplifying assumption that $m = 1$ and that the vector \mathbf{y} is a column vector containing all of our response variables. Though the basis expansion function allows us to capture non-linear relationships in data we call this a linear model because the relationship between the basis-expanded features and the predictions is linear w.r.t. the parameter vector. It is helpful to recall that we have computed the maximum likelihood estimate of this model as:

$$\theta^{\text{MAP}} = (\mathbf{X}^\top \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I}_d)^{-1} \mathbf{X}^\top \mathbf{y}$$

2.2 Bayes Theorem

In our last lecture, we introduced Bayes theorem as critical component to a probabilistic learning paradigm:

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}$$

We used this formula to discuss conjugate priors and derive the posterior for Bernoulli random variables. In this lecture, we will deal with the slightly more complex setting of conditional density estimation where we would like to ultimately reason about $p(y|\mathbf{x}^*, \theta)$ where \mathbf{x}^* is some unseen feature vector.

3 Bayesian Linear Regression

In any Bayesian machine learning model, there are two critical ingredients that we must carefully consider: the likelihood and the prior. To start, we will assume that we have a single dimensional response/label.

Likelihood The likelihood in supervised learning has the critical task of modelling the observational noise of a process. Throughout our exposition of linear models, we have focused on Gaussian noise (and will continue to do so). Our likelihood with respect to a given observation (x, y) can be expressed as:

$$p(y|x, \theta) = \mathcal{N}(\phi(x)\theta, \sigma^2)$$

We can write out the full equation in matrix form for all of our observations (and plugging into the Gaussian density equation):

$$p(\mathbf{y}|\mathbf{X}, \theta) = \frac{1}{(2\pi\sigma^2)^{N/2}} \cdot \exp\left(-\frac{(\mathbf{y} - \phi(\mathbf{X})\theta)^\top (\mathbf{y} - \phi(\mathbf{X})\theta)}{2\sigma^2}\right) \quad \mu = \phi(x)\theta \quad \sigma^2 = \sigma^2$$

We will make the simplifying assumption that σ is known.

Prior As we have seen in our previous exercises the conjugate prior for the mean of a Gaussian with known variance is a Gaussian, so we know that in this case it will be analytically convenient to select a Gaussian prior:

$$p(\theta) = \mathcal{N}(0, \tau^2 \mathbf{I})$$

where \mathbf{I} is the $k \times k$ identity matrix, our prior is an isotropic Gaussian with variance τ^2 . We have actually seen in previous lectures that we can express this distribution as:

$$p(\theta) = \frac{1}{(2\pi\tau^2)^{k/2}} \cdot \exp\left(-\frac{\theta^\top \theta}{2\tau^2}\right) \quad \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}}, \text{ need } x \rightarrow \theta, \mu=0, \sigma^2=\tau^2$$

We will now explore two ways of deriving the Bayesian posterior for our model.

3.1 Deriving the Posterior via Densities

A concrete way to compute the posterior distribution use the method of equating coefficients. The first step is to simply plug in the densities and use algebra to simplify until we arrive at a form that looks approximately like the posterior form we expect (from conjugacy). Here, let us consider computing the posterior by taking the log of both sides of our unnormalized posterior:

$$\log(p(\theta|\mathcal{D})) \propto \log(p(\theta|\mathcal{D})) + \log(p(\theta)) \quad (1)$$

$$p(\theta|\mathcal{D}): \text{const} e^{\frac{(\mathbf{y} - \phi(\mathbf{X})\theta)^\top (\mathbf{y} - \phi(\mathbf{X})\theta)}{2\sigma^2}} = \frac{1}{2\sigma^2} ((\mathbf{y} - \phi(\mathbf{X})\theta)^\top (\mathbf{y} - \phi(\mathbf{X})\theta)) + \frac{1}{2\tau^2} (\theta^\top \theta) \quad (2)$$

$$p(\theta): \text{const} e^{\frac{\theta^\top \theta}{2\tau^2}}, \tau \text{ is } \sigma \text{ for this} = \frac{1}{2\sigma^2} \left(\theta^\top \phi(\mathbf{X})^\top \phi(\mathbf{X}) \theta - 2\mathbf{y}^\top \phi(\mathbf{X}) \theta + \mathbf{y}^\top \mathbf{y} \right) + \frac{1}{2\tau^2} \theta^\top \theta \quad (3)$$

$$= \frac{1}{2\sigma^2} \theta^\top \phi(\mathbf{X})^\top \phi(\mathbf{X}) \theta + \frac{2}{2\sigma^2} \mathbf{y}^\top \phi(\mathbf{X}) \theta - \frac{1}{2\sigma^2} \mathbf{y}^\top \mathbf{y} + \frac{1}{2\tau^2} \theta^\top \theta \quad (4)$$

$$\propto \tau^2 \theta^\top \phi(\mathbf{X})^\top \phi(\mathbf{X}) \theta + 2\tau^2 \mathbf{y}^\top \phi(\mathbf{X}) \theta - \tau^2 \mathbf{y}^\top \mathbf{y} + \sigma^2 \theta^\top \theta \quad (5)$$

While perhaps unintuitive at first, we can see that the above equation is quadratic in θ which is the variable our posterior density is over. Thus, we know we ought to be able to proceed with equating coefficients with a Gaussian. To start, we group terms that look like $\theta^\top \mathbf{M} \theta$, and pull out the θ s:

$$\text{group } \theta^\top \theta, \theta = \tau^2 \theta^\top \phi(\mathbf{X})^\top \phi(\mathbf{X}) \theta + \sigma^2 \theta^\top \theta + 2\tau^2 \mathbf{y}^\top \phi(\mathbf{X}) \theta + \tau^2 \mathbf{y}^\top \mathbf{y} \quad (6)$$

$$= \theta^\top (\tau^2 \phi(\mathbf{X})^\top \phi(\mathbf{X}) + \sigma^2 \mathbf{I}) \theta + 2\tau^2 \mathbf{y}^\top \phi(\mathbf{X}) \theta - \tau^2 \mathbf{y}^\top \mathbf{y} \quad (7)$$

Before doing anymore algebra, recall the reason we chose a Gaussian prior is because we knew it was the conjugate prior to a Gaussian posterior from our previous exercises. And we know that we can write any Gaussian distribution as:

$$\log(\mathcal{N}(\theta; \mu, \Sigma)) \propto (\theta - \mu)^\top \Sigma^{-1} (\theta - \mu) \quad (8)$$

$$= \theta^\top \Sigma^{-1} \theta - 2\theta^\top \Sigma^{-1} \mu + \mu^\top \Sigma^{-1} \mu \quad (9)$$

$$= \theta^\top \Sigma^{-1} \theta - 2\theta^\top \Sigma^{-1} \mu + \text{const.} \quad (10)$$

So, now we have two equations, the one above (Equation (10)) is that of the log of a Gaussian, and our previous Equation (7) which is the log of our posterior which we also know is Gaussian. So what we would like to do is somehow get the two equal to each other. First, we can drop the constant in equation (10) because we only have proportionality. Then, we can see that the first term in (7) looks a lot like the first term in (10). So it is natural to first set:

$$\Sigma^{-1} = (\tau^2 \phi(\mathbf{X})^\top \phi(\mathbf{X}) + \sigma^2 \mathbf{I})$$

. Finally, we want to take care of the next term where we have:

$$2\theta^\top \Sigma^{-1} \mu = 2\tau^2 \mathbf{y}^\top \phi(\mathbf{X}) \theta$$

$$2\tau^2 \mathbf{y}^\top \phi(\mathbf{X}) \theta = 2\mathbf{M}^\top \boldsymbol{\varepsilon}^{-1} \theta$$

$$\mathbf{M}^\top: \tau^2 \mathbf{y}^\top \phi(\mathbf{X}) \theta (\boldsymbol{\varepsilon}^{-1} \theta)^{-1} = \tau^2 \mathbf{y}^\top \phi(\mathbf{X}) \theta \theta^\top \boldsymbol{\varepsilon} = \tau^2 \mathbf{y}^\top \phi(\mathbf{X}) \boldsymbol{\Sigma}$$

$$\boldsymbol{\mu}: \boldsymbol{\varepsilon}^\top \phi(\mathbf{X})^\top \mathbf{y}$$

and we know the final parameter that we must find equality for is the mean, μ . Notice that we do not need to worry about the $\mathbf{y}^\top \mathbf{y}$ term as this is a constant with respect to μ . Luckily, we already have Σ^{-1} in the terms of equation (7), so we can simply solve for μ to get that:

$$\begin{aligned} 2\theta^\top \Sigma^{-1} \mu &= 2\tau^2 \mathbf{y}^\top \phi(\mathbf{X}) \theta \\ \mu &= \tau^2 \Sigma \phi(\mathbf{X})^\top \mathbf{y} \end{aligned}$$

So now that we have solved for the mean and covariance, we can say that the posterior of our Bayesian linear regression model is:

$$\mathcal{N}(\theta; \tau^2 \Sigma \phi(\mathbf{X})^\top \mathbf{y}, (\tau^2 \phi(\mathbf{X})^\top \phi(\mathbf{X}) + \sigma^2 \mathbf{I})^{-1})$$

3.1.1 Interpreting our Posterior Values

Let us take a closer look at the posterior that we have computed. In particular, let's reason about if the expressions we have come up with for the mean and covariance matrix of our Bayesian linear regression make intuitive sense. This will not only deepen our understanding of the quantities we have just derived, but will also let us understand if there are conditions under which our inference is not well defined. Beginning with the posterior covariance matrix (in no small part because the mean relies on the covariance):

$$\Sigma = (\frac{1}{\sigma^2} \phi(\mathbf{X})^\top \phi(\mathbf{X}) + \sigma^2 \mathbf{I})^{-1}$$

We know that not all matrices are invertible, yet our covariance matrix relies on the matrix inverse of a particular quantity. We can start with the fact that $\sigma^2 \mathbf{I}$ and $\tau^2 \phi(\mathbf{X})^\top \phi(\mathbf{X})$ are positive semi-definite matrices. However, not all PSD matrices are invertible. It is a useful exercise, and is covered in lecture, to prove that the covariance we have computed is indeed invertible. Next, let's reason about the mean:

$$\mu = \tau^2 \Sigma \phi(\mathbf{X})^\top \mathbf{y}$$

One thing we know about normal distributions is that their mode, the value maximizing the p.d.f. coincides with the mean of the distribution. In previous lectures, we introduced the concept of a *maximum a posteriori* (MAP) estimate, which is the parameter which has the highest probability when we assume a prior distribution over our weights. So here, we would like to check that the quantity of our mean is equivalent to the MAP estimate we computed in previous lectures. Expanding our equation for the mean we have:

$$\begin{aligned} \mu &= \tau^2 \Sigma \phi(\mathbf{X})^\top \mathbf{y} \\ &= \tau^2 (\tau^2 \phi(\mathbf{X})^\top \phi(\mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \phi(\mathbf{X})^\top \mathbf{y} \\ &= (\phi(\mathbf{X})^\top \phi(\mathbf{X}) + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} \phi(\mathbf{X})^\top \mathbf{y} \end{aligned}$$

Now, if we return to lecture 5, we will see this exactly the same as the quantity we derived before. With both of these, we have done a pretty good sanity check that our posterior makes sense.

3.2 Deriving the Posterior via Joint Gaussian

Another method for arriving at the posterior distribution is to start by modelling the joint distribution of θ and \mathbf{y} as a Gaussian distribution. It may not be clear exactly how that helps us, but recall that the product of densities that we reasoned about above was easy to interpret as a joint distribution. Moreover, there are many probability rules that allow us to manipulate joint distributions to compute their marginal and conditional distributions. Gaussian distributions have particularly nice rules for marginalizing and conditioning. We state them below:

Theorem 3.1. Marginalization Given a Gaussian random variable:

$$p(a,b) = \mathcal{N}\left(\begin{bmatrix} a \\ b \end{bmatrix}; \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}, \begin{bmatrix} \Sigma_{a,a} & \Sigma_{a,b} \\ \Sigma_{b,a} & \Sigma_{b,b} \end{bmatrix}\right)$$

$\mathcal{N}(a)$ $\mathcal{N}(a,b)$

The marginal distribution of a is given by:

$$p(a) = \mathcal{N}(a; \mu_a, \Sigma_{a,a})$$

Theorem 3.2. Conditioning Given a Gaussian random variable:

$$\mathcal{N}\left(\begin{bmatrix} a \\ b \end{bmatrix}; \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}, \begin{bmatrix} \Sigma_{a,a} & \Sigma_{a,b} \\ \Sigma_{b,a} & \Sigma_{b,b} \end{bmatrix}\right)$$

The conditional distribution of $p(a|b)$ is given by:

$$\begin{aligned} p(a|b) &= \mathcal{N}(a; \mu_{a|b}, \Sigma_{a|b}), \\ \mu_{a|b} &= \mu_a + \Sigma_{a,b} \Sigma_{b,b}^{-1} (b - \mu_b) \\ \Sigma_{a|b} &= \Sigma_{a,a} - \Sigma_{a,b} \Sigma_{b,b}^{-1} \Sigma_{b,a} \end{aligned}$$

$p(a|b)$ $p(b)$

Now that we have established these theorems, we will show how we can use them to compute the Bayesian linear regression posterior. We start by writing out the joint distribution:

$$p(\theta, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \theta \\ \mathbf{y} \end{bmatrix}; \begin{bmatrix} \mathbb{E}_{p(\theta)}[\theta] \\ \mathbb{E}_{\mathbf{y}}[\mathbf{y}] \end{bmatrix}, \begin{bmatrix} \mathbb{V}[\theta], \mathbb{C}[\theta, \mathbf{y}] \\ \mathbb{C}[\mathbf{y}, \theta], \mathbb{V}[\mathbf{y}] \end{bmatrix}\right)$$

Ultimately, as we show in our lecture slides, we can compute the values for each entry in the joint Gaussian above and plug into our conditioning theorem to arrive at the posterior as:

$$p(\theta|\mathbf{y}) = \mathcal{N}\left(\theta; \Phi(X)^\top [\Phi(X)\Phi(X)^\top + \sigma^2 \mathbf{I}_N]^{-1} \mathbf{y}, \mathbf{I}_M - \Phi(X)^\top [\Phi(X)\Phi(X)^\top + \sigma^2 \mathbf{I}_N]^{-1} \Phi(X)\right)$$

Notice, that this form of the posterior is different than the one we arrived at via equating coefficients. In your exercises, you will be asked to prove that these two are indeed identical.

4 Posterior Predictive Distribution

We know that when we make predictions about a new, unseen point x^* in linear regression we simply plug the value into our model, i.e., multiply by our parameters, and the result is a value that we call \hat{y} , our prediction. However, in our Bayesian linear regression model we do not have a single set of parameters, instead we have a continuous distribution over parameters, which makes the idea of “plugging” unseen values into our model a little bit less clear. However, the key thing to understand is that in a Bayesian framework (almost) everything is treated as a probability distribution including our prediction. So what is the exact form of the probability distribution that we seek? Well we know we want the probability distribution to be over the space of outputs y and we know that we have access to the unseen value x^* and our previous dataset, \mathcal{D} , so we can write out the distribution we are interested in as: $p(y|x^*, \mathcal{D})$. This is known as the posterior predictive distribution. To compute this, we of course can use our basic rules of probability to arrive at an equation for this distribution in terms of quantities we already know:

$$p(y|x^*, \mathcal{D}) = \int_{\theta \in \Theta} p(y|x^*, \theta) p(\theta|\mathcal{D}) d\theta$$

$p(x) = \int p(x|y) p(y)$

Notice that this uses slightly odd notation as x^* is not a random variable. However, we include it on the conditioning side to indicate that the distribution over outputs is taken with respect to the input x^* . The first is just our likelihood and the second is our posterior distribution, both of which we know. So we can plug in our densities and then think about how to attack this integral:

$$\begin{aligned} p(y|x^*, \mathcal{D}) &= \int_{\theta \in \Theta} p(y|x^*, \theta) p(\theta|\mathcal{D}) d\theta \\ &= \int_{\theta \in \Theta} \mathcal{N}(\phi(x^*)^\top \theta, \tau^2) \mathcal{N}(\theta; \mu, \Sigma) d\theta \\ &\propto \int_{\theta \in \Theta} \exp\left(\frac{1}{\tau^2}(y - \phi(x^*)^\top \theta)^2\right) \exp\left(-\frac{1}{2}(\theta - \mu)^\top \Sigma^{-1}(\theta - \mu)\right) d\theta \\ &= \int_{\theta \in \Theta} \exp\left(\frac{1}{\tau^2}(y^2 - 2(\theta^\top \phi(x^*))y + (\theta^\top \phi(x^*))^2)\right) \exp\left(-\frac{1}{2}(\theta - \mu)^\top \Sigma^{-1}(\theta - \mu)\right) d\theta \\ &= \int_{\theta \in \Theta} \exp\left(\frac{1}{\tau^2}(y^2 - 2(\theta^\top \phi(x^*))y + (\theta^\top \phi(x^*))^2) - \frac{1}{2}(\theta^\top \Sigma \theta - 2\theta^\top \Sigma^{-1}\mu + \mu^\top \Sigma^{-1}\mu)\right) d\theta \\ &\propto \int_{\theta \in \Theta} \exp\left(\frac{1}{\tau^2}(y^2 - 2(\theta^\top \phi(x^*))y + (\theta^\top \phi(x^*))^2) - \frac{1}{2}(\theta^\top \Sigma \theta - 2\theta^\top \Sigma^{-1}\mu)\right) d\theta \end{aligned}$$

We will pause the derivation here and make a key observation: the above expression is quadratic in y and so, as in our derivation of the posterior, we can complete the square to arrive at the closed form normal distribution for the posterior predictive. We do not complete the derivation here and instead leave it as an exercise problem.

A Lecture 9: Bayesian Inference, Continued

Question 1 (Vector Ordering in Gaussians). Consider a joint Gaussian density on a vector \mathbf{z} that can be split up as $\mathbf{z} = [\mathbf{x}^\top, \mathbf{y}^\top]^\top$ (this notation denotes stacking):

$$p(\mathbf{z}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}; \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{C} \end{bmatrix}\right). \quad (11)$$

Now consider the permuted vector $\mathbf{z}' = [\mathbf{y}^\top, \mathbf{x}^\top]^\top$. Show that the Gaussian distribution of $p(\mathbf{z}')$ has a mean with rows, and a covariance matrix with swapped rows and columns.

Hint: Notice that \mathbf{z}' is a permuted version of \mathbf{z} . We can write this mathematically as $\mathbf{z}' = \mathbf{P}\mathbf{z}$, where \mathbf{P} is a permutation matrix:

$$\mathbf{P} = \begin{bmatrix} 0 & \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix}. \quad (12)$$

Note: This is an important skill! Notice that the Gaussian conditioning formula in the formula sheet is only written in one way.

Question 2 (Woodbury Identity). We saw that the two different ways of deriving the Gaussian posterior, gave two different results:

$$\begin{aligned} p(\boldsymbol{\theta}|\mathbf{y}) &= \mathcal{N}\left(\boldsymbol{\theta}; \Phi(X)^\top [\Phi(X)\Phi(X)^\top + \sigma^2 \mathbf{I}_N]^{-1} \mathbf{y}, \quad \mathbf{I}_M - \Phi(X)^\top [\Phi(X)\Phi(X)^\top + \sigma^2 \mathbf{I}_N]^{-1} \Phi(X)\right), \\ p(\boldsymbol{\theta}|\mathbf{y}) &= \mathcal{N}\left(\boldsymbol{\theta}; \left[\frac{1}{\sigma^2} \Phi(X)^\top \Phi(X) + \mathbf{I}_M\right]^{-1} \frac{1}{\sigma^2} \Phi(X)^\top \mathbf{y}, \quad \left[\frac{1}{\sigma^2} \Phi(X)^\top \Phi(X) + \mathbf{I}_M\right]^{-1}\right). \end{aligned} \quad (13)$$

Apply the Woodbury identity to $\left[\frac{1}{\sigma^2} \Phi(X)^\top \Phi(X) + \mathbf{I}_M\right]^{-1}$ to show that the two solutions are equal.

Note: This is an excellent exercise for your matrix algebra skills. The main difficulty that most uninitiated have, is that multiplication no longer commutes. You need to be careful that you consistently pre/post multiply matrices.

Question 3 (BLR Predictive). For a Bayesian Linear Regression model:

1. Find $p(\mathbf{y}^*, \mathbf{y})$. $y^* = \text{unseen}$
2. Find $p(\mathbf{y}^*|\mathbf{y})$.

Question 1 (Vector Ordering in Gaussians). Consider a joint Gaussian density on a vector \mathbf{z} that can be split up as $\mathbf{z} = [\mathbf{x}^\top, \mathbf{y}^\top]^\top$ (this notation denotes stacking):

$$p(\mathbf{z}) = \mathcal{N}\left(\underbrace{\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}}_{\mathbf{z}}; \underbrace{\begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}}_{\mathbf{z}}, \underbrace{\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{C} \end{bmatrix}}_{\mathbf{z}}\right). \quad (11)$$

Now consider the permuted vector $\mathbf{z}' = [\mathbf{y}^\top, \mathbf{x}^\top]^\top$. Show that the Gaussian distribution of $p(\mathbf{z}')$ has a mean with rows, and a covariance matrix with swapped rows and columns.

Hint: Notice that \mathbf{z}' is a permuted version of \mathbf{z} . We can write this mathematically as $\mathbf{z}' = \mathbf{P}\mathbf{z}$, where \mathbf{P} is a permutation matrix:

$$\mathbf{P} = \begin{bmatrix} 0 & \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix}. \quad (12)$$

Note: This is an important skill! Notice that the Gaussian conditioning formula in the formula sheet is only written in one way.

$$\begin{array}{c} \mathbb{R}^E \\ \downarrow \\ \mathbf{z}: \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \end{array} \quad \mathbf{x} \in \mathbb{R}^D, \mathbf{y} \in \mathbb{R}^E \quad \mathbf{z}: \begin{pmatrix} -\mathbf{x} - \\ -\mathbf{y} - \end{pmatrix} \quad \mathbf{z}': \begin{pmatrix} \mathbf{y}' \\ \mathbf{x}' \end{pmatrix} : \begin{pmatrix} \mathbf{F} \mathbf{y} - \\ \mathbf{G} \mathbf{x} - \end{pmatrix}$$

$$p(\mathbf{z}): \underbrace{\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}}_{\mathbb{R}^{D+E}}; \underbrace{\begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}}_{\mathbb{R}^{D+E}}; \underbrace{\begin{pmatrix} \tilde{\mathbf{A}} & \mathbf{B} \\ \mathbf{B}^\top & \tilde{\mathbf{C}} \end{pmatrix}}_{\in \mathbb{R}^{(D+E) \times (D+E)}}$$

$$\text{Show } p(\mathbf{z}') = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ \mathbf{x} \end{bmatrix}; \begin{bmatrix} \mathbf{b} \\ \mathbf{a} \end{bmatrix}, \begin{bmatrix} \mathbf{C} & \mathbf{B}^\top \\ \mathbf{B} & \mathbf{A} \end{bmatrix}\right)$$

$$\mathbf{z}' = \mathbf{P}\mathbf{z}: \begin{bmatrix} 0 & \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{y} \\ \mathbf{x} \end{bmatrix}$$

$$\text{If } \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \rightarrow \mathbf{A}\mathbf{x} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top) \quad \text{From Q1 sheet 6}$$

$$\text{Here } \mathbf{A} = \mathbf{P} \therefore \boldsymbol{\mu} \rightarrow \mathbf{P}\boldsymbol{\mu} = \begin{bmatrix} 0 & \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix} = \begin{bmatrix} \mathbf{b} \\ \mathbf{a} \end{bmatrix}$$

$$\boldsymbol{\Sigma} \rightarrow \mathbf{P}\boldsymbol{\Sigma}\mathbf{P}^\top = \begin{bmatrix} 0 & \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{C} \end{bmatrix} \begin{bmatrix} 0 & \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix} = \begin{bmatrix} \mathbf{B}^\top & \mathbf{C} \\ \mathbf{A} & \mathbf{B} \end{bmatrix} \begin{bmatrix} 0 & \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix} = \begin{bmatrix} \mathbf{C} & \mathbf{B}^\top \\ \mathbf{B} & \mathbf{A} \end{bmatrix}$$

Question 2 (Woodbury Identity). We saw that the two different ways of deriving the Gaussian posterior, gave two different results:

$$p(\boldsymbol{\theta}|\mathbf{y}) = \mathcal{N}\left(\boldsymbol{\theta}; \Phi(X)^\top [\Phi(X)\Phi(X)^\top + \sigma^2 \mathbf{I}_N]^{-1} \mathbf{y}, \quad \mathbf{I}_M - \Phi(X)^\top [\Phi(X)\Phi(X)^\top + \sigma^2 \mathbf{I}_N]^{-1} \Phi(X)\right),$$

$$p(\boldsymbol{\theta}|\mathbf{y}) = \mathcal{N}\left(\boldsymbol{\theta}; \left[\frac{1}{\sigma^2} \Phi(X)^\top \Phi(X) + \mathbf{I}_M\right]^{-1} \frac{1}{\sigma^2} \Phi(X)^\top \mathbf{y}, \quad \left[\frac{1}{\sigma^2} \Phi(X)^\top \Phi(X) + \mathbf{I}_M\right]^{-1}\right). \quad (13)$$

Apply the Woodbury identity to $\left[\frac{1}{\sigma^2} \Phi(X)^\top \Phi(X) + \mathbf{I}_M\right]^{-1}$ to show that the two solutions are equal.

Note: This is an excellent exercise for your matrix algebra skills. The main difficulty that most uninitiated have, is that multiplication no longer commutes. You need to be careful that you consistently pre/post multiply matrices.

Woodbury Identity

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}$$

$$\begin{pmatrix} \mathbf{I}_M + \frac{1}{\sigma^2} \Phi(X)^\top \Phi(X) \\ A \quad U \quad C \quad V \end{pmatrix}^{-1}$$

$$\begin{pmatrix} \mathbf{I}_M + \frac{1}{\sigma^2} \Phi(X)^\top \Phi(X) \\ A \quad U \quad C \quad V \end{pmatrix}^{-1} \quad A = \mathbf{I}_M \quad C = \frac{1}{\sigma^2} \mathbf{I}_N \rightarrow C^{-1} = \sigma^2 \mathbf{I}_N$$

$$\mathbf{I}_M - \mathbf{I}_M \Phi(X)^\top (\Phi(X) \mathbf{I}_M \Phi(X)^\top + \sigma^2 \mathbf{I}_N)^{-1} \Phi(X) \mathbf{I}_M$$

$$= \mathbf{I}_M - \Phi(X)^\top (\Phi(X) \Phi(X)^\top + \sigma^2 \mathbf{I}_N)^{-1} \Phi(X)$$

$$(\mathbf{I}_M + \frac{1}{\sigma^2} \Phi(X)^\top \Phi(X))^{-1} \frac{1}{\sigma^2} \Phi(X)^\top \mathbf{y}$$

$$\left(\mathbf{I}_M - \Phi(X)^\top (\Phi(X) \Phi(X)^\top + \sigma^2 \mathbf{I}_N)^{-1} \Phi(X) \right) \frac{1}{\sigma^2} \Phi(X)^\top \mathbf{y}$$

$$= \frac{1}{\sigma^2} \Phi(X)^\top \mathbf{y} - \frac{1}{\sigma^2} \Phi(X)^\top \mathbf{Z}^{-1} \Phi(X) \Phi(X)^\top \mathbf{y}$$

$$= \frac{1}{\sigma^2} \Phi(X)^\top (\mathbf{y} - \mathbf{Z}^{-1} \Phi(X) \Phi(X)^\top \mathbf{y})$$

$$= \frac{1}{\sigma^2} \Phi(X)^\top (\mathbf{I} - \mathbf{Z}^{-1} \Phi(X) \Phi(X)^\top) \mathbf{y}$$

$$= \frac{1}{\sigma^2} \Phi(X)^\top \mathbf{Z}^{-1} (\mathbf{Z} - \Phi(X) \Phi(X)^\top) \mathbf{y}$$

$$= \frac{1}{\sigma^2} \Phi(X)^\top \mathbf{Z}^{-1} (\Phi(X) \Phi(X)^\top + \sigma^2 \mathbf{I}_N - \Phi(X) \Phi(X)^\top) \mathbf{y}$$

$$= \frac{1}{\sigma^2} \Phi(X)^\top \mathbf{Z}^{-1} \sigma^2 \mathbf{I}_N \mathbf{y}$$

$$= \Phi(X)^\top \mathbf{Z}^{-1} \mathbf{y}$$

$$= \Phi(X)^\top (\Phi(X) \Phi(X)^\top + \sigma^2 \mathbf{I}_N)^{-1} \mathbf{y}$$

∩

Question 3 (BLR Predictive). For a Bayesian Linear Regression model:

1. Find $p(\mathbf{y}^*, \mathbf{y})$. $y^* = \text{unseen}$

2. Find $p(\mathbf{y}^* | \mathbf{y})$.

$$p(y(x^*), D) = \int p(y(x^*), \theta) p(\theta | D) d\theta$$

1. $p(y^*, y) \propto p(y | y^*) p(y^*)$

$$y = \phi(\theta) + \varepsilon \quad \varepsilon \sim N(0, \sigma^2 I)$$

$$y \sim N(\kappa\theta, \sigma^2)$$

$$y^* = \phi(x^*)\theta + \varepsilon$$