

Tutorial Linear Regression

Suggested Reading:

Murphy, chapter 7

OR

Bishop, chapter 3

1 Linear regression

Problem 1: Show that the matrix

$$\Phi(\Phi^T \Phi)^{-1} \Phi^T$$

takes any vector and projects it onto the space spanned by the columns of Φ . Use this result to show that the least square solution for linear regression corresponds to an orthogonal projection of the vector \mathbf{T} (denoted by \mathbf{Z} in class!) onto the manifold S as shown in Figure 1. There, the subspace S is spanned by the basis functions $\phi_j(\mathbf{x})$ in which each basis function is viewed as a vector φ_j of length N with elements $\phi_j(\mathbf{x}_n)$. (Hint: You might want consider what $\Phi(\Phi^T \Phi)^{-1} \Phi^T$ resembles, e.g. how does it relate to the maximum likelihood solution for linear regression.)

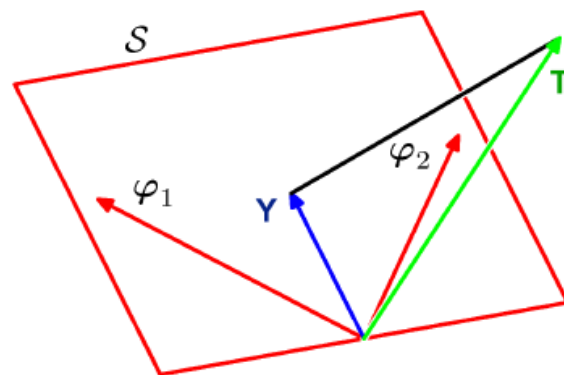


Figure 1: The projection property of $\Phi(\Phi^T \Phi)^{-1} \Phi^T$.

If we set $\Theta = I$ in the section on weighted linear regression, we see that for the standard maximum likelihood problem we want to solve $\Phi^T \Phi \mathbf{w} = \Phi^T \mathbf{z}$ which is equivalent to $\Phi^T (\Phi \mathbf{w} - \mathbf{z}) = 0$. For the solution \mathbf{w}_{MLE} one can deduce from the last equation that (i) $\Phi \mathbf{w}_{ML}$ is a vector in S and (ii) $\Phi \mathbf{w}_{ML} - \mathbf{z}$ (the *error* vector between \mathbf{z} and $\Phi \mathbf{w}_{ML}$) is orthogonal to S . Thus $\Phi \mathbf{w}_{MLE}$ is the orthogonal projection of \mathbf{T} on S . is Φ , thus the whole matrix resembles a projection into a space spanned by the columns of Φ . To show that the least square solution corresponds to an *orthogonal* projection consider the dot product of the projection and the difference vector between \mathbf{t} and the projection, that is

$$(\mathbf{t} - \Phi(\Phi^T \Phi)^{-1} \Phi^T \mathbf{t})^T (\Phi(\Phi^T \Phi)^{-1} \Phi^T \mathbf{t})$$

which is

$$\mathbf{t}^T \Phi(\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} - \mathbf{t}^T \Phi(\Phi^T \Phi)^{-1} \Phi^T \Phi(\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} = 0$$

2 Ridge regression

Problem 2: Using singular value decomposition of the design matrix $\Phi = \mathbf{U}\mathbf{D}\mathbf{V}^T$ show that the output on the training set fitted with the ridge regression solution $\hat{\mathbf{w}}^{ridge}$ can be written as

$$\sum_j \left(\frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j \mathbf{u}_j^T \right) \mathbf{z}$$

where \mathbf{u}_j are the columns of \mathbf{U} , d_j the elements of \mathbf{D} and λ the cost factor of the ℓ_2 regularization. What is the interpretation of this formula?

Based on the SVD of Φ , we can write (the trick is here to rewrite λI as $\lambda V V^T$ and factor matrices out (and remembering that $AB^{-1} = B^{-1}A^{-1}$)):

$$\hat{\mathbf{w}}^{ridge} = V(D^2 + \lambda I)^{-1} D U^T \mathbf{z}$$

Then

$$\Phi \hat{\mathbf{w}}^{ridge} = \Phi(\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{z} = U D (D^2 + \lambda I)^{-1} D U^T \mathbf{z}$$

First observation, $D^2 + \lambda I$ is a diagonal matrix, with $d_j^2 + \lambda$ on the diagonal, d_j are the singular values of Φ . Therefore, its inverse is again a diagonal, with $1/d_j^2 + \lambda$ on the diagonal. And therefore $D(D^2 + \lambda I)^{-1} D$ is also a diagonal matrix (product of diagonal matrices), with $\frac{d_j^2}{d_j^2 + \lambda}$ on the diagonal. This matrix gets multiplied from the right to U , i.e. it only scales the columns. Finally the product of two matrices of the form AB^T can be written as the sum of the outer product of the respective columns of A and B .

Now let's put the formula in words. First, $\mathbf{u}_j^T \mathbf{y}$ computes the representation of \mathbf{y} with respect to the orthonormal basis U , and then reconstructs \mathbf{y} in this basis, however with the coordinates *shrunk* ($\lambda > 0$ and thus $\frac{d_j^2}{d_j^2 + \lambda} < 1$). A greater amount of shrinkage is applied to the coordinates with smaller singular values. (What does a small singular value mean? We will later discuss that the SVD of Φ is another way of expressing the *principal components* of the variables in Φ . These are directions in the space spanned by the training examples in which the training data varies, small singular values are directions in which the training data varies very little. Hence, ridge regression shrinks those directions most. The implicit assumption (or justification for this behaviour) is that the output will vary most with those directions that vary most.)

3 Multi-output linear regression

Problem 3: In class, we only considered functions of the form $f : \mathbb{R}^n \rightarrow \mathbb{R}$. What about the general case of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$? For linear regression with multiple outputs, write down the loglikelihood formulation and derive the MLE of the parameters.

The observation \mathbf{z}_i is a vector with $\mathbf{z}_i \sim \mathcal{N}(\mathbf{W}\mathbf{x}_i, \Sigma)$, $\mathbf{W} \in R^{m \times n}$, $\Sigma \in R^{m \times m}$, covariance Σ is known and fixed for all possible observations. For n i.i.d observed pairs $(\mathbf{x}_i, \mathbf{z}_i)$, the likelihood is $\prod_i \mathcal{N}(\mathbf{W}\mathbf{x}_i, \Sigma)$, and thus the negative log-likelihood is $\text{const} + \frac{1}{2} \sum_i (\mathbf{z}_i - \mathbf{W}\mathbf{x}_i)^T \Sigma^{-1} (\mathbf{z}_i - \mathbf{W}\mathbf{x}_i)$. Let $\mathbf{L}\mathbf{L}^T = \Sigma$ (the cholesky decomposition of Σ), so $\Sigma^{-1} = \mathbf{L}^{-1}\mathbf{L}^{-T}$. Using this decomposition write $\mathbf{L}^{-T}(\mathbf{z}_i - \mathbf{W}\mathbf{x}_i) = (\mathbf{L}^{-T}\mathbf{z}_i - \mathbf{L}^{-T}\mathbf{W}\mathbf{x}_i) = (\hat{\mathbf{z}}_i - \hat{\mathbf{W}}\mathbf{x}_i)$. Using this transformation, the negative log-likelihood now becomes $\text{const} + \frac{1}{2} \sum_i (\hat{\mathbf{z}}_i - \hat{\mathbf{W}}\mathbf{x}_i)^T (\hat{\mathbf{z}}_i - \hat{\mathbf{W}}\mathbf{x}_i)$. Using similar reasoning to the lecture (or Problem 2 above), we can rewrite this as $(\Phi\hat{\mathbf{W}} - \hat{\mathbf{Z}})^T (\Phi\hat{\mathbf{W}} - \hat{\mathbf{Z}})$. Note that $\hat{\mathbf{Z}}$ is a matrix that has the vectors $\hat{\mathbf{z}}_i$ as its rows. Matrix calculus (derivative of the negative log-likelihood with respect to $\hat{\mathbf{W}}$) then gives us $\hat{\mathbf{W}}_{MLE} = (\Phi^T\Phi)^{-1}\Phi^T\hat{\mathbf{Z}}$. So these are m single least square problems for every *column* of $\hat{\mathbf{Z}}$. Finally, transforming back $\hat{\mathbf{W}}_{MLE}$ gives $\mathbf{W}_{MLE} = \mathbf{L}^T\hat{\mathbf{W}}_{MLE}$.

4 Bayesian Linear Regression

Problem 4: We have seen that, as the size of a data set increases, the uncertainty associated with the posterior distribution over model parameters decreases. Prove the following matrix identity

$$(M + vv^T)^{-1} = M^{-1} - \frac{(M^{-1}v)(v^T M^{-1})}{1 + v^T M^{-1}v}$$

and, using it, show that the uncertainty $\sigma_N^2(x)$ associated with the bayesian linear regression function given by

$$\sigma_N^2(x) = \frac{1}{\beta} + \phi(x)^T S_N \phi(x) \tag{1}$$

(where S_N is the covariance of the posterior $p(w|\mathcal{D})$. (see e.g. equations 3.53 and 3.54 in Bishop)) satisfies

$$\sigma_{N+1}^2(x) \leq \sigma_N^2(x) \tag{2}$$

You may want to use

$$\Phi_{N+1}^T \Phi_{N+1} = \Phi_N^T \Phi_N + \phi(x_{N+1})\phi(x_{N+1})^T$$

Starting from

$$(M + vv^T)^{-1} = M^{-1} - \frac{(M^{-1}v)(v^T M^{-1})}{1 + v^T M^{-1}v}$$

multiply

$$M + vv^T$$

on both sides, giving

$$I = I + M^{-1}vv^T - \left(\frac{M^{-1}vv^T + M^{-1}vv^T M^{-1}vv^T}{1 + v^T M^{-1}v} \right)$$

$v^T M^{-1}v$ is a scalar:

$$M^{-1}vv^T + M^{-1}vv^T M^{-1}vv^T = (1 + v^T M^{-1}v)M^{-1}vv^T$$

To prove eq. (1), we first need the definition of σ_N^2 again and also that of S_N (Φ_N denotes the design matrix containing the first N observations).

$$\begin{aligned} \sigma_N^2(\mathbf{x}) &= \frac{1}{\beta} + \phi(\mathbf{x})^T S_N \phi(\mathbf{x}) \\ S_N^{-1} &= \alpha I + \beta \Phi_N^T \Phi_N \end{aligned}$$

We now use the following identity:

$$\Phi_N^T \Phi_N = \sum_{k=1}^N \phi(\mathbf{x}_k) \phi(\mathbf{x}_k)^T$$

so we can write

$$\Phi_{N+1}^T \Phi_{N+1} = \Phi_N^T \Phi_N + \phi(\mathbf{x}_{N+1}) \phi(\mathbf{x}_{N+1})^T$$

and hence

$$S_{N+1}^{-1} = \alpha I + \beta \Phi_{N+1}^T \Phi_{N+1} = \alpha I + \Phi_N^T \Phi_N + \phi(\mathbf{x}_{N+1}) \phi(\mathbf{x}_{N+1})^T = S_N^{-1} + \phi(\mathbf{x}_{N+1}) \phi(\mathbf{x}_{N+1})^T$$

Substituting this into $\sigma_{N+1}^2(\mathbf{x})$ and using the above matrix identity we get:

$$\begin{aligned}
\sigma_{N+1}^2(\mathbf{x}) &= \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_{N+1} \phi(\mathbf{x}) \\
&= \frac{1}{\beta} + \phi(\mathbf{x})^T (\mathbf{S}_N^{-1} + \phi(\mathbf{x}_{N+1})\phi(\mathbf{x}_{N+1})^T)^{-1} \phi(\mathbf{x}) \\
&= \frac{1}{\beta} + \phi(\mathbf{x})^T \left(\mathbf{S}_N - \frac{\mathbf{S}_N \phi(\mathbf{x}_{N+1})\phi(\mathbf{x}_{N+1})^T \mathbf{S}_N}{1 + \phi(\mathbf{x}_{N+1})^T \mathbf{S}_N \phi(\mathbf{x}_{N+1})} \right) \phi(\mathbf{x}) \\
&= \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}) - \frac{\phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}_{N+1})\phi(\mathbf{x}_{N+1})^T \mathbf{S}_N \phi(\mathbf{x})}{1 + \phi(\mathbf{x}_{N+1})^T \mathbf{S}_N \phi(\mathbf{x}_{N+1})}
\end{aligned}$$

The nominator of the last fraction $(\phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}_{N+1}))^2 \geq 0$. Also (with $\alpha, \beta \geq 0$)

$$\phi(\mathbf{x}_{N+1})^T \mathbf{S}_N \phi(\mathbf{x}_{N+1}) = \phi(\mathbf{x}_{N+1})^T \alpha \mathbf{I} \phi(\mathbf{x}_{N+1}) + \phi(\mathbf{x}_{N+1})^T \beta \mathbf{\Phi}_N^T \mathbf{\Phi}_N \phi(\mathbf{x}_{N+1}) \geq 0$$

so we have proven eq. (1).

Problem 5: We know that the posterior for a linear regression algorithm with a likelihood defined by $p(Z \mid W, \beta) = \prod_{n=1}^N \mathcal{N}(Z_n \mid W^T \Phi(X_n), \beta^{-1})$ and prior given by $p(w) = \mathcal{N}(W \mid M_0, \mathbf{S}_0)$ is

$$p(W \mid Z) = \mathcal{N}(W \mid M_N, \mathbf{S}_N)$$

where

$$\begin{aligned} M_N &= \mathbf{S}_N (S_0^{-1} M_0 + \beta \Phi^T Z) \\ S_N^{-1} &= S_0^{-1} + \beta \Phi^T \Phi \end{aligned}$$

Let's assume β is a known constant. Verify that this is the form of the posterior we would derive.

From Bayes' theorem we have

$$p(W \mid Z) \propto p(Z \mid W)p(W),$$

where the factors on the r.h.s. are given by the equations listed above. Writing this out in full, we get

$$\begin{aligned} p(W \mid Z) &\propto \left[\prod_{n=1}^N \mathcal{N}(Z_n \mid W^T \Phi(X_n), \beta^{-1}) \right] \mathcal{N}(w \mid M_0, \mathbf{S}_0) \\ &\propto \exp\left(\frac{-\beta}{2}(Z - \Phi W)^T(Z - \Phi W)\right) \exp\left(\frac{-1}{2}(W - M_0)^T \mathbf{S}_0^{-1}(W - M_0)\right) \\ &= \exp\left(\frac{-1}{2}(W^T(\mathbf{S}_0^{-1} + \beta \Phi \Phi^T)W - \beta W^T \Phi^T Z + \beta Z^T Z + M_0^T \mathbf{S}_0^{-1}W - W^T \mathbf{S}_0^{-1}M_0 + M_0^T \mathbf{S}_0^{-1}M_0)\right) \\ &= \exp\left(\frac{-1}{2}(W^T(\mathbf{S}_0^{-1} + \beta \Phi^T \Phi)W - (\mathbf{S}_0^{-1}M_0 + \beta \Phi^T Z)^T W - W^T(\mathbf{S}_0^{-1}M_0 + \beta \Phi^T Z) + \beta Z^T Z + M_0^T \mathbf{S}_0^{-1}M_0)\right) \\ &= \exp\left(\frac{-1}{2}(w - M_N)^T \mathbf{S}_N^{-1}(W - M_N)\right) \exp\left(\frac{-1}{2}(\beta Z^T Z + M_0^T \mathbf{S}_0^{-1}M_0 - M_N^T \mathbf{S}_N^{-1}M_N)\right) \end{aligned}$$

where we used the definitions for M_N and \mathbf{S}_N^{-1} from above when completing the square in the last step. The first exponential corresponds to the posterior, unnormalized Gaussian distribution over W , while the second exponential is independent of W and hence can be absorbed into the normalization factor.

5 Online Learning

Problem 6: Suppose we are using a linear basis function model where the posterior distribution is given by $p(W \mid Z) = \mathcal{N}(W \mid M_N, \mathbf{S}_N)$ and we have already observed N data points. That means that this posterior can be regarded as the prior for the next observation. By considering an additional data point (X_{N+1}, z_{N+1}) , and by completing the square in the exponential, show that the resulting posterior distribution is again given by the posterior mentioned above, but with \mathbf{S}_N replaced by \mathbf{S}_{N+1} and M_N replaced by M_{N+1} .

Combining the prior

$$p(W) = \mathcal{N}(W \mid M_N, \mathbf{S}_N)$$

and the likelihood

$$p(Z_{N+1} \mid X_{N+1}, W) = \left(\frac{\beta}{2\pi} \right)^{\frac{1}{2}} \exp \left(\frac{-\beta}{2} (Z_{N+1} - W^T \Phi_{N+1})^2 \right)$$

where $\Phi_{N+1} = \Phi(X_{N+1})$, we obtain a posterior of the form

$$p(W \mid Z_{N+1}, X_{N+1}, M_N, \mathbf{S}_N) \propto \exp \left(-\frac{1}{2} (W - M_N)^T \mathbf{S}_N^{-1} (W - M_N) - \frac{1}{2} \beta (Z_{N+1} - W^T \Phi_{N+1})^2 \right).$$

We can expand the argument of the exponential, omitting the $-1/2$ factors, as follows

$$\begin{aligned} & (W - M_N)^T \mathbf{S}_N^{-1} (W - M_N) + \beta (Z_{N+1} - W^T \Phi_{N+1})^2 \\ &= W^T \mathbf{S}_N^{-1} W - 2W^T \mathbf{S}_N^{-1} M_N + \beta W^T \Phi_{N+1}^T \Phi_{N+1} W - 2\beta W^T \Phi_{N+1} Z_{N+1} + \text{const} \\ &= W^T (\mathbf{S}_N^{-1} + \beta \Phi_{N+1} \Phi_{N+1}^T) W - 2W^T (\mathbf{S}_N^{-1} M_N + \beta \Phi_{N+1} Z_{N+1}) + \text{const}, \end{aligned}$$

where const denotes remaining terms independent of W . From this we can read off the desired result directly,

$$p(W \mid Z_{N+1}, X_{N+1}, M_N, \mathbf{S}_N) = \mathcal{N}(W \mid M_{N+1}, \mathbf{S}_{N+1}),$$

with

$$\mathbf{S}_{N+1}^{-1} = \mathbf{S}_N^{-1} + \beta \Phi_{N+1} \Phi_{N+1}^T$$

and

$$M_{N+1} = \mathbf{S}_{N+1} (\mathbf{S}_N^{-1} M_N + \beta \Phi_{N+1} Z_{N+1}).$$