

# Machine Learning: Exercises for Block II (Linear Regression)

Isabel Valera

## Exercise 1: Orthogonal

- i) Show that the matrix  $\Phi(\Phi^\top\Phi)^{-1}\Phi^\top$  takes any vector  $v$  and projects it onto the space spanned by the columns of  $\Phi$ .
- ii) Use this result to show that the least-squares solution given in equation 2 corresponds to an orthogonal projection of the vector  $t$  onto the manifold  $S$  as shown in figure 1.

$$\Phi(\Phi^\top\Phi)^{-1}\Phi^\top \quad (1)$$

$$w_{ML} = (\Phi^\top\Phi)^{-1}\Phi^\top t \quad (2)$$

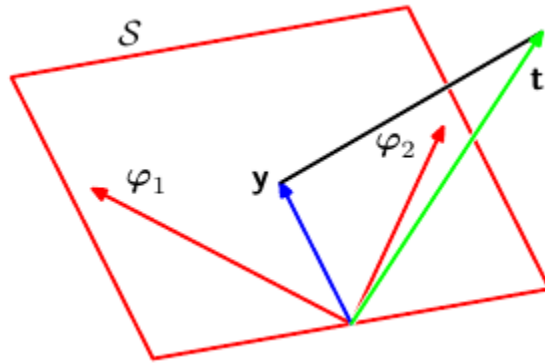


Figure 1

## Exercise 2: Sample weights

Consider a data set in which each data point  $y_n$  is associated with a weighting factor  $r_n > 0$ , so that the sum-of-squares loss function is as in equation 3.

- i) Find an expression for the solution  $w^*$  that minimizes this error function.
- ii) Given an interpretation of the weighted sum-of-squares error function in terms of data dependent noise variance.

- iii) Given an interpretation of the weighted sum-of-squares error function in terms of replicated data points.

$$L(\mathbf{X}, \mathbf{y}, \mathbf{w}) = \frac{1}{2} \sum_{n=1}^N r_n (y_n - \mathbf{w}^\top \varphi(\mathbf{x}_n))^2 \quad (3)$$

### Exercise 3: Independent noise

Consider the linear model given in equation 4 together with the sum-of-squares loss function given in 5. Now suppose that Gaussian noise  $\varepsilon_i$  with zero mean and variance  $\sigma^2$  is added independently to each of the input variables  $x_i$ . By making use of  $\mathbb{E}[\varepsilon_i] = 0$  and  $\mathbb{E}[\varepsilon_i \varepsilon_j] = \delta_{ij} \sigma^2$  show that minimizing  $L(\mathbf{X}, \mathbf{y}, \mathbf{w})$  averaged over the noise distribution is equivalent to minimizing the sum of squares loss for noise-free input variables with the addition of a weight-decay regularization term, in which the bias parameter  $w_0$  is omitted from the regularizer.

$$f(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i \quad (4)$$

$$L(\mathbf{X}, \mathbf{y}, \mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (f(\mathbf{x}_n, \mathbf{w}) - y_n)^2 \quad (5)$$

### Exercise 4: Linear basis functions

Consider a linear basis function regression model for a multivariate target variable  $\mathbf{y}$  having a Gaussian distribution as given in 6 where  $f(\mathbf{x}, \mathbf{W}) = \mathbf{W}^\top \varphi(\mathbf{x})$  together with a training data set comprising input basis vectors  $\varphi(\mathbf{x}_n)$  and corresponding target vectors  $\mathbf{y}_n$ , with  $n = 1, \dots, N$ .

$$p(\mathbf{y} | \mathbf{W}, \Sigma) = \text{Gauss}(\mathbf{y} | f(\mathbf{x}, \mathbf{W}), \Sigma) \quad (6)$$

- i) Show that the maximum likelihood solution  $\mathbf{W}_{ML}$  for the parameter matrix  $\mathbf{W}$  has the property that each column is given by an expression as in equation 2, which was the solution for an isotropic noise distribution. Note that this is independent of the covariance matrix  $\Sigma$ .
- ii) Show that the maximum likelihood solution for  $\Sigma$  is given by

$$\Sigma = \frac{1}{N} \sum_{n=1}^N (\mathbf{y}_n - \mathbf{W}_{ML}^\top \varphi(\mathbf{x}_n)) (\mathbf{y}_n - \mathbf{W}_{ML}^\top \varphi(\mathbf{x}_n))^\top$$

### Exercise 5: Gauss-Markov theorem

- i) Let  $\theta = \alpha^\top \mathbf{w}$  be a linear combination of the parameters  $\mathbf{w}$ . Prove the Gauss-Markov theorem: the least squares estimate of  $\theta = \alpha^\top \mathbf{w}$  has variance no bigger than that of any other linear unbiased estimate  $\alpha^\top \mathbf{w}$ .
- ii) The matrix inequality  $B \preceq A$  holds if  $A - B$  is positive semidefinite. Show that if  $\hat{\Sigma}$  is the covariance matrix of the least squares estimate of  $\mathbf{w}$  and  $\tilde{\Sigma}$  is the covariance matrix of any other linear unbiased estimate, then  $\hat{\Sigma} \preceq \tilde{\Sigma}$ .

## Exercise 6: Ridge regression

- i) Show that the ridge regression estimate is the mean (and mode) of the posterior distribution, under a Gaussian prior  $\mathbf{w} \sim \text{Gauss}(0, \tau \mathbf{I})$  and Gaussian sampling model  $\mathbf{y} \sim \text{Gauss}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I})$ .
- ii) Find the relationship between the regularization parameter  $\lambda$  in the ridge formula 7, and the variances  $\tau$  and  $\sigma^2$ .

$$\sum_{i=1}^N \left( y_i - w_0 - \sum_{j=1}^p x_{ij} w_j \right)^2 + \lambda \sum_{j=1}^p w_j^2 \quad (7)$$

## Exercise 7: Gaussian parameters

Assume  $y_i = \text{Gauss}(w_0 + \mathbf{x}_i^\top \mathbf{w}, \sigma^2)$ ,  $i = 1, 2, \dots, N$ , and the parameters  $w_j, j = 1, \dots, p$  are each distributed as  $\text{Gauss}(0, \tau^2)$ , independently of one another. Assuming  $\sigma^2$  and  $\tau^2$  are known, and  $w_0$  is not governed by a prior (or has flat improper prior), show that the (minus) log-posterior density of  $\mathbf{w}$  is proportional to  $\sum_{i=1}^N \left( y_i - w_0 - \sum_j x_{ij} w_j \right)^2 + \lambda \sum_{j=1}^p w_j^2$  where  $\lambda = \frac{\sigma^2}{\tau^2}$

## Exercise 8: Artificial data

Show that the ridge regression estimates can be obtained by ordinary least squares regression on an augmented data set. We augment the centered matrix  $\mathbf{X}$  with  $p$  additional rows  $\sqrt{\lambda} \mathbf{I}$ , and augment  $\mathbf{y}$  with  $p$  zeros. By introducing artificial data having response value zero, the fitting procedure is forced to shrink the coefficients towards zero. This is related to the idea of *hints* due to Abu-Mostafa (1995), where model constraints are implemented by adding artificial data examples that satisfy them.

## References

- [1] C. M. Bishop. *Pattern recognition and machine learning*. springer, 2006.
- [2] J. Friedman, T. Hastie, R. Tibshirani, et al. *The elements of statistical learning*, volume 1. Springer series in statistics New York, 2001.