

## Machine Learning Worksheet 5

### Linear Regression

## 1 Probability Theory

**Problem 1.** Let  $X$  have a continuous cdf  $F_X(x)$ . Define the random variable  $Y$  as  $Y = F_X(X)$ . Assuming that  $F_X(x)$  is strictly increasing, how is  $Y$  distributed? Show your work.

**Problem 2.** Show that the sum of two independent Gaussian random variables ( $\mathbf{X}_1$  and  $\mathbf{X}_2$ ) is Gaussian. Some of the properties of Gaussians mentioned in the lecture can help.

**Problem 3.** Let  $Z = (X, Y)$  be a bivariate normal distributed random variable. Furthermore, let  $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$  and  $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$ . Assume that  $\rho(X, Y) = 0$ . Show that in this case  $X$  and  $Y$  are independent.

## 2 Weighted Linear Regression

Consider a linear regression problem in which we want to “weight” different training examples differently. Specifically, suppose we want to minimize

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \theta_n (z_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2$$

**Problem 4.** We already worked out what happens for the case where all the weights  $\theta_n$  are the same. In this problem, we will generalize some of those ideas to the weighted setting, and also implement the locally weighted linear regression algorithm.

1. Show that  $E(\mathbf{w})$  can also be written

$$E(\mathbf{w}) = (\mathbf{z} - \Phi \mathbf{w})^T \Theta (\mathbf{z} - \Phi \mathbf{w}) \tag{1}$$

for an appropriate diagonal matrix  $\Theta$ , and where  $\Phi$  and  $\mathbf{z}$  are as defined in class. State clearly what  $\Theta$  is.

2. Now let all the  $\theta_n$  equal 1. By differentiating Eq. 1 with respect to  $\mathbf{w}$ , derive the normal equations for the least squares problem, as given in class.
3. Generalize the normal equations to the case of arbitrary  $\theta_n$ .
4. Suppose we have a training set  $(\mathbf{x}_n, z_n)$ ;  $n = 1, \dots, N$  of  $N$  independent examples, but in which the  $z_n$  were observed with differing variances. Specifically, suppose that

$$p(z_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}(z_n | \mathbf{w}^T \Phi(\mathbf{x}_n), \sigma_n^2)$$

where the  $\sigma_n$  are fixed, known, constants. Show that finding the maximum likelihood estimate of  $\mathbf{w}$  reduces to solving a weighted linear regression problem. State clearly what the  $\theta_n$  are in terms of the  $\sigma_n$ .

5. With *ordinary* linear regression it may be a good idea to *rescale* the *columns* of the design matrix – in particular when using nonlinear basis function expansions (e.g. like polynomial expansion). Using the normal equations, prove that rescaling the design matrix does not change the predicted values for some test dataset.

### 3 Basisfunctions

**Problem 5.** Show that the tanh function and the logistic sigmoid function are related by

$$\tanh(x) = 2\sigma(2x) - 1$$

Thus, show that a general linear combination of logistic sigmoid functions of the form

$$y(x, \mathbf{w}) = w_0 + \sum_{j=1}^M w_j \sigma\left(\frac{x - \mu_j}{s}\right)$$

is equivalent to a linear combination of tanh functions of the form

$$y(x, \mathbf{u}) = u_0 + \sum_{j=1}^M u_j \tanh\left(\frac{x - \mu_j}{2s}\right)$$

and find expressions to relate the new parameters  $\{u_0, \dots, u_M\}$  to the original parameters  $\{w_0, \dots, w_M\}$ .

### 4 Ridge regression

**Problem 6.** Show that the following holds: The ridge regression estimates can be obtained by ordinary least squares regression on an augmented dataset: Augment the design matrix  $\Phi$  with  $p$  additional rows  $\sqrt{\lambda}\mathbf{I}$  and augment  $\mathbf{z}$  with  $p$  zeros.

**Problem 7.** 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

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

where  $\mathbf{u}_j$  are the columns of  $\mathbf{U}$ ,  $d_j$  the elements of  $\mathbf{D}$  and  $\lambda$  the cost factor of the  $\ell_2$  regularization. What is the interpretation of this formula?

## 5 Multi-output linear regression

**Problem 8.** 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.

## 6 Bayesian Linear Regression

**Problem 9.** ★ We have seen that, as the size of a data set increases, the uncertainty associated with the posterior distribution over model parameters decreases (see tower equalities). 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  $\sigma_N^2(\mathbf{x})$  associated with the bayesian linear regression function given by eq. (26) on the slides satisfies

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