

c) The aim of this exercise is to deepen your understanding of regression and relate it to other statistical measures such as covariance and correlation.

Consider two data series,  $X = (x_1, x_2, \dots, x_n)$  and  $Y = (y_1, y_2, \dots, y_n)$ , with  $\bar{X} = 0$  and  $\bar{Y} = 0$ . We use linear regression (ordinary least squares) to regress  $Y$  against  $X$  (without fitting any intercept), as in  $Y = aX + \epsilon$  where  $\epsilon$  denotes a series of error terms.

Derive the value of the regression coefficient  $a$  in terms of the standard deviations  $\sigma_X$  and  $\sigma_Y$  and the correlation  $\rho_{XY}$  between the two data series.

From class we know that  $a = (X^T X)^{-1} X^T y$ . In our one dimensional case this yields

$$a = \left( \sum_{i=1}^n x_i^2 \right)^{-1} \sum_{i=1}^n x_i y_i \quad (1.1)$$

We recognize

$$\sum_{i=1}^n x_i^2 \stackrel{\bar{x}=0}{=} \sum_{i=1}^n (x_i - \bar{x})^2 = N \cdot \text{Var}(X) \quad (1.2)$$

and

$$\sum_{i=1}^n x_i y_i \stackrel{\substack{\bar{x}=0 \\ \bar{y}=0}}{=} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = N \cdot \text{Cov}(X, Y) \quad (1.3)$$

Plugging in the results (1.2) and (1.3) into (1.1) yields

$$a = (N \cdot \text{Var}(X))^{-1} \cdot N \cdot \text{Cov}(X, Y) = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$



d) In this exercise we will show that adding noise to a data point is the same as adding regularization to your mean squared error cost function.

Given a training dataset  $D = \{x_i, y_i\}_{i=1}^n$  where  $x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$  and the mean squared error cost function  $f(x, y; w) = \frac{1}{n} \sum_{i=1}^n (y_i - \langle w, x_i \rangle)^2$ .

Consider adding a noise term  $\epsilon_i$  where  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  to each data sample in the training set  $x_i$  such that  $\mathbb{E}[\epsilon_i] = 0$  and  $\mathbb{E}[\epsilon_i \epsilon_j] = \sigma^2$ . We will assume  $n \rightarrow \infty$ . Thus, we write our cost function  $f(x, y; w) = \frac{1}{n} \sum_{i=1}^n (y_i - \langle w, x_i \rangle)^2$  as  $f(x, y; w) = \mathbb{E}[(y_i - \langle w, x_i \rangle)^2]$ . Prove that the following equation holds:

$$\mathbb{E}[(y_i - \langle w, x_i + \epsilon_i \rangle)^2] = \mathbb{E}[(y_i - \langle w, x_i \rangle)^2] + \sigma^2 \sum_{i=1}^d w_i^2$$

$$\begin{aligned} \mathbb{E}[(y_i - \langle w, x_i + \epsilon_i \rangle)^2] &= \mathbb{E}[(y_i - \langle w, x_i \rangle - \langle w, \epsilon_i \rangle)^2] \quad | \text{Linearity of inner product} \\ &= \mathbb{E}[(y_i - \langle w, x_i \rangle)^2 - 2(y_i - \langle w, x_i \rangle) \langle w, \epsilon_i \rangle + \langle w, \epsilon_i \rangle^2] \\ &= \mathbb{E}[(y_i - \langle w, x_i \rangle)^2] + \mathbb{E}[\langle w, \epsilon_i \rangle^2] \quad | \langle w, \epsilon_i \rangle = \epsilon_i^T w \end{aligned}$$

Now we have to make a small explanatory detour:

$$\langle w, \epsilon_i \rangle^2 = \langle w, \epsilon_i \rangle \cdot \langle w, \epsilon_i \rangle$$

$$= \underbrace{(w^T \epsilon_i)}_{\mathbb{R}} \cdot \underbrace{(w^T \cdot \epsilon_i)}_{\mathbb{R}}$$

$$= w^T \cdot (\epsilon_i \epsilon_i^T) \cdot w$$

$$= w^T \begin{pmatrix} \epsilon_{i1} & \epsilon_{i1} & \dots & \epsilon_{i1} & \epsilon_{id} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \epsilon_{id} & \epsilon_{i1} & \dots & \epsilon_{id} & \epsilon_{id} \end{pmatrix} w \quad | \text{Since matrix multiplication is associative.}$$

$$= \sum_{j=1}^d \sum_{k=1}^d w_j \cdot \epsilon_{ij} \epsilon_{ik} \cdot w_k$$

From  $\mathbb{E}[\epsilon_i \epsilon_i^T] = \sigma^2 I$  we know that  $\mathbb{E}[\epsilon_{ij} \epsilon_{ik}] = \sigma^2$  if  $k=j$  and 0 otherwise. Therefore we get

$$\mathbb{E}[\langle w, \epsilon_i \rangle^2] = \sum_{j=1}^d \sum_{k=1}^d w_j \cdot \mathbb{E}[\epsilon_{ij} \epsilon_{ik}] \cdot w_k = \sigma^2 \cdot \sum_{j=1}^d w_j^2$$



- a) Since we want good generalization (test) performance on data even though we learn only with available (train) data therefore, model capacity must be appropriate to the true complexity of the task. Higher capacity for a simpler task can lead to overfitting and low capacity for a complex task can lead to underfitting. Therefore when looking for an optimal capacity we generally make a Bias and Variance trade-off. More conclusively, the expected test mean square error (MSE), for a given value  $x$ , can always be decomposed into the sum of three fundamental quantities: the variance of  $\hat{f}(x)$ , the squared bias of  $\hat{f}(x)$  and the variance of the error terms  $\varepsilon$ .  
(Reading ISLR: section 2.2.2 The Bias-Variance Trade-off)

$$\mathbb{E}[(y - \hat{f}(x))^2] = \text{Var}(\hat{f}(x)) + [\text{Bias}(\hat{f}(x))]^2 + \text{Var}(\varepsilon).$$

Give explanation of the terms Bias and Variance and how does they relate to Overfitting and Underfitting? (max 150 words)

- **Variance** refers to the amount by which  $\hat{f}$  would change if we estimated it using a different training data set. Since the training data are assumed to be sampled from a probability distribution, different training data sets will result in a different  $\hat{f}$ .
- **Bias** refers to the error that is introduced by approximating a real-life problem, which may be extremely complicated, by a much simpler model. In mathematical terms bias is given as

$$\text{Bias}(\hat{f}(x)) = \mathbb{E}[\hat{f}] - f$$

Where  $f$  is the function we are trying to approximate.

- Under the condition that the overall error remains constant, increasing the variance enables the model to fit more complex functions (i.e increased capacity). Thus, the model is able to fit training observations very well (i.e low MSE). Due to the fact that the model is incentivized by the loss function to minimize its MSE, this can lead to overfitting. Otherwise, increasing the bias reduces the capacity and potentially leads to underfitting.

b) Ridge Regression:

In the slide "Fitting with Regularization" a regularization term is introduced to the least squares loss to avoid over/underfitting. This is also known as ridge regression.

$$J(w) = MSE_{train} + \lambda w^T w$$

Here  $w$  are the weights and  $\lambda$  is the regularization parameter. Provide a closed form solution for  $w$  which minimizes this loss. Please show all steps in the solution.

We have  $MSE_{train} = \frac{1}{m} \|Xw - y\|_2^2$ , thus

$$J(w) = \frac{1}{m} \|Xw - y\|_2^2 + \lambda w^T w$$

We now take the gradient of  $J(w)$

$$\nabla_w J(w) = \frac{1}{m} \cdot 2X^T(Xw - y) + 2\lambda w$$

We then set the gradient to zero

$$\nabla_w J(w) \stackrel{!}{=} 0$$

$$\Leftrightarrow X^T(Xw - y) = m\lambda w$$

$$\Leftrightarrow X^T X w - X^T y = m\lambda w$$

$$\Leftrightarrow (X^T X - m\lambda I)w = X^T y$$

$$\text{Thus } \hat{w} = (X^T X - m\lambda I)^{-1} X^T y$$

Authors:

- Name: Mhd Jawad Al Rahwanji - MatrikelNummer: 7038980 - Email: mhal00002@stud.uni-saarland.de
- Name: Christian Singer - MatrikelNummer: 7039059 - Email: chsi00002@stud.uni-saarland.de

### Exercise 3.1 - Eigenvalue Decomposition

- a) Yes,  $M$  is symmetrical. Eigenvalue decomposition is possible.
- b)  $M$  is singular because  $\det(M) = 0$ . One of its eigenvalues is 0.
- c) Using the characteristic equation,  $\det(A - \lambda I) = \lambda(\lambda(4 - \lambda) - 3) = 0$ . After solving it we get  $\lambda_1 = 0$ ,  $\lambda_2 = 1$  and  $\lambda_3 = 3$ . The corresponding eigenvectors would be:  $v_1 = k(1, 1, 1)^T$ ,  $v_2 = (0, 0, 0)^T$  and  $v_3 = k(1, \frac{-1}{2}, 1)^T$ .

### Exercise 3.2 - Linear Regression

- a) In the case of predicting illness, Accuracy is how often the model made a correct prediction. Precision is how many of its positive predictions were ill. Recall (Sensitivity) reports how many of the ill were identified. F1 combines Precision and Recall.

It's a trade-off. Precision is preferred to reduce false alarms. Recall is preferred in sensitive tasks. Accuracy is used with balanced class distributions and the classes are of equal importance. Whereas F1 is used in cases when the class distribution is unbalanced.

- b) Elaborations:
  - a) One unit change in  $X_1$  adds 10 to  $Y$  which may cause the described effect but not at all times.
  - b) The 3rd term has a very small contribution to  $Y$  due to its logarithmic nature and plus the value of the intercept. A 1 unit change wouldn't cause a 50% change in  $Y$ .
  - c) One 100% change in  $X_2$  in the range  $[0, 1]$  subtracts from  $Y$  which may cause the described effect but not at all times.
  - d) A company with \$0 raised, \$1 initial stock value and \$0 debt would be valued at 1 million dollars the year after but that wouldn't make sense.

Authors:

- Name: Mhd Jawad Al Rahwanji - MatrikelNummer: 7038980 - Email: mhal00002@stud.uni-saarland.de
- Name: Christian Singer - MatrikelNummer: 7039059 - Email: chsi00002@stud.uni-saarland.de

### Exercise 3.4 - Validation set and Cross Validation

- a) The advantage cross validation (CV) has over a holdout set is allowing the model to train on multiple train-test splits. That gives us a better idea about the performance by utilizing all of the available data. That is particularly useful when the available dataset is small. Also, with CV we eliminate split dependency resulting in a more realistic performance score for the model. But bear in mind that CV is computationally more expensive.
- b) The process would go as follows:
  - For each hyperparameter (order of the polynomial = 1, 5, 9), we train linear regression using the augmented features with a  $k$ -fold cross validation (CV) ( $k = 5$ ).
  - Each time (for example, order = 1) we train 5 linear regressions each trained on different train-test splits. The 5 models are scored (using MSE for example) and then their scores are averaged to one final score for the current parameter.
  - As for the means of comparison, we sort the final scores achieved by the different hyperparameters (the averaged scores of the 5 linear regressions we trained).
  - We pick the hyperparameter with the best overall final score.
- c) Both runs will score 35. In LOOCV, we train  $n$  models ( $n = \#$  of observations). Each iteration we train on  $n-1$  observations and test on the remaining observation. As a result, no matter how we shuffle the dataset we will end up with the same final score.