**CISPA** HELMHOLTZ CENTER FOR INFORMATION SECURITY

UNIVERSITÄT DES SAARLANDES

**Deadline:** Thursday, December 15, 2022, 16:00

Before solving the exercises, read the instructions on the course website.

- For each theoretical problem, submit a single `pdf` file that contains your answer to the respective problem. This file may be a scan of your (legible) handwriting.

- For each practical problem, submit a single `zip` file that contains

  - the completed jupyter notebook (`.ipynb`) file,
  - any necessary files required to reproduce your results, and
  - a `pdf` report generated from the jupyter notebook that shows all your results.

- For the bonus question, submit a single `zip file` that contains

  - a `pdf` file that includes your answers to the theoretical part,
  - the completed jupyter notebook (`.ipynb`) file for the practical component,
  - any necessary files required to reproduce your results, and
  - a `pdf` report generated from the jupyter notebook that shows your results.

- Every team member has to submit a signed Code of Conduct.

**Problem 1** (T, 14 Points).     **Cross-Validation.**

1. [$4pts$] Explain the impact of the value for $k$ in $k$-fold cross validation. Where does $k$-fold CV fit in between the validation set approach and LOOVC and what is the advantage of using it?

2. [$2pts$] For the hat matrix $H$, defined as

$$H = X(X^T X)^{-1} X^T \ ,$$

the diagonal element $h_i = H_{(ii)}$ is called the leverage. What is the meaning of the leverage $h_i$ for sample $i$? Consequently, what effect on model estimation does removing a sample with high leverage from the dataset have?

3. [$8pts$] Prove that for linear and polynomial least squares regression, the LOOCV estimate for the test MSE can be calculated as

$$\mathrm{CV}_{(n)} = \frac{1}{n} \sum_{i=1}^{n} \left( \frac{y_i - \hat{y}_i}{1 - h_i} \right)^2 \ .$$

*Hint:* First, properly understand and define all variables that are relevant for the situation where you set sample $i$ aside.

**Problem 2** (T, 6 Points).     **The Bootstrap.**
We will now derive the probability that a given observation is part of a bootstrap sample of size $n$.
Suppose that we obtain a bootstrap sample from a set of $n$ observations.

1. [$2pts$] What is the probability that the first bootstrap observation is not the $j$th observation from the original sample? Justify your answer.

2. [$2pts$] Argue that the probability that the $j$th observation is not in the bootstrap sample is $(1 - 1/n)^n$.

3. [$2pts$] Comment on the behavior of the above probabilities with increasing sample size $n$.

**Elements of Machine Learning, WS 2022/2023**
Jilles Vreeken and Aleksandar Bojchevski
Exercise Sheet #3: *Generalization*

C I S P A
HELMHOLTZ CENTER FOR
INFORMATION SECURITY

UNIVERSITÄT
DES
SAARLANDES

---

**Problem 3** (T, 10 Points).    **Correlated Variables.**
It is well-known that ridge regression tends to give similar coefficient values to correlated variables, whereas the Lasso may give quite different coefficient values to correlated variables. We will now explore this property in a very simple setting. Suppose that $n = 2$, $p = 2$, $x_{11} = x_{12}, x_{21} = x_{22}$. Furthermore, suppose that $y_1 + y_2 = 0$, $x_{11} + x_{21} = 0$ and $x_{12} + x_{22} = 0$, so that the estimated intercept in a least squares, ridge regression, or lasso model is zero, $\hat{\beta}_0 = 0$.

1. [*2pts*] Write out the ridge regression optimization problem explicitly in this setting.

2. [*2pts*] What does the solution space for ordinary least squares regression look like here?

3. [*4pts*] Show that in this example, the ridge coefficient estimates satisfy $\hat{\beta}_1 = \hat{\beta}_2$.

4. [*2pts*] Explain how this example connects to the statement that ridge regression tends to give similar coefficient values to correlated variables.

**Problem 4** (P, 20 Points).    **Cross-Validation in the Wild.**
Go through *5.3 Lab: Cross-Validation and the Bootstrap* (ISLR p.212–218) as well as *6.5.2 Lab: Ridge Regression and the Lasso* (ISLR p. 274–278). The objective of this programming exercise is to predict the logarithm of the prostate specific antigen (PSA) level using several predictors. You can find the data *prostate_train.csv, prostate_test.csv* in the CMS.

1. [*1pts*] Load the data. Make sure to split the training, respectively test, dataset into one dataset containing the target column **lpsa**, and another one containing the remaining features.

2. [*4pts*] Use `sklearn.linear_model.Ridge` to fit a ridge regression model on the training data and predict **lpsa** from all other features. Repeat this for different regularization parameters $\lambda$ and plot the resulting coefficient values in relation to $\lambda$ (cf. Figure 6.4, p. 238, ISLR). What can you observe?

3. [*4pts*] Use `sklearn.model_selection.KFold` to perform 10-fold cross-validation on the training data and determine the optimal value for $\lambda$ in the ridge regression model. Report the average train and test MSE for this $\lambda$, i.e., the average of the train and test MSEs on the holdout datasets.

4. [*4pts*] Now use `sklearn.linear_model.Lasso` to fit Lasso models for **lpsa**, and plot the Lasso coefficients in relation to the regularization parameter $\lambda$ (cf. Figure 6.6, p. 242, ISLR). What can you say in comparison to the plot in part 2 of this exercise? Make at least two observations.

5. [*5pts*] Perform 10-fold cross-validation on the training data to determine the optimal value for $\lambda$ in Lasso. Again report average train and test MSEs for this $\lambda$. How many and which features are used? Compare this to the coefficients determined for ridge regression in part 3 of this exercise.

6. [*2pts*] Compare the performance, in terms of MSE, of the best models generated in parts 3. and 5. Which model would you choose and why? What alternative model could have been used?

**Problem 5** (Bonus).    **Coefficient Uncertainty.**
This bonus problem has both theoretical and practical parts.

1. **Theoretical**. Consider the regularized regression problem

$$\underset{\beta}{\text{minimize}} \quad (y - X\beta)^T(y - X\beta) + \lambda|\beta|^q \ , \tag{5.1}$$

   which becomes ridge regression for $q = 2$, respectively Lasso for $q = 1$. While minimizing the above objective gives us a *point estimate* of $\beta$, you have seen in the lecture how a validation approach such as bootstrapping can show us a *distribution* over possible values of $\beta$ (e.g. Fig. 5.10 in ISLR).

   However, bootstrapping can only be done after model estimation. In this exercise, we will investigate what happens if we include a distribution over $\beta$ directly in our model.

   (a) We consider a linear model with independent, Gaussian distributed samples of $Y$, that is, $p(y_i \mid x_i, \beta_i) \sim \mathcal{N}(y_i \mid x_n^T\beta, \sigma^2)$ for $i = 0, .., n$. Write out $p(Y \mid X, \beta)$.

   (b) We now also place a distribution over $\beta$. For now, we assume a Gaussian distribution, where $p(\beta) \sim \mathcal{N}(0, 1/\lambda)$ for some parameter $\lambda$. Consider

$$\hat{\beta} = \underset{\beta}{\max} \ \log\Big(p(\beta) \ p(Y \mid X, \beta)\Big) \ . \tag{5.2}$$

   Simplify the above expression. Where is the maximum attained?

   (c) Compare the expression in (b) with ridge regression, i.e. Eq. (5.1) where $q = 2$. Can you draw a connection between accounting for the coefficient distribution and regularization?

   (d) How would you need to modify (b) to arrive at ordinary least squares regression, Eq. (5.1) with $\lambda = 0$? Explain what this means for the distribution over $\beta$.

   (e) How would you need to modify (b) to arrive at the Lasso regression, Eq. (5.1) with $q = 1$? Explain what this means for the distribution over $\beta$.

2. **Practical**. We will now investigate this problem further on synthetic data.

   (a) Write a Python function computing the expression you found for Eq. (5.2).

   (b) Solve the expression in Eq. (5.2) for $\beta$. *Useful function:* `scipy.optimize`.

   (c) Now, we sample different values for $\beta$ from the Gaussian distribution we assumed. Create a histogram over these values, and also indicate the estimate $\hat{\beta}$ that you found in part (b).

   (d) Repeat steps (b) and (c) for increasing numbers of datapoints, $n = 10, n = 100$ and $n = 1000$. Can you observe a trend?