

# Computational Statistics and Data Analysis (MVComp2)

## Exercise 8

Lecturer Tristan Berau

Semester Wi23/24

Due Dec. 14, 2023, 23:59

### 1 Partial derivative of the residual sum of squares (2 points)

Recall the residual sum of squares,  $\text{RSS} = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$ .

Show that

$$\begin{aligned}\frac{\partial}{\partial w_k} \text{RSS}(\mathbf{w}) &= a_k w_k - c_k \\ a_k &= 2\|x_{:,k}\|^2 \\ c_k &= 2 \sum_{i=1}^n x_{ik} (y_i - \mathbf{w}_{-k}^\top \mathbf{x}_{i,-k}) = 2\mathbf{x}_{:,k}^\top \mathbf{r}_k,\end{aligned}$$

where  $\mathbf{w}_{-k}$  corresponds to  $\mathbf{w}$  without component  $k$ ,  $\mathbf{x}_{i,-k}$  is  $\mathbf{x}_i$  without component  $k$ , and  $\mathbf{r}_k = \mathbf{y} - \mathbf{w}_{-k}^\top \mathbf{x}_{:, -k}$  is the residual due to using all the features except feature  $k$ .

**Hint:** Partition the weights into those involving  $k$  and those not involving  $k$ .

### 2 Lasso regression (2 points)

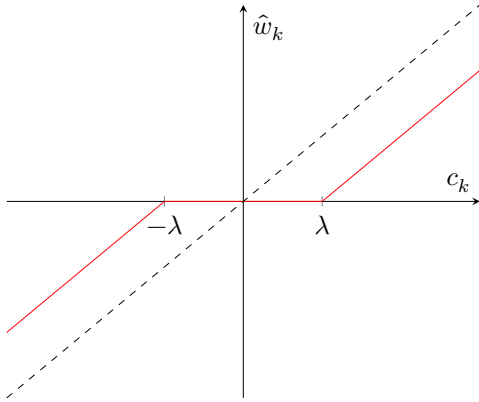
Recall the result of problem 1, where the RSS provides the smooth part of the lasso objective

$$\mathcal{L}(\mathbf{w}) = \text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1.$$

- Find the optimal parameters,  $\hat{w}_k$  as a function of  $c_k$ . You should find a relation as illustrated in the Figure.
- Instead of using the optimal parameters derived in (a), it helps to approximate it with a so-called soft-threshold function

$$\text{SoftThreshold}(x, \delta) = \text{sign}(x) \max(|x| - \delta, 0).$$

Sketch  $\hat{w}_d = \text{SoftThreshold}(\frac{c_d}{a_d}, \frac{\lambda}{a_d})$  and compare it with both the optimal parameters derived in (a), as well as the vanilla RSS from problem 1. Interpret how the two Lasso methods affect the optimal parameters.



**Hint:** We can generalize the notion of a derivative for certain functions with local discontinuities. Such functions are called subdifferentiable. A famous example is the absolute value function,  $f(x) = |x|$ . Its subdifferential is given by

$$\partial f(x) = \begin{cases} -1, & \text{if } x < 0 \\ [-1, 1], & \text{if } x = 0 \\ +1, & \text{if } x > 0 \end{cases}$$

where the notation  $[-1, 1]$  means any value between -1 and 1 inclusive.

### 3 Prostate cancer (6 points)

We will analyze a prostate cancer dataset, which you can download here: [prostate.csv](#). It contains 97 datapoints, each with the following features:

Variable	Description
lcavol	(log) Cancer volume
lweight	(log) Weight
age	Patient age
lbph	(log) Vening Prostatic Hyperplasia
svi	Seminal Vesicle Invasion
lcp	(log) Capsular Penetration
gleason	Gleason score
pgg45	Percent of Gleason score 4 or 5
lpsa	(log) Prostate Specific Antigen

The objective is to build a predictive model for `lpsa`. We will denote the vector of labels of `lpsa`,  $\mathbf{y}$ . The rest of the dataset will consist of the feature matrix,  $\mathbf{X}$ , made of the abovementioned variables, as well as an extra column  $[1, \dots, 1]^T$  to account for the intercept of your linear models.

- Plot the correlation between `lpsa` and each one of the three features that correlate the strongest with `lpsa`.
- Check that your dataset matrix,  $\mathbf{X}$ , is full rank.
- Split your data into a training and a test set. Keep 30 datapoints for the test set. (Feel free to use existing libraries, like `scikit-learn`.)

- (d) Implement yourself an ordinary least squares solver. Do not use existing statistics / machine learning libraries (though feel free to use linear-algebra libraries). To simplify your task, consider filling in the following template:

```
import numpy as np
from dataclasses import dataclass, field

@dataclass
class OrdinaryLeastSquares:
    model_params: np.ndarray = field(init=False)

    def __post_init__(self):
        self.model_params = None

    def fit(self, X_train_: np.ndarray, y_train_: np.ndarray) -> None:
        # self.model_params = ... (fit your model)
        pass

    def predict(self, X_test_: np.ndarray) -> np.ndarray:
        # make predictions for X_test_
        pass

    def rmse(self, X_test_: np.ndarray, y_test_: np.ndarray) -> float:
        # compute the root-mean-squared-error
        pass
```

Print out the set of optimal parameters for your training set, as well as the RMSE of the test set.

- (e) Implement yourself ridge regression. Consider extending the following template:

```
@dataclass
class RidgeRegression(OrdinaryLeastSquares):
    ridge_penalty: float

    def fit(self, X_train_: np.ndarray, y_train_: np.ndarray) -> None:
        # self.model_params = ... (fit your model)
        pass
```

where `ridge_penalty` is the coefficient in front of the  $\ell_2$  regularization term. Exclude the intercept from the regularization procedure. Plot the values of all model parameters (except the intercept) as a function of the regularization penalty term. Have the penalty term vary from 0 to 100. On a separate figure, plot the RMSE of the test set as a function of the penalty term. Indicate at what penalty value your best model is.

- (f) Implement yourself lasso regression. Use the objective function and soft threshold function from Problem 2 (b) to iteratively update your model parameters until convergence (controlled by `num_iter` in the template code below)

```
@dataclass
class LassoRegression(OrdinaryLeastSquares):
    lasso_penalty: float
```

```
def fit(self, X_train_: np.ndarray, y_train_: np.ndarray, num_iter: int = 1000) -> None:
    # self.model_params = ... (fit your model)
    pass
```

where `lasso_penalty` is the coefficient in front of the  $\ell_1$  regularization term. Exclude the intercept from the regularization procedure. Plot the values of all model parameters (except the intercept) as a function of the regularization penalty term. Have the penalty term vary from  $10^{-1}$  to  $10^2$ . On a separate figure, plot the RMSE of the test set as a function of the penalty term. Indicate at what penalty value your best model is.

- (g) Compare your sets of best parameters across the three solvers, and comment on the impact of ridge and lasso regression.