

# Homework 2, weeks 36 and 37

## Data Analysis and Machine Learning FYS-STK3155/FYS4155

Department of Physics, University of Oslo, Norway

Sep 1, 2020

\*

### Exercise 1: Adding Ridge and Lasso Regression

This exercise is a continuation of exercise 3 from exercise set 1 (week 35). We will use the same function to generate our data set, still staying with a simple function  $y(x)$  which we want to fit using linear regression, but now extending the analysis to include the Ridge and the Lasso regression methods. You can use the code under the Regression as an example on how to use the Ridge and the Lasso methods, see the [regression slides](#)).

We will thus again generate our own dataset for a function  $y(x)$  where  $x \in [0, 1]$  and defined by random numbers computed with the uniform distribution. The function  $y$  is a quadratic polynomial in  $x$  with added stochastic noise according to the normal distribution  $\mathcal{N}(t, \infty)$ .

The following simple Python instructions define our  $x$  and  $y$  values (with 100 data points).

```
x = np.random.rand(100)
y = 2.0+5*x*x+0.1*np.random.randn(100)
```

```
paragraphparagraph>paragraph>-0.5em
```

a) Write your own code for the Ridge method (see chapter 3.4 of Hastie *et al.*, equations (3.43) and (3.44)) and compute the parametrization for different values of  $\lambda$ . Compare and analyze your results with those from exercise 3. Study the dependence on  $\lambda$  while also varying the strength of the noise in your expression for  $y(x)$ .

**Solution.** The code here allows you to perform your own Ridge calculation and perform calculations for various values of the regularization parameter  $\lambda$ . This program can easily be extended upon.

```
import os
import numpy as np
```

```

import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

def R2(y_data, y_model):
    return 1 - np.sum((y_data - y_model) ** 2) / np.sum((y_data - np.mean(y_data)) ** 2)
def MSE(y_data, y_model):
    n = np.size(y_model)
    return np.sum((y_data - y_model) ** 2) / n

# A seed just to ensure that the random numbers are the same for every run.
# Useful for eventual debugging.
np.random.seed(3155)

x = np.random.rand(100)
y = 2.0 + 5 * x * x + 0.1 * np.random.randn(100)

# number of features p (here degree of polynomial
p = 3
# The design matrix now as function of a given polynomial
X = np.zeros((len(x), p))
X[:, 0] = 1.0
X[:, 1] = x
X[:, 2] = x * x
# We split the data in test and training data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
scaler = StandardScaler()
scaler.fit(X_train)
X_train_scaled = scaler.transform(X_train)
X_test_scaled = scaler.transform(X_test)

# matrix inversion to find beta
OLSbeta = np.linalg.inv(X_train.T @ X_train) @ X_train.T @ y_train
print(OLSbeta)
# and then make the prediction
ytildeOLS = X_train @ OLSbeta
print("Training R2 for OLS")
print(R2(y_train, ytildeOLS))
print("Training MSE for OLS")
print(MSE(y_train, ytildeOLS))
ypredictOLS = X_test @ OLSbeta
print("Test R2 for OLS")
print(R2(y_test, ypredictOLS))
print("Test MSE OLS")

```

```

print(MSE(y_test,ypredictOLS))

# Repeat now for Ridge regression and various values of the regularization parameter
I = np.eye(p,p)
# Decide which values of lambda to use
nlambdas = 20
MSEPredict = np.zeros(nlambdas)
MSETrain = np.zeros(nlambdas)
lambdas = np.logspace(-4, 1, nlambdas)
for i in range(nlambdas):
    lmb = lambdas[i]
    Ridgebeta = np.linalg.inv(X_train.T @ X_train+lmb*I) @ X_train.T @ y_train
    # and then make the prediction
    ytildeRidge = X_train @ Ridgebeta
    ypredictRidge = X_test @ Ridgebeta
    MSEPredict[i] = MSE(y_test,ypredictRidge)
    MSETrain[i] = MSE(y_train,ytildeRidge)
# Now plot the results
plt.figure()
plt.plot(np.log10(lambdas), MSETrain, label = 'MSE Ridge train')
plt.plot(np.log10(lambdas), MSEPredict, 'r--', label = 'MSE Ridge Test')
plt.xlabel('log10(lambda)')
plt.ylabel('MSE')
plt.legend()
plt.show()

```

paragraphparagraph>paragraph>-0.5em

b) Repeat the above but using the functionality of **Scikit-Learn**. Compare your code with the results from **Scikit-Learn**. Remember to run with the same random numbers for generating  $x$  and  $y$ .

**Solution.** To use **scikit-learn** with Ridge, we simply need to add the relevant function **Ridge()**, as done in the code here.

```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
import sklearn.linear_model as skl

def R2(y_data, y_model):
    return 1 - np.sum((y_data - y_model) ** 2) / np.sum((y_data - np.mean(y_data)) ** 2)
def MSE(y_data,y_model):
    n = np.size(y_model)
    return np.sum((y_data-y_model)**2)/n

```

```

# A seed just to ensure that the random numbers are the same for every run.
# Useful for eventual debugging.
np.random.seed(3155)

x = np.random.rand(100)
y = 2.0+5*x*x+0.1*np.random.randn(100)

# number of features p (here degree of polynomial
p = 3
# The design matrix now as function of a given polynomial
X = np.zeros((len(x),p))
X[:,0] = 1.0
X[:,1] = x
X[:,2] = x*x
# We split the data in test and training data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
scaler = StandardScaler()
scaler.fit(X_train)
X_train_scaled = scaler.transform(X_train)
X_test_scaled = scaler.transform(X_test)

# matrix inversion to find beta
OLSbeta = np.linalg.inv(X_train.T @ X_train) @ X_train.T @ y_train
print(OLSbeta)
# and then make the prediction
ytildeOLS = X_train @ OLSbeta
print("Training R2 for OLS")
print(R2(y_train,ytildeOLS))
print("Training MSE for OLS")
print(MSE(y_train,ytildeOLS))
ypredictOLS = X_test @ OLSbeta
print("Test R2 for OLS")
print(R2(y_test,ypredictOLS))
print("Test MSE OLS")
print(MSE(y_test,ypredictOLS))

# Repeat now for Ridge regression and various values of the regularization parameter
I = np.eye(p,p)
# Decide which values of lambda to use
nlambdas = 100
MSEPredict = np.zeros(nlambdas)
MSEPredictSKL = np.zeros(nlambdas)
MSETrain = np.zeros(nlambdas)
lambdas = np.logspace(-4, 0, nlambdas)

```

```

for i in range(nlambdas):
    lmb = lambdas[i]
    # add ridge
    clf_ridge = skl.Ridge(alpha=lmb).fit(X_train, y_train)
    yridge = clf_ridge.predict(X_test)
    Ridgebeta = np.linalg.inv(X_train.T @ X_train+lmb*I) @ X_train.T @ y_train
    # and then make the prediction
    ytildeRidge = X_train @ Ridgebeta
    ypredictRidge = X_test @ Ridgebeta
    MSEPredict[i] = MSE(y_test,ypredictRidge)
    MSEPredictSKL[i] = MSE(y_test,yridge)
    MSETrain[i] = MSE(y_train,ytildeRidge)
#then plot the results
plt.figure()
plt.plot(np.log10(lambdas), MSETrain, label = 'MSE Ridge train')
plt.plot(np.log10(lambdas), MSEPredict, 'r--', label = 'MSE Ridge Test')
plt.plot(np.log10(lambdas), MSEPredictSKL, 'g--', label = 'MSE Ridge sickit-learn Test')
plt.xlabel('log10(lambda)')
plt.ylabel('MSE')
plt.legend()
plt.show()

```

paragraphparagraph>paragraph>-0.5em

c) Our next step is to study the variance of the parameters  $\beta_1$  and  $\beta_2$  (assuming that we are parameterizing our function with a second-order polynomial). We will use standard linear regression and the Ridge regression. You can now opt for either writing your own function or using **Scikit-Learn** to find the parameters  $\beta$ . From your results calculate the variance of these parameters (recall that this is equal to the diagonal elements of the matrix  $(\hat{X}^T \hat{X}) + \lambda \hat{I})^{-1}$ ). Discuss the results of these variances as functions of  $\lambda$ . In particular, try to link your discussion with the discussion in Hastie *et al.* and their figures 3.10 and 3.11. **Scikit-Learn** may not provide the variance of the parameters  $\beta$ . This needs to be checked. With your own code you can however do so.

### Solution.

```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
import sklearn.linear_model as skl

def R2(y_data, y_model):
    return 1 - np.sum((y_data - y_model) ** 2) / np.sum((y_data - np.mean(y_data)) ** 2)
def MSE(y_data,y_model):

```

```

n = np.size(y_model)
return np.sum((y_data-y_model)**2)/n

# A seed just to ensure that the random numbers are the same for every run.
# Useful for eventual debugging.
np.random.seed(3155)

x = np.random.rand(100)
y = 2.0+5*x*x+0.1*np.random.randn(100)

# number of features p (here degree of polynomial
p = 3
# The design matrix now as function of a given polynomial
X = np.zeros((len(x),p))
X[:,0] = 1.0
X[:,1] = x
X[:,2] = x*x
# We split the data in test and training data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
scaler = StandardScaler()
scaler.fit(X_train)
X_train_scaled = scaler.transform(X_train)
X_test_scaled = scaler.transform(X_test)

# matrix inversion to find beta
OLSbeta = np.linalg.inv(X_train.T @ X_train) @ X_train.T @ y_train
print(OLSbeta)
# The variance is given by the inverse of the matrix  $X^T X$ 
print(np.linalg.inv(X_train.T @ X_train))

# Repeat now for Ridge regression and various values of the regularization parameter
I = np.eye(p,p)
# Decide which values of lambda to use
nlambda = 10
MSEPredict = np.zeros(nlambda)
MSEPredictSKL = np.zeros(nlambda)
MSETrain = np.zeros(nlambda)
lambda = np.logspace(-4, 0, nlambda)
for i in range(nlambda):
    lmb = lambda[i]
    Ridgebeta = np.linalg.inv(X_train.T @ X_train+lmb*I) @ X_train.T @ y_train
    print(np.linalg.inv(X_train.T @ X_train+lmb*I))

```

paragraphparagraph>paragraph>-0.5em

d) Repeat the previous step but add now the Lasso method, see equation (3.53) of Hastie *et al.*. Discuss your results and compare with standard regression and the Ridge regression results. You can write your own code or use the functionality of **scikit-learn**. We recommend the latter since we have not yet discussed how to solve the Lasso equations numerically. Also, you do not need to compute the variance of the parameters  $\beta$  but you can extract their values and study their behavior as functions of the regularization parameter  $\lambda$ .

### Solution.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
import sklearn.linear_model as skl
#from sklearn.linear_model import LinearRegression, Ridge, Lasso
def R2(y_data, y_model):
    return 1 - np.sum((y_data - y_model) ** 2) / np.sum((y_data - np.mean(y_data)) ** 2)
def MSE(y_data, y_model):
    n = np.size(y_model)
    return np.sum((y_data - y_model) ** 2) / n

# A seed just to ensure that the random numbers are the same for every run.
# Useful for eventual debugging.
np.random.seed(3155)

x = np.random.rand(100)
y = 2.0 + 5 * x * x + 0.1 * np.random.randn(100)

# number of features p (here degree of polynomial
p = 3
# The design matrix now as function of a given polynomial
X = np.zeros((len(x), p))
X[:, 0] = 1.0
X[:, 1] = x
X[:, 2] = x * x
# We split the data in test and training data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
scaler = StandardScaler()
scaler.fit(X_train)
X_train_scaled = scaler.transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

```

# matrix inversion to find beta
OLSbeta = np.linalg.inv(X_train.T @ X_train) @ X_train.T @ y_train
print(OLSbeta)
# and then make the prediction
ytildeOLS = X_train @ OLSbeta
print("Training R2 for OLS")
print(R2(y_train,ytildeOLS))
print("Training MSE for OLS")
print(MSE(y_train,ytildeOLS))
ypredictOLS = X_test @ OLSbeta
print("Test R2 for OLS")
print(R2(y_test,ypredictOLS))
print("Test MSE OLS")
print(MSE(y_test,ypredictOLS))

# Repeat now for Ridge regression and various values of the regularization parameter
I = np.eye(p,p)
# Decide which values of lambda to use
nlambdas = 100
MSEPredictLasso = np.zeros(nlambdas)
MSEPredictRidge = np.zeros(nlambdas)
lambdas = np.logspace(-4, 0, nlambdas)
for i in range(nlambdas):
    lmb = lambdas[i]
    # add ridge
    clf_ridge = skl.Ridge(alpha=lmb).fit(X_train, y_train)
    clf_lasso = skl.Lasso(alpha=lmb).fit(X_train, y_train)
    yridge = clf_ridge.predict(X_test)
    ylasso = clf_lasso.predict(X_test)
    MSEPredictLasso[i] = MSE(y_test,ylasso)
    MSEPredictRidge[i] = MSE(y_test,yridge)
#then plot the results
plt.figure()
plt.plot(np.log10(lambdas), MSEPredictRidge, 'r--', label = 'MSE Ridge Test')
plt.plot(np.log10(lambdas), MSEPredictLasso, 'g--', label = 'MSE Lasso Test')
plt.xlabel('log10(lambda)')
plt.ylabel('MSE')
plt.legend()
plt.show()

```

paragraphparagraph>paragraph>-0.5em

e) Finally, using **Scikit-Learn** or your own code, compute also the mean square error, a risk metric corresponding to the expected value of the squared



(quadratic) error defined as

$$MSE(\hat{y}, \tilde{y}) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2,$$

and the  $R^2$  score function. If  $\tilde{y}_i$  is the predicted value of the  $i$ -th sample and  $y_i$  is the corresponding true value, then the score  $R^2$  is defined as

$$R^2(\hat{y}, \tilde{y}) = 1 - \frac{\sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2}{\sum_{i=0}^{n-1} (y_i - \bar{y})^2},$$

where we have defined the mean value of  $\hat{y}$  as

$$\bar{y} = \frac{1}{n} \sum_{i=0}^{n-1} y_i.$$

Discuss these quantities as functions of the variable  $\lambda$  in the Ridge and Lasso regression methods.

**Solution.** These results can all be studied with the codes we have above. These scores are included in the codes above.

\*

Exercise 2: Normalizing our data

A much used approach before starting to train the data is to preprocess our data. Normally the data may need a rescaling and/or may be sensitive to extreme values. Scaling the data renders our inputs much more suitable for the algorithms we want to employ.

**Scikit-Learn** has several functions which allow us to rescale the data, normally resulting in much better results in terms of various accuracy scores. The **StandardScaler** function in **Scikit-Learn** ensures that for each feature/predictor we study the mean value is zero and the variance is one (every column in the design/feature matrix). This scaling has the drawback that it does not ensure that we have a particular maximum or minimum in our data set. Another function included in **Scikit-Learn** is the **MinMaxScaler** which ensures that all features are exactly between 0 and 1. The

The **Normalizer** scales each data point such that the feature vector has a euclidean length of one. In other words, it projects a data point on the circle (or sphere in the case of higher dimensions) with a radius of 1. This means every data point is scaled by a different number (by the inverse of its length). This normalization is often used when only the direction (or angle) of the data matters, not the length of the feature vector.

The **RobustScaler** works similarly to the **StandardScaler** in that it ensures statistical properties for each feature that guarantee that they are on the same scale. However, the **RobustScaler** uses the median and quartiles, instead of mean

and variance. This makes the RobustScaler ignore data points that are very different from the rest (like measurement errors). These odd data points are also called outliers, and might often lead to trouble for other scaling techniques.

It also common to split the data in a **training** set and a **testing** set. A typical split is to use 80% of the data for training and the rest for testing. This can be done as follows with our design matrix  $X$  and data  $y$  (remember to import **scikit-learn**)

```
# split in training and test data
X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2)
```

Then we can use the standard scaler to scale our data as

```
scaler = StandardScaler()
scaler.fit(X_train)
X_train_scaled = scaler.transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

In this exercise we want you to compute the MSE for the training data and the test data as function of the complexity of a polynomial, that is the degree of a given polynomial. We want you also to compute the  $R^2$  score as function of the complexity of the model for both training data and test data. You should also run the calculation with and without scaling.

One of the aims is to reproduce Figure 2.11 of [Hastie et al.](#) We will also use Ridge and Lasso regression.

Our data is defined by  $x \in [-3, 3]$  with a total of for example 100 data points.

```
np.random.seed()
n = 100
maxdegree = 14
# Make data set.
x = np.linspace(-3, 3, n).reshape(-1, 1)
y = np.exp(-x**2) + 1.5 * np.exp(-(x-2)**2) + np.random.normal(0, 0.1, x.shape)
```

where  $y$  is the function we want to fit with a given polynomial.

paragraphparagraph>paragraph>-0.5em

a) Write a first code which sets up a design matrix  $X$  defined by a fifth-order polynomial. Scale your data and split it in training and test data.

**Solution.**

```
import matplotlib.pyplot as plt
import numpy as np
from sklearn.linear_model import LinearRegression, Ridge, Lasso
from sklearn.preprocessing import PolynomialFeatures
from sklearn.model_selection import train_test_split
from sklearn.pipeline import make_pipeline
```

```

np.random.seed(2018)
n = 50
maxdegree = 5
# Make data set.
x = np.linspace(-3, 3, n).reshape(-1, 1)
y = np.exp(-x**2) + 1.5 * np.exp(-(x-2)**2) + np.random.normal(0, 0.1, x.shape)
TestError = np.zeros(maxdegree)
TrainError = np.zeros(maxdegree)
polydegree = np.zeros(maxdegree)
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
scaler = StandardScaler()
scaler.fit(X_train)
x_train_scaled = scaler.transform(x_train)
x_test_scaled = scaler.transform(x_test)

for degree in range(maxdegree):
    model = make_pipeline(PolynomialFeatures(degree=degree), LinearRegression(fit_intercept=False))
    clf = model.fit(x_train_scaled, y_train)
    y_fit = clf.predict(x_train_scaled)
    y_pred = clf.predict(x_test_scaled)
    polydegree[degree] = degree
    TestError[degree] = np.mean( np.mean((y_test - y_pred)**2) )
    TrainError[degree] = np.mean( np.mean((y_train - y_fit)**2) )

plt.plot(polydegree, TestError, label='Test Error')
plt.plot(polydegree, TrainError, label='Train Error')
plt.legend()
plt.show()

```

paragraphparagraph>paragraph>-0.5em

b) Perform an ordinary least squares and compute the means squared error and the  $R^2$  factor for the training data and the test data, with and without scaling.

**Solution.** This requires a simple extension to the above code where you simply add a statement calling the  $R^2$  function included in the same code.

paragraphparagraph>paragraph>-0.5em

c) Add now a model which allows you to make polynomials up to degree 15. Perform a standard OLS fitting of the training data and compute the MSE and  $R^2$  for the training and test data and plot both test and training data MSE and  $R^2$  as functions of the polynomial degree. Compare what you see with Figure 2.11 of Hastie et al. Comment your results. For which polynomial degree do you find an optimal MSE (smallest value)?

**Solution.** Here you simply need to change the degree of the polynomial in the above code to  $n = 15$ .

paragraphparagraph>paragraph>-0.5em

d) Repeat part (2c) but now using Ridge regressions with various hyperparameters  $\lambda$ . Make the same plots for the optimal  $\lambda$  value for each polynomial degree. Compare these results with those from the standard OLS approach.

**Solution.** Here you need to add for example the same loop over the parameters  $\lambda$  as you did in the first exercise, that is add

```
nlambdas = 100
MSEPredictRidge = np.zeros(nlambdas)
lambdas = np.logspace(-4, 0, nlambdas)
for i in range(nlambdas):
    lmb = lambdas[i]
    # add ridge
    clf_ridge = skl.Ridge(alpha=lmb).fit(X_train_scaled, y_train)
```

The plotting functionality of the first exercise can be reused here as well.