
Assignment IV: Bayes' Theorem and Ridge Regression

Exercises in Machine Learning (190.013), SS2022
Stefan Nehl¹

¹stefan-christopher.nehl@stud.unileoben.ac.at, MNR: 00935188, Montanuniversität Leoben, Austria

May 9, 2022

In the fourth assignment, I had to solve three different tasks. First, calculate the probability for a positive test result with the Bayes' Theorem, next describe the ridge regression and derive the weight update for with the least squares regression and last implement the ridge regression. The implementation of the ridge regression also includes testing the model and plotting its results.

1 Bayes' Theorem

The Bayes' Theorem is a mathematical formula which describes the probability of an event. Furthermore, it is used for calculating conditional probabilities. (Joyce, 2021)

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}$$

(Rueckert, 2022)

I used this formula to calculate the probability to be infected with SARS CoV2 and having a positive test result of an antigen test. Let $A \in [\text{infected}, \text{non-infected}]$ the event, which defines if a person is infected or not, and $B \in [+, -]$ the event, which defines the result of the antigen test.

1.1 Implementation

First, I created the following variables with the values.

Table 1: Variables and Values

name	value
<i>populationAustria</i>	9095538
<i>activeCases</i>	441098
<i>covTestSensitivity</i>	0.971
<i>covTestSpecific</i>	0.995

First, I set the variable for $p(+|inf)$ to the value of *covTestSensitivity* and the variable for $p(-|nInf)$ to the value of *covTestSpecific*. Next, I calculated the value for $p(inf)$, $p(nInf)$ and stored the values in the variables *pInfected* and *pNotInfected*.

$$p(inf) = \text{activeCases} / \text{populationAustria}$$

$$p(nInf) = 1 - p(inf)$$

The variable $p(inf)$ defines the value for the probability to be infected with covid and $p(nInf)$ not. Furthermore, the abbreviation for infected is *inf* and for non infected *nInf*. The abbreviation for having a positive test result is $+$ and for a negative test result $-$. Next, I initialized the following variables and calculated their values with the following formulas.

$$p(nInf \& -) = p(nInf) * p(-|non-infected)$$

$$p(-) = p(inf \& -) + p(nInf \& -)$$

$$p(nInf \& +) = p(nInf) - p(nInf \& -)$$

$$p(+) = p(inf \& +) + p(nInf \& +)$$

$$p(-|inf) = \frac{p(inf \& -)}{p(-)}$$

$$p(+|nInf) = \frac{p(nInf \& +)}{p(nInf)}$$

Last, I used the Bayes' Theorem to calculate the $p(infected|+)$ value.

$$p(inf|+) = \frac{p(+|inf) * p(inf)}{p(+)}$$

1.2 Result and Conclusion

The results of the calculation is displayed in Table 2. The result for $p(inf|+)$ is $0.630525 \approx 63.05\%$. Which means there is a 63.05% chance to be infected with covid and get a positive test result. The implemented code can be found in the appendix of this paper.

Table 2: Results

name	value
$p(- inf)$	00.15%
$p(+ nInf)$	02.90%
$p(inf)$	04.85%
$p(nInf)$	95.15%
$p(+)$	07.47%
$p(inf +)$	63.05%

2 Ridge Regression

Ridge regression is used for parameter estimation to address the collinearity problem in multiple linear regression. (McDonald, 2009) The Ridge Regression adds the quadratic regularization term $\frac{\lambda}{2} (\omega^T \omega)$ to the objective J_{LS} . (Rueckert, 2022)

2.1 Derivation of the Least Squares Solution

For the derivation of the least squares I defined the vectors $\mathbf{y} \in \mathbb{R}$, $\mathbf{A} \in \mathbb{R}^{n \times M}$ and $\omega \in \mathbb{R}^M$ where M is the dimension and n the number of samples.

$$\frac{\partial J_{LS}}{\partial \omega} = \frac{\partial}{\partial \omega} \{1/2\sigma^{-2}(\mathbf{y} - \mathbf{A}\omega)^T(\mathbf{y} - \mathbf{A}\omega)\}$$

After the partial deviation we receive the following equation.

$$\frac{\partial J_{LS}}{\partial \omega} = 1/2\sigma^{-2}(-2\mathbf{y}^T \mathbf{A} + 2\omega^T \mathbf{A}^T \mathbf{A})$$

I set this equation to zero to calculate the least square solution.

$$\begin{aligned} \frac{\partial J_{LS}}{\partial \omega} &= 0, \\ 1/2\sigma^{-2}(-2\mathbf{y}^T \mathbf{A} + 2\omega^T \mathbf{A}^T \mathbf{A}) &= 0, \\ 1/2\sigma^{-2}(-2\mathbf{y}^T \mathbf{A} + 2\omega^T \mathbf{A}^T \mathbf{A}) &= 0, \\ -\mathbf{y}^T \mathbf{A} + \omega^T \mathbf{A}^T \mathbf{A} \omega &= 0, \\ \omega &= (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}. \end{aligned}$$

(Rueckert, 2022) Important here is, that the matrix \mathbf{A} has a full rank and is invertible. If this is not the case I would use the Moore–Penrose inverse which is described in the following formula.

$$\mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T.$$

(Rueckert, 2022)

3 Implementation of Ridge Regression

The last task was to implement the ridge regression for a given dataset. The dataset includes longitude and latitude of a map with the corresponding temperature data.

3.1 Import Data and Implementation

For the implementation I first created the abstract class *Regression* which includes the abstract methods *importData*, *generateTrainingSubset*, *computeLinearRidgeRegression*, *testModel*, *computeError*, *plotError*, *plotHeatMap*, *computeMeanOfError*. Then I created the class *RidgeRegression* which implements those methods and has the parameter *trainStep* which indicates the size of the training data. For importing the data I used the *scipy.io* library which is able to read *MatLab* files and load this data. I used only the first dataset of the time series data to create the model. The method *generateTrainingSubset* creates the training data with the parameter *trainStep*. The parameter *trainStep* defines the size of the training set. For example, if the value is set to 4 every 4. values is used for the calculation of the weight values.

3.2 Ridge Regression

For the Implementation of the Ridge Regression calculation I added the method *computeLinearRidgeRegression* and made the calculation. I followed the formula from chapter 2 and used the *numpy* library to do the matrix calculations. I added a helper method which I used to create the feature vector for the calculation. The method is named *createFeatureVector* and takes the parameter *x*. The parameter *x* is the vector of the current y-value. In our case it's a two dimensional vector with the longitude and latitude. I augmented this vector and created the new vector *featureVector* which is a three dimensional vector. The first dimension contains a 1, the second the latitude and the third the longitude. This vector was then returned from the method. After the creation of the feature vectors, the method *computeLinearRidgeRegression* finishes its calculations and returns the weight vector. The weight vector is used in the method *testModel* where the vector is multiplied with the given test values in the dataset.

3.3 Calculate Error

After testing the model I calculated the error between the model and the provided test data. For the error computed I added the method *computeError*, which takes a list of y values, stored in the variable *yStar*, to compute the difference between the calculated result and the parameter *yStar*. Next, I sorted the values descending for plotting and created a *panda data frame* with the *pandas* library. Last, I computed also the mean of the error values with *numpy*.

3.4 Plotting

For plotting I used the *seaborn* library for the heatmap and the descending bar chart and *matplotlib* for the error differences between different lambdas and train

steps. Both libraries, *seaborn* and *matplotlib* are working great together, because *seaborn* is build on top of *matplotlib*. Because of the native support of *pandas* with the *seaborn* I used the data frame type for preparing the data. For the descending error plot I just created a data frame and passed the value to the function *barplot* in the *seaborn* library. The heat map was a little bit more difficult to create. I created a data frame and used the *pivot* function and sorted the values descending in respect of the longitude value. The sort operation was used, because of the ascending longitude values in the dataset. The prepared data frame was passed to the *heatmap* function of the *seaborn* library.

3.5 Results

The error of the model between the is displayed in Figure 1. This figure shows a descending error plot with a λ from 0.1. Raising λ does not made any changes in my model. Even raising the value to 50, displayed in Figure 2, only changes the value slightly. The train step was at 4 with both plots. The small changes of

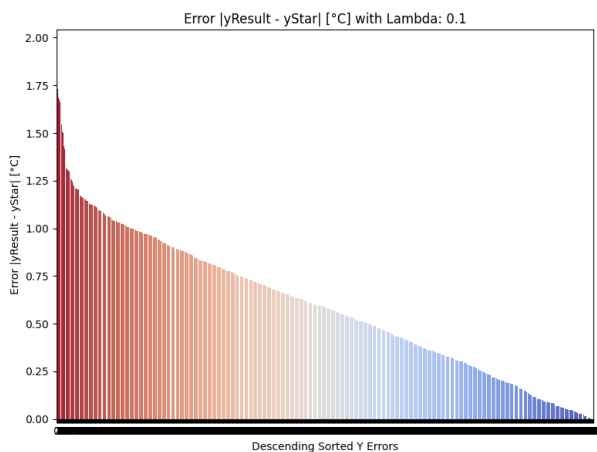


Figure 1: Descending error plot with $\lambda:0.1$

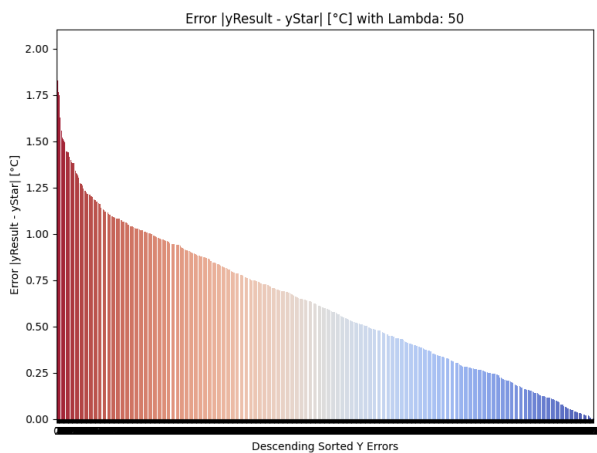


Figure 2: Descending error plot with $\lambda:50$

the error with a different λ is also displayed in Figure 3

and 4. The highest error is exactly on the same position in the heat map, only the value itself changed slightly. After I made some tests with the different λ values

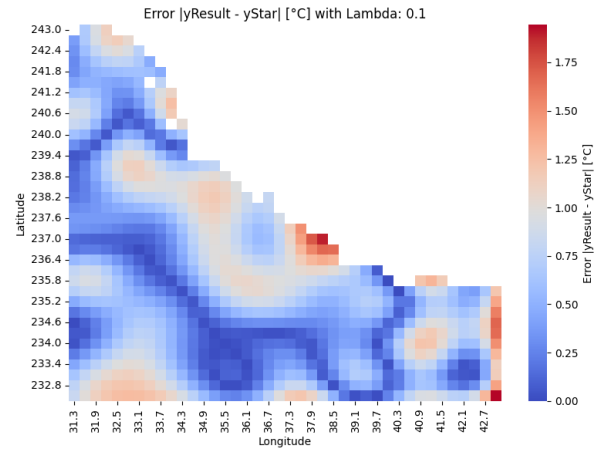


Figure 3: Heatmap with $\lambda:0.1$

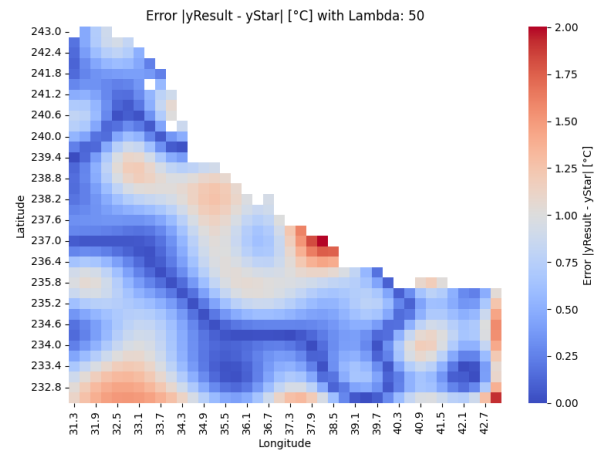


Figure 4: Heatmap with $\lambda:50$

I also started to adjust the train step. I reduced the train step to 1 and repeated the tests with different λ values. I displayed the results in Figure 5. Figure 5 shows, that the changes of the train step and λ values only made small adjustment to the results itself. Also the mean values of every tests has only small changes in the value itself. The result of the mean values is displayed in Table 3.

3.6 Conclusion

I had some difficulties with the Implementation of the ridge regression. Especially the augmentation of the values gave me some headache. I even tried to use a *Gaussian Basis Function* with the *Polynomial Basis Function* to create a feature vector. However, all the tries didn't result in a better model. Furthermore, the results were much worse than the straight forward implementation of the feature vector with 1, Latitude

and Longitude. One explanation of the error could be the small size of the data set. A larger data set could improve the accuracy of the model. The code of the implementation is in the appendix of this paper. Small side note, I had to remove the correct temperature unit from the code because it resulted in an *UTF-8* issue.

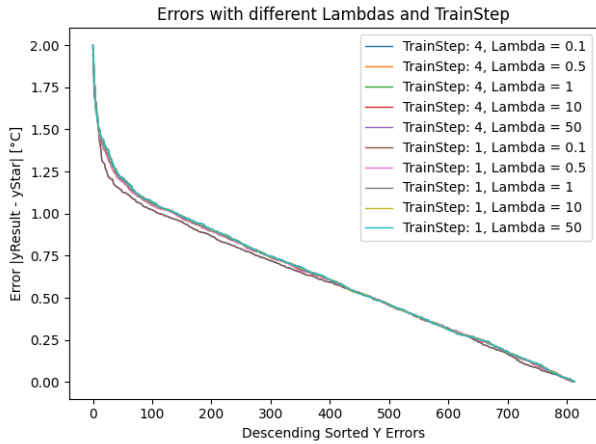


Figure 5: Descending error plot with different lambdas and train steps

Bibliography

- Joyce, James (2021). "Bayes' Theorem". In: *The Stanford Encyclopedia of Philosophy*. Ed. by Edward N. Zalta. Fall 2021. Metaphysics Research Lab, Stanford University.
- McDonald, Gary C. (2009). "Ridge regression". In: *WIREs Computational Statistics* 1.1, pp. 93–100. doi: <https://doi.org/10.1002/wics.14>. eprint: <https://wires.onlinelibrary.wiley.com/doi/pdf/10.1002/wics.14>. URL: <https://wires.onlinelibrary.wiley.com/doi/abs/10.1002/wics.14>.
- Rueckert, Elmar (2022). *An Introduction to Probabilistic Machine Learning*. Elmar Rueckert.

Table 3: Mean Errors

Train Step	λ	Mean
4	0.1	0.5938
4	0.5	0.6117
4	1.0	0.6156
4	10.0	0.6195
4	50.0	0.6198
1	0.1	0.5938
1	0.5	0.6117
1	1.0	0.6156
1	10.0	0.6195
1	50.0	0.6198

APPENDIX

```
1 import numpy as np
2
3 populationAustria = 9095538
4 activeCases = 441098
5 covTestSensitivity = 0.971
6 covTestSpecific = 0.995
7
8 #A...[infected, not-infected]
9 #B...[positive, negative]
10
11 pInfected = activeCases / populationAustria
12 pNotInfected = 1 - pInfected
13
14 pPositiveInfected = covTestSensitivity
15 pNegativeNotInfected = covTestSpecific
16
17 pInfectedAndPositive = pInfected * pPositiveInfected
18 pInfectedAndNegative = pInfected - pInfectedAndPositive
19
20 pNotInfectedAndNegative = pNotInfected * pNegativeNotInfected
21 pNegative = pInfectedAndNegative + pNotInfectedAndNegative
22
23 pNotInfectedAndPositive = pNotInfected - pNotInfectedAndNegative
24 pPositive = pInfectedAndPositive + pNotInfectedAndPositive
25
26 pNegativeInfected = pInfectedAndNegative / pNegative
27 pPositiveNotInfected = pNotInfectedAndPositive / pNotInfected
28
29 #BayesTheorem
30 pInfectedPositive = pPositiveInfected * pInfected / pPositive
31
32 print("p(-|inf): " + str(pNegativeInfected))
33 print("p(+|nInf): " + str(pPositiveNotInfected))
34 print("p(inf): " + str(pInfected))
35 print("p(nInf): " + str(pNotInfected))
36 print("p(+): " + str(pPositive))
37 print("p(inf|+): " + str(pInfectedPositive))
```

```

1 import sys
2 import seaborn as sbr
3 import matplotlib.pyplot as plt
4
5 sys.path.insert(0, '../modules')
6 from RidgeRegression import RidgeRegression
7
8 def doTest(lambdaValue, trainStep, plot=False):
9     ridgeRegression = RidgeRegression(100)
10    print("import data")
11    ridgeRegression.importData()
12    print("generate Trainingset")
13    ridgeRegression.generateTrainingSubset()
14    print("Train")
15    weightVector = ridgeRegression.computeLinearRidgeRegression(lambdaValue)
16    ridgeRegression.testModel(weightVector)
17    print("Calculate Error")
18    yTest = ridgeRegression.getYTestData()
19    ridgeRegression.computeError(yTest)
20    ridgeRegression.computeMeanOfError()
21    print("Mean Error with Lambda " + str(lambdaValue) + ": " + str(
        ridgeRegression.getMeanError()))
22
23    if plot:
24        print("plot error")
25        ridgeRegression.plotError()
26        print("plot heatmap")
27        ridgeRegression.plotHeatMap()
28
29    return ridgeRegression.getDescSortedError()
30
31 firstTrainSetErrors = []
32
33 trainStep = 4
34 firstTrainSetErrors.append((0.1, trainStep, (doTest(0.1, trainStep, plot=True))))
35 firstTrainSetErrors.append((0.5, trainStep, (doTest(0.5, trainStep))))
36 firstTrainSetErrors.append((1, trainStep, (doTest(1, trainStep))))
37 firstTrainSetErrors.append((10, trainStep, (doTest(10, trainStep))))
38 firstTrainSetErrors.append((50, trainStep, (doTest(50, trainStep, plot=True))))
39
40 trainStep = 1
41 firstTrainSetErrors.append((0.1, trainStep, (doTest(0.1, trainStep))))
42 firstTrainSetErrors.append((0.5, trainStep, (doTest(0.5, trainStep))))
43 firstTrainSetErrors.append((1, trainStep, (doTest(1, trainStep))))
44 firstTrainSetErrors.append((10, trainStep, (doTest(10, trainStep))))
45 firstTrainSetErrors.append((50, trainStep, (doTest(50, trainStep))))
46
47 for testSet in firstTrainSetErrors:
48     lambdaValue = testSet[0]
49     trainStep = testSet[1]
50     errors = testSet[2]
51     x = errors["values"]
52     y = errors["error"]
53     plt.plot(x, y, linewidth=1, label=f"TrainStep: {trainStep}, Lambda = {
        lambdaValue}")
54
55 plt.title("Errors with different Lambdas and TrainStep")
56 plt.xlabel("Descending Sorted Y Errors")
57 plt.ylabel("Error |yResult - yStar| [C]")

```

```
58 plt.tight_layout()  
59 plt.legend()  
60 plt.show()
```



```

1 import matplotlib.pyplot
2 import scipy.io as sio
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import matplotlib.ticker as ticker
6 import seaborn as sbr
7 import pandas as pd
8
9 from inference import Regression
10 from GaussDistribution import GaussDistribution
11
12 class RidgeRegression(Regression):
13
14     def __init__(self, trainStep):
15         self.trainStep = trainStep
16         self.hasGenerated1DTestData = False
17
18     def importData(self):
19         dataDictionary = sio.loadmat("AssignmentIV_data_set.mat")
20
21         tempTimeData = dataDictionary.get("TempField")
22         self.tempData = np.array(tempTimeData[:, :, 0])
23
24         longData = dataDictionary.get("LongitudeScale") #x-value: Long
25         latData = dataDictionary.get("LatitudeScale") #y-value: Lat
26
27         self.x_test = np.array(dataDictionary.get("x_test")) # Testdata from
28         self.y_test = np.array(dataDictionary.get("y_test")) # Testdata from
29         DataSet
30
31         longLatData = []
32         arrayValues = []
33
34         for y in range(len(self.tempData)):
35             for x in range(len(self.tempData[y])):
36                 value = self.tempData[y,x]
37                 if np.isnan(value):
38                     continue
39
40                 arrayValues.append(value)
41                 longLatData.append((latData[y][0], longData[x][0])) #[0] needed
42                 because of strange import of long/latData
43
44         self.inputValues = np.array(longLatData)
45         self.outputValues = np.array(arrayValues)
46         self.numberOfSamples = len(self.inputValues)
47
48     def generateTrainingSubset(self):
49         self.trainSubsetInput = np.array(self.inputValues[0:len(self.inputValues)
50             ):self.trainStep])
51         self.trainSubsetOutput = np.array(self.outputValues[0:len(self.
52             outputValues):self.trainStep])
53
54     def createFeatureVector(self, x):
55         featureVector = []
56         featureVector.append(1)
57
58         for i in range(len(x)):

```

```

55         newXVector = x[i]
56         featureVector.append(newXVector)
57
58     return featureVector
59
60 def computeLinearRidgeRegression(self, lambdaValue):
61     self.lambdaValue = lambdaValue
62     X = np.vstack([self.createFeatureVector(x) for x in self.
63                     trainSubsetInput]))
64     Y = np.vstack([y for y in self.trainSubsetOutput]))
65
66     XT = np.transpose(X)
67     XTX = np.matmul(XT, X) + self.lambdaValue * np.identity(X.shape[1])
68     self.weightVector = np.matmul(np.linalg.inv(XTX), XT), Y
69     return self.weightVector
70
71 def testModel(self, weight):
72     self.yResult = np.hstack([x @ self.weightVector for x in self.x_test])
73
74 def getYTestData(self):
75     return self.y_test
76
77 def computeError(self, yStar):
78     self.yError = np.transpose(abs(self.yResult - yStar))
79     reversedArray = np.flip(np.sort(self.yError, 0))
80     self.errorDataFrame = pd.DataFrame({
81         'values': range(len(reversedArray)),
82         'error': [error[0] for error in reversedArray]
83     })
84
85 def getDescSortedError(self):
86     return self.errorDataFrame
87
88 def computMeanOfError(self):
89     self.meanError = np.mean(self.yError)
90
91 def getMeanError(self):
92     return self.meanError
93
94 def plotError(self):
95     plt.figure(figsize=(8, 6))
96     errorPlot = sbr.barplot(data=self.errorDataFrame, x="values", y="error",
97                             palette="coolwarm_r")
98
99     plt.xlabel("Descending Sorted Y Errors")
100    plt.ylabel("Error |yResult - yStar| [C]")
101    plt.title("Error |yResult - yStar| [C] with Lambda: " + str(self.
102                lambdaValue))
103    plt.tight_layout()
104    matplotlib.pyplot.show()
105
106 def plotHeatMap(self):
107     tempErrorData = \
108     {
109         'Lat': [round(long,2) for long in self.x_test[:, 1]],
110         'Long': [round(lat,2) for lat in self.x_test[:, 2]],
111         'error': [error[0] for error in self.yError]
112     }

```

```
111 tempDataFrame = pd.DataFrame(tempErrorData)
112 tempDataFrame = tempDataFrame.pivot("Long", "Lat", "error")
113 reversedTempErrorData = tempDataFrame.sort_values(("Long"), ascending=
    False)
114
115 def fmt(x, y):
116     return '{:,.2f}'.format(x)
117
118 plt.figure(figsize=(8,6))
119 errorHeatMap = sbr.heatmap(reversedTempErrorData, vmin=0.0, cmap="
    coolwarm", cbar_kws={"label": "Error |yResult - yStar| [C]"})
120 ax = errorHeatMap.axes
121
122 plt.xlabel("Longitude")
123 plt.ylabel("Latitude")
124 plt.title("Error |yResult - yStar| [C] with Lambda: " + str(self.
    lambdaValue))
125 plt.tight_layout()
126 matplotlib.pyplot.show()
```