

Precision Preprocessing and Hyperparameter Optimization & A KNN-Based Framework for College Retention Prediction

Samson Rozansky, Brian Zhang

William G. Enloe High School

IB Computer Science HL 4B

March 8th, 2024

§1 Introduction

§1.1 Research Objective

College is one of the most effective tools that facilitates class mobility, but 24% of students drop out in their first year of college. We want to reduce this amount by attempting to create a model that can detect “at-risk” students and give them the appropriate resources needed for a better future. This in turn will enhance economic opportunities for thousands of students and result in a more productive economy overall.

Most universities begin reaching out to students when they start failing courses. We believe this addresses the problem too late. We want to be very proactive at addressing these problems before their GPA start taking hits. This is because many graduate programs, internships, REU's, and many more consider GPA as a critical factor for deciding whether they should admit one student over another. So being proactive is one of the main reasons we began our inquiry.

§1.2 Research Questions

- What is the best distance to use to determine students that are “at-risk”?
- What is the most optimal ”K” amount for classifying students as ”at-risk”?
- What are some attributes we can throw away to improve accuracy?
- What are some attributes that are important to detecting “at-risk” students?

§2 Data

Our data came from the Capacitação da Administração Pública, or Public Administration Testing in Portugal. It came with 37 different attributes, as shown below:

Feature	Data Type and Description
Curricular units 1st sem (grade)	Integer - Grade average in the 1st semester (0 to 20).

Curricular units 1st sem (without evaluations)	Integer - Number of units without evaluations in 1st semester.
Curricular units 2nd sem (credited)	Integer - Number of credited units in 2nd semester.
Curricular units 2nd sem (enrolled)	Integer - Number of units enrolled in 2nd semester.
Curricular units 2nd sem (evaluations)	Integer - Number of evaluations in 2nd semester.
Curricular units 2nd sem (approved)	Integer - Number of units approved in 2nd semester.
Curricular units 2nd sem (grade)	Integer - Grade average in 2nd semester (0 to 20).
Curricular units 2nd sem (without evaluations)	Integer - Number of units without evaluations in 2nd semester.
Unemployment rate	Continuous - Unemployment rate (%).
Inflation rate	Continuous - Inflation rate (%).
GDP	Continuous - Gross Domestic Product.
Target	Categorical - Three-category classification (dropout, enrolled, graduate).
Marital Status	Integer - Marital Status (1 – single, 2 – married, 3 – widower, etc.).
Application mode	Integer - Various application mode codes.
Application order	Integer - Application order (0 - first choice; 9 - last choice).
Course	Integer - Various course codes.
Daytime/evening attendance	Integer - 1 – daytime, 0 - evening.
Previous qualification	Integer - Education Level and related qualifications.
Previous qualification (grade)	Continuous - Grade of previous qualification (0 to 200).
Nationality	Integer - Nationality codes.
Mother's qualification	Integer - Education Level of mother.

Father's qualification	Integer - Education Level of father.
Mother's occupation	Integer - Occupation of mother.
Father's occupation	Integer - Occupation of father.
Admission grade	Continuous - Admission grade (0 to 200).
Displaced	Integer - 1 - yes, 0 - no.
Educational special needs	Integer - 1 - yes, 0 - no.
Debtor	Integer - 1 - yes, 0 - no.
Tuition fees up to date	Integer - 1 - yes, 0 - no.
Gender	Integer - Gender (1 - male, 0 - female).
Scholarship holder	Integer - 1 - yes, 0 - no.
Age at enrollment	Integer - Age of student at enrollment.
International	Integer - 1 - yes, 0 - no.
Curricular units 1st sem (credited)	Integer - Number of units credited in 1st semester.
Curricular units 1st sem (enrolled)	Integer - Number of units enrolled in 1st semester.
Curricular units 1st sem (evaluations)	Integer - Number of evaluations in 1st semester.
Curricular units 1st sem (approved)	Integer - Number of units approved in 1st semester.

Table 1: List of features, data types, and descriptions.

The data refer to records of students enrolled between the academic years 2008/2009 to 2018/2019. These include data from 17 undergraduate degrees such as design, education, nursing, journalism, management, social service, and technologies.

Shown below is an image of the correlation heatmap of the dataset, where a brighter (more yellow) color indicates stronger positive correlation and a darker (more dark purple) color indicates stronger negative correlation.

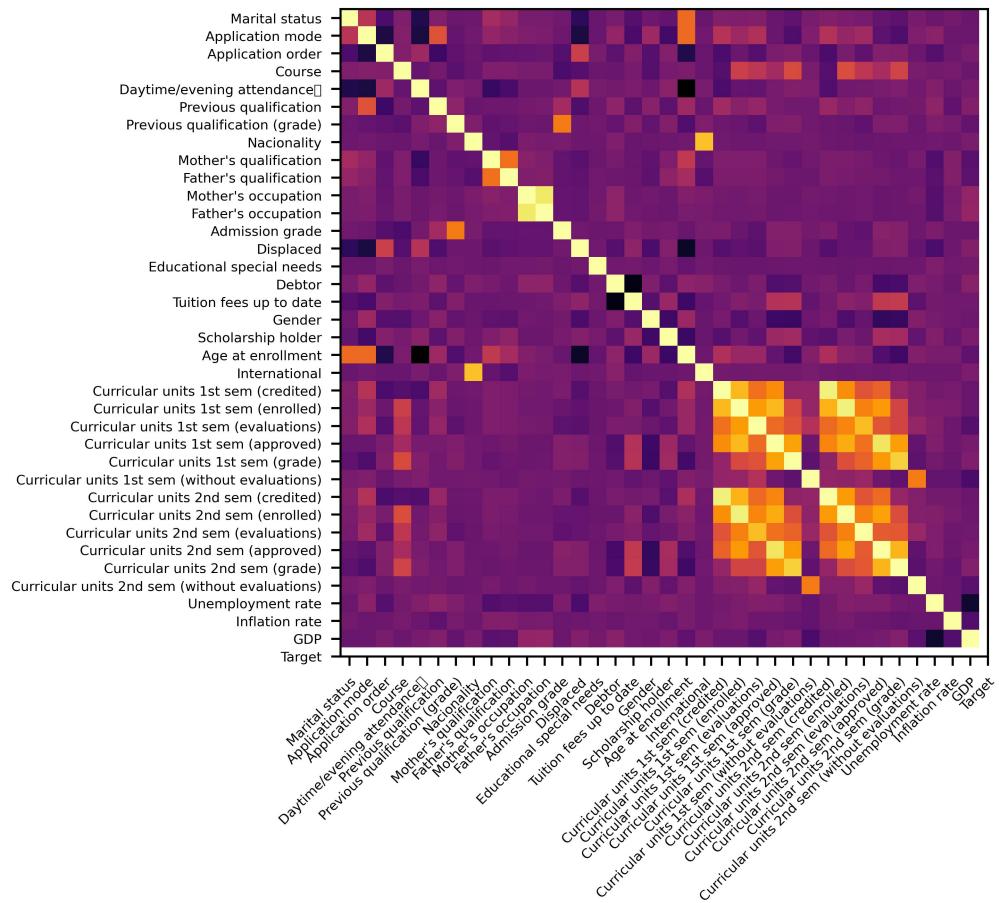


Figure 1: Correlation heatmap for the dataset. Brighter colors indicate stronger positive correlation. Dimmer colors indicate stronger negative correlation.

§3 Preprocessing

Before the data in the dataset can be used in models, we must prepare it for analysis by cleaning, transforming, and formatting it in a way that makes it suitable for the models. Raw data often is incomplete, has errors, and is inconsistent, and this can lead to biased models or incompatability with certain models. Data preprocessing is the process of dealing with the issues with raw data, and preparing it for use in models. This ensures that machine learning models can be used and are reliable.

The first error we ran into, was that the data used semi-colons as delimiters, when we need commas, since CSV stands for Comma Separated Values.

The second error we ran into, was that the data contained students who were currently enrolled. This is not useful for the exploration since we are trying to predict whether a student will drop out or not, and currently enrolled students would make no sense as a target value

The last error we ran into, was that the data varied a lot, and would benefit if it were normalized

§3.1 Normalization

There are large variations in the magnitude of data, with magnitudes of data ranging from nearly 10,000 to less than 1. Shown below is a box plot of the data: This difference in magnitude of data can confound the distance metrics, as the model may believe that the features with magnitude 10,000 are far more important than the features with magnitude 1.

In order to deal with this, the data was normalized based on the z-scores of each data point in each feature. For a given feature, let μ be the average of all data points of the feature and σ be the standard deviation of all data points of the feature. Then, the formula for the z-score of some arbitrary data point x in the feature is as follows:

$$z = \frac{x - \mu}{\sigma}.$$

Each data point was then converted into its z-score and a new dataset was created, with the following box plot. This yields a normalized range of values for the model

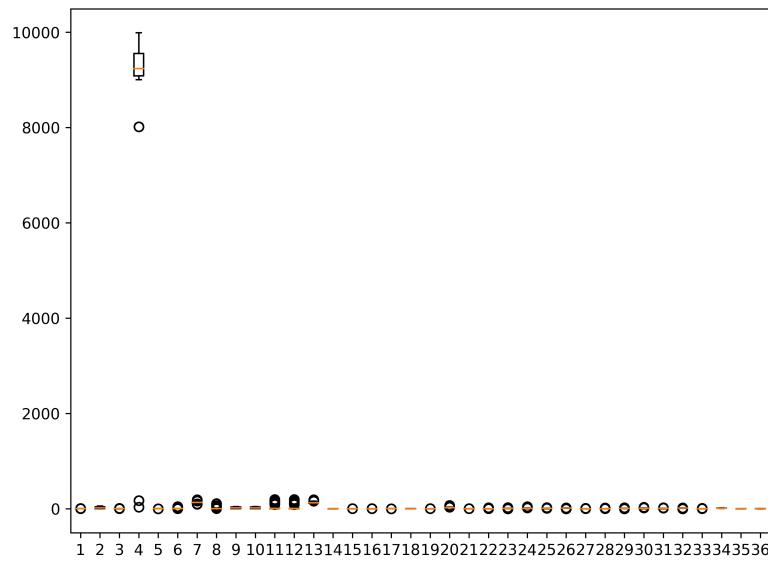


Figure 2: Boxplot of raw data prior to normalization.

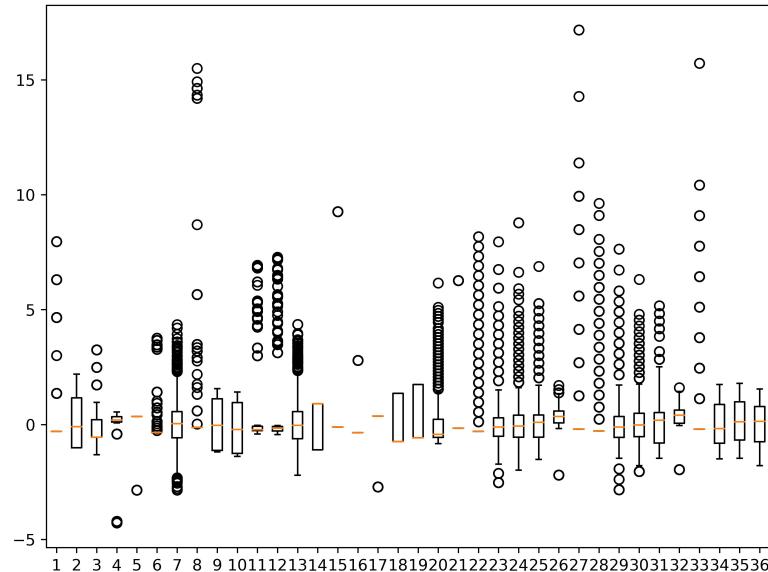


Figure 3: Boxplot of data after z-score normalization.

to be trained on. In addition to this, the normalized dataset has the benefit of each feature having mean $\mu = 0$ and standard deviation $\sigma = 1$.

§3.2 Recursive Feature Elimination

Recursive feature elimination (RFE) is a method used to find the optimal set of features within a dataset. This method works recursively by following these steps:

1. In the dataset with n features, the model is trained on each distinct group of $n - 1$ features.
2. The accuracy for each of the groups of $n - 1$ features is recorded.
3. The group of $n - 1$ features with the highest accuracy is kept and becomes the new dataset.
4. Steps 1-3 are repeated until there is 1 feature left in the dataset.
5. The set of features with the highest accuracy out of all the tests is determined.

§3.2.1 Benefits

Some benefits of this are that we can see which features we can ignore and remove them to improve overall accuracy all while using a relatively simple method and algorithm. These features may be confounding the model by providing unnecessary information.

The lower amount of features also means that there is less data for the model to be trained on, thereby decreasing training time and processing resources required.

§3.2.2 Limitations

Some limitations of this are that we may not be able to create the most accurate group of features using RFE. Removing two features with smaller accuracy increases might result in greater accuracy gains than removing the one feature that provided the greatest accuracy increase.

In order to test every possible group of features that can be used, the time complexity is of magnitude $O(2^N)$, which is far too time-consuming for feasible testing.

Name	Formula
L1	$\sum x_i - y_i $
L2	$\sqrt{\sum (x_i - y_i)^2}$
Cosine	$1 - \cos \theta = 1 - \frac{X \cdot Y}{\ X\ \ Y\ }$

Table 2: Distance metrics and formulae

Name	k	Accuracy
Cosine	16	0.8705
L1	10	0.8788
L2	6	0.8664

 Table 3: Best values of k for each of the models

§4 Model Training

We used three different distance metrics to compare the entries in the model.

For each of the distances, we tested k from 1 to 20, and determined which k had the greatest accuracy. The optimal values are below: Shown below is a graph of the k value and the accuracy of the model for predicting the training dataset and the testing dataset.

§5 Results

The following figures present the results of our models: Out of the three distance

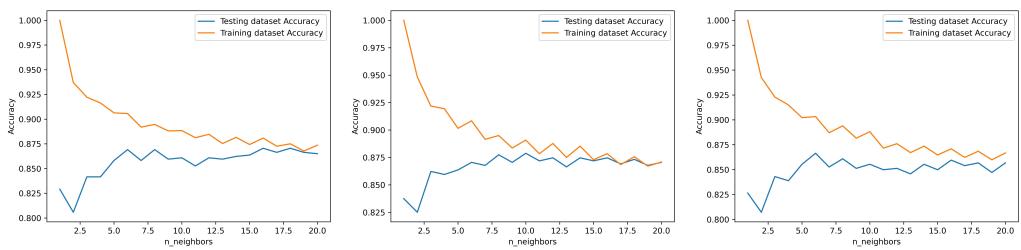


Figure 4: Graph of the k value and accuracy for the three models (Cosine, L1, L2 from left to right)

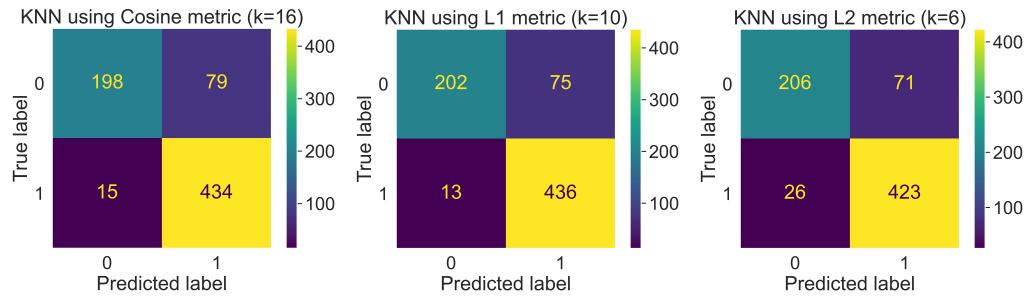


Figure 5: The confusion Matrices for Cosine, L1, and L2 distance metrics

Model	Accuracy	Recall	Precision
Cosine	0.8705	0.9666	0.8460
L1	0.8788	0.9710	0.8532
L2	0.8664	0.9421	0.8563

Table 4: Various model metrics compared to distance metric

metrics, L1 had the greatest accuracy with 87.88%. Shown below is a detailed breakdown of the metrics of the model.

The ROC metric is to check the fit of your model. It is a plot of the True Positive Rate against the False Positive Rate. The model performance is determined by looking at the area under the ROC curve. The model with the best ROC is L1.

Afterwards, RFE was used to determine the optimal set of features. Results are shown:

Metric	Accuracy	Number of features
Cosine	0.9045	12
L1	0.9105	14
L2	0.9105	13

Table 5: Detailed results after RFE

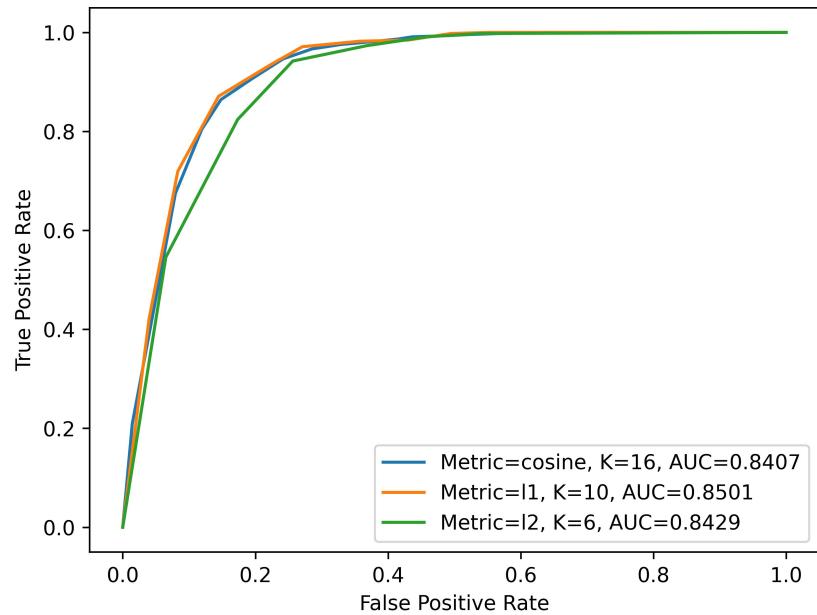


Figure 6: ROC curve with AUC

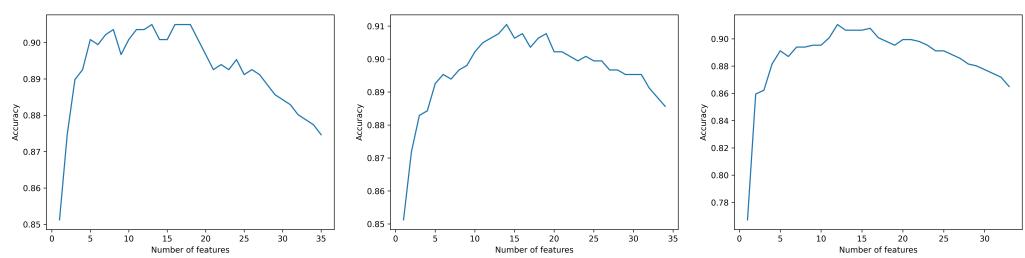


Figure 7: Graphs of Recursive Feature Elimination with the three different distance metrics

§6 Conclusion

Predicting whether people will drop out or not through our KNN-based model and preprocessing techniques can have a significant positive impact. By identifying at-risk students early, universities can provide targeted support before issues escalate to dropouts. This proactive approach increases student success rates and unlocks greater economic opportunities. Ultimately, this data-driven research highlights the potential of predictive analytics in building more equitable and supportive higher education systems.

In future studies about this topic, instead of sourcing it only from the Portuguese "Capacitação da Administração", we want to source the information from around the world including countries that are in different geological regions, cultural norms, and much more. This would let us make sure we aren't being influenced by these factors, and we could see which factors influence whether a student will drop out or not.

§7 Works Cited and Links

Hanson, Melanie. "College Dropout Rate [2023]: by Year + Demographics." Education Data Initiative, 29 October 2023, <https://educationdata.org/college-dropout-rates>. Accessed 11 March 2024.

"Predict Students' Dropout and Academic Success." UCI Machine Learning Repository, <https://archive.ics.uci.edu/dataset/697/predict+students+dropout+and+academic+success>. Accessed 11 March 2024.

Rose, Marley. "College Dropout Rate in the U.S." BestColleges.com, <https://www.bestcolleges.com/research/college-dropout-rate/>. Accessed 11 March 2024.

The GitHub repo <https://github.com/samson-rozansky/IB-Computer-Science-KNN>

§A Code

§A.1 Preprocessing

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pathlib
4 import pandas as pd
5
6 ROOT = pathlib.Path(__file__).parent.parent.resolve().joinpath("data")
7
8 data = pd.read_csv(ROOT.joinpath("data_raw.csv"))
9
10 data_normal = data.copy()
11 for column in data_normal.columns[:36]:
12     data_normal[column] = (data_normal[column] - data_normal[
13         column].mean()) / data_normal[column].std(ddof=0)
14
15 ret = []
16 for index, row in data_normal.iterrows():
17     if (row["Target"] != "Enrolled"):
18         if (row["Target"] == "Graduate"):
19             row["Target"] = 1
20         else:
21             row["Target"] = 0
22     ret.append(row)
23
24 df = pd.DataFrame(ret)
25 df.reset_index()
26 df.to_csv(ROOT.joinpath("data_normal.csv"), index = False)

```

§A.2 Boxplot Raw

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pathlib
4 import pandas as pd
5

```

```
6 ROOT = pathlib.Path(__file__).parent.parent.resolve().joinpath("data")
7 OUTPUT = pathlib.Path(__file__).parent.parent.resolve().joinpath("figs")
8
9 data = pd.read_csv(ROOT.joinpath("data_raw.csv"))
10
11 fig = plt.figure()
12 ax = fig.add_axes([0, 0, 1, 1])
13 bp = ax.boxplot(data.iloc[:, :36])
14 plt.savefig(OUTPUT.joinpath("boxplot_raw.jpg"), bbox_inches = "tight", transparent = True, dpi = 600)
```

§A.3 Boxplot Normalized

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pathlib
4 import pandas as pd
5
6 ROOT = pathlib.Path(__file__).parent.parent.resolve().joinpath("data")
7 OUTPUT = pathlib.Path(__file__).parent.parent.resolve().joinpath("figs")
8
9 DATA_FILE = ROOT.joinpath("data_normal.csv")
10
11 if not DATA_FILE.is_file():
12     import preprocess
13
14 df = pd.read_csv(DATA_FILE)
15
16 fig = plt.figure()
17 ax = fig.add_axes([0, 0, 1, 1])
18 bp = ax.boxplot(df.iloc[:, :36])
19 plt.savefig(OUTPUT.joinpath("boxplot_normal.jpg"), bbox_inches = "tight", transparent = True, dpi = 600)
```

§A.4 Correlation Heatmap

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import pathlib
5
6 ROOT = pathlib.Path(__file__).parent.parent.resolve().joinpath("data")
7 OUTPUT = pathlib.Path(__file__).parent.parent.resolve().joinpath("figs")
8
9 DATA_FILE = ROOT.joinpath("data_raw.csv")
10
11 df = pd.read_csv(DATA_FILE)
12
13 correlations = np.round(df.corr().to_numpy(), decimals=2)
14
15 #plt.style.use('https://github.com/dhaitz/matplotlib-stylesheets/
16 #               raw/master/pitayasmoothie-dark.mplstyle')
16 #plt.style.use(['dark_background'])
17
18 plt.rcParams.update({'font.size': 5})
19
20 fig, ax = plt.subplots()
21 im = ax.imshow(correlations, cmap = 'inferno')
22
23 # Show all ticks and label them with the respective list entries
24 ax.set_xticks(np.arange(len(list(df.columns))), labels=list(df.
25 columns))
25 ax.set_yticks(np.arange(len(list(df.columns))), labels=list(df.
26 columns))
26
27 # Rotate the tick labels and set their alignment.
28 plt.setp(ax.get_xticklabels(), rotation=45, ha="right",
29           rotation_mode="anchor")
30
31 # Loop over data dimensions and create text annotations.
31 """
32 for i in range(len(list(df.columns))):
```

```
33     for j in range(len(list(df.columns))):  
34         text = ax.text(j, i, correlations[i, j], ha="center", va="center", color="aqua")  
35     """  
36  
37 fig.tight_layout()  
38  
39 plt.savefig(OUTPUT.joinpath("heatmap.jpg"), bbox_inches = "tight",  
            transparent = True, dpi = 600)
```

§A.5 Hyperparameter Optimization

```
1 from sklearn.neighbors import KNeighborsClassifier  
2 from sklearn.model_selection import train_test_split  
3  
4 import os  
5 import numpy as np  
6 import matplotlib.pyplot as plt  
7 import pathlib  
8 import pandas as pd  
9 import seaborn as sns  
10  
11 from sklearn import metrics  
12 from sklearn.metrics import precision_score, accuracy_score,  
    confusion_matrix, ConfusionMatrixDisplay  
13  
14 ROOT = pathlib.Path(__file__).parent.parent.resolve().joinpath("data")  
15 OUTPUT = pathlib.Path(__file__).parent.parent.resolve().joinpath("figs\\distance")  
16  
17 DATA_FILE = ROOT.joinpath("data_normal.csv")  
18  
19 if not DATA_FILE.is_file():  
20     import preprocess  
21  
22 data = pd.read_csv(DATA_FILE)  
23  
24 # Create feature and target arrays
```



```

60     open(write_file, 'x')
61     with open(write_file, 'w') as f:
62         f.write('Evaluation Metrics for ' + name[m] + '\n')
63         f.write('Accuracy: ' + str(metrics.accuracy_score(y_test,
64             predictions)) + '\n')
65         f.write('Recall: ' + str(metrics.recall_score(y_test,
66             predictions)) + '\n')
67         f.write('F1 Score: ' + str(metrics.f1_score(y_test,
68             predictions)) + '\n')
69         f.write('Precision: ' + str(metrics.precision_score(y_test,
70             predictions)) + '\n')
71
72     # Confusion Matrix
73     cm = confusion_matrix(y_test, predictions, labels = None)
74     #sns.set(font_scale=2)
75     disp = ConfusionMatrixDisplay(confusion_matrix = cm,
76         display_labels = None)
77     disp.plot()
78     plt.grid(False)
79     plt.title("KNN using " + name[m] + " metric (k=" + str(best_k)
80     + ")")
81     #plt.savefig(OUTPUT.joinpath(name[m] + "_confusion_matrix.jpg")
82     #, bbox_inches = "tight", transparent = True, dpi = 600)
83     plt.close()
84
85     # graph of k vs accuracy
86     plt.plot(neighbors, test_accuracy, label = name[m] + ' Testing
87     dataset Accuracy')
88     plt.plot(neighbors, train_accuracy, label = name[m] + '
89     Training dataset Accuracy')
90
91     plt.legend()
92     plt.xlabel('n_neighbors')
93     plt.ylabel('Accuracy')
94     #plt.savefig(OUTPUT.joinpath(m + "_neighbors_vs_accuracy.jpg")
95     #, bbox_inches = "tight", transparent = True, dpi = 600)
96     #plt.close()
97
98     plt.savefig(OUTPUT.joinpath("neighbors_vs_accuracy.jpg"),

```

```
bbox_inches = "tight", transparent = True, dpi = 600)
```

§A.6 ROC Curve

```
1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.model_selection import train_test_split
3
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pathlib
7 import pandas as pd
8 import seaborn as sns
9
10 from sklearn import metrics
11 from sklearn.metrics import precision_score, accuracy_score,
12     confusion_matrix, ConfusionMatrixDisplay
13
14 ROOT = pathlib.Path(__file__).parent.parent.resolve().joinpath("data")
15
16 OUTPUT = pathlib.Path(__file__).parent.parent.resolve().joinpath("figs")
17
18 DATA_FILE = ROOT.joinpath("data_normal.csv")
19
20
21 data = pd.read_csv(DATA_FILE)
22
23 # Create feature and target arrays
24 X = data.drop(columns = ['Target'])
25 y = data['Target']
26
27 # Split into training and test set
28 X_train, X_test, y_train, y_test = train_test_split(X, y,
29     test_size = 0.2, random_state=42)
30
31 metrics_list = ['cosine', 'l1', 'l2']
32 k = [16, 10, 6]
```

```

32
33 for i in range(3):
34     model = KNeighborsClassifier(n_neighbors=k[i], metric=
35         metrics_list[i])
36     model.fit(X_train, y_train)
37     y_pred = model.predict(X_test)
38     y_pred_proba = model.predict_proba(X_test)[:, 1]
39     fpr, tpr, _ = metrics.roc_curve(y_test, y_pred_proba)
40     auc = round(metrics.roc_auc_score(y_test, y_pred), 4)
41
42     plt.plot(fpr, tpr, label = "Metric=" + metrics_list[i] + ", K=" +
43         " " + str(k[i]) + ", AUC=" + str(auc))
44
45 plt.ylabel('True Positive Rate')
46 plt.xlabel('False Positive Rate')
47 plt.legend()
48 plt.savefig(OUTPUT.joinpath("roc_curve.jpg"), bbox_inches = "tight"
49             , transparent = True, dpi = 600)

```

§A.7 Recursive Feature Elimination

```

1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.model_selection import train_test_split
3
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pathlib
7 import pandas as pd
8 import seaborn as sns
9
10 from sklearn import metrics
11 from sklearn.metrics import precision_score, accuracy_score,
12     confusion_matrix, ConfusionMatrixDisplay
13
14 ROOT = pathlib.Path(__file__).parent.parent.resolve().joinpath("data")
15 OUTPUT = pathlib.Path(__file__).parent.parent.resolve().joinpath("figs\\rfe")

```

```
16 DATA_FILE = ROOT.joinpath("data_normal.csv")
17
18 if not DATA_FILE.is_file():
19     import preprocess
20
21 data = pd.read_csv(DATA_FILE)
22
23 features = []
24 accuracies = []
25
26 def RFE(df, m, k):
27     best_acc = 0
28     feature = 'temp'
29     for column in df:
30         if (column == 'Target'):
31             continue
32         temp = df.copy()
33         temp.drop(columns = [column], inplace = True)
34         X = temp.drop(columns = ['Target'])
35         y = temp['Target']
36
37         X_train, X_test, y_train, y_test = train_test_split(X, y,
38             test_size = 0.2, random_state=42)
39         model = KNeighborsClassifier(n_neighbors=k, metric=m)
40         model.fit(X_train, y_train)
41         y_pred = model.predict(X_test)
42         cur_acc = metrics.accuracy_score(y_test, y_pred)
43         if (cur_acc > best_acc):
44             feature = column
45             best_acc = cur_acc
46
47         df.drop(columns = [feature], inplace = True)
48         features.append(feature)
49         accuracies.append(best_acc)
50
51         if (len(df.columns) > 2):
52             print("Number eliminated: " + str(len(features)))
53             print("Number remaining: " + str(len(df.columns)) + "\n")
```

```
53     RFE(df.copy(), m, k)
54
55     return
56
57
58 metrics_list = ['cosine', 'l1', 'l2']
59 k = [16, 10, 6]
60
61 for i in range(3):
62     RFE(data, metrics_list[i], k[i])
63     write_file = OUTPUT.joinpath(metrics_list[i] + "_RFE.txt")
64     if not write_file.is_file():
65         open(write_file, 'x')
66         open(write_file, 'w').close()
67
68     for j in range(len(features)):
69         with open(write_file, 'a') as f:
70             f.write("Removed: " + features[j] + "\nAccuracy: " +
71 str(accuracies[j]) + "\n\n")
72
73     plt.plot(np.arange(len(features), 0, -1), accuracies)
74     plt.xlabel('Number of features')
75     plt.ylabel('Accuracy')
76
77     plt.savefig(OUTPUT.joinpath(metrics_list[i] + "_RFE.jpg"),
78                 bbox_inches = "tight", transparent = True, dpi = 600)
79     plt.close()
80
81     features = []
82     accuracies = []
```