

Machine Learning Assignment 2

Supervised Learning Classification using K-Nearest Neighbor with KDTree Algorithm

by:

Muhammad Furqon Fahlevi 1301194214

Muhammad Ilham Mubarak 1301194276

A. Problem Formulation

Classification refers to a predictive modeling problem where a class label is predicted for a given example of input data. In this assignment, we are going to predict whether customers are interested in buying a new vehicle or not based on customer data at the dealer.

B. Data Exploration and Preparation

1. Import Library

In this step, we import all of the necessary library to complete the assignment

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
from scipy import stats
from math import sqrt
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.decomposition import PCA
from collections import Counter
import heapq
```

2. Import Datasets

i. Data Train

```
[ ] # Import dataset (kendaraan_train.csv) into new dataframe named data_train
data_train = pd.read_csv("https://raw.githubusercontent.com/furqonfahlevi/classification-ml/main/kendaraan_train.csv")
data_train.head()
```

	id	Jenis_Kelamin	Umur	SIM	Kode_Daerah	Sudah_Asuransi	Umur_Kendaraan	Kendaraan_Rusak	Premi	Kanal_Penjualan	Lama_Berlangganan	Tertarik
0	1	Wanita	30.0	1.0	33.0	1.0	< 1 Tahun	Tidak	28029.0	152.0	97.0	0
1	2	Pria	48.0	1.0	39.0	0.0	> 2 Tahun	Pernah	25800.0	29.0	158.0	0
2	3	NaN	21.0	1.0	46.0	1.0	< 1 Tahun	Tidak	32733.0	160.0	119.0	0
3	4	Wanita	58.0	1.0	48.0	0.0	1-2 Tahun	Tidak	2630.0	124.0	63.0	0
4	5	Pria	50.0	1.0	35.0	0.0	> 2 Tahun	NaN	34857.0	88.0	194.0	0

ii. Data Test

```
[ ] # Import dataset (kendaraan_test.csv) into new dataframe named data_test
data_test = pd.read_csv("https://raw.githubusercontent.com/furqonfahlevi/classification-ml/main/kendaraan_test.csv")
data_test.head()
```

	Jenis_Kelamin	Umur	SIM	Kode_Daerah	Sudah_Asuransi	Umur_Kendaraan	Kendaraan_Rusak	Premi	Kanal_Penjualan	Lama_Berlangganan	Tertarik
0	Wanita	49	1	8	0	1-2 Tahun	Pernah	46963	26	145	0
1	Pria	22	1	47	1	< 1 Tahun	Tidak	39624	152	241	0
2	Pria	24	1	28	1	< 1 Tahun	Tidak	110479	152	62	0
3	Pria	46	1	8	1	1-2 Tahun	Tidak	36266	124	34	0
4	Pria	35	1	23	0	1-2 Tahun	Pernah	26963	152	229	0

iii. Data Shape

Shape function return the shape of an array. Shape of an array is a tuple with the number of elements per axis (dimension). In our case, data train has 285831 lines and 12 columns, data test has 47639 lines and 11 columns.

```
[ ] print("Data Train Shape", data_train.shape)
    print("Data Test Shape", data_test.shape)

Data Train Shape (285831, 12)
Data Test Shape (47639, 11)
```

3. Pre-processing for Datasets

i. Drop Unnecessary Column

Drop unnecessary column											
<pre>[] del data_train['id'] data_train.head()</pre>											
	Jenis_Kelamin	Umur	SIM	Kode_Daerah	Sudah_Asuransi	Umur_Kendaraan	Kendaraan_Rusak	Premi	Kanal_Penjualan	Lama_Berlangganan	Tertarik
0	Wanita	30.0	1.0	33.0	1.0	< 1 Tahun	Tidak	28029.0	152.0	97.0	0
1	Pria	48.0	1.0	39.0	0.0	> 2 Tahun	Pemah	25800.0	29.0	158.0	0
2	NaN	21.0	1.0	46.0	1.0	< 1 Tahun	Tidak	32733.0	160.0	119.0	0
3	Wanita	58.0	1.0	48.0	0.0	1-2 Tahun	Tidak	2630.0	124.0	63.0	0
4	Pria	50.0	1.0	35.0	0.0	> 2 Tahun	NaN	34857.0	88.0	194.0	0

ii. Duplicated Data

Duplicated data leads to inaccurate results, we tried to figure it out with dropping them.

Duplicated data

```
[ ] # Check duplicate data
    print("Duplicated data:", data_train.duplicated().sum())
    print("Duplicated data:", data_test.duplicated().sum())
    data_train.drop_duplicates(inplace=True)
    data_test.drop_duplicates(inplace=True)
```

```
Duplicated data: 169
Duplicated data: 3
```

```
[ ] print("Duplicated data:", data_train.duplicated().sum())
    print("Duplicated data:", data_test.duplicated().sum())
```

```
Duplicated data: 0
Duplicated data: 0
```

iii. Encoder Data

Categorical data are variables that contain label values rather than numeric values. In this case, we tried to convert a numerical variable to an ordinal variable. Each categorical column assigned with an integer value.

Encoder Data

```
[8] categorical = ["Jenis_Kelamin", "Kendaraan_Rusak", "Umur_Kendaraan"]
```

```
[9] data_train = pd.get_dummies(data_train, columns=categorical)  
data_train.sample(3)
```

	Umur	Kode_Daerah	Sudah_Asuransi	Premi	Kanal_Penjualan	Lama_Berlangganan	Tertarik	Jenis_Kelamin_Pria	Jenis_Kelamin_Wanita	Kendaraan_Rusak_Pernah	Kendaraan_Rusak_Tidak	Umur_Kendaraan_1-2 Tahun	Umur_Kendaraan_< 1 Tahun	Umur_Kendaraan_> 2 Tahun
47489	26	1	46	0	2620	196	288	0	1	0	0	0	1	0
22869	31	1	8	0	61098	152	272	1	1	0	1	0	0	1
13426	30	1	21	1	32211	152	107	0	0	1	0	1	0	1

```
[10] data_test = pd.get_dummies(data_test, columns=categorical)  
data_test.sample(3)
```

	Umur	Kode_Daerah	Sudah_Asuransi	Premi	Kanal_Penjualan	Lama_Berlangganan	Tertarik	Jenis_Kelamin_Pria	Jenis_Kelamin_Wanita	Kendaraan_Rusak_Pernah	Kendaraan_Rusak_Tidak	Umur_Kendaraan_1-2 Tahun	Umur_Kendaraan_< 1 Tahun	Umur_Kendaraan_> 2 Tahun
47489	26	1	46	0	2620	196	288	0	1	0	0	0	1	0
22869	31	1	8	0	61098	152	272	1	1	0	1	0	0	1
13426	30	1	21	1	32211	152	107	0	0	1	0	1	0	1

iv. Check NaN

We check whether there is a NaN data or not in every column. After checked, we tried to fill every NaN value with 0, mean and median. The data that are filled with 0 because the data only consist of 0 and 1.

Check NaN

```
[11] data_train.isna().sum()
```

```
Umur          14199  
SIM           14484  
Kode_Daerah   14291  
Sudah_Asuransi 14229  
Premi         14510  
Kanal_Penjualan 14297  
Lama_Berlangganan 13926  
Tertarik      0  
Jenis_Kelamin_Pria 0  
Jenis_Kelamin_Wanita 0  
Kendaraan_Rusak_Pernah 0  
Kendaraan_Rusak_Tidak 0  
Umur_Kendaraan_1-2 Tahun 0  
Umur_Kendaraan_< 1 Tahun 0  
Umur_Kendaraan_> 2 Tahun 0  
dtype: int64
```

```
[12] data_train["SIM"].fillna(0, inplace=True)  
data_train["Sudah_Asuransi"].fillna(0, inplace=True)  
data_train["Umur"].fillna(data_train["Umur"].mean(), inplace=True)  
data_train["Premi"].fillna(data_train["Premi"].median(), inplace=True)  
data_train["Kanal_Penjualan"].fillna(data_train["Kanal_Penjualan"].mean(), inplace=True)  
data_train["Lama_Berlangganan"].fillna(data_train["Lama_Berlangganan"].mean(), inplace=True)  
data_train["Kode_Daerah"].fillna(data_train["Kode_Daerah"].mean(), inplace=True)
```

```
[13] data_train.isna().sum()
```

```

Umur      0
SIM       0
Kode_Daerah  0
Sudah_Asuransi  0
Premi     0
Kanal_Penjualan  0
Lama_Berlangganan  0
Tertarik  0
Jenis_Kelamin_Pria  0
Jenis_Kelamin_Wanita  0
Kendaraan_Rusak_Pernah  0
Kendaraan_Rusak_Tidak  0
Umur_Kendaraan_1-2 Tahun  0
Umur_Kendaraan_< 1 Tahun  0
Umur_Kendaraan_> 2 Tahun  0
dtype: int64

```

```
[14] data_test.isna().sum()
```

```

Umur      0
SIM       0
Kode_Daerah  0
Sudah_Asuransi  0
Premi     0
Kanal_Penjualan  0
Lama_Berlangganan  0
Tertarik  0
Jenis_Kelamin_Pria  0
Jenis_Kelamin_Wanita  0
Kendaraan_Rusak_Pernah  0
Kendaraan_Rusak_Tidak  0
Umur_Kendaraan_1-2 Tahun  0
Umur_Kendaraan_< 1 Tahun  0
Umur_Kendaraan_> 2 Tahun  0
dtype: int64

```

v. Data Standardization

Standardization is scaling technique where the values are centered around the mean with a unit standard deviation. This means that the mean of the attribute becomes zero and the resultant distribution has a unit standard deviation.

Data Standarization

```

[40] scaler = MinMaxScaler()
      columns = ["Umur", "Kode_Daerah", "Premi", "Kanal_Penjualan", "Lama_Berlangganan"]
      data_train[columns] = scaler.fit_transform(data_train[columns])
      data_train.sample(3)

```

	Umur	Kode_Daerah	Sudah_Asuransi	Premi	Kanal_Penjualan	Lama_Berlangganan	Tertarik	Jenis_Kelamin_Pria	Jenis_Kelamin_Wanita	Kendaraan_Rusak_Pernah	Kendaraan_Rusak_Tidak	Umur_Kendaraan_1-2 Tahun	Umur_Kendaraan_< 1 Tahun	Umur_Kendaraan_> 2 Tahun
238578	0.305846	1.0	0.558462	0.0	0.363507	0.758059	0.754325	0	1	1	0	1	0	0
198917	0.303077	1.0	0.230769	0.0	0.046341	0.154321	0.062744	1	1	0	1	0	1	0
381921	0.015385	1.0	0.113846	0.0	0.014780	0.381481	0.100346	0	0	1	1	0	0	1

```
[41] data_test[colums] = scaler.fit_transform(data_test[colums])
data_test.sample(3)
```

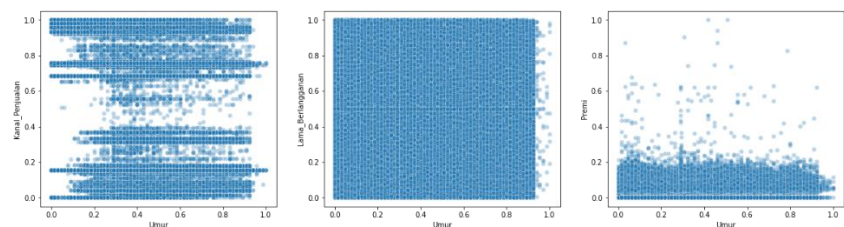
Umur	Salah_Jumlah	Salah_Jumlah	Premi	Kanal_Penjualan	Lama_Berlangganan	Tertarik	Salah_Kelamin_Pria	Salah_Kelamin_Kanita	Kendaraan_Kecil_Personal	Kendaraan_Kecil_Tidak	Umur_Kendaraan_1_2 Tahun	Umur_Kendaraan_3_5 Tahun	Umur_Kendaraan_6_2 Tahun
54902	0.401538	1	0.557602	0	0.960915	0.154021	0.416605	0	1	1	0	1	0
27128	0.500000	1	0.208482	0	0.852944	0.060247	0.164487	0	1	1	0	1	0
18272	0.420109	1	0.153846	0	0.916827	0.726259	0.217363	0	1	1	0	1	0

vi. Check Outliers

Outlier is an observation point that is distant from other observations. We tried to check continuous data whether there is an outlier or not.

Check outliers

```
[17] fig, axes = plt.subplots(nrows=1, ncols=3, figsize=(20, 5))
sns.scatterplot(data=data_train, x="Umur", y="Kanal_Penjualan", ax=axes[0], alpha=0.3)
sns.scatterplot(data=data_train, x="Umur", y="Lama_Berlangganan", ax=axes[1], alpha=0.3)
sns.scatterplot(data=data_train, x="Umur", y="Premi", ax=axes[2], alpha=0.3)
plt.show()
```

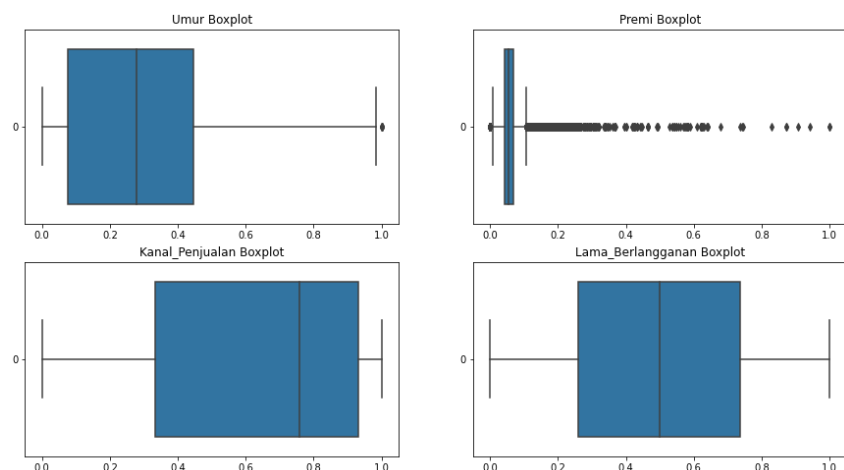


```
[18] fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(15, 8))

axes[0, 0].title.set_text("Umur Boxplot")
axes[0, 1].title.set_text("Premi Boxplot")
axes[1, 0].title.set_text("Kanal_Penjualan Boxplot")
axes[1, 1].title.set_text("Lama_Berlangganan Boxplot")

sns.boxplot(data=data_train["Umur"], ax=axes[0, 0], orient="horizontal")
sns.boxplot(data=data_train["Premi"], ax=axes[0, 1], orient="horizontal")
sns.boxplot(data=data_train["Kanal_Penjualan"], ax=axes[1, 0], orient="horizontal")
sns.boxplot(data=data_train["Lama_Berlangganan"], ax=axes[1, 1], orient="horizontal")
plt.show()
```

As we can see from the boxplot below, “Umur” and “Premi” column have outliers.



vii. Z-Score Method

Z-score is measured in terms of standard deviations from the mean. Any z-score greater than 3 or less than -3 is considered to be an

outlier (threshold). In this case, we tried to find z-score from the chosen columns.

Z-Score method

```
[42] columns = ['Kanal_Penjualan', 'Premi', 'Umur', 'Lama_Berlangganan']
      z = np.abs(stats.zscore(data_train[columns]))
      threshold = 3

      data_train_new = data_train[(z < threshold).all(axis=1)]
      data_train_new.sample(3)
```

Umur	Kode_Sarah	Sudah_Asuranci	Premi	Kanal_Penjualan	Lama_Berlangganan	Tertarik	jenis_kelamin_Pria	jenis_kelamin_Wanita	Kendaraan_Rusak_Perusak	Kendaraan_Rusak_Tidak	Umur_Kendaraan_1_2 Tahun	Umur_Kendaraan_3_5 Tahun	Umur_Kendaraan_6_9 Tahun	Umur_Kendaraan_10 Tahun
198673	0.546154	1.0	0.507632	1.0	0.357985	0.932099	0.495247	0	1	0	0	1	0	1
138118	0.503231	1.0	0.533462	0.0	0.352376	0.154321	0.086886	1	1	0	1	0	0	0
91378	0.288958	1.0	0.546154	0.0	0.386491	0.944444	0.540203	1	0	1	1	0	1	0

```
[43] columns = ['Kanal_Penjualan', 'Premi', 'Umur', 'Lama_Berlangganan']
      z = np.abs(stats.zscore(data_test[columns]))
      threshold = 3

      data_test_new = data_test[(z < threshold).all(axis=1)]
      data_test_new.sample(3)
```

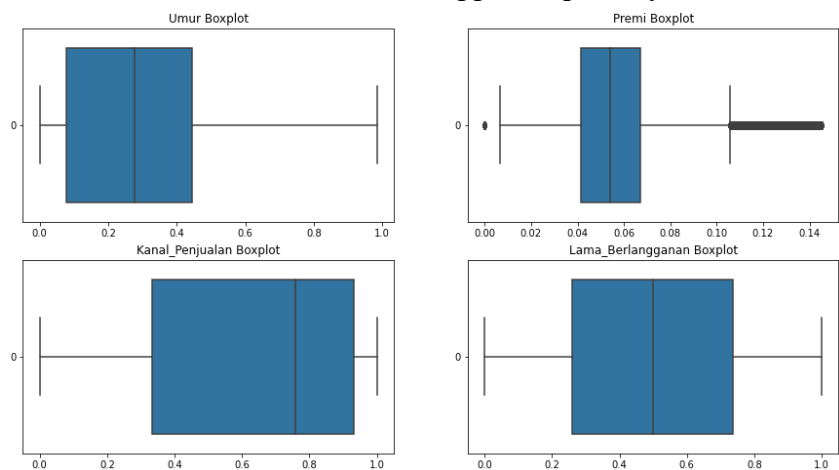
Umur	Kode_Sarah	Sudah_Asuranci	Premi	Kanal_Penjualan	Lama_Berlangganan	Tertarik	jenis_kelamin_Pria	jenis_kelamin_Wanita	Kendaraan_Rusak_Perusak	Kendaraan_Rusak_Tidak	Umur_Kendaraan_1_2 Tahun	Umur_Kendaraan_3_5 Tahun	Umur_Kendaraan_6_9 Tahun	Umur_Kendaraan_10 Tahun
94073	0.400000	1	0.368015	1	0.000000	0.154321	0.902699	0	1	0	0	1	1	0
11679	0.522677	1	0.133846	0	0.000000	0.154321	0.301038	0	0	1	1	0	1	0
19872	0.600000	1	0.538462	0	0.040564	0.154321	0.433986	0	1	0	1	0	1	0

```
[21] fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(15, 8))

      axes[0, 0].title.set_text("Umur Boxplot")
      axes[0, 1].title.set_text("Premi Boxplot")
      axes[1, 0].title.set_text("Kanal_Penjualan Boxplot")
      axes[1, 1].title.set_text("Lama_Berlangganan Boxplot")

      sns.boxplot(data=data_train_new["Umur"], ax=axes[0, 0], orient="horizontal")
      sns.boxplot(data=data_train_new["Premi"], ax=axes[0, 1], orient="horizontal")
      sns.boxplot(data=data_train_new["Kanal_Penjualan"], ax=axes[1, 0], orient="horizontal")
      sns.boxplot(data=data_train_new["Lama_Berlangganan"], ax=axes[1, 1], orient="horizontal")
      plt.show()
```

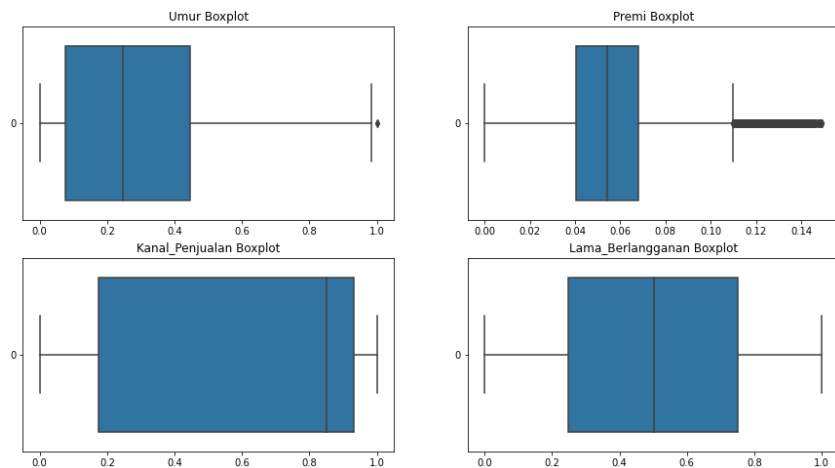
After doing the z-score method, we can see from the boxplot below that most of the outliers almost disappear especially the “Premi”.



```
[22] fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(15, 8))

axes[0, 0].title.set_text("Umur Boxplot")
axes[0, 1].title.set_text("Premi Boxplot")
axes[1, 0].title.set_text("Kanal_Penjualan Boxplot")
axes[1, 1].title.set_text("Lama_Berlangganan Boxplot")

sns.boxplot(data=data_test_new["Umur"], ax=axes[0, 0], orient="horizontal")
sns.boxplot(data=data_test_new["Premi"], ax=axes[0, 1], orient="horizontal")
sns.boxplot(data=data_test_new["Kanal_Penjualan"], ax=axes[1, 0], orient="horizontal")
sns.boxplot(data=data_test_new["Lama_Berlangganan"], ax=axes[1, 1], orient="horizontal")
plt.show()
```



viii. Dimensional Reduction (PCA)

Dimension reduction eliminates noisy data dimensions and thus and improves accuracy in classification, in addition to reduced computational cost.

Dimensial Reduction (PCA)

```
[ ] X = data_train_new.drop(columns=["Tertarik"])
    y = data_train_new[["Tertarik"]]

X_test_data = data_test_new.drop(columns=["Tertarik"])
y_test_data = data_test_new[["Tertarik"]]
```

```
[ ] pca = PCA(n_components=2)
    principal_df = pd.DataFrame(data = pca.fit_transform(X))
    principal_test = pd.DataFrame(data = pca.fit_transform(X_test_data))
```

```
[ ] print(f"Data Train New Shape: {principal_df.shape}")
    print(f"Data Test New Shape: {principal_test.shape}")
```

```
Data Train New Shape: (284032, 2)
Data Test New Shape: (47383, 2)
```

We can see from the image above; those datasets have been reduced from 15 columns into 2 columns. This are very important because PCA creates a smaller data in volume and has the same analytical results as the original representation.

C. Modeling

1. Data Splitting

In our case, we tried to split the data with Pareto Principle, or it also called 80/20 rule, 80% of effects come from 20% of causes.

Data Splitting

```
[ ] X = data_train_new.drop(columns=["Tertarik"])
    y = data_train_new[["Tertarik"]]

[ ] X_train, X_test, y_train, y_test = train_test_split(principal_df, y, test_size = 0.2, random_state = 42)

[ ] print(f"X_train shape: {X_train.shape} | X_test shape: {X_test.shape}")
    print(f"y_train shape: {y_train.shape} | y_test shape: {y_test.shape}")

X_train shape: (227225, 2) | X_test shape: (56807, 2)
y_train shape: (227225, 1) | y_test shape: (56807, 1)
```

2. Evaluation Metrics

```
1 class Metrics:
2     def __init__(self, prediction, test):
3
4         if isinstance(test, pd.DataFrame):
5             test = test.values
6
7         if isinstance(prediction, pd.DataFrame):
8             prediction = prediction.values
9
10        self.prediction = np.array(prediction)
11        self.test = np.array(test).reshape(-1)
12
13        unique = set(self.prediction)
14        matrix = [list() for x in range(len(unique))]
15        for i in range(len(unique)):
16            matrix[i] = [0 for x in range(len(unique))]
17        lookup = dict()
18        for i, value in enumerate(unique):
19            lookup[value] = i
20        for i in range(len(self.prediction)):
21            x = lookup[self.prediction[i]]
22            y = lookup[self.test[i]]
23            matrix[y][x] += 1
24        matrix[0][1], matrix[1][0] = matrix[1][0], matrix[0][1]
25
26        self.tn = matrix[0][0]
27        self.tp = matrix[1][1]
28        self.fn = matrix[1][0]
29        self.fp = matrix[0][1]
30        self.matrix = matrix
31
32    def accuracy_metric(self):
33        correct = 0
34        for i in range(len(self.prediction)):
35            if self.prediction[i] == self.test[i]:
36                correct += 1
37        return correct / float(len(self.prediction)) * 100.0
38
39    def confusion_matrix(self):
40        return self.matrix
41
42    def f1_score(self):
43        p = self.precision_score()
44        r = self.recall_score()
45        return 2 * ((p * r) / (p + r)) / 100
46
47    def precision_score(self):
48        return (self.tp / ((self.tp) + (self.fp))) * 100
49
50    def recall_score(self):
51        return (self.tp / ((self.tp) + (self.fn))) * 100
52
53    def visualize_confusion_matrix(self):
54        cf = self.confusion_matrix()
55        sns.heatmap(cf,
56                    xticklabels=["Negative (0)", "Positive (1)"],
57                    yticklabels=["Negative (0)", "Positive (1)"],
58                    annot=True, fmt='d',
59                    cmap="YlGnBu", linewidths=10)
60        plt.show()
61
62    def show(self):
63        self.visualize_confusion_matrix()
64        print("Accuracy with K-NN: {:.2f}%".format(self.accuracy_metric()))
65        print("Precision with K-NN: {:.2f}%".format(self.precision_score()))
66        print("Recall with K-NN: {:.2f}%".format(self.recall_score()))
67        print("F1-Score with K-NN: {:.2f}%".format(self.f1_score()))
```


i. Confusion Matrix

Confusion matrix is a summary of prediction results on a classification problem. It shows the ways in which your classification model is confused when it makes predictions. We create ours based on image below:

		Predicted classes	
		Negative 0	Positive 1
Actual classes	Negative 0	TN	FP
	Positive 1	FN	TP

Therefore, we applied into our code:

```
1 def __init__(self, prediction, test):
2
3     if isinstance(test, pd.DataFrame):
4         test = test.values
5
6     if isinstance(prediction, pd.DataFrame):
7         prediction = prediction.values
8
9     self.prediction = np.array(prediction)
10    self.test = np.array(test).reshape(-1)
11
12    unique = set(self.prediction)
13    matrix = [list() for x in range(len(unique))]
14    for i in range(len(unique)):
15        matrix[i] = [0 for x in range(len(unique))]
16    lookup = dict()
17    for i, value in enumerate(unique):
18        lookup[value] = i
19    for i in range(len(self.prediction)):
20        x = lookup[self.prediction[i]]
21        y = lookup[self.test[i]]
22        matrix[y][x] += 1
23    matrix[0][1], matrix[1][0] = matrix[1][0], matrix[0][1]
24
25    self.tn = matrix[0][0]
26    self.tp = matrix[1][1]
27    self.fn = matrix[1][0]
28    self.fp = matrix[0][1]
29    self.matrix = matrix
```

ii. Accuracy Score

Classification accuracy is the ratio of number of correct predictions to the total number of input samples. We need to find is predict

data same as test data. If they are equal, they will be assigned into “correct” variable and then we calculate based on formula below.

$$\text{Accuracy} = \frac{\text{number of correct predictions}}{\text{total number of predictions made}}$$

We applied this formula into our case to find the classification accuracy.

```
1 def accuracy_metric(self):
2     correct = 0
3     for i in range(len(self.prediction)):
4         if self.prediction[i] == self.test[i]:
5             correct += 1
6     return correct / float(len(self.prediction)) * 100.0
```

iii. Precision Score

Precision is the ratio of correctly predicted positive observations to the total predicted positive observations, with formula:

$$\text{Precision} = \frac{\text{True Positive}}{(\text{True Positive} + \text{False Positive})}$$

```
1 def precision_score(self):
2     return (self.tp / ((self.tp) + (self.fp))) * 100
```

iv. Recall Score

Recall is the ratio of correctly predicted positive observations to all observations in actual class, with formula:

$$\text{Recall} = \frac{\text{True Positive}}{(\text{True Positive} + \text{False Negative})}$$

```
1 def recall_score(self):
2     return (self.tp / ((self.tp) + (self.fn))) * 100
```

v. F1 Score

F1 score is defined as the harmonic mean of precision and recall. Combining the precision and recall metrics into a single metric, F1 score work well on imbalanced data. With formula:

$$F1\ Score = 2 * \frac{(Precision * Recall)}{(Precision + Recall)}$$

```
1 def f1_score(self):
2     p = self.precision_score()
3     r = self.recall_score()
4     return (2 * ((p * r) / (p + r))) / 100
```

3. KDTree

KDTree is a generalization of binary search tree that stores points in k-dimensional space. This means that KDTree can be used to store an array of points in the Cartesian plane. Early stages of starting KDTree are define the nodes first. In our case, every node contains points/value, label (Tertarik) with binary value and distance. Every node contains left and right node to connect between other nodes.

```
kNN
[ ] class KNode:
    def __init__(self, points, y, left = None, right = None, distance = np.inf):
        self.points = points
        self.y = y
        self.left = left
        self.right = right
        self.distance = distance
    def __lt__(self, other):
        return self.distance < other.distance
```

```

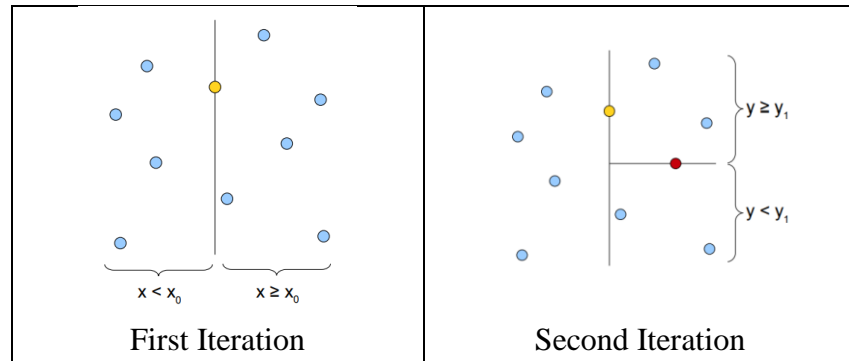
1 class KDTree:
2     def __init__(self, k = 2, p = 2):
3         self.tree = None
4         self.k = k
5         self.depth = 0
6         self.heap = []
7         self.p = p
8
9     def fit(self, X, y):
10        self.X = X
11        self.y = y
12
13        if isinstance(X, pd.DataFrame):
14            X = X.values
15
16        if isinstance(y, pd.DataFrame):
17            y = y.values
18
19        self.tree = self._construct_tree(np.array(X), np.array(y).reshape(-1), 0)
20
21    def _construct_tree(self, points, y, depth):
22        if len(points) == 0:
23            return None
24
25        k = len(points[0])
26        axis = depth % k
27
28        sort_by_axis = np.argsort(points[:, axis])
29        sorted_points = points[sort_by_axis]
30        sorted_y = y[sort_by_axis]
31        mid = len(sorted_points) // 2
32
33        return KNode(
34            sorted_points[mid],
35            sorted_y[mid],
36            self._construct_tree(sorted_points[:mid], sorted_y[:mid], depth + 1),
37            self._construct_tree(sorted_points[mid + 1:], sorted_y[mid + 1:], depth + 1)
38        )
39
40    def nearest_neighbour_search(self, query_point):
41        k = len(query_point)
42        heap = []
43
44        def search(node, depth):
45            if node == None:
46                return
47
48            nonlocal heap
49
50            d = np.linalg.norm(query_point - node.points, ord = self.p)
51            node.distance = -d
52
53            if len(heap) < self.k:
54                heapq.heappush(heap, node)
55            else:
56                heapq.heappushpop(heap, node)
57
58            axis = depth % k
59
60            if query_point[axis] < node.points[axis]:
61                close, other = node.left, node.right
62            else:
63                close, other = node.right, node.left
64
65            search(close, depth + 1)
66
67            delta = abs(query_point[axis] - node.points[axis])
68            nearest = abs(heap[0].distance)
69            isFull = len(heap) > self.k
70
71            if len(heap) < self.k or delta < nearest:
72                search(other, depth + 1)
73
74        search(self.tree, 0)
75        return heap
76
77    def predict(self, X_test):
78        results = []
79        if isinstance(X_test, pd.DataFrame):
80            X_test = X_test.values
81
82        for test in X_test:
83            result = self.nearest_neighbour_search(test)
84            predict_values = [item.y for item in result]
85            counter = Counter(predict_values)
86            results.append(counter.most_common(1)[0][0])
87        return results

```

i. Tree Construction

The tree construction follows the basic tree construction of a binary search tree. On every iteration, the algorithm will find the median of our sorted list of points at a certain axis. Then it will recursively

split the data by doing the same steps and changing the axis. The simplified process can be seen in the image below.



As for the implementation, it can be seen in the image below

```

1 def _construct_tree(self, points, y, depth):
2     if len(points) == 0:
3         return None
4
5     k = len(points[0])
6     axis = depth % k
7
8     sort_by_axis = np.argsort(points[:, axis])
9     sorted_points = points[sort_by_axis]
10    sorted_y = y[sort_by_axis]
11    mid = len(sorted_points) // 2
12
13    return KNode(
14        sorted_points[mid],
15        sorted_y[mid],
16        self._construct_tree(sorted_points[:mid], sorted_y[:mid], depth + 1),
17        self._construct_tree(sorted_points[mid + 1:], sorted_y[mid + 1:], depth + 1)
18    )

```

ii. k-Nearest Neighbor Search

The next concept of our KDTree is the searching method. As a new test point appears, the algorithm will traverse the tree recursively in a binary manner. For each visit, a Priority Queue will record the distance, thus storing all the best distance up until the lead node. After finding the leaf node, it will check whether there is a smaller distance. If it holds true, then the algorithm will backtrack to the correct points. The implementation for traversing the tree can be seen below.

```

1  def nearest_neighbour_search(self, query_point):
2      k = len(query_point)
3      heap = []
4
5      def search(node, depth):
6          if node == None:
7              return
8
9          nonlocal heap
10
11         d = np.linalg.norm(query_point - node.points, ord = self.p)
12         node.distance = -d
13
14         if len(heap) < self.k:
15             heapq.heappush(heap, node)
16         else:
17             heapq.heappushpop(heap, node)
18
19         axis = depth % k
20
21         if query_point[axis] < node.points[axis]:
22             close, other = node.left, node.right
23         else:
24             close, other = node.right, node.left
25
26         search(close, depth + 1)
27
28         delta = abs(query_point[axis] - node.points[axis])
29         nearest = abs(heap[0].distance)
30         isFull = len(heap) > self.k
31
32         if len(heap) < self.k or delta < nearest:
33             search(other, depth + 1)
34
35     search(self.tree, 0)
36     return heap

```

D. Evaluation

1. Creating k-NN Model

We create k-NN Model with $K = 5$ and $p = 2$ (Euclidean Distance). Here we pass in our training and test data which are 'X_train', and 'y_train'. In our case, 'y_pred' is the result of our predictions.

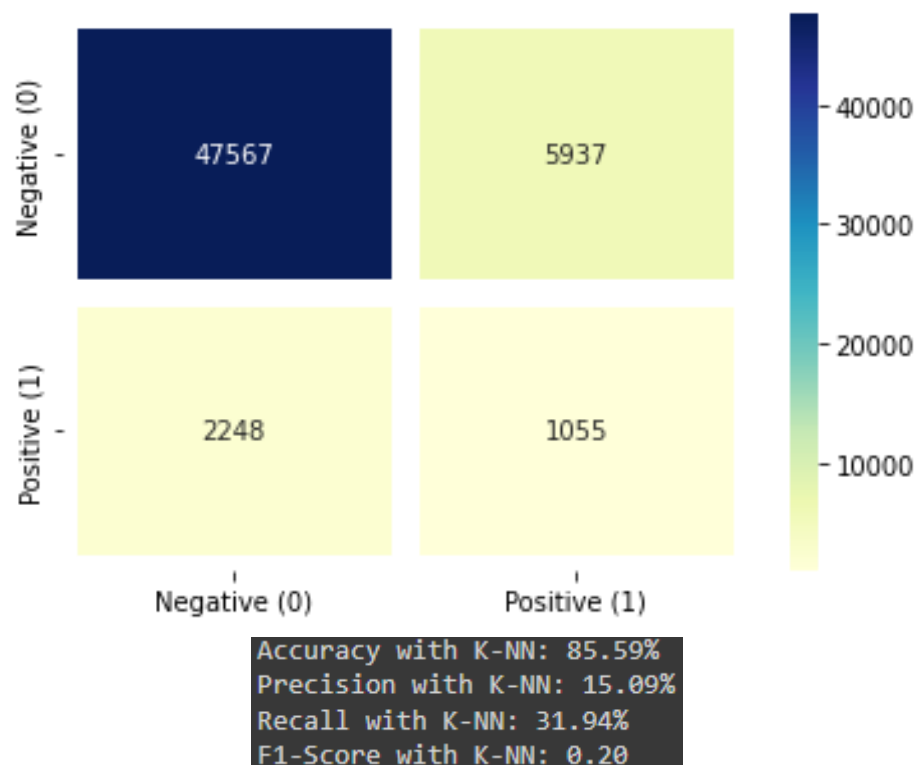
```

[ ] knn = KDTree(k = 5, p = 2)
    knn.fit(X_train, y_train)
    y_pred = knn.predict(X_test)

```

2. Default Result

We can see from the image below that, the metrics has a high accuracy. But somehow, the F1 Score is quite small. Since F1 score is designed work well on imbalanced data, we trying to reach the best F1 Score as much as possible.



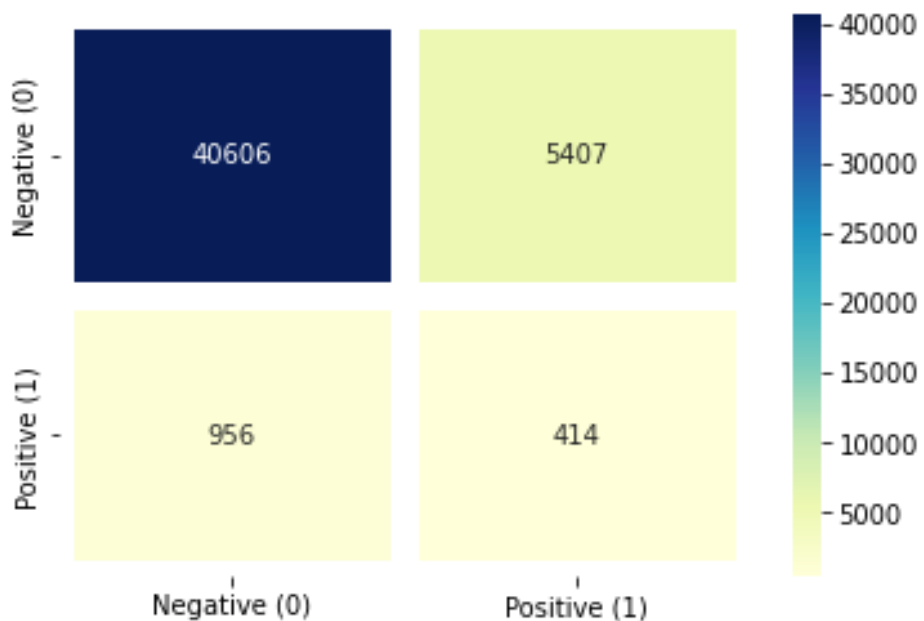
3. Validation Test

With validation test, we change the data frame of our test data with `principal_test` (data test that has gone through the PCA process).

```
[ ] y_pred_test = knn.predict(principal_test)
```

4. Validation Result

We can see after validation test is done; F1 Score is still not high enough. This can be affected by imbalanced data.



```
Accuracy with K-NN: 86.57%
Precision with K-NN: 7.11%
Recall with K-NN: 30.22%
F1-Score with K-NN: 0.12
```

E. Experiment

In this experiment, we will try the pre-processing step. We know that our original data is imbalance, which we will try to oversample and undersample the dataset to balance it. We will also try to check the k values and check the accuracy.

1. Oversampling and Undersampling

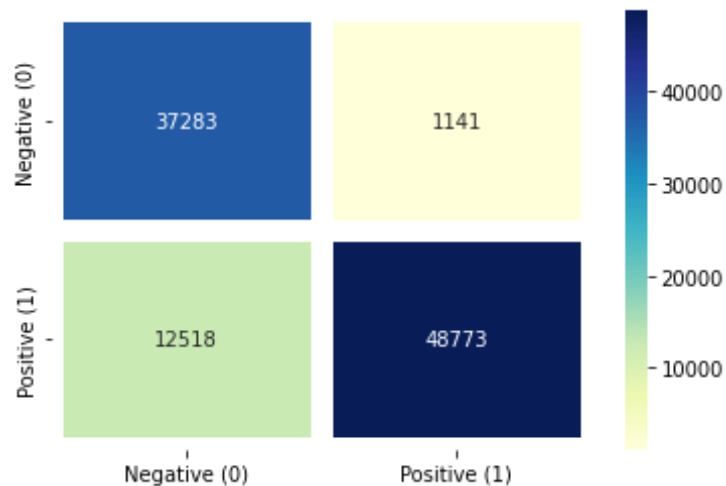
Imbalanced data refers to those types of datasets where the target class has an uneven distribution of observations. Example of a classification problem where the distribution of examples across the known classes is biased or skewed. This can cause a poor predictive performance. We trying to figure it out with oversampling and undersampling.

i. Oversampling

Oversampling is duplicating samples from the minority class. We call RandomOverSampler library into our case.

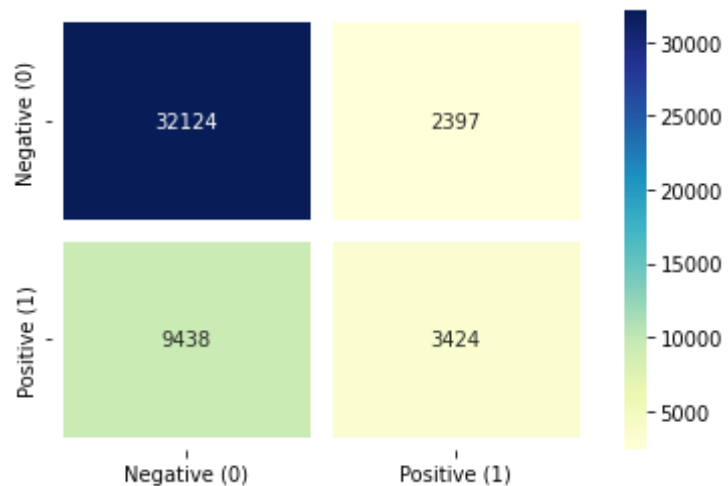
- **Original Oversampling**

With similar K and P or distance, we got:



```
Accuracy with K-NN: 86.30%
Precision with K-NN: 97.71%
Recall with K-NN: 79.58%
F1-Score with K-NN: 0.88
```

- **Validation Test Oversampling**



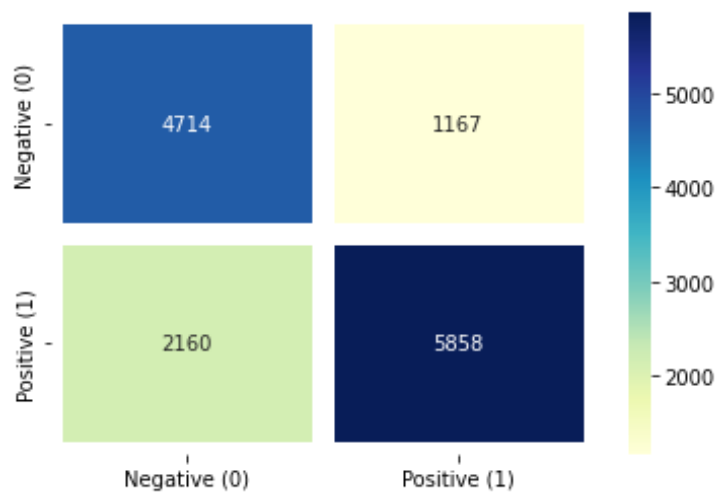
```
Accuracy with K-NN: 75.02%
Precision with K-NN: 58.82%
Recall with K-NN: 26.62%
F1-Score with K-NN: 0.37
```

ii. Undersampling

Undersampling is duplicating samples from the majority class. We call RandomUnderSampler library into our case.

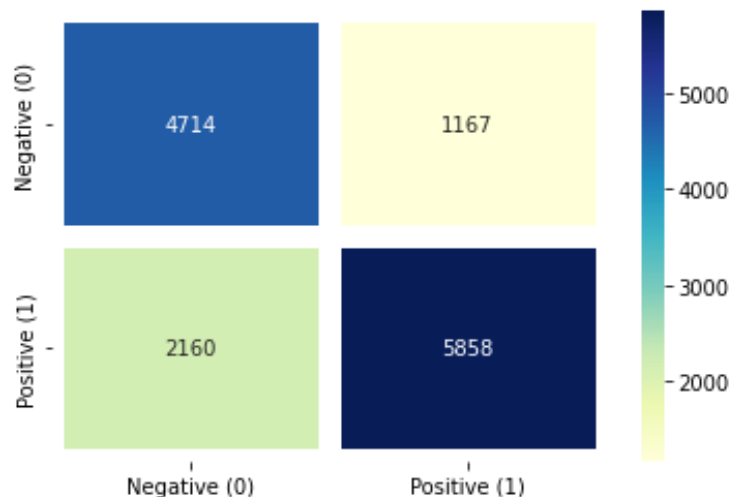
- **Original Undersampling**

With similar K and P or distance, we got:



```
Accuracy with K-NN: 76.06%
Precision with K-NN: 83.39%
Recall with K-NN: 73.06%
F1-Score with K-NN: 0.78
```

- **Validation Test Undersampling**



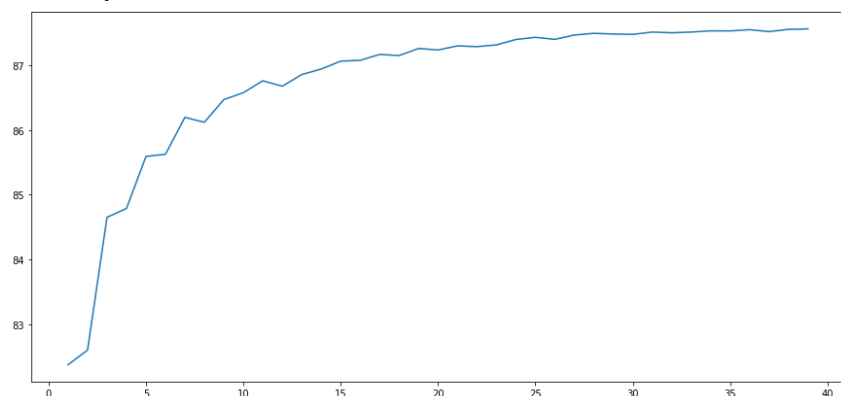
```
Accuracy with K-NN: 72.94%
Precision with K-NN: 75.38%
Recall with K-NN: 27.81%
F1-Score with K-NN: 0.41
```

2. Accuracy Score vs K-Value

In this stage, we tried to figure out the trend between accuracy score and K-Value.

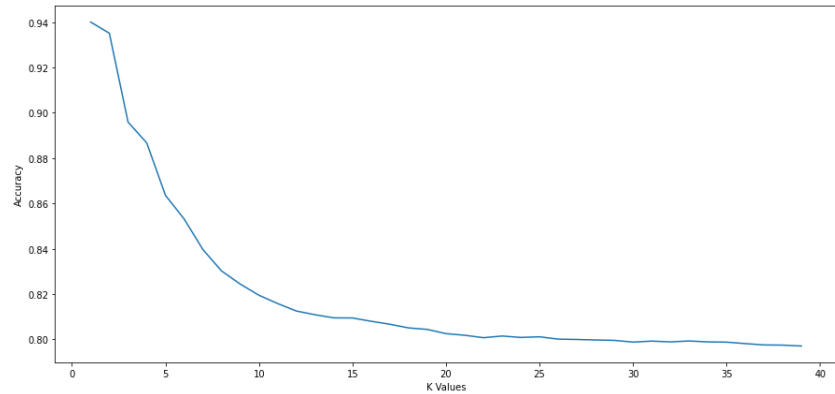
i. Original

We can see from original /without over and under sampling, the larger K makes larger accuracy. This is expected since our original dataset is populated more with the '0' label. When the k is increased, it would discover more '1' which gradually improve the accuracy.



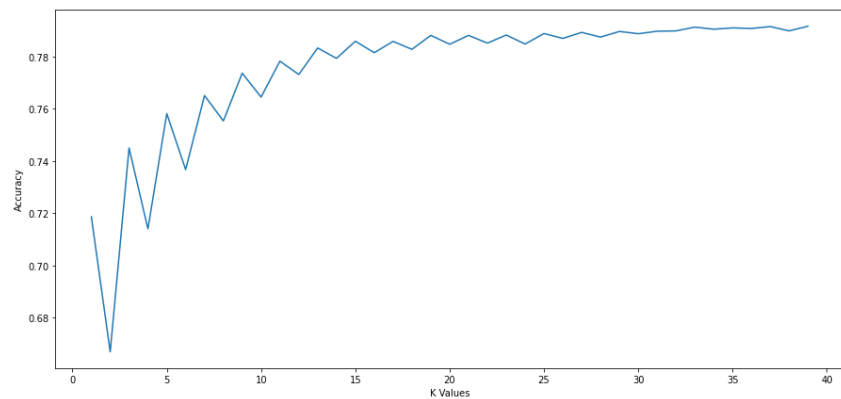
ii. Oversampling

This is trend after oversampling. We can see that the larger K makes accuracy is decreasing. This happens because our dataset is balance. When increasing the k values, it will found more the wrong neighbor value, thus decreasing the accuracy.



iii. Undersampling

This is trend after undersampling. We can see that the larger K, the larger also the accuracy, but not as much as the original data. We can also clearly see that the trend forms spikes.



F. Conclusion

In this assignment, we are given the task to do classification using machine learning. Therefore, we use K-Nearest Neighbor for our model with KDTree implementation. This reduces the time taken when doing the search compared to brute force. We also found out that our original data is not balance which means that there are more ‘Non-Tertarik’ class than the ‘Tertarik’ which result our model to perform poorly with the high-accuracy but low F1-score. To improve our model, we pre-process our dataset to balance the label class. Our method to balance the dataset is by using oversampling and undersampling. We found that model that train using the oversampled dataset performs better than the original and the undersampling one. In conclusion, our best model is the one with the oversampled dataset and K of 5. The result of our evaluation is that our model has a F1-score of 0.78 and accuracy of 72.06% when predicting the test data.