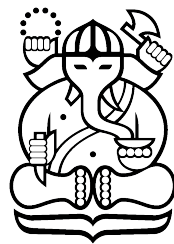


# Implementasi Algoritma KNN dan Naive Bayes

Dibuat sebagai tugas besar Mata Kuliah Inteligensi Buatan - IF3170

Dibuat oleh:

Henry Anand Septian Radityo	13521004
Matthew Mahendra	13521007
Hidayatullah Wildan Ghaly Buchary	13521015
Ahmad Nadil	13521024



Program Studi Teknik Informatika  
Institut Teknologi Bandung  
Bandung  
2023

# Daftar Isi

<b>1</b>	<b>Implementasi Algoritma KNN</b>	<b>2</b>
1.1	Algoritma KNN . . . . .	2
1.2	Implementasi . . . . .	2
<b>2</b>	<b>Implementasi Algoritma Naive Bayes</b>	<b>4</b>
2.1	Algoritma Naive Bayes . . . . .	4
2.2	Implementasi . . . . .	4
<b>3</b>	<b>Perbandingan Hasil Prediksi Algoritma dengan Pustaka Scikit-Learn</b>	<b>6</b>
3.1	Perbandingan KNN . . . . .	6
3.1.1	Hasil Prediksi KNN . . . . .	6
3.1.2	Perbandingan dengan Pustaka Scikit . . . . .	6
3.2	Scikit Naive Bayes . . . . .	7
3.2.1	Hasil Prediksi Naive Bayes . . . . .	7
3.2.2	Perbandingan dengan Pustaka Scikit . . . . .	7
3.3	Insight dari Implementasi Scikit KNN dan Scikit Naive Bayes . . . . .	8
<b>4</b>	<b>Pemrosesan Data untuk Submisi Kaggle</b>	<b>9</b>
4.1	Pemrosesan . . . . .	9
<b>A</b>	<b>Pembagian Tugas</b>	<b>11</b>
<b>B</b>	<b>Github dan Notebook</b>	<b>12</b>

# Bab 1

## Implementasi Algoritma KNN

### 1.1 Algoritma KNN

Algoritma K-Nearest Neighbors (KNN) adalah algoritma *supervised learning* yang biasanya digunakan untuk melakukan klasifikasi kelas terhadap suatu instans data. Algoritma akan membuat suatu model berdasarkan data latih yang sudah ada. Model yang dibentuk akan digunakan untuk dihitung jarak antar atributnya dengan suatu data tes untuk kemudian dihitung jaraknya yang paling kecil.

Untuk semua instans data latih yang selisihnya dengan data tes sudah dihitung, diurutkan dari yang terkecil hingga terbesar. Dari pengurutan tersebut, diambil  $k$  instans data latih untuk dilihat label kelasnya. Label kelas akan dihitung sesuai kategori, e.g. 2 ya 3 tidak 1 ragu. Setelahnya akan diambil label dengan jumlah terbanyak dari instans yang telah diambil. Label tersebut yang akan digunakan sebagai klasifikasinya.

### 1.2 Implementasi

Implementasi algoritma KNN menggunakan bahasa pemrograman Python dengan pendekatan pemrograman berorientasi objek. Metode fit digunakan untuk mempersiapkan model dan metode predict digunakan untuk melakukan prediksi. Semua input menggunakan tipe dataframe. Perhitungan jarak menggunakan Euclidean Distance dengan perhitungan sebagai berikut,

$$d = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

Terdapat konfigurasi untuk proses pelabelan suatu instans. Metode pertama yaitu "uniform" menggunakan metode seperti pada penjelasan algoritma. Metode kedua yaitu "distance" menghitung berdasarkan bobot invers dari jarak suatu data. Dengan perhitungan metode kedua, bisa jadi suatu data yang memiliki jarak yang lebih dekat namun jumlah yang lebih sedikit menjadi yang terpilih karena pembobotan invers dari jarak.

Implementasi algoritma dapat dilihat pada listing 1.1

Listing 1.1: Implementasi Algoritma KNN

```
1 import numpy as np
2 class KNN_Model():
3     # Constructor
4     # Define the number of neighbors to use in prediction
```

```

5  # Default neighbors = 5
6  def __init__(self, k=3, dist=1, weights="uniform"):
7      self.k = k
8      self.dist = dist
9      self.weights = weights
10
11  # Fit Method
12  # Fit the training dataset to the model
13  # X are the features, Y is the target
14  def fit(self, X, y):
15      self.features = X.columns
16      self.X_train = X.to_numpy()
17      self.y_train = y.to_numpy()
18
19  # Predict Method
20  # Prediction using KNN Algorithm
21  # X is the data to be predicted
22  # Return Predicted Data
23  def __predict(self, X):
24      if (self.dist == 1):
25          distances = [self.__euclidean_distance(X, train) for train
in self.X_train]
26      else:
27          # other distance e.g. : Manhattan
28          distances = [self.__manhattan_distance(X, train) for train
in self.X_train]
29      k_indices = np.argsort(distances)[:self.k]
30      k_nearest_labels = [self.y_train[i] for i in k_indices]
31      most_common = None
32      if(self.weights == "uniform"):
33          most_common = np.bincount(k_nearest_labels)
34      elif(self.weights == "distance"):
35          weights = 1 / (np.array(distances) + 1e-10)
36          chosen_weights = [weights[i] for i in k_indices]
37          most_common = np.bincount(k_nearest_labels, weights=
chosen_weights)
38      return np.argmax(most_common)
39
40
41  def predict(self, X):
42      X_ = X.to_numpy()
43      result = [self.__predict(test) for test in X_]
44      return np.array(result)
45
46  # Distance Method for KNN
47  def __euclidean_distance(self, p1, p2):
48      p1_numeric = np.array(p1, dtype=float)
49      p2_numeric = np.array(p2, dtype=float)
50      return np.sqrt(np.sum((p1_numeric - p2_numeric) ** 2))
51
52  def __manhattan_distance(self, p1, p2):
53      p1_numeric = np.array(p1, dtype=float)
54      p2_numeric = np.array(p2, dtype=float)
55      return np.abs(p1_numeric - p2_numeric)

```

## Bab 2

# Implementasi Algoritma Naive Bayes

### 2.1 Algoritma Naive Bayes

Naive Bayes adalah algoritma klasifikasi probabilistik yang berdasarkan pada Teorema Bayes. Algoritma ini sangat efektif dalam kasus-kasus di mana dimensi fitur sangat besar, seperti dalam pemrosesan bahasa alami atau analisis teks. Dalam implementasi Naive Bayes yang telah saya buat, setiap fitur dalam dataset dianggap independen dari fitur lainnya, sebuah asumsi yang disebut "naive". Meskipun asumsi ini tampaknya terlalu sederhana atau "naif", dalam banyak kasus, algoritma ini masih menunjukkan performa yang sangat baik.

Algoritma ini bekerja dengan menghitung probabilitas posterior dari setiap kelas berdasarkan input. Ini dilakukan dengan mengalikan probabilitas priori kelas dengan probabilitas dari masing-masing fitur yang diberikan kelas tersebut. Kemudian, kelas dengan probabilitas tertinggi dipilih sebagai hasil prediksi. Probabilitas dihitung menggunakan perhitungan sebagai berikut,

$$P(y|x_1, x_2, \dots, x_n) = P(y) \times P(x_1|y) \times P(x_2|y) \times \dots \times P(x_n|y)$$

### 2.2 Implementasi

Implementasi algoritma Naive Bayes menggunakan bahasa pemrograman Python dengan pendekatan pemrograman berorientasi objek. Metode fit digunakan untuk mempersiapkan model dan metode predict digunakan untuk melakukan prediksi. Metode fit melakukan proses perhitungan  $P(y)$  sebagai probabilitas dari target serta menghitung mean dan varian dari fitur. Perhitungan mean dan fitur digunakan karena pendekatan Naive Bayes yang digunakan adalah menggunakan distribusi gauss. Perhitungan  $P(x_i|y)$  adalah sebagai berikut,

$$P(x|y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} e^{-\frac{(x-\mu_y)^2}{2\sigma_y^2}}$$

Listing 2.1: Implementasi Algoritma Naive Bayes

```
1 class NaiveBayes:
2     def fit(self, X, y):
3         n_samples, n_features = X.shape
```

```

4         self._classes = np.unique(y)
5         n_classes = len(self._classes)
6
7         # calculate mean, var, and prior for each class
8         self._mean = np.zeros((n_classes, n_features), dtype=np.float64
9     )
10        self._var = np.zeros((n_classes, n_features), dtype=np.float64)
11        self._priors = np.zeros(n_classes, dtype=np.float64)
12
13        for idx, c in enumerate(self._classes):
14            X_c = X[y == c]
15            self._mean[idx, :] = X_c.mean(axis=0)
16            self._var[idx, :] = X_c.var(axis=0)
17            self._priors[idx] = X_c.shape[0] / float(n_samples)
18
19    def predict(self, X):
20        y_pred = [self._predict(test) for test in np.array(X)]
21        return np.array(y_pred)
22
23    def _predict(self, test):
24        posteriors = []
25
26        # calculate posterior probability for each class
27        for idx, c in enumerate(self._classes):
28            prior = np.log(self._priors[idx])
29            posterior = np.sum(np.log(self._pdf(idx, test)))
30            posterior = posterior + prior
31            posteriors.append(posterior)
32
33        # return class with the highest posterior
34        return self._classes[np.argmax(posteriors)]
35
36    def _pdf(self, class_idx, x):
37        mean = self._mean[class_idx]
38        var = self._var[class_idx]
39        numerator = np.exp(-((x - mean) ** 2) / (2 * var + 1e-10))
40        denominator = np.sqrt(2 * np.pi * var)
41        return numerator / denominator

```

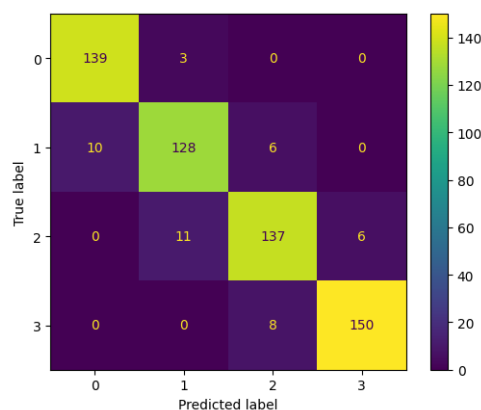
## Bab 3

# Perbandingan Hasil Prediksi Algoritma dengan Pustaka Scikit-Learn

### 3.1 Perbandingan KNN

#### 3.1.1 Hasil Prediksi KNN

Hasil prediksi menggunakan KNN mendapatkan skor 0.93 untuk metrik F1, Precision, dan recall. *Preprocessing* yang dilakukan akan dijelaskan pada bab berikutnya. Confusion matrix yang dihasilkan dapat dilihat pada gambar 3.1.



Gambar 3.1: Confusion Matrix untuk Prediksi KNN

#### 3.1.2 Perbandingan dengan Pustaka Scikit

Implementasi pustaka Scikit-learn `KNeighborsClassifier` menghasilkan hasil yang serupa dengan implementasi algoritma K-Nearest Neighbors (KNN) oleh kelompok kami ketika menggunakan nilai  $k$  yang sama. Perbandingan ini dengan cara membandingkan akurasi. Pustaka Scikit-learn `KNeighborsClassifier` ini menawarkan fleksibilitas dalam pemilihan metode perhitungan jarak, termasuk opsi untuk menggunakan jarak Euclidean, Manhattan, atau Minkowski. Selain itu, Scikit-learn `KNeighborsClassifier` memungkinkan penggunaan bobot dalam perhitungan, seperti bobot yang seragam dimana semua tetangga

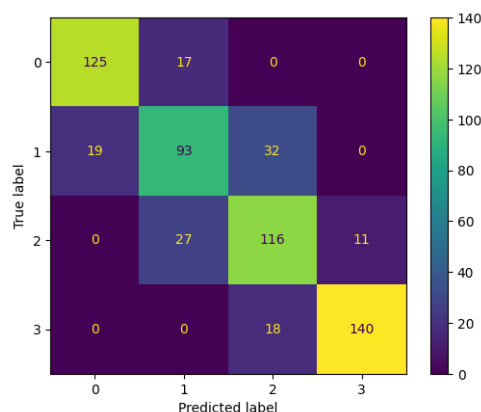
dianggap sama pentingnya, atau bobot berdasarkan jarak, di mana tetangga yang lebih dekat memiliki pengaruh yang lebih besar dalam menentukan klasifikasi. Dalam pengimplementasian Scikit-learn pada KNN, kami memberikan input berupa `n_neighbors` bernilai  $k$  yang sama dengan KNN yang kami buat dan `weights` bernilai "distance" pada fungsi `KNeighborsClassifier`. Setelah itu dilakukan fit dengan masukan `X_train` dan `Y_train`. Hasil prediksi menggunakan data validasi menghasilkan nilai yang sama dengan prediksi dengan KNN yang dibuat. Karena itu, dapat disimpulkan bahwa algoritma yang dibuat adalah sama dengan algoritma pada pustaka Scikit.

## 3.2 Scikit Naive Bayes

Implementasi pustaka Scikit-learn `GaussianNB` menghasilkan hasil yang sama dengan implementasi algoritma Naive Bayes oleh kelompok kami. Perbandingan ini dengan cara membandingkan akurasi. Pustaka `GaussianNB` di Scikit-learn fokus pada penerapan Naive Bayes untuk fitur yang diasumsikan mengikuti distribusi Gaussian (normal), dan tidak menyediakan metode perhitungan jarak yang dapat dipilih. Berbeda dengan beberapa algoritma lain yang mungkin menawarkan opsi menggunakan jarak secara murni atau menggunakan bobot, `GaussianNB` bekerja berdasarkan prinsip dasar Naive Bayes yang mengasumsikan independensi antar fitur dan menggunakan probabilitas untuk mengklasifikasikan data. Dalam pengimplementasian Scikit-learn pada Naive Bayes, kami memberikan input berupa input kosong pada fungsi `GaussianNB`. setelah itu dilakukan fit dengan masukan `X_train` dan `Y_train`.

### 3.2.1 Hasil Prediksi Naive Bayes

Hasil Prediksi Naive Bayes mendapatkan skor akurasi sebesar 0.79 dengan precision dan recall sebesar 0.79 sehingga mendapatkan F1-score sebesar 0.79. *Preprocessing* yang dilakukan akan dijelaskan pada bab berikutnya. Confusion matrix yang dihasilkan dapat dilihat pada gambar 3.2.



Gambar 3.2: Confusion Matrix untuk Prediksi Naive Bayes

### 3.2.2 Perbandingan dengan Pustaka Scikit

Implementasi pustaka Scikit-learn untuk model Naive bayes menghasilkan akurasi, precision, dan recall yang sama dengan implementasi model dari kelompok kami. Hal ini



membuktikan bahwa implementasi dari naive bayes memberikan hasil yang positif dikarenakan parameter yang diberikan sama dengan default dari pustaka scikit. Data yang dimasukkan untuk kedua model juga merupakan data yang diproses dengan pipeline yang sama sehingga dipastikan hasil prediksi dari model kami dan pustaka merupakan hasil yang sama.

### **3.3 Insight dari Implementasi Scikit KNN dan Scikit Naive Bayes**

Dalam analisis kami, penggunaan Scikit-learn `KNeighborsClassifier` dan `GaussianNB` menunjukkan kemampuan yang efektif dalam mengklasifikasikan data, mirip dengan implementasi algoritma yang kami kembangkan sendiri. Khusus untuk `KNeighborsClassifier`, keunggulannya terletak pada fleksibilitas dalam pemilihan metode perhitungan jarak dan penggunaan bobot. Faktor ini sangat penting karena memungkinkan penyesuaian algoritma sesuai dengan sifat data yang spesifik. Dengan menggunakan bobot berdasarkan jarak, dimana tetangga yang lebih dekat memiliki pengaruh yang lebih besar, kami dapat meningkatkan akurasi klasifikasi. Hal ini menunjukkan bahwa Scikit-learn tidak hanya menyediakan alat yang efisien tetapi juga fleksibel, memungkinkan penyesuaian yang dapat meningkatkan performa model tergantung pada kebutuhan dataset tertentu. Hasil prediksi menggunakan data validasi menghasilkan nilai yang sama dengan prediksi dengan Naive Bayes yang dibuat. Karena itu, dapat disimpulkan bahwa algoritma yang dibuat adalah sama dengan algoritma pada pustaka Scikit.

Implementasi `GaussianNB` pada Scikit-learn menawarkan pendekatan yang berbeda. Algoritma ini mengasumsikan bahwa fitur mengikuti distribusi Gaussian, yang merupakan asumsi yang kuat namun seringkali efektif dalam berbagai kasus. Tidak adanya kebutuhan untuk pemilihan metode perhitungan jarak atau penggunaan bobot menunjukkan simplicitas dari algoritma Naive Bayes, yang mengandalkan probabilitas bersyarat dan asumsi independensi antar fitur. Hal ini dapat sangat berguna dalam kasus-kasus di mana hubungan antara fitur tidak terlalu rumit atau ketika distribusi data cukup normal. Kemudahan implementasi dan efisiensi komputasi dari `GaussianNB` menjadikannya pilihan yang menarik untuk banyak masalah klasifikasi, terutama ketika diperlukan solusi yang cepat dan sederhana.

## Bab 4

# Pemrosesan Data untuk Submisi Kaggle

### 4.1 Pemrosesan

Pemrosesan yang dilakukan adalah membuat pipeline untuk membantu pembersihan data. Pipeline yang dilakukan adalah sebagai berikut:

1. Menghapus data outlier khususnya di atribut 'fc'
2. Mengganti nilai invalid pada semua atribut dengan median
3. Melakukan filter atribut sehingga menyisakan beberapa atribut yaitu battery\_power, int\_memory, mobile\_wt, px\_height, px\_width, ram, sc\_h dan kolom target

Dilakukan penghapusan outlier untuk membantu model dalam melakukan training lalu untuk semua nilai invalid seperti pada kolom sc\_w terdapat beberapa nilai 0. Hal ini mustahil apabila dilihat dari metadata yang diberikan. Selanjutnya tahap terakhir dari pipeline adalah menghilangkan semua atribut yang memiliki korelasi yang kurang dan persebaran yang kurang merata sehingga model seperti knn akan lebih baik dalam mencari suatu jarak dari antara kolom test dan kolom train.

Listing 4.1: Pipeline

```
1 class OutlierRemoval(BaseEstimator, TransformerMixin):
2     def __init__(self, columns=None):
3         self.columns = columns
4         self.lower_bounds = {}
5         self.upper_bounds = {}
6
7     def fit(self, X, y=None):
8         if self.columns is None:
9             self.columns = X.columns
10
11         for column in self.columns:
12             q3 = X[column].quantile(q=0.75)
13             q1 = X[column].quantile(q=0.25)
14             IQR = q3 - q1
15             lower_bound = q1 - 1.5 * IQR
16             upper_bound = q3 + 1.5 * IQR
17
18             self.lower_bounds[column] = lower_bound
```

```

19         self.upper_bounds[column] = upper_bound
20
21     return self
22
23     def transform(self, X):
24         for column in self.columns:
25             X = X[(X[column] >= self.lower_bounds[column]) & (X[column]
26                 <= self.upper_bounds[column])]
27             return X
28
29 class InvalidReplacement(BaseEstimator, TransformerMixin):
30     def __init__(self):
31         return None
32
33     def fit(self, X, y=None):
34         self.median = X['px_height'].median()
35         self.columns = X.columns
36         return self
37
38     def transform(self, X):
39         X_ = X.copy()
40         for col in self.columns:
41             if col != 'price_range':
42                 X_[col] = X_[col].replace(0, X_[col].median())
43         return X_
44
45 class DropColumn(BaseEstimator, TransformerMixin):
46     def __init__(self, column):
47         self.column = column
48
49     def fit(self, X, y=None):
50         return self
51
52     def transform(self, X):
53         X_ = X.drop(columns = self.column)
54         return X_
55
56 pipeline = Pipeline([
57     ('removing outliers', OutlierRemoval(columns=numeric_columns)),
58     ('invalid data replacement', InvalidReplacement()),
59     ('feature selection', DropColumn(non_numeric_columns + ['fc'] +
60         remove_col))
61 ])

```

# Lampiran A

## Pembagian Tugas

NIM	Pekerjaan
13521004	KNN, Submisi Kaggle
13521007	KNN, Scikit, Submisi Kaggle
13521015	Naive Bayes, Submisi Kaggle
13521024	Naive Bayes, Submisi Kaggle

Tabel A.1: Pembagian Kerja

# Lampiran B

## Github dan Notebook

- Repository Github: [https://github.com/MHEN2606/Tubes2\\_melihata](https://github.com/MHEN2606/Tubes2_melihata)
- Kaggle Notebook: <https://www.kaggle.com/code/matthewmahendra/tubes2-melihata/notebook>