

Laporan Tugas Besar 1
IF3170 Inteligensi Artifisial
Implementasi Algoritma Pembelajaran Mesin
Semester I Tahun 2024/2025



Disusun Oleh:
Mohammad Nugraha Eka Prawira (13522001)
Muhammad Al Thariq Fairuz (13522027)
Randy Verdian (13522067)
Emery Fathan Zwageri (13522079)

Program Studi Teknik Informatika
Sekolah Teknik Elektro dan Informatika
Institut Teknologi Bandung
2024

Bab 1	
Implementasi KNN.....	3
Bab 2	
Implementasi Naive-Bayes.....	5
Bab 3	
Implementasi ID3.....	7
Bab 4	
Tahap Cleaning dan Preprocessing.....	11
4.1. Cleaning.....	11
4.2. Preprocessing.....	11
Bab 5	
Perbandingan Hasil Prediksi.....	14
Bab 6	
Kontribusi Kelompok.....	17
Bab 7	
Referensi.....	18

Bab 1

Implementasi KNN

```
KNN.py
1 import numpy as np
2 from sklearn.base import BaseEstimator, ClassifierMixin # Import these so we can use cross_val_score to evaluate our model, we take classifiermixin since we are building a classifier
3
4 class KNeighbors(BaseEstimator, ClassifierMixin): # Inherit from BaseEstimator and ClassifierMixin
5     """
6     Create an instance of the K-Nearest Neighbors algorithm
7     """
8
9     def __init__(self, k: int = 3, metrics: str = 'euclidean'):
10         """
11         Constructor for the K-Nearest Neighbors algorithm
12         """
13         self.k = k
14         self.metrics = metrics
15
16     def fit(self, X : np.ndarray, y : np.ndarray):
17         """
18         Fit the training data to the model
19         """
20         self.X_train = X
21         self.y_train = y
22
23     def euclidean_distance(self, x1 : np.ndarray, x2 : np.ndarray) -> float:
24         """
25         Compute the Euclidean distance between two points
26         """
27         return np.sqrt(np.sum((x1 - x2)**2))
28
29     def manhattan_distance(self, x1 : np.ndarray, x2 : np.ndarray) -> float:
30         """
31         Compute the Manhattan distance between two points
32         """
33         return np.sum(np.abs(x1 - x2))
34
35     def minkowski_distance(self, x1 : np.ndarray, x2 : np.ndarray, p : int) -> float :
36         """
37         Compute the Minkowski distance between two points
38         """
39         return np.sum(np.abs(x1 - x2)**p)**(1/p)
40
41     def _predict(self, x : np.ndarray, p : int = 2) -> int:
42         """
43         Predict the class label for a single sample
44         """
45
46         # Compute distances between x and all examples in the training set
47         if self.metrics == 'euclidean':
48             distances = [self.euclidean_distance(x, x_train) for x_train in self.X_train]
49         elif self.metrics == 'manhattan':
50             distances = [self.manhattan_distance(x, x_train) for x_train in self.X_train]
51         elif self.metrics == 'minkowski':
52             distances = [self.minkowski_distance(x, x_train, p) for x_train in self.X_train]
53         else:
54             raise ValueError('Invalid metrics')
55
56         # Sort by distance and return indices of the first k neighbors
57         k_indices = np.argsort(distances)[:self.k] # Return the indices of the k nearest neighbors from the lowest distance
58
59         # Extract the labels of the k nearest neighbor training samples
60         k_nearest_labels = [self.y_train[i] for i in k_indices]
61
62         # Return the most common class label
63         most_common, counts = np.unique(k_nearest_labels, return_counts=True)
64         return most_common[np.argmax(counts)]
65
66     def predict(self, X : np.ndarray, p : int = 2) -> np.ndarray:
67         """
68         Predict the class labels for a set of samples
69         """
70         y_pred = [self._predict(x, p) for x in X]
71         return np.array(y_pred)
```

Snipped

Kelas KNeighbors merupakan implementasi algoritma KNN dari *scratch* menggunakan python. Kelas ini menurunkan kelas dari BaseEstimator, ClassifierMixin agar kompatibel dengan cross_val_score dari scikit-learn sebagai fungsi untuk mengukur kemampuan model.

Kelas ini memiliki dua parameter, yaitu jumlah tetangga dan metrik yang digunakan, metrik yang tersedia ada euclidean, manhattan, dan minkowski. Algoritma ini bekerja dengan menghitung jarak antara data yang akan diprediksi dengan semua data yang ada pada *training* data dengan menggunakan metrik yang telah dipilih sebelumnya. Setelah perhitungan jarak dilakukan, akan diambil indeks dari k tetangga terdekat dan kemudian akan dihitung kelas terbanyak yang dimiliki oleh tetangga tersebut, kelas yang paling banyak muncul di antara tetangga-tetangga ini akan menjadi kelas yang diprediksi untuk data baru.

Bab 2

Implementasi Naive-Bayes

```
GNB.py

1 import numpy as np
2 from sklearn.base import BaseEstimator, ClassifierMixin
3
4 class GaussianNaiveBayes(BaseEstimator, ClassifierMixin):
5     """
6     Create an instance of the Gaussian Naive Bayes algorithm
7     """
8
9     def fit(self, X : np.ndarray, y : np.ndarray):
10        """
11        Fit the training data to the model, calculate the mean, variance, and prior for each class
12        """
13
14        # Check for all unique values (classes) in the label
15        self.classes = np.unique(y)
16        self.epsilon = 1e-9 # Prevent division by zero
17
18        # Store the mean, variance, and prior for each class
19        self.mean = {}
20        self.var = {}
21        self.priors = {}
22
23        # Calculate the mean, variance, and prior for each class
24        for cls in self.classes:
25            X_c = X[y == cls]
26            self.mean[cls] = X_c.mean(axis=0)
27            self.var[cls] = X_c.var(axis=0)
28            self.priors[cls] = X_c.shape[0] / X.shape[0]
29
30        # Calculate the probability of a feature vector x in a given class
31        def _pdf(self, cls : int, x : np.ndarray) -> float:
32            """
33            Calculate the probability of a feature vector x in a given class
34            """
35
36            mean = self.mean[cls]
37            var = self.var[cls] + self.epsilon
38            numerator = np.exp(-(x - mean) ** 2 / (2 * var))
39            denominator = np.sqrt(2 * np.pi * var)
40
41            # Error handling: Add epsilon to denominator where it is zero
42            denominator = np.where(denominator == 0, self.epsilon, denominator)
43
44            return numerator / denominator
45
46        def predict_instance(self, x : np.ndarray) -> int:
47            """
48            Predict the class label for a single sample
49            """
50
51            posteriors = []
52
53            for cls in self.classes:
54                prior = np.log(self.priors[cls])
55                pdf_values = self._pdf(cls, x)
56
57                # Adding epsilon before taking the log to avoid log(0)
58                pdf_values = np.where(pdf_values == 0, self.epsilon, pdf_values)
59
60                # Calculate the posterior and handle -inf values
61                log_pdf_values = np.log(pdf_values)
62                log_pdf_values[np.isnan(log_pdf_values)] = -1e9 # Replace -inf with a very large negative number
63
64                posterior = np.sum(log_pdf_values)
65                posterior = prior + posterior
66                posteriors.append(posterior)
67
68            return self.classes[np.argmax(posteriors)]
69
70        def predict(self, X : np.ndarray) -> np.ndarray:
71            """
72            Predict the class label for all of the samples
73            """
74
75            y_pred = [self.predict_instance(x) for x in X]
76            return np.array(y_pred)
```

Snipped

Kelas GaussianNaiveBayes merupakan implementasi dari algoritma Naive Bayes dengan asumsi Gaussian (distribusi normal) yang dibuat dari awal menggunakan python. Kelas ini juga menurunkan kelas dari BaseEstimator dan ClassifierMixin untuk kompatibilitas dengan fungsi cross_val_score dari scikit-learn.

Algoritma ini pertama kali mengidentifikasi semua kelas unik yang ada dalam data. Untuk setiap kelas, algoritma menghitung:

- Mean (rata-rata): Menghitung nilai rata-rata setiap fitur untuk data yang termasuk dalam kelas tersebut
- Variance (varians): Menghitung sebaran data untuk setiap fitur dalam kelas tersebut
- Prior: Menghitung probabilitas awal setiap kelas, yaitu jumlah data dalam kelas tersebut dibagi total jumlah data

Kemudian, akan dihitung probabilitas sebuah nilai fitur termasuk dalam suatu kelas menggunakan distribusi normal. Rumus yang digunakan adalah rumus distribusi normal:

$$(1/\sqrt{2\pi\sigma^2}) * e^{-(x - \mu)^2/2\sigma^2}$$

Dimana μ adalah mean, σ^2 adalah variance, epsilon ($1e-9$) ditambahkan untuk menghindari pembagian dengan nol

Untuk setiap kelas, akan dihitung posterior probability dengan rumus Bayes dan digunakan logaritma untuk menghindari underflow karena perkalian angka-angka kecil. Prior probability diubah ke dalam bentuk logaritma dan PDF juga dihitung untuk setiap fitur dan dijumlahkan (dalam bentuk log). Kelas dengan posterior probability tertinggi dipilih sebagai hasil prediksi.

Bab 3

Implementasi ID3

```
1 import numpy as np
2 from sklearn.base import BaseEstimator, ClassifierMixin
3
4 class Node:
5     """
6     Create a node for the Decision Tree
7     """
8
9     def __init__(self, feature_index: int = None, threshold: float = None, left: 'Node' = None, right: 'Node' = None, info_gain: float = None, value: int = None):
10         """
11         Constructor for the Node class
12         """
13         self.feature_index = feature_index
14         self.threshold = threshold
15         self.left = left
16         self.right = right
17         self.info_gain = info_gain
18         self.value = value
19
20 class DecisionTree(BaseEstimator, ClassifierMixin):
21     """
22     Create an instance of the Decision Tree algorithm
23     """
24     def __init__(self, min_samples_split: int = 4, max_depth: int = 4, mode: str = "gini"):
25         """
26         Constructor for the Decision Tree model
27         """
28         self.min_samples_split = min_samples_split
29         self.max_depth = max_depth
30         self.root = None
31         self.mode = mode
32
33     def split(self, dataset: np.ndarray, feature_index: int, threshold: float) -> tuple[np.ndarray, np.ndarray]:
34         """
35         Split the dataset into two parts based on the feature and threshold
36         """
37         dataset_left = np.array([row for row in dataset if row[feature_index] <= threshold])
38         dataset_right = np.array([row for row in dataset if row[feature_index] > threshold])
39         return dataset_left, dataset_right
40
41     def entropy(self, y: np.ndarray) -> float:
42         """
43         Calculate the entropy of the dataset
44         """
45         class_labels = np.unique(y)
46         entropy = 0
47         for cls in class_labels:
48             p = len(y[y == cls]) / len(y)
49             entropy += -p * np.log2(p)
50
51         return entropy
52
53     def gini_index(self, y: np.ndarray) -> float:
54         """
55         Calculate the Gini index of a tree node
56         """
57         class_labels = np.unique(y)
58         gini = 0
59         for cls in class_labels:
60             p = len(y[y == cls]) / len(y)
61             gini += p * (1 - p)
62
63         return gini
64
65     def calculate_leaf_value(self, y: np.ndarray) -> int:
66         """
67         Calculate the leaf node value
68         """
69         Y = list(y)
70         return max(Y, key=Y.count)
71
72     def information_gain(self, parent_node: np.ndarray, left_child: np.ndarray, right_child: np.ndarray, mode: str = "gini") -> float:
73         """
74         Calculate the information gain of a split
75         """
76
77         left_weight = len(left_child) / len(parent_node)
78         right_weight = len(right_child) / len(parent_node)
79
80         if mode == "gini":
81             gain = self.gini_index(parent_node) - (left_weight * self.gini_index(left_child) + right_weight * self.gini_index(right_child))
82         else:
83             gain = self.entropy(parent_node) - (left_weight * self.entropy(left_child) + right_weight * self.entropy(right_child))
84
85         return gain
```

Snipped

```

87 def get_best_split(self, dataset: np.ndarray, num_features: int) -> dict:
88     """
89     Find the best split for the dataset based on the information gain
90     """
91     best_split= {}
92     max_info_gain= -float('inf')
93
94     for feature_index in range (num_features):
95         feature_values= dataset[:, feature_index]
96         possible_thresholds= np.unique(feature_values)
97
98         for thresholds in possible_thresholds:
99             dataset_left, dataset_right= self.split(dataset, feature_index, thresholds)
100
101             if len(dataset_left) > 0 and len(dataset_right) > 0:
102                 y, y_left, y_right= dataset[:, -1], dataset_left[:, -1], dataset_right[:, -1]
103
104                 current_info_gain= self.information_gain(y, y_left, y_right, self.mode)
105
106                 if current_info_gain > max_info_gain:
107                     best_split['feature_index']= feature_index
108                     best_split['threshold']= thresholds
109                     best_split['data_left']= dataset_left
110                     best_split['data_right']= dataset_right
111                     best_split['info_gain']= current_info_gain
112                     max_info_gain= current_info_gain
113
114     return best_split
115
116 def construct_tree(self, dataset: np.ndarray, current_depth: int = 0) -> 'Node':
117     """
118     Construct the decision tree recursively
119     """
120     X, Y= dataset[:, :-1], dataset[:, -1]
121     num_samples, num_features= X.shape
122
123     if num_samples >= self.min_samples_split and current_depth <= self.max_depth:
124         best_split= self.get_best_split(dataset, num_features)
125
126         if best_split['info_gain'] > 0:
127             left_subtree= self.construct_tree(best_split['data_left'], current_depth+1)
128             right_subtree= self.construct_tree(best_split['data_right'], current_depth+1)
129             return Node(best_split['feature_index'], best_split['threshold'], left_subtree, right_subtree, best_split['info_gain'])
130
131     # Calculate leaf node
132     leaf_value= self.calculate_leaf_value(Y)
133
134     return Node(value= leaf_value)
135
136 # def print(self, tree= None, indent= " "):
137
138 #     if tree is None:
139 #         tree= self.root
140
141 #     if tree.value is not None:
142 #         print(tree.value)
143
144 #     else:
145 #         print(f"{tree.feature_index}, {tree.threshold} {tree.info_gain}")
146 #         print(f"{indent}left: ", end= "")
147 #         self.print(tree.left, indent + indent)
148 #         print(f"{indent}right: ", end= "")
149 #         self.print(tree.right, indent + indent)
150
151 def fit(self, X: np.ndarray, Y: np.ndarray):
152     """
153     Fit the training data to the model
154     """
155     Y= np.reshape(Y, (-1, 1))
156     dataset= np.concatenate((X, Y), axis= 1)
157     self.root= self.construct_tree(dataset)
158
159 def _predict(self, x: np.ndarray, tree: 'Node') -> int:
160     """
161     Predict the class label for a single sample
162     """
163
164     if tree.value is not None:
165         return tree.value
166
167     feature_val= x[tree.feature_index]
168
169     if feature_val <= tree.threshold:
170         return self._predict(x, tree.left)
171     else:
172         return self._predict(x, tree.right)
173
174 def predict(self, X: np.ndarray) -> np.ndarray:
175     """
176     Predict the class labels for a set of samples
177     """
178     return [self._predict(x, self.root) for x in X]

```


ID3 adalah algoritma yang membangun model prediksi dalam bentuk struktur pohon. Implementasi ini menggunakan pendekatan top-down (dari atas ke bawah) dengan dua metrik impurity yang bisa dipilih: Gini Index atau Entropy. Namun, karena yang diminta adalah ID3 model yang basisnya hanya menerima metrik impurity entropy, maka pada tugas besar kali ini, pelatihan dan pengujian model hanya akan menggunakan metrik impurity entropy saja. Model ini dibangun dengan meniru model CART yang ada pada *library* scikit-learn.

Implementasi ini mendefinisikan kelas Node yang merepresentasikan setiap titik keputusan dalam pohon. Kemudian, kelas DecisionTree: Kelas utama ini mewarisi dari BaseEstimator dan ClassifierMixin untuk kompatibilitas dengan cross_val_score scikit-learn. Parameter inisialisasinya meliputi:

- `min_samples_split`: Jumlah minimum sampel yang dibutuhkan untuk melakukan split
- `max_depth`: Kedalaman maksimum pohon
- `mode`: Metrik yang digunakan (hanya entropy yang akan digunakan pada tugas besar kali ini)

ID3 bekerja dengan membangun struktur pohon keputusan melalui proses *training* yang berurutan. Pertama, algoritma menerima data *training* yang terdiri dari fitur dan label kelas. Proses pembangunan pohon dimulai dari akar (root) dan berlanjut secara rekursif ke bawah.

Pada setiap node, algoritma mencari cara terbaik untuk memisahkan data menjadi dua kelompok. Ini dilakukan dengan mengevaluasi setiap fitur dan mencoba berbagai nilai threshold untuk menemukan pemisahan yang menghasilkan information gain tertinggi. Information gain menunjukkan seberapa baik sebuah pemisahan dapat memurnikan data, semakin tinggi nilainya, semakin baik pemisahan tersebut.

Pengujian model ini menggunakan metrik impurity entropy, yaitu dengan cara mengukur ketidakpastian dalam data. Metrik ini membantu algoritma menentukan seberapa baik sebuah pemisahan dalam memisahkan kelas-kelas yang berbeda.

Proses pembuatan cabang baru akan terus berlanjut sampai salah satu kondisi pemberhentian terpenuhi:

- Jumlah sampel di node kurang dari minimum yang ditentukan (`min_samples_split`)
- Kedalaman pohon mencapai batas maksimum (`max_depth`)
- Tidak ada lagi pemisahan yang menghasilkan information gain positif
- Semua sampel di node tersebut sudah berasal dari kelas yang sama (homogen)

Ketika melakukan prediksi untuk data baru, algoritma mulai dari node akar dan mengikuti jalur ke bawah berdasarkan nilai-nilai fitur data tersebut. Pada setiap node internal (bukan daun), algoritma membandingkan nilai fitur dengan threshold node tersebut. Jika nilainya lebih kecil atau sama dengan threshold, data akan diarahkan ke cabang kiri; jika lebih besar, akan diarahkan ke cabang kanan.

Proses ini berlanjut sampai mencapai node daun (leaf node). Label kelas yang tersimpan di node daun tersebut menjadi prediksi untuk data input. Label ini ditentukan berdasarkan kelas mayoritas dari data *training* yang berakhir di node tersebut selama proses *training*.

Bab 4

Tahap Cleaning dan Preprocessing

4.1. Cleaning

Pada tahap cleaning kita melakukan imputasi pada *missing values* karena hasil dari EDA kami datanya sangat *noisy* dan kelas cenderung bernilai 'normal'. Awalnya kami membuat beberapa imputer kami sendiri, tetapi tidak perform baik sehingga kami menggunakan imputer yang biasa saja. Ini beberapa Imputer yang kami gunakan:

1. SimpleImputer(Strategy='most_frequent')
Imputer ini kami gunakan untuk impute missing value pada fitur kategorikal
2. SimpleImputer(Strategy='median')
Imputer ini kami gunakan untuk impute missing value pada fitur numerikal

4.2. Preprocessing

Pada tahap ini kami membuat pipeline dengan method Pipeline dan ColumnTransformer dari scikit learn. Kami mendefinisikan dua pipeline, satu untuk kolom numerik satu untuk kategorikal

```
num_cols= list(df_train.select_dtypes(include=np.number))

cat_cols = list(df_train.select_dtypes(include='object'))
all_cols = list(df_train.columns)

cat_transformer = Pipeline(steps = [
    ('imputer', SimpleImputer(strategy='most_frequent')),
    ('OHE', OneHotEncoder(handle_unknown='ignore', sparse_output = False)),
])

num_transformer = Pipeline(steps = [
    ('imputer', SimpleImputer(strategy='median')),
    ('scaler', StandardScaler()),
])

preprocessor = ColumnTransformer(
    transformers = [
        ('cat_transformer', cat_transformer, cat_cols),
        ('num_transformer', num_transformer, num_cols)
    ]
)
```

```
▶ X_train_transformed = preprocessor.fit_transform(X_train)
   X_val_transformed = preprocessor.transform(X_val)

[ ] le= LabelEncoder()
    y_train= le.fit_transform(y_train)
    y_val= le.transform(y_val)
```

Beberapa teknik feature encoder yang kami gunakan:

1. One Hot Encoder
Kami menggunakan One Hot Encoder untuk encoding data kategorikal karena kebanyakan dari data kategorikal tidak memiliki arti urutan, karena itu lah kami menggunakan OHE ketimbang Ordinal Encoder. Sparse = false karena kita ingin keluarannya numpy array bukan sparse matrix, kekurangannya adalah lebih boros memory.
2. StandardScaler
Kami menggunakan ini untuk meng scale datanya agar berada pada skala yang sama yaitu mean =0 variance= 1.
3. LabelEncoder
Digunakan untuk encoding target menjadi nilai integer ordinal

Kami meng compile Pipeline tadi dengan menggunakan ColumnTransformer agar masing masing pipeline bisa diterapkan pada kolom yang sesuai. Setelah dicompile kita fit_transform pada data training dan transform pada data validasi.

Kami juga mengimplementasikan AutoEncoder dengan tensorflow dan juga menggunakan PCA, tetapi tidak memberikan hasil yang baik sehingga tidak kami gunakan untuk submisi.

```

class AutoEncoders(Model):
    def __init__(self, output_units):
        super().__init__()
        self.encoder = Sequential(
            [
                Dense(156, activation = 'relu'),
                Dense(128, activation = 'relu'),
            ]
        )
        self.decoder = Sequential(
            [
                Dense(156, activation= "relu"),
                Dense(output_units, activation="tanh")
            ]
        )
    def call(self, inputs):
        encoded = self.encoder(inputs)
        decoded = self.decoder(encoded)
        return decoded

[ ] ae = AutoEncoders(len(X_train_scaled.columns))

[ ] ae.compile(
    loss='mae',
    metrics=['mae'],
    optimizer='adam'
)

history = ae.fit(
    X_train_scaled,
    X_train_scaled,
    epochs=10,
    batch_size=32,
    validation_data=(X_test_scaled, X_test_scaled)
)

```

Bab 5

Perbandingan Hasil Prediksi

1. KNN

Scratch

```
[ ] from KNN import KNNeighbours
```

```
knn_data= X_train_transformed[:10000]
knn_data_label= y_train[:10000]
```

```
knn_test= X_val_transformed[:1000]
knn_test_label= y_val[:1000]
```

```
[ ] knn_scratch= KNNeighbours(k=5)
knn_scratch.fit(knn_data,knn_data_label)
y_pred= knn_scratch.predict(knn_test)
```

```
f1_macro_scratch = f1_score(knn_test_label,y_pred,average ='macro')
print(f'F1 Macro Average: {f1_macro_scratch}')
```

```
precision    recall  f1-score   support
```

```
0      0.20      0.08      0.12        12
1      0.00      0.00      0.00        11
2      0.33      0.31      0.32       74
3      0.58      0.74      0.65      189
4      0.46      0.49      0.48        90
5      0.99      0.98      0.99       235
6      0.90      0.87      0.89      313
7      0.69      0.49      0.58        69
8      0.50      0.14      0.22         7
```

```
accuracy          0.75      1000
macro avg         0.52      0.46      0.47      1000
weighted avg      0.74      0.75      0.74      1000
```

Scikit-learn

```
knn= KNeighborsClassifier(
    n_neighbors=5,
    metric='euclidean'
)
```

```
knn.fit(knn_data, knn_data_label)
y_pred= knn.predict(knn_test)
```

```
f1_macro_scratch = f1_score(knn_test_label, y_pred,average ='macro')
print(f'F1 Macro Average: {f1_macro_scratch}')
```

```
precision    recall  f1-score   support
```

```
0      0.20      0.08      0.12        12
1      0.00      0.00      0.00        11
2      0.31      0.30      0.31       74
3      0.58      0.73      0.64      189
4      0.46      0.49      0.48        90
5      0.99      0.98      0.99       235
6      0.90      0.87      0.89      313
7      0.69      0.49      0.58        69
8      0.50      0.14      0.22         7
```

```
accuracy          0.74      1000
macro avg         0.52      0.45      0.47      1000
weighted avg      0.74      0.74      0.74      1000
```

Perbandingan antara KNN dari Scratch dan KNN dari Scikit-learn menunjukkan bahwa meskipun keduanya menggunakan konsep dasar yang sama, yaitu menghitung jarak dan memilih tetangga terdekat untuk prediksi, Scikit-learn lebih efisien dan cepat berkat optimasi yang sudah dilakukan, termasuk penggunaan multi-threading. Hal ini menghasilkan akurasi yang lebih tinggi dan waktu eksekusi yang lebih cepat pada Scikit-learn, terutama saat menangani dataset besar. Sementara itu, KNN dari Scratch memberikan pemahaman lebih mendalam tentang algoritma, namun memerlukan optimasi lebih lanjut untuk meningkatkan performa, seperti mengurangi waktu eksekusi yang lebih lama dibandingkan dengan implementasi built-in dari Scikit-learn.

2. Naive Bayes

Scratch	Scikit-learn
---------	--------------

<pre>[] from GNB import GaussianNaiveBayes gnb_scratch= GaussianNaiveBayes() gnb_scratch.fit(X_train_transformed, y_train) y_pred= gnb_scratch.predict(X_val_transformed) f1_macro_scratch = f1_score(y_val,y_pred,average = 'macro') print(f'F1 Macro Average: {f1_macro_scratch}')</pre> <p>F1 Macro Average: 0.15689196303232086</p>	<pre>gnb= GaussianNB() gnb.fit(X_train_transformed, y_train) y_pred= gnb.predict(X_val_transformed) f1_macro_sklearn = f1_score(y_val,y_pred,average = 'macro') print(f'F1 Macro Average: {f1_macro_sklearn}')</pre> <p>F1 Macro Average: 0.21145610600630813</p>
--	--

Keduanya memiliki hasil F1 Macro Average yang berbeda dimana scikit-learn memiliki skor yang lebih tinggi dibanding dari scratch yaitu 0.211 berbanding dengan 0.156. Perbandingan antara Gaussian Naïve Bayes (GNB) dari Scratch dan GNB dari Scikit-learn menunjukkan bahwa keduanya menggunakan pendekatan yang serupa dalam mengklasifikasikan data, yaitu dengan menghitung statistik penting untuk tiap kelas dan menggunakan distribusi Gaussian untuk memprediksi kelas dengan probabilitas tertinggi. Namun, GNB dari Scikit-learn memiliki keunggulan dalam hal akurasi karena menangani edge-case dengan lebih baik, seperti pembagian dengan nol menggunakan epsilon dan penanganan underflow pada probabilitas yang sangat kecil. Selain itu, GNB dari Scikit-learn lebih efisien dalam waktu eksekusi dan lebih stabil, terutama pada dataset besar dengan banyak fitur. Untuk meningkatkan model yang dibuat dari scratch, diperlukan optimasi, penanganan edge-case yang lebih baik, serta pemilihan teknik encoding yang lebih sesuai, seperti label encoding, untuk menghindari masalah curse of dimensionality yang menghambat pelatihan model.

3. ID3

Scratch	Scikit-learn
<pre>[] from DecisionTree import DecisionTree dt_data= X_train_transformed[:1000] dt_data_label= y_train[:1000] dt_test= X_val_transformed[:100] dt_test_label= y_val[:100] [] dt_scratch= DecisionTree(mode='entropy') dt_scratch.fit(dt_data, dt_data_label) y_pred= dt_scratch.predict(dt_test) f1_macro_scratch = f1_score(dt_test_label,y_pred,average = 'macro') print(f'F1 Macro Average: {f1_macro_scratch}')</pre> <p>F1 Macro Average: 0.4820763561551082</p>	<pre>dt= DecisionTreeClassifier(criterion = "entropy", max_depth= None, min_samples_leaf=4) dt.fit(dt_data, dt_data_label) y_pred= dt.predict(dt_test) f1_macro_sklearn = f1_score(dt_test_label,y_pred,average = 'macro') print(f'F1 Macro Average: {f1_macro_sklearn}')</pre> <p>F1 Macro Average: 0.4704355314133676</p>

Kedua algoritma memiliki F1 Macro Average yang hampir sama yaitu Scratch

bernilai 0.48 sedangkan scikit-learn bernilai 0.47. Perbandingan antara Decision Tree dari Scratch dan Decision Tree dari Scikit-learn menunjukkan bahwa keduanya memiliki prinsip dasar yang sama. Namun, Decision Tree dari Scikit-learn memiliki keunggulan dalam hal akurasi karena optimisasi yang lebih baik (walaupun pada kasus ini scikit-learn memperoleh skor lebih kecil, kemungkinan besar karena hyperparameter default yang ada pada scikit-learn selain yang di-*define* menyebabkan model scikit learn *overfit*), termasuk penanganan kesalahan presisi floating-point dan pruning yang lebih efisien, serta manajemen hyperparameter yang lebih efektif. Selain itu, model dari Scikit-learn lebih efisien dalam waktu eksekusi, sementara model dari scratch cenderung lebih lambat dan kurang optimal.

Bab 6

Kontribusi Kelompok

Mohammad Nugraha Eka Prawira (13522001)	Pre Processing, Laporan
Muhammad Al Thariq Fairuz (13522027)	Model from scratch, pre processing, Laporan
Randy Verdian (13522067)	EDA, Pre processing, Laporan
Emery Fathan Zwageri (13522079)	AutoEncoders, preprocessing, kaggle gaming, Laporan

Bab 7

Referensi

[KNN](#)

[GNB](#)

[DTL](#)

[KNeighborsClassifier — scikit-learn 1.6.0 documentation](#)

[Gaussian Naive Bayes - GeeksforGeeks](#)

[GaussianNB — scikit-learn 1.6.0 documentation](#)

[Decision Tree - GeeksforGeeks](#)

[DecisionTreeClassifier — scikit-learn 1.6.0 documentation](#)

<https://www.analyticsvidhya.com/blog/2021/06/dimensionality-reduction-using-autoencoders-in-python/>

[Autoencoders -Machine Learning - GeeksforGeeks](#)