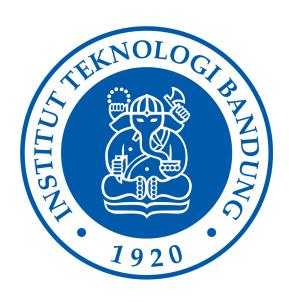
Laporan Tugas Besar 2 IF3170 Inteligensi Artifisial Implementasi Algoritma Pembelajaran Mesin Semester I Tahun 2024/2025



Disusun oleh:

1.	Maulvi Ziadinda Maulana	13522122
2.	Ahmad Rafi Maliki	13522137
3.	Nicholas Reymond Sihite	13522144
4.	Albert Ghazaly	13522150

INSTITUT TEKNOLOGI BANDUNG 2024

Daftar Isi

Daftar Isi	2
Bagian 1. Deskripsi Permasalahan	3
Bagian 2. Implementasi Algoritma	4
2.1. K-Nearest Neighbour (KNN)	4
2.2. Gaussian Naive Bayes (NB)	6
2.3. Iterative Dichotomiser 3 (ID3)	11
Bagian 3. Pembersihan dan Prapemrosesan Data	17
3.1. Pembersihan Data	17
3.1.1. Penanganan Nilai Kosong	17
3.1.2. Penanganan Nilai Outlier	20
3.1.3. Penanganan Nilai Duplikat	20
3.1.4. Rekayasa Fitur	20
3.2. Prapemprosesan Data	22
3.2.1. Feature Scaling.	22
3.2.2. Feature Encoding.	23
3.2.3. Handling Imbalanced Dataset	25
3.2.4. Normalisasi Data	27
3.2.5. Pengurangan Dimensi Data	27
3.3. Pipeline per Algoritma	29
3.3.1. Pipeline k-Nearest Neighbour	29
3.3.2. Pipeline Gaussian Naive-Bayes	30
3.3.3. Pipeline Iterative Dichotomiser 3	30
Bagian 4. Percobaan dan Analisis	32
4.1. Percobaan	32
4.2. Analisis	38
Pembagian Tugas	39
Referensi	40

Bagian 1. Deskripsi Permasalahan



Gambar 1. Ilustrasi Pembelajaran Mesin

(sumber: https://mvjce.edu.in/blog/future-of-machine-learning-trends-applications/)

Pembelajaran mesin merupakan salah satu cabang dari kecerdasan buatan yang memungkinkan sistem untuk belajar dari data dan membuat prediksi atau keputusan tanpa diprogram secara eksplisit.

Dataset UNSW-NB15 adalah kumpulan data lalu lintas jaringan yang mencakup berbagai jenis serangan siber dan aktivitas normal. Pada tugas ini, mahasiswa diminta untuk mengimplementasikan algoritma pembelajaran mesin yang telah dipelajari di kuliah, yaitu KNN, Gaussian Naive-Bayes, dan ID3 pada dataset UNSW-NB15.

Bagian 2. Implementasi Algoritma

2.1. K-Nearest Neighbour (KNN)

Algoritma ini diimplementasikan dalam sebuah kelas yang bernama KNN dengan 3 atribut wajib dan 1 atribut opsional, yaitu: *integer* 'k' untuk menyimpan banyaknya *nearest neighbour* yang diinginkan; *string* 'distance_metric' untuk menyimpan jenis metrik pengukuran distance; DataFrame 'data' untuk menyimpan data *train*; dan atribut opsional *integer* 'p' untuk menyimpan nilai p apabila menggunakan metrik Minkowski Distance.

Konfigurasi algoritma (pengaturan nilai 'k', 'distance_metric', dan 'p') dilakukan pada konstruktor kelas. *Konfigurasi* data *train* dilakukan pada metode 'fit' agar sesuai dengan *pipeline* yang dibuat. Untuk memprediksi kelas sebuah DataFrame, kelas ini menggunakan sebuah metode yang bernama 'predict' yang menerima parameter 'data_to_predict' berupa DataFrame.

Proses yang terjadi pada metode ini adalah sebagai berikut: memisahkan *feature* target, yaitu attack_cat, dengan *feature* lainnya dari data *train*; mengambil nilai-nilai yang ada pada DataFrame (*train* dan *test*) sehingga tipe datanya menjadi np.ndarray; melakukan kalkulasi jarak antara data-data pada data *train* dan *test* dengan bantuan fungsi 'cdist' dari Scipy; mengambil sebanyak k data dari data *train* dengan 'distance' terdekat terhadap setiap data pada data *test*; mengambil kelas mayoritas dari setiap k data tersebut; dan mengembalikan sebuah *array* yang berisi prediksi kelas data *test*. Implementasi kelas KNN dapat dilihat pada potongan kode berikut.

```
class KNN(BaseEstimator, ClassifierMixin):
    def __init__(self) -> None:
        self.k = 5
        self.distance_metric = "3"

        if ((self.distance_metric == "minkowski") or
(self.distance_metric == "3")):
            self.p = 10

    def calculateManhattanDistance(self, test_data: np.ndarray,
train_data: np.ndarray):
        distance = cdist(test_data, train_data, metric='cityblock')
        return distance

    def calculateEuclideanDistance(self, test_data: np.ndarray,
train_data: np.ndarray):
        distance = cdist(test_data, train_data, metric='euclidean')
        return distance
```

```
def calculateMinkowskiDistance(self, test data: np.ndarray,
train data: np.ndarray):
       distance = cdist(test data, train data, metric='minkowski',
p=self.p)
       return distance
    def calculateDistance(self, test data: np.ndarray, train data:
np.ndarray):
        if ((self.distance metric == "manhattan") or
(self.distance metric == "1")):
            return self.calculateManhattanDistance(test data, train data)
        elif ((self.distance metric == "euclidean") or
(self.distance metric == "2")):
            return self.calculateEuclideanDistance(test data, train data)
        elif ((self.distance metric == "minkowski") or
(self.distance metric == "3")):
            return self.calculateMinkowskiDistance(test data, train data)
    def fit(self, X, y):
       X = pd.DataFrame(X)
        y = pd.DataFrame(y)
        train data = pd.concat([X, y], axis=1)
        sample fraction = min(1, 10000 / len(train data))
        train data = train data.groupby('attack cat',
group keys=False).apply(lambda x:
x.sample(frac=sample fraction)).reset index(drop=True)
        self.data = train data
    def predict(self, data to predict: pd.DataFrame):
        class predictions = []
        features = self.data.columns.difference(["attack cat"])
        train data = self.data[features].values
        train labels = self.data['attack cat'].values
        test data = data to predict[features].values
        distances = self.calculateDistance(test data, train data)
        for i in range(len(test data)):
            k nearest indices = distances[i].argsort()[:self.k]
            k_nearest_labels = train_labels[k_nearest_indices]
            unique, counts = np.unique(k nearest labels,
return counts=True)
            class predictions.append(unique[np.argmax(counts)])
        return class predictions
```

2.2. Gaussian Naive Bayes (NB)

Gaussian Naive Bayes merupakan metode pembelajaran mesin variasi dari Naive Bayes Classifier yang digunakan untuk memprediksi kelas dari data yang fiturnya bertipe

continuous. Sama seperti NB yang standar, cara kerja algoritma ini adalah menghitung probabilitas posterior dari tiap kelas label yang diperoleh dari hasil perkalian probabilitas kondisional untuk tiap fitur dan menjadikan kelas dengan probabilitas tertinggi sebagai hasil dari prediksi. Bedanya adalah pada *Gaussian Naive Bayes*, nilai probabilitas posteriornya diperoleh dari hasil kalkulasi yang melibatkan mean dan juga standar deviasi dari data latih.

$$p(C_k \mid \mathbf{x}) = rac{p(C_k) \ p(\mathbf{x} \mid C_k)}{p(\mathbf{x})}$$

Gambar 2. Probabilitas Posterior Naive Bayes

Probabilitas posterior dibaca sebagai probabilitas kelas C_k terjadi bila fitur bernilai x. Persamaan inilah yang akan menjadi landasan untuk menghitung probabilitas gabungan untuk membentuk model Naive Bayes. Probabilitas prior $p(C_k)$ merupakan peluang kejadian data memiliki kelas C_k . Likelihood $p(x|C_k)$ merupakan peluang kejadian data memiliki fitur x jika diketahui kelas C_k . Evidence p(x) merupakan peluang kejadian data memiliki fitur x.

$$egin{aligned} p(C_k \mid x_1, \dots, x_n) & \propto \ p(C_k, x_1, \dots, x_n) \ & = p(C_k) \ p(x_1 \mid C_k) \ p(x_2 \mid C_k) \ p(x_3 \mid C_k) \ \cdots \ & = p(C_k) \prod_{i=1}^n p(x_i \mid C_k) \,, \end{aligned}$$

Gambar 3. Probabilitas Gabungan Naive Bayes

Berikut merupakan cara untuk menghitung nilai probabilitas gabungan dari sebuah kelas. Nilai probabilitas gabungan diperoleh dengan cara mengalikan probabilitas prior dan juga probabilitas kondisional tiap fitur menggunakan *chain rule*. Hal tersebut dilakukan untuk semua kelas yang mungkin dan kelas yang memiliki nilai probabilitas tertinggi dijadikan sebagai hasil dari prediksi.

Seperti yang sudah dijelaskan sebelumnya, perbedaan *Gaussian Naive Bayes* terletak di cara menghitung probabilitas kondisional. Pada kasus ini tidak mungkin menghitung nilai probabilitas kondisional dengan cara yang sudah dijelaskan untuk tiap fitur karena fiturnya tidak bernilai diskrit, tetapi *continuous*.

$$p(x=v\mid C_k) = rac{1}{\sqrt{2\pi\sigma_k^2}}\,e^{-rac{(v-\mu_k)^2}{2\sigma_k^2}}$$

Gambar 4. Probabilitas Kondisional Gaussian Naive Bayes

Pada *Gaussian Naive Bayes*, nilai probabilitas kondisional diperoleh dengan menggunakan formula pada Gambar 4. Formula merupakan formula dari distribusi normal (*gaussian*) yang menggunakan nilai mean dan juga standar deviasi dari fitur terkait. Perhitungan probabilitas gabungan kemudian dilanjutkan seperti halnya Naive Bayes yang standar.

Pada implementasi kode, ada sedikit penyesuaian untuk cara perhitungan probabilitas gabungungan. Karena nilai standar deviasi tidak jarang sengat kecil (diantara nol dan satu), perhitungan perkalian diubah menjadi penjumlahan dalam basis logaritma untuk menghindari terjadinya error akibat perkalian dengan angka yang sangat kecil.

```
class CustomGaussianNB(BaseEstimator, ClassifierMixin):
   def init (self):
       self.smoothing = 1e-9
       self.label encoder = LabelEncoder()
       self.model = {}
       self.classes = None
       self.class priors = None
       self.fitted = False
   def fit(self, X, y):
       if self.fitted:
           return self
       if isinstance(X, pd.DataFrame):
           X = X.to numpy()
       if isinstance(y, pd.Series):
            y = y.to numpy()
       y encoded = self.label encoder.fit transform(y)
        self.classes = self.label encoder.classes
       class counts = np.bincount(y encoded)
       self.class priors = np.log(class counts / len(y encoded))
        self.model = {}
        for label in np.unique(y encoded):
            class data = X[y encoded == label]
            self.model[label] = {}
            for i in range(X.shape[1]):
                feature data = class data[:, i]
                mean = np.mean(feature data)
                std = np.std(feature data)
                std = max(std, self.smoothing)
                self.model[int(label)][int(i)] = {'mean': mean, 'std':
std}
       self.fitted = True
       return self
```

```
def predict(self, X):
        if isinstance(X, pd.DataFrame):
            X = X.to numpy()
        predictions = []
        for row in X:
            log likelihoods = {}
            predicted label encoded = None
            max log likelihood = float('-inf')
            for label in self.model.keys():
                log likelihood = self.class priors[label]
                for i in range(len(row)):
                    mean = self.model[label][i]['mean']
                    std = self.model[label][i]['std']
                    coefficient = -np.log(std * np.sqrt(2 * np.pi))
                    exponent = -0.5 * ((row[i] - mean) / std) ** 2
                    log likelihood += coefficient + exponent
                if predicted label encoded == None or log likelihood >
max log likelihood:
                    max log likelihood = log likelihood
                    predicted label encoded = label
                log likelihoods[label] = log likelihood
            predicted class =
self.label encoder.inverse transform([predicted label encoded])[0]
            predictions.append(predicted class)
        return np.array(predictions)
    def save model(self, model path, download=False):
        if not self.fitted:
            raise ValueError("Model is not fitted. Please fit the model
before saving.")
        model data = {
            'class priors': self.class priors.tolist(),
            'classes': self.classes.tolist(),
            'model': {},
        for label in self.model:
            model data['model'][str(label)] = {}
            for i in self.model[label]:
                model data['model'][str(label)][str(i)] = {
```

```
'mean': str(self.model[label][i]['mean']),
                    'std': str(self.model[label][i]['std'])
                }
        with open (model path, 'w') as f:
            json.dump(model data, f, indent=4)
        print(f"\nModel saved to: {model path}")
        if download:
            files.download(model path)
    def load model (self, model path):
        with open (model_path, 'r') as f:
            model data = json.load(f)
        self.class priors = np.array(model data['class priors'])
        self.classes = np.array(model data['classes'])
        self.model = model data['model']
        new model = {}
        for label str, features in self.model.items():
            label = int(label str)
            new model[label] = {}
            for i str, feature data in features.items():
                i = int(i str)
                feature data['mean'] = float(feature data['mean'])
                feature data['std'] = float(feature data['std'])
                new model[label][i] = feature data
        self.model = new model
        self.label encoder.classes =
np.array(np.array(model data['classes']))
        self.fitted = True
```

Potongan kode di atas merupakan implementasi kami terhadap Gaussian Naive Bayes dengan mekanisme seperti yang sudah dijelaskan sebelumnya. Implementasi kami wujudkan dalam bentuk kelas CustomGaussianNB yang meng-inherit kelas BaseEstimator dan ClassifierMixin dari pustaka sklearn. Kelas ini memiliki lima buah atribut:

1. smoothing

Yaitu nilai "penghalus" untuk mengganti nilai standar deviasi yang bernilai nol akan diubah menjadi nilai smoothing ini yang merupakan nilai sangat kecil yaitu 10^-9.

2. label encoder

Yaitu objek dari pustaka sklearn yang akan digunakan untuk mengkonversi nilai label dari *string* ke *integer*.

3. model

Yaitu sebuah data bertipe *dictionary* untuk menyimpan hasil pembelajaran.

4. classes

Yaitu sebuah *list* yang menyimpan semua kelas label yang mungkin.

5. class priors

Yaitu sebuah *list* yang menyimpan nilai probabilitas prior dari tiap kelas.

6. fitted

Yaitu sebuah variabel bertipe *boolean* untuk menandakan apakah objek sudah belajar atau belum. Nilai atribut ini dapat berubah setelah melakukan pembelajaran atau setelah melakukan pemuatan hasil pembelajaran dari *external file*.

Selain atribut yang disebutkan di atas, kelas ini juga memiliki lima buah metode:

1. init (self)

Metode konstruktor.

2. fit(self, X, y)

Metode ini merupakan metode yang dipanggil untuk melakukan pembelajaran. Terdapat dua buah argumen yaitu *X* yang merupakan *list of tuples* yang tiap *tuple* nya merupakan nilai fitur dari data dan juga *y* yang merupakan kelas label dari *tuple* yang sejajar. Method ini akan mempopulasi nilai dari atribut *models*.

Setelah metode ini dijalankan akan dihasilkan *models* yang merupakan sebuah *dictionary* dengan data nilai mean dan standar deviasi untuk setiap fitur dari setiap label.

3. predict(self, X)

Metode ini merupakan metode yang dipanggil untuk melakukan prediksi kelas dari sebuah dataset. Terdapat sebuah argumen *X* yang merupakan *list of tuples* yang tiap *tuple* nya merupakan nilai fitur dari data. Metode ini mengembalikan list label kelas hasil prediksi.

4. save model(self, model path, download=False)

Metode ini digunakan untuk menyimpan model hasil pembelajaran ke *external* file.

5. load model(self, model path)

Metode ini digunakan untuk memuat model hasil pembelajaran dari external file.

2.3. Iterative Dichotomiser 3 (ID3)

Algoritma *Iterative Dichotomiser 3* (ID3) merupakan salah satu jenis algoritma yang digunakan untuk membuat *decision tree* dari sebuah dataset. Algoritma ini bekerja dengan memanfaatkan beberapa parameter atau variabel penilaian (*measurement*), yakni *Entropy* (*H*) dan *Information Gain* (*IG*). *Entropy* adalah suatu pengukuran (*measurement*) ketidakpastian pada suatu data atau set.

Dalam mengukur Entropy, jika kita memiliki data atau set S, Entropy dari set tersebut H(S) dapat dihitung:

$$\mathrm{H}(S) = \sum_{x \in X} -p(x) \log_2 p(x)$$

Gambar 5 Perhitungan Entropy pada Set S

Penjelasan dari persamaan tersebut adalah Entropy dari set *S* adalah penjumlahan dari negatif proporsi kelas *x* dari seluruh set *S* dikali *Log* 2 dari proporsi kelas *x* dari seluruh set *S*. Hasil dari perhitungan tersebut (*Entropy*) akan merepresentasikan ketidakpastian dan digunakan sebagai salah satu komponen dalam menghitung *Information Gain* (*IG*). *Information Gain* sendiri merupakan komponen yang digunakan untuk menentukan atribut terbaik yang akan digunakan sebagai *node* dalam tree.

Dalam menghitung nilai *Information Gain*, jika kita memiliki atribut A pada sebuah set atau data S, *Information Gain* dari atribut tersebut IG(S, A), dapat dihitung:

$$IG(S,A) = \mathrm{H}(S) - \sum_{t \in T} p(t) \mathrm{H}(t) = \mathrm{H}(S) - \mathrm{H}(S|A).$$

Gambar 6 Perhitungan Entropy pada Atribut A

Penjelasan dari persamaan tersebut adalah *Information Gain* dari atribut A adalah hasil perhitungan Entropy S, H(S), dikurang penjumlahan dari perkalian proporsi dari kelas t dan total data di set S dan Entropy dari kelas t. Persamaan tersebut dapat juga dinyatakan sebagai Entropy dari set S dikurang Entropy S terhadap atribut A.

```
class ID3DecisionTree(BaseEstimator, ClassifierMixin):
    def __init__(self):
        self.tree: Optional[Dict[str, Any]] = None
        self.label: str = ""

    def calculate_entropy(self, data: pd.DataFrame) -> float:
        _, counts = np.unique(data.iloc[:, -1], return_counts=True)
        probabilities = counts / counts.sum()
        return -np.sum(probabilities * np.log2(probabilities + 1e-9))

    def calculate_information_gain(self, data: pd.DataFrame, feature:
    str) -> float:
```

```
total entropy = self.calculate entropy(data)
        unique values = data[feature].unique()
        weighted entropy = 0
        for value in unique values:
            subset = data[data[feature] == value]
            weighted entropy += (len(subset) / len(data)) *
self.calculate entropy(subset)
        return total entropy - weighted entropy
    def find_best_splits(self, data: pd.DataFrame, feature: str,
max splits: int) -> List[float]:
        sorted data = data.sort values(feature)
        unique values = sorted data[feature].unique()
        print(f"Unique value length of column '{feature}' is
{len(unique values)}")
        if len(unique values) <= 1:
            return []
        max splits = min(max splits, len(unique values) - 1)
        split points = []
        last split = None
        for i in range(1, max splits + 1):
            quantile = i / (max splits + 1)
            split value = np.quantile(unique values, quantile)
            if last split is None or not np.isclose(split value,
last split):
                split points.append(split value)
                last split = split value
        split points = sorted(set(split points))
        if len(split points) < 1:
            return []
        return split points
    def calculate information gain on split(self, data: pd.DataFrame,
left: pd.DataFrame, right: pd.DataFrame) -> float:
        total entropy = self.calculate entropy(data)
        left weight = len(left) / len(data)
        right_weight = len(right) / len(data)
        return total_entropy - (left_weight *
self.calculate entropy(left) + right weight *
self.calculate entropy(right))
    def convert numeric to categorical(self, data: pd.DataFrame) ->
pd.DataFrame:
```

```
for column in data.columns[:-1]:
            if pd.api.types.is numeric dtype(data[column]):
                print(f"Processing numeric column: {column}")
                unique values = data[column].unique()
                if len(unique values) < 20:
                    max split = len(unique values)
                else:
                    \max \text{ split} = \inf(\text{len}(\text{data}) / 100)
                min val = data[column].min()
                max val = data[column].max()
                split_points = self.find_best_splits(data, column,
max splits=max split)
                if split points:
                    bins = [min val] + split points + [max val]
                    bins = sorted(set(bins))
                    labels = [f'<= {split points[0]}'] + \</pre>
                             [f'{split_points[i]} - {split_points[i + 1]}'
for i in range(len(split points) - 1)] + \
                             [f'> {split points[-1]}']
                    data[column] = pd.cut(data[column], bins=bins,
labels=labels, include lowest=True)
        return data
    def best feature to split(self, data: pd.DataFrame) -> str:
        gains = {feature: self.calculate information gain(data, feature)
for feature in data.columns[:-1]}
        return max(gains, key=gains.get)
    def build tree(self, data: pd.DataFrame) -> Dict[str, Any]:
        target = data.columns[-1]
        if len(data[target].unique()) == 1:
            return data[target].iloc[0]
        if len(data.columns) == 1:
            return data[target].mode().iloc[0]
        best feature = self.best feature to split(data)
        print(f"Feature {best feature} processed")
        tree = {best feature: {}}
        majority class = data[target].mode().iloc[0]
        for value in data[best feature].unique():
            subset = data[data[best feature] ==
value].drop(columns=[best feature])
            tree[best feature][value] = self.build tree(subset)
        tree[best feature]["unknown"] = majority class
```

```
return tree
   def fit(self, X: pd.DataFrame, y: pd.Series) -> Dict[str, Any]:
       self.label = 'attack cat'
       if isinstance(X, np.ndarray):
            X = pd.DataFrame(X)
            print(f"Converted X to pandas DataFrame. Type of X:
{type(X)}; Shape of X: {X.shape}")
        if len(X) != len(y):
            raise ValueError("The number of rows in X and the number of
elements in y must be the same.")
       data = X.copy()
       print(f"Target column 'attack_cat':\n{y.head()}")
       try:
           data[self.label] = y.values
            print(f"Data after adding target column:\n{data.head()}")
        except Exception as e:
            print(f"Error occurred while adding target column: {e}")
       processed data = self.convert numeric to categorical(data)
       print("Building tree")
       columns = [col for col in processed data.columns if col !=
self.label]
       columns.append(self.label)
       processed data = processed data[columns]
       print(f"Columns after reordering: {processed data.columns}")
        self.tree = self.build tree(processed data)
   def predict sample(self, tree: dict, sample: pd.Series) -> Any:
       while isinstance(tree, dict):
            feature = next(iter(tree))
            value = sample.get(feature)
            if isinstance (value, pd. Series):
                value = value.iloc[0]
            if pd.isna(value):
                return "default class"
            branches = tree[feature]
            found match = False
            for key in branches:
                if isinstance(key, str) and ' - ' in key:
                    lower bound, upper bound = key.split(' - ')
                    lower bound = float(lower bound)
                    upper bound = float(upper bound)
```

```
if isinstance(value, (int, float)):
                        if lower bound <= value <= upper bound:
                            tree = branches[key]
                             found match = True
                            break
                    else:
                        continue
                elif isinstance(key, str) and any(op in key for op in
['<=', '>=', '<', '>', '==', '!=']):
                    operator, threshold = key.split(' ', 1)
                    threshold = float(threshold) if operator not in
['==', '!='] else threshold
                    if isinstance(value, (int, float)):
                        if operator == '<=':</pre>
                            if value <= threshold:</pre>
                                tree = branches[key]
                                 found match = True
                                break
                        elif operator == '<':
                             if value < threshold:
                                 tree = branches[key]
                                 found match = True
                                break
                        elif operator == '>=':
                            if value >= threshold:
                                 tree = branches[key]
                                 found match = True
                                break
                        elif operator == '>':
                            if value > threshold:
                                 tree = branches[key]
                                 found match = True
                                break
                        elif operator == '==':
                            if value == threshold:
                                 tree = branches[key]
                                 found match = True
                                break
                        elif operator == '!=':
                            if value != threshold:
                                 tree = branches[key]
                                 found match = True
                                break
                    else:
                        continue
                elif isinstance(value, (str, int, float)):
                    if value == key:
                        tree = branches[key]
                        found match = True
                        break
                else:
```

```
continue
            if not found match:
                if "unknown" in branches:
                    # return "Normal"
                    return branches["unknown"]
                else:
                    return "default class"
        return tree
   def predict(self, df: pd.DataFrame) -> pd.Series:
     if isinstance(df, np.ndarray):
          df = pd.DataFrame(df)
          print(f"Converted df to pandas DataFrame. Type of df:
{type(df)}; Shape of df: {df.shape}")
     predictions = df.apply(lambda row: self.predict sample(self.tree,
row), axis=1)
     return predictions
   def save_model(self, file_path: str) -> None:
       if self.tree is None:
            raise ValueError("Model is not trained yet, cannot save.")
       with open(file_path, 'w') as f:
            json.dump(self.tree, f, indent=4)
   def load model(self, file path: str) -> None:
        with open(file_path, 'r') as f:
            self.tree = json.load(f)
```

Bagian 3. Pembersihan dan Prapemrosesan Data

3.1. Pembersihan Data

Pembersihan data adalah langkah awal yang paling penting dalam mempersiapkan dataset untuk melatih sebuah model Artificial Intelligence. Data awal seringkali mengandung berbagai kesalahan termasuk nilai kosong, nilai yang duplikat, dan data yang inkonsisten. Proses pembersihan data meliputi menangani nilai kosong dengan imputasi atau penghapusan; menangani outlier dengan menghapus atau mentransformasi nilai; memvalidasi data agar konsisten dan sesuai tipe yang diharapkan; menghapus data duplikat yang dapat mempengaruhi hasil pelatihan; serta melakukan pengubahan fitur, seperti normalisasi, pengubahan skala nilai, atau encoding, untuk meningkatkan relevansi data dalam model.

3.1.1. Penanganan Nilai Kosong

Untuk menangani nilai kosong, kami melakukan *Missing Value Analysis* terlebih dahulu. Berdasarkan analisis pada tabel, beberapa kolom memiliki persentase nilai yang hilang sekitar 5%, dengan tipe data yang beragam, termasuk numerik (*float* dan *integer*) serta kategori (nominal). Kolom-kolom numerik, seperti *dur*, *sbytes*, dan *dbytes*, memiliki nilai yang hilang yang signifikan namun tidak dominan, sedangkan kolom kategori, seperti statedan service, juga menunjukkan pola kehilangan serupa. Berikut adalah hasil lengkap dari *Missing Value Analysis* yang kami lakukan.

Tabel 1. Hasil Missing Value Analysis

No	Name	Missing Count	Missing Percentage	Туре
0	state	8805	5.021644	nominal
1	dur	8722	4.974307	float
2	sbytes	8561	4.882486	integer
3	dbytes	8869	5.058144	integer
4	sttl	8825	5.03305	integer
5	dttl	8654	4.935526	integer
6	sloss	8794	5.01537	integer
7	dloss	8978	5.120308	integer
8	service	8791	5.013659	nominal
9	sload	8786	5.010808	float
10	dload	8837	5.039894	float
11	spkts	8654	4.935526	integer

12	dpkts	8686	4.953776	integer
14	is_sm_ips_ports	8746	4.987995	binary
15	ct_state_ttl	8635	4.92469	integer
16	ct_flw_http_mthd	8647	4.931533	integer
17	is_ftp_login	8647	4.931533	binary
18	ct_ftp_cmd	8842	5.042745	integer
19	ct_srv_src	8851	5.047878	integer
20	ct_srv_dst	8774	5.003964	integer
21	ct_dst_ltm	8738	4.983432	integer
22	ct_src_ltm	8823	5.031909	integer
23	ct_src_dport_ltm	8775	5.004534	integer
24	ct_dst_sport_ltm	8788	5.011948	integer
25	ct_dst_src_ltm	8895	5.072972	integer
26	swin	8740	4.984573	integer
27	dwin	8779	5.006815	integer
28	stcpb	8672	4.945791	integer
29	dtcpb	8803	5.020503	integer
30	smean	8788	5.011948	integer
31	dmean	8855	5.050159	integer
32	trans_depth	8785	5.010237	integer
33	response_body_len	8791	5.013659	integer
34	proto	8826	5.03362	nominal
35	sjit	8738	4.983432	float
36	djit	8846	5.045027	float
37	sinpkt	8707	4.965752	float
38	dinpkt	8734	4.981151	float
39	teprtt	8836	5.039323	float
40	synack	8736	4.982292	float
41	ackdat	8595	4.901877	float

Kami merancang kelas *CustomImputer* untuk menangani nilai yang kosong pada data. Untuk kolom numerik, kami mengganti nilai dengan rata-rata (mean) agar tidak mengubah distribusi data secara signifikan. Sedangkan untuk kolom kategori, kami mengganti nilai yang kosong dengan nilai yang paling sering muncul. Selain itu, kami juga mengubah nilai '-' menjadi NaN untuk memudahkan pengolahan. Dengan pendekatan ini, kami memastikan bahwa dataset yang dihasilkan bebas dari nilai hilang, konsisten, dan siap untuk analisis atau modeling lebih lanjut.

```
class CustomImputer(BaseEstimator, TransformerMixin):
    def __init__(self, numeric_columns=None,
categorical columns=None):
        self.numeric columns = numeric columns
        self.categorical columns = categorical columns
        self.numeric imputer = SimpleImputer(strategy='mean')
        self.categorical imputer =
SimpleImputer(strategy='most frequent')
    def fit(self, X, y=None):
        if self.numeric columns is None:
            self.numeric columns =
X.select_dtypes(include=['int64', 'float64']).columns.tolist()
        if self.categorical columns is None:
            self.categorical columns =
X.select dtypes(include=['object', 'category']).columns.tolist()
        if len(self.numeric columns) > 0: # Check if
numeric columns is not empty
            self.numeric imputer .fit(X[self.numeric columns])
        if len(self.categorical columns) > 0: # Check if
categorical columns is not empty
self.categorical imputer .fit(X[self.categorical columns])
        return self
    def transform(self, X):
        X transformed = X.copy()
        X transformed.replace('-', np.nan, inplace=True)
        if len(self.numeric columns) > 0: # Check if
numeric columns is not empty
            X transformed[self.numeric columns] =
self.numeric imputer .transform(X transformed[self.numeric columns]
        if len(self.categorical columns) > 0: # Check if
categorical columns is not empty
            X transformed[self.categorical columns] =
self.categorical imputer .transform(X transformed[self.categorical
```

```
columns])
    return X_transformed

def fit_transform(self, X, y=None):
    return self.fit(X, y).transform(X)
```

3.1.2. Penanganan Nilai Outlier

Dalam analisis dataset ini, kami memutuskan untuk tidak melakukan penanganan outlier karena mempertimbangkan kasus yang berkaitan dengan serangan siber. Kami berpikir bahwa *outlier*, yang biasanya dianggap sebagai data yang menyimpang, justru dapat memiliki makna yang signifikan dalam kasus ini. Dalam serangan siber, outlier sering kali merepresentasikan aktivitas yang tidak normal atau pola serangan yang berbeda dari lalu *lintas* jaringan yang biasa. Misalnya, lonjakan nilai pada kolom seperti *sbytes* atau *dbytes* dapat mencerminkan anomali seperti serangan *Denial of Service* (DoS).

Dengan tidak menghapus atau mengubah outlier, kami memastikan bahwa informasi penting yang dapat mengindikasikan serangan tidak hilang. Penanganan yang salah terhadap outlier berisiko menghilangkan pola penting yang dapat membantu dalam mendeteksi atau mengklasifikasikan serangan. Oleh karena itu, pendekatan ini kami ambil untuk menjaga integritas data dan memastikan bahwa analisis atau model klasifikasi yang dibangun dapat menangkap anomali dengan lebih baik, yang merupakan inti dari deteksi serangan siber.

3.1.3. Penanganan Nilai Duplikat

Menghapus duplikat data merupakan langkah penting untuk menjaga integritas data dan memastikan analisis yang akurat. Duplikasi dapat menyebabkan bias dalam model *Artificial Intelligence*, meningkatkan risiko *overfitting*, dan memperbesar ukuran dataset secara tidak perlu, sehingga meningkatkan biaya komputasi. Fungsi *removeDuplicates* menggunakan metode *drop_duplicates*() untuk menghapus entri yang sama dalam dataset, memastikan data lebih bersih, efisien, dan konsisten untuk analisis lebih lanjut.

```
def removeDuplicates(df):
    df = df.drop_duplicates()
    return df
```

3.1.4. Rekayasa Fitur

Rekayasa fitur atau *Feature engineering* adalah proses membuat, mengubah, atau mengurangi fitur dalam dataset untuk meningkatkan kemampuan model

Artificial Intelligence dalam belajar pola dan membuat prediksi yang akurat. Berikut adalah kelas FeatureEngineer yang kami buat.

```
class FeatureEngineer(BaseEstimator, TransformerMixin):
    def __init__(self, numeric columns=None,
categorical columns=None):
        self.numeric columns = numeric_columns
        self.categorical columns = categorical columns
        self.columns to drop = []
    def fit(self, X, y=None):
        # Ensure target column is passed as y
        if y is None:
            raise ValueError("The target variable 'y' must be
provided for the feature engineering process.")
        # Step 1: Identify numeric columns with low variance
        if len(self.numeric columns) > 0:
            numeric data = X[self.numeric columns]
            variances = numeric data.var()
            low variance columns = variances[variances <</pre>
0.01].index.tolist()
        else:
            low variance columns = []
        # Step 2: Remove irrelevant columns
        irrelevant columns = ['id']
        # Step 3: Feature importance analysis (numeric vs
categorical target)
        if len(self.numeric columns) > 0:
            numeric data = X[self.numeric columns]
            anova_f_values, _ = f_classif(numeric_data, y)
            low importance columns = [
                col for col, score in zip(self.numeric columns,
anova f values) if score < 0.1
        else:
            low importance columns = []
        # Step 4: Redundancy analysis (highly correlated numeric
features)
        if len(self.numeric columns) > 0:
            corr matrix = X[self.numeric columns].corr().abs()
            upper triangle =
corr matrix.where (np.triu (np.ones (corr matrix.shape),
k=1).astype(bool))
            redundant columns = [column for column in
upper triangle.columns if any(upper triangle[column] > 0.9)]
            redundant columns = []
        # Combine all columns to drop
        self.columns_to_drop = list(set(
            low_variance_columns + irrelevant_columns +
low importance columns + redundant columns
```

```
return self

def transform(self, X):
    # Drop identified columns from X
    X_transformed = X.drop(columns=self.columns_to_drop,
errors='ignore')
    return X_transformed

def fit_transform(self, X, y=None):
    return self.fit(X, y).transform(X)
```

Kelas *FeatureEngineer* yang kami buat ini bertujuan untuk mengidentifikasi fitur yang tidak relevan, memiliki variansi rendah, atau redundan, sekaligus memastikan bahwa fitur yang digunakan memiliki hubungan signifikan dengan fitur target yaitu *attack cat*.

Langkah-langkah dalam kelas ini digunakan untuk mengoptimalkan relevansi fitur, mengurangi risiko *overfitting*, dan meningkatkan efisiensi pelatihan model. Identifikasi fitur dengan variansi rendah dilakukan karena fitur tersebut tidak memberikan informasi yang cukup untuk membedakan target, sehingga akan menyederhanakan model. Kolom yang tidak relevan, seperti 'id', dihapus karena tidak memiliki kontribusi terhadap prediksi model. Kami juga menganallisis pentingnya fitur menggunakan ANOVA F-test mengevaluasi hubungan antara fitur numerik dan target untuk mengidentifikasi fitur dengan pengaruh kecil. Selain itu, analisis korelasi digunakan untuk menghapus fitur yang sangat berkorelasi (>0.9) guna mengurangi redundansi informasi. Pendekatan berbasis statistik dan korelasi ini memastikan dataset berkualitas lebih baik sebelum modeling.

3.2. Prapemprosesan Data

3.2.1. Feature Scaling

Feature scaling adalah teknik prapemrosesan data yang bertujuan untuk menyamakan skala dari fitur-fitur dalam dataset sehingga algoritma *Artificial Intelligence* dapat bekerja secara efektif. Skala fitur yang berbeda dapat mempengaruhi kinerja model, terutama untuk algoritma yang sensitif terhadap skala seperti *K-Nearest Neighbors* (KNN). Beberapa metode umum dalam *feature scaling* adalah normalisasi (Min-Max Scaling), standardisasi (Z-score Scaling), *robust scaling* (berbasis IQR), dan *Log Transformation*.

Kelas *FeatureScaler* yang kami buat dirancang untuk mengotomatiskan proses *scaling* pada fitur numerik dalam dataset. Dengan parameter scaler, tiap algoritma dapat menentukan metode *scaling* yang diinginkan, seperti *StandardScaler*, *MinMaxScaler*, atau lainnya dari library sklearn. Jika tidak diberikan, fungsi secara

default menggunakan *StandardScaler*, yang merupakan metode standarisasi berbasis Z-score. Kelas ini akan memastikan bahwa fitur numerik dalam dataset memiliki skala yang seragam. Berikut adalah implementasi kelas yang kami buat.

3.2.2. Feature Encoding

Data yang digunakan untuk pembelajaran mesin seringkali memiliki tipe data yang berbeda. Tipe data string merupakan tipe data yang tidak cocok jika digunakan untuk melakukan perhitungan. Maka dari itu, fitur dengan tipe data ini kami *encode* menjadi bentuk numerik yang dapat lebih mudah diolah. Beberapa metode *encoding* yang kami implementasikan adalah *one-hot encoding*, *label encoding*, *dan target encoding*. Berikut adalah kelas *FeatureEncoder* yang kami buat.

```
class FeatureEncoder(BaseEstimator, TransformerMixin):
   def init (self, categorical columns=None,
encoding strategy='onehot'):
        self.categorical columns = categorical columns
        self.encoding strategy = encoding strategy
       self.encoders = {}
   def fit(self, X, y=None):
       X = X.copy()
        if self.encoding strategy not in ['target', 'label',
'onehot'l:
            raise ValueError("encoding strategy must be one of
'target', 'label', or 'onehot'")
        # Automatically detect categorical columns if not specified
        if self.categorical columns is None:
            self.categorical_columns =
X.select dtypes(include=['object', 'category']).columns.tolist()
```

```
# Fit the encoder for each column
        for col in self.categorical columns:
            if self.encoding strategy == 'onehot':
                enc = OneHotEncoder(sparse=False,
handle unknown='ignore')
            elif self.encoding strategy == 'label':
                enc = LabelEncoder()
            elif self.encoding strategy == 'target':
                if y is None:
                    raise ValueError("Target encoding requires the
target variable y.")
                enc = TargetEncoder()
            if self.encoding strategy == 'label':
                enc.fit(X[col].astype(str))
            else:
                enc.fit(X[[col]].astype(str), y)
            self.encoders[col] = enc # Save the fitted encoder for
the column
       return self
    def transform(self, X):
       X = X.copy()
        for col, enc in self.encoders.items():
            if self.encoding strategy == 'label':
                transformed =
pd.Series(enc.transform(X[col].astype(str)), index=X.index,
name=col)
                X[col] = transformed
            else:
                transformed = enc.transform(X[[col]].astype(str))
                # Convert the transformed array or sparse matrix
into a DataFrame
                if hasattr(transformed, "toarray"): # Handle
sparse matrices
                    transformed = pd.DataFrame(
                        transformed.toarray(),
                        columns=[f"{col}_{i}" for i in
range(transformed.shape[1])],
                        index=X.index,
                elif isinstance(transformed, np.ndarray): # Handle
numpy arrays
                    transformed = pd.DataFrame(
                        transformed,
                        columns=[f"{col} {i}" for i in
range(transformed.shape[1])],
                        index=X.index,
                elif isinstance(transformed, pd.DataFrame): #
Handle DataFrames
```

3.2.3. Handling Imbalanced Dataset

Dataset yang tidak seimbang artinya sebagian besar dari *tuple* data memiliki kelas yang sama dan ada kelas yang hanya diwakili oleh sedikit *tuple*. Hal ini dapat mempengaruhi proses pembelajaran yang mengakibatkan model hanya bagus dalam memprediksi kelas tertentu. Ada dua metode utama untuk mengatasi masalah ini, *oversampling* dan *undersampling*.

Oversampling dilakukan ketika ingin menyeimbangkan data dengan kelas minoritas. Hal ini dapat dicapai dengan menggunakan kelas SMOTE dari pustaka sklearn. Kelas ini mampu menciptakan data sintetis dari data minoritas yang ada sehingga frekuensi data dengan kelas minoritas tadi dapat menyaingi kelas mayoritas. Metode ini dapat mengurangi overfitting terhadap kelas mayoritas dan sering dipakai ketika ada ketidak seimbangan data yang signifikan.

Undersampling dilakukan ketika ingin menyeimbangkan data dengan kelas mayoritas. Hal ini dapat dicapai dengan kelas RandomUnderSampler dari pustaka sklearn. Kelas ini bekerja dengan cara mengurangi data dengan kelas mayoritas pada dataset secara acak sehingga jumlahnya bisa setara dengan kelas minoritas. Keuntungan metode ini adalah mudah dilakukan dan juga mengurangi beban komputasi karena data yang diolah akan lebih sedikit. Metode ini sering digunakan ketika dataset terlalu besar dan perlu penyeimbangan. Berikut adalah implementasi dari SamplingHandler yang kami buat.

```
self.target column = target column
        self.strategy = strategy
        self.sampling ratio = sampling ratio
        self.random state = random state
        self.smote = SMOTE(sampling strategy=sampling ratio,
random state=random state)
    def fit(self, X, y=None):
        # SMOTE requires X and y, so it needs to be initialized
with proper data during fitting.
        if self.strategy in ['oversample', 'both'] and y is not
None:
            self.smote.fit(X, y)
        return self
    def transform(self, X, y=None):
        if y is None:
            raise ValueError("y (target column) must be provided
during transform.")
        X, y = X.copy(), y.copy()
        majority_class = y.value counts().idxmax()
        minority class = y.value counts().idxmin()
        # Perform undersampling
        if self.strategy in ['undersample', 'both']:
            majority data = X[y == majority class]
            minority data = X[y == minority class]
            n majority samples = int(len(minority data) /
self.sampling ratio)
            majority data resampled = resample(
                majority data,
                replace=False,
                n_samples=n majority samples,
                random state=self.random state
            )
            X = pd.concat([majority data resampled, minority data])
            y = pd.concat([y.loc[majority data resampled.index],
y[y == minority class]])
        # Perform oversampling
        if self.strategy in ['oversample', 'both']:
            X, y = self.smote.fit resample(X, y)
        # Shuffle the data to ensure randomness
        X = pd.DataFrame(X).reset index(drop=True)
        y = pd.Series(y).reset index(drop=True)
        return X, y
    def fit transform(self, X, y=None):
        return self.fit(X, y).transform(X, y)
```

3.2.4. Normalisasi Data

Normalisasi data adalah proses untuk mengubah nilai fitur ke dalam skala tertentu (misalnya, antara 0 dan 1) atau ke distribusi tertentu seperti distribusi normal. Hal ini berguna untuk memastikan bahwa model atau proses yang mengasumsikan normalitas data dapat bekerja dengan baik. Selain itu, normalisasi juga membantu mengurangi pengaruh besar kecilnya skala fitur (*magnitude effect*) dalam model. Berikut adalah implementasi dari *DataNormalizer* yang kami buat.

```
class DataNormalizer(BaseEstimator, TransformerMixin):
def __init__(self, feature_range=(0, 1),
numerical_columns=None, scaler=StandardScaler()):
        self.feature range = feature range
        self.numeric columns = numerical columns
        self.method = StandardScaler()
        self.scaler = scaler
    def fit(self, X, y=None):
        X = X.copy()
        # Identify numeric columns if not provided
        if self.numeric columns is None:
             self.numeric columns =
X.select dtypes(include=['int64', 'float64']).columns
        # Fit the scaler to the selected numeric columns
        self.scaler.fit(X[self.numeric columns])
        return self
    def transform(self, X):
        X = X.copy()
        if self.scaler is None:
            raise ValueError("The normalizer must be fitted before
transforming the data.")
        # Apply normalization to numeric columns
        X[self.numeric columns] =
self.scaler.transform(X[self.numeric columns])
        return X
    def fit transform(self, X, y=None):
        return self.fit(X, y).transform(X)
```

3.2.5. Pengurangan Dimensi Data

Mengurangi dimensi data dapat mempercepat proses pembelajaran model serta membuat keterikatan nilai antarfitur semakin kuat. Setelah melalui *pipeline* pemrosesan ini, seluruh fitur yang ada sudah berubah menjadi fitur baru hasil pengolahan fitur aslinya.

Pada implementasi kode, kami menggunakan kelas PCA dari pustaka sklearn. Kelas ini dapat dikostumisasi dengan cara memberi *key argument* "n_components" yang bernilai *integer* pada konstruktornya. Nilai ini menyatakan jumlah fitur yang ingin dihasilkan setelah melewati komponen *pipeline* ini. Kelas dari pustaka ini mampu membuat fitur baru hasil reduksi dimensi secara otomatis.

```
class DimensionalityReduction(BaseEstimator, TransformerMixin):
    def init (self, method='pca', n components=2,
random state=None):
        if method not in ['pca', 'tsne', 'autoencoder']:
            raise ValueError ("Method must be one of 'pca', 'tsne',
or 'autoencoder'.")
        self.method = method
        self.n components = n components
        self.random state = random state
        self.model = None
    def fit(self, X, y=None):
       X = X.copy()
        if self.method == 'pca':
            self.model = PCA(n components=self.n components,
random state=self.random state)
            self.model.fit(X)
       elif self.method == 'tsne':
            # t-SNE does not require fitting; it works directly in
transform.
            self.model = TSNE(n_components=self.n_components,
random state=self.random state)
        elif self.method == 'autoencoder':
            self.model = self. build autoencoder(X.shape[1],
self.n components)
            self.model.fit(X, X, epochs=50, batch size=32,
verbose=0)
       return self
    def transform(self, X):
       X = X.copy()
        if self.model is None:
            raise ValueError("The model must be fitted before
transforming the data.")
        if self.method == 'pca':
            return self.model.transform(X)
        elif self.method == 'tsne':
           return self.model.fit transform(X)
        elif self.method == 'autoencoder':
            encoder = self.model.get layer('encoder')
            return encoder.predict(X)
```

```
def fit transform(self, X, y=None):
        return self.fit(X, y).transform(X)
    def build autoencoder(self, input dim, latent dim):
        # Encoder
        input layer = layers.Input(shape=(input dim,))
        encoded = layers.Dense(64, activation='relu')(input layer)
        encoded = layers.Dense(32, activation='relu') (encoded)
        latent = layers.Dense(latent dim, activation='relu',
name='encoder') (encoded)
        # Decoder
        decoded = layers.Dense(32, activation='relu')(latent)
        decoded = layers.Dense(64, activation='relu')(decoded)
        output layer = layers.Dense(input dim,
activation='linear') (decoded)
        # Autoencoder Model
        autoencoder = Model(inputs=input layer,
outputs=output layer)
        autoencoder.compile(optimizer='adam', loss='mse')
        return autoencoder
```

3.3. *Pipeline* per Algoritma

Setiap algoritma membutuhkan pembersihan data, tetapi prapemrosesan data yang belum tentu sama. Oleh sebab itu, diperlukan *pipeline* yang berbeda-beda untuk setiap algoritma. Berikut merupakan spesifikasinya.

3.3.1. Pipeline k-Nearest Neighbour

Pemrosesan data yang digunakan pada algoritma kNN adalah *imputer* dan *encoder*. *Imputer* digunakan untuk mengganti *value* yang hilang agar data dapat lebih akurat diprediksi oleh kNN. *Encoder* digunakan untuk mengubah data kategorikal menjadi numerikal agar fitur-fitur datanya seragam. Awalnya, akan digunakan *scaling* untuk mengubah skala fitur-fitur yang sudah dibuat menjadi numerikal menjadi 0-1 agar perhitungan *distance* lebih *scalable*. Namun, setelah dilakukan eksperiman, ternyata *scaling* membuat hasil prediksi menjadi lebih buruk.

3.3.2. Pipeline Gaussian Naive-Bayes

Setelah dilakukan eksperimen, diperoleh konfigurasi pipeline yang paling optimal untuk Gaussian Naive Bayes adalah sebagai berikut. *Imputer* adalah *pipeline* yang wajib digunakan karena berfungsi untuk mengatasi nilai yang hilang pada data. Pada kasus ini digunakan kelas SimpleImputer dari pustaka sklean dengan strategi mean untuk tipe data numerik dan juga strategi modus untuk tipe data kategorikal. *Encoder* juga wajib digunakan untuk model ini karena penting bagi model ini untuk data yang diolah bertipe data numerik. Sayangnya, masih banyak fitur yang bertipe data kategorikal dan perlu di-*encode* terlebih dahulu. *Scaler* dan *Normalizer* perlu digunakan juga agar nilai data numerik rentangnya seragam untuk berbagai fitur agar hasil perhitungannya lebih seimbang. *Dim Reduction* juga digunakan agar beban komputasi lebih rendah akibat berkurangnya dimensi data.

```
pipe = Pipeline([
            ('imputer', CustomImputer(
              numeric columns=numeric columns,
              categorical columns=categorical columns filtered
            )),
             ('encoder', FeatureEncoder(
                categorical columns=categorical columns filtered,
                encoding strategy='target'
            )),
             ('scaler', FeatureScaler(
                scaler=MinMaxScaler(),
                numeric columns=numeric columns
            ('normalizer', DataNormalizer(
                numerical columns=numeric columns
            )),
            ('dim reduction', DimensionalityReduction(
                method='pca',
                n components=5,
                random state=42
            )),
            ('classifier', classifier)
        ])
```

3.3.3. Pipeline Iterative Dichotomiser 3

Melalui beberapa hipotesis dan hasil eksperimen, kita menyeleksi beberapa kriteria untuk digunakan dalam pipeline. Yang pertama, *imputer*, ini sudah jelas kami implementasikan karena banyaknya nilai kosong pada data set baik *training* maupun testing. Oleh karena itu, kami implementasikan *imputer* pada *pipeline*. Selain itu, kami juga menggunakan *encoder* untuk meng-*encode* fitur-fitur kategorik menjadi numerik agar data menjadi seragam. Selanjutnya, data tersebut bisa dijalankan algoritma ID3 menggunakan *classifier*.

Bagian 4. Percobaan dan Analisis

4.1. Percobaan

Tabel 2. Perbandingan Hasil Prediksi Algoritma terhadap Hasil Prediksi Pustaka

No.	Percobaan dan Prediksi					
	Data yang digunakan: data n Rasio pemisahan <i>train</i> : <i>vali</i> Spesifikasi algoritma KNN: k = 5	<i>idation</i> = 66 : 3		i data <i>trc</i>	iin	
	distance_metric = euclidear	1				
	Model yang dibuat sendiri					
	Prediksi Algoritma KNN	Prediction Accuracy:		ty Eval	luation	:
			precision	recall	f1-score	support
1.		Analysis	0.07	0.05	0.06	442
1.		Backdoor	0.02	0.01	0.01	386
		DoS	0.29	0.33	0.31	2712
		Exploits	0.39	0.51	0.44	7383
		Fuzzers	0.34	0.29	0.31	4020
		Generic	0.95	0.97	0.96	8844
		Normal	0.79	0.77	0.78	12382
		Reconnaissance	0.65	0.34	0.45	2320
		Shellcode	0.17	0.03	0.05	250
		Worms	0.00	0.00	0.00	29
		accuracy			0.64	38768
		macro avg	0.37	0.33	0.34	38768
		weighted avg	0.64	0.64	0.63	38768

Со	nfusi	ion Ma	atrix	:							
]]	21	9	137	196	20	11	39	9	0		0]
[10	3	127	196	18	6	24	2	. 0		0]
[67	54	882	1263	140	76	171	55	4		0]
[113	84	1231	3748	832	135	1076	161	2		1]
ſ	41	19	274	1549	1154	61	830	75	17		0]
1	7	6	54	112		8575	49	12			1]
1	28	9		1622	893		9526	93			0]
1	28	13	192		200	59	285	785			0]
_t	20	2	132	69	75	21	50	11			0]
, .											
L	0	0	1	7	12	0	9	0	0		0]]
CI	lassif	icatio		ort: ecisio	n ı	recall	f1-s	core	supp	ort	
	A	\nalys:	is	0.0	10	0.00		0.00		660	
		Backdoo	or	0.0	3	0.00 0.72		0.00 0.06		660 576	
	В	Backdoo	or oS)3)3	0.00		0.00	4	660	
	B	Backdoo Do Exploi Fuzze	or oS ts rs	0.0 0.0 0.4 0.2)3)3 !3 !9	0.00 0.72 0.01 0.05 0.66		0.00 0.06 0.01 0.08 0.40	4 11 6	660 576 047 .020	
	B	Backdoo Do Exploit Fuzzen Gener:	or oS ts rs ic	0.0 0.0 0.4 0.2 0.9)3)3 !3 !9	0.00 0.72 0.01 0.05 0.66 0.95		0.00 0.06 0.01 0.08 0.40 0.95	4 11 6 13	660 576 047 020 001	
Re	B E	Backdoo Do Exploi Fuzze	or oS ts rs ic al	0.0 0.0 0.4 0.2	13 13 13 19 15	0.00 0.72 0.01 0.05 0.66		0.00 0.06 0.01 0.08 0.40	4 11 6 13 18	660 576 047 .020	
Re	E Econna	Backdoo Do Exploit Fuzzen Gener: Norma aissand	or oS ts rs ic al ce de	0.6 0.4 0.2 0.9 0.7 0.6	93 13 19 15 19 14	0.00 0.72 0.01 0.05 0.66 0.95 0.70 0.00		0.00 0.06 0.01 0.08 0.40 0.95 0.74 0.00	4 11 6 13 18	660 576 947 920 9001 3200 3480 3462	
Re	E Econna	Backdoo Do Exploit Fuzzen Gener: Norma	or oS ts rs ic al ce de	0.0 0.4 0.2 0.9 0.7	93 13 19 15 19 14	0.00 0.72 0.01 0.05 0.66 0.95 0.70		0.00 0.06 0.01 0.08 0.40 0.95 0.74 0.00	4 11 6 13 18	660 576 047 020 0001 200 480	
Re	B E econna Sh	Backdoo Do Exploit Fuzzen Gener: Norma issand ellcoo Worn	or oS ts rs ic al ce de ms	0.6 0.4 0.2 0.9 0.7 0.6 0.6	93 13 13 19 15 19 10 10	0.00 0.72 0.01 0.05 0.66 0.95 0.70 0.00		0.00 0.06 0.01 0.08 0.40 0.95 0.74 0.00 0.00	4 11 6 13 18 3	660 576 947 920 9001 3200 4480 4462 374 43	
Re	B Econna Sh a ma	Backdoo Do Exploid Fuzzen Gener: Norma Dissand Worn	or oS ts rs ic al ce de ms	0.6 0.4 0.2 0.9 0.7 0.6	13 13 13 19 15 19 10 10	0.00 0.72 0.01 0.05 0.66 0.95 0.70 0.00		0.00 0.06 0.01 0.08 0.40 0.95 0.74 0.00 0.00	4 11 6 13 18 3	660 576 947 920 9001 3200 3480 3462 374 43	
	B E econna Sh a ma weigh	Backdoo Do Exploit Fuzzen Gener: Norma issand ellcoo Worn accurad	or oS ts rs ic al ce de ms cy vg	0.6 0.4 0.2 0.9 0.7 0.6 0.6	13 13 13 19 15 19 10 10	0.00 0.72 0.01 0.05 0.66 0.95 0.70 0.00 0.00		0.00 0.06 0.01 0.08 0.40 0.95 0.74 0.00 0.00	4 11 6 13 18 3	660 576 047 020 6001 2200 4480 4462 374 43	
	econna Sh ma weigh nfusic 0	Backdoo Do Exploit Fuzzer Gener: Norma Dissandellcoo Worr Accurac Accurac Dissandellcoo Worr Accurac A	or oS ts rs ic al ce de ms cy vg vg	0.6 0.4 0.2 0.9 0.7 0.6 0.6 0.2	13 13 13 19 15 19 10 10 10 10 10 10 10 10 10 10 10 10 10	0.00 0.72 0.01 0.05 0.66 0.95 0.70 0.00 0.00	6 6	0.00 0.06 0.01 0.08 0.40 0.95 0.74 0.00 0.00 0.52 0.51	4 11 6 13 18 3 57 57	660 576 047 020 0001 1200 1480 4462 374 43 863 863 863	0]
Co	econna Sh a ma weigh weigh	Backdoo Do Exploit Fuzzer Gener Norma Aissandellcoo Worn Accurac Accur	or oS ts rs ic al ce de ms cy vg vg	0.6 0.4 0.2 0.9 0.7 0.6 0.6 0.2 0.5	13 13 13 19 15 19 10 10 10 10 10 10 10 10 10 10 10 10 10	0.00 0.72 0.01 0.05 0.66 0.95 0.70 0.00 0.00	6 6	0.00 0.06 0.01 0.08 0.40 0.95 0.74 0.00 0.00 0.52 0.51	4 11 6 13 18 3 57 57 57	660 576 .047 .020 .0001 .200 .4480 .4462 .374 .43 .863 .863 .863	0]
Co	econna Sh ma weigh nfusic 0	Backdoo Do Exploit Fuzzer Gener: Norma sissandellcoo Worn accurade accurade accurade accurade accurade accurade accurade 448 448 445	or oS ts rs ic al ce de ms cy vg rix: (0.6 0.4 0.2 0.9 0.7 0.6 0.6 0.6 0.5	23 23 29 25 25 26 26 26 27 27 27 27 27 27 27 27 27 27 27 27 27	0.00 0.72 0.01 0.05 0.66 0.95 0.70 0.00 0.00	6 6	0.00 0.06 0.01 0.08 0.40 0.95 0.74 0.00 0.00 0.52 0.22 0.51	4 11 6 13 18 3 57 57	660 576 047 020 0001 1200 1480 4462 374 43 863 863 863	
Co	econna Sh ma weigh nfusio 0 0 0 0	Sackdoo Do Exploid Fuzzer Gener: Norma Sissand Word Accurac accura accura accura accura accura accura accura accura accura acc	or oS ts rs ic cal cce dde ms cy vg vg rix: () () () () () () () () () () () () ()	0.6 0.4 0.2 0.9 0.7 0.6 0.6 0.2 0.5 2: 2 2: 1887 504 21:	23 23 29 25 26 26 26 26 27 27 28 28 29 29 20 21 21 21 22 23 24 24 24 24 24 24 24 24 24 24 24 24 24	0.00 0.72 0.01 0.05 0.66 0.90 0.00 0.00 0.31 0.52	6 6 8 2 9 27 9 27 9 8	0.00 0.06 0.01 0.08 0.40 0.95 0.74 0.00 0.00 0.52 0.52 0.51	4 11 6 13 18 3 57 57 57 57	660 576 .047 .020 .001 .2200 .4480 .4462 .374 .43 .863 .863 .863	0] 0] 0] 0]
Co	econna Sh ma weigh nfusic 0 0 0 0	Sackdoo Do Exploid Fuzzer Gener: Norma Aissand Hellcoo Worr Accurac Ac	or oS ts rs ic al ce de ms vg vg vg	0.6 0.4 0.2 0.9 0.7 0.6 0.6 0.2 0.5 2: 2 186 7 56 4 21:	23 23 29 25 26 26 26 27 27 28 29 29 20 20 21 21 21 21 21 21 21 21 21 21 21 21 21	0.00 0.72 0.01 0.05 0.66 0.90 0.00 0.00 0.31 0.52	6 6 6 8 27 9 27 9 8 8 22 7 9	0.00 0.06 0.01 0.08 0.40 0.95 0.74 0.00 0.00 0.02 0.52 0.51	4 11 6 13 18 3 57 57 57 57 57	660 576 .047 .020 .0001 .4480 .4462 .374 .43 .863 .863 .863	0] 0] 0] 0] 0]
Co	econna Sh ma weigh nfusio 0 0 0 0	Backdoo Do Exploit Fuzzer Gener: Norma Dissandellcoo Worn Accurac Dorn Mat 448 415 2920 3771 1359 455 1342	or oS ts rs ic al ce de ms vg vg vg	0.6 0.4 0.2 0.9 0.7 0.6 0.6 0.2 0.5 0.2 2.5 2.1 2.1 3.1 3.1 3.1 3.1 3.1	53 63 63 64 66 66 67 67 68 68 69 69 60 60 60 60 60 60 60 60 60 60	0.00 0.72 0.01 0.05 0.66 0.95 0.00 0.00 0.31 0.52 7 7 7 7 7 7 5 6 6 19 6 19 6 34 8	6 6 6 6 6 9 27 9 27 9 27 9 7 9 7 9 7 9 7 9 7 9 7	0.00 0.06 0.01 0.08 0.40 0.95 0.74 0.00 0.00 0.02 0.52 0.51	4 11 6 13 18 3 57 57 57 57	660 576 .047 .020 .001 .2200 .4480 .4462 .374 .43 .863 .863 .863	0] 0] 0] 0]
Co	econna Sh ma weigh nfusic 0 0 0 0 0	Backdoo Do Exploid Fuzzer Genera Norma dissanda dissanda dissanda dissanda dissanda Worr accurada discurada dissanda Horr accurada dissanda Horr accurada dissanda Horr accurada dissanda dissanda Horr accurada dissanda d	or oS ts rs ic al ce de ms vg vg vg	0.6 0.4 0.2 0.9 0.7 0.6 0.6 0.2 0.5 0.2 2.5 2.1 2.1 3.1 3.1	5 11 5 9 6 5 7 384 1 320 8 320	0.00 0.72 0.01 0.05 0.66 0.95 0.00 0.00 0.00 0.52 7 7 7 7 7 5 5 6 6 19 6 19 6 19 19 19 19 19 19 19 19 19 19 19 19 19	6 6 6 6 6 9 27 9 27 9 27 9 7 9 7 9 7 9 7 9 7 9 7	0.00 0.06 0.01 0.08 0.40 0.95 0.74 0.00 0.00 0.00 0.52 0.22 0.51	4 11 6 13 18 3 57 57 57 57 57	660 576 .047 .020 .0001 .2200 .4462 .374 .43 .863 .863 .863	0] 0] 0] 0] 0]

Prediksi Algoritma ID3

Prediction Validity Evaluation: Accuracy: 0.7230 Classification Report: precision recall f1-score support 0.06 0.03 0.04 660 Backdoor 0.11 0.07 0.09 DoS 4047 0.24 Exploits 0.64 0.57 11020 0.49 0.48 6001 Generic 0.99 0.97 0.98 13200 Normal 0.85 0.85 0.85 18480 0.80 0.67 Shellcode 0.35 0.19 Worms 0.11 0.02 0.04 0.72 0.42 57863 0.41 0.72 57863 0.46 macro avg weighted avg 57863

Con	fusio	n Matr	ix:								
11	23	12	140	420	14	2	22	26	1	0]	
[17	40	132	296	35	3	37	16	0	0]	
[115	111	968	2250	257	19	203	117	6	1]	
[161	153	1347	7076	1053	29	993	190	15	3]	
[14	15	211	1415	2807	17	1375	113	33	1]	
[1	8	34	184	38	12861	65	7	2	0]	
[11	12	89	1176	1341	22	15698	102	26	3]	
[26	14	192	735	91	3	82	2313	6	0]	
[0	0	6	104	101	1	93	21	48	0]	
[0	0	0	31	5	0	6	0	0	1]]	

Model scikit-learn

Prediksi Algoritma KNN

knn_classifier = KNeighborsClassifier(n_neighbors=5)

Running with KNN Classifier: Prediction Validity Evaluation: Accuracy: 0.6346 Classification Report: precision recall f1-score support Analysis 0.10 0.07 0.08 660 Backdoor 0.03 0.01 0.02 576 0.33 DoS 0.28 0.31 4947 Exploits 0.39 0.50 0.44 11020 6001 Fuzzers 0.35 0.31 0.33 Generic 0.96 0.96 0.96 13200 Normal 0.79 0.76 0.78 18480 Reconnaissance 0.68 0.32 0.44 3462 Shellcode 0.30 0.07 0.12 374 Worms 0.00 0.00 0.00 accuracy 0.63 57863 macro avg 0.39 0.33 57863 0.35 weighted avg 0.64 0.63 0.63 57863 Confusion Matrix: 200 275 29 29 0] 1918 228 247 0] 62 1904 5508 1279 148 1700 0] 101 0] 2264 1192 0] 0] 109 218 60 12708 2485 1350 47 225 86 14121 10 0] 1120 343 50 442 0] 60 øj]

Prediksi Algoritma GNB

gaussian_nb_classifier = GaussianNB()

Running with Gaussian Naive Bayes Classifier: Prediction Validity Evaluation: Accuracy: 0.4145 Classification Report: precision recall f1-score support Analysis 0.00 0.00 0.00 660 Backdoor 0.00 0.00 0.00 576 DoS 0.05 0.00 0.00 4047 Exploits 11020 0.61 0.02 0.04 6001 Fuzzers 0.25 0.34 0.29 Generic 0.44 0.95 0.60 13200 Normal 0.82 0.44 0.57 18480 Reconnaissance 0.09 0.28 0.14 3462 Shellcode 0.00 0.00 0.00 374 Worms 0.01 0.02 0.01 accuracy 0.41 57863 macro avg 0.23 0.51 0.21 0.17 57863 weighted avg 0.51 0.41 0.37 57863 Confusion Matrix: 84 455 0 103 0] 0] 7] 20] 0 48 452 54 0 0 0 20 57 383 3131 367 96 0 0 25 249 2186 4464 875 3201 0 57 2026 1805 239 1874 0] 0 107 12567 414 108 1] 0 0 35 2596 3876 3588 131] 8177 0 4 614 1753 129 962 0 0 0 95 170 17 92 øj 0 0 1]] 19

		Duranian with transform	
		Running with ID3 Classifier:	
		Prediction Validity Evaluation: Accuracy: 0.7803	
		Classification Report:	
		precision recall f1-score support	
		Analysis 0.14 0.16 0.15 660 Backdoor 0.15 0.14 0.14 576	
		DoS 0.32 0.35 0.33 4047	
		Exploits 0.67 0.67 0.67 11020 Fuzzers 0.65 0.64 0.64 6001	
		Generic 0.98 0.98 0.98 13200	
		Normal 0.91 0.90 0.90 18480 Reconnaissance 0.78 0.75 0.76 3462	
		Shellcode 0.45 0.45 0.45 374	
		Worms 0.35 0.30 0.33 43	
		accuracy 0.78 57863	
		macro avg 0.54 0.53 0.54 57863 weighted avg 0.78 0.78 0.78 57863	
		Confusion Matrix:	
			0] 0]
		[208 142 14 0 3 1835 175 47 47 164 24	2]
			l6] 0]
		[6 5 68 129 18 12944 17 10 2	1]
			0] 1]
		[0 4 17 58 80 2 31 8 170	4]
		[0 0 2 23 2 0 0 3 0 1	[3]]
	Data yang digunakan: data	test yang di-submit ke Kaggle	
	Spesifikasi algoritma KNN k = 5	1:	
	$\begin{vmatrix} k-3 \\ distance metric = manhati$	tan	
	distance_metric - mannati	un	
	Model yang dibuat sendiri		
	Prediksi Algoritma KNN	knn_prediction_manhattan.csv	
		Complete · 13522144 Nicholas Reymond Sihite · 16m ago	8689
2.			
	Prediksi Algoritma GNB	predicted_label.csv 0.188	279
		Complete · Ahmad Rafi Maliki · 15h ago	,,,,
	Prediksi Algoritma ID3	ID3_Prediction6.csv	9533
		Complete · Albert150 · 3h ago	
	Model scikit-learn		
	Prediksi Algoritma KNN	Tidak dapat diuji karena dilarang pada spesifikasi	
1			

Prediksi Algoritma GNB	Tidak dapat diuji karena dilarang pada spesifikasi
Prediksi Algoritma ID3	Tidak dapat diuji karena dilarang pada spesifikasi

4.2. Analisis

Berdasarkan hasil kedua percobaan, algoritma yang paling baik dalam memprediksi attack_cat data (*validation* dan *test*) adalah algoritma *Iterative Dichotomiser 3* dengan skor prediksi *validation 72.3*% dan *test 39.53*%. Setelah itu, algoritma terbaik kedua adalah *k-Nearest Neighbour* dengan skor prediksi *validation 63.3*% dan *test 28.69*%. Algoritma terbaik ketiga adalah *Gaussian Naive Bayes* dengan skor prediksi *validation* adalah 52.44% dan *test 18.879*%. Hasil yang didapat dengan menggunakan model dari *library* 'scikit' juga tidak berbeda dengan hasil yang didapat dengan menggunakan model yang dibuat sendiri, yaitu urutan keakuratan modelnya adalah ID3 - KNN - GNB.

Selain itu, berdasarkan hasil percobaan terhadap algoritma KNN, tidak ada perbedaan signifikan antara hasil prediksi ketiga metrik pengukuran *distance*. Selain itu, tidak ada juga perbedaan signifikan antara 5 *neighbour* dengan 10 *neighbour*.

Pembagian Tugas

Tabel 3. Kontribusi Anggota

NIM - Nama	Pekerjaan
13522122 - Maulvi Ziadinda Maulana	Pembersihan dan prapemrosesan data
13522137 - Ahmad Rafi Maliki	Algoritma Gaussian Naive-Bayes
13522144 - Nicholas Reymond Sihite	Algoritme k-Nearest Neighbour
13522150 - Albert Ghazaly	Algoritma Iterative Dichotomiser 3

Referensi

Elmenshawii, F. 2022. *KNN From Scratch* di https://www.kaggle.com/code/fareselmenshawii/knn-from-scratch (diakses 15 Desember 2024).

Wikipedia kontributor. 2024. *Klasifikasi Naive Bayes* di https://id.wikipedia.org/wiki/Klasifikasi Naive Bayes (diakses 15 Desember 2024).