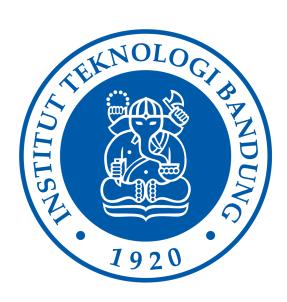
Laporan Tugas Besar 2 IF3170 Inteligensi Artifisial Implementasi Algoritma Pembelajaran Mesin Semester I Tahun 2024/2025



Disusun oleh:

1.	Maulana Muhamad Susetyo	13522127
2.	Andi Marihot Sitorus	13522138
3.	Muhammad Dzaki Arta	13522149
4.	Muhammad Rasheed Qais Tandiung	13522158

INSTITUT TEKNOLOGI BANDUNG 2024

Implementasi Algoritma

1. K-Nearest Neighbor (KNN)

Algoritma K-Nearest Neighbor (KNN) adalah salah satu metode klasifikasi yang paling sederhana dan intuitif dalam machine learning. KNN bekerja dengan cara mengklasifikasikan data baru berdasarkan kedekatannya dengan data yang ada di dalam dataset. Kelas KNN dalam implementasi ini dirancang untuk menyediakan berbagai fungsi yang mendukung proses pelatihan model, prediksi, serta penyimpanan dan pemuatan model. Kelas ini memiliki beberapa fungsi utama yang akan dijelaskan secara rinci berikut.

Fungsi fit

Fungsi ini digunakan untuk melatih model dengan data latih. Fungsi ini menerima dua parameter: X_train, yang merupakan fitur atau atribut data latih, dan y_train, yang merupakan label kelas dari data latih. Fungsi ini akan menyimpan data latih dalam variabel instansi sehingga dapat digunakan dalam fungsi lain, seperti prediksi dan perhitungan jarak.

```
def fit(self, X_train, y_train):
    self.X_train = X_train
    self.y_train = y_train
```

Fungsi predict

Fungsi ini menerima data uji (X_test) dan memprediksi kelas untuk setiap baris data uji. Untuk setiap baris data uji, fungsi ini akan memanggil fungsi get_neighbours untuk mencari tetangga terdekat dari data uji berdasarkan metrik jarak yang dipilih (seperti Euclidean, Manhattan, atau Minkowski). Setelah mendapatkan tetangga terdekat, fungsi ini akan menghitung mayoritas kelas yang ada di antara tetangga tersebut dan mengembalikan kelas yang paling sering muncul sebagai prediksi.

```
def predict(self, X_test):
    preds = []
    for test_row in X_test:
        nearest_neighbours = self.get_neighbours(test_row)
        unique, counts = np.unique(nearest_neighbours,
return_counts=True)
        majority = unique[np.argmax(counts)]
        preds.append(majority)
    return np.array(preds)
```

Fungsi get_neighbours

Fungsi ini digunakan untuk mencari tetangga terdekat dari suatu data uji. Fungsi ini menghitung jarak antara baris data uji dengan setiap baris data latih menggunakan berbagai metrik yang dapat dipilih: Euclidean, Manhattan, atau Minkowski. Setelah menghitung jarak, fungsi ini mengurutkan jarak tersebut dan memilih k tetangga terdekat

untuk data uji. Fungsi ini kemudian mengembalikan kelas dari tetangga-tetangga tersebut untuk digunakan dalam perhitungan mayoritas kelas pada fungsi predict.

```
def get neighbours(self, test row):
        distances = list()
        for (train row, train class) in zip(self.X train,
self.y train):
            if self.metric == 'euclidean':
                dist =
np.sqrt(np.sum((train row-test row)**2))
            elif self.metric == 'manhattan':
                np.sum(np.abs(train row-test row))
            elif self.metric == 'minkowski':
np.sum(np.abs(train row-test row)**self.p)**(1/self.p)
            else:
                raise NameError('Parameter metrics should be
euclidean, manhattan and minkowski')
            distances.append((dist, train class))
        distances.sort(key=lambda x: x[0])
        neighbours = list()
        for i in range(self.k):
            neighbours.append(distances[i][1])
        return neighbours
```

Fungsi save_to_file

Fungsi ini digunakan untuk menyimpan model KNN yang telah dilatih ke dalam sebuah file. Fungsi ini menyimpan informasi tentang nilai k (jumlah tetangga yang digunakan), metric (jenis metrik jarak), serta data latih (X_train dan y_train). Data tersebut disimpan dalam format JSON agar mudah diakses kembali di lain waktu.

```
def save_to_file(self, filename):
    model_data = {
        'k': self.k,
        'metric': self.metric,
        'X_train': self.X_train.tolist(),
        'y_train': self.y_train.tolist()
    }
    with open(filename, 'w') as f:
        json.dump(model_data, f)
    print(f"Model saved to {filename}")
```

Fungsi load_from_file

Fungsi ini digunakan untuk memuat model KNN yang telah disimpan sebelumnya dari sebuah file. Fungsi ini akan membaca file JSON yang berisi data model, kemudian

mengembalikan nilai k, metric, dan data latih yang diperlukan untuk melanjutkan prediksi dengan model yang sama.

```
def load_from_file(self, filename):
    with open(filename, 'r') as f:
        model_data = json.load(f)
    self.k = model_data['k']
    self.metric = model_data['metric']
    self.X_train = np.array(model_data['X_train'])
    self.y_train = np.array(model_data['y_train'])
    print(f"Model loaded from {filename}")
    return self
```

2. Naive Bayes

Algoritma *Naive Bayes* merupakan salah satu metode untuk perancangan model pada *machine learning*, yang menggunakan pendekatan probabilistik untuk melakukan prediksi kelas, yaitu dengan menentukan kelas sebuah input berdasarkan probabilitas yang terobservasi dari training data.

Implementasi *Naive Bayes* pada tugas besar ini merupakan tipe *Gaussian Naive Bayes*, yaitu versi Naive Bayes yang fitur numeriknya diasumsikan mengikuti distribusi Gaussian.

Salah satu aspek utama pada algoritma *Naive Bayes* adalah mengasumsikan independensi antar fitur (naive). Hal ini dapat menjadi aspek yang memperkuat atau memperlemah model dalam prediksi, tergantung data yang dimiliki.

Fungsi fit

```
def fit(self, X_train, y_train):
    self.classes = np.unique(y_train)

for c in self.classes:
    X_c = X_train[y_train == c]

    self.class_mean[c] = X_c.mean(axis=0)
    self.class_var[c] = X_c.var(axis=0) + 1e-9
    self.class_prior[c] = X_c.shape[0] /
X_train.shape[0]
```

Fungsi ini akan menerima X_train dan y_train, lalu untuk setiap kelas akan menetapkan mean, variansi, dan probabilitas dari kelas tersebut.

Fungsi predict

```
def predict(self, X_test):
    preds = []
    for test_row in X_test:
```

```
posteriors = []
    for c in self.classes:
        prior_log = np.log(self.class_prior[c])
        likelihood_log =

np.sum(np.log(self._gaussian_likelihood(test_row, c)))
        posterior = prior_log + likelihood_log
        posteriors.append(posterior)
        preds.append(self.classes[np.argmax(posteriors)])

return np.array(preds)
```

Fungsi prediksi pada implementasi Gaussian Naive Bayes akan menghitung untuk setiap kelas c nilai posterior, yaitu penjumlahan prior probability dari kelas c dengan likelihood untuk setiap atribut sebuah baris pada setiap fitur, yaitu L(atribut | c).

Fungsi prediksi lalu akan mengembalikan kelas dengan nilai posterior tertinggi.

Fungsi _gaussian_likelihood

```
def _gaussian_likelihood(self, x, class_label):
    mean = self.class_mean[class_label]
    var = self.class_var[class_label]
    numerator = np.exp(-0.5 * ((x - mean) ** 2) / var)
    denominator = np.sqrt(2 * np.pi * var)
    return numerator / denominator
```

Fungsi ini akan menghitung likelihood sebuah nilai x pada *gaussian distribution* yang diberikan oleh mean dan var yang sudah disimpan model untuk kelas class_label.

Fungsi save_to_file

```
def save_to_file(self, filename):
    with open(filename, 'wb') as f:
        pickle.dump(self, f)

@staticmethod
def load_from_file(filename):
    with open(filename, 'rb') as f:
        return pickle.load(f)
```

3. ID3

Kelas ID3 memiliki sebuah parameter yaitu max_depth untuk kedalaman maksimalnya. Dalam kelas ID3 terdapat fungsi fit dan predict yang akan digunakan pada pipeline. Terdapat juga fungsi-fungsi yang digunakan untuk membuat decision tree menggunakan algoritma ID3.

Fungsi entropy

Menghitung entropy variabel target

```
def entropy(self, y):
    _, y = np.unique(y, return_inverse=True)
    counts = np.bincount(y)
    probabilities = counts / len(y)
    return -np.sum(probabilities * np.log2(probabilities +
1e-9))

self.y_train = y_train
```

Fungsi information_gain

Menghitung information gain dari sebuah kolom

```
def information gain (self, X, y, feature index,
threshold=None):
        """Calculate information gain for a feature."""
        total entropy = self.entropy(y)
        # Split on threshold if provided
        if threshold is None:
            values, counts = np.unique(X[:, feature index],
return counts=True)
            weighted entropy = 0
            for i, value in enumerate(values):
                subset y = y[X[:, feature index] == value]
                weighted entropy += (counts[i] / len(y)) *
self.entropy(subset y)
        else:
            # Binary split for numerical feature
            left indices = X[:, feature index] <= threshold</pre>
            right indices = X[:, feature index] > threshold
            left entropy = self.entropy(y[left indices]) if
np.any(left indices) else 0
            right entropy = self.entropy(y[right indices]) if
np.any(right indices) else 0
            left weight = np.sum(left indices) / len(y)
            right weight = np.sum(right indices) / len(y)
            weighted entropy = left weight * left entropy +
right weight * right entropy
        return total entropy - weighted entropy
```

Fungsi best_split

Mencari kolom dengan information gain terbaik dan menambahkan threshold jika kolom merupakan numerik

```
def best split(self, X, y):
        best feature = None
        best gain = -np.inf
        best threshold = None
        for feature index in range(X.shape[1]):
            if np.issubdtype(X[:, feature index].dtype,
np.number): # Numerical feature
                thresholds = np.unique(X[:, feature index])
                for threshold in thresholds:
                    gain = self.information gain(X, y,
feature index, threshold)
                    if gain > best gain:
                        best gain = gain
                        best feature = feature index
                        best threshold = threshold
            else: # Categorical feature
                gain = self.information gain(X, y,
feature index)
                if gain > best gain:
                    best gain = gain
                    best feature = feature index
                    best threshold = None
        return best feature, best threshold
```

Fungsi build_tree

Membuat decision tree menggunakan fungsi-fungsi yang telah dibuat

```
def build_tree(self, X, y, feature_names, depth=0):
    # Stopping condition
    if len(np.unique(y)) == 1 or (self.max_depth is not
None and depth >= self.max_depth):
        return np.bincount(y).argmax()

if X.shape[1] == 0:
        return np.bincount(y).argmax()

# Find the best feature to split on
        best_feature_index, best_threshold =
self.best_split(X, y)
        best_feature_name = feature_names[best_feature_index]

if best_threshold is not None:
        best_feature_name = f"{best_feature_name} <=
{best_threshold}"

tree = {best_feature_name: {}}</pre>
```

```
# Handle splitting for numerical or categorical
features
        if best threshold is None:
            feature values = np.unique(X[:,
best feature index])
            for value in feature values:
                indices = X[:, best feature index] == value
                subtree = self.build tree(
                    np.delete(X[indices], best feature index,
axis=1),
                    y[indices],
                    [name for i, name in
enumerate(feature names) if i != best feature index],
                    depth + 1
                tree[best feature name][value] = subtree
            # Numerical split: create <= threshold and >
threshold branches
            left indices = X[:, best feature index] <=</pre>
best threshold
            right indices = X[:, best feature index] >
best threshold
            subtree left = self.build tree(
                np.delete(X[left indices], best feature index,
axis=1),
                y[left_indices],
                [name for i, name in enumerate(feature names)
if i != best feature index],
                depth + 1
            subtree right = self.build tree(
                np.delete(X[right indices],
best feature index, axis=1),
                y[right_indices],
                [name for i, name in enumerate(feature names)
if i != best feature index],
                depth + 1
            tree[best feature name]["Yes"] = subtree left
            tree[best feature name]["No"] = subtree right
        return tree
```

```
def fit(self, X, y):
    if isinstance(X, pd.DataFrame):
        feature_names = X.columns.tolist()
    else:
        feature_names = [str(i) for i in
range(X.shape[1])]

X = np.array(X)
y = np.array(y)
self.tree = self.build_tree(X, y, feature_names)
return self
```

Fungsi predict

Menggunakan decision tree untuk menentukan nilai dari variabel target

```
def predict instance(self, instance, tree):
        if not isinstance(tree, dict):
            return tree
        feature = next(iter(tree))
       feature name, threshold = feature.split(" <= ") if " <=</pre>
" in feature else (feature, None)
        feature value = instance[feature name] if
isinstance(instance, pd.Series) else
instance[int(feature name)]
        if threshold is not None:
            threshold = float(threshold)
            branch = "Yes" if feature value <= threshold else</pre>
"No"
        else:
            branch = feature value
        if branch not in tree[feature]:
            return None # Handle unseen feature values
        return self.predict instance(instance,
tree[feature][branch])
    def predict(self, X):
        if isinstance(X, pd.DataFrame):
            predictions = [self.predict instance(row,
self.tree) for _, row in X.iterrows()]
        else:
            predictions = [self.predict instance(row,
self.tree) for row in X]
        return np.array(predictions)
```

Fungsi save_to_file

```
def save_to_file(self, filename):
    if self.tree is None:
        raise ValueError("Model is not trained yet.")
    with open(filename, 'w') as f:
        json.dump(self.tree, f)
    print(f"Model saved to {filename}")
```

Fungsi load_from_file

```
def load_from_file(self, filename):
    with open(filename, 'r') as f:
        self.tree = json.load(f)
    print(f"Model loaded from {filename}")
    return self
```

Data Cleaning dan Preprocessing

1. Cleaning

Berikut adalah fungsi clean_data yang melakukan pembersihan data sebelum dilakukan preprocessing

```
def clean_data(df_real):
    df = df_real.copy()

    df = gk_lama_imputer(df)
    df = remove_outliers(df)
    df = remove_duplicates(df)
    df = remove_columns(df)

return df
```

fungsi gk_lama_imputer mengganti nilai null dengan nilai:

- median pada data numerikal
- modus pada data kategorikal

```
def gk_lama_imputer(df):
   numeric_imputer = SimpleImputer(strategy='mean')
   categorical_imputer =
SimpleImputer(strategy='most_frequent')

   df[Numerical_Columns] =
   numeric_imputer.fit_transform(df[Numerical_Columns])
   df[Categorical_Columns] =
   categorical_imputer.fit_transform(df[Categorical_Columns])
   return df
```

fungsi remove_outliers membuang nilai yang tidak sesuai dari nilai biner dan 99th percentile pada data numerik yang memiliki nilai diluar batas normal (ditentukan oleh penulis)

```
def remove_outliers_num(df_real, num_cols, lower_percentile=0,
    upper_percentile=0.99):

    df = df_real.copy()

    mask = pd.Series(True, index=df.index)

    for col in num_cols:
        lower_limit = df[col].quantile(lower_percentile)
        upper_limit = df[col].quantile(upper_percentile)
        mask &= ((df[col].isna()) | ((df[col] >= lower_limit) &
```

```
(df[col] <= upper limit)))</pre>
  return df[mask]
def remove outliers cat(df real, cat cols):
  df = df real.copy()
  for col in cat cols:
    most common = df[col].value counts().nlargest(2).index
    df_cleaned = df[df[col].isin(most_common)]
  return df
def remove outliers(df real, rightmost percent=0.01):
  df = df real.copy()
  rightmost columns = [
      "trans depth", "response body len", "sbytes", "dbytes",
"sloss", "dloss", "sjit", "djit", "sinpkt", "dinpkt"
 most common columns = ["is ftp login", "swin", "dwin"]
  df cleaned col = remove outliers cat(df,
most common columns)
  df cleaned = remove outliers num(df cleaned col,
rightmost columns)
  return df cleaned
```

Fungsi remove_duplicates menghapus row yang sama persis

```
def remove_duplicates(df):
    result = df.drop_duplicates()
    return df
```

Fungsi remove_column menghapus kolom yang redundan/tidak diperlukan. Pada kasus ini, setelah melakukan eksperimen, ditemukan nilai akurasi lebih tinggi jika kolom dengan korelasi tinggi TIDAK dihapus.

```
def remove_columns(df_real):
    df = df_real.copy()

# abs_corr = corr.abs()
# upper_triangle = np.triu(np.ones(abs_corr.shape), k=1)
# high_corr_pairs = [
# (abs_corr.index[i], abs_corr.columns[j])
```

```
for i, j in zip(*np.where((abs corr > 0.9) &
(upper triangle == 1)))
  # 1
  # qroups = []
  # seen = set()
  # for f1, f2 in high corr pairs:
        if f1 not in seen and f2 not in seen:
            groups.append({f1, f2})
            seen.update({f1, f2})
  #
        else:
            for group in groups:
                if f1 in group or f2 in group:
                    group.update({f1, f2})
                    seen.update({f1, f2})
                    break
  # groups = [list(group) for group in groups]
  # target correlations = abs corr['attack cat']
  # columns to drop = []
  # for group in groups:
     lowest corr column = target correlations[group].idxmin()
      columns to drop.append(lowest corr column)
  # df = df.drop(columns=columns to drop)
  # columns_to_keep = []
  # for group in groups:
  # highest corr column =
target_correlations[group].idxmax()
     columns to keep.append(highest corr column)
  # ungrouped columns = [
       col for col in abs corr.columns if col not in seen and
col != 'attack cat'
  # columns to keep.extend(ungrouped columns)
 # columns to keep.append('attack cat')
 # df = df[columns to keep]
  return df.drop('label', axis = 1)
```

2. Preprocessing

Berikut adalah pipeline untuk mendapatkan prediksi dari model

```
def pipeliner(model):
  train set clean = clean data(train set)
  val set clean = clean data(val set)
 Categorical Columns for Pipelining =
train set clean.select dtypes(include=['object',
'bool']).columns.tolist()
 Numerical Columns for Pipelining =
train set clean.select dtypes(include=['int64','float64']).columns.tolis
t()
 Categorical Columns for Pipelining = [col for col in
Categorical Columns for Pipelining if col not in ["attack cat", "id",
"label"]]
 Numerical Columns for Pipelining = [col for col in
Numerical Columns for Pipelining if col not in ["attack cat", "id",
"label"]]
  numerical transformer = Pipeline(steps=[
    ('scaler', PowerTransformer(method='yeo-johnson')),
    ('pca',
PCA(n components=Numerical Columns for Pipelining. len ()))
  categorical transformer = Pipeline(steps=[
    ('onehot', OneHotEncoder(handle unknown='ignore',
sparse output=False))
 1)
  preprocessor = ColumnTransformer(
    transformers=[
      ('num', numerical transformer, Numerical Columns for Pipelining),
      ('cat', categorical transformer,
Categorical Columns for Pipelining)
    remainder='passthrough'
  pipeline = Pipeline(steps=[
        ('preprocessor', preprocessor),
        ('model', model)
    1)
  X train = train set clean.drop('attack cat', axis=1)
 y train = train set clean['attack cat']
 X val = val set clean.drop('attack cat', axis=1)
  y val = val set clean['attack cat']
  label encoder = LabelEncoder()
```

```
y_train = label_encoder.fit_transform(y_train)
y_val = label_encoder.transform(y_val)

pipeline.fit(X_train, y_train)
y_pred = pipeline.predict(X_val)

print("Validation Accuracy:", accuracy_score(y_val, y_pred))
```

Preprocessing yang dilakukan adalah sebagai berikut:

- Untuk kolom numerikal, pertama dilakukan scaler menggunakan metode 'yeo-johnson'. Seharusnya bisa menggunakan 'box-cox' tetapi kami mendapatkan error. Kemudian dilakukan PCA dengan jumlah kolom yang sama dengan jumlah kolom sekarang karena menghasilkan akurasi paling tinggi, lebih tinggi daripada melakukan prediksi tanpa melakukan PCA
- Untuk kolom kategorikal, dilakukan OneHotEncoding untuk mengubah format data menjadi machine-learnable.
- Terakhir untuk variabel target, dilakukan LabelEncoding agar dapat digunakan oleh ID3

Perbandingan Model From Scratch dan Library

1. K-Nearest Neighbor (KNN)

Dikarenakan waktu eksekusi yang lama, kami memutuskan untuk melakukan uji coba melakukan 10.000 data pertama.

```
print(f"Model loaded from {filename}")
                return self
        model = KNN(n_neighbors=5, metric='euclidean')
        pipeliner(model)
     √ 1m 34.8s
     {'dpkts', 'synack', 'spkts', 'ct_dst_src_ltm'}
     (7564, 40)
     {'dpkts', 'synack', 'spkts', 'ct_dst_src_ltm'}
     (1903, 40)
     Validation Accuracy: 0.4025223331581713
D ~
        from sklearn.neighbors import KNeighborsClassifier
        model = KNeighborsClassifier(n_neighbors=71, weights='distance')
        # Uncomment kalo mau dirun, ini lama bgt
        pipeliner(model)
     ✓ 0.5s
     {'dpkts', 'synack', 'spkts', 'ct_dst_src_ltm'}
     (7564, 40)
     {'dpkts', 'synack', 'spkts', 'ct_dst_src_ltm'}
     (1903, 40)
     Validation Accuracy: 0.4030478192327903
```

2. Naive Bayes

```
[58] model = NaiveBayes()
pipeliner(model)

{'synack', 'spkts', 'ct_dst_src_ltm', 'dloss'}
(132685, 40)
{'synack', 'spkts', 'ct_dst_src_ltm', 'dloss'}
(33178, 40)
<ipython-input-57-f8dc687ced5b>:24: RuntimeWarning: divide by zero encountered in log likelihood_log = np.sum(np.log(self._gaussian_likelihood(test_row, c)))
Validation Accuracy: 0.10943999035505456

[45] from sklearn.naive_bayes import GaussianNB
model = GaussianNB()
pipeliner(model)

**Yalidation Accuracy: 0.6133883898969197
```

Hasil yang didapatkan model yang sudah dibuat memiliki validation accuracy 0.109, berada di bawah nilai validation accuracy menggunakan model scikit-learn. Hal ini mungkin dapat diatribusikan ke algoritmanya yang banyak melakukan perhitungan matematis. Perhitungan ini memerlukan perhitungan yang presisi yang mungkin terhilangkan pada kode model yang dibuat. Perhitungan yang presisi mungkin dapat dilakukan dengan menggunakan library dan metode yang lebih kompleks.

Selain itu, Naive Bayes merupakan model yang sangat sensitif terhadap dependensi data. Model ini mengasumsikan bahwa semua fitur independen terhadap yang lainnya. Jika terdapat beberapa fitur yang saling berkolerasi dan dependen, hal ini dapat membuat kinerja model menjadi lebih buruk.

3. ID3

Berikut adalah hasil perbandingan implementasi ID3 dengan algoritma dari library namun hanya menggunakan 5000 data pertama dikarenakan durasi yang terlalu lama.

```
return sett
D ~
        model = ID3(max_depth=5)
        pipeliner(model)
[38] V 2m 46.1s
   {'ct_dst_src_ltm', 'dpkts', 'spkts', 'synack'}
    (3783, 40)
    {'ct_dst_src_ltm', 'dpkts', 'spkts', 'synack'}
    (945, 40)
    Validation Accuracy: 0.6952380952380952
D ~
       from sklearn.tree import DecisionTreeClassifier
       model = DecisionTreeClassifier(criterion='entropy')
        pipeliner(model, max_depth=5)
    ✓ 0.6s
    {'ct_dst_src_ltm', 'dpkts', 'spkts', 'synack'}
    (3783, 40)
     {'ct_dst_src_ltm', 'dpkts', 'spkts', 'synack'}
    (945, 40)
    Validation Accuracy: 0.6994708994708995
```

Pembagian Tugas

NIM	Tugas
13522127	Data Cleaning and Preprocessing
13522138	ID3 Implementation
13522149	K-Nearest Neighbour
13522158	Naive Bayes

Referensi

https://medium.com/geekculture/step-by-step-decision-tree-id3-algorithm-from-scratch-in-python-no-fancy-library-4822bbfdd88f