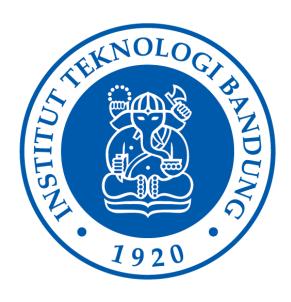
LAPORAN TUGAS BESAR II IF3170 INTELEGENSI BUATAN

IMPLEMENTASI ALGORITMA KNN dan Naive-Bayes



DISUSUN OLEH:

13521044 RACHEL GABRIELA CHEN

13521046 JEFFREY CHOW

13521094 ANGELA LIVIA ARUMSARI

13521134 RINALDY ADIN

PROGRAM STUDI TEKNIK INFORMATIKA SEKOLAH TEKNIK ELEKTRO DAN INFORMATIKA INSTITUT TEKNOLOGI BANDUNG 2023

DAFTAR ISI

DAFTAR ISI	
DAFTAR TABEL	
IMPLEMENTASI KNN	
IMPLEMENTASI NAIVE-BAYES	6
PERBANDINGAN HASIL IMPLEMENTASI DAN PUSTAKA	9
SUBMISI KAGGLE	11
KONTRIBUSI KELOMPOK	13
REPOSITORY	14

DAFTAR TABEL

Tabel 1. Atribut Kelas KNeighborsClassifier	4
Tabel 2. Method Kelas KNeighborsClassifier	
Tabel 3. Atribut Kelas NaiveBayesGaussian	6
Tabel 4. Method Atribut Kelas NaiveBayesGaussian	6
Tabel 5. Atribut Kelas NaiveBayesCategorical	7
Tabel 6. Method Kelas NaiveBayesCategorical	8
Tabel 7. Performance Score KNN	9
Tabel 8. Performance Score Naive Bayes	

IMPLEMENTASI KNN

Algoritma K-Nearest Neighbors (KNN) diimplementasikan dalam kelas KNeighborsClassifier. Daftar atribut dari *class* ini adalah sebagai berikut.

Tabel 1. Atribut Kelas KNeighborsClassifier

Nama Atribut	Keterangan
n_neighbors	Jumlah <i>nearest neighbor</i> yang diambil dalam melakukan prediksi
metric	Pilihan <i>metric</i> yang digunakan untuk menghitung <i>distance</i> antara 2 titik: - 'manhattan': Perhitungan jarak dengan <i>manhattan distance</i> - 'euclidean': Perhitungan jarak dengan <i>euclidean distance</i>
X_train	List yang berisi nilai-nilai X yang digunakan untuk <i>training</i>
y_train	List yang berisi nilai-nilai target(y) yang digunakan untuk <i>training</i>
weight	Pilihan weight yang digunakan untuk melakukan prediksi: - 'distance': weight dari setiap titik adalah invers dari jaraknya ke titik yang akan diprediksi - 'uniform': setiap titik memiliki weight yang sama

Daftar method dari class KNeighborsClassifier adalah sebagai berikut.

Tabel 2. Method Kelas KNeighborsClassifier

Nama Method	Keterangan	
KNeighborsClassifier(n_nei ghbors, metric, weight)	Konstruktor kelas KNeighborsClassifier dengan argumen n_neighbors, metric, dan weight sebagai parameter yang digunakan untuk prediksi	
fit(X_train, y_train)	Fungsi yang digunakan untuk <i>set</i> X_train dan y_train yang digunakan untuk <i>training</i>	
predict(X_test)	Fungsi yang digunakan untuk melakukan prediksi nilai y untuk setiap elemen pada X_test. Fungsi ini akan memanggil get_nearest_neighbors() untuk mendapatkan daftar neighbor terdekat. Kemudian, fungsi ini akan menghitung target dengan memilih target dengan weight terbesar. Jika parameter weight pada kelas adalah 'uniform', maka akan dipilih target dengan count terbanyak pada list nearest neighbors. Jika parameter weight pada kelas adalah 'distance', maka akan dipilih target dengan sum invers jarak terbesar pada list nearest neighbors.	

score(X_test, y_test)	Fungsi ini adalah fungsi untuk menghitung akurasi hasil fungsi predict(X_test) terhadap y_test
get_nearest_neighbors(X_test)	Fungsi ini merupakan helper function untuk mendapatkan nearest neighbors dari setiap X pada X_test. Fungsi ini menggunakan helper function calculate_test_to_train_distance(X_test) untuk mendapatkan distance titik test terhadap setiap titik train. Kemudian hasilnya di-sort secara ascending. Fungsi ini kemudian mengembalikan n_neighbors pertama pada list tersebut.
calculate_test_to_train_ distance(X_test)	Fungsi ini merupakan fungsi untuk menghitung jarak setiap X pada X_test terhadap setiap titik X pada X_train. Jika metric adalah 'manhattan', fungsi ini akan memanggil helper function calculate_manhattan_distance, jika 'euclidean', akan memanggil helper function calculate_euclidean_distance

IMPLEMENTASI NAIVE-BAYES

Algoritma Naive-Bayes adalah sebuah algoritma klasifikasi probabilistik yang didasari oleh *Bayes' Theorem*. Algoritma ini memiliki asumsi bahwa semua *features* yang digunakan untuk klasifikasi independen satu sama lain. Algoritma ini membedakan pemrosesan fitur kategorikal dan numerik sehingga diimplementasikan dalam dua kelas, yaitu NaiveBayesGaussian dan NaiveBayesCategorical.

Kelas NaiveBayesGaussian menangani fitur numerik dengan daftar atribut sebagai berikut.

Tabel 3. Atribut Kelas NaiveBayesGaussian

Nama Atribut	Keterangan
numerical_columns	List yang berisi nilai-nilai seluruh kolom x numerik yang digunakan untuk <i>training</i>
y_train	List yang berisi nilai-nilai target (y) yang digunakan untuk training
target_probs	Dictionary dengan pasangan " key : $value$ " dengan key merupakan nilai unik dari kolom target (y) dan $value$ bernilai $P(key)$
target_values	List ini berisi nilai-nilai unik dari kolom target (y)

Kelas NaiveBayesGaussian dilengkapi dengan method sebagai berikut.

Tabel 4. Method Atribut Kelas NaiveBayesGaussian

Nama Method	Keterangan	
NaiveBayesGaussian()	Konstruktor dari kelas NaiveBayesGaussian yang melakukan inisialisasi pada seluruh atribut kelas.	
fit(numerical_columns, y_train)	Fungsi akan melakukan <i>set</i> pada numerical_columns dan y_train pada atribut kelas sesuai dengan parameter yang diberikan.	
	Fungsi ini juga akan melakukan kalkulasi untuk mengisi atribut target_probs dan conditional_probs. Target_probs akan dikalkulasi dengan menghitung jumlah tiap nilai unik dan membaginya dengan banyak data pada kolom target.	
predict_proba(x_test_numeric al)	Fungsi ini akan menginisialisasi perhitungan probabilitas setiap nilai unik target untuk setiap baris data numerik yang diberikan pada x_test_numeric, yaitu baris-baris data yang	

	ingin diprediksi nilai y nya. Perhitungan probabilitas dilakukan dengan memanggil fungsigaussian_proba terhadap masing-masing baris data.	
gaussian_proba(x_test_nu merical)	Fungsi ini bertujuan untuk menghitung probabilitas dari setiap nilai unik target untuk x_test_numeric yang berupa baris data masukan yang ingin diprediksi. Probabilitas dari setiap nilai unik dikembalikan dalam bentuk <i>list</i> .	
	Perhitungan dilakukan dengan mencari nilai $P(x \mid y)$ untuk setiap kolom untuk setiap nilai target y yang mungkin. Nilai x merupakan nilai dari suatu kolom pada x_test_numeric. Nilai $P(x \mid y)$ dihitung dengan,	
	$P(x y) = \frac{1}{\sqrt{2\pi\sigma^2}} exp(-\frac{(x-\mu)^2}{2\sigma^2})$	
	$\sqrt{2\pi\sigma^2}$ 2σ Untuk menghindari perkalian nilai $P(x \mid y)$ yang terlalu kecil, dicari nilai logaritma dari nilai probabilitas,	
	$log(P(x y)) = -0.5 log(2\pi\sigma^2) - \frac{(x-\mu)^2}{2\sigma^2}$	
	Dengan σ^2 bernilai dari varians dari semua nilai pada kolom yang dicari pada x_train yang memiliki nilai target y , serta μ merupakan nilai rata-rata dari nilai pada kolom yang dicari pada x_train yang memiliki nilai target y .	
	Probabilitas dari nilai target y untuk x_test_numeric dapat dihitung dengan,	
	$P(y x1, x2, x3) = P(y) \prod P(xi y)$	
	Dengan dihitung secara logaritma,	
	$log(P(y x1,x2,)) = log(P(y)) + \sum log(P(xi \mid y))$	
	dengan nilai $P(y)$ ditemukan dari atribut target_probs pada kelas.	
	Nilai probabilitas dari setiap nilai y dalam bentuk logaritma akan dihitung eksponennya untuk dikembalikan.	

Kelas ${\tt NaiveBayesCategorical}$ menangani fitur kategorikal dengan daftar atribut sebagai berikut.

Tabel 5. Atribut Kelas NaiveBayesCategorical

Nama Atribut	Keterangan
categorical_columns	List yang berisi nilai-nilai seluruh kolom x kategorikal yang digunakan untuk <i>training</i>

y_train	List yang berisi nilai-nilai target (y) yang digunakan untuk training
target_probs	Dictionary dengan pasangan "key: value" dengan key merupakan nilai unik dari kolom target (y) dan value bernilai $P(key)$
conditional_probs	List yang berisi dictionary dengan pasangan "a-b: value" dengan a merupakan nilai unik dari kolom pada categorical_columns, b merupakan nilai unik dari kolom target (y), dan value merupakan nilai dari $P(a \mid b)$

Kelas NaiveBayesCategorical dilengkapi dengan method sebagai berikut.

Tabel 6. Method Kelas NaiveBayesCategorical

Nama Method	Keterangan	
NaiveBayesCategorical()	Konstruktor dari kelas NaiveBayesCategorical yang melakukan inisialisasi pada seluruh atribut kelas.	
fit(categorical_columns, y_train)	Fungsi akan melakukan <i>set</i> pada categorical_columns dan y_train pada atribut kelas sesuai dengan parameter yang diberikan.	
	Fungsi ini juga akan melakukan kalkulasi untuk mengisi atribut target_probs dan conditional_probs. Target_probs akan dikalkulasi dengan menghitung jumlah tiap nilai unik dan membaginya dengan banyak data pada kolom target. Conditional_probs berupa "a-b: value" akan dikalkulasi dengan menghitung banyak data pada salah satu kolom categorical_columns yang bernilai a ketika kolom target bernilai b dibagi dengan banyaknya data pada kolom target yang bernilai b. Pada perhitungan conditional_probs akan dilakukan normalisasi.	
predict_proba(x_test)	Fungsi ini akan melakukan perhitungan probabilitas masing-masing baris data terhadap klasifikasi kolom target. Nilai probabilitas masing-masing klasifikasi kolom target dilakukan dengan rumus berikut,	
	$P(y x1, x2, x3) = P(y) \prod P(xi y)$	
	Hal tersebut diimplementasi dengan mendapatkan nilai $P(xi \mid y)$ dari conditional_probs, serta nilai $P(y)$ yang didapatkan dari target_probs. Untuk hasil yang lebih optimal, probabilitas akan dinormalisasi dengan pembagian terhadap jumlah nilai probabilitas pada baris data tersebut.	

PERBANDINGAN HASIL IMPLEMENTASI DAN PUSTAKA

Untuk melakukan perbandingan *performance* antara hasil implementasi dan pustaka, digunakan beberapa metrik evaluasi kinerja klasifikasi. Metrik-metrik ini termasuk *Accuracy Average*, F1 *Macro Average*, F1 *Micro Average*, *Precision Macro Average*, *Precision Micro Average*, *Recall Macro Average*, dan *Recall Micro Average*.

Berikut adalah *performance score* dari algoritma KNN yang dihasilkan oleh hasil implementasi dan pustaka:

Tabel 7. Performance Score KNN

Matrile	Performance Score	
Metrik	Implementasi	Pustaka
Accuracy Average	0.94	0.94
F1 Macro Average	0.94	0.94
F1 Micro Average	0.94	0.94
Precision Macro Average	0.94	0.94
Precision Micro Average	0.94	0.94
Recall Macro Average	0.94	0.94
Recall Micro Average	0.94	0.94

Dari tabel *performance score* terhadap algoritma KNN diatas, *performance score* untuk algoritma KNN hasil implementasi dan pustaka sama yang menunjukkan bahwa hasil implementasi sudah baik dan benar. Keduanya menunjukkan *performance score* di angka 0.94 yang menunjukkan process *tuning* dan *training* model sudah baik.

Berikut adalah *performance score* dari algoritma Naive Bayes yang dihasilkan oleh hasil implementasi dan pustaka:

Tabel 8. Performance Score Naive Bayes

Metrik	Performance Score		
	Implementasi	Pustaka	
Accuracy Average	0.7817	0.7817	
F1 Macro Average	0.7809	0.7809	
F1 Micro Average	0.7817	0.7817	
Precision Macro Average	0.7814	0.7814	

Precision Micro Average	0.7817	0.7817
Recall Macro Average	0.7806	0.7806
Recall Micro Average	0.7817	0.7817

Dari tabel *performance score* algoritma Naive-Bayes di atas, algoritma yang diimplementasikan memiliki *performance score* yang sama. Hal ini membuktikkan bahwa implementasi dari Naive Bayes yang diterapkan sudah baik dan benar. Hal ini juga dapat dilihat dari prediksi kedua implementasi yang menghasilkan hasil yang serupa. Akurasi yang bernilai 0.78 menandakan bahwa pada kasus ini algoritma Naive-Bayes tidak sebaik algoritma KNN dalam melakukan prediksi klasifikasi. Hal ini dapat disebabkan oleh asumsi pada Naive Bayes yang menganggap seluruh *features* independen.

SUBMISI KAGGLE

Submisi Kaggle dilakukan dengan memanfaatkan implementasi algoritma KNN yang memiliki akursi lebih tinggi. Untuk melakukan submisi dilakukan pemrosesan menggunakan parameter yang telah diperoleh melalui eksperimen menggunakan scikit-learn library dan eksperimen lebih lanjut. Berikut merupakan pemrosesan menggunakan parameter yang didapatkan melalui eksperimen menggunakan scikit-learn library.

```
from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import f_classif
from lib.knn import KNeighborsClassifier

df_train = pd.concat([pd.read_csv('data/data_train.csv'), pd.read_csv('data/data_validation.csv')])
df_test = pd.read_csv('data/data_test.csv')
X_train = df_train.drop('price_range', axis=1)
y_train = df_train['price_range'].values
X_test = df_test

X_train_selected, X_test_selected, selected_features = get_x_test_best_features(N_FEATURES_SELECTED,
X_train, y_train, X_test)
X_test_selected = X_test_selected.values
KNN = KNeighborsClassifier(n_neighbors=BEST_N_NEIGHBORS, metric=BEST_METRIC, weight=BEST_WEIGHT)
KNN.fit(X_train_selected,y_train)
predictions = df_test['price_range'] = KNN.predict(X_test_selected)
df_test[['id','price_range']].to_csv("submission.csv", index=False)
```

Berikut merupakan pemrosesan menggunakan parameter yang didapatkan melalui eksperimen lebih lanjut.

```
from lib.knn import KNeighborsClassifier

df_train = pd.concat([pd.read_csv('data/data_train.csv'), pd.read_csv('data/data_validation.csv')])

df_test = pd.read_csv('data/data_test.csv')

X_train = df_train.drop('price_range', axis=1)

y_train = df_train['price_range'].values

X_test = df_test

X_train_selected, X_test_selected, _ = get_x_test_best_features(4, X_train, y_train, X_test)

X_test_selected = X_test_selected.values

WeightedKNN = KNeighborsClassifier(n_neighbors=2, metric='euclidean', weight='distance')

WeightedKNN.fit(X_train_selected, y_train)

weighted_predictions = df_test['price_range'] = WeightedKNN.predict(X_test_selected)

df_test[['id','price_range']].to_csv("submission3.csv", index=False)
```

Pemrosesan untuk melakukan submisi Kaggle secara lebih jelas adalah sebagai berikut.

- 1. Mendefinisikan df_train yang berupa data dari data_train.csv dan data_validation.csv
- Mendefinisikan df_test yang berupa data dari data_test.csv

- 3. Mendefinisikan X_train yang berupa df_train tanpa kolom target yaitu kolom "price_range"
- 4. Mendefinisikan y_train kolom target "price_range" pada df_train
- 5. Mencari X_train_selected, X_test_selected, dan selected_features melalui fungsi get_x_test_best_features untuk mencari fitur paling baik yang akan digunakan untuk melakukan *training*
- 6. Membuat model KNN dengan KNeighborsClassifier dengan parameter jumlah neighbor, metric, dan weight
- 7. Melakukan prediksi menggunakan model KNN dengan X_test_selected yang telah didapatkan
- 8. Melakukan export kolom id dan price range ke csv untuk melakukan submisi

Berdasarkan kedua submisi dan eksperimen yang dilakukan, didapatkan parameter yang paling optimal yaitu fitur terbaik yang dipilih empat, jumlah *neighbor* dua, *metric manhattan*, dan *weight distance*.

KONTRIBUSI KELOMPOK

Fitur	Kontributor
KNN	13521044, 13521046
Naive-Bayes	13521094, 13521134

REPOSITORY

Seluruh implementasi lengkap dapat diakses melalui repository berikut https://github.com/liviaarumsari/Tubes2-whoosh