

BAB II

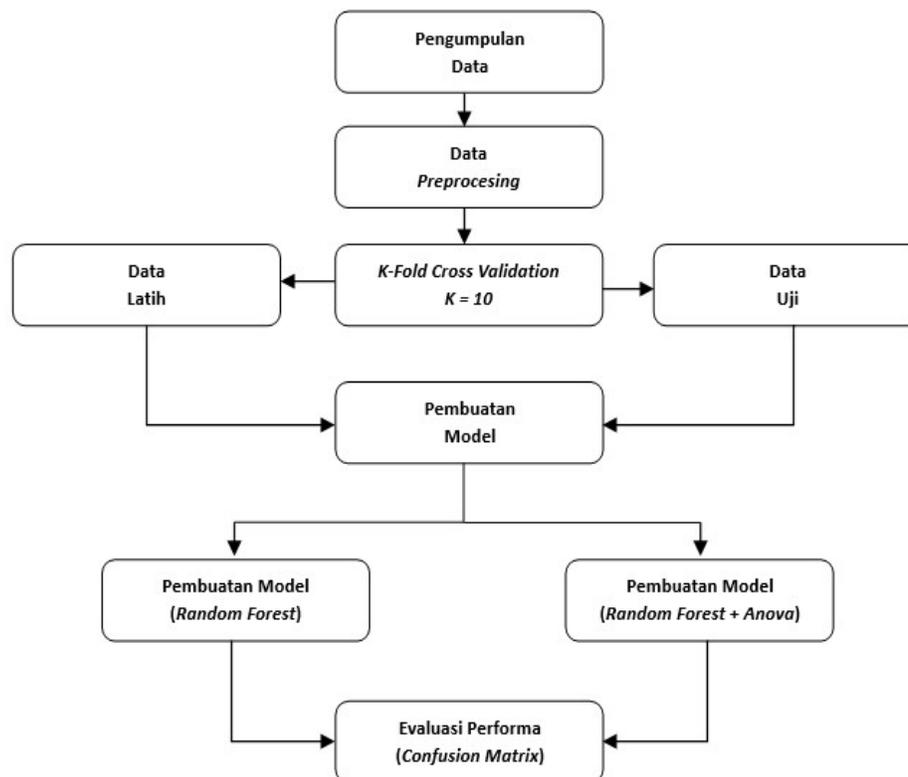
METODOLOGI PENELITIAN

2.1 Objek Penelitian

Stunting merupakan masalah kesehatan yang terjadi pada anak-anak. Di kota Samarinda kasus Stunting terbilang cukup tinggi yaitu sebesar 25,3% (Cindy Mutia Annur, 2023). Oleh karena itu, objek yang digunakan dalam penelitian ini adalah data Stunting pada tahun 2023 yang didapat dari 26 puskesmas di lingkup Dinas kesehatan Kota Samarinda. Penelitian dilakukan di Dinas Kesehatan Kota Samarinda yang berlokasi di Jl. Milono No.1, 75122, Bugis, Kec. Samarinda Kota, Kota Samarinda, Kalimantan Timur 76112.

2.2 Teknik Analisis Data

Penelitian ini bertujuan untuk mengembangkan sebuah metodologi yang terstruktur dan terorganisir dengan baik untuk memastikan bahwa proses penelitian dilaksanakan secara konsisten dan teratur. Teknik analisis data yang akan digunakan telah dirancang untuk mengikuti alur kerja yang sistematis, yang akan diuraikan dalam alur penelitian berikut :



Gambar 2. 1 Alur Penelitian

2.2.1 Pengumpulan data

Tahapan awal pada penelitian ini adalah melakukan perhaman pada data, dimulai dari pengumpulan data lalu dilanjutkan dengan proses data preparation atau persiapan data kasus Stunting. Data yang didapat memiliki 23 kolom dengan 22 atribut dan 1 atribut sebagai label (TB/U). Dapat dilihat pada Tabel 2.1 adalah data yang didapat dari Dinas kesehatan Kota Samarinda.

Tabel 2.1 Atribut Data Dinas Kesehatan Kota Samarinda

No	Atribut	Type Data	Keterangan
1	Nama	<i>String</i>	Nama
2	JK	<i>String</i>	Jenis Kelamin
3	BB Lahir	<i>Integer</i>	Berat Badan Lahir
4	TB Lahir	<i>Integer</i>	Tinggi Badan Lahir
5	Tgl Lahir	<i>Date</i>	Tanggal Lahir
6	Nama Ortu	<i>String</i>	Nama Orang Tua
7	Provinsi	<i>String</i>	Provinsi
8	Kab/Kota	<i>String</i>	Kabupaten atau Kota
9	Kec	<i>String</i>	Kecamatan
10	Puskesmas	<i>String</i>	Lokasi Puskesmas
11	Posyandu	<i>String</i>	Lokasi Posyandu
12	RT	<i>Integer</i>	Lokasi RT
13	Total Pengukuran	<i>Integer</i>	Total Pengukuran
14	Tanggal Pengukuran	<i>Integer</i>	Tanggal Pengukuran
15	Berat	<i>Integer</i>	Berat Badan
16	Tinggi	<i>Integer</i>	Tinggi Badan
17	BB/U	<i>numeric</i>	Berat Badan Menurut Umur
18	ZS BB/U	<i>numeric</i>	Z Score Berat Badan menurut Umur
19	TB/U	<i>numeric</i>	Tinggi Badan menurut Umur
20	ZS TB/U	<i>numeric</i>	Z Score Tinggi Badan menurut Umur
21	BB/TB	<i>numeric</i>	Berat Badan menurut Tinggi Badan
22	ZS BB/TB	<i>numeric</i>	Z Score Berat Badan menurut Tinggi Badan

2.2.2 Data Pre-Processing

Pada Tahapan ini akan memiliki 4 tahapan, yaitu data *selection*, data *cleaning*, data *Transformation* dan data *Balancing*. Data yang digunakan berasal dari data Dinas Kesehatan Kota Samarinda dengan proses persiapan sebagai berikut :

a) Data Selection

Pada tahapan ini, akan dilakukan pengambilan data dengan memilih fitur atau atribut yang relevan untuk digunakan pada penelitian sedangkan fitur yang dianggap tidak relevan akan dihilangkan, seperti pada tabel 2.2 menampilkan data awal yang didapatkan dari Dinas Kesehatan Kota Samarinda memiliki 23 kolom, lalu didapatkan 11 kolom yang dianggap tidak relevan untuk digunakan dalam mengklasifikasikan stunting pada anak, setelah dilakukan proses seleksi maka didapatkan 13 atribut yang dijadikan fitur dan 1 atribut (TB/U) sebagai target atau kelas. Dapat dilihat pada Tabel 2.2 adalah data yang telah diseleksi.

Tabel 2.2 Data Selection

No	Atribut	Keterangan	Tipe Data
1	Nama	Nama	String
2	JK	Jenis Kelamin	String
3	Berat	Berat Badan	Integer
4	Tinggi	Tinggi Badan	Integer
5	LiLA	Lingkar Lengan Atas	Integer
6	BB/U	Berat Badan Menurut Umur	numeric
7	ZS BB/U	Z Score Berat Badan menurut Umur	numeric
8	Tanggal Pengukuran	Tanggal Pengukuran	Integer
9	ZS TB/U	Z Score Tinggi Badan menurut Umur	numeric
10	BB/TB	Berat Badan menurut Tinggi Badan	numeric
11	ZS BB/TB	Z Score Berat Badan menurut Tinggi Badan	numeric
12	Naik Berat B	Naik Berat Badan	String
13	Jml Vit A	Jumlah Vitamin A	integer
14	TB/U	Tinggi Badan menurut Umur	numeric

b) Data Cleaning

Pada Tahapan ini data yang digunakan akan dilakukan proses pembersihan dengan menghapus data yang memiliki nilai #N/A (*no value is available*) atau data yang tidak memiliki nilai serta yang kembar atau terduplikasi. Proses ini akan menggunakan bahasa pemrograman *Python* serta library *Pandas*, dengan fungsi *dropna()* untuk menghapus baris yang memiliki nilai hilang dan fungsi *drop_duplicates()* untuk menghapus baris yang merupakan duplikat.

```

missing_values = data.isna().sum()

# Menampilkan jumlah missing values untuk setiap kolom
print("Jumlah missing values untuk setiap kolom:")
print(missing_values)

```

Gambar 2. 2 Kode Menampilkan Nilai Kosong

Berikut adalah tabel 2.3 yang berisi tentang tiap parameter yang digunakan dalam kode tersebut beserta penjelasan fungsi tiap parameter.

Tabel 2.3 Parameter Kode Menampilkan Nilai Kosong

Parameter	Keterangan
data	Ini adalah variabel yang merujuk pada DataFrame pandas. DataFrame adalah struktur data dua dimensi yang bisa menyimpan data dalam bentuk tabel dengan baris dan kolom.
isna()	Fungsi isna() dari pandas digunakan untuk mendeteksi nilai-nilai yang hilang dalam DataFrame. Fungsi ini mengembalikan DataFrame baru di mana setiap entri adalah True jika entri asli adalah nilai yang hilang (NaN, None, atau NaT), dan False jika tidak.
sum()	Metode sum() digunakan di sini setelah isna() untuk menghitung jumlah True di setiap kolom dari DataFrame yang dihasilkan oleh isna(). Ini efektif menghitung jumlah nilai yang hilang dalam setiap kolom dari DataFrame asli data.

```

▶ # Menghitung jumlah data sebelum penghapusan duplikat
total_data_sebelum = len(data)

# Mengurutkan data berdasarkan kolom "Nama" dan "Tanggal Pengukuran" secara menurun
data = data.sort_values(by=['Nama', 'Tanggal Pengukuran'], ascending=[True, False])

# Menghapus duplikat berdasarkan kolom "Nama" dan mempertahankan yang pertama (yang memiliki tanggal pengukuran terbaru)
data = data.drop_duplicates(subset='Nama', keep='first')

#Menghapus Nilai Kosong Tiap Kolom
data.dropna(inplace=True)

# Menghitung jumlah data setelah penghapusan duplikat
total_data_sesudah = len(data)

print(f"Total data sebelum penghapusan duplikat: {total_data_sebelum}")
print(f"Total data setelah penghapusan duplikat: {total_data_sesudah}")

```

Gambar 2.3 Kode Menghapus Data Terduplikasi Dan Data Kosong

Berikut adalah tabel 2.4 yang berisi tentang tiap parameter yang digunakan dalam kode tersebut beserta penjelasan fungsi tiap parameter.

Tabel 2.4 Parameter Kode Menghapus Data Terduplikasi Dan Data Kosong

Fungsi	Parameter	Keterangan
<i>len()</i>	: data	data adalah <i>DataFrame</i> yang sedang diproses. Fungsi <i>len()</i> di sini digunakan untuk menghitung jumlah baris dalam <i>DataFrame</i> .
<i>sort_values()</i>	: <i>by</i>	Parameter <i>by</i> menerima list dari nama-nama kolom yang akan dijadikan dasar pengurutan. Di sini, pengurutan dilakukan berdasarkan kolom "Nama" dan "Tanggal Pengukuran".
	: <i>ascending</i>	Parameter <i>ascending</i> menerima list dari nilai <i>boolean</i> yang menentukan urutan pengurutan untuk setiap kolom yang ditentukan di <i>by</i> . Nilai <i>True</i> untuk "Nama" berarti pengurutan secara <i>ascending</i> (A-Z), dan <i>False</i> untuk "Tanggal Pengukuran" berarti pengurutan secara <i>descending</i> (tanggal terbaru ke terlama).
<i>drop_duplicates()</i>	: <i>subset</i>	Parameter <i>subset</i> menentukan kolom yang akan digunakan untuk identifikasi duplikat. Di sini, 'Nama' adalah kolom yang digunakan, yang berarti duplikat akan dicari berdasarkan nama.
	: <i>keep</i>	Parameter <i>keep</i> menentukan duplikat mana yang akan dipertahankan. <i>first</i> menunjukkan bahwa hanya baris pertama dari setiap duplikat yang akan dipertahankan, yang di kasus ini adalah entri dengan 'Tanggal Pengukuran' terbaru karena pengurutan sebelumnya.
<i>print()</i>	: <i>f-string</i>	<i>f-string</i> digunakan untuk memformat string, menyisipkan nilai dari variabel langsung ke dalam <i>string</i> yang dicetak. <i>total_data_sebelum</i> dan <i>total_data_sesudah</i> adalah variabel yang nilai-nya disisipkan dan dicetak.

c) Data Transformation

Pada tahapan ini dilakukan dengan proses konversi nilai-nilai atribut yang kategorikal menjadi bentuk numerik, tahapan ini perlu dilakukan karena dalam penggunaan library sklearn, hanya bisa menerima atribut dengan tipe numerik. Beberapa atribut yang ditransformasi meliputi jenis kelamin, naik berat badan, berat badan menurut umur, dan berat badan menurut tinggi badan. Proses ini akan menggunakan library *scikit-learn* dengan fungsi *LabelEncoder*. *Label Encoder* tersebut digunakan untuk mengubah teks atau data kategorikal menjadi data numerik dalam satu kolom data secara otomatis (M. M. Mafa'atih, 2020)

```
from sklearn.preprocessing import LabelEncoder

# Buat instance LabelEncoder
encoder = LabelEncoder()

# Kolom yang perlu di-encode
columns_to_encode = ['JK', 'BB/U', 'BB/TB', 'Naik Berat Badan']

# Loop melalui setiap kolom yang perlu di-encode dan terapkan LabelEncoder
for column in columns_to_encode:
    data[column] = encoder.fit_transform(data[column])
```

Gambar 2.4 Kode Tranformasi Data Pada Kolom Atribut

```
# Ubah TB/U
data['TB/U'] = data['TB/U'].replace({'Normal': 0, 'Tinggi': 0, 'Pendek': 1, 'Sangat Pendek': 1})
```

Gambar 2.5 Kode Tranformasi Data Pada Kolom Label

Berikut adalah tabel 2.5 yang berisi tentang tiap parameter yang digunakan dalam kode tersebut beserta penjelasan fungsi tiap parameter.

Tabel 2.5 Parameter Kode Tranformasi

Fungsi	Parameter	Keterangan
<i>LabelEncoder()</i>	-	<i>LabelEncoder</i> adalah kelas dari <i>sklearn.preprocessing</i> yang digunakan untuk mengkonversi label kategori menjadi label numerik.
<i>fit_transform()</i>	<i>data[column]</i>	Metode <i>fit_transform()</i> digunakan untuk menyesuaikan encoder dan mengubah label menjadi bentuk <i>numerik</i> . Parameter <i>data[column]</i> merupakan kolom spesifik dari <i>DataFrame</i> data yang sedang diolah.
<i>replace()</i>	{'Normal': 0, 'Tinggi': 0, 'Pendek': 1, 'Sangat Pendek': 1}	Metode <i>replace()</i> digunakan untuk mengganti nilai-nilai dalam sebuah <i>Series</i> atau <i>DataFrame</i> . Parameter ini merupakan <i>dictionary</i> yang menentukan nilai lama dan baru; dalam hal ini, 'Normal' dan 'Tinggi' diganti dengan 0, 'Pendek' dan 'Sangat Pendek' diganti dengan 1.

d) Data Balancing

Pada tahapan terakhir preprocessing, ketidakseimbangan kelas akan terjadi ketika jumlah atribut mengalami ketidakseimbangan antara kelas mayoritas dan kelas minoritas. Sehingga Perlu melakukan data balancing menggunakan modul *python imblearn.over_sampling* dengan mengimpor fungsi *SMOTE (sampling strategy)* untuk melakukan *oversampling*.

```
x = data.drop(['TB/U'],axis=1)
y = data['TB/U']
from imblearn.over_sampling import SMOTE
sm = SMOTE(random_state=42, sampling_strategy=1)
X_imp, y_imp = sm.fit_resample(x, y)
```

Gambar 2.6 Kode Data Balancing SMOTE

Berikut adalah tabel 2.6 yang berisi tentang tiap parameter yang digunakan dalam kode tersebut beserta penjelasan fungsi tiap parameter.

Tabel 2.6 Parameter Kode Data Balancing SMOTE

Fungsi	Parameter	Keterangan
<i>drop()</i>	<i>['TB/U']</i>	Parameter ini menentukan nama kolom atau daftar kolom yang akan dihapus dari <i>DataFrame</i> data. Di sini, kolom 'TB/U' akan dihapus.
	<i>axis=1</i>	<i>Axis</i> 1 mereferensikan kolom. Dengan <i>axis=1</i> , operasi <i>drop()</i> akan menghapus kolom yang ditentukan.
<i>SMOTE</i>	<i>random_state=42</i>	Memastikan bahwa hasil <i>oversampling</i> konsisten di berbagai eksekusi dengan menyediakan <i>seed</i> untuk generator bilangan acak. 42 adalah nilai yang sering digunakan dalam contoh untuk tujuan <i>reproducibilitas</i> .
	<i>sampling_strategy=1</i>	Menentukan strategi <i>sampling</i> yang diinginkan. Nilai 1 mengindikasikan bahwa kelas minoritas akan <i>dirsample</i> hingga jumlahnya sama dengan kelas mayoritas.
<i>fit_resample()</i>	<i>x, y</i>	Metode <i>fit_resample()</i> digunakan untuk mengajarkan model <i>SMOTE</i> pada data input (<i>x</i>) dan target (<i>y</i>), kemudian mengembalikan dataset yang telah di-resample. <i>x</i> adalah <i>DataFrame</i> yang dihasilkan setelah menghapus kolom 'TB/U', dan <i>y</i> adalah nilai target yang berasal dari kolom 'TB/U'.

2.2.3 Pembagian data

Sebelum memulai proses pemodelan, dataset yang akan digunakan harus dibagi menjadi dua bagian utama, yaitu data latih dan data uji. Data latih berfungsi sebagai dasar untuk pembangunan model, sementara data uji digunakan untuk mengevaluasi kinerja model yang telah dibuat. Dalam penelitian ini, akan menggunakan teknik *K-Fold Cross Validation*, Hal ini dapat membantu mengidentifikasi

$$l(y) = \underset{c}{\operatorname{argmax}} \left(\sum_{n=1}^N I_{hn}(y) = C \right) \quad (2.1)$$

- a. Langkah Pertama sampel dataset tersebut dimanfaatkan untuk membentuk pohon keputusan atau *tree* ke- i ($i = 1, 2, \dots, k$). Dalam tahap pembuatan *tree* ini, digunakan metode *Classification and Regression Tree* (CART), di mana gini impurity diaplikasikan sebagai kriteria untuk menentukan pembagian pada setiap node dalam *tree* (Yoga Siswa, 2023) .

```
#Kode Inisialisasi Random Forests
from sklearn.ensemble import RandomForestClassifier
clf = RandomForestClassifier(max_depth=1, random_state=42)
```

Gambar 2.7 Kode Inisialisasi Model *Random Forest*

- b. Langkah kedua melakukan pembagian data menggunakan *K-Fold Cross Validation* dengan $K = 10$.

```
# Kode kfold Cross Validation
from sklearn.model_selection import cross_val_score, StratifiedKFold
cv = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)
```

Gambar 2.8 Kode Pembagian Data Latih dan Data Uji

Berikut adalah tabel 2.8 yang berisi tentang tiap parameter yang digunakan dalam kode tersebut beserta penjelasan fungsi tiap parameter:

Tabel 2.8 Parameter Pembagian Data Latih dan Data Uji

Parameter	Keterangan
<i>n_splits</i>	Jumlah lipatan (fold) yang akan dibuat. Dalam kasus ini, nilai <i>n_splits</i> = 10 berarti akan memiliki 10 lipatan untuk cross-validation.
<i>shuffle</i>	Menentukan apakah data akan diacak sebelum dibagi menjadi lipatan. Jika <i>shuffle</i> = True, data akan diacak sebelum dibagi, sehingga setiap lipatan akan memiliki sampel yang berbeda-beda. Jika <i>shuffle</i> = false, data akan dibagi berdasarkan urutan aslinya.
<i>random_state</i>	Menentukan seed untuk generator nomor acak yang digunakan saat mengacak data. Penggunaan nilai <i>random_state</i> yang sama akan menghasilkan pembagian yang sama saat melakukan pengacakan.

- c. Langkah ketiga dalam pemodelan algoritma Random Forest melibatkan penggunaan *library scikit-learn* untuk mengimplementasikan model. Pada tahap ini, akan menggunakan metode *k-fold cross-validation* untuk membagi dataset menjadi data latih dan uji, serta untuk mengukur performa model secara objektif. Pertama, membuat model *Random Forest* menggunakan *RandomForestClassifier* dari *scikit-learn*. Selanjutnya, melakukan perulangan melalui setiap *fold* yang dihasilkan oleh *k-fold cross-validation*. Pada setiap perulangan, dataset dibagi menjadi data latih dan uji, dan model *Random Forest* dilatih menggunakan data latih. Setelah dilatih, model tersebut digunakan untuk melakukan prediksi pada data uji.

Skor evaluasi akurasi dihitung menggunakan metrik yang sesuai dari library *sklearn.metrics* dan disimpan untuk setiap *fold*. Selain itu, matriks kebingungan (confusion matrix) dari setiap *fold* juga dihitung menggunakan fungsi *confusion_matrix* dari *scikit-learn*. Hasil prediksi dan skor evaluasi ini kemudian disimpan untuk digunakan dalam analisis lebih lanjut tentang performa model pada dataset yang diberikan.

```

from sklearn.model_selection import cross_val_score, StratifiedKFold
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, f1_score
import numpy as np

# Perform cross-validation on the Random Forest classifier with K=10
clf = RandomForestClassifier(max_depth=1, random_state=42)

# Kode kFold Cross Validation
cv = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)

# Lists to store evaluation metrics for each fold
conf_matrices = []
accuracies = []
precisions = []
recalls = []
f1_scores = []

for train_index, test_index in cv.split(X_imb, y_imb):
    X_train, X_test = X_imb.iloc[train_index], X_imb.iloc[test_index]
    y_train, y_test = y_imb.iloc[train_index], y_imb.iloc[test_index]

    # Train the model
    clf.fit(X_train, y_train)

    # Predict on the test set
    y_pred = clf.predict(X_test)

    # Calculate evaluation metrics
    conf_matrices.append(confusion_matrix(y_test, y_pred))
    accuracies.append(accuracy_score(y_test, y_pred))

```

Gambar 2.9 Kode Permodelan Algoritma *Random Forest*

Berikut adalah tabel 2.9 yang berisi tentang tiap parameter yang digunakan dalam kode tersebut beserta penjelasan fungsi tiap parameter:

Tabel 2.9 Parameter Permodelan Algoritma *Random Forest*

Parameter	Keterangan
<i>clf</i>	Model <i>Random Forest</i> yang digunakan untuk melakukan klasifikasi.
<i>train_index</i>	Indeks data untuk data latih.
<i>test_index</i>	Indeks data untuk data uji.
<i>accuracies</i>	List yang akan menyimpan skor akurasi dari setiap <i>fold</i> .
<i>confusion_matrices</i>	List yang akan menyimpan <i>confusion matrix</i> dari setiap <i>fold</i> .
<i>X_train</i>	Subset atribut untuk training.
<i>X_test</i>	Subset atribut untuk testing.
<i>y_train</i>	Subset label untuk training..
<i>y_test</i>	Subset label untuk testing.
<i>y_pred</i>	Hasil prediksi dari subset test.

- d. Langkah keempat Setelah proses *K-fold cross-validation* selesai dilakukan dan model *Random Forest* telah diuji pada setiap *fold*, skor evaluasi seperti akurasi dari setiap *fold* dapat ditampilkan. Dalam pengulangan menggunakan *enumerate* dan *zip*, akurasi dari masing-masing *fold* diproses dan ditsmpilkan. Selain itu, rata-rata akurasi dari semua *fold* juga dihitung dengan menggunakan fungsi *np.mean()* dari *NumPy*. Hasil rata-rata skor evaluasi akurasi dari seluruh *fold* juga dicetak. Terakhir, rata-rata dari *confusion matrix* dari semua *fold*

dihitung menggunakan fungsi `np.mean()` dan ditampilkan. *Confusion matrix*) rata-rata ini memberikan gambaran tentang seberapa baik model dapat mengklasifikasikan setiap kelas pada data uji.

```

import numpy as np
# Print the evaluation metrics for each fold
for i, (conf_matrix, accuracy) in enumerate(zip(conf_matrices, accuracies), 1):
    print(f"Fold-{i}: Accuracy={accuracy}")

# Print the average evaluation metrics
print('')
print('Hasil Rata Rata')
print(f"Average Accuracy: {np.mean(accuracies)}")

# Calculate and print the average confusion matrix
avg_conf_matrix = sum(conf_matrices)
print("Average Confusion Matrix:")
print(avg_conf_matrix)

```

Gambar 2. 10 Kode Hasil akurasi dari setiap *Fold*

Berikut adalah tabel 2.10 yang berisi tentang tiap parameter yang digunakan dalam kode tersebut beserta penjelasan fungsi tiap parameter

Tabel 2. 10 Parameter Hasil Skor Dari Setiap *Fold*

Parameter	Keterangan
<i>np</i>	Alias untuk <i>numpy</i> , sebuah library <i>Python</i> yang digunakan untuk melakukan operasi numerik pada <i>array</i> , matriks, dan struktur data numerik lainnya.
<i>enumerate</i>	Fungsi <i>Python</i> yang digunakan untuk mengambil elemen dari sebuah iterable (seperti <i>list</i> , <i>tuple</i> , atau <i>string</i>) bersamaan dengan indeksnya dalam bentuk <i>tuple</i> .
<i>zip</i>	Fungsi <i>Python</i> yang digunakan untuk menggabungkan elemen-elemen yang bersesuaian dari dua atau lebih <i>iterable</i> (<i>list</i> , <i>tuple</i> , atau <i>array</i>) menjadi <i>tuple-tuple</i> baru
<i>accuracies</i>	List yang menyimpan nilai akurasi untuk setiap <i>fold</i> dalam <i>cross-validation</i>
<i>mean</i>	Fungsi <i>numpy</i> yang digunakan untuk menghitung rata-rata dari sebuah array atau list
<i>axis</i>	Parameter yang digunakan dalam fungsi <i>np.mean</i> untuk menentukan sumbu (<i>axis</i>) mana yang akan dihitung rata-ratanya
<i>avg_cm</i>	Nilai rata-rata dari <i>confusion matrix</i> untuk semua <i>fold</i> dalam <i>cross-validation</i>

- b) Tahap seleksi fitur menggunakan metode *ANOVA* dengan cara menentukan atribut yang memiliki bobot atau pengaruh signifikan terhadap variabel yang akan diklasifikasi. Nilai F yang lebih besar menunjukkan bahwa ada perbedaan yang signifikan antara kelompok, karena menunjukkan bahwa variansi antar kelompok lebih besar dibandingkan dengan variansi dalam kelompok. Sebaliknya,

nilai F yang rendah menunjukkan bahwa perbedaan antar kelompok tidak signifikan, menunjukkan bahwa kelompok-kelompok tersebut memiliki variansi yang serupa. Berikut tahapan kerja dari seleksi fitur *Analysis of Variance* (ANOVA) adalah sebagai berikut :

- a. Semua fitur dipilih pada dataset Stunting
- b. Fungsi fitur target dari scikit-learn dihitung menggunakan *ANOVA F-Score* untuk setiap fitur. Berikut adalah rumus untuk menghitung *ANOVA*:

$$F = \frac{\text{variance between groups}}{\text{variance within groups}} \quad (2.2)$$

$$\text{Variance between groups} = \frac{\sum_i^n n_i (\bar{Y}_i - \bar{Y})^2}{(k - 1)} \quad (2.3)$$

$$\text{Variance within groups} = \frac{\sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y})^2}{(n - k)} \quad (2.4)$$

Keterangan :

- n_i = Jumlah sampel di grup ke-i
- \bar{Y}_i = Rata-rata sampel di grup ke-i
- \bar{Y} = Rata-rata total sampe
- k = Jumlah grup
- n = Jumlah total sampel

- c. Hasil dari pengujian akan digunakan dalam pemilihan fitur yang mempunyai pengaruh tinggi dan varian terendah dipilih dalam percobaan ini lalu akan diuji dengan `SelectKbest()`; k mewakili jumlah fitur yang digunakan untuk dataset akhir
- d. Jumlah fitur dengan pengaruh tertinggi akan digunakan untuk membuat berbagai subset fitur.

```

from sklearn.feature_selection import SelectKBest, f_classif

X = data

# Menggunakan SelectKBest dengan ANOVA untuk memilih fitur terbaik
selector = SelectKBest(score_func=f_classif)
X_selected = selector.fit_transform(X, y)

# Mengambil nilai F dan nama fitur dari SelectKBest
f_values = selector.scores_
features = X.columns

# Membuat DataFrame dari fitur dan nilai F
import pandas as pd
df_f = pd.DataFrame({'Feature': features, 'F_Value': f_values})

# Sorting DataFrame berdasarkan nilai F secara ascending
df_sorted = df_f.sort_values('F_Value', ascending=False)

# Membuat diagram bar untuk menampilkan nilai F dari semua fitur
plt.figure(figsize=(12, 6))
plt.bar(df_sorted['Feature'], df_sorted['F_Value'], color='skyblue')
plt.xlabel('Fitur')
plt.ylabel('Nilai F')
plt.title('Perbandingan Nilai F dari Semua Fitur (Ascending)')
plt.xticks(rotation=90)
plt.show()

```

Gambar 2. 11 Kode Seleksi Fitur *ANOVA*

Berikut adalah tabel 2.11 yang berisi tentang tiap parameter yang digunakan dalam kode tersebut beserta penjelasan fungsi tiap parameter

Tabel 2.11 Parameter Seleksi Fitur *ANOVA*

Fungsi	Parameter	Keterangan
<i>SelectKBest</i>	<i>score_func=f_classif</i>	<i>SelectKBest</i> adalah metode untuk seleksi fitur yang memilih K fitur terbaik berdasarkan nilai statistik tertentu. <i>score_func = f_classif</i> menunjukkan penggunaan <i>F-test ANOVA</i> untuk tugas klasifikasi, yang mengukur pengaruh tiap fitur terhadap <i>variabilitas</i> antara kelas.
<i>fit_transform()</i>	<i>X, y</i>	Metode <i>fit_transform()</i> digunakan untuk mengajarkan <i>selector</i> pada data X (fitur) dan y (target) dan mengubah X menjadi subset fitur yang terpilih berdasarkan skor tertinggi.
<i>DataFrame()</i>	{'Feature': <i>features</i> , 'F_Value': <i>f_values</i> }	Membuat <i>DataFrame</i> dari dua <i>arrays</i> , <i>features</i> yang berisi nama fitur dari X dan <i>f_values</i> yang berisi nilai F dari setiap fitur setelah dihitung oleh <i>SelectKBest</i> .
<i>sort_values()</i>	'F_Value', <i>ascending=False</i>	Mengurutkan <i>DataFrame</i> berdasarkan kolom <i>F_Value</i> secara <i>descending</i> untuk melihat fitur dengan pengaruh paling tinggi terlebih dahulu.

2.2.5 Evaluasi

Ditahapan evaluasi akan dilakukan pengukuran akurasi dari algoritma yang digunakan dengan kualitas data training serta akan diuji dengan menggunakan teknik *Confusion Matrix*.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \times 100\% \quad (2.5)$$

Keterangan :

- TP (*True Positive*) : Jumlah data point berlabel *yes* yang nilainya diidentifikasi benar.
- TP (*True Negative*) : Jumlah data point berlabel *no* yang nilainya diidentifikasi salah.
- FP (*False Positive*) : Jumlah data point berlabel *yes* yang nilai sebenarnya diidentifikasi salah
- FN (*False Negative*) : Jumlah data point berlabel *no* yang nilai sebenarnya teridentifikasi benar