

Laporan pengerjaan Pertemuan 5

Langkah 1 — Muat Data

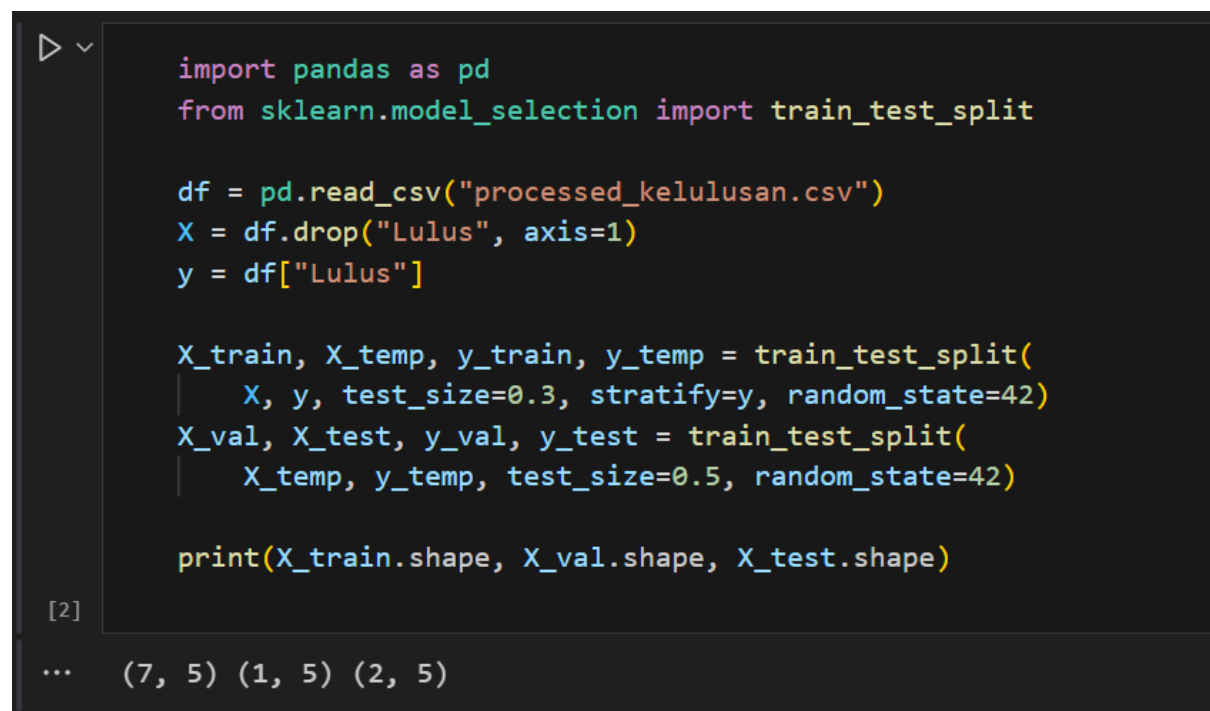
1. pakai file split dari Pertemuan 4, lalu copy paste dari lembar kerja ke vscode

```
import pandas as pd

X_train = pd.read_csv("X_train.csv")
X_val   = pd.read_csv("X_val.csv")
X_test  = pd.read_csv("X_test.csv")
y_train = pd.read_csv("y_train.csv").squeeze("columns")
y_val   = pd.read_csv("y_val.csv").squeeze("columns")
y_test  = pd.read_csv("y_test.csv").squeeze("columns")

print(X_train.shape, X_val.shape, X_test.shape)
```

2. Hasilnya seperti berikut ini



```
import pandas as pd
from sklearn.model_selection import train_test_split

df = pd.read_csv("processed_kelulusan.csv")
X = df.drop("Lulus", axis=1)
y = df["Lulus"]

X_train, X_temp, y_train, y_temp = train_test_split(
    X, y, test_size=0.3, stratify=y, random_state=42)
X_val, X_test, y_val, y_test = train_test_split(
    X_temp, y_temp, test_size=0.5, random_state=42)

print(X_train.shape, X_val.shape, X_test.shape)
```

[2]

... (7, 5) (1, 5) (2, 5)

Langkah 2 — Baseline Model & Pipeline

Bangun baseline terstandar menggunakan Logistic Regression + pipeline preprocessing.

1. Copy paste coding dari lembar kerja ke vscode dan dapat hasilnya seperti ini

```
from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import StandardScaler
from sklearn.impute import SimpleImputer
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import f1_score, classification_report

num_cols = X_train.select_dtypes(include="number").columns

pre = ColumnTransformer([
    ("num", Pipeline([("imp", SimpleImputer(strategy="median")),
                     ("sc", StandardScaler())]), num_cols),
], remainder="drop")

logreg = LogisticRegression(max_iter=1000, class_weight="balanced", random_state=42)
pipe_lr = Pipeline([("pre", pre), ("clf", logreg)])

pipe_lr.fit(X_train, y_train)
y_val_pred = pipe_lr.predict(X_val)
print("Baseline (LogReg) F1(val):", f1_score(y_val, y_val_pred, average="macro"))
print(classification_report(y_val, y_val_pred, digits=3))
```

[3]

```
... Baseline (LogReg) F1(val): 1.0
      precision    recall  f1-score   support

      1         1.000      1.000      1.000         1

 accuracy         1.000
 macro avg         1.000
 weighted avg         1.000
```

Baseline berfungsi sebagai acuan. Peningkatan harus dibuktikan dengan metrik, bukan asumsi.

Langkah 3 — Model Alternatif (Random Forest)

Copy paste coding dari lembar kerja ke vscode dan dapat hasilnya seperti ini

```
from sklearn.ensemble import RandomForestClassifier

rf = RandomForestClassifier(
    n_estimators=300, max_features="sqrt", class_weight="balanced", random_state=42
)
pipe_rf = Pipeline([("pre", pre), ("clf", rf)])

pipe_rf.fit(X_train, y_train)
y_val_rf = pipe_rf.predict(X_val)
print("RandomForest F1(val):", f1_score(y_val, y_val_rf, average="macro"))
```

```
RandomForest F1(val): 1.0
```

Langkah 4 — Validasi Silang & Tuning Ringkas

Copy paste coding dari lembar kerja ke vscode dan dapat hasilnya seperti ini

```
from sklearn.model_selection import StratifiedKFold, GridSearchCV

skf = StratifiedKFold(n_splits=3, shuffle=True, random_state=42)
param = {
    "clf__max_depth": [None, 12, 20, 30],
    "clf__min_samples_split": [2, 5, 10]
}
gs = GridSearchCV(pipe_rf, param_grid=param, cv=skf,
                  scoring="f1_macro", n_jobs=-1, verbose=1)
gs.fit(X_train, y_train)
print("Best params:", gs.best_params_)
print("Best CV F1:", gs.best_score_)

best_rf = gs.best_estimator_
y_val_best = best_rf.predict(X_val)
print("Best RF F1(val):", f1_score(y_val, y_val_best, average="macro"))
```

```
Fitting 3 folds for each of 12 candidates, totalling 36 fits
Best params: {'clf__max_depth': None, 'clf__min_samples_split': 2}
Best CV F1: 1.0
Best RF F1(val): 1.0
```

Langkah 5 — Evaluasi Akhir (Test Set)

1. Copy paste code dari lembar kerja ke vscode

```
from sklearn.metrics import confusion_matrix, roc_auc_score, precision_recall_curve, roc_curve
import matplotlib.pyplot as plt

final_model = best_rf # atau pipe_lr jika baseline lebih baik
y_test_pred = final_model.predict(X_test)

print("F1(test):", f1_score(y_test, y_test_pred, average="macro"))
print(classification_report(y_test, y_test_pred, digits=3))
print("Confusion matrix (test):")
print(confusion_matrix(y_test, y_test_pred))

# ROC-AUC (jika ada predict_proba)
if hasattr(final_model, "predict_proba"):
    y_test_proba = final_model.predict_proba(X_test)[:,1]
    try:
        print("ROC-AUC(test):", roc_auc_score(y_test, y_test_proba))
    except:
        pass
    fpr, tpr, _ = roc_curve(y_test, y_test_proba)
    plt.figure(); plt.plot(fpr, tpr); plt.xlabel("FPR"); plt.ylabel("TPR"); plt.title("ROC (test)")
    plt.tight_layout(); plt.savefig("roc_test.png", dpi=120)
```

2. Hasilnya seperti ini

```
F1(test): 1.0
      precision    recall  f1-score   support

     0       1.000      1.000      1.000         2

 accuracy          1.000
 macro avg          1.000
 weighted avg       1.000

Confusion matrix (test):
[[2]]
ROC-AUC(test): nan
d:\machine_learning\.venv\lib\site-packages\sklearn\metrics\classification.py:534: UserWarning: A single label
warnings.warn(
d:\machine_learning\.venv\lib\site-packages\sklearn\metrics\ranking.py:424: UndefinedMetricWarning: Only one cl
warnings.warn(
d:\machine_learning\.venv\lib\site-packages\sklearn\metrics\ranking.py:1201: UndefinedMetricWarning: No positiv
warnings.warn(
```

