

```

import pandas as pd
import numpy as np
from scipy.io import arff
from sklearn.model_selection import KFold
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.metrics import accuracy_score
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from xgboost import XGBClassifier
from sklearn.neural_network import MLPClassifier
from scipy.spatial.distance import jensenshannon
from scipy.stats import wasserstein_distance
from tabulate import tabulate
import sys
import os

# Set environment variable to allow CPU for large dataset
os.environ["TABPFN_ALLOW_CPU_LARGE_DATASET"] = "1"

# Add local TabPFGen path
sys.path.insert(0, './src/tabpfngen_backup')
from tabpfngen import TabPFGen
import warnings
warnings.filterwarnings("ignore")

# =====
# Step 1: Load the Dataset
# =====
file_path = r'C:\Users\Manthan Goyal\Desktop\Team-Project\TabPFGen\datasets\dermatology.arff'
data, meta = arff.loadarff(file_path)
df = pd.DataFrame(data)

# Decode byte strings
df = df.applymap(lambda x: x.decode('utf-8') if isinstance(x, bytes) else x)

# =====
# Step 2: Preprocessing
# =====
df.replace('?', np.nan, inplace=True)

# Convert all columns to numeric where possible
for column in df.columns:
    df[column] = pd.to_numeric(df[column], errors='coerce')

# Drop columns that are entirely NaN
df.dropna(axis=1, how='all', inplace=True)
df.dropna(inplace=True)

# Encode target variable
target_col = 'class'
le = LabelEncoder()
df[target_col] = le.fit_transform(df[target_col])

# Separate features and target
X = df.drop(columns=[target_col])
y = df[target_col]

# Standardize the features
scaler = StandardScaler()
X = pd.DataFrame(scaler.fit_transform(X), columns=X.columns)

# =====
# Step 3: Cross-Validation and Evaluation
# =====
kf = KFold(n_splits=2, shuffle=True, random_state=42)
classifiers = {
    'LR': LogisticRegression(max_iter=1000),
    'RF': RandomForestClassifier(),
    'XG BOOST': XGBClassifier(eval_metric='logloss', use_label_encoder=False),
    'MLP': MLPClassifier(max_iter=500)
}

# Storage for results
results = []

# Perform 3 Repeats of 2-Fold Cross-Validation
for repeat in range(1, 4):
    for fold, (train_index, test_index) in enumerate(kf.split(X), 1):
        # Split the data
        train_X, test_X = X.iloc[train_index], X.iloc[test_index]
        train_y, test_y = y.iloc[train_index], y.iloc[test_index]

        # Generate Synthetic Data
        generator = TabPFGen(n_sgld_steps=100)
        X_synth, y_synth = generator.generate_classification(
            train_X.to_numpy(),
            train_y.to_numpy(),
            int(0.5 * len(train_X)) #50% OF THE DATA
        )

        # Evaluate each classifier
        for name, clf in classifiers.items():
            clf.fit(train_X, train_y)
            real_acc = accuracy_score(test_y, clf.predict(test_X))
            results.append([repeat, fold, name, real_acc])

        # JSD and Wasserstein calculations
        real_dist = np.bincount(train_y, minlength=len(np.unique(train_y))) / len(train_y)
        synth_dist = np.bincount(y_synth, minlength=len(np.unique(train_y))) / len(y_synth)
        jsd_value = jensenshannon(real_dist, synth_dist) if len(real_dist) == len(synth_dist) else np.nan
        wd_value = wasserstein_distance(np.sort(train_y), np.sort(y_synth))

        # Store distance metrics
        results.append([repeat, fold, "JSD", jsd_value])
        results.append([repeat, fold, "WD", wd_value])

# =====
# Step 4: Format the Output
# =====
# Create a new DataFrame to match the required structure
models = ['LR', 'RF', 'MLP', 'XG BOOST', 'JSD', 'WD']
columns = ['R1-F1', 'R1-F2', 'R2-F1', 'R2-F2', 'R3-F1', 'R3-F2', 'AVERAGE']
output_df = pd.DataFrame(index=models, columns=columns)

# Fill the DataFrame with values
for repeat in range(1, 4):
    for fold in range(1, 3):
        col_name = f'R{repeat}-F{fold}'

        for model in models:
            value = [x[3] for x in results if x[0] == repeat and x[1] == fold and x[2] == model]
            if value:
                output_df.at[model, col_name] = value[0]

```

```
# Calculate the AVERAGE for each row
output_df['AVERAGE'] = output_df.iloc[:, :-1].apply(pd.to_numeric, errors='coerce').mean(axis=1)

# =====
# Step 5: Display in Terminal
# =====
print("\n    Final Cross-Validation Summary:\n")
print(tabulate(output_df, headers='keys', tablefmt='grid'))
```