# Machine Learning Project 1
# Scikit-Learn and Cross Validation Project

Lakshmi Chakradhar Vijayarao

Net ID: lxv230000

## Introduction

In this project, I built classification models to check if the mushroom is edible or poisonous using the mushroom_m i x e d _ 50000.csv dataset. I used three machine learning algorithms: Decision Tree Classifier, Random Forest Classifier and KNN Classifier. The models were evaluated using grid search and cross-validation, and the final model was saved for accuracy competition.

## Dataset

The dataset contains 50,000 samples with 20 features. The target variable has two classes: e (edible mushroom) and p (poisonous mushroom).

## Model Building

I used the scikit-learn library to implement two models:

- Decision Tree Classifier
- Random Forest Classifier
- KNN Classifier

### Cross-Validation and Grid Search

I performed grid search with 5-fold cross-validation. The parameter grids for each model were as follows:

| Model | Parameter | Values |
|---|---|---|
| Decision Tree | max depth | [10, 20, 30] |
| | min samples split | [2, 5, 10] |
| | min samples leaf | [1, 2, 4] |
| Random Forest | n estimators | [50, 100] |
| | max depth | [10, 15] |
| | min samples split | [5, 10] |
| | min samples leaf | [1, 2, 4] |
| | max features | ['sqrt', 'log2'] |
| KNN | neighbors | [3. 5. 7] |
| | Weights. | ['uniform' , 'distance'] |

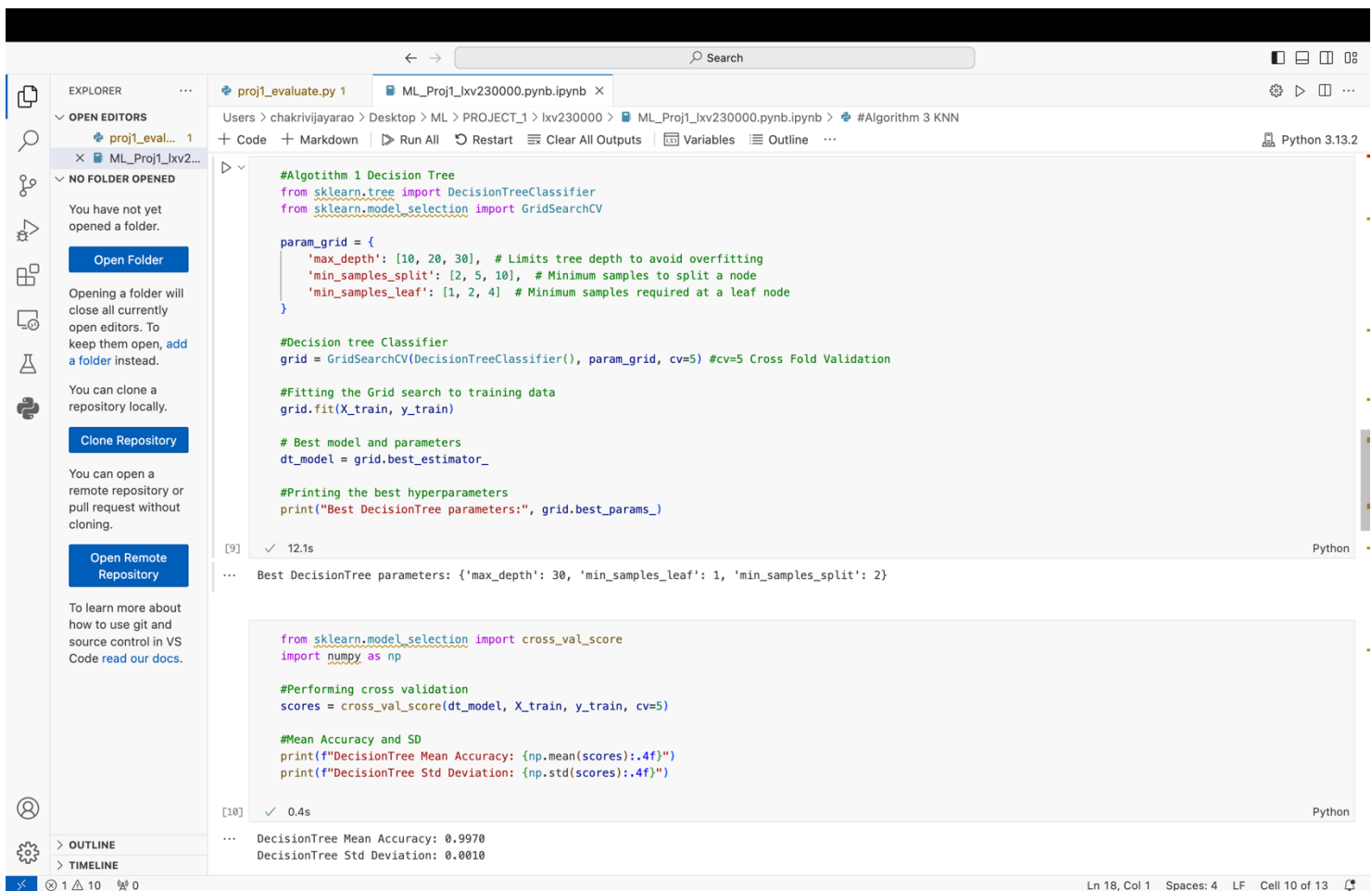Table 1: Grid search parameters for Decision Tree, Random Forest and KNN

# Documentation

The approach for this project can be summarized in the following steps:

- **Data Loading**: I began by loading the dataset using pandas. The dataset was read into a DataFrame, where the features (X) were separated from the target variable (y).

- **Train-Test Split**: I split the data into training and testing sets using train test split from scikit-learn. This ensured that our model would be evaluated on unseen data, providing a fair assessment of its performance.

- **Model Definition**: I defined three classifiers: DecisionTreeClassifier, RandomForestClassifier and KNNClassifier.

- **Parameter Tuning**: To find the best hyperparameters, we employed Grid Search with cross-validation. This involved defining a parameter grid for each model, which allowed me to systematically explore different combinations of parameters to find the optimal settings for our classifiers.

- **Model Evaluation**: After fitting the models with the training data, we evaluated their performance using accuracy scores calculated on the test data. I also recorded the mean accuracy and standard deviation from the cross-validation results.

- **Model Saving**: The best performing model based on accuracy was saved for future use. I utilized Python's pickle module for serialization.

# Results

## Decision Tree:

## Random Forest:



```python
#Algorithm 2 Random Forest
from sklearn.ensemble import RandomForestClassifier

param_grid = {
    'n_estimators': [50, 100],  # Number of trees to balance accuracy and execution time
    'max_depth': [10, 15],  # Reduces complexity and overfitting
    'min_samples_split': [5, 10],  # Ensures nodes have enough samples before splitting
    'min_samples_leaf': [2, 4],  # Avoids overfitting on small splits
    'max_features': ['sqrt', 'log2']  # Reduces dimensionality impact
}

#Random Forest Classifier
grid = GridSearchCV(RandomForestClassifier(), param_grid, cv=5)
#Fitting the Grid search to training data
grid.fit(X_train, y_train)

# Best model and parameters
rf_model = grid.best_estimator_

#Printing the best hyperparameters
print("Best RandomForest parameters:", grid.best_params_)
```
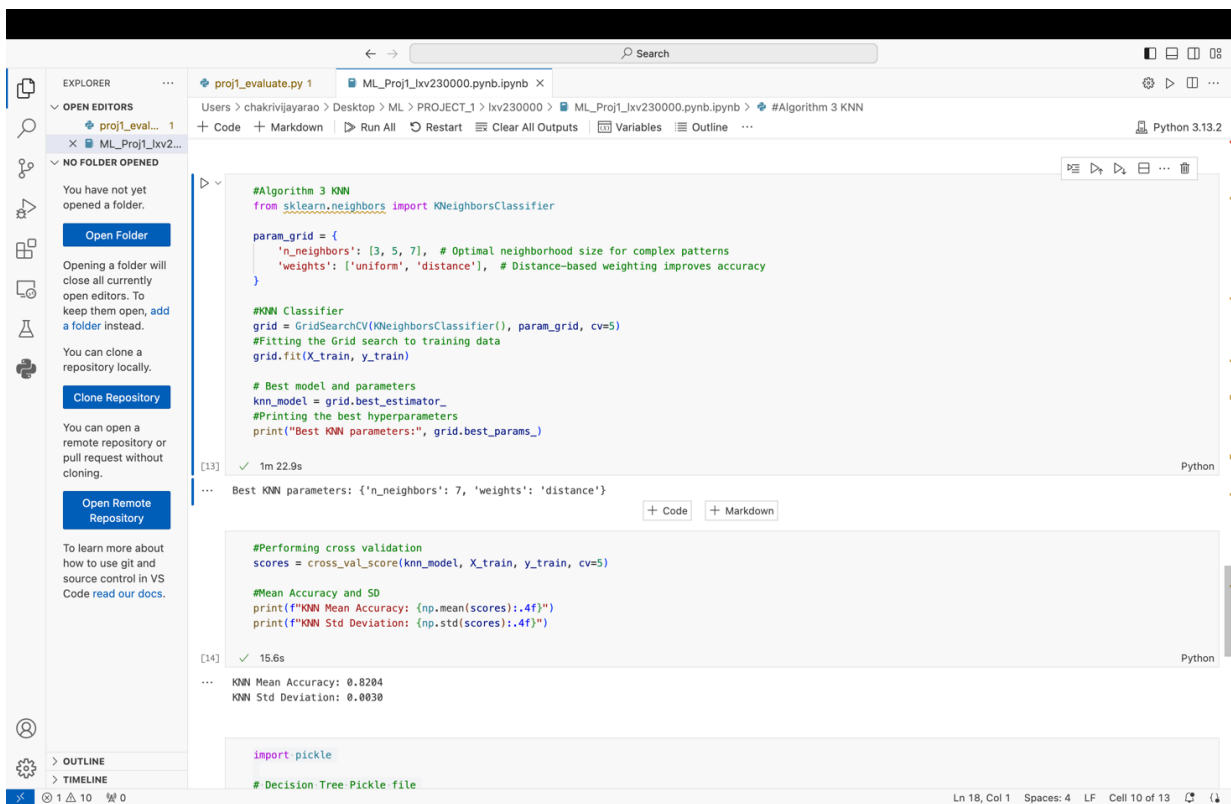
[11] ✓ 2m 56.6s

Best RandomForest parameters: {'max_depth': 15, 'max_features': 'log2', 'min_samples_leaf': 2, 'min_samples_split': 5, 'n_estimators': 100}

```python
#Performing cross validation
scores = cross_val_score(rf_model, X_train, y_train, cv=5)

#Mean Accuracy and SD
print(f"RandomForest Mean Accuracy: {np.mean(scores):.4f}")
print(f"RandomForest Std Deviation: {np.std(scores):.4f}")
```

[12] ✓ 7.8s

RandomForest Mean Accuracy: 0.9996
RandomForest Std Deviation: 0.0003

## KNN:



```python
#Algorithm 3 KNN
from sklearn.neighbors import KNeighborsClassifier

param_grid = {
    'n_neighbors': [3, 5, 7],  # Optimal neighborhood size for complex patterns
    'weights': ['uniform', 'distance'],  # Distance-based weighting improves accuracy
}

#KNN Classifier
grid = GridSearchCV(KNeighborsClassifier(), param_grid, cv=5)
#Fitting the Grid search to training data
grid.fit(X_train, y_train)

# Best model and parameters
knn_model = grid.best_estimator_
#Printing the best hyperparameters
print("Best KNN parameters:", grid.best_params_)
```

[13] ✓ 1m 22.9s

Best KNN parameters: {'n_neighbors': 7, 'weights': 'distance'}

```python
#Performing cross validation
scores = cross_val_score(knn_model, X_train, y_train, cv=5)

#Mean Accuracy and SD
print(f"KNN Mean Accuracy: {np.mean(scores):.4f}")
print(f"KNN Std Deviation: {np.std(scores):.4f}")
```

[14] ✓ 15.6s

KNN Mean Accuracy: 0.8204
KNN Std Deviation: 0.0030

```python
import pickle

# Decision Tree Pickle file
```

**The results of grid search with cross-validation are summarized below:**

| Hyperparameter | Decision Tree Classifier | Random Forest Classifier | KNN Classifier |
|---|---|---|---|
| max depth | 30 | 15 | - |
| max features | - | Log2 | - |
| min samples leaf | 1 | 2 | - |
| min samples split | 2 | 5 | - |
| n estimators | - | 100 | - |
| N neighbours | - | - | 7 |
| weights | - | - | distance |

Table 2: Best Parameters Found for Decision Tree and Random Forest Classifiers

| Model | Mean Accuracy | Accuracy Std. Dev. |
|---|---|---|
| **Decision Tree** | 0.9970 | 0.0010 |
| **Random Forest** | 0.9996 | 0.0003 |
| **KNN** | 0.8204 | 0.0030 |

Table 3: Cross-Validation and Test Accuracy Results

Evalvuation Program Output:



## **Conclusion**

The Random Forest model achieved the highest accuracy of 99.96%, and it was saved as proj1_chosen_model.pkl.