

```
In [1]: 1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import train_test_split, GridSearchCV
4 import matplotlib.pyplot as plt
```

```
In [2]: 1 data = pd.read_csv("C:/Users/MANOJ/Documents/week5/nursery/nursery.data", names=[ 'parents', 'has_nur
2 data.head()
```

Out[2]:

	parents	has_nurs	form	children	housing	finance	social	health	class
0	usual	proper	complete	1	convenient	convenient	nonprob	recommended	recommend
1	usual	proper	complete	1	convenient	convenient	nonprob	priority	priority
2	usual	proper	complete	1	convenient	convenient	nonprob	not_recom	not_recom
3	usual	proper	complete	1	convenient	convenient	slightly_prob	recommended	recommend
4	usual	proper	complete	1	convenient	convenient	slightly_prob	priority	priority

```
In [3]: 1 data['class'].unique() # 5 class dataset
```

Out[3]: array(['recommend', 'priority', 'not\_recom', 'very\_recom', 'spec\_prior'],  
dtype=object)

```
In [4]: 1 data['class'].value_counts()
```

Out[4]: not\_recom 4320  
priority 4266  
spec\_prior 4044  
very\_recom 328  
recommend 2  
Name: class, dtype: int64

```
In [5]: 1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 12960 entries, 0 to 12959
Data columns (total 9 columns):
#   Column      Non-Null Count  Dtype
---  -
0   parents    12960 non-null  object
1   has_nurs   12960 non-null  object
2   form       12960 non-null  object
3   children   12960 non-null  object
4   housing    12960 non-null  object
5   finance    12960 non-null  object
6   social     12960 non-null  object
7   health     12960 non-null  object
8   class      12960 non-null  object
dtypes: object(9)
memory usage: 911.4+ KB
```

```
In [6]: 1 X = data.drop('class',axis=1)
2 Y = data['class']
```

```
In [7]: 1 Y[Y == 'spec_prior'] = 'recommend'
2 Y[Y == 'very_recom'] = 'recommend'
3
```

```
In [8]: 1 Y.value_counts()
```

Out[8]: recommend 4374  
not\_recom 4320  
priority 4266  
Name: class, dtype: int64

```
In [9]: 1 # To store the test score of 5 iterations in order to compute the mean and variance of testing acc
2 DecisionTree_categorical = []
3 DecisionTree_oneHot = []
4 Log_Regression = []
5 K_Nearest_Neighbors = []
```

The task is conducted 5 times, therefore the data would be split into train, validation and test 5 times randomly.

```
In [47]: 1 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=42)
```

```
In [48]: 1 x_train, x_val, y_train, y_val = train_test_split(X_train, Y_train, test_size=0.1, random_state=42)
```

To fit the Decision Tree with Categorical data

```
In [49]: 1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn import tree
3 from sklearn.preprocessing import LabelEncoder
```

```
In [50]: 1 encoder = LabelEncoder()
2 for i in X.columns: # Encoding the categorical values of features
3     x_train[i] = encoder.fit_transform(x_train[i])
4     x_val[i] = encoder.fit_transform(x_val[i])
5     X_test[i] = encoder.fit_transform(X_test[i])
```

```
In [51]: 1 model = DecisionTreeClassifier().fit(x_train,y_train)
```

```
In [52]: 1 parameters = {'max_depth':[i for i in range(1,11)],
2               'criterion':['entropy','gini'],
3               'splitter':['best','random'],
4               'min_samples_split': [i for i in range(2,11)],
5               'min_samples_leaf': [i for i in range(1,11)]}
```

#### The parameters considered for tuning in Decision Trees:

1. The max depth parameter is searched in range (1,10) as it covers shallow trees and trees of considerable depth.
2. The criterion is chosen to be either gini or entropy which is based on gini index and information gain at each node respectively.
3. The splitter is chosen as either best or random. best chooses the best split at each node based on criterion and random chooses a best split among a random subset of features.
4. min\_samples\_split range from (2,10), determining the optimum number of samples required to split at a node. Higher value prevents overfitting therefore it is searched till the value 10.
5. min\_samples\_leaf is searched between (1,10) which specifies the minimum no of samples to be present at each leaf node.

```
In [53]: 1 C = GridSearchCV(model, param_grid=parameters, cv=5)
2 C.fit(x_val,y_val)
3 Model_categorical = C.best_estimator_
4 print(f"Best Decision Tree parameters: {C.best_params_}")
5
6 Model_categorical.fit(x_train, y_train)
7 DecisionTree_categorical.append(Model_categorical.score(X_test,Y_test))
```

```
Best Decision Tree parameters: {'criterion': 'gini', 'max_depth': 9, 'min_samples_leaf': 1, 'min_samples_split': 6, 'splitter': 'best'}
```

```
In [54]: 1 DecisionTree_categorical
```

```
Out[54]: [0.9347993827160493,
0.9621913580246914,
0.9621913580246914,
0.9533179012345679,
0.9533179012345679]
```

```
In [55]: 1 A = [round(i*100,2) for i in DecisionTree_categorical]
2 DecisionTree_categorical,A
```

```
Out[55]: ([0.9347993827160493,
0.9621913580246914,
0.9621913580246914,
0.9533179012345679,
0.9533179012345679],
[93.48, 96.22, 96.22, 95.33, 95.33])
```

```
In [56]: 1 dt_mean_categorical = np.mean(A)
2 dt_var_categorical = np.var(A)
3 print("Mean of Test Performance:",round(dt_mean_categorical,2))
4 print("Variance of Test Performance:",round(dt_var_categorical,2))
```

Mean of Test Performance: 95.32  
Variance of Test Performance: 1.0

To fit Decision Tree with one hot encoded labels

```
In [58]: 1 for i in range(5):
2     X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=42)
3     x_train, x_val, y_train, y_val = train_test_split(X_train, Y_train, test_size=0.1, random_state=42)
4
5     # To get One hot encoded Labels
6     x_train_onehot = pd.get_dummies(x_train)
7     x_val_onehot = pd.get_dummies(x_val)
8     X_test_onehot = pd.get_dummies(X_test)
9
10    model_oneHot = DecisionTreeClassifier().fit(x_train_onehot,y_train)
11
12    parameters = {'max_depth':[i for i in range(1,11)],
13                  'criterion':['entropy','gini'],
14                  'splitter':['best','random'],
15                  'min_samples_split': [i for i in range(2,11)],
16                  'min_samples_leaf': [i for i in range(1,11)]}
17
18    C = GridSearchCV(model_oneHot, param_grid=parameters, cv=5)
19    C.fit(x_val_onehot,y_val)
20    Model_oneHot = C.best_estimator_
21    print(f"Best Decision Tree parameters: {C.best_params_}")
22
23    Model_oneHot.fit(x_train_onehot, y_train)
24    DecisionTree_oneHot.append(Model_oneHot.score(X_test_onehot,Y_test))
```

Best Decision Tree parameters: {'criterion': 'entropy', 'max\_depth': 10, 'min\_samples\_leaf': 1, 'min\_samples\_split': 2, 'splitter': 'best'}

Best Decision Tree parameters: {'criterion': 'entropy', 'max\_depth': 10, 'min\_samples\_leaf': 1, 'min\_samples\_split': 2, 'splitter': 'best'}

C:\Users\MANOJ\AppData\Roaming\Python\Python310\site-packages\numpy\ma\core.py:2820: RuntimeWarning: invalid value encountered in cast

\_data = np.array(data, dtype=dtype, copy=copy,

Best Decision Tree parameters: {'criterion': 'entropy', 'max\_depth': 10, 'min\_samples\_leaf': 1, 'min\_samples\_split': 2, 'splitter': 'random'}

Best Decision Tree parameters: {'criterion': 'entropy', 'max\_depth': 10, 'min\_samples\_leaf': 1, 'min\_samples\_split': 2, 'splitter': 'random'}

Best Decision Tree parameters: {'criterion': 'gini', 'max\_depth': 9, 'min\_samples\_leaf': 1, 'min\_samples\_split': 6, 'splitter': 'best'}

```
In [60]: 1 B = [round(i*100,2) for i in DecisionTree_oneHot]
2 DecisionTree_oneHot,B
```

```
Out[60]: ([0.9594907407407407,
0.9583333333333334,
0.9591049382716049,
0.9571759259259259,
0.9571759259259259],
[95.95, 95.83, 95.91, 95.72, 95.72])
```

```
In [61]: 1 dt_mean_oneHot = np.mean(B)
2 dt_var_oneHot = np.var(B)
3 print("Mean of Test Performance:",round(dt_mean_oneHot,2))
4 print("Variance of Test Performance:",round(dt_var_oneHot,2))
```

Mean of Test Performance: 95.83  
Variance of Test Performance: 0.01

### K-Nearest Neighbors

```
In [62]: 1 from sklearn.neighbors import KNeighborsClassifier
```

```
In [77]: 1 K_Nearest_Neighbors = []
2 for i in range(5):
3     X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=42)
4     x_train, x_val, y_train, y_val = train_test_split(X_train, Y_train, test_size=0.1, random_state=42)
5
6     # To get One hot encoded Labels
7     x_train_onehot = pd.get_dummies(x_train)
8     x_val_onehot = pd.get_dummies(x_val)
9     X_test_onehot = pd.get_dummies(X_test)
10
11     Model3 = KNeighborsClassifier()
12     Model3.fit(x_train_onehot,y_train)
13
14     parameter_knn={'n_neighbors' : [i for i in range(2,np.random.randint(5,15))]}
15     # Searching for the most optimal number of Neighbors
16     K = GridSearchCV(Model3, param_grid=parameter_knn, cv=5)
17     K.fit(x_val_onehot,y_val)
18     Model_KNN = K.best_estimator_
19     Model_KNN.fit(x_train_onehot,y_train)
20     print(f"Best KNN parameters: {K.best_params_}")
21
22     K_Nearest_Neighbors.append(Model_KNN.score(X_test_onehot,Y_test))
23
24
```

Best KNN parameters: {'n\_neighbors': 4}  
Best KNN parameters: {'n\_neighbors': 8}  
Best KNN parameters: {'n\_neighbors': 10}  
Best KNN parameters: {'n\_neighbors': 6}  
Best KNN parameters: {'n\_neighbors': 10}

```
In [79]: 1 k = [round(i*100,2) for i in K_Nearest_Neighbors]
2 K_Nearest_Neighbors,k
```

```
Out[79]: ([0.9120370370370371,
0.9483024691358025,
0.9629629629629629,
0.9402006172839507,
0.9629629629629629],
[91.2, 94.83, 96.3, 94.02, 96.3])
```

```
In [80]: 1 k_mean = np.mean(k)
2 k_var = np.var(k)
3 print("Mean of Test Performance:",round(k_mean,2))
4 print("Variance of Test Performance:",round(k_var,2))
```

Mean of Test Performance: 94.53  
Variance of Test Performance: 3.54

### Logistic Regression

```
In [81]: 1 from sklearn.linear_model import LogisticRegression
```

```

In [114]: 1 Log_Regression = []
          2 for i in range(5):
          3     X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=42)
          4     x_train, x_val, y_train, y_val = train_test_split(X_train, Y_train, test_size=0.1, random_state=42)
          5
          6     # To get One hot encoded Labels
          7     x_train_onehot = pd.get_dummies(x_train)
          8     x_val_onehot = pd.get_dummies(x_val)
          9     X_test_onehot = pd.get_dummies(X_test)
          10
          11     Model2 = LogisticRegression().fit(x_train_onehot,y_train)
          12     parameters_logistic = {'penalty': ['l1'],
          13                           'solver': ['saga'],
          14                           'C': [np.random.uniform(0.02,0.9)], 'max_iter':[4000]}
          15
          16     L = GridSearchCV(Model2, param_grid=parameters_logistic, cv=5)
          17     L.fit(x_val_onehot,y_val)
          18     Model_Logistic = L.best_estimator_
          19     Model_Logistic.fit(x_train_onehot, y_train)
          20     print(f"Best Logistic Regression parameters: {L.best_params_}")
          21
          22     Log_Regression.append(Model_Logistic.score(X_test_onehot,Y_test))
          23

```

```

Best Logistic Regression parameters: {'C': 0.0481045829569674, 'max_iter': 4000, 'penalty': 'l1', 'solver': 'saga'}
Best Logistic Regression parameters: {'C': 0.45728212719036515, 'max_iter': 4000, 'penalty': 'l1', 'solver': 'saga'}
Best Logistic Regression parameters: {'C': 0.3522040644630245, 'max_iter': 4000, 'penalty': 'l1', 'solver': 'saga'}
Best Logistic Regression parameters: {'C': 0.2298302051294932, 'max_iter': 4000, 'penalty': 'l1', 'solver': 'saga'}
Best Logistic Regression parameters: {'C': 0.7035688089332067, 'max_iter': 4000, 'penalty': 'l1', 'solver': 'saga'}

```

The parameters to tune in Logistic Regression:

1. Penalty is set to L1
2. solver is set to saga as it supports both L1 penalty and is suitable for Multi class problem. It is algorithm used in the optimisation problem.
3. Looking for the best lambda(coefficient of the regularization term)

```

In [101]: 1 Log_Regression

```

```

Out[101]: [0.9128086419753086,
           0.9170524691358025,
           0.9166666666666666,
           0.9170524691358025,
           0.9147376543209876]

```

```

In [102]: 1 l = [round(i*100,2) for i in Log_Regression]
          2 Log_Regression,l

```

```

Out[102]: ([0.9128086419753086,
            0.9170524691358025,
            0.9166666666666666,
            0.9170524691358025,
            0.9147376543209876],
           [91.28, 91.71, 91.67, 91.71, 91.47])

```

```

In [103]: 1 l_mean = np.mean(l)
          2 l_var = np.var(l)
          3 print("Mean of Test Performance:",round(l_mean,2))
          4 print("Variance of Test Performance:",round(l_var,2))

```

```

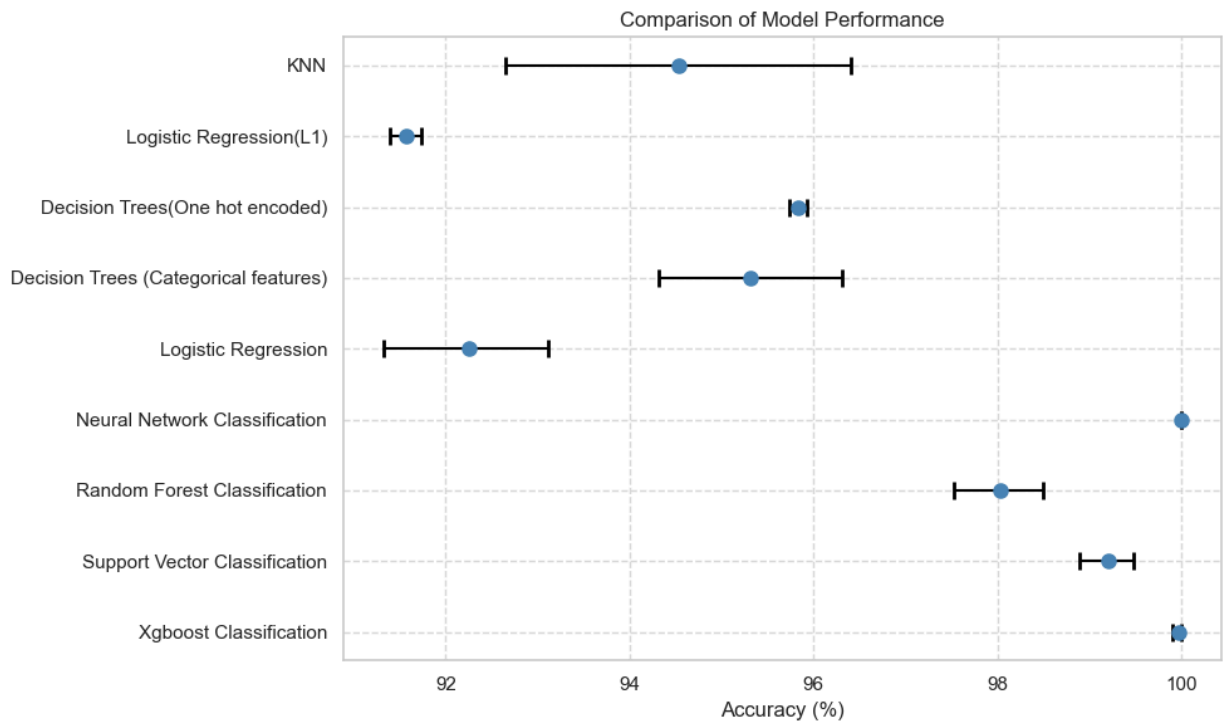
Mean of Test Performance: 91.57
Variance of Test Performance: 0.03

```

```

In [115]: 1 # Plotting the data
2 models = ['Xgboost Classification', 'Support Vector Classification', 'Random Forest Classification',
3           'Neural Network Classification', 'Logistic Regression',
4           'Decision Trees (Categorical features)',
5           'Decision Trees(One hot encoded)', 'Logistic Regression(L1)', 'KNN']
6
7 # Define mean accuracy values and error ranges for each model
8 means = [99.969, 99.198, 98.025, 100, 92.253, dt_mean_categorical, dt_mean_oneHot, l_mean, k_mean
9 lower_error = [0.062, 0.309, 0.494, 0.0, 0.926, np.sqrt(dt_var_categorical), np.sqrt(dt_var_oneHot),
10               np.sqrt(l_var), np.sqrt(k_var)]
11 upper_error = [0.031, 0.277, 0.463, 0.0, 0.864, np.sqrt(dt_var_categorical),
12               np.sqrt(dt_var_oneHot), np.sqrt(l_var), np.sqrt(k_var)]
13
14 error = [lower_error, upper_error]
15
16 fig, ax = plt.subplots(figsize=(10, 6))
17
18 ax.errorbar(means, models, xerr=error, fmt='o', color='steelblue', ecolor='black', capsize=5, marker='o')
19
20 ax.set_xlabel('Accuracy (%)')
21 ax.set_title('Comparison of Model Performance')
22
23 ax.grid(True, linestyle='--', alpha=0.7)
24
25 plt.tight_layout()
26 plt.show()
27

```



In [ ]: 1

In [ ]: 1

In [ ]: 1

In [ ]: 1

In [ ]: 1

