EXPERIMENT NO 8 Random Forest Model on the MNIST data. **IMPORTING LIBRARIES** In [1]: import numpy as np import pandas as pd import matplotlib.pyplot as plt import seaborn as sns from sklearn.model_selection import train_test_split from sklearn.ensemble import RandomForestClassifier from sklearn.metrics import classification report, confusion matrix from sklearn.model_selection import RandomizedSearchCV, GridSearchCV from sklearn.model_selection import cross_val_score READING THE DATASET In [2]: train = pd.read csv('/DOCUMENTS/COLLEGE/CLASSES/EXPERIMENT NO 8/train.csv/train.csv') test = pd.read_csv('/DOCUMENTS/COLLEGE/CLASSES/EXPERIMENT_NO_8/test.csv/test.csv') In [3]: train.head() label pixel0 pixel1 pixel4 pixel5 pixel6 pixel7 pixel8 ... pixel774 pixel775 pixel776 pixel777 pixel778 pixel779 r Out[3]: pixel2 pixel3 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 2 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 3 4 0 4 0 0 0 0 0 0 0 5 rows × 785 columns In [4]: test.head() pixel0 pixel7 pixel8 pixel9 ... pixel774 pixel775 pixel776 pixel777 Out[4]: pixel1 pixel2 pixel3 pixel4 pixel5 pixel6 pixel778 pixel779 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 2 0 0 0 0 0 0 0 0 0 0 0 0 0 0 3 0 4 5 rows × 784 columns In [5]: # Seperate the target and independant variables df x = train.iloc[:, 1:]df y = train.iloc[:, 0] In [6]: # function to print dataset def print_image(row, df): temp = df.iloc[row, :].values temp = temp.reshape(28, 28).astype('uint8') plt.imshow(temp) In [7]: print_image(3, df_x) 5 10 15 20 25 10 15 25 **EXPLORARTORY DATA ANALYSIS (EDA)**

C:\Users\HP\anaconda3\lib\site-packages\seaborn\ decorators.py:36: FutureWarning: Pass the following variable a s a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arg

Ó

25

Ó

uments without an explicit keyword will result in an error or misinterpretation.

In [8]: # Check the frequency of each number df_y.value_counts().sort_index() 4132 Out[8]:

4684

4177

4351

4072

3795

4137

4401

4063 4188

Name: label, dtype: int64

<AxesSubplot:xlabel='label', ylabel='count'>

1 2 3

UNDERSTANDING THE TEST DATA

plt.subplot(10, 10, i + 1)

plt.figure(figsize=(12, 10))

print image(i, test)

for i in range(100):

5 label

Ó

X_train, X_test, y_train, y_test = train_test_split(df_x,

25 Ò

bootstrap=True, class_weight=None, criterion='gini', max_depth=None, max features='auto', max_leaf_nodes=None,

 $n_{jobs=1}$,

verbose=0,

precision recall f1-score support

 0.96
 0.98
 0.97

 0.99
 0.99
 0.99

 0.96
 0.97
 0.96

 0.96
 0.94
 0.95

 0.96
 0.96
 0.96

 0.97
 0.99
 0.98

 0.97
 0.96
 0.97

 0.95
 0.95
 0.95

 0.95
 0.95
 0.95

 0.95
 0.95
 0.95

oob_score=False, random_state=None,

warm_start=False)

min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=1, min_samples_split=2

min_weight_fraction_leaf=0.0,

756 841

899

812

8400

0.96

0]

0.96 0.96 0.96

25 Ó 25 Ò

df_y,

25

test size=0.2, random_state=0)

Ò

sns.countplot(df_y)

warnings.warn(

4000

3000

2000

1000

0

20 0

20 0

20

0

0

0

In [11]:

In [12]:

Out[12]:

In [13]:

In [14]:

In [15]:

In [16]:

Out[16]:

MODEL PREPRATION

Split the dataset

print(X_train.shape) print(X_test.shape)

(33600, 784) (8400, 784)

MODEL CREATION

training random Forest

rf.fit(X_train, y_train)

pred = rf.predict(X_test) print("Classification Report")

print("Confusion Report")

Classification Report

5

8

accuracy

macro avg weighted avg

Confusion Report

Cross validation

rf.fit(df_x, df_y)

0.963904761904762 The accuracy has marginally increased to 96.21%

print(np.mean(score))

Predicting on test data pred = rf.predict(test)

ImageId Label

27996

27997

27998

27999

28000

---X----X----X----

28000 rows × 2 columns

27995

27996

27997

27998

27999

RandomForestClassifier(n jobs=1)

print(classification_report(y_test, pred))

print(confusion_matrix(y_test, pred))

0.97 0.95

[[793 0 0 2 2 3 6 0 7 [0 947 7 1 1 2 1 0 2 [6 1 834 5 3 0 3 4 4

USING THE CROSS VALIDATION

rf = RandomForestClassifier(n_estimators=100)

score = cross_val_score(rf, df_x, df_y)

pred = pd.Series(pred, name="Label")

2

pd.concat([pd.Series(range(1, 28001), name="ImageId"), pred], axis=1)

rf = RandomForestClassifier(n_estimators=100,

2

3

4

5

6

7

8

In [9]:

Out[9]:

In [10]: