22/11/2024, 13:06	Random_Forest_Classifier - Colab
# IMPORTANT: SOME KAGGLE DATA SOURCES ARE PRIVATE # RUN THIS CELL IN GROER TO IMPORT YOUR KAGGLE DATA SOURCES. import kagglehub kagglehub.login()	
# IMPORTANT: RUN THIS CELL IN ORDER TO IMPORT YOUR KAGGLE DATA SOURCES, # THEN FELL FREE TO DELETE THIS CELL. # NOTE: THIS NOTEONE ENVIRONMENT DIFFERS FROM KAGGLE'S PYTHON # ENVIRONMENT SO THERE MAY BE MISSING LIBRARIES USED BY YOUR # NOTEONE. Kagglesvij24_bank_01_path = kagglehub.dataset_download('kagglesvij24/bank-01') kagglesvij24_bank_full_version1_path = kagglehub.dataset_download('kagglesvij24/bank-full-version1') print('Data source import complete.')	
# This Python 3 environment comes with many balged analytics libraries installed # It is defined by the kagale (python Octor image: https://github.com/kagale/docker-python # For example, here's several helpful packages to lead import numny as no # linear algebra import numny as no # linear algebra import pands as no # data processing, CSV file I/O (e.g. pd.read_csv) # Input data files are available in the read-only "/input/" directory # For example, running this (by clicking run or pressing Shift-Enter) will list all files under the input directory # sport os for dirrame, filenames in os.walk('haggle/input'): for filename in filenames: print(os.path.join(dirname, filename)) # You can write up to 2668 to the current directory (/haggle/norking/) that gets preserved as output when you create a version using "Save & Run All" # You can also mitte temporary files to //kaggle/norking/), but they won't be saved outside of the current session	
/kaggle/input/bank-full-version1/bank-full.csv Bank Marketing Campaign- Data Pre-Processing and Model Deployment	
i □ Mage.png i □ Mage.png	
Data Pre-Processing	
Steps of preprocessing of data Import necessary library Read Dataset sanity check of dataStep Exploratory Data Analysis (EDA) Missing Value findings Outliers findings Duplicate Findings Normalization	

https://colab.research.google.com/#fileId=https://sac/Joogle.com/#fileId=https://sac/Joogle.com/#fileId=https://sac/Joogle.com/#fileId=https://sac/Joogle.com/#fileId=https://sac/Joogle.com/#fileId=https://sac/Joogle.com/#fileId=https://sac/Joogle.com/#fileId=https://sac/Joogle.com/#fileId=https://sac/Joogle.com/#fileId=https://sac/Joogle.com/#fileId=https://sac/Joogle.com//sac/

Exploratory Data Analysis Using Pandas for basic statistics, summary, and descriptive analysis.

Create histograms,boxplots,scatter plots, and other visualization to understand data distribution and relationships.
 Identify outliers and anomalies that migth affect analysis.

Importing Necessary Libraries

import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns

Reading Dataset

Bank_data = pd.read_csv("/kaggle/input/bank-full-version1/bank-full.csv") import pandas as pd Bank_data = pd.read_csv('/kaggle/input/bank-full-version1/bank-full.csv')

Index(['sl. no', 'age', 'job', 'marital', 'education', 'default', 'balance',
 'housing', 'loan', 'contact', 'day', 'month', 'duration', 'campaign',
 'pdays', 'previous', 'poutcome', 'y'],
 dtype='object')

0 1 58 management married tertiary no 2143 yes no unknown 5 may 261 1 -1 0 unknown no 1 2 44 technician single secondary no 29 yes no unknown 5 may 151 1 -1 0 unknown no 2 3 33 entrepreneur married secondary no 2 yes yes unknown 5 may 76 1 -1 0 unknown no 3 4 47 blue-collar married unknown no 1506 yes no unknown 5 may 92 1 -1 0 unknown no

Bank_data sl. no age job marital education default balance housing loan contact day month duration campaign pdays previous poutcome y 0 1 58 management married tertiary no 2143 yes no unknown 5 may 261 1 -1 0 unknown no 1 2 44 technician single secondary no 29 yes no unknown 5 may 151 1 -1 0 unknown no 2 3 33 entrepreneur married secondary no 2 yes yes unknown 5 may 76 1 -1 0 unknown no 3 4 47 blue-collar married unknown no 1506 yes no unknown 5 may 92 1 -1 0 unknown no 4 5 33 unknown single unknown no 1 no no unknown 5 may 198 1 -1 0 unknown no

4 5 33 unknown single unknown no 1 no no unknown 5 may 198 1 -1 0 unknown no

45207 45208 71 retired divorced primary no 1729 no no cellular 17 nov 456 2 -1 0 unknown yes **45208** 45209 72 retired married secondary no 5715 no no cellular 17 nov 1127 5 184 3 success yes **45209** 45210 57 blue-collar married secondary no 668 no no telephone 17 nov 508 4 -1 0 unknown no **45210** 45211 37 entrepreneur married secondary no 2971 no no cellular 17 nov 361 2 188 11 other no 45211 rows × 18 columns Bank_data.isna().sum() age 8
job 8
job 8
marital 8
education 8
default 8
balance 8
housing 8
loan 8
contact 8
day 8
month 6
duration 8
duration 9
pdays 8
previous 9
poutcome 8
y tope: int64

Bank_data.shape → (45211, 18) Bank_data.info()

Sanity Check

count 45211.000000 45211.000000 45211.000000 45211.000000 45211.000000 45211.000000 45211.000000 45211.000000
 mean
 22606.000000
 40.936210
 1362.272058
 15.806419
 258.163080
 2.763841
 40.197828
 0.580323

 std
 13051.435847
 10.618762
 3044.765829
 8.322476
 257.527812
 3.098021
 100.128746
 2.303441
 min 1.000000 18.000000 -8019.000000 1.000000 0.000000 1.000000 -1.000000 0.000000 **25**% 11303.500000 33.000000 72.000000 8.000000 103.000000 1.000000 -1.000000 0.000000 **50%** 22606.000000 39.000000 448.000000 16.000000 180.000000 2.000000 -1.000000 0.000000 **75%** 33908.500000 48.000000 1428.000000 21.000000 319.000000 3.000000 -1.000000 0.000000
 max
 45211.000000
 95.000000
 102127.000000
 31.000000
 4918.000000
 63.000000
 871.000000
 275.000000
 for col in Bank_data.select_dtypes(include='object').columns: print (col)
print(Bank_data[col].unique())

 $\overrightarrow{\exists}$ sl. no age balance day duration campaign pdays previous

print(Bank_data[col].unique())

job
 ['management' 'technician' 'entrepreneur' 'blue-collar' 'unknown'
 'retired' 'admin.' 'services' 'self-employed' 'unemployed' 'housemaid'
 'student']
 marital
 ['married' 'single' 'divorced']
 education
 ['tertiary' 'secondary' 'unknown' 'primary']
 default
 ['no' 'yes']
 housing
 ['yes' 'no']
 loan
 ['no' 'yes']
 contact
 ['unknown' 'cellular' 'telephone']
 month
 ['may' 'jun' 'jul' 'aug' 'oct' 'nov' 'dec' 'jan' 'feb' 'mar' 'apr' 'sep']
 poutcome
 ['unknown' 'fallure' 'other' 'success']
 y
 ['no' 'yes'] y ['no' 'yes'] import numpy as np

Identify features with missing values features_na = [feature for feature in Bank_data.columns if Bank_data[feature].isnull().sum() > 0] # Print the percentage of missing values for each feature with missing values
if features_na:
 for feature in features_na:
 print(f'{feature}: {np.round(Bank_data[feature].isnull().mean(), 4) * 100}% missing values')
else:
 print('No missing value found')

→ No missing value found print(column,Bank_data[column].nunique()) sl. no 45211
age 77
job 12
marital 3
education 4
default 2
balance 7168
housing 2
loan 2
contact 3
day 31
month 12
duration 1573
campaign 48
pdays 559
previous 41
poutcome 4
y 2

Explore categorical features

Identify discrete numerical features discrete_features = [feature for feature in numerical_features if len(Bank_data[feature].unique()) < 25] # Print the number of discrete numerical features
print('Discrete variables count: {}'.format(len(discrete_features))) # Display the discrete numerical features
print(discrete_features)

Print the number of continuous numerical features
print('Continuous numerical features count: ()'.format(len(continuous_numerical_features)))

Continuous numerical features count: 8
['sl. no', 'age', 'balance', 'day', 'duration', 'campaign', 'pdays', 'previous'] numerical_features=[feature for feature in Bank_data.columns if ((Bank_data[feature].dtypes != '0') and (feature not in ['y']))]
print('Number of numerical variables:'), len ('numerical_features')
Bank_data[numerical_features].head() The Number of numerical variables:
sl. no age job marital education default balance housing loan contact day month duration campaign pdays previous poutcome 0 1 58 management married tertiary no 2143 yes no-unknown 5 may 261 1 -1 0 unknown 1 2 44 technician single secondary no 29 yes no unknown 5 may 151 1 -1 0 unknown 2 3 33 entrepreneur married secondary no 2 yes yes unknown 5 may 76 1 -1 0 unknown

3 4 47 blue-collar married unknown no 1506 yes no unknown 5 may 92 1 -1 0 unknown 4 5 33 unknown single unknown no 1 no no unknown 5 may 198 1 -1 0 unknown

Identify continuous numerical features continuous_numerical_feature in numerical_features if feature not in discrete_features and feature != 'y']

categorical_features = [feature for feature in Bank_data.columns if (Bank_data[feature].dtypes == 'object') and (feature not in ['deposit'])] categorical_features

Bank_data.drop(['default'], inplace=True) sl. no age job marital education default balance housing loan contact day month duration campaign pdays previous poutcome y 0 1 58 management married tertiary no 2143 yes no unknown 5 may 261 1 -1 0 unknown no 1 2 44 technician single secondary no 29 yes no unknown 5 may 151 1 -1 0 unknown no 2 3 33 entrepreneur married secondary no 2 yes yes unknown 5 may 76 1 -1 0 unknown no 3 4 47 blue-collar married unknown no 1506 yes no unknown 5 may 92 1 -1 0 unknown no 4 5 33 unknown single unknown no 1 no no unknown 5 may 198 1 -1 0 unknown no **45206** 45207 51 technician married tertiary no 825 no no cellular 17 nov 977 3 -1 0 unknown yes 45207 45208 71 retired divorced primary no 1729 no no cellular 17 nov 456 2 -1 0 unknown yes **45208** 45209 72 retired married secondary no 5715 no no cellular 17 nov 1127 5 184 3 success yes 45209 45210 57 blue-collar married secondary no 668 no no telephone 17 nov 508 4 -1 0 unknown no 45210 45211 37 entrepreneur married secondary no 2971 no no cellular 17 nov 361 2 188 11 other no 45211 rows × 18 columns Bank_data = Bank_data.drop_duplicates()
Bank_data

2 3 33 entrepreneur married secondary no 2 yes yes unknown 5 may 76 1 -1 0 unknown no 3 4 47 blue-collar married unknown no 1506 yes no unknown 5 may 92 1 -1 0 unknown no 4 5 33 unknown single unknown no 1 no no unknown 5 may 198 1 -1 0 unknown no 45206 45207 51 technician married tertiary no 825 no no cellular 17 nov 977 3 -1 0 unknown yes 45207 45208 71 retired divorced primary no 1729 no no cellular 17 nov 456 2 -1 0 unknown yes **45208** 45209 72 retired married secondary no 5715 no no cellular 17 nov 1127 5 184 3 success yes **45210** 57 blue-collar married secondary no 668 no no telephone 17 nov 508 4 -1 0 unknown no **45210** 45211 37 entrepreneur married secondary no 2971 no no cellular 17 nov 361 2 188 11 other no age job education educatio

sl. no age job marital education default balance housing loan contact day month duration campaign pdays previous poutcome y 0 1 58 management married tertiary no 2143 yes no unknown 5 may 261 1 -1 0 unknown no 1 2 44 technician single secondary no 29 yes no unknown 5 may 151 1 -1 0 unknown no

Bank_data=pd.read_csv("/kaggle/input/bank-full-version1/bank-full.csv") Descriptive Statistics of the Numerical Column sl. no age balance day duration campaign pdays previous count 45211.000000 45211.000000 45211.000000 45211.000000 45211.000000 45211.000000 45211.000000

Display the correlation matrix
print(correlation_matrix)

Exploratory Data Analysis (EDA)

50% 22606.00000 39.00000 448.00000 16.00000 180.00000 2.00000 -1.00000 0.000000 **75**% 33908.500000 48.000000 1428.000000 21.000000 319.000000 3.000000 -1.000000 0.000000
 max
 45211.000000
 95.000000
 102127.000000
 31.000000
 4918.000000
 63.000000
 871.000000
 275.000000
 # correlation with heatmap to interpret the relation and multicolliniarity # Select numerical columns
numerical_columns = Bank_data.select_dtypes(include='number').columns

 mean
 22606.000000
 40.936210
 1362.272058
 15.806419
 258.163080
 2.763841
 40.197828
 0.580323

 std
 13051.435847
 10.618762
 3044.765829
 8.322476
 257.527812
 3.098021
 100.128746
 2.303441

 min
 1.000000
 18.000000
 -8019.000000
 1.000000
 0.000000
 1.000000
 -1.000000
 0.000000

\$1. no 1.000000 0.014973 0.07559 0.061465 0.013031 0.102884 age 0.014973 1.000000 0.007573 1.000000 0.007573 0.007569 0.061465 0.013031 0.102884 age 0.014973 1.000000 0.007783 -0.009120 0.00468 0.004760 balance 0.075639 0.097783 1.000000 0.004503 0.021560 -0.014578 day -0.061465 -0.093120 0.004503 1.0000000 0.030206 0.162490 duration 0.013031 0.000468 0.021560 -0.030206 0.000000 0.030206 0.1600000 0.004579 campaign -0.162884 0.004760 -0.014578 0.162490 -0.084570 1.0000000 pdays 0.437729 -0.023788 0.004373 0.09344 -0.001565 -0.088628 previous 0.271098 0.001288 0.06674 -0.051710 0.001203 -0.032855 pdays previous \$1. no 0.437729 0.271098 age -0.023758 0.001288 balance 0.003435 0.016674 day -0.093044 -0.051710 duration -0.001565 0.001203 campaign -0.088628 -0.032855 pdays 1.000000 0.454820 previous 0.454820 1.000000 # correlation with heatmap to interpret the relation and multicolliniarity import seaborn as sns import matplotlib.pyplot as plt # Plot the correlation matrix using a heatmap
plt.figure(figsize=(12, 8))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
plt.title('Correlation Matrix')

print(Bank_data.columns) Index(['sl.no', 'age', 'job', 'marital', 'education', 'default', 'balance', 'housing', 'loan', 'contact', 'day', 'month', 'duration', 'campaign', 'pdays', 'previous', 'poutcome', 'y'], dtype='object') import numpy as np lw, up = whisker(Bank_data['duration'])
print(f'Lower whisker: {lw}')
print(f'Upper whisker: {up}')

https://colab.research.google.com/#fileId=https%3A//storage.googleapis.com/kaggle-colab-exported-notebooks/random-forest-classifier-843cf9a241122/auto/storage/goog4_request%26X-Goog-Signature%3D78c6635c496240f26278a1f010fcefbda3381f739dcc523de696ca11fd224da818310a891d1ca4d8b9e4a835a65fba48dedcd61b96eed3d78897f732... 2/3

newdataframe1.rename(columns = {'y-new':'deposited?'}, inplace = True)

DataFrame For Model Training

2/11/2024, 13:06 # Model Training Dataframe newdataframe1	Random_Forest_Classifier - Colab

# Verify the column has been dropped print(newdataframe1.columns) Show hidden output	
* Model Training Dataframe -Final newdataframe1	
1. 0 Age Datase day duration companies of the companies o	
from sklearw.andsl_selection import train_test_split # Prepape features and target * modatafremed['deposited'], axis+1 y = modatafremed['deposited'], axis+1 y = modatafremed['deposited'], axis+1 # Split the data into training and testing sets **Xtrain_Xtext_X, Ytrain_/ ytest_split(x, y, test_slice*0.3, random_state*42) print(f'Training data shape: (%_train_shape), (%_train_shape)) Training data shape: (%_train_shape), (%_train_shape)) Training data shape: (%_train_shape), (%_train_shape)) f=train = %_train_copy() of_train = %_train_copy() of_train = %_train_copy() of_train_feposited'] = %_train	
Salamon Sala	
> Non-deposited?: 88.34 % deposited?: 11.66 % x_train=df_train.drop(['deposited?'],axis=1) y_train=df_train['deposited?']	

y_train 23	
THE PROPERTY OF THE PROPERTY O	
<pre>target_names=['No Deposited', 'Deposited'] from sklearn.metrics import precision_score, recall_score, fi_score, accuracy_score def evaluation_metrics(y_test, y_pre, target_names): # scores print("Accuracy :", accuracy_score(y_test, y_pre)) print("Precision_score(y_test, y_pre, pos_label='yes')) print("Precision_score(y_test, y_pre, pos_label='yes')) print("Fi Score :", fi_score(y_test, y_pre, pos_label='yes')) x_train</pre>	
10747 10748 38 Lank day devention capality plays provided job-blanc-callar	
Signature Sign	
Note	
<pre>model = RandomorestClassifier(random_state=2) model.fxt(_mrin,rrain) #predictions y_pre = model.predict(_xtest) evaluation_metrics(_y_test, y_pre, target_names)</pre>	
The Silvan monits (speet (monitoriticalization process, press, pr	
plt.tere() plt.gend() plt.grid(True) plt.show()	

precision recall f1-score support

No Deposited 8.93 8.97 8.95 11966
Deposited 8.68 8.49 8.57 1598

accuracy 9.81 8.73 8.76 13564
weighted avg 8.81 8.73 8.76 13564
Accuracy: 8.9128575641403716
AUC: 8.9421902942776749

Confusion Matrix for RandomForestClassifier