

```
1 !pip install pandas
2 !pip install numpy
3
```

Requirement already satisfied: pandas in /usr/local/lib/python3.11
 Requirement already satisfied: numpy>=1.23.2 in /usr/local/lib/pyt
 Requirement already satisfied: python-dateutil>=2.8.2 in /usr/loc
 Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/pyt
 Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/pyt
 Requirement already satisfied: six>=1.5 in /usr/local/lib/python3
 Requirement already satisfied: numpy in /usr/local/lib/python3.11

```
1 import pandas as pd
2 import numpy as np
```

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Mounted at /content/drive

```
1 ds_path = f"/content/drive/MyDrive/Intellipaat-sem2/ML Algo/Assignment 1/Mushroom.csv"
```

```
1 df_mushroom = pd.read_csv(ds_path)
```

1. Data Exploration

```
1 # 1. The first few rows of the dataset
2 print("df_mushroom.head()")
3 print(df_mushroom.head())
```

df_mushroom.head()

	class	cap-shape	cap-surface	cap-color	bruises	odor	gill-attachment	gill-spacing	gill-size	gill-color	stalk-surface-below-ring	stalk-color-above-ring	stalk-color-below-ring	veil-type	veil-color	ring-number	ring-type	spore-print-color	population	habitat
0	p	x	s	n	t	p		c	n	k	...	w	w	w	p	o	p	k	s	u
1	e	x	s	y	t	a		c	b	k	...	w	w	w	p	o	p	n	n	g
2	e	b	s	w	t	l		c	b	n	...	w	w	w	p	o	p	n	n	m
3	p	x	y	w	t	p		c	n	n	...	w	w	w	p	o	p	k	s	u
4	e	x	s	g	f	n		w	b	k	...	w	w	w	p	o	e	n	a	g

[5 rows x 23 columns]

Resources X

You are not subscribed. [Learn more](#)

You currently have zero compute units available. Resources offered free of charge are not guaranteed. Purchase more units [here](#).

At your current usage level, this runtime may last up to 2 hours 30 minutes.

[Manage sessions](#)

Want more memory and disk space? [X](#)

[Upgrade to Colab Pro](#)

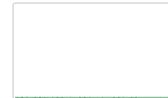
Python 3 Google Compute Engine backend (GPU)

Showing resources from 8:09 AM to 9:25 AM

System RAM
2.1 / 12.7 GB



GPU RAM
0.0 / 15.0 GB



Disk
31.1 / 112.6 GB



```
3 print("Dataset Info:")
4 print(df_mushroom.info())
```

→ Dataset Info:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8124 entries, 0 to 8123
Data columns (total 23 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   class            8124 non-null    object  
 1   cap-shape        8124 non-null    object  
 2   cap-surface      8124 non-null    object  
 3   cap-color        8124 non-null    object  
 4   bruises          8124 non-null    object  
 5   odor             8124 non-null    object  
 6   gill-attachment  8124 non-null    object  
 7   gill-spacing     8124 non-null    object  
 8   gill-size        8124 non-null    object  
 9   gill-color       8124 non-null    object  
 10  stalk-shape     8124 non-null    object  
 11  stalk-root       8124 non-null    object  
 12  stalk-surface-above-ring 8124 non-null    object  
 13  stalk-surface-below-ring 8124 non-null    object  
 14  stalk-color-above-ring 8124 non-null    object  
 15  stalk-color-below-ring 8124 non-null    object  
 16  veil-type       8124 non-null    object  
 17  veil-color      8124 non-null    object  
 18  ring-number     8124 non-null    object  
 19  ring-type       8124 non-null    object  
 20  spore-print-color 8124 non-null    object  
 21  population      8124 non-null    object  
 22  habitat         8124 non-null    object  
dtypes: object(23)
memory usage: 1.4+ MB
None
```

```
1 # 3. Look for missing values
2 print("\nMissing Values in Each Column:")
3 print(df_mushroom.isnull().sum())
4
5 # Get the columns with missing values
6 columns_with_nulls = df_mushroom.columns[df_mushroom.isnull().any()]
7
8 # Print the column names with missing values
9 print("Columns with missing values:")
10 print(columns_with_nulls)
```

→ Missing Values in Each Column:

```
class              0
cap-shape          0
cap-surface        0
cap-color          0
bruises           0
odor              0
gill-attachment   0
gill-spacing      0
gill-size          0
gill-color         0
stalk-shape       0
stalk-root         0
stalk-surface-above-ring 0
stalk-surface-below-ring 0
stalk-color-above-ring 0
stalk-color-below-ring 0
veil-type          0
veil-color         0
ring-number        0
ring-type          0
spore-print-color 0
population         0
habitat            0
dtype: int64
```

```
Columns with missing values:  
Index([], dtype='object')
```

```
1 # 4. unique values in each column  
2 print("\nUnique Values in Each Column:")  
3 for column in df_mushroom.columns:  
4     print(f"{column}: {df_mushroom[column].nunique()} unique values")
```

↳

```
Unique Values in Each Column:  
class: 2 unique values  
cap-shape: 6 unique values  
cap-surface: 4 unique values  
cap-color: 10 unique values  
bruises: 2 unique values  
odor: 9 unique values  
gill-attachment: 2 unique values  
gill-spacing: 2 unique values  
gill-size: 2 unique values  
gill-color: 12 unique values  
stalk-shape: 2 unique values  
stalk-root: 5 unique values  
stalk-surface-above-ring: 4 unique values  
stalk-surface-below-ring: 4 unique values  
stalk-color-above-ring: 9 unique values  
stalk-color-below-ring: 9 unique values  
veil-type: 1 unique values  
veil-color: 4 unique values  
ring-number: 3 unique values  
ring-type: 5 unique values  
spore-print-color: 9 unique values  
population: 6 unique values  
habitat: 7 unique values
```

```
1 # 5. Distribution of the target variable ('class')  
2 print("\nTarget Variable [ class ] Distribution:")  
3 print(df_mushroom['class'].value_counts())
```

↳

```
Target Variable [ class ] Distribution:  
class  
e    4208  
p    3916  
Name: count, dtype: int64
```

```
1 # 6. Basic statistics  
2 print("\nBasic Statistics:")  
3 print(df_mushroom.describe(include='all'))
```

↳

```
Basic Statistics:  
      class cap-shape cap-surface cap-color bruises odor gill-a  
count   8124     8124     8124     8124     8124     8124  
unique    2        6        4       10        2        9  
top      e        x        y        n        f        n  
freq    4208    3656    3244    2284    4748    3528  
  
      gill-spacing gill-size gill-color ... stalk-surface-below  
count      8124     8124     8124    ...  
unique      2        2       12    ...  
top        c        b        b    ...  
freq     6812    5612    1728    ...  
  
      stalk-color-above-ring stalk-color-below-ring veil-type ve  
count          8124           8124           8124  
unique          9             9             1  
top            w             w             p  
freq     4464           4384           8124  
  
      ring-number ring-type spore-print-color population habitat  
count      8124     8124     8124     8124     8124  
unique      3         5         9         6         7
```

top freq	o 7488	p 3968	w 2388	v 4040	d 3148
-------------	-----------	-----------	-----------	-----------	-----------

[4 rows x 23 columns]

```
1 # 7. Check the distribution of values for each categorical column
2 for column in df_mushroom.columns:
3     print(f"\nValue Counts for {column}:")
4     print(df_mushroom[column].value_counts())
```

→ START-COLOR-PRINTING

```
w    4384
p    1872
g    576
n    512
b    432
o    192
e    96
c    36
y    24
Name: count, dtype: int64
```

```
Value Counts for veil-type:
veil-type
p    8124
Name: count, dtype: int64
```

```
Value Counts for veil-color:
veil-color
w    7924
n    96
o    96
y    8
Name: count, dtype: int64
```

```
Value Counts for ring-number:
ring-number
o    7488
t    600
n    36
Name: count, dtype: int64
```

```
Value Counts for ring-type:
ring-type
p    3968
e    2776
l    1296
f    48
n    36
Name: count, dtype: int64
```

```
Value Counts for spore-print-color:
spore-print-color
w    2388
n    1968
k    1872
h    1632
r    72
u    48
o    48
y    48
b    48
Name: count, dtype: int64
```

```
Value Counts for population:
population
v    4040
y    1712
s    1248
```

```
1 print(df_mushroom.shape)
```

→ (8124, 23)

```
1 # Remove data with any missing information for now
2 null_columns = len(columns_with_nulls)
3 print(null_columns)
4 if (null_columns > 0):
5   df_mushroom = df_mushroom.dropna()
```

→ 0

```
1 print(df_mushroom.shape)
```

→ (8124, 23)

```
1 # Get summary stats for the categorical features
2 df_mushroom.describe(include = ['O'])
```

→

	class	cap-shape	cap-surface	cap-color	bruises	odor	gill-attachment	gill-spacing	gill-width	gill-thickness	stalk-shape	stalk-root	stalk-surface-above-ring	stalk-surface-below-ring	stalk-color-above-ring	stalk-color-below-ring	veil-color	ring-number	ring-type	spore-print-color	population	habitat
count	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	
unique	2	6	4	10	2	9	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	
top	e	x	y	n	f	n	f	n	n	n	n	n	n	n	n	n	n	n	n	n	n	
freq	4208	3656	3244	2284	4748	3528	7914	7914	7914	7914	7914	7914	7914	7914	7914	7914	7914	7914	7914	7914	7914	

4 rows × 23 columns

```
1 df_mushroom.describe()
```

→

	class	cap-shape	cap-surface	cap-color	bruises	odor	gill-attachment	gill-spacing	gill-width	gill-thickness	stalk-shape	stalk-root	stalk-surface-above-ring	stalk-surface-below-ring	stalk-color-above-ring	stalk-color-below-ring	veil-color	ring-number	ring-type	spore-print-color	population	habitat
count	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	8124	
unique	2	6	4	10	2	9	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	
top	e	x	y	n	f	n	f	n	n	n	n	n	n	n	n	n	n	n	n	n	n	
freq	4208	3656	3244	2284	4748	3528	7914	7914	7914	7914	7914	7914	7914	7914	7914	7914	7914	7914	7914	7914	7914	

4 rows × 23 columns

```
1 # Create the feature/flag for eatable / non eatable
2 df_mushroom['class'] = df_mushroom['class'].apply(lambda x: 1 if x
```

```
1 df_mushroom.head()
```



	class	cap-shape	cap-surface	cap-color	bruises	odor	gill-attachment	gill-spacer
0	0	x	s	n	t	p	f	
1	1	x	s	y	t	a	f	
2	1	b	s	w	t	l	f	
3	0	x	y	w	t	p	f	
4	1	x	s	g	f	n	f	

5 rows × 23 columns

```
1 # Check the event rate
2 df_mushroom['class'].value_counts()/len(df_mushroom)
```



count

class	count
1	0.517971
0	0.482029

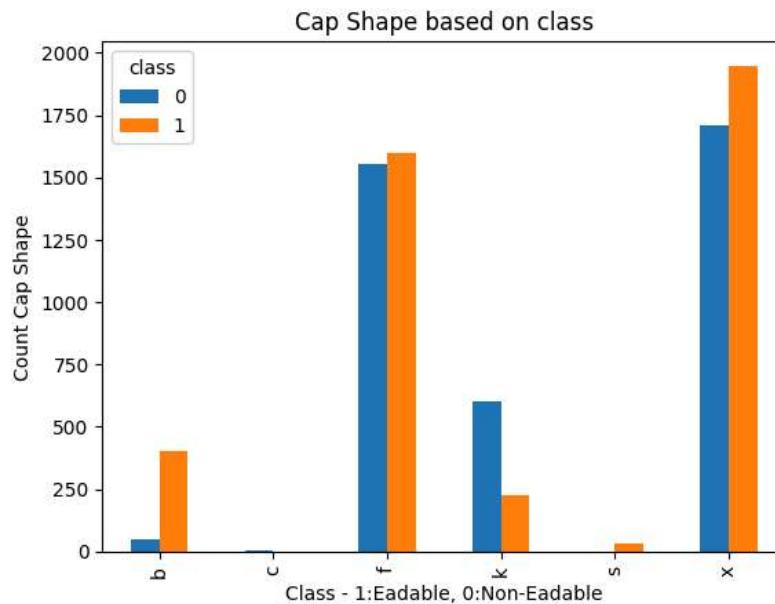
dtype: float64

```
1 !pip install matplotlib
2 !pip install seaborn
```

```
→ Requirement already satisfied: matplotlib in /usr/local/lib/python
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/pyt
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/li
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/li
Requirement already satisfied: numpy>=1.23 in /usr/local/lib/pyth
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/
Requirement already satisfied: pillow>=8 in /usr/local/lib/python
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib
Requirement already satisfied: python-dateutil>=2.7 in /usr/local
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3
Requirement already satisfied: seaborn in /usr/local/lib/python3.
Requirement already satisfied: numpy!=1.24.0,>=1.20 in /usr/local
Requirement already satisfied: pandas>=1.2 in /usr/local/lib/pyth
Requirement already satisfied: matplotlib!=3.6.1,>=3.4 in /usr/lo
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/pyt
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/li
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/li
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/
Requirement already satisfied: pillow>=8 in /usr/local/lib/python
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib
Requirement already satisfied: python-dateutil>=2.7 in /usr/local
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/pyt
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/p
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3
```

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 pd.crosstab(df_mushroom['cap-shape'],df_mushroom['class']).plot(kind
5 plt.title('Cap Shape based on class')
6 plt.xlabel('Class - 1:Eadable, 0:Non-Eadable')
7 plt.ylabel('Count Cap Shape')
```

→ Text(0, 0.5, 'Count Cap Shape')

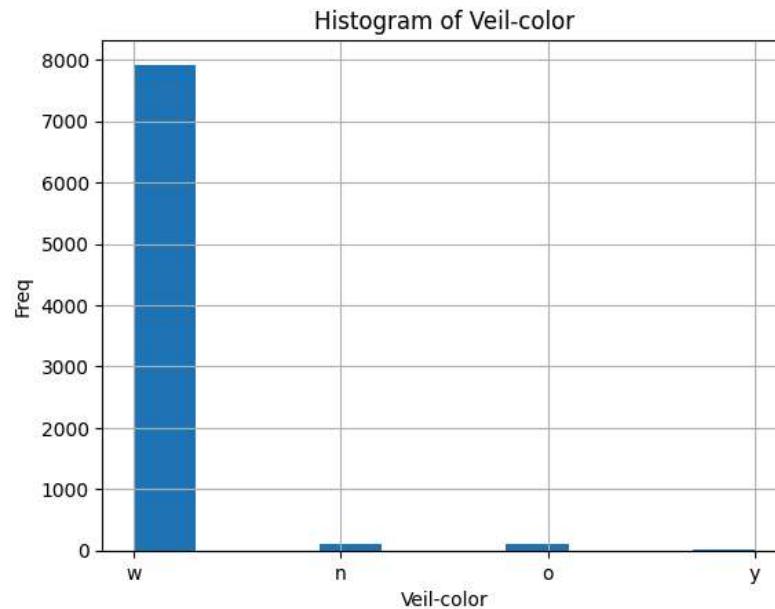


```

1 #veil-color: brown=n,orange=o,white=w,yellow=y
2
3 df_mushroom['veil-color'].hist()
4 plt.title('Histogram of Veil-color')
5 plt.xlabel('Veil-color')
6 plt.ylabel('Freq')
7

```

→ Text(0, 0.5, 'Freq')



```

1 # Convert the categorical features
2 # Creating dummies for all these variables
3 d_cap_shape = pd.get_dummies(df_mushroom['cap-shape'], prefix='cap'
4 d_cap_surface = pd.get_dummies(df_mushroom['cap-surface'], prefix=
5 d_cap_color = pd.get_dummies(df_mushroom['cap-color'], prefix='cap'
6 d_bruises = pd.get_dummies(df_mushroom['bruises'], prefix='bruises'
7 d_odor = pd.get_dummies(df_mushroom['odor'], prefix='odor', dtype=i
8 d_gill_attachment = pd.get_dummies(df_mushroom['gill-attachment'],

```

```

9 d_gill_spacing = pd.get_dummies(df_mushroom['gill-spacing'], prefix='gill_')
10 d_gill_size = pd.get_dummies(df_mushroom['gill-size'], prefix='gill_')
11 d_gill_color = pd.get_dummies(df_mushroom['gill-color'], prefix='gill_')
12 d_stalk_shape = pd.get_dummies(df_mushroom['stalk-shape'], prefix='stalk_')
13 d_stalk_root = pd.get_dummies(df_mushroom['stalk-root'], prefix='stalk_')
14 d_stalk_surface_above_ring = pd.get_dummies(df_mushroom['stalk-surface-above-ring'], prefix='stalk_')
15 d_stalk_surface_below_ring = pd.get_dummies(df_mushroom['stalk-surface-below-ring'], prefix='stalk_')
16 d_stalk_color_above_ring = pd.get_dummies(df_mushroom['stalk-color-above-ring'], prefix='stalk_')
17 d_stalk_color_below_ring = pd.get_dummies(df_mushroom['stalk-color-below-ring'], prefix='stalk_')
18 d_veil_type = pd.get_dummies(df_mushroom['veil-type'], prefix='veil_')
19 d_veil_color = pd.get_dummies(df_mushroom['veil-color'], prefix='veil_')
20 d_ring_number = pd.get_dummies(df_mushroom['ring-number'], prefix='ring_')
21 d_ring_type = pd.get_dummies(df_mushroom['ring-type'], prefix='ring_')
22 d_spore_print_color = pd.get_dummies(df_mushroom['spore-print-color'], prefix='spore_')
23 d_population = pd.get_dummies(df_mushroom['population'], prefix='pop_')
24 d_habitat = pd.get_dummies(df_mushroom['habitat'], prefix='habitat_')
25
26

```

```

1 # Create the final dataset with all the relevant features - both discrete and continuous
2 feature_x_cont = ['class','cap-shape','cap-surface','cap-color','bruises','odor','gill-spacing','gill-size','gill-color','stalk-shape','stalk-root','stalk-surface-above-ring','stalk-surface-below-ring','stalk-color-above-ring','stalk-color-below-ring','veil-type','veil-color','ring-number','ring-type','spore-print-color','population','habitat']
3
4 df_mushroom_cont = df_mushroom[feature_x_cont]
5
6 # Creating the Final data with all the relevant fields and Dep Var
7 df_mushroom_new = pd.concat([df_mushroom['class'],d_cap_shape,d_ca

```

```
1 df_mushroom_new.shape
```

→ (8124, 118)

```
1 df_mushroom_new.head()
```

	class	cap-shape_b	cap-shape_c	cap-shape_f	cap-shape_k	cap-shape_s	cap-shape_x	sl
0	0	0	0	0	0	0	0	1
1	1	0	0	0	0	0	0	1
2	1	1	0	0	0	0	0	0
3	0	0	0	0	0	0	0	1
4	1	0	0	0	0	0	0	1

5 rows × 118 columns

```
1 df_mushroom_new.info()
```

→ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 8124 entries, 0 to 8123
Columns: 118 entries, class to habitat_w
dtypes: int64(118)
memory usage: 7.3 MB

▼ Train Test Split

Double-click (or enter) to edit

```

1 #X = df_mushroom_new.drop('class', axis=1)
2 #y = df_mushroom_new['class']
3

```

```

4 X = df_mushroom.drop('class', axis=1)
5 y = df_mushroom['class']
6
7 X = X.map(lambda x: ord(x))

1 from sklearn.model_selection import train_test_split

1 X_train,X_test,y_train,y_test=train_test_split(X,y,train_size=0.75

```

✓ Model Building

```

1 import sklearn.metrics as metrics
2 from sklearn.linear_model import LogisticRegression

1 from sklearn.datasets import make_classification
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.model_selection import GridSearchCV, train_test_split
4 from sklearn.metrics import accuracy_score

1 # Step 3: Define the model
2 model = RandomForestClassifier(random_state=42)
3
4 # Step 4: Define the hyperparameter grid
5 param_grid = {
6     'n_estimators': [50, 100, 150],           # Number of trees in the forest
7     'max_depth': [None, 10, 20],             # Maximum depth of the trees
8     'min_samples_split': [2, 5, 10],        # Minimum number of samples required to split an internal node
9     'min_samples_leaf': [1, 2, 4]           # Minimum number of samples required at each leaf node
10 }
11
12 # Step 5: Set up GridSearchCV
13 grid_search = GridSearchCV(
14     estimator=model,
15     param_grid=param_grid,
16     cv=5,                      # 5-fold cross-validation
17     scoring='accuracy',       # Performance metric
18     verbose=1,                 # Print progress
19     n_jobs=-1                  # Use all available CPU cores
20 )
21
22 # Step 6: Perform the grid search
23 grid_search.fit(X_train, y_train)
24
25 # Step 7: Print the best hyperparameters and score
26 print("Best Parameters:", grid_search.best_params_)
27 print("Best Cross-Validation Accuracy:", grid_search.best_score_)
28
29 # Step 8: Evaluate the model on the test set
30 best_model = grid_search.best_estimator_
31 y_pred = best_model.predict(X_test)
32 test_accuracy = accuracy_score(y_test, y_pred)
33 print("Test Set Accuracy:", test_accuracy)

```

→ Fitting 5 folds for each of 81 candidates, totalling 405 fits
 Best Parameters: {'max_depth': None, 'min_samples_leaf': 1, 'min_...
 Best Cross-Validation Accuracy: 1.0
 Test Set Accuracy: 1.0

```
1 best_model.fit(X_train, y_train)
```

```
RandomForestClassifier
RandomForestClassifier(n_estimators=50, random_state=42)

1 best_model.score(X_train,y_train)
→ 1.0
```

Model Prediction

```
1 y_train_pred=best_model.predict(X_train)
```

```
1 y_test_pred=best_model.predict(X_test)
```

Model Evaluation

```
1 from sklearn.metrics import confusion_matrix,accuracy_score,recall_
```

```
1 print('Train Accuracy:',accuracy_score(y_train,y_train_pred))
2 print('Test Accuracy:',accuracy_score(y_test,y_test_pred))
3 print('Train Recall:',recall_score(y_train,y_train_pred))
4 print('Test Recall:',recall_score(y_test,y_test_pred))
5 print('Train Precision:',precision_score(y_train,y_train_pred))
6 print('Test Precision:',precision_score(y_test,y_test_pred))
7 print('Train F1 Score:',f1_score(y_train,y_train_pred))
8 print('Test F1 Score:',f1_score(y_test,y_test_pred))
```

```
→ Train Accuracy: 1.0
  Test Accuracy: 1.0
  Train Recall: 1.0
  Test Recall: 1.0
  Train Precision: 1.0
  Test Precision: 1.0
  Train F1 Score: 1.0
  Test F1 Score: 1.0
```

```
1 print(classification_report(y_train,y_train_pred))
2 print("xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx")
3 print(classification_report(y_test,y_test_pred))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	2925
1	1.00	1.00	1.00	3168
accuracy			1.00	6093
macro avg	1.00	1.00	1.00	6093
weighted avg	1.00	1.00	1.00	6093
xx				
	precision	recall	f1-score	support
0	1.00	1.00	1.00	991
1	1.00	1.00	1.00	1040
accuracy			1.00	2031
macro avg	1.00	1.00	1.00	2031
weighted avg	1.00	1.00	1.00	2031

```
1 from sklearn.tree import DecisionTreeClassifier
```

```
1 DT_gcv=DecisionTreeClassifier()
```

```
1 param_grid = {'max_features': ['sqrt', 'log2'],
2                 'ccp_alpha': [1, 0.1, .01, .001, 0.0001, 0.00001],
3                 'max_depth' : [4,5,6,7,8,9,10,11,12,15],
4                 'criterion' :['gini', 'entropy', 'log_loss'],
5                 'splitter' :['best', 'random'],
6                 'min_samples_split': [2, 5, 10],
7                 'min_samples_leaf': [1, 2, 4]
8             }
```

```
1 model_cv=GridSearchCV(DT_gcv,param_grid,scoring="f1",verbose=2,cv=
```

```
1 model_cv.fit(X_train,y_train)
```

→ Fitting 5 folds for each of 6480 candidates, totalling 32400 fits

```
GridSearchCV
  best_estimator_:
    DecisionTreeClassifier
      DecisionTreeClassifier
```

```
1 best_model=model_cv.best_estimator_
2 best_param=model_cv.best_params_
```

```
1 print(best_param)
```

→ { 'ccp_alpha': 0.001, 'criterion': 'entropy', 'max_depth': 10, 'ma

```
1 best_model.fit(X_train,y_train)
```

```
1 y_train_pred_cv=best_model.predict(X_train)
```

```
1 y_test_pred_cv=best_model.predict(X_test)
```

```
1 print(classification_report(y_train,y_train_pred_cv))
2 print("XXXXXXXXXXXXXXXXXXXXXXXXXXXXXX")
3 print(classification_report(y_test,y_test_pred_cv))
```

→

	precision	recall	f1-score	support
0	1.00	1.00	1.00	2925
1	1.00	1.00	1.00	3168
accuracy			1.00	6093
macro avg	1.00	1.00	1.00	6093
weighted avg	1.00	1.00	1.00	6093

XXXXXXXXXXXXXXXXXXXXXXXXXXXXXX

	precision	recall	f1-score	support
0	1.00	1.00	1.00	991
1	1.00	1.00	1.00	1040
accuracy			1.00	2031
macro avg	1.00	1.00	1.00	2031
weighted avg	1.00	1.00	1.00	2031

Linear Regression is a basic supervised machine learning algorithm where the relationship between the dependent variable and one or more independent variables is modeled by a straight line (linear equation).

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.model_selection import GridSearchCV
3
4 # Define the model
5 model = LogisticRegression( ) # Increase max_iter as needed
6
7
```

```
1 model.fit(X_train,y_train)

→ /usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_log
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as
<https://scikit-learn.org/stable/modules/preprocessing.html>
Please also refer to the documentation for alternative solver opt
https://scikit-learn.org/stable/modules/linear_model.html#log

```
n_iter_i = _check_optimize_result(
    LogisticRegression( ⓘ ⓘ)
LogisticRegression()
```

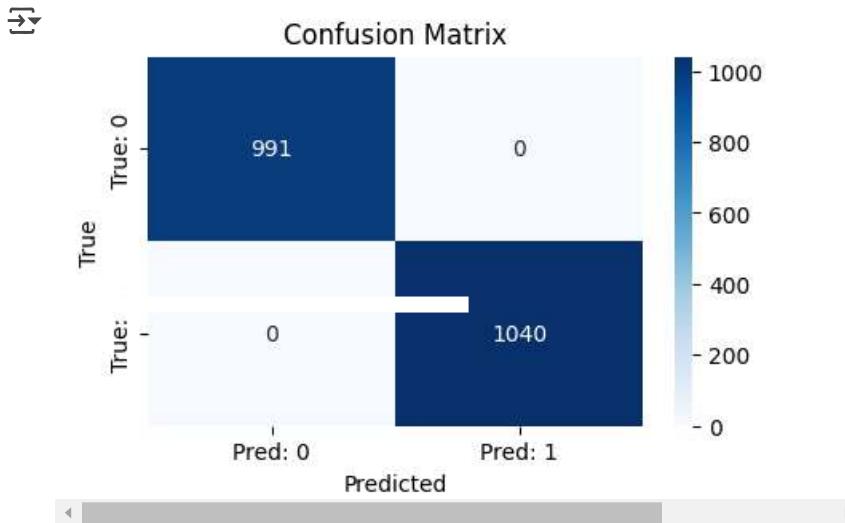
```
1 y_train_pred = model.predict(X_test)
2 y_test_pred = model.predict(X_test)
3
4

1 print('Train Accuracy:',accuracy_score(y_train,y_train_pred))
2 print('Test Accuracy:',accuracy_score(y_test,y_test_pred))
3 print('Train Recall:',recall_score(y_train,y_train_pred))
4 print('Test Recall:',recall_score(y_test,y_test_pred))
5 print('Train Precision:',precision_score(y_train,y_train_pred))
6 print('Test Precision:',precision_score(y_test,y_test_pred))
7 print('Train F1 Score:',f1_score(y_train,y_train_pred))
8 print('Test F1 Score:',f1_score(y_test,y_test_pred))
9 print('Train roc_auc_score:',roc_auc_score(y_train,y_train_pred))
10 print('Test roc_auc_score:',roc_auc_score(y_test,y_test_pred))
```

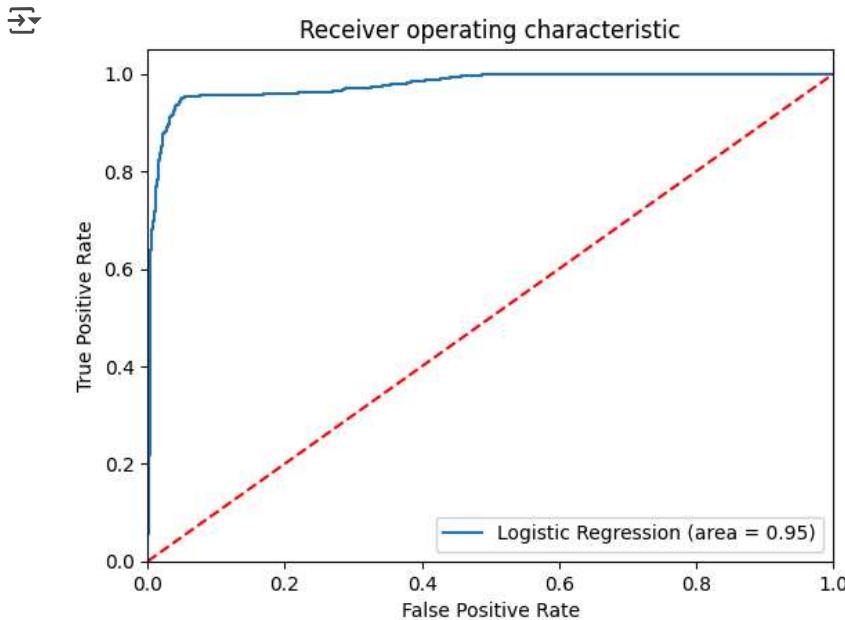
```
→ Train Accuracy: 1.0
Test Accuracy: 1.0
Train Recall: 1.0
Test Recall: 1.0
Train Precision: 1.0
Test Precision: 1.0
Train F1 Score: 1.0
Test F1 Score: 1.0
Train roc_auc_score: 1.0
Test roc_auc_score: 1.0
```

```
1 from sklearn.metrics import confusion_matrix
2
3 # Generate confusion matrix
4 cm = confusion_matrix(y_test,y_test_pred)
5
6 # Plot confusion matrix as a heatmap
7 plt.figure(figsize=(5, 3))
8 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['P'
9 plt.title('Confusion Matrix')
10 plt.xlabel('Predicted')
```

```
11 plt.ylabel('True')
12 plt.show()
```



```
1 # a graph that shows how well a model performs across all possible
2 # against the false positive rate (FPR)
3 from sklearn.metrics import roc_curve
4 logit_roc_auc = roc_auc_score(y_test, model.predict(X_test))
5 fpr, tpr, thresholds = roc_curve(y_test, model.predict_proba(X_te
6 plt.figure()
7 plt.plot(fpr, tpr, label='Logistic Regression (area = %0.2f)' % lo
8 plt.plot([0, 1], [0, 1], 'r--')
9 plt.xlim([0.0, 1.0])
10 plt.ylim([0.0, 1.05])
11 plt.xlabel('False Positive Rate')
12 plt.ylabel('True Positive Rate')
13 plt.title('Receiver operating characteristic')
14 plt.legend(loc="lower right")
15 plt.savefig('Log_ROC')
16 plt.show()
```



✓ Hyperparameter tuning with GridSearchCV

```

1 from sklearn.metrics import classification_report# Define the para
2 param_grid = {
3     'C': [0.01, 0.1, 1, 10, 100], # Inverse of regularization str
4     'solver': ['lbfgs', 'liblinear', 'saga'], # Solvers
5     'l1_ratio': [0, 0.5, 1], # Used only with elasticnet
6 }
7
8 # Initialize GridSearchCV
9 grid_search = GridSearchCV(estimator=model,
10                           param_grid=param_grid,
11                           cv=5,
12                           scoring='accuracy',
13                           verbose=1,
14                           n_jobs=-1)
15
16 # Fit the model
17 grid_search.fit(X_train, y_train)
18
19 # Best parameters and model
20 print("Best Parameters:", grid_search.best_params_)
21 best_model = grid_search.best_estimator_
22
23 # Evaluate the model
24 y_pred = best_model.predict(X_test)
25 print("\nClassification Report:\n", classification_report(y_test, :

```

→ Fitting 5 folds for each of 45 candidates, totalling 225 fits
 Best Parameters: {'C': 10, 'l1_ratio': 0, 'solver': 'liblinear'}

	precision	recall	f1-score	support
0	1.00	1.00	1.00	991
1	1.00	1.00	1.00	1040
accuracy			1.00	2031
macro avg	1.00	1.00	1.00	2031
weighted avg	1.00	1.00	1.00	2031

/usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_log
 warnings.warn(

```
1 best_model.fit(X_train,y_train)
```

→ /usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_log
 warnings.warn(
 ▾ LogisticRegression
 LogisticRegression(C=10, l1_ratio=0, solver='liblinear')

```
1 y_test_tred = model.predict(X_test)
```

✓ SVM

```

1 from sklearn.svm import SVC
2 svm = SVC(kernel='rbf',random_state=0,probability=True)
3 svm.fit(X_train,y_train)
```

→ ▾ SVC
 SVC(probability=True, random_state=0)

```

1 y_test_pred = svm.predict(X_test)

1 y_train_pred = svm.predict(X_train)

1 print('Train Accuracy:',accuracy_score(y_train,y_train_pred))
2 print('Test Accuracy:',accuracy_score(y_test,y_test_pred))
3 print('Train Recall:',recall_score(y_train,y_train_pred))
4 print('Test Recall:',recall_score(y_test,y_test_pred))
5 print('Train Precision:',precision_score(y_train,y_train_pred))
6 print('Test Precision:',precision_score(y_test,y_test_pred))
7 print('Train F1 Score:',f1_score(y_train,y_train_pred))
8 print('Test F1 Score:',f1_score(y_test,y_test_pred))
9 print('Train roc_auc_score:',roc_auc_score(y_train,y_train_pred))
10 print('Test roc_auc_score:',roc_auc_score(y_test,y_test_pred))

```

```

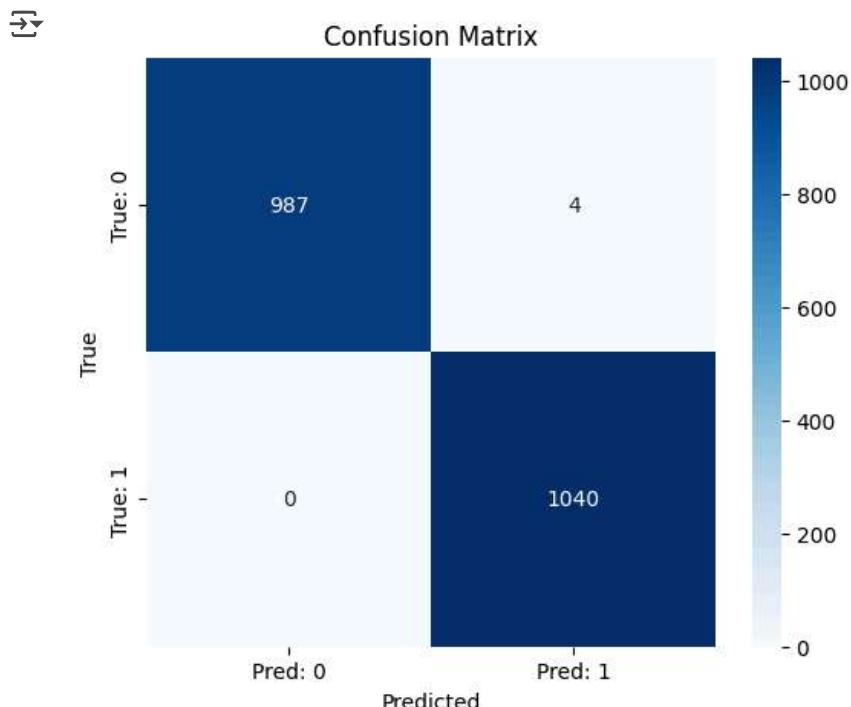
[4]: Train Accuracy: 0.995240439849007
      Test Accuracy: 0.9980305268340719
      Train Recall: 0.9996843434343434
      Test Recall: 1.0
      Train Precision: 0.9912363067292644
      Test Precision: 0.9961685823754789
      Train F1 Score: 0.9954424013829954
      Test F1 Score: 0.9980806142034548
      Train roc_auc_score: 0.9950558469308469
      Test roc_auc_score: 0.9979818365287588

```

```

1 from sklearn.metrics import confusion_matrix
2
3 # Generate confusion matrix
4 cm = confusion_matrix(y_test, y_test_pred)
5
6 # Plot confusion matrix as a heatmap
7 plt.figure(figsize=(6, 5))
8 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['Pred: 0', 'Pred: 1'],
9 plt.title('Confusion Matrix')
10 plt.xlabel('Predicted')
11 plt.ylabel('True')
12 plt.show()

```



```

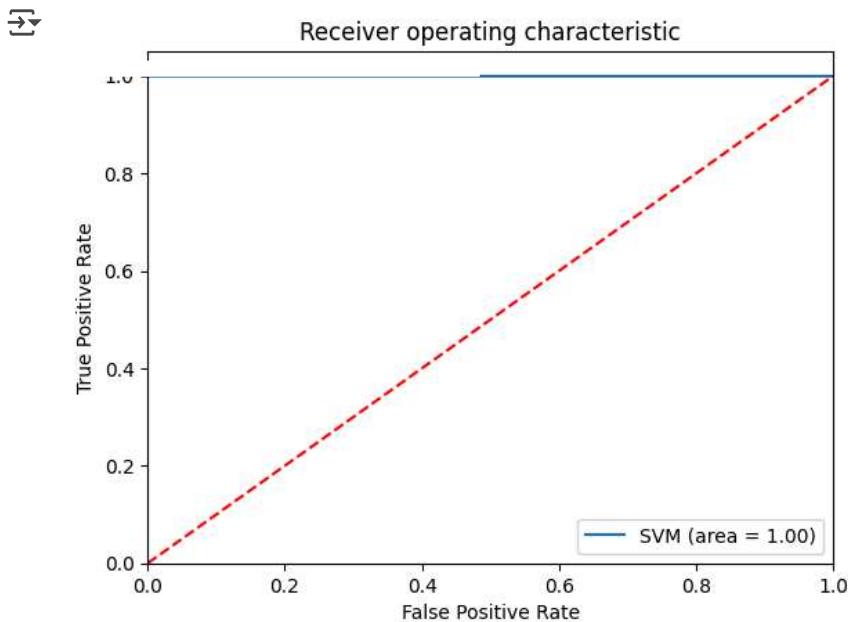
1 # a graph that shows how well a model performs across all possible
2 # against the false positive rate (FPR)

```

```

3 from sklearn.metrics import roc_curve
4 logit_roc_auc = roc_auc_score(y_test, svm.predict(X_test))
5 fpr, tpr, thresholds = roc_curve(y_test, svm.predict_proba(X_test))
6 plt.figure()
7 plt.plot(fpr, tpr, label='SVM (area = %0.2f)' % logit_roc_auc)
8 plt.plot([0, 1], [0, 1], 'r--')
9 plt.xlim([0.0, 1.0])
10 plt.ylim([0.0, 1.05])
11 plt.xlabel('False Positive Rate')
12 plt.ylabel('True Positive Rate')
13 plt.title('Receiver operating characteristic')
14 plt.legend(loc="lower right")
15 plt.savefig('Log_ROC')
16 plt.show()

```



```

1 # Define the parameter grid for GridSearchCV
2 param_grid = {
3     'C': [0.1, 1, 10, 100],  # Regularization parameter
4     'kernel': ['linear', 'poly', 'rbf', 'sigmoid'],  # Kernel types
5 }
6
7 # Initialize GridSearchCV
8 grid_search_svm = GridSearchCV(estimator=svm,
9                                 param_grid=param_grid,
10                                cv=5,
11                                scoring='accuracy',
12                                verbose=1,
13                                n_jobs=-1)
14
15 # Fit the model
16 grid_search_svm.fit(X_train, y_train)
17
18 # Best parameters and model
19 print("Best Parameters:", grid_search_svm.best_params_)
20 best_model_svm = grid_search_svm.best_estimator_
21
22 # Evaluate the model
23 y_test_pred = best_model_svm.predict(X_test)
24 print("\nClassification Report:\n", classification_report(y_test, y_
25
...

```

... Fitting 5 folds for each of 16 candidates, totalling 80 fits

*Conclusion : Three Models are tried *

1. RandomForestClassifier
2. DecisionTreeClassifier
3. LogisticRegression
4. SVM

[Change runtime type](#)