

Comment Share

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

Q You are exploring the wilderness of Mushroom, a land populated by a plethora of diverse fauna and flora. In particular, Mushroomia is known for its unparalleled variety in mushrooms. However, not all the mushrooms in Mushroomia are edible. As you make your way through Mushroomia, you would like to know which mushrooms are edible, in order to forage for supplies for your daily mushroom soup.

You have access to:

- Shroomster Pro Max™ - a state of the art data collection device, developed by Mushroomia, that allows you to collect various data points about any mushroom you encounter in the wild
- The National Archives on Mushrooms - a dataset collected over the years by the government of Mushroomia

To address this problem, you decide to use the skills you learnt in CSM148 and train machine learning models on the The National Archives on Mushrooms in order to use your Shroomster Pro Max™ to determine whether the mushrooms you encounter on your adventure can be added to your daily mushroom soup.

This project will be more unstructured than the previous two projects in order to allow you to experience how data science problems are solved in practice. There are two parts to this project: a Jupyter Notebook with your code (where you explore, visualize, process your data and train machine learning models) and a report (where you explain the various choices you make in your implementation and analyze the final performance of your model).

1. Loading and Viewing Data

```
[1] import numpy as np
import seaborn as sns # used for plot interactive graph.

[2] # Delete this cell if not using google colab
from google.colab import files
files.upload()

[3] # Only keep features that shroomster can detect
# remove class feature, because shroomster cannot detect
# labels + concat['class']
labels = pd.read_csv('mushroom_train.csv', delimiter=';')
test = pd.read_csv('mushroom_test.csv', delimiter=';')

[4] concert = pd.concat([train, test])
concert
```

	cap-diameter	cap-shape	cap-surface	cap-color	does-bruise-or-bled	gill-attachment	gill-spacing	gill-color	stem-height	stem-width	stem-root	stem-surface	stem-color	veil-type	veil-color	has-ring	ring-type	spore-print-color	habitat	season	
0	p	15.26	x	g	o	r	e	Nan	w	16.95	—	s	y	w	u	w	t	g	Nan	d	w
1	p	16.60	x	g	o	r	e	Nan	w	17.99	—	s	y	w	u	w	t	g	Nan	d	w
2	p	14.07	x	g	o	r	e	Nan	w	17.00	—	s	y	w	u	w	t	g	Nan	d	w
3	p	14.17	f	h	e	r	e	Nan	w	15.77	—	s	y	w	u	w	t	p	Nan	d	w
4	p	14.64	x	h	o	r	e	Nan	w	16.53	—	s	y	w	u	w	t	p	Nan	d	w
...
10851	e	52.41	o	y	y	r	p	Nan	y	5.47	—	Nan	k	k	Nan	Nan	f	f	Nan	d	u
10852	e	54.81	o	y	y	r	p	Nan	y	6.67	—	Nan	k	k	Nan	Nan	f	f	Nan	d	s
10853	e	49.95	o	y	y	r	p	Nan	y	6.43	—	Nan	k	n	Nan	Nan	f	f	Nan	d	u
10854	e	53.16	o	y	y	r	p	Nan	y	6.99	—	Nan	k	k	Nan	Nan	f	f	Nan	d	s
10855	e	49.78	o	y	y	r	p	Nan	y	5.77	—	Nan	k	k	Nan	Nan	f	f	Nan	d	u

61069 rows × 21 columns

2. Splitting Data into Features and Labels

```
[5] # Only keep features that shroomster can detect
# remove class feature, because shroomster cannot detect
labels = pd.concat(['class'])
labels
```

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
0	0	0	0	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	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

61069 rows × 21 columns

3. Data Exploration and Visualization (Combined 3/4/5)

```
[7] # Hint: we have done this in both project 1 & 2
import matplotlib.pyplot as plt
from matplotlib import style
style.use('fivethirtyeight')
# Import matplotlib.pyplot as plt plotting package
```

```
[8] from sklearn.preprocessing import OrdinalEncoder
```

```
[9] # convert objects to numerical types
concert['cap-diameter'] = concert['cap-diameter'].apply(np.float64)
concert['cap-shape'] = concert['cap-shape'].apply(np.float64)
concert['stem-width'] = concert['stem-width'].apply(np.float64)

# split datasets
train_nf = concat[labels]
test_nf = concat[labels].iloc[50000:]

# convert objects to categorical
ordinal_encoder = OrdinalEncoder()
concat[['cap-shape', 'cap-surface', 'cap-color', 'does-bruise-or-bled', 'gill-attachment', 'gill-spacing', 'stem-color', 'veil-type', 'veil-color', 'has-ring', 'ring-type', 'spore-print-color', 'habitat', 'season']] = ordinal_encoder.fit_transform(concat[['cap-shape', 'cap-surface', 'cap-color', 'does-bruise-or-bled', 'gill-attachment', 'gill-spacing', 'stem-color', 'veil-type', 'veil-color', 'has-ring', 'ring-type', 'spore-print-color', 'habitat', 'season']])

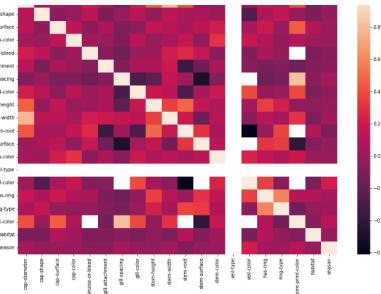
# train display
train_display = concat[labels].iloc[:50000]
```

61069 rows × 21 columns

```
[10] train_display.hist(bins=10, figsize=(20,15))
plt.show()
```

```
[11] fig, ax = plt.subplots(figsize=(15,10))
sns.heatmap(train_display.corr(), ax=ax)
```

```
[12] cax = plt.colorbar()
```



5. Data Augmentation (Creating at least 2 New Features)

```
[14] # Similar to Project 1 and 2.

every similar features: stem-width, cap-diameter; spore-print-color, gill-spacing;
remove_features = ['veil-type', 'veil-color', 'spore-print-color', 'cap-diameter', 'gill-spacing']
#remove_features = []
#train_df.drop(remove_features, axis=1, inplace=True)

#train_df['cap-diam-div-stem-width'] = train_df['cap-diameter']/(train_df['stem-width']+1)
#train_df['cap-diam-mul-stem-width'] = train_df['cap-diameter']*(train_df['stem-width']+1)
#train_df['cap-diam-div-item-height'] = train_df['cap-diameter']/(train_df['stem-height']+1)
#train_df['cap-diam-mul-item-height'] = train_df['cap-diameter']*(train_df['stem-height']+1)
#train_df['item-width-div-item-height'] = train_df['item-width']/(train_df['item-height']+1)
#train_df['item-width-mul-item-height'] = train_df['item-width']*(train_df['item-height']+1)
#train_df['item-height-div-item-height'] = train_df['item-height']/(train_df['item-width']+1)
#train_df['item-height-mul-item-width'] = train_df['item-height']*(train_df['item-width']+1)

#python:Input:17-94efc97250011: SettingWithCopyWarning
#A value is trying to be set on a copy of a slice from a DataFrame.
#Try using .loc[row\_indexer,col\_indexer]: value instead.

see the caeavts in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
train_df['cap-diam-div-stem-width'] = train_df['cap-diameter']/(train_df['stem-width']+1)
#train_df['cap-diam-mul-item-height'] = train_df['cap-diameter']*(train_df['item-height']+1)
#train_df['item-width-div-item-height'] = train_df['item-width']/(train_df['item-height']+1)
#train_df['item-width-mul-item-height'] = train_df['item-width']*(train_df['item-height']+1)

#python:Input:17-94efc97250011: SettingWithCopyWarning
#A value is trying to be set on a copy of a slice from a DataFrame.
#Try using .loc[row\_indexer,col\_indexer]: value instead.

see the caeavts in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
train_df['cap-diam-mul-item-height'] = train_df['cap-diameter']*(train_df['item-height']+1)
#train_df['item-width-div-item-height'] = train_df['item-width']/(train_df['item-height']+1)
#train_df['item-width-mul-item-height'] = train_df['item-width']*(train_df['item-height']+1)

#python:Input:17-94efc97250011: SettingWithCopyWarning
#A value is trying to be set on a copy of a slice from a DataFrame.
#Try using .loc[row\_indexer,col\_indexer]: value instead.

see the caeavts in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
train_df['item-width-div-item-height'] = train_df['item-width']/(train_df['item-height']+1)
#train_df['item-width-mul-item-height'] = train_df['item-width']*(train_df['item-height']+1)

#python:Input:17-94efc97250011: SettingWithCopyWarning
#A value is trying to be set on a copy of a slice from a DataFrame.
#Try using .loc[row\_indexer,col\_indexer]: value instead.

see the caeavts in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
train_df['item-width-mul-item-height'] = train_df['item-width']*(train_df['item-height']+1)
#train_df['item-height-div-item-width'] = train_df['item-height']/(train_df['item-width']+1)

#python:Input:17-94efc97250011: SettingWithCopyWarning
#A value is trying to be set on a copy of a slice from a DataFrame.
#Try using .loc[row\_indexer,col\_indexer]: value instead.

see the caeavts in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
train_df['item-height-div-item-width'] = train_df['item-height']/(train_df['item-width']+1)
#train_df['item-height-mul-item-width'] = train_df['item-height']*(train_df['item-width']+1)

#python:Input:17-94efc97250011: SettingWithCopyWarning
#A value is trying to be set on a copy of a slice from a DataFrame.
#Try using .loc[row\_indexer,col\_indexer]: value instead.

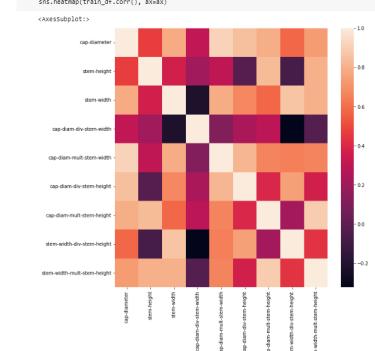
#train_df['item-height-mul-item-width'] = train_df['item-height']*(train_df['item-width']+1)

#train_df['cap-diameter-cap-surface-cap-color-doess-bruis-or-bled-gill-attachment-gill-spacing-gill-color-stem-height-stem-width-ring-type-spore-print-color-habitat-season-cap-diam-div-stem-width-cap-diam-mul-stem-width-cap-diam-div-item-height-cap-diam-mul-item-height-item-width-div-item-height-item-width-mul-item-height']

0   15.26   x   g   o   f   t   e   NaN   w   16.95   17.09   ...   g   NaN   d   w   0.43580   260.794   0.850139   258.6570   0.920809   268.6755
1   16.60   x   g   o   f   t   e   NaN   w   17.89   18.19   ...   g   NaN   d   u   0.865034   301.9540   0.874144   298.6340   0.957873   327.2381
2   14.07   x   g   o   f   t   e   NaN   w   17.00   17.74   ...   g   NaN   d   w   0.780000   249.6018   0.746404   250.4460   0.843617   315.7270
3   14.17   f   h   e   f   t   e   NaN   w   15.77   15.98   ...   p   NaN   d   w   0.834511   226.4366   0.844961   223.4690   0.962092   252.0046
4   14.64   x   h   o   f   t   e   NaN   w   16.53   17.20   ...   p   NaN   d   w   0.804396   251.8000   0.835140   241.9992   0.861175   254.1860
...   ...
#train_df['cap-diameter-cap-surface-cap-color-doess-bruis-or-bled-gill-attachment-gill-spacing-gill-color-stem-height-stem-width-ring-type-spore-print-color-habitat-season-cap-diam-div-stem-width-cap-diam-mul-stem-width-cap-diam-div-item-height-cap-diam-mul-item-height-item-width-div-item-height-item-width-mul-item-height']

50208  1.18   s   s   y   f   f   f   f   f   3.93   6.22   ...   f   NaN   d   a   0.163435   7.3396   0.239351   4.6374   1.261663   24.4446
50209  1.27   f   s   y   f   f   f   f   f   3.18   5.43   ...   f   NaN   d   a   0.197512   6.8961   0.303828   4.0386   1.299043   17.2074
50210  1.27   s   s   y   f   f   f   f   f   3.06   6.37   ...   f   NaN   d   u   0.172320   0.0899   0.261317   4.9022   1.310700   24.5882
50211  1.24   f   s   y   f   f   f   f   f   3.58   5.44   ...   f   NaN   d   u   0.192547   6.7455   0.271939   4.4144   1.162982   19.3664
50212  1.17   s   s   y   f   f   f   f   f   3.25   5.45   ...   f   NaN   d   u   0.181395   6.3705   0.275294   3.8025   1.262353   17.7125
```

(14) Fig. 4 x = plt.subplots(Figsize=(9,9), axes=True)



```
[17] # keep cap-diameter-div-stem-width and cap-diameter-div-item-height because they are the most unique
remove_features = ['cap-diam-mul-stem-width', 'cap-diam-mul-item-height', 'item-width-div-item-height', 'item-width-mul-item-height']
test_df['cap-diam-div-stem-width'] = test_df['cap-diameter']/(test_df['stem-width']+1)
test_df['cap-diam-div-item-height'] = test_df['cap-diameter']/(test_df['item-height']+1)

#python:Input:17-2d316449:11: SettingWithCopyWarning
#A value is trying to be set on a copy of a slice from a DataFrame.
#Try using .loc[row\_indexer,col\_indexer]: value instead.

see the caeavts in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
#test_df['cap-diam-div-item-height'] = test_df['cap-diameter']/(test_df['item-height']+1)

#python:Input:17-2d316449:11: SettingWithCopyWarning
#A value is trying to be set on a copy of a slice from a DataFrame.
#Try using .loc[row\_indexer,col\_indexer]: value instead.

see the caeavts in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
test_df['cap-diam-div-item-height'] = test_df['cap-diameter']/(test_df['item-height']+1)
```

4. Data Processing

```
[18] from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.compose import ColumnTransformer, make_column_transformer
```

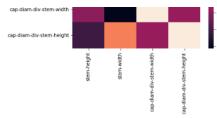
```
[19] # 4.1. Convert the 'class' column into labels: 'p' (poisonous) -> 0, 'e' (edible) -> 1
converted_labels = labels.replace('p', 0, 'e': 1)
```

```
[20] # 4.2. You can drop columns if you see fit
#drop_gill_color due to high correlation
#consider dropping veil color
features = train_df.drop(remove_features, axis=1)
test_df = test_df.drop(remove_features, axis=1)
test_df = test_df.drop(['habitat'], axis=1)
```

(11) sns.heatmap(train_x.corr())

features are pretty good





[11] # 3. See any incomplete data? We learned how to deal with them in project 1.
`train_X.isnull().sum()

```
cap-shape          0
cap-surface        12298
cap-color          0
bruise            0
does-bruise-or-bleed  0
gill-attachment    7766
gill-color          0
stem-color          0
stem-width          42388
stem-length         36939
stem-root           0
ring-ring           0
ring-size          1745
habitat             0
season              0
cap-dim-div-stem-width  0
cap-dim-div-stem-height 0
dtype: int64
```

[12] multiple_features = ['stem-height']
`train_X = train_X.dropna(subset=features, axis=1)
test_X = test_X.dropna(subset=features, axis=1)
train_X.dropna()`

	cap-shape	cap-surface	cap-color	does-bruise-or-bleed	gill-attachment	gill-color	stem-height	stem-width	stem-surface	stem-color	has-ring	ring-type	habitat	season	cap-dim-div-item-width	cap-dim-div-stem-height
0	x	g	o	f	e	w	16.95	17.09	y	w	t	g	d	w	0.043560	0.050139
1	x	g	o	f	e	w	17.99	16.19	y	w	t	g	d	w	0.060524	0.071444
2	x	g	o	f	e	w	17.80	17.74	y	w	t	g	d	w	0.750000	0.764044
3	f	h	e	f	e	w	15.77	15.88	y	w	t	p	d	w	0.834511	0.844961
4	x	h	o	f	e	w	16.93	17.20	y	w	t	p	d	w	0.043496	0.035140
...
48443	f	k	n	f	p	n	3.83	5.94	k	n	t	f	d	u	0.940922	1.351967
48444	s	k	n	f	p	n	4.37	7.18	k	n	t	f	d	u	0.707824	1.008346
48445	f	k	y	f	p	n	3.72	5.76	k	n	t	f	d	u	0.059763	0.944915
48446	s	k	y	f	p	n	4.41	5.35	k	n	t	f	d	u	0.070866	1.022181
48447	f	k	n	f	p	n	3.43	5.08	k	n	t	f	d	u	0.028847	1.137686

14264 rows × 16 columns

[14] X_train, X_test, y_train, y_test = train_X, test_X, converted_labels.iloc[50011], converted_labels.iloc[50011]

```
[15] print(X_train.shape, y_train.shape)
print(X_test.shape, y_test.shape)
(S5011, 16) (50011, 1)
(50011, 16) (50011, 1)
```

[16] numerical_features = ['stem-height', 'stem-width', 'cap-dim-div-stem-width', 'cap-dim-div-stem-height']
categorical_features = ['cap-shape', 'cap-color', 'does-bruise-or-bled', 'gill-color', 'has-ring', 'habitat', 'season']

num_pipeline = Pipeline([
 ('std_scaler', StandardScaler())
])
full_pipeline = ColumnTransformer([
 ("num", num_pipeline, numerical_features),
 ("cat", OneHotEncoder(categories='auto', handle_unknown='ignore'), categorical_features)
])
prep_train = full_pipeline.fit_transform(X_train).toarray()
prep_text = full_pipeline.transform(X_test).toarray()

6. Logistic Regression & Statistical Hypothesis Testing

```
[17] from sklearn.linear_model import LogisticRegression
from sklearn import metrics
logreg = LogisticRegression(solver='liblinear')
logreg.fit(prep_train, y_train)
```

```
[18] predicted_log = logreg.predict(prep_text)
print("With %s %%" % (accuracy_, metrics.accuracy_score(y_text, predicted_log)))
print("With %s %%" % (precision_, metrics.precision_score(y_text, predicted_log)))
print("With %s %%" % (recall_, metrics.recall_score(y_text, predicted_log)))
print("With %s %%" % (f1, metrics.f1_score(y_text, predicted_log)))
```

```
Accuracy: 0.437639
Precision: 0.518287
Recall: 0.344711
F1: 0.389527
```

[19] import statsmodels.api as sm
stats_data = train_df[['stem-height', 'stem-width', 'cap-dim-div-stem-width', 'cap-dim-div-stem-height']]
stats_data

	stem-height	stem-width	cap-dim-div-stem-width	cap-dim-div-stem-height
0	16.95	17.09	0.845660	0.050139
1	17.99	18.19	0.865604	0.071444
2	17.80	17.74	0.750000	0.764044
3	15.77	15.88	0.834511	0.844961
4	16.93	17.20	0.043496	0.035140
...
50208	3.93	6.22	0.163435	0.239351
50209	3.16	5.43	0.197512	0.303528
50210	3.88	6.37	0.172320	0.261517
50211	3.58	5.44	0.192547	0.271930
50212	3.25	5.45	0.181395	0.275294

50213 rows × 4 columns

```
[20] y_train
```

	0	1	2	3	4	...
50208	0	0	0	0	0	...
50209	0	0	0	0	0	...
50210	0	0	0	0	0	...
50211	0	0	0	0	0	...
50212	0	0	0	0	0	...

50213 rows × 5 columns

```
[21] sm_X = sm.add_constant(stats_data)
logit_mod = sm.Logit(y_train, sm_X)
logit_mod = logit_mod.fit()
print(logit_mod.summary())
optimization terminated successfully.
    Current function value: 0.644947
    Iterations 5
    Logit Regression Results
```

	coef	std err	z	P> z	[0.05	0.95]
Intercept	4.1238	1.3561	3.056	0.002	0.128	0.975
stem-height	0.0365	0.005	7.527	<0.001	0.026	0.064
stem-width	-0.0001	0.000	-0.024	0.981	-0.001	0.000
cap-dim-div-stem-width	0.186	0.054	3.458	<0.001	0.023	0.353
cap-dim-div-stem-height	0.1527	0.041	3.725	<0.001	0.026	0.405

7. Dimensionality Reduction using PCA

```
[22] # PCA: https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html
prep_train.shape
```

(50213, 6)

[23] from sklearn import decomposition

```
[24] pca = decomposition.PCA(n_components=5)
pca_train = pca.fit_transform(prep_train)
pca_train.shape
```

(50213, 5)

[25] logreg.fit(pca_train, y_train)

```
[26] LogisticRegression
```

(LogisticRegression(solver='liblinear'))

```
[27] logreg.fit(pca_train, y_train)
predicted_pca = logreg.predict(pca_test)
```

```
print("With %s %%" % (accuracy_, metrics.accuracy_score(y_text, predicted_pca)))
print("With %s %%" % (precision_, metrics.precision_score(y_text, predicted_pca)))
print("With %s %%" % (recall_, metrics.recall_score(y_text, predicted_pca)))
print("With %s %%" % (f1, metrics.f1_score(y_text, predicted_pca)))
```

accuracy_

precision_

recall_

f1

▼ 8. Experiment with any 2 other models (Non-Ensemble)

```
[37] # Model: https://scikit-learn.org/stable/modules/supervised_learning.html
from sklearn.neighbors import KNeighborsClassifier
neighbors = KNeighborsClassifier(n_neighbors=1000)
neighbors.fit(x_train, y_train)
neighbors.predict(neighbors.predict(x_test))
print("1000 NN % (% accuracy)", metrics.accuracy_score(y_test, neighbors.predict()))
print("1000 NN % (% precision)", metrics.precision_score(y_test, neighbors.predict()))
print("1000 NN % (% recall)", metrics.recall_score(y_test, neighbors.predict()))
print("1000 NN % (% f1", metrics.f1_score(y_test, neighbors.predict()))

Accuracy: 0.611551
Precision: 0.606000
Recall: 0.616100
F1: 0.609331

[38] from sklearn.tree import DecisionTreeClassifier
tree = DecisionTreeClassifier(max_depth=2)
tree.fit(x_train, y_train)
tree.predict(tree.predict(x_test))
print("2 DTC % (% accuracy)", metrics.accuracy_score(y_test, tree.predict()))
print("2 DTC % (% precision)", metrics.precision_score(y_test, tree.predict()))
print("2 DTC % (% recall)", metrics.recall_score(y_test, tree.predict()))
print("2 DTC % (% f1", metrics.f1_score(y_test, tree.predict()))

Accuracy: 0.571942
Precision: 0.656607
Recall: 0.551329
F1: 0.592333
```

▼ 9. Experiment with 1 Ensemble Method

```
[39] # Ensemble method: https://scikit-learn.org/stable/modules/ensemble.html
from sklearn.ensemble import RandomForestClassifier
forest = RandomForestClassifier(n_estimators=100)
forest.fit(x_train, y_train)
forest.predict(forest.predict(x_test))
print("100 RF % (% accuracy)", metrics.accuracy_score(y_test, forest.predict()))
print("100 RF % (% precision)", metrics.precision_score(y_test, forest.predict()))
print("100 RF % (% recall)", metrics.recall_score(y_test, forest.predict()))
print("100 RF % (% f1", metrics.f1_score(y_test, forest.predict()))

Accuracy: 0.591194
Precision: 0.620556
Recall: 0.486467
F1: 0.534044
```

▼ 10. Cross-Validation & Hyperparameter Tuning for All 3 Models

```
[40] # Cross-validation: https://scikit-learn.org/stable/modules/cross_validation.html
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score

kf = KFold(n_splits = 10, random_state=42, shuffle=True)

knn_clf = KNeighborsClassifier()
knn_clf.fit(x_train, y_train)
cv_knn = cross_val_score(knn_clf, x_train, y_train, cv=kf)
print("10 CV Accuracy: %.2f%% (%.2f%%)" % (cv_knn.mean()*100.0, cv_knn.std() * 100.0))

tree_clf = DecisionTreeClassifier()
tree_clf.fit(x_train, y_train, cv=kf)
print("Decision Tree CV Accuracy: %.2f%% (%.2f%%)" % (tree_clf.score(x_train, y_train) * 100.0, tree_clf.cv_results_.mean("std_err")))

forest_clf = RandomForestClassifier()
forest_clf.fit(x_train, y_train, cv=kf)
print("Random Forest CV Accuracy: %.2f%% (%.2f%%)" % (forest_clf.score(x_train, y_train) * 100.0, forest_clf.cv_results_.mean("std_err")))

cv_knn = cv_knn.mean()
Decision Tree CV Accuracy: 90.65
Decision Forest CV Accuracy: 95.215

[41] # Hyperparameter tuning: https://scikit-learn.org/stable/modules/grid_search.html
from sklearn.model_selection import GridSearchCV

knn_params = [
    {"n_neighbors": [5, 55, 555], "weights": ["uniform", "distance"], "algorithm": ["brute", "kd_tree", "ball_tree", "auto"]}
]
knn_clf_grid = GridSearchCV(estimator=knn_clf, param_grid=knn_params, cv=kf, verbose=True)
knn_clf_grid.fit(x_train, y_train)
print("KNN Grid Best Estimator: ", knn_clf_grid.best_estimator_.get_params())
fitting 10 folds for each of 32 candidates, totalling 320 fits
{'Algorithm': 'brute', 'LeafSize': 50, 'Metric': 'minkowski', 'n_neighbors': 5, 'p': 1, 'Weights': 'distance'}

[42] tree_params = [
    {"max_depth": [2, 5, 8], "criterion": ["gini", "entropy"], "log_loss": []}
]
tree_clf_grid = GridSearchCV(estimator=tree_clf, param_grid=tree_params, cv=kf, verbose=True)
tree_clf_grid.fit(x_train, y_train)
print("Tree Grid Best Estimator: ", tree_clf_grid.best_estimator_.get_params())
fitting 10 folds for each of 3 candidates, totalling 90 fits
{'Criterion': 'gini', 'LeafSize': 50, 'Metric': 'minkowski', 'n_neighbors': None, 'n_jobs': None, 'n_neighbors': 5, 'p': 1, 'Weights': 'distance'}

[43] forest_params = [
    {"n_estimators": [1, 5, 10], "max_features": [None, 1, 3]}
]
forest_clf_grid = GridSearchCV(estimator=forest_clf, param_grid=forest_params, cv=kf, verbose=True)
forest_clf_grid.fit(x_train, y_train)
print("Forest Grid Best Estimator: ", forest_clf_grid.best_estimator_.get_params())
fitting 10 folds for each of 6 candidates, totalling 60 fits
{'n_estimators': 1, 'MaxDepth': 8, 'ClassWeight': None, 'Criterion': 'gini', 'MaxDepth': 8, 'MaxFeatures': None, 'MinInurityDecrease': 0.0, 'MinSamplesLeaf': 1, 'MinSamplesSplit': 2, 'MinWeightFractionLeaf': 0.0, 'RandomState': None, 'Splitter': 'best'}
```

▼ 11. Report Final Results

```
[44] # 1. K-Nearest Neighbors Classifier
final_knn = KNeighborsClassifier(algorithm='brute', leaf_size=30, metric='minkowski', metric_params=None, n_jobs=None, n_neighbors=5, p=2, weights='distance')
final_knn.fit(x_train, y_train)
final_knn.predict(final_knn.predict(x_test))
print("Final KNN % (% accuracy)", metrics.accuracy_score(y_test, final_knn.predict()))
print("Final KNN % (% precision)", metrics.precision_score(y_test, final_knn.predict()))
print("Final KNN % (% recall)", metrics.recall_score(y_test, final_knn.predict()))
print("Final KNN % (% f1", metrics.f1_score(y_test, final_knn.predict()))

Accuracy: 0.48455
Precision: 0.548515
Recall: 0.480924
F1: 0.464309

[45] final_tree = DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini', max_depth=None, max_features=None, max_leaf_nodes=None,
    min_impurity_decrease=0.0, min_samples_leaf=1, min_samples_split=1, min_weight_fraction_leaf=0.0,
    presort=False, splitter='best')
final_tree.fit(x_train, y_train)
print("Final Tree % (% accuracy)", metrics.accuracy_score(y_test, final_tree.predict()))
print("Final Tree % (% precision)", metrics.precision_score(y_test, final_tree.predict()))
print("Final Tree % (% recall)", metrics.recall_score(y_test, final_tree.predict()))
print("Final Tree % (% f1", metrics.f1_score(y_test, final_tree.predict()))

Accuracy: 0.590440
Precision: 0.630492
Recall: 0.516264
F1: 0.549980

[46] final_forest = RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None, criterion='gini', max_depth=None, max_features=None,
    max_leaf_nodes=None, min_impurity_decrease=0.0, min_samples_leaf=1, min_samples_split=1, min_weight_fraction_leaf=0.0, n_estimators=10, n_jobs=None,
    oob_score=False, random_state=None)
final_forest.fit(x_train, y_train)
final_forest.predict(final_forest.predict(x_test))
print("Final Forest % (% accuracy)", metrics.accuracy_score(y_test, final_forest.predict()))
print("Final Forest % (% precision)", metrics.precision_score(y_test, final_forest.predict()))
print("Final Forest % (% recall)", metrics.recall_score(y_test, final_forest.predict()))
print("Final Forest % (% f1", metrics.f1_score(y_test, final_forest.predict()))

Accuracy: 0.545965
Precision: 0.609116
Recall: 0.527399
F1: 0.485437
```