



## ▼ This project needs get\_dummies and CountVectorizer / TfidfVectorizer before or after tts

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 %matplotlib inline
```

## ▼ Load the Amazon Review data

```
In [2]: 1 file_path = r"D:\LIUZHICHENG\Udemy\Machine Learning\Classification Bootcamp\ML Classification Package\6. Decision Trees and Random Forests\Amazon Reviews.csv"
2 df = pd.read_csv(filepath_or_buffer=file_path, sep="\t")
3 df.head()
```

Out[2]:

	rating	date	variation	verified_reviews	feedback
0	5	31-Jul-18	Charcoal Fabric	Love my Echo!	1
1	5	31-Jul-18	Charcoal Fabric	Loved it!	1
2	4	31-Jul-18	Walnut Finish	Sometimes while playing a game, you can answer...	1
3	5	31-Jul-18	Charcoal Fabric	I have had a lot of fun with this thing. My 4 ...	1
4	5	31-Jul-18	Charcoal Fabric	Music	1

```
In [3]: 1 df.isnull().sum()
```

Out[3]:

	rating	date	variation	verified_reviews	feedback
0	0	0	0	0	0
1	0	0	0	0	0
2	0	0	0	0	0
3	0	0	0	0	0
4	0	0	0	0	0

## ► EDA

[...]

## ▼ Data Preprocessing

```
In [9]: 1 df.drop(columns=["date", "rating"], inplace=True)
2 # get dummy for "variation"
3 df = pd.get_dummies(data=df, columns=["variation"], drop_first=True) # drop_first=True -> dummy trap
```

```
In [10]: 1 # CountVectorizer
2 from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
3 cv = CountVectorizer()
4 review_cv = cv.fit_transform(df["verified_reviews"])
```

```
In [11]: 1 review_cv.shape
```

Out[11]: (3150, 4044)

```
In [12]: 1 print(cv.get_feature_names())
kie', 'telephone', 'television', 'tell', 'telling', 'tells', 'temp', 'temperature', 'temps', 'tempting', 'ten', 'tend', 'tends', 'terminology', 'terrible', 'terrific', 'test', 'tested', 'testing', 'texas', 'text', 'texts', 'tg', 'tge', 'than', 'thank', 'thanks', 'that', 'thats', 'the', 'theater', 'theecho', 'their', 'theirs', 'them', 'themes', 'themselves', 'then', 'theories', 'there', 'therefore', 'thermostat', 'these', 'thestand', 'thete', 'they', 'thick', 'thing', 'things', 'think', 'thinking', 'third', 'this', 'thongs', 'thorough', 'thoroughly', 'those', 'thou', 'though', 'thought', 'thoughts', 'thousands', 'three', 'thrilled', 'through', 'throughout', 'throw', 'thrown', 'thru', 'thu', 'thumb', 'thumbs', 'thunderstorm', 'thursday', 'ti', 'tickled', 'tiempo', 'tiens', 'ties', 'til', 'till', 'time', 'timer', 'timers', 'times', 'timing', 'tin', 'ting', 'tinker', 'tinkering', 'tinny', 'tiny', 'tipping', 'tips', 'tired', 'title', 'tivo', 'to', 'toda', 'today', 'toddler', 'together', 'toilet', 'told', 'tomorrow', 'tomy', 'ton', 'tones', 'tons', 'tony', 'too', 'took', 'too1', 'tools', 'tooth', 'top', 'topic', 'tosca', 'total', 'totally', 'totally', 'tou', 'touch', 'touching', 'touted', 'toward', 'towards', 'tower', 'town', 'toy', 'tp', 'track', 'traditional', 'traffic', 'trailer', 'trailers', 'trained', 'trainees', 'training', 'transferring', 'travel', 'traveling', 'travelling', 'traves', 'treadmill', 'treat', 'treble', 'trek', 'tremendous', 'trending', 'trial', 'tricks', 'tricky', 'tried', 'tries', 'trigger', 'trip', 'trivia', 'trouble', 'trouble shooting', 'troublesome', 'troubling', 'true', 'truly', 'trust', 'try', 'trying', 'tube', 'tubi', 'tune', 'tunein', 'tunes', 'turn', 'turned', 'turning', 'turns', 'tv', 'tvs', 'tweeter', 'twice', 'twist', 'twitter', 'two', 'ty', 'type', 'typed', 'types', 'typical', 'typically', 'typing', 'titulos', 'udefulness', 'ugly', 'uhyour', 'ummm', 'un', 'unable', 'unacceptable', 'unavailable', 'unbelievable', 'uncle', 'under', 'underestimated', 'understand', 'understanding', 'understands', 'understood', 'unexpected', 'unfortunately', 'unhappy', 'unhelpful', 'unico', 'unimportant', 'uninstall', 'unique', 'unit', 'units', 'unless', 'unlikely', 'unlimited', 'unlocking', 'unnannounced', 'unnecessary', 'unobtrusive', 'unplug', 'unplugged', 'unresponsive', 'unsettling', 'untapped', 'until', 'unusable', 'unused', 'unwitty', 'unwrapped', 'up', 'upcoming', 'update', 'updated', 'updates', 'updating', 'upgrade', 'upgraded', 'upgrades', 'upgrading',
```

```
In [13]: 1 type(review_cv)
```

Out[13]: scipy.sparse.csr.csr\_matrix

```
In [14]: 1 review_cv.toarray()
```

Out[14]: array([[0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 ...,
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0]], dtype=int64)

```
In [15]: 1 df = df.drop(columns="verified_reviews")
2 df.head()
```

Out[15]: . . . variation Black variation Black variation Black variation Charcoal variation Configuration: variation Heather variation Oak

feedback	Dot	Plus	Show	Spot	Fabric	Fire TV Stick	Gray Fabric	Finish
0	1	0	0	0	0	1	0	0
1	1	0	0	0	0	1	0	0
2	1	0	0	0	0	0	0	0
3	1	0	0	0	0	1	0	0
4	1	0	0	0	0	1	0	0

In [16]: 1 pd.DataFrame(review\_cv.toarray())

Out[16]:

0	1	2	3	4	5	6	7	8	9	...	4034	4035	4036	4037	4038	4039	4040	4041	4042	4043
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
2	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	1	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
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
3145	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
3146	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
3147	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
3148	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
3149	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0

3150 rows × 4044 columns

In [17]: 1 df = pd.concat(objs=[df, pd.DataFrame(review\_cv.toarray())], axis=1)  
2 df

Out[17]:

feedback	variation_Black_Dot	variation_Black_Plus	variation_Black_Show	variation_Black_Spot	variation_Charcoal_Fabric	variation_Configuration_Fire_TV_Stick	variation_Heather_Gray_Fabric	variation_C_Finish
0	1	0	0	0	0	1	0	0
1	1	0	0	0	0	1	0	0
2	1	0	0	0	0	0	0	0
3	1	0	0	0	0	1	0	0
4	1	0	0	0	0	1	0	0
...	...	...	...	...	...	...	...	...
3145	1	1	0	0	0	0	0	0
3146	1	1	0	0	0	0	0	0
3147	1	1	0	0	0	0	0	0
3148	1	0	0	0	0	0	0	0
3149	1	1	0	0	0	0	0	0

3150 rows × 4060 columns

In [18]: 1 X = df.drop(columns="feedback").to\_numpy()  
2 X

Out[18]: array([[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],  
[1, 0, 0, ..., 0, 0, 0]], dtype=int64)

In [19]: 1 y = df["feedback"].to\_numpy()  
2 y

Out[19]: array([1, 1, 1, ..., 1, 1, 1], dtype=int64)

## Model Training

CountVectorizer before tts

In [20]: 1 # Decision Trees  
2 from sklearn.model\_selection import train\_test\_split  
3 # stratify=y  
4 X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2, random\_state=0, stratify=y)  
5  
6 from sklearn.tree import DecisionTreeClassifier  
7 clf\_dt = DecisionTreeClassifier()  
8 clf\_dt.fit(X\_train, y\_train)  
9  
10 from sklearn.metrics import confusion\_matrix, accuracy\_score, classification\_report  
11  
12 predictions = clf\_dt.predict(X\_test)  
13 print(classification\_report(y\_test, predictions))  
14 print(accuracy\_score(y\_test, predictions))  
15 print(confusion\_matrix(y\_test, predictions))

	precision	recall	f1-score	support
0	0.59	0.53	0.56	51
1	0.96	0.97	0.96	579
accuracy			0.93	630
macro avg	0.77	0.75	0.76	630
weighted avg	0.93	0.93	0.93	630

0.9317460317460318

```
[[ 27  24]
 [ 19 560]]
```

```
In [21]: # Random Forest
1 from sklearn.model_selection import train_test_split
2 # stratify=y
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0, stratify=y)
4
5 from sklearn.ensemble import RandomForestClassifier
6 clf_rf = RandomForestClassifier()
7 clf_rf.fit(X_train, y_train)
8
9 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
10
11 predictions = clf_rf.predict(X_test)
12 print(classification_report(y_test, predictions))
13 print(accuracy_score(y_test, predictions))
14 print(confusion_matrix(y_test, predictions))
```

	precision	recall	f1-score	support
0	1.00	0.25	0.41	51
1	0.94	1.00	0.97	579
accuracy			0.94	630
macro avg	0.97	0.63	0.69	630
weighted avg	0.94	0.94	0.92	630

0.9396825396825397  
[[ 13 38]
 [ 0 579]]

## Model Training

### CountVectorizer after tts

```
In [22]: # read csv
1 df = pd.read_csv(filepath_or_buffer=file_path, sep="\t")
2 # drop columns, "rating", "date"
3 df.drop(columns=["rating", "date"], inplace=True)

In [23]: # get_dummies for "variation"
1 df = pd.get_dummies(data=df, columns=["variation"], drop_first=True)

In [24]: X = df.drop(columns="feedback")
1 y = df["feedback"]

In [25]: from sklearn.model_selection import train_test_split
1 # stratify=y
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0, stratify=y)

In [26]: X_train.shape
Out[26]: (2520, 16)

In [27]: X_test.shape
Out[27]: (630, 16)

In [28]: from sklearn.feature_extraction.text import CountVectorizer
1 cv = CountVectorizer()
2 X_train_review_cv = cv.fit_transform(raw_documents=X_train["verified_reviews"])
3 X_test_review_cv = cv.transform(raw_documents=X_test["verified_reviews"])

In [29]: X_train_review_cv
Out[29]: <2520x3696 sparse matrix of type '<class 'numpy.int64'>'  
with 48733 stored elements in Compressed Sparse Row format>

In [30]: X_train_review_cv.shape
Out[30]: (2520, 3696)

In [31]: X_test_review_cv.shape
Out[31]: (630, 3696)

In [32]: X_train_review_cv.toarray()
Out[32]: array([[0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 ...,
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0]], dtype=int64)

In [33]: X_test_review_cv.toarray()
Out[33]: array([[0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 ...,
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0]], dtype=int64)

In [34]: X_train = X_train.drop(columns="verified_reviews").reset_index() # reset_index() for X_train
1 # concatenate X_train and "review" sparse matrix
2 X_train = pd.concat(objs=[X_train, pd.DataFrame(X_train_review_cv.toarray())], axis=1)
3
4 X_train
```

```
Out[34]:
```

index	variation_Black	variation_Black	variation_Black	variation_Black	variation_Charcoal	variation_Configuration:	variation_Heather	variation_Oak
0	0	0	0	0	0	Eira TV Chair	Grey Fabric	Finch

	Dot	Plus	Show	Spot	Twink	Fire TV Stick	Gray Fabric	Oak Finish
0	1655	0	0	0	0	0	0	0
1	177	0	0	0	0	0	0	1
2	2381	0	0	0	0	1	0	0
3	273	0	0	0	0	0	1	0
4	1344	0	0	0	1	0	0	0
...	...	...	...	...	...	...	...	...
2515	2295	0	0	0	0	1	0	0
2516	1276	0	0	0	0	0	0	0
2517	3113	1	0	0	0	0	0	0
2518	367	0	0	0	0	0	0	0
2519	3007	1	0	0	0	0	0	0

2520 rows × 3712 columns

```
In [35]: X_test = X_test.drop(columns="verified_reviews").reset_index() # reset_index() for X_test
# concatenate X_test and "review" sparse matrix
X_test = pd.concat(objs=[X_test, pd.DataFrame(X_test_review_cv.toarray())], axis=1)
X_test
```

Out[35]:

	index	variation_Black Dot	variation_Black Plus	variation_Black Show	variation_Black Spot	variation_Charcoal Fabric	variation_Configuration: Fire TV Stick	variation_Heather Gray Fabric	variation_Oak Finish
0	876	0	0	0	0	1	0	0	0
1	1172	0	0	0	0	0	0	0	0
2	149	0	0	0	0	1	0	0	0
3	232	0	0	0	0	0	0	0	0
4	1562	0	0	1	0	0	0	0	0
...	...	...	...	...	...	...	...	...	...
625	2509	1	0	0	0	0	0	0	0
626	898	0	0	0	0	1	0	0	0
627	2737	1	0	0	0	0	0	0	0
628	2210	0	0	0	0	0	1	0	0
629	1957	0	1	0	0	0	0	0	0

630 rows × 3712 columns

```
In [36]: # Decision Trees
from sklearn.tree import DecisionTreeClassifier
clf_dt = DecisionTreeClassifier()
clf_dt.fit(X_train, y_train)

from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
predictions = clf_dt.predict(X_test)
print(classification_report(y_test, predictions))
print(accuracy_score(y_test, predictions))
print(confusion_matrix(y_test, predictions))
```

	precision	recall	f1-score	support
0	0.66	0.53	0.59	51
1	0.96	0.98	0.97	579
accuracy			0.94	630
macro avg	0.81	0.75	0.78	630
weighted avg	0.93	0.94	0.94	630

0.9396825396825397  
[[ 27 24]  
 [ 14 565]]

```
In [37]: # Random Forest
from sklearn.ensemble import RandomForestClassifier
clf_rf = RandomForestClassifier()
clf_rf.fit(X_train, y_train)

from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
predictions = clf_rf.predict(X_test)
print(classification_report(y_test, predictions))
print(accuracy_score(y_test, predictions))
print(confusion_matrix(y_test, predictions))
```

	precision	recall	f1-score	support
0	0.93	0.25	0.40	51
1	0.94	1.00	0.97	579
accuracy			0.94	630
macro avg	0.93	0.63	0.68	630
weighted avg	0.94	0.94	0.92	630

0.9380952380952381  
[[ 13 38]  
 [ 1 578]]

## Model Training

### TfidfVectorizer before tts

```
In [38]: # read csv
df = pd.read_csv(filepath_or_buffer=file_path, sep="\t")
```

```

3 # drop columns, "rating", "date"
4 df.drop(columns=["rating", "date"], inplace=True)
5
6 # get_dummies for "variation"
7 df = pd.get_dummies(data=df, columns=["variation"], drop_first=True)
8
9 from sklearn.feature_extraction.text import TfidfVectorizer
10
11 tv = TfidfVectorizer()
12 review_tv = tv.fit_transform(raw_documents=df["verified_reviews"])
13
14 df = pd.concat(objs=[df, pd.DataFrame(review_tv.toarray())], axis=1)
15 df.drop(columns="verified_reviews", inplace=True)
16
17 X = df.drop(columns="feedback")
18 y = df["feedback"]
19
20 from sklearn.model_selection import train_test_split
21 # stratify=y
22 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0, stratify=y)

```

In [39]:

```

1 # Decision Trees
2 from sklearn.tree import DecisionTreeClassifier
3
4 clf_dt = DecisionTreeClassifier()
5 clf_dt.fit(X_train, y_train)
6
7 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
8
9 predictions = clf_dt.predict(X_test)
10 print(classification_report(y_test, predictions))
11 print(accuracy_score(y_test, predictions))
12 print(confusion_matrix(y_test, predictions))

precision    recall   f1-score   support
          0       0.55      0.59      0.57       51
          1       0.96      0.96      0.96      579

accuracy                           0.93      630
macro avg       0.75      0.77      0.76      630
weighted avg    0.93      0.93      0.93      630

0.926984126984127
[[ 30  21]
 [ 25 554]]

```

In [40]:

```

1 # Random Forest
2 from sklearn.ensemble import RandomForestClassifier
3
4 clf_rf = RandomForestClassifier()
5 clf_rf.fit(X_train, y_train)
6
7 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
8
9 predictions = clf_rf.predict(X_test)
10 print(classification_report(y_test, predictions))
11 print(accuracy_score(y_test, predictions))
12 print(confusion_matrix(y_test, predictions))

precision    recall   f1-score   support
          0       1.00      0.25      0.41       51
          1       0.94      1.00      0.97      579

accuracy                           0.94      630
macro avg       0.97      0.63      0.69      630
weighted avg    0.94      0.94      0.92      630

0.9396825396825397
[[ 13  38]
 [ 0 579]]

```

## Model Training

### TfidfVectorizer after tts

In [41]:

```

1 # read csv
2 df = pd.read_csv(filepath_or_buffer=file_path, sep="\t")
3 # drop columns, "rating", "date"
4 df.drop(columns=["rating", "date"], inplace=True)
5
6 # get_dummies for "variation"
7 df = pd.get_dummies(data=df, columns=["variation"], drop_first=True)
8
9 X = df.drop(columns="feedback")
10 y = df["feedback"]
11
12 from sklearn.model_selection import train_test_split
13 # stratify=y
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0, stratify=y)
15
16 from sklearn.feature_extraction.text import TfidfVectorizer
17
18 tv = TfidfVectorizer()
19 X_train_review_tv = tv.fit_transform(raw_documents=X_train["verified_reviews"])
20 X_test_review_tv = tv.transform(raw_documents=X_test["verified_reviews"])
21
22 X_train = X_train.drop(columns="verified_reviews").reset_index() # reset_index() for X_train
23 # concatenate X_train and "review" sparse matrix
24 X_train = pd.concat(objs=[X_train, pd.DataFrame(X_train_review_tv.toarray())], axis=1)
25
26 X_test = X_test.drop(columns="verified_reviews").reset_index() # reset_index() for X_test
27 # concatenate X_test and "review" sparse matrix
28 X_test = pd.concat(objs=[X_test, pd.DataFrame(X_test_review_tv.toarray())], axis=1)

```

In [42]:

```

1 # Decision Trees
2 from sklearn.tree import DecisionTreeClassifier
3
4 clf_dt = DecisionTreeClassifier()

```

```

5 clf_dt.fit(X_train, y_train)
6
7 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
8
9 predictions = clf_dt.predict(X_test)
10 print(classification_report(y_test, predictions))
11 print(accuracy_score(y_test, predictions))
12 print(confusion_matrix(y_test, predictions))

```

	precision	recall	f1-score	support
0	0.62	0.47	0.53	51
1	0.95	0.97	0.96	579
accuracy			0.93	630
macro avg	0.78	0.72	0.75	630
weighted avg	0.93	0.93	0.93	630

0.933333333333333  
[[ 24 27]  
[ 15 564]]

In [43]:

```

1 # Random Forest
2 from sklearn.ensemble import RandomForestClassifier
3
4 clf_rf = RandomForestClassifier()
5 clf_rf.fit(X_train, y_train)
6
7 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
8
9 predictions = clf_rf.predict(X_test)
10 print(classification_report(y_test, predictions))
11 print(accuracy_score(y_test, predictions))
12 print(confusion_matrix(y_test, predictions))

```

	precision	recall	f1-score	support
0	0.89	0.31	0.46	51
1	0.94	1.00	0.97	579
accuracy			0.94	630
macro avg	0.92	0.66	0.72	630
weighted avg	0.94	0.94	0.93	630

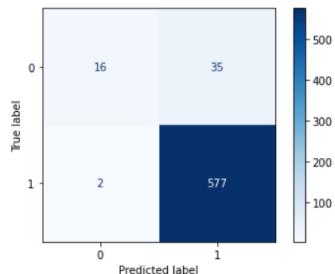
0.9412698412698413  
[[ 16 35]  
[ 2 577]]

In [44]:

```

1 from sklearn.metrics import plot_confusion_matrix
2
3 display = plot_confusion_matrix(estimator=clf_rf, X=X_test, y_true=y_test, cmap="Blues", values_format=".3g");

```



In [45]:

```

1 display.confusion_matrix

```

Out[45]: array([[ 16, 35],  
[ 2, 577]], dtype=int64)