

```
import seaborn as sns
import pandas as pd
import matplotlib.pyplot as plt
```

```
df=pd.read_csv('survey_lung_cancer.csv')
```

```
df.head(3)
```

	GENDER	AGE	SMOKING	YELLOW_FINGERS	ANXIETY	PEER_PRESSURE	CHRONIC DISEASE	FATIGUE
0	M	69	1	2	2	1	1	2
1	M	74	2	1	1	1	2	2
2	F	59	1	1	1	2	1	2

Next steps:

Generate code with df

New interactive sheet

```
df.tail()
```

	GENDER	AGE	SMOKING	YELLOW_FINGERS	ANXIETY	PEER_PRESSURE	CHRONIC DISEASE	FATIGUE
304	F	56	1	1	1	2	2	2
305	M	70	2	1	1	1	1	2
306	M	58	2	1	1	1	1	1
307	M	67	2	1	2	1	1	2
308	M	62	1	1	1	2	1	2

```
df.sample()
```

	GENDER	AGE	SMOKING	YELLOW_FINGERS	ANXIETY	PEER_PRESSURE	CHRONIC DISEASE	FATIGUE
195	M	69	1	2	2	1	1	1

```
df.shape
```

(309, 16)



df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 309 entries, 0 to 308
Data columns (total 16 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   GENDER                                309 non-null    object
1   AGE                                   309 non-null    int64
2   SMOKING                              309 non-null    int64
3   YELLOW_FINGERS                       309 non-null    int64
4   ANXIETY                              309 non-null    int64
5   PEER_PRESSURE                        309 non-null    int64
6   CHRONIC DISEASE                      309 non-null    int64
7   FATIGUE                             309 non-null    int64
8   ALLERGY                              309 non-null    int64
9   WHEEZING                             309 non-null    int64
10  ALCOHOL CONSUMING                    309 non-null    int64
11  COUGHING                             309 non-null    int64
12  SHORTNESS OF BREATH                  309 non-null    int64
13  SWALLOWING DIFFICULTY                309 non-null    int64
14  CHEST PAIN                           309 non-null    int64
15  LUNG_CANCER                          309 non-null    object
dtypes: int64(14), object(2)
memory usage: 38.8+ KB
```

df.describe()

	AGE	SMOKING	YELLOW_FINGERS	ANXIETY	PEER_PRESSURE	CHRONIC DISEASE
count	309.000000	309.000000	309.000000	309.000000	309.000000	309.000000
mean	62.673139	1.563107	1.569579	1.498382	1.501618	1.504854
std	8.210301	0.496806	0.495938	0.500808	0.500808	0.500787
min	21.000000	1.000000	1.000000	1.000000	1.000000	1.000000
25%	57.000000	1.000000	1.000000	1.000000	1.000000	1.000000
50%	62.000000	2.000000	2.000000	1.000000	2.000000	2.000000
75%	69.000000	2.000000	2.000000	2.000000	2.000000	2.000000
max	87.000000	2.000000	2.000000	2.000000	2.000000	2.000000

df.dtypes

0

GENDER	object
AGE	int64
SMOKING	int64
YELLOW_FINGERS	int64
ANXIETY	int64
PEER_PRESSURE	int64
CHRONIC_DISEASE	int64
FATIGUE	int64
ALLERGY	int64
WHEEZING	int64
ALCOHOL_CONSUMING	int64
COUGHING	int64
SHORTNESS_OF_BREATH	int64
SWALLOWING_DIFFICULTY	int64
CHEST_PAIN	int64
LUNG_CANCER	object

dtype: object

```
df.isnull().sum()
```

```

0
GENDER 0
AGE 0
SMOKING 0
YELLOW_FINGERS 0
ANXIETY 0
PEER_PRESSURE 0
CHRONIC DISEASE 0
FATIGUE 0
ALLERGY 0
WHEEZING 0
ALCOHOL CONSUMING 0
COUGHING 0
SHORTNESS OF BREATH 0
SWALLOWING DIFFICULTY 0
CHEST PAIN 0
LUNG_CANCER 0

dtype: int64

```

```

df['LUNG_CANCER'].value_counts()

count
LUNG_CANCER
YES      270
NO       39

dtype: int64

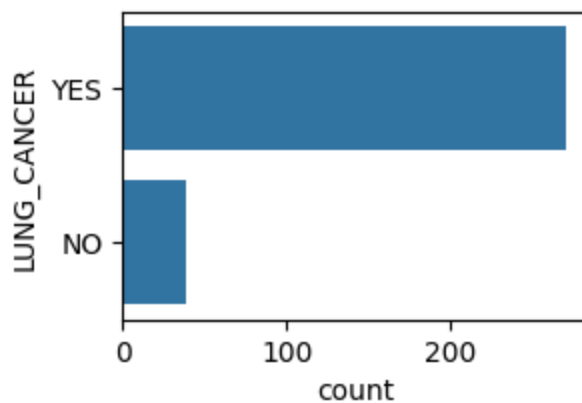
```

```

plt.figure(figsize=(3,2))
sns.countplot(df['LUNG_CANCER'])

```

<Axes: xlabel='count', ylabel='LUNG_CANCER'>



encoding

x and y me break

train test split

standard scaler

mode train

✓ encoding

Double-click (or enter) to edit

```
from sklearn.preprocessing import LabelEncoder  
le=LabelEncoder()  
df['GENDER']=le.fit_transform(df['GENDER'])  
df['LUNG_CANCER']=le.fit_transform(df['LUNG_CANCER'])
```

df

	GENDER	AGE	SMOKING	YELLOW_FINGERS	ANXIETY	PEER_PRESSURE	CHRONIC DISEASE	FATIGUE
0	1	69	1	2	2	1	1	2
1	1	74	2	1	1	1	2	2
2	0	59	1	1	1	2	1	2
3	1	63	2	2	2	1	1	1
4	0	63	1	2	1	1	1	1
...
304	0	56	1	1	1	2	2	2
305	1	70	2	1	1	1	1	2
306	1	58	2	1	1	1	1	1
307	1	67	2	1	2	1	1	2
308	1	62	1	1	1	2	1	2

309 rows × 16 columns

Next steps:

[Generate code with df](#)

[New interactive sheet](#)

```
x = df.drop('LUNG_CANCER', axis=1)
y = df['LUNG_CANCER']
```

```
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=42)
```

```
from sklearn.preprocessing import StandardScaler
sc=StandardScaler()
x_train=sc.fit_transform(x_train)
x_test=sc.transform(x_test)
```

```
from sklearn.svm import SVC
model=SVC()
model.fit(x_train,y_train)
```

▼ SVC ⓘ ?
SVC()

```
model.score(x_train,y_train)*100,model.score(x_test,y_test)*100
```

```
(94.73684210526315, 96.7741935483871)
```

```
y_pred=model.predict(x_test)
```

```
from sklearn.metrics import confusion_matrix  
cm=confusion_matrix(y_test,y_pred)
```

```
cm
```

```
array([[ 1,  1],  
       [ 1, 59]])
```

```
from sklearn.metrics import accuracy_score, classification_report  
print("Accuracy:",accuracy_score(y_test,y_pred)*100)
```