

Import libraries

```
In [2]: import pandas as pd
import numpy as np
from sklearn import preprocessing
import matplotlib.pyplot as plt
import matplotlib inline
import seaborn as sns

df = pd.read_csv('bank-additional-full.csv',sep=';')
term_deposits = df.copy()
# Name a grasp of how our data looks.
df.head()
```

```
Out[2]:
```

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	...	campaign
0	56	housemaid	married	basic4y	no	no	no	telephone	may	mon	...	
1	57	services	married	highschool	unknown	no	no	telephone	may	mon	...	
2	37	services	married	highschool	no	yes	no	telephone	may	mon	...	
3	40	admin.	married	basic6y	no	no	no	telephone	may	mon	...	
4	56	services	married	highschool	no	no	yes	telephone	may	mon	...	

5 rows x 21 columns

```
In [3]: df.describe()
```

```
Out[3]:
```

	age	duration	campaign	pdays	previous	emp.var.rate	cons.price.idx	cons
count	41188.000000	41188.000000	41188.000000	41188.000000	41188.000000	41188.000000	41188.000000	4118
mean	40.02406	258.285010	2.567593	962.475454	0.172963	0.081886	93.575664	-4
std	10.42125	259.279249	2.770014	186.910907	0.439401	1.570960	0.578840	-4
min	17.00000	0.000000	1.000000	0.000000	0.000000	-3.400000	92.201000	-5
25%	32.00000	102.000000	1.000000	999.000000	0.000000	-1.800000	93.075000	-4
50%	38.00000	180.000000	2.000000	999.000000	0.000000	1.100000	93.749000	-4
75%	47.00000	319.000000	3.000000	999.000000	0.000000	1.400000	93.994000	-3
max	98.00000	4918.000000	56.000000	999.000000	7.000000	1.400000	94.767000	-2

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 41188 entries, 0 to 41187
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  ---                ---
0   age                    41188 non-null  int64
1   job                    41188 non-null  object
2   marital                41188 non-null  object
3   education              41188 non-null  object
4   default                41188 non-null  object
5   housing                41188 non-null  object
6   loan                   41188 non-null  object
7   contact                41188 non-null  object
8   month                  41188 non-null  object
9   day_of_week            41188 non-null  object
10  duration                41188 non-null  int64
11  campaign                41188 non-null  int64
12  pdays                   41188 non-null  int64
13  previous                41188 non-null  object
14  poutcome                41188 non-null  object
15  emp.var.rate            41188 non-null  float64
16  cons.price.idx           41188 non-null  float64
17  cons.conf.idx            41188 non-null  float64
18  euribor3m                41188 non-null  float64
19  nr.employed              41188 non-null  float64
20  y                        41188 non-null  object
dtypes: float64(5), int64(5), object(11)
memory usage: 6.6+ MB
```

```
In [5]: df.isnull().sum()
```

```
age          0
job           0
marital       0
education     0
default       0
housing       0
loan          0
contact       0
month         0
day_of_week   0
duration      0
campaign      0
pdays        0
previous      0
poutcome      0
emp.var.rate   0
cons.price.idx 0
cons.conf.idx  0
euribor3m      0
nr.employed    0
y              0
dtype: int64
```

Fortunately, there are no missing values. If there were missing values we will have to fill them with the median, mean or mode. I tend to use the median but in this scenario there is no need to fill any missing values. This will definitely make our job easier!

```
In [6]: df.nunique()
```

```
Out[6]:
```

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	...	campaign
age	78											
job	12											
marital	4											
education	8											
default	3											
housing	3											
loan	3											
contact	2											
month	10											
day_of_week	5											
duration	1544											
campaign	27											
pdays	27											
previous	8											
poutcome	3											
emp.var.rate	10											
cons.price.idx	26											
cons.conf.idx	26											
euribor3m	11											
nr.employed	11											
y	2											
dtype:	int64											

Now,lets read the test data

```
In [13]: test=pd.read_csv('bank-additional.csv',sep=';')

In [14]: test.head()
```

```
Out[14]:
```

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	...	campaign
0	30	blue-collar	married	basic9y	no	yes	no	cellular	may	fri	...	
1	39	services	married	high.school	no	no	no	telephone	jun	fri	...	
2	25	services	single	high.school	no	yes	no	telephone	may	wed	...	
3	38	services	married	basic9y	no	unknown	unknown	telephone	jun	fri	...	
4	47	admin.	married	university.degree	no	yes	no	cellular	nov	mon	...	

5 rows x 21 columns

Concating train and test data

```
In [15]: data=pd.concat((df,test))
data.replace(['basic6y','basic4y','basic9y'],'basic',inplace=True)

In [16]: data.head()
```

```
Out[16]:
```

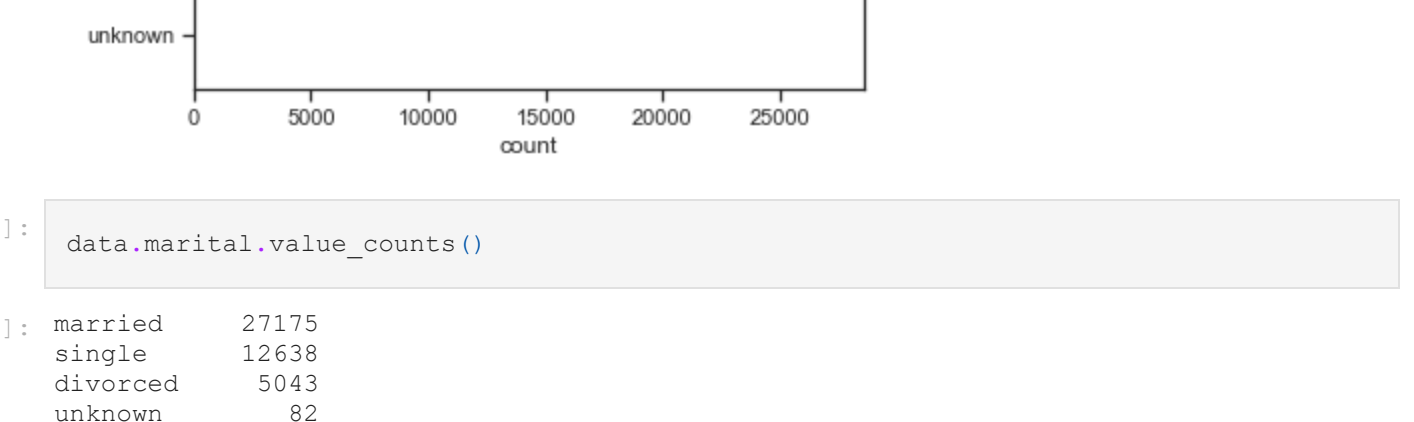
	age	job	marital	education	default	housing	loan	contact	month	day_of_week	...	campaign
0	56	housemaid	married	basic	no	no	no	telephone	may	mon	...	
1	57	services	married	highschool	unknown	no	no	telephone	may	mon	...	
2	37	services	married	highschool	no	yes	no	telephone	may	mon	...	
3	40	admin.	married	basic	no	no	no	telephone	may	mon	...	
4	56	services	married	highschool	no	no	yes	telephone	may	mon	...	

5 rows x 21 columns

Exploratory data analysis

```
In [17]: sns.set(style="ticks", color_codes=True)
sns.countplot(y='job', data=data)
```

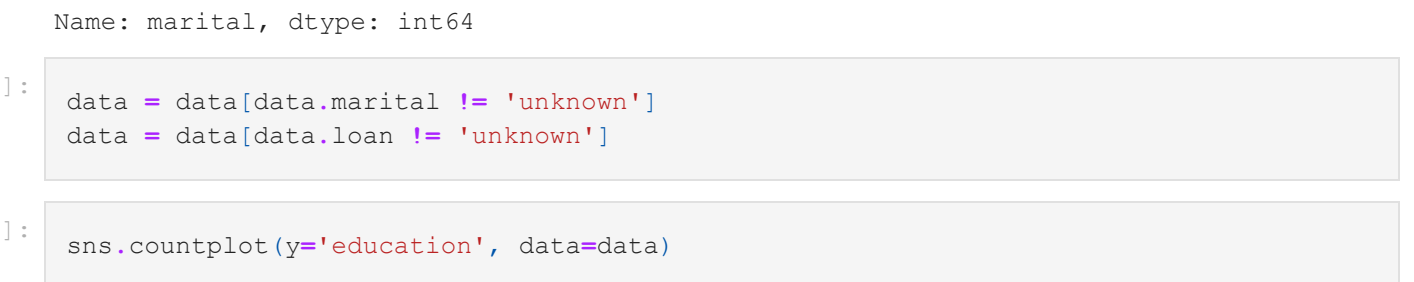
```
Out[17]:
```



```
In [18]: data = data[data.job != 'unknown']

In [19]: sns.countplot(y='marital', data=data)
```

```
Out[19]:
```



```
In [20]: data.marital.value_counts()
```

```
Out[20]:
```

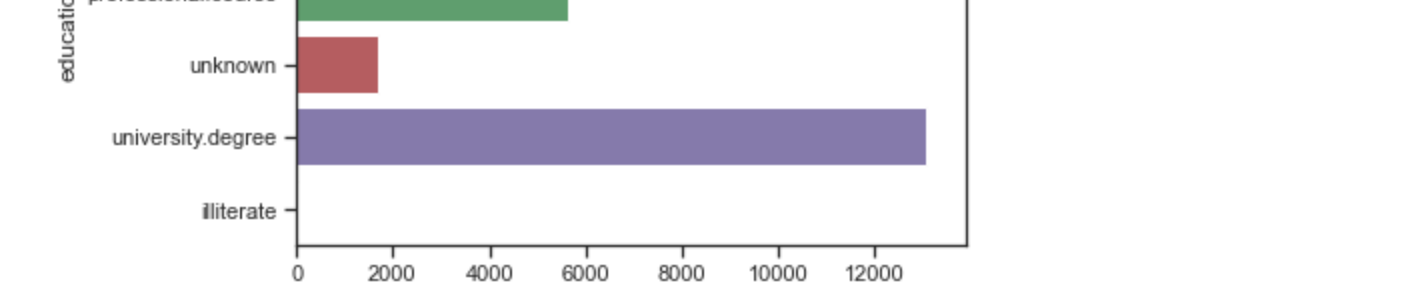
marital	count
married	27175
single	12638
divorced	5043
unknown	82

Name: marital, dtype: int64

```
In [21]: data = data[data.marital != 'unknown']
data = data[data.loan != 'unknown']

In [22]: sns.countplot(y='education', data=data)
```

```
Out[22]:
```



```
In [23]: data = data[data.education != 'illiterate']

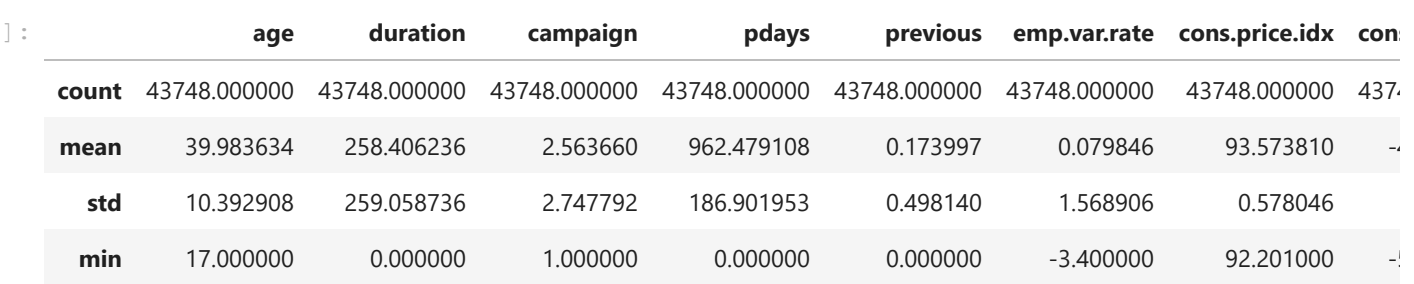
In [24]: data.describe()
```

```
Out[24]:
```

	age	duration	campaign	pdays	previous	emp.var.rate	cons.price.idx	cons
count	43748.000000	43748.000000	43748.000000	43748.000000	43748.000000	43748.000000	43748.000000	437
mean	39.983634	258.406236	2.563660	962.479108	0.173997	0.079846	93.573810	-4
std	10.392908	259.058736	2.747792	186.901953	0.498140	1.568906	0.578046	-4
min	17.00000	0.000000	1.000000	0.000000	0.000000	-3.400000	92.201000	-5
25%	32.00000	102.000000	1.000000	999.000000	0.000000	-1.800000	93.075000	-4
50%	38.00000	180.000000	2.000000	999.000000	0.000000	1.100000	93.749000	-4
75%	47.00000	320.000000	3.000000	999.000000	0.000000	1.400000	93.994000	-3
max	98.00000	4918.000000	43.000000	999.000000	7.000000	1.400000	94.767000	-2

```
In [25]: sns.countplot(y='y', data=data)
```

```
Out[25]:
```



It is clear that the dataset is imblanced

Data preprocessing

Firstly,convert the categorical features into numeric using label encoding.

In label encoding in Python, we replace the categorical value with a numeric value between 0 and the number of classes minus 1. If the categorical variable value contains 5 distinct classes, we use (0, 1, 2, 3, and 4).

```
In [26]: def categorize(df):
new_df = df.copy()
le = preprocessing.LabelEncoder()

new_df['job'] = le.fit_transform(new_df['job'])
new_df['marital'] = le.fit_transform(new_df['marital'])
new_df['education'] = le.fit_transform(new_df['education'])
new_df['default'] = le.fit_transform(new_df['default'])
new_df['housing'] = le.fit_transform(new_df['housing'])
new_df['month'] = le.fit_transform(new_df['month'])
new_df['loan'] = le.fit_transform(new_df['loan'])
new_df['contact'] = le.fit_transform(new_df['contact'])
new_df['day_of_week'] = le.fit_transform(new_df['day_of_week'])
new_df['poutcome'] = le.fit_transform(new_df['poutcome'])
new_df['y'] = le.fit_transform(new_df['y'])

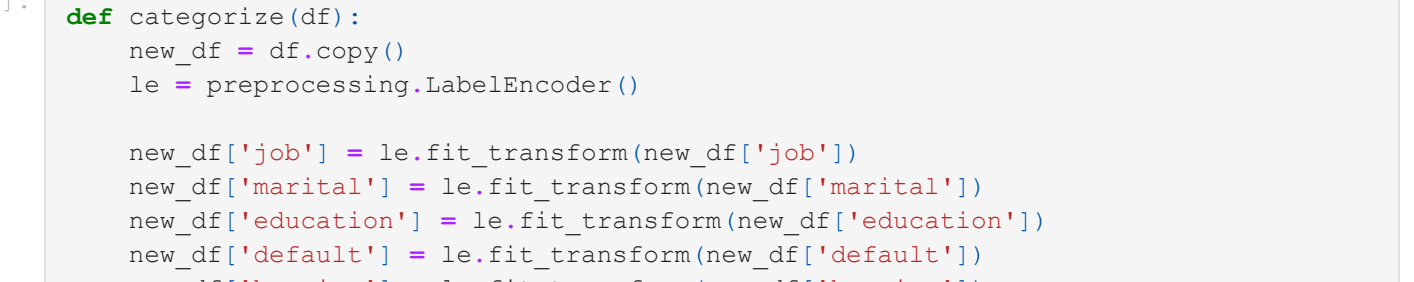
return new_df

In [27]: data = categorize(data)
data = data.convert_objects(convert_numeric=True)
```

Checking for outliers using boxplots

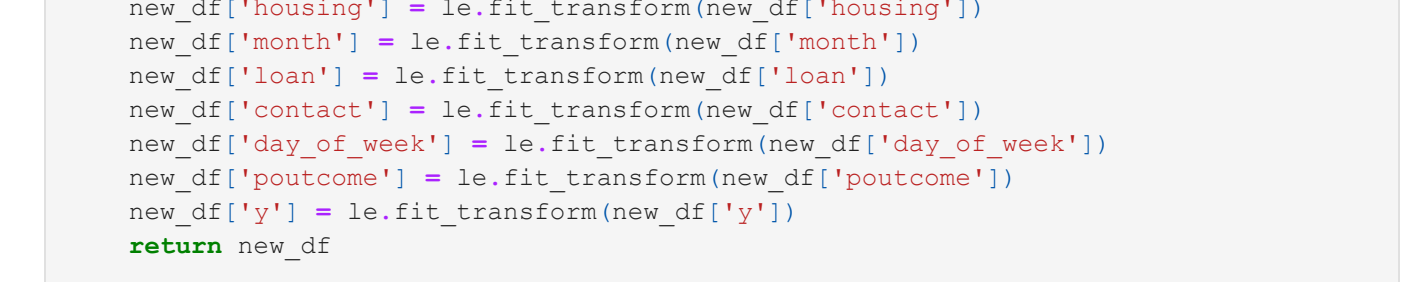
```
In [28]: sns.boxplot(x='y', y='duration', data=data)
```

```
Out[28]:
```



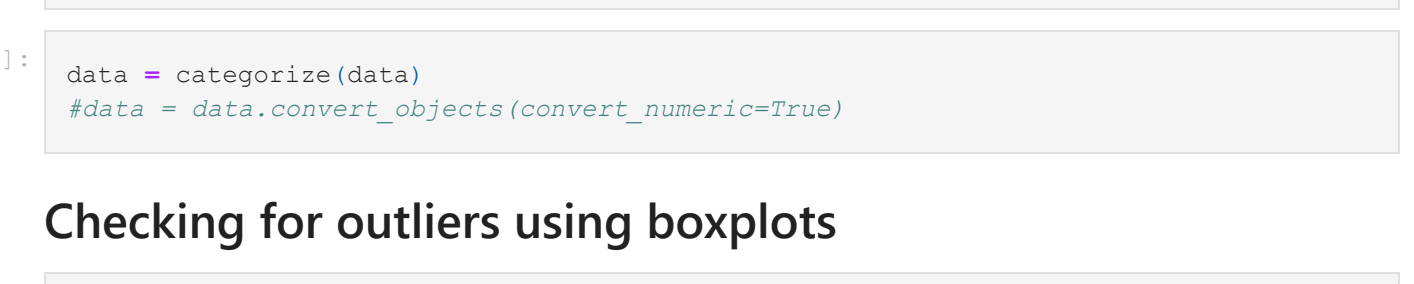
```
In [29]: sns.boxplot(x='y', y='education', data=data)
```

```
Out[29]:
```



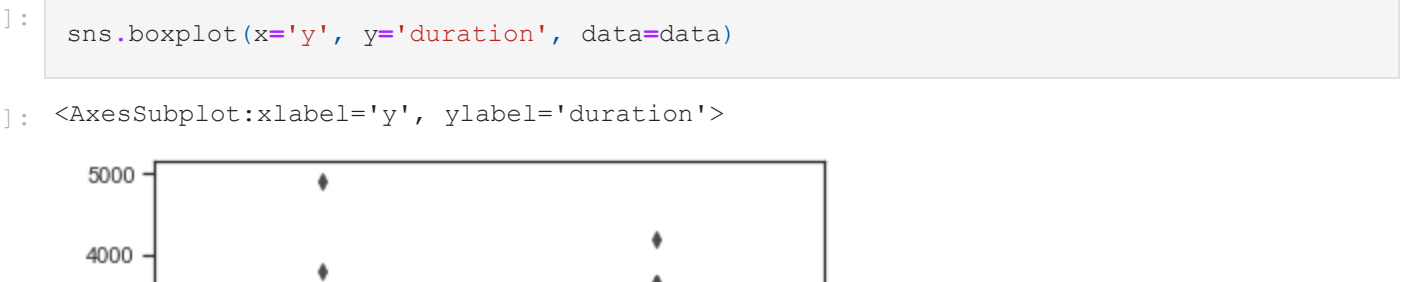
```
In [30]: sns.boxplot(x='y', y='housing', data=data)
```

```
Out[30]:
```



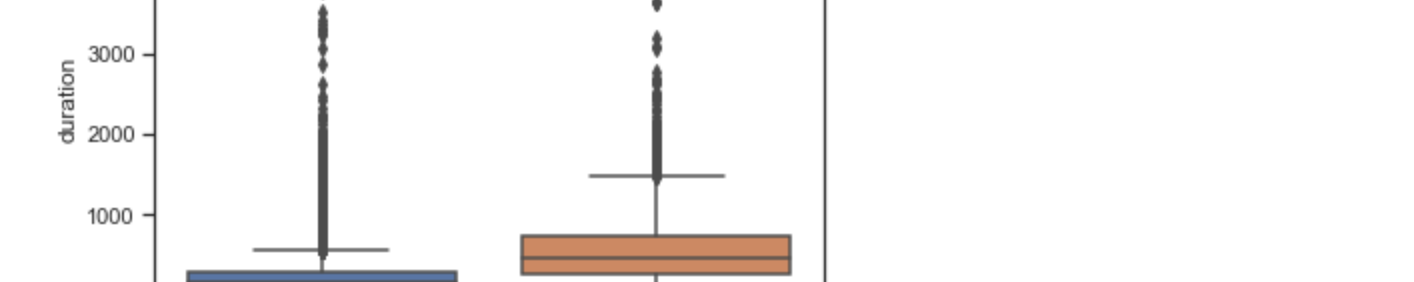
```
In [31]: sns.boxplot(x='y', y='age', data=data)
```

```
Out[31]:
```



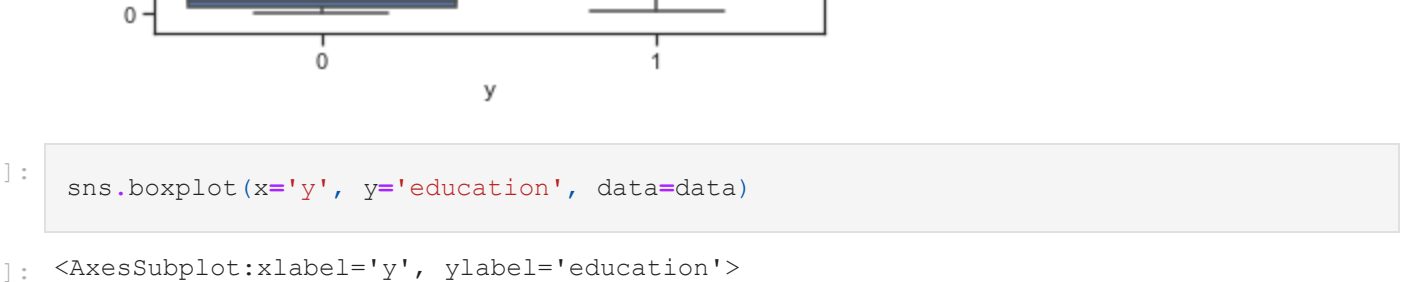
```
In [32]: sns.boxplot(x='y', y='job', data=data)
```

```
Out[32]:
```



```
In [33]: sns.boxplot(x='y', y='campaign', data=data )
```

```
Out[33]:
```



Removing outliers

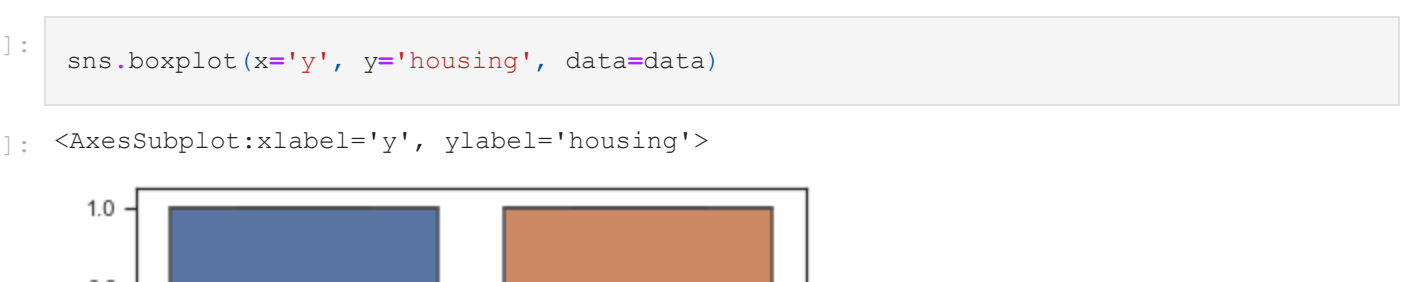
```
In [34]: def remove_outliers(df, column, minimum, maximum):
col_values = df[column].values
return df[np.where(np.logical_and(col_values<minimum, col_values>maximum), col_values==True)]

In [35]: min_val = data["duration"].min()
min_val = 1500
data = remove_outliers(df=data, column='duration', minimum=min_val, maximum=max_val)
max_val = 80
data = remove_outliers(df=data, column='age', minimum=min_val, maximum=max_val)
min_val = data["campaign"].min()
max_val = 6
data = remove_outliers(df=data, column='campaign', minimum=min_val, maximum=max_val)
```

Dropping less meaningful columns

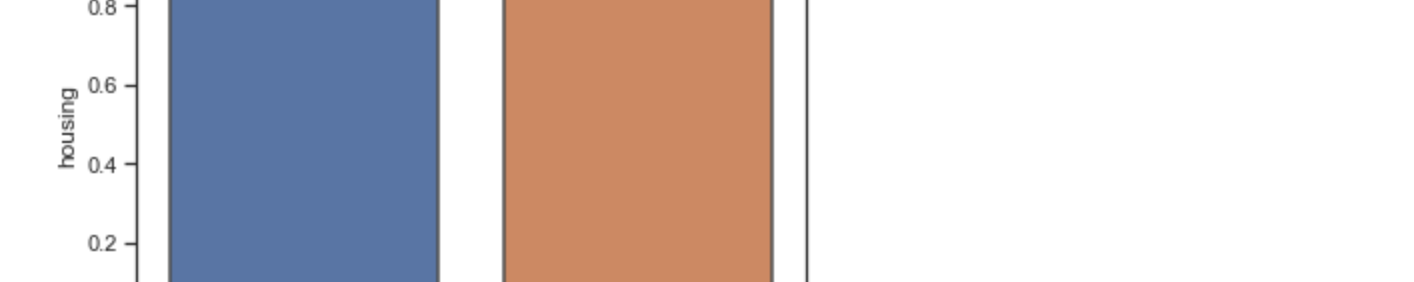
```
In [36]: sns.countplot(x='education',hue='y',data=data)
```

```
Out[36]:
```



```
In [37]: sns.countplot(x='default',hue='y',data=data)
```

```
Out[37]:
```

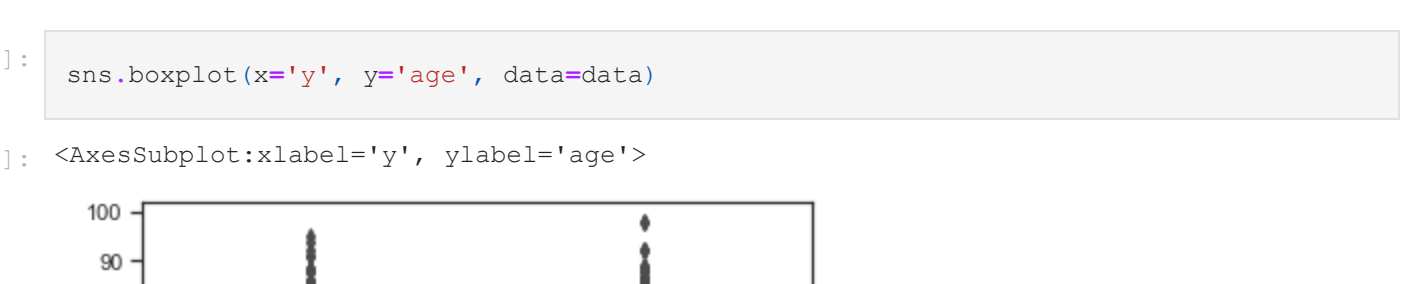


It is skewed to 0. So We can drop this

```
In [38]: data = data.drop('default',axis=1)

In [39]: sns.countplot(x='poutcome',hue='y',data=data)
```

```
Out[39]:
```

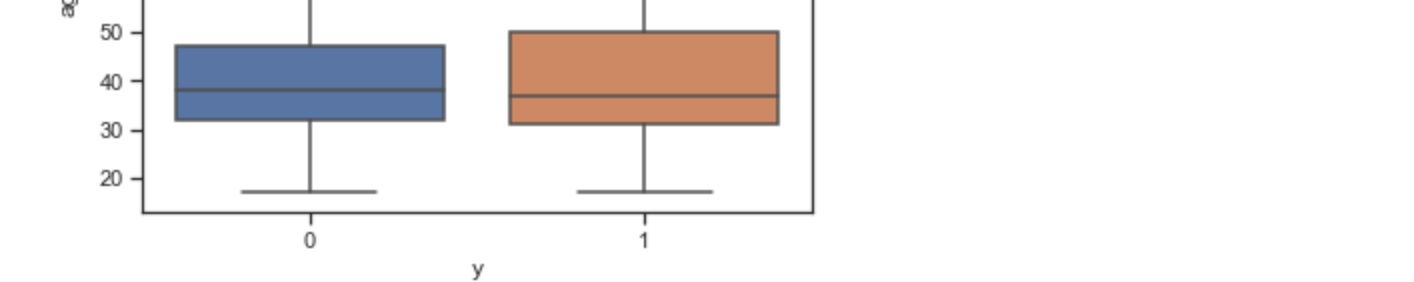


So many non existent values. We can drop this

```
In [40]: data = data.drop('poutcome',axis=1)


In [41]: sns.countplot(x='loan',hue='y',data=data)
```

```
Out[41]:
```



```
In [42]: sns.countplot(x='contact',hue='y',data=data)
```

```
Out[42]:
```



```
In [43]: data =data.drop('contact',axis=1)

In [44]: data = data.drop(['emp.var.rate','cons.price.idx','cons.conf.idx','euribor3m','nr.employed',
In [45]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 43748 entries, 0 to 41187
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  ---                ---
0   age                    43748 non-null  float64
1   job                    43748 non-null  int32
2   marital                43748 non-null  int32
3   education              43748 non-null  int32
4   housing                43748 non-null  int32
5   loan                   43748 non-null  int32
6   month                  43748 non-null  int32
7   day_of_week            43748 non-null  int32
8   duration                43748 non-null  float64
9   campaign                43748 non-null  float64
10  pdays                   43748 non-null  int64
11  previous                43748 non-null  int64
12  y                        43748 non-null  int32
dtypes: float64(3), int32(8), int64(2)
memory usage: 4.3 MB
```

```
In [46]: data.head()
```

```
Out[46]:
```

	age	job	marital	education	housing	loan	month	day_of_week	duration	campaign	pdays	previous
0	56.0	3	1	0	0	0	6	1	261.0	1.0	999	0
1	57.0	7	1	1	0	0	6	1	149.0	1.0	999	0
2	37.0	7	1	1	1	0	6	1	226.0	1.0	999	0
3	40.0	0	1	0	0	0	6	1	151.0	1.0	999	0
4	56.0	7	1	1	0	1	6	1	307.0	1.0	999	0

Splitting into train and test data

import all necessary libraries

```
In [47]: from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.decomposition import PCA
from sklearn import metrics
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
```

```
In [48]: X = data.drop('y',axis = 1).values
y = data['y'].values
X_train, X_test, Y_train, Y_test = train_test_split(X, y, test_size=0.25, random_state=1)
```

```
In [49]: scaler = StandardScaler()
scaler.fit(X_train)
X_train = scaler.fit_transform(X_train)
```

```
In [50]: pca = PCA(n_components=10)
pca.fit(X_train)
X_train = pca.fit_transform(X_train)
```

```
In [51]: X_train.shape
```

```
Out[51]:
```

(32811, 10)

Building different Models and validating using 10 fold cross validation

```
In [52]: models = []
models.append(('LR', LogisticRegression()))
models.append(('LDA', LinearDiscriminantAnalysis()))
models.append(('KNN', KNeighborsClassifier()))
models.append(('Decision-Tree', DecisionTreeClassifier()))
models.append(('Gaussian', GaussianNB()))
models.append(('SVM', SVC()))
models.append(('RandomForest', RandomForestClassifier(max_depth = 8, n_estimators = 120)))
models.append(('Ada', AdaBoostClassifier(n_estimators = 120)))
```

```
In [53]: results = []
names = []
for name, model in models:
kfold = model_selection.KFold(n_splits=10, random_state=42)
cv_results = model_selection.cross_val_score(model, X_train, Y_train, cv=kfold, scoring='accuracy')
results.append(cv_results)
names.append(name)
msg = '{}: {:.2f} {:.2f} {:.2f}'.format(name, cv_results.mean(), cv_results.std(), cv_results.min())
print(msg)
```

```
C:\Users\Threshun Duvvuru\anaconda3\lib\site-packages\sklearn\model_selection\_split.py:293: FutureWarning: Setting a random_state has no effect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuffle=True.
warnings.warn(
LR: 0.903595662058859
```

```
C:\Users\Threshun Duvvuru\anaconda3\lib\site-packages\sklearn\model_selection\_split.py:293: FutureWarning: Setting a random_state has no effect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuffle=True.
warnings.warn(
LDA: 0.8985707323442396
```

```
C:\Users\Threshun Duvvuru\anaconda3\lib\site-packages\sklearn\model_selection\_split.py:293: FutureWarning: Setting a random_state has no effect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuffle=True.
warnings.warn(
KNN: 0.8959804395183542
```

```
C:\Users\Threshun Duvvuru\anaconda3\lib\site-packages\sklearn\model_selection\_split.py:293: FutureWarning: Setting a random_state has no effect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuffle=True.
```