

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
from pandas.api.types import CategoricalDtype
import matplotlib.pyplot as plt
from sklearn.preprocessing import LabelEncoder
from xgboost import XGBClassifier
from sklearn.metrics import accuracy_score
from sklearn.model_selection import GridSearchCV
```

```
In [2]: train_set = pd.read_csv('http://archive.ics.uci.edu/ml/machine-learning-databases/adult/adult.data', header = None)
test_set = pd.read_csv('http://archive.ics.uci.edu/ml/machine-learning-databases/adult/adult.test', skiprows = 1, header = None)
col_labels = ['age', 'workclass', 'fnlwgt', 'education', 'education_num', 'marital_status', 'occupation', 'relationship',
              'race', 'sex', 'capital_gain', 'capital_loss', 'hours_per_week',
              'native_country', 'wage_class']
train_set.columns = col_labels
test_set.columns = col_labels
```

```
In [3]: print('Train size:: {}'.format(train_set.shape))
print('Test size:: {}'.format(test_set.shape))
```

Train size:: (32561, 15)
Test size:: (16281, 15)

```
In [4]: train_set.head()
```

Out[4]:

	age	workclass	fnlwgt	education	education_num	marital_status	occupation	relationship
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife

In [5]: `test_set.head()`

Out[5]:

	age	workclass	fnlwgt	education	education_num	marital_status	occupation	relationship
0	25	Private	226802	11th	7	Never-married	Machine-op-inspct	Own-child
1	38	Private	89814	HS-grad	9	Married-civ-spouse	Farming-fishing	Husband
2	28	Local-gov	336951	Assoc-acdm	12	Married-civ-spouse	Protective-serv	Husband
3	44	Private	160323	Some-college	10	Married-civ-spouse	Machine-op-inspct	Husband
4	18	?	103497	Some-college	10	Never-married	?	Own-child

In [6]: `train_set.describe()`

Out[6]:

	age	fnlwgt	education_num	capital_gain	capital_loss	hours_per_week
count	32561.000000	3.256100e+04	32561.000000	32561.000000	32561.000000	32561.000000
mean	38.581647	1.897784e+05	10.080679	1077.648844	87.303830	40.437456
std	13.640433	1.055500e+05	2.572720	7385.292085	402.960219	12.347429
min	17.000000	1.228500e+04	1.000000	0.000000	0.000000	1.000000
25%	28.000000	1.178270e+05	9.000000	0.000000	0.000000	40.000000
50%	37.000000	1.783560e+05	10.000000	0.000000	0.000000	40.000000
75%	48.000000	2.370510e+05	12.000000	0.000000	0.000000	45.000000
max	90.000000	1.484705e+06	16.000000	99999.000000	4356.000000	99.000000

In [7]: `train_set.isna().sum()`

Out[7]:

```
age          0
workclass    0
fnlwgt       0
education    0
education_num 0
marital_status 0
occupation   0
relationship 0
race         0
sex          0
capital_gain 0
capital_loss 0
hours_per_week 0
native_country 0
wage_class   0
dtype: int64
```

```
In [8]: train_set.describe(include=["object"]).T
```

Out[8]:

	count	unique	top	freq
workclass	32561	9	Private	22696
education	32561	16	HS-grad	10501
marital_status	32561	7	Married-civ-spouse	14976
occupation	32561	15	Prof-specialty	4140
relationship	32561	6	Husband	13193
race	32561	5	White	27816
sex	32561	2	Male	21790
native_country	32561	42	United-States	29170
wage_class	32561	2	<=50K	24720

```
In [9]: test_set.describe(include=["object"]).T
```

Out[9]:

	count	unique	top	freq
workclass	16281	9	Private	11210
education	16281	16	HS-grad	5283
marital_status	16281	7	Married-civ-spouse	7403
occupation	16281	15	Prof-specialty	2032
relationship	16281	6	Husband	6523
race	16281	5	White	13946
sex	16281	2	Male	10860
native_country	16281	41	United-States	14662
wage_class	16281	2	<=50K	12435

```
In [10]: # Check if duplicate data exist
check_dup = train_set.duplicated().any()
print("Are there any duplicated values in data? ", check_dup)

if check_dup:
    train_set = train_set.drop_duplicates()
else:
    print("There are not duplicated values in data.")

check_dup_test = test_set.duplicated().any()
print("Are there any duplicated values in data? ", check_dup)

if check_dup:
    test_set = test_set.drop_duplicates()
else:
    print("There are not duplicated values in data.")
```

```
Are there any duplicated values in data?  True
Are there any duplicated values in data?  True
```

```
In [11]: object_columns=train_set.select_dtypes(include=["object"]).columns
for i in range(len(object_columns)):
    print("----- {}-----".format(object_columns[i]))
    print(train_set[object_columns[i]].value_counts())
```

```
----- workclass-----
Private                22673
Self-emp-not-inc       2540
Local-gov              2093
?                      1836
State-gov              1298
Self-emp-inc           1116
Federal-gov            960
Without-pay            14
Never-worked           7
Name: workclass, dtype: int64
----- education-----
HS-grad                10494
Some-college           7282
Bachelors              5353
Masters                1722
Assoc-voc              1382
11th                   1175
Assoc-acdm             1067
10th                   933
7th-8th                645
Prof-school            576
9th                    514
12th                   433
Doctorate              413
5th-6th                332
1st-4th                166
Preschool              50
Name: education, dtype: int64
----- marital_status-----
Married-civ-spouse     14970
Never-married          10667
Divorced               4441
Separated              1025
Widowed                993
Married-spouse-absent  418
Married-AF-spouse      23
Name: marital_status, dtype: int64
----- occupation-----
Prof-specialty         4136
Craft-repair           4094
Exec-managerial        4065
Adm-clerical           3768
Sales                  3650
Other-service          3291
Machine-op-inspct      2000
?                      1843
Transport-moving       1597
Handlers-cleaners      1369
Farming-fishing        992
Tech-support           927
Protective-serv        649
Priv-house-serv        147
Armed-Forces           9
Name: occupation, dtype: int64
----- relationship-----
Husband                13187
```

```

Not-in-family      8292
Own-child          5064
Unmarried          3445
Wife               1568
Other-relative     981
Name: relationship, dtype: int64
----- race-----
White              27795
Black              3122
Asian-Pac-Islander 1038
Amer-Indian-Eskimo 311
Other              271
Name: race, dtype: int64
----- sex-----
Male              21775
Female            10762
Name: sex, dtype: int64
----- native_country-----
United-States      29153
Mexico             639
?                  582
Philippines        198
Germany            137
Canada             121
Puerto-Rico        114
El-Salvador         106
India              100
Cuba                95
England             90
Jamaica             81
South              80
China               75
Italy               73
Dominican-Republic 70
Vietnam            67
Guatemala          62
Japan               62
Poland              60
Columbia           59
Taiwan             51
Haiti              44
Iran               43
Portugal           37
Nicaragua          34
Peru               31
Greece             29
France             29
Ecuador            28
Ireland            24
Hong               20
Trinidad&Tobago    19
Cambodia           19
Laos               18
Thailand            18
Yugoslavia         16
Outlying-US(Guam-USVI-etc) 14
Hungary            13

```

```

Honduras      13
Scotland      12
Holand-Netherlands  1
Name: native_country, dtype: int64
----- wage_class-----
<=50K      24698
>50K       7839
Name: wage_class, dtype: int64

```

```

In [12]: train_set.replace([' ?'], np.nan, inplace=True)
         test_set.replace([' ?'], np.nan, inplace=True)

```

```

In [13]: train_set.isnull().sum()

```

```

Out[13]: age      0
         workclass  1836
         fnlwgt    0
         education  0
         education_num  0
         marital_status  0
         occupation  1843
         relationship  0
         race      0
         sex      0
         capital_gain  0
         capital_loss  0
         hours_per_week  0
         native_country  582
         wage_class    0
         dtype: int64

```

```

In [14]: test_set.isnull().sum()

```

```

Out[14]: age      0
         workclass  963
         fnlwgt    0
         education  0
         education_num  0
         marital_status  0
         occupation  966
         relationship  0
         race      0
         sex      0
         capital_gain  0
         capital_loss  0
         hours_per_week  0
         native_country  274
         wage_class    0
         dtype: int64

```



```
In [15]: object_columns=train_set.select_dtypes(include=["object"]).columns
         for i in range(len(object_columns)):
             print("----- {}-----".format(object_columns[i]))
             print(train_set[object_columns[i]].value_counts())
```

```
----- workclass-----
Private                22673
Self-emp-not-inc       2540
Local-gov              2093
State-gov              1298
Self-emp-inc           1116
Federal-gov            960
Without-pay            14
Never-worked           7
Name: workclass, dtype: int64
----- education-----
HS-grad                10494
Some-college           7282
Bachelors              5353
Masters                1722
Assoc-voc              1382
11th                   1175
Assoc-acdm             1067
10th                   933
7th-8th                645
Prof-school            576
9th                    514
12th                   433
Doctorate              413
5th-6th                332
1st-4th                166
Preschool              50
Name: education, dtype: int64
----- marital_status-----
Married-civ-spouse     14970
Never-married          10667
Divorced               4441
Separated              1025
Widowed                993
Married-spouse-absent  418
Married-AF-spouse      23
Name: marital_status, dtype: int64
----- occupation-----
Prof-specialty         4136
Craft-repair           4094
Exec-managerial        4065
Adm-clerical           3768
Sales                  3650
Other-service          3291
Machine-op-inspct      2000
Transport-moving       1597
Handlers-cleaners      1369
Farming-fishing        992
Tech-support           927
Protective-serv        649
Priv-house-serv        147
Armed-Forces           9
Name: occupation, dtype: int64
----- relationship-----
Husband                13187
Not-in-family          8292
Own-child              5064
```

```

Unmarried          3445
Wife                1568
Other-relative     981
Name: relationship, dtype: int64
----- race-----
White              27795
Black              3122
Asian-Pac-Islander 1038
Amer-Indian-Eskimo 311
Other              271
Name: race, dtype: int64
----- sex-----
Male              21775
Female            10762
Name: sex, dtype: int64
----- native_country-----
United-States      29153
Mexico              639
Philippines        198
Germany            137
Canada             121
Puerto-Rico        114
El-Salvador        106
India              100
Cuba                95
England             90
Jamaica             81
South              80
China              75
Italy              73
Dominican-Republic 70
Vietnam            67
Japan              62
Guatemala          62
Poland             60
Columbia           59
Taiwan             51
Haiti              44
Iran               43
Portugal           37
Nicaragua          34
Peru               31
France            29
Greece            29
Ecuador           28
Ireland           24
Hong              20
Cambodia           19
Trinidad&Tobago    19
Thailand           18
Laos              18
Yugoslavia         16
Outlying-US(Guam-USVI-etc) 14
Honduras           13
Hungary            13
Scotland           12
Holand-Netherlands 1

```

Name: native_country, dtype: int64

----- wage_class-----

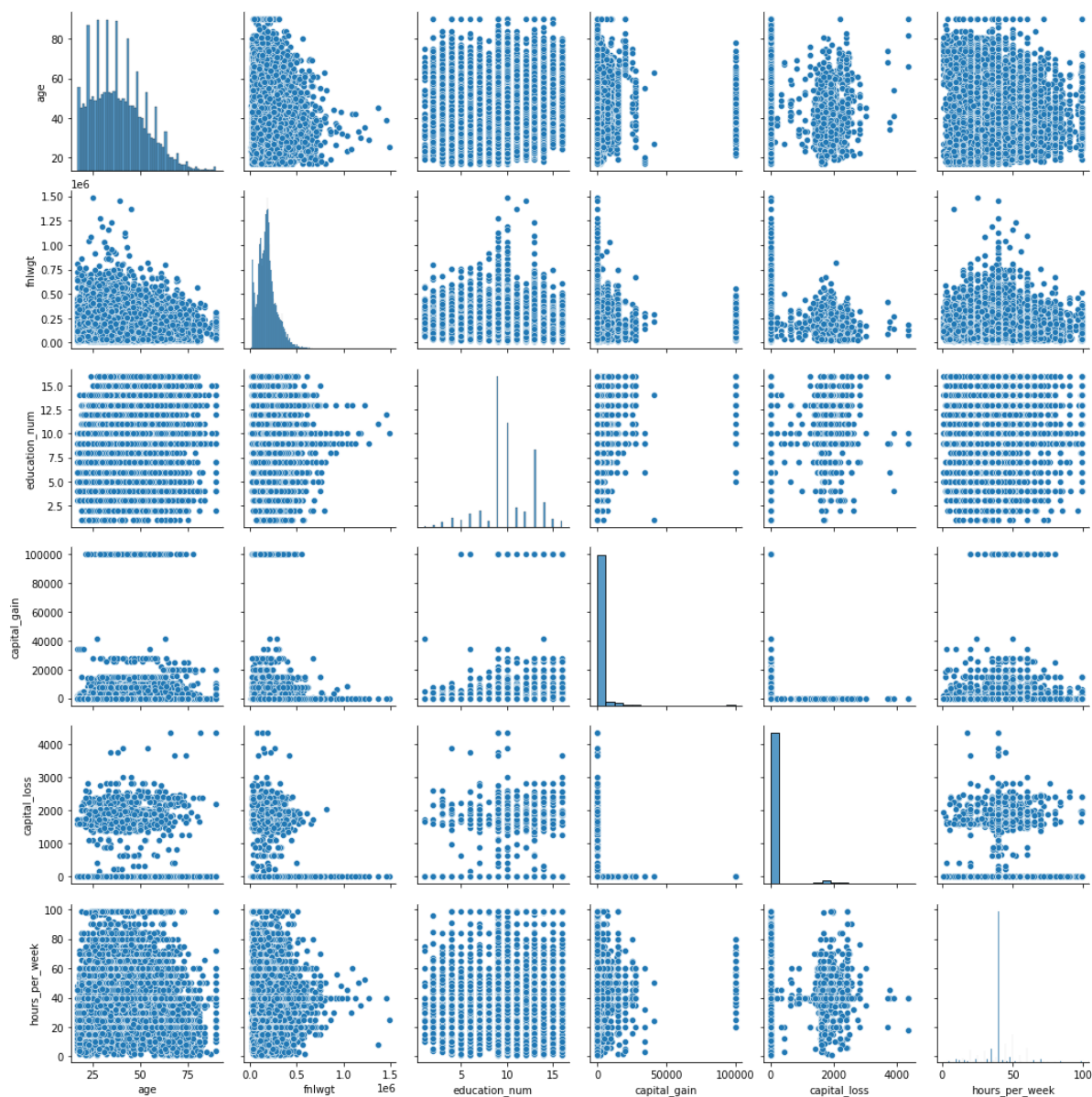
<=50K 24698

>50K 7839

Name: wage_class, dtype: int64

In [16]: sns.pairplot(train_set)

Out[16]: <seaborn.axisgrid.PairGrid at 0x254902a9488>



```
In [17]: train_set.education.value_counts()
```

```
Out[17]: HS-grad      10494
Some-college  7282
Bachelors    5353
Masters      1722
Assoc-voc    1382
11th         1175
Assoc-acdm   1067
10th         933
7th-8th      645
Prof-school  576
9th          514
12th         433
Doctorate    413
5th-6th      332
1st-4th      166
Preschool    50
Name: education, dtype: int64
```

```
In [18]: train_set.education_num.value_counts()
```

```
Out[18]: 9      10494
10     7282
13     5353
14     1722
11     1382
7      1175
12     1067
6       933
4       645
15      576
5       514
8       433
16      413
3       332
2       166
1        50
Name: education_num, dtype: int64
```

```
In [19]: ## As education and education_num are same, we are dropping education column
train_set.drop(columns=['education'], inplace=True)
test_set.drop(columns=['education'], inplace=True)
```

```
In [20]: train_set["education_num"]=train_set["education_num"].astype(CategoricalDtype(
ordered=True)) #ordinal data.
test_set["education_num"]=test_set["education_num"].astype(CategoricalDtype(ordered=True))
train_set["education_num"].head()
```

```
Out[20]: 0    13
1    13
2     9
3     7
4    13
Name: education_num, dtype: category
Categories (16, int64): [1 < 2 < 3 < 4 ... 13 < 14 < 15 < 16]
```

```
In [21]: train_set.education_num.value_counts()
```

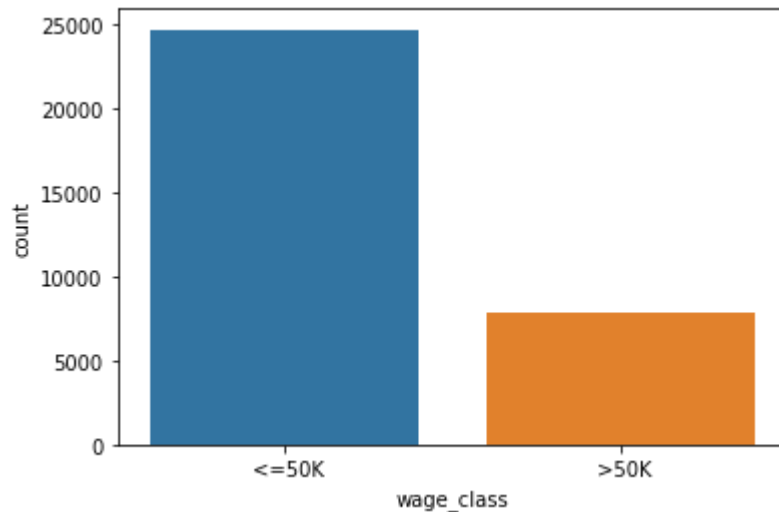
```
Out[21]: 9    10494
10    7282
13    5353
14    1722
11    1382
7     1175
12    1067
6     933
4     645
15    576
5     514
8     433
16    413
3     332
2     166
1       50
Name: education_num, dtype: int64
```

```
In [22]: sns.countplot(train_set.wage_class)
```

C:\Users\Urvi\AppData\Roaming\Python\Python37\site-packages\seaborn_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

FutureWarning

```
Out[22]: <AxesSubplot:xlabel='wage_class', ylabel='count'>
```



```
In [23]: train_set.skew()
```

```
Out[23]: age          0.557663  
fnlwgt        1.447703  
capital_gain  11.949403  
capital_loss   4.592702  
hours_per_week 0.228759  
dtype: float64
```

```
In [24]: plt.figure(figsize=(30,15))

plt.subplot(321)
sns.boxplot(train_set["age"])
plt.xticks(rotation=90)

plt.subplot(322)
sns.boxplot(train_set["capital_loss"])
plt.xticks(rotation=90)

plt.subplot(323)
sns.boxplot(train_set["capital_gain"])
plt.xticks(rotation=90)

plt.subplot(324)
sns.boxplot(train_set["fnlwt"])
plt.xticks(rotation=90)

plt.subplot(325)
sns.boxplot(train_set["hours_per_week"])
plt.xticks(rotation=90)

plt.subplots_adjust(hspace=0.5)
plt.show()
```


C:\Users\Urvi\AppData\Roaming\Python\Python37\site-packages\seaborn_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

FutureWarning

C:\Users\Urvi\AppData\Roaming\Python\Python37\site-packages\seaborn_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

FutureWarning

C:\Users\Urvi\AppData\Roaming\Python\Python37\site-packages\seaborn_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

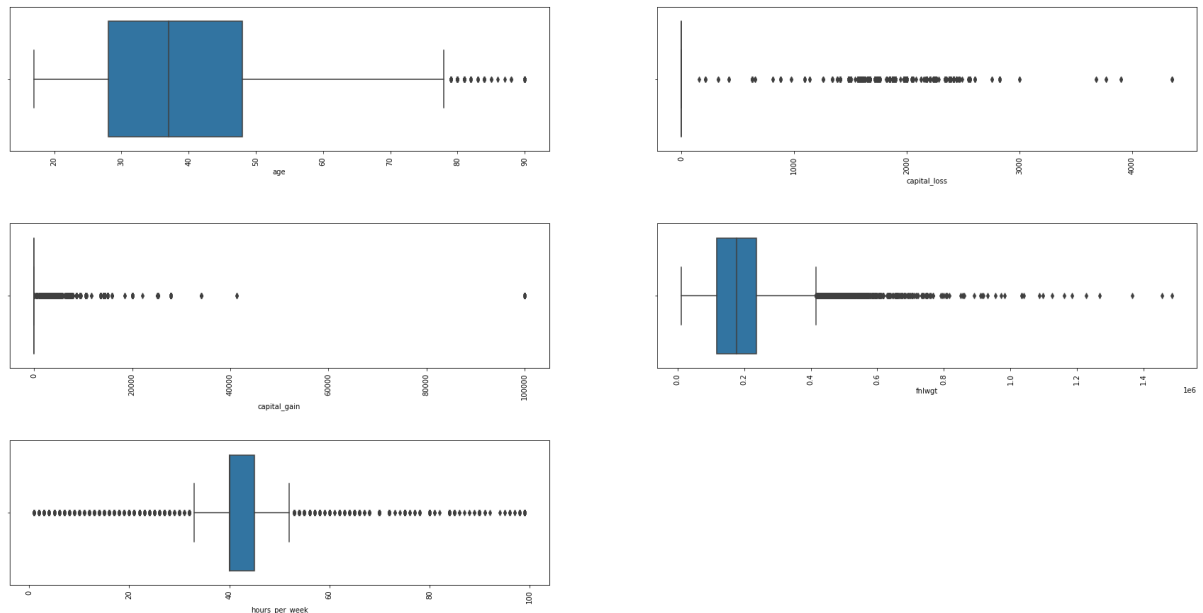
FutureWarning

C:\Users\Urvi\AppData\Roaming\Python\Python37\site-packages\seaborn_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

FutureWarning

C:\Users\Urvi\AppData\Roaming\Python\Python37\site-packages\seaborn_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

FutureWarning

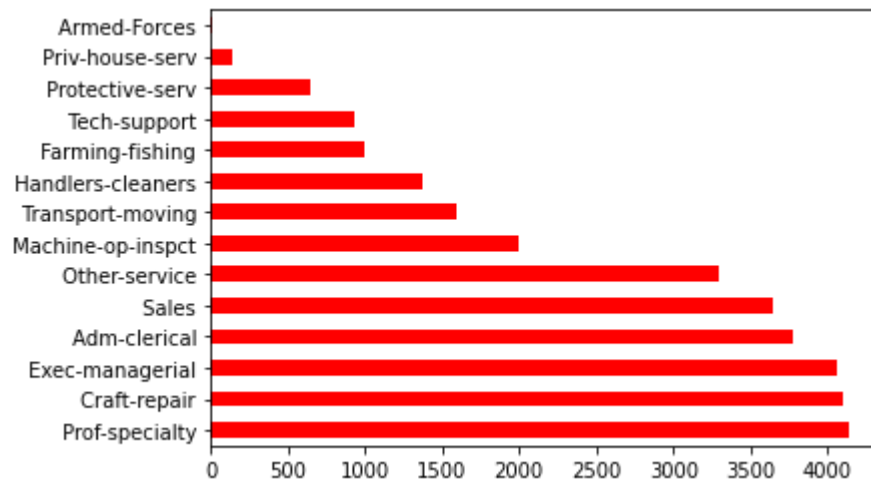


```
In [25]: numeric_columns=list(train_set.select_dtypes(include=["int64"]).columns) #numeric columns
numeric_columns
```

```
Out[25]: ['age', 'fnlwgt', 'capital_gain', 'capital_loss', 'hours_per_week']
```

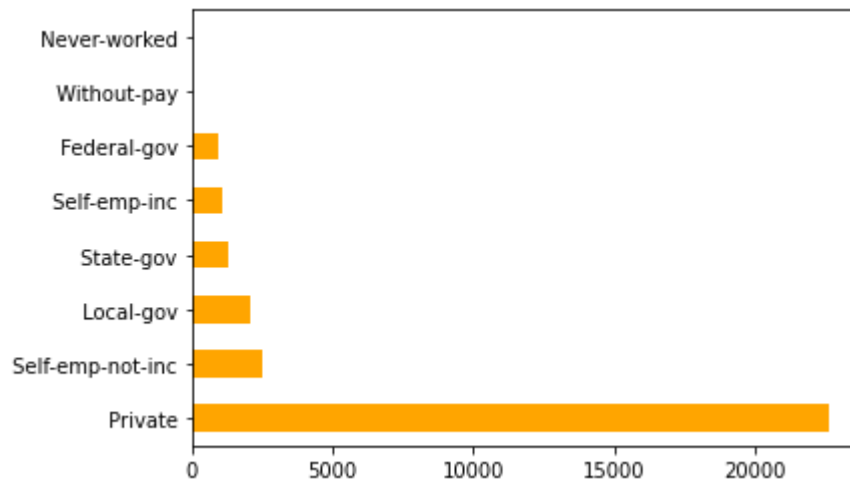
```
In [26]: train_set["occupation"].value_counts().plot.barh(color="red")
```

Out[26]: <AxesSubplot:>



```
In [27]: train_set["workclass"].value_counts().plot.barh(color="orange")
```

Out[27]: <AxesSubplot:>



```
In [28]: train_set.native_country.value_counts()
```

```
Out[28]: United-States      29153
Mexico      639
Philippines 198
Germany     137
Canada      121
Puerto-Rico 114
El-Salvador 106
India       100
Cuba        95
England     90
Jamaica     81
South       80
China       75
Italy       73
Dominican-Republic 70
Vietnam     67
Japan       62
Guatemala   62
Poland      60
Columbia    59
Taiwan      51
Haiti       44
Iran        43
Portugal    37
Nicaragua   34
Peru        31
France      29
Greece      29
Ecuador     28
Ireland     24
Hong        20
Cambodia    19
Trinidad&Tobago 19
Thailand     18
Laos        18
Yugoslavia  16
Outlying-US(Guam-USVI-etc) 14
Honduras    13
Hungary     13
Scotland    12
Holand-Netherlands 1
Name: native_country, dtype: int64
```

```
In [29]: for i in ["occupation", "workclass", "native_country"]:
          train_set[i].fillna(train_set[i].mode()[0], inplace=True)
          test_set[i].fillna(test_set[i].mode()[0], inplace=True)
```

```
In [30]: train_set.isnull().sum()
```

```
Out[30]: age                0  
workclass            0  
fnlwgt              0  
education_num       0  
marital_status      0  
occupation          0  
relationship        0  
race                0  
sex                 0  
capital_gain        0  
capital_loss        0  
hours_per_week      0  
native_country      0  
wage_class          0  
dtype: int64
```

```
In [31]: test_set.isnull().sum()
```

```
Out[31]: age                0  
workclass            0  
fnlwgt              0  
education_num       0  
marital_status      0  
occupation          0  
relationship        0  
race                0  
sex                 0  
capital_gain        0  
capital_loss        0  
hours_per_week      0  
native_country      0  
wage_class          0  
dtype: int64
```

```
In [32]: #Assigning the numeric values to the string type variables
number = LabelEncoder()
train_set['workclass'] = number.fit_transform(train_set['workclass'])
train_set['marital_status'] = number.fit_transform(train_set['marital_status'])
train_set['occupation'] = number.fit_transform(train_set['occupation'])
train_set['relationship'] = number.fit_transform(train_set['relationship'])
train_set['race'] = number.fit_transform(train_set['race'])
train_set['sex'] = number.fit_transform(train_set['sex'])
train_set['native_country'] = number.fit_transform(train_set['native_country'])
train_set['wage_class'] = number.fit_transform(train_set['wage_class'])

test_set['workclass'] = number.fit_transform(test_set['workclass'])
test_set['marital_status'] = number.fit_transform(test_set['marital_status'])
test_set['occupation'] = number.fit_transform(test_set['occupation'])
test_set['relationship'] = number.fit_transform(test_set['relationship'])
test_set['race'] = number.fit_transform(test_set['race'])
test_set['sex'] = number.fit_transform(test_set['sex'])
test_set['native_country'] = number.fit_transform(test_set['native_country'])
test_set['wage_class'] = number.fit_transform(test_set['wage_class'])
```

```
In [33]: train_set['age_bin'] = pd.cut(train_set['age'], 20)
test_set['age_bin'] = pd.cut(test_set['age'], 20)
```

```
In [34]: train_set['hours-per-week_bin'] = pd.cut(train_set['hours_per_week'], 10)
test_set['hours-per-week_bin'] = pd.cut(test_set['hours_per_week'], 10)
```

```
In [35]: train_set[['wage_class', 'age']].groupby(['wage_class'], as_index=False).mean()
.sort_values(by='age', ascending=False)
```

Out[35]:

	wage_class	age
1	1	44.250925
0	0	36.787392

```
In [36]: train_set = train_set.apply(LabelEncoder().fit_transform)
train_set.head()
```

Out[36]:

	age	workclass	fnlwgt	education_num	marital_status	occupation	relationship	race	sex	c
0	22	6	2671	12	4	0	1	4	1	
1	33	5	2926	12	2	3	0	4	1	
2	21	3	14086	8	0	5	1	4	1	
3	36	3	15336	6	2	5	0	2	1	
4	11	3	19355	12	2	9	5	2	0	

```
In [37]: test_set = test_set.apply(LabelEncoder().fit_transform)
test_set.head()
```

Out[37]:

	age	workclass	fnlwgt	education_num	marital_status	occupation	relationship	race	sex	c
0	8	3	8931	6	4	6	3	2	1	
1	21	3	1888	8	2	4	0	4	1	
2	11	1	11540	11	2	10	0	4	1	
3	27	3	5146	9	2	6	0	2	1	
4	1	3	2450	9	4	9	3	4	0	

```
In [38]: train_x = train_set.drop(columns=['wage_class', 'age', 'hours_per_week'])
test_x = test_set.drop(columns=['wage_class', 'age', 'hours_per_week'])
train_y = train_set.wage_class
test_y = test_set.wage_class
```

XGBoost Classifier

```
In [39]: # fit model no training data
model = XGBClassifier(objective='binary:logistic')
model.fit(train_x, train_y)
```

Out[39]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1, colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1, importance_type='gain', interaction_constraints='', learning_rate=0.300000012, max_delta_step=0, max_depth=6, min_child_weight=1, missing=nan, monotone_constraints='()', n_estimators=100, n_jobs=0, num_parallel_tree=1, objective='binary:logistic', random_state=0, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1, tree_method='exact', validate_parameters=1, verbosity=None)

```
In [40]: # cheking training accuracy
y_pred = model.predict(train_x)
predictions = [round(value) for value in y_pred]
accuracy = accuracy_score(train_y, predictions)
accuracy
```

Out[40]: 0.9019270369118234

```
In [41]: # cheking initial test accuracy
y_pred = model.predict(test_x)
predictions = [round(value) for value in y_pred]
accuracy = accuracy_score(test_y, predictions)
accuracy
```

Out[41]: 0.854939788645859

```
In [42]: param_grid={
          'learning_rate': [1, 0.5, 0.1, 0.01, 0.001],
          'max_depth': [3, 5, 10, 20],
          'n_estimators': [10, 50, 100, 200]
        }
```

```
In [43]: grid= GridSearchCV(XGBClassifier(objective='binary:logistic'), param_grid)
```

```
In [44]: grid.fit(train_x, train_y)
```

```
Out[44]: GridSearchCV(cv=None, error_score=nan,
                      estimator=XGBClassifier(base_score=None, booster=None,
                                              colsample_bylevel=None,
                                              colsample_bynode=None,
                                              colsample_bytree=None, gamma=None,
                                              gpu_id=None, importance_type='gain',
                                              interaction_constraints=None,
                                              learning_rate=None, max_delta_step=None,
                                              max_depth=None, min_child_weight=None,
                                              missing=nan, monotone_constraints=None,
                                              n_es...
                                              random_state=None, reg_alpha=None,
                                              reg_lambda=None, scale_pos_weight=None,
                                              subsample=None, tree_method=None,
                                              validate_parameters=None, verbosity=None),
                      iid='deprecated', n_jobs=None,
                      param_grid={'learning_rate': [1, 0.5, 0.1, 0.01, 0.001],
                                   'max_depth': [3, 5, 10, 20],
                                   'n_estimators': [10, 50, 100, 200]},
                      pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
                      scoring=None, verbose=0)
```

```
In [45]: grid.best_params_
```

```
Out[45]: {'learning_rate': 0.1, 'max_depth': 5, 'n_estimators': 200}
```

```
In [46]: # Create new model using the same parameters
          new_model=XGBClassifier(learning_rate= 0.1, max_depth= 5, n_estimators= 200)
          new_model.fit(train_x, train_y)
```

```
Out[46]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                      colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,
                      importance_type='gain', interaction_constraints='',
                      learning_rate=0.1, max_delta_step=0, max_depth=5,
                      min_child_weight=1, missing=nan, monotone_constraints='()',
                      n_estimators=200, n_jobs=0, num_parallel_tree=1,
                      objective='binary:logistic', random_state=0, reg_alpha=0,
                      reg_lambda=1, scale_pos_weight=1, subsample=1,
                      tree_method='exact', validate_parameters=1, verbosity=None)
```

```
In [47]: y_pred_new = new_model.predict(test_x)
         predictions_new = [round(value) for value in y_pred_new]
         accuracy_new_xgboost = accuracy_score(test_y, predictions_new)
         accuracy_new_xgboost
```

Out[47]: 0.8597321209142296

```
In [48]: new_model.get_booster().get_score(importance_type="gain")
```

```
Out[48]: {'relationship': 128.6458774135325,
          'education_num': 41.5177175380451,
          'capital_gain': 37.861184960465046,
          'capital_loss': 18.21865458632269,
          'age_bin': 15.75366145925717,
          'occupation': 9.55463262210465,
          'hours-per-week_bin': 12.643608384179403,
          'fnlwgt': 3.129041437650815,
          'marital_status': 38.846798557440955,
          'sex': 9.694870202,
          'race': 4.3695650333863645,
          'workclass': 5.375567432009281,
          'native_country': 3.365511446216217}
```

Features relationship, education_num, marital_status, capital_gain, capital_loss, age_bin, sex, hours-per-week_bin has more importance as compared to other variables

Project Done By: Urvi Gadda

mailto:urvigada96@gmail.com

```
In [ ]:
```