

▼ Census Income EDA

Prepared By:

Bankole Ayoade

GitHub: <https://github.com/acaaattunde2012>

LinkedIn: <https://www.linkedin.com/in/bankole-ayoade-fca-acti-cipfa-161406a7/>

▼ About Census Income Dataset

The dataset comprises 32561 rows and 15 columns as follows:

- income_class
 - | 50K, <=50K.
- age: continuous.
- workclass: Private, Self-emp-not-inc, Self-emp-inc, Federal-gov, Local-gov, State-gov, Without-pay, Never-worked.
- fnlwgt: continuous.
- education: Bachelors, Some-college, 11th, HS-grad, Prof-school, Assoc-acdm, Assoc-voc, 9th, 7th-8th, 12th, Masters, 1st-4th, * 10th, Doctorate, 5th-6th, Preschool.
- education-num: continuous.
- marital-status: Married-civ-spouse, Divorced, Never-married, Separated, Widowed, Married-spouse-absent, Married-AF-spouse.
- occupation: Tech-support, Craft-repair, Other-service, Sales, Exec-managerial, Prof-specialty, Handlers-cleaners, Machine-op-inspct, Adm-clerical, Farming-fishing, Transport-moving, Priv-house-serv, Protective-serv, Armed-Forces.
- relationship: Wife, Own-child, Husband, Not-in-family, Other-relative, Unmarried.
- race: White, Asian-Pac-Islander, Amer-Indian-Eskimo, Other, Black.
- sex: Female, Male.
- capital-gain: continuous.
- capital-loss: continuous.
- hours-per-week: continuous.
- native-country: United-States, Cambodia, England, Puerto-Rico, Canada, Germany, Outlying-US(Guam-USVI-etc), India, Japan, Greece, South, China, Cuba, Iran, Honduras, Philippines, Italy, Poland, Jamaica, Vietnam, Mexico, Portugal, Ireland, France, Dominican-Republic, Laos, Ecuador, Taiwan, Haiti, Columbia, Hungary, Guatemala, Nicaragua, Scotland, Thailand, Yugoslavia, El-Salvador, Trinadad&Tobago, Peru, Hong, Holand-Netherlands.

▼ Importing Dependencies

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import warnings
6 import plotly.express as px
7 import plotly.graph_objects as go
8 import plotly.io as pio
9 pio.templates.default = "plotly_white"
10
11 warnings.filterwarnings("ignore")
12
13 %matplotlib inline
```

▼ Data Ingestion:

```
1 url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/adult/adult.data'
```

```
1 #income_df=pd.read_csv(url, header=None)
2 income_df=pd.read_csv("adultdata.csv", header=None)
```

```
1 income_df.head(5)
```

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	0	40	United-States	<=50K
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	13	United-States	<=50K
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	40	United-States	<=50K
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0	13	United-States	<=50K

```
1 # Insert column names
2 col_name = ["age", "workclass", "fnlwgt", "education", "education-num", "marital_status", "occupation", "relationship", "race",
3             "sex", "capital_gain", "capital_loss", "hours_per_week", "native_country", "income_class"]
4 income_df.columns = col_name
5 income_df.head(5)
```

	age	workclass	fnlwgt	education	education-num	marital_status	occupation	relationship	race	sex	capital_gain	capital_loss	hours_per_week	native_country	income_class
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	0	16	United-States	<=50K
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	13	United-States	<=50K
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	13	United-States	<=50K
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0	13	United-States	<=50K
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	13	United-States	<=50K

```
1 # Make a copy of the dataset
2 income = income_df.copy()
```

▼ Data Profiling

```
1 # Check the first 5 rows
2 income.head(5)
```

	age	workclass	fnlwgt	education	education-num	marital_status	occupation	relationship	race	sex	capital_gain	capital_loss	hours_per_week	native_country	income_class
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	0	16	United-States	<=50K
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	13	United-States	<=50K
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	13	United-States	<=50K
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0	13	United-States	<=50K
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	13	United-States	<=50K

```
1 # Check the last 5 rows
2 income.tail(5)
```

	age	workclass	fnlwgt	education	education-num	marital-status	occupation	relationship	race	sex	capital_gain	capital_loss
32556	27	Private	257302	Assoc-acdm	12	Married-civ-spouse	Tech-support	Wife	White	Female	0	0
32557	40	Private	154374	HS-grad	9	Married-civ-spouse	Machine-op-inspect	Husband	White	Male	0	0

```

1 # Insert column names
2 col_name = ["age", "workclass", "fnlwgt", "education", "education-num", "marital_status", "occupation", "relationship", "race",
3             "sex", "capital_gain", "capital_loss", "hour_per_week", "native_country", "income_class"]
4 income.columns = col_name
5 income.head(5)

```

	age	workclass	fnlwgt	education	education-num	marital_status	occupation	relationship	race	sex	capital_gain	capital_loss
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	(
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	(
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	(
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	(
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	(

```

1 # Show the no of rows and columns
2 income.shape

```

(32561, 15)

```

1 # Show more info about the dataset
2 income.info()

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32561 entries, 0 to 32560
Data columns (total 15 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   age              32561 non-null   int64  
 1   workclass        32561 non-null   object  
 2   fnlwgt           32561 non-null   int64  
 3   education        32561 non-null   object  
 4   education-num    32561 non-null   int64  
 5   marital-status   32561 non-null   object  
 6   occupation       32561 non-null   object  
 7   relationship     32561 non-null   object  
 8   race              32561 non-null   object  
 9   sex               32561 non-null   object  
 10  capital_gain     32561 non-null   int64  
 11  capital_loss     32561 non-null   int64  
 12  hour_per_week    32561 non-null   int64  
 13  native_country   32561 non-null   object  
 14  income_class     32561 non-null   object  
dtypes: int64(6), object(9)
memory usage: 3.7+ MB

```

```

1 # Checking missing values
2 income.isnull().sum()

```

```

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
hour_per_week 0
native_country 0

```

```
income_class      0
dtype: int64
```

```
1 # Checking missing values
2 income.isna().sum()
```

```
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
hour_per_week    0
native_country   0
income_class     0
dtype: int64
```

```
1 # Checking for duplicates
2 income[income.duplicated()]
```

	age	workclass	fnlwgt	education	education-num	marital_status	occupation	relationship	race	sex	capital_gain	capital_loss
4881	25	Private	308144	Bachelors	13	Never-married	Craft-repair	Not-in-family	White	Male	0	
5104	90	Private	52386	Some-college	10	Never-married	Other-service	Not-in-family	Asian-Pac-Islander	Male	0	
8171	21	Private	250051	Some-college	10	Never-married	Prof-farm	Own-child	White	Female	0	

```
1 income.drop_duplicates(keep='first', inplace=True)
```

```
2 income[income.duplicated()]
```

	age	workclass	fnlwgt	education	education-num	marital_status	occupation	relationship	race	sex	capital_gain	capital_loss
4881	25	Private	308144	Bachelors	13	Never-married	Craft-repair	Not-in-family	White	Male	0	

```
1 # Checking for unique values for each column
```

```
2 for col in income.columns:
```

```
3     print(f"===== {col} =====")
```

```
4     print(income[col].unique())
```

```
5     print("=====")
```

```
6
```

```
=====age=====
[39 50 38 53 28 37 49 52 31 42 30 23 32 40 34 25 43 54 35 59 56 19 20 45
 22 48 21 24 57 44 41 29 18 47 46 36 79 27 67 33 76 17 55 61 70 64 71 68
 66 51 58 26 60 90 75 65 77 62 63 80 72 74 69 73 81 78 88 82 83 84 85 86
 87]

=====workclass=====
[' State-gov' ' Self-emp-not-inc' ' Private' ' Federal-gov' ' Local-gov'
 ' ?' ' Self-emp-inc' ' Without-pay' ' Never-worked']

=====fnlwgt=====
[ 77516  83311 215646 ... 34066 84661 257302]

=====education=====
[' Bachelors' ' HS-grad' ' 11th' ' Masters' ' 9th' ' Some-college'
 ' Assoc-acdm' ' Assoc-voc' ' 7th-8th' ' Doctorate' ' Prof-school'
 ' 5th-6th' ' 10th' ' 1st-4th' ' Preschool' ' 12th']

=====education-num=====
[13 9 7 14 5 10 12 11 4 16 15 3 6 2 1 8]

=====marital_status=====
[' Never-married' ' Married-civ-spouse' ' Divorced'
 ' Married-spouse-absent' ' Separated' ' Married-AF-spouse' ' Widowed']

=====occupation=====
[' Adm-clerical' ' Exec-managerial' ' Handlers-cleaners' ' Prof-specialty'
 ' Other-service' ' Sales' ' Craft-repair' ' Transport-moving'
 ' Farming-fishing' ' Machine-op-inspct' ' Tech-support' ' ?'
 ' Protective-serv' ' Armed-Forces' ' Priv-house-serv']

=====relationship=====
[' Not-in-family' ' Husband' ' Wife' ' Own-child' ' Unmarried'
 ' Other-relative']

=====race=====
[' White' ' Black' ' Asian-Pac-Islander' ' Amer-Indian-Eskimo' ' Other']

=====sex=====
[' Male' ' Female']

=====capital_gain=====
[ 2174    0 14084  5178  5013  2407 14344 15024  7688 34095  4064  4386
 7298 1409  3674 1055  3464 2050 2176  594 20051  6849 4101 1111
 8614 3411 2597 25236 4650 9386 2463 3103 10605 2964 3325 2580
 3471 4865 99999 6514 1471 2329 2105 2885 25124 10520 2202 2961
 27828 6767 2228 1506 13550 2635 5556 4787 3781 3137 3818 3942
 914 401 2829 2977 4934 2062 2354 5455 15020 1424 3273 22040
 4416 3908 10566 991 4931 1086 7430 6497 114 7896 2346 3418
 3432 2907 1151 2414 2290 15831 41310 4588 2538 3456 6418 1848
 3887 5721 9562 1455 2036 1831 11678 2936 2993 7443 6360 1797
 1173 4687 6723 2009 6097 2653 1639 18481 7978 2387 5060]

=====capital_loss=====
[ 0 2042 1408 1902 1573 1887 1719 1762 1564 2179 1816 1980 1977 1876
 1340 2206 1741 1485 2339 2415 1380 1721 2051 2377 1669 2352 1672 653
 2392 1504 2001 1590 1651 1628 1848 1740 2002 1579 2258 1602 419 2547
 2174 2205 1726 2444 1138 2238 625 213 1539 880 1668 1092 1594 3004]


```

```
1 # Separate both the dataset columns into numeric and categorical columns
```

```
2 numeric_cols = [feature for feature in income.columns if income[feature].dtype != 'O']
```

```

3 categorical_cols =[feature for feature in income.columns if income[feature].dtype == 'O']
1 print('Categorical Features')
2 print(categorical_cols)
3 print('=====')
4 print('Numeric Features')
5 print(numeric_cols)

Categorical Features
['workclass', 'education', 'marital_status', 'occupation', 'relationship', 'race', 'sex', 'native_country', 'income_class']
=====
Numeric Features
['age', 'fnlwgt', 'education-num', 'capital_gain', 'capital_loss', 'hour_per_week']

1 # Removing trailing spaces from all the categorical columns
2 for col in categorical_cols:
3     income[col]=income[col].str.strip()
4     print(income[col].unique())

['State-gov' 'Self-emp-not-inc' 'Private' 'Federal-gov' 'Local-gov' '?'
 'Self-emp-inc' 'Without-pay' 'Never-worked']
['Bachelors' 'HS-grad' '11th' 'Masters' '9th' 'Some-college' 'Assoc-acdm'
 'Assoc-voc' '7th-8th' 'Doctorate' 'Prof-school' '5th-6th' '10th'
 '1st-4th' 'Preschool' '12th']
['Never-married' 'Married-civ-spouse' 'Divorced' 'Married-spouse-absent'
 'Separated' 'Married-AF-spouse' 'Widowed']
['Adm-clerical' 'Exec-managerial' 'Handlers-cleaners' 'Prof-specialty'
 'Other-service' 'Sales' 'Craft-repair' 'Transport-moving'
 'Farming-fishing' 'Machine-op-inspct' 'Tech-support' '?'
 'Protective-serv' 'Armed-Forces' 'Priv-house-serv']
['Not-in-family' 'Husband' 'Wife' 'Own-child' 'Unmarried' 'Other-relative']
['White' 'Black' 'Asian-Pac-Islander' 'Amer-Indian-Eskimo' 'Other']
['Male' 'Female']
['United-States' 'Cuba' 'Jamaica' 'India' '?' 'Mexico' 'South'
 'Puerto-Rico' 'Honduras' 'England' 'Canada' 'Germany' 'Iran'
 'Philippines' 'Italy' 'Poland' 'Columbia' 'Cambodia' 'Thailand' 'Ecuador'
 'Laos' 'Taiwan' 'Haiti' 'Portugal' 'Dominican-Republic' 'El-Salvador'
 'France' 'Guatemala' 'China' 'Japan' 'Yugoslavia' 'Peru'
 'Outlying-US(Guam-USVI-etc)' 'Scotland' 'Trinidad&Tobago' 'Greece'
 'Nicaragua' 'Vietnam' 'Hong' 'Ireland' 'Hungary' 'Holand-Netherlands']
['<=50K' '>50K']

1 # Cleaning up the some categorical columns by replacing Special xter with Most frequent
2 income['workclass'] = income['workclass'].astype(str)
3 income['occupation'] = income['occupation'].astype(str)
4 income['native_country'] = income['native_country'].astype(str)
5 income['education']=income['education'].astype(str)
6 income['workclass'] = income['workclass'].str.replace('?', "Private")
7 income['occupation'] = income['occupation'].str.replace('?', "Prof-specialty")
8 income['native_country'] = income['native_country'].str.replace('?', "United-States")
9 income['education'].replace({'11th':'HS-grad", "9th":"HS-grad", "7th-8th":"HS-grad", "1st-4th":"HS-grad", "5th-6th":"HS-grad", "10th":"H

1 income['workclass'].unique()

array(['State-gov', 'Self-emp-not-inc', 'Private', 'Federal-gov',
       'Local-gov', 'Self-emp-inc', 'Without-pay', 'Never-worked'],
      dtype=object)

1 income['occupation'].unique()

array(['Adm-clerical', 'Exec-managerial', 'Handlers-cleaners',
       'Prof-specialty', 'Other-service', 'Sales', 'Craft-repair',
       'Transport-moving', 'Farming-fishing', 'Machine-op-inspct',
       'Tech-support', 'Protective-serv', 'Armed-Forces',
       'Priv-house-serv'], dtype=object)

1 income['education'].unique()

array(['Bachelors', 'HS-grad', 'Masters', 'Some-college', 'Assoc-acdm',
       'Assoc-voc', 'Doctorate', 'Prof-school', 'Preschool'], dtype=object)

1 income['native_country'].unique()

array(['United-States', 'Cuba', 'Jamaica', 'India', 'Mexico', 'South',
       'Puerto-Rico', 'Honduras', 'England', 'Canada', 'Germany', 'Iran',
       'Philippines', 'Italy', 'Poland', 'Columbia', 'Cambodia',
       'Thailand', 'Ecuador', 'Laos', 'Taiwan', 'Haiti', 'Portugal',
       'Dominican-Republic', 'El-Salvador', 'France', 'Guatemala',
       'China', 'Japan', 'Yugoslavia', 'Peru',
       'Outlying-US(Guam-USVI-etc)', 'Scotland', 'Trinidad&Tobago',
       'Greece', 'Nicaragua', 'Vietnam', 'Hong', 'Ireland', 'Hungary',
       'Holand-Netherlands'], dtype=object)

```

```
1 # Statistical analysis about the dataset
2 income.describe()
```

	age	fnlwgt	education-num	capital_gain	capital_loss	hour_per_week
count	32537.000000	3.253700e+04	32537.000000	32537.000000	32537.000000	32537.000000
mean	38.585549	1.897808e+05	10.081815	1078.443741	87.368227	40.440329
std	13.637984	1.055565e+05	2.571633	7387.957424	403.101833	12.346889
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	2.369930e+05	12.000000	0.000000	0.000000	45.000000
max	90.000000	1.484705e+06	16.000000	99999.000000	4356.000000	99.000000

```
1 income.describe(include='all').T
```

	count	unique	top	freq	mean	std	min	25%	50%	75%	max
age	32537.0	NaN		NaN	38.585549	13.637984	17.0	28.0	37.0	48.0	90.0
workclass	32537	8	Private	24509	NaN	NaN	NaN	NaN	NaN	NaN	NaN
fnlwgt	32537.0	NaN		NaN	189780.848511	105556.471009	12285.0	117827.0	178356.0	236993.0	1484705.0
education	32537	9	HS-grad	14692	NaN	NaN	NaN	NaN	NaN	NaN	NaN
education-num	32537.0	NaN		NaN	10.081815	2.571633	1.0	9.0	10.0	12.0	16.0
marital-status	32537	7	Married-civ-spouse	14970	NaN	NaN	NaN	NaN	NaN	NaN	NaN
occupation	32537	14	Prof-specialty	5979	NaN	NaN	NaN	NaN	NaN	NaN	NaN
relationship	32537	6	Husband	13187	NaN	NaN	NaN	NaN	NaN	NaN	NaN
race	32537	5	White	27795	NaN	NaN	NaN	NaN	NaN	NaN	NaN
sex	32537	2	Male	21775	NaN	NaN	NaN	NaN	NaN	NaN	NaN
capital_gain	32537.0	NaN		NaN	1078.443741	7387.957424	0.0	0.0	0.0	0.0	99999.0
capital_loss	32537.0	NaN		NaN	87.368227	403.101833	0.0	0.0	0.0	0.0	4356.0
hour_per_week	32537.0	NaN		NaN	40.440329	12.346889	1.0	40.0	40.0	45.0	99.0
native_country	32537	41	United-States	29735	NaN	NaN	NaN	NaN	NaN	NaN	NaN
income_class	32537	2	<=50K	24698	NaN	NaN	NaN	NaN	NaN	NaN	NaN

```
1 # Checking the correlation among the numeric columns
2 income[numeric_cols]
```

	age	fnlwgt	education-num	capital_gain	capital_loss	hour_per_week
0	39	77516	13	2174	0	40
1	50	83311	13	0	0	13
2	38	215646	9	0	0	40
3	53	234721	7	0	0	40
4	28	338409	13	0	0	40
...
32556	27	257302	12	0	0	38
32557	40	154374	9	0	0	40
32558	58	151910	9	0	0	40
32559	22	201490	9	0	0	20
32560	52	287927	9	15024	0	40

32537 rows × 6 columns

```
1 # Correlation among the numeric columns
2 income[numeric_cols].corr()
```

	age	fnlwgt	education-num	capital_gain	capital_loss	hour_per_week
age	1.000000	-0.076447	0.036224	0.077676	0.057745	0.068515
fnlwgt	-0.076447	1.000000	-0.043388	0.000429	-0.010260	-0.018898
education-num	0.036224	-0.043388	1.000000	0.122664	0.079892	0.148422
capital_gain	0.077676	0.000429	0.122664	1.000000	-0.031639	0.078408
capital_loss	0.057745	0.010260	0.079892	0.031639	1.000000	0.051020

```
1 # Checking missing values again after columns cleaning
```

```
2 income.isnull().sum()
```

```
3 #income['workclass'].isnull().sum()
```

```
4 #income[income['workclass'] == str(np.nan)]
```

```
5
```

```
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
hour_per_week 0
native_country 0
income_class  0
dtype: int64
```

```
1 for col in categorical_cols:
```

```
2     print(f"{col} ")
```

```
3     print(f"{income[col].value_counts(normalize=True)*100}")
```

```
4     print("====")
```

```
workclass
Private      75.326551
Self-emp-not-inc 7.806497
Local-gov    6.432677
State-gov    3.989304
Self-emp-inc 3.429941
Federal-gov 2.950487
Without-pay   0.043028
Never-worked 0.021514
Name: workclass, dtype: float64
=====
```

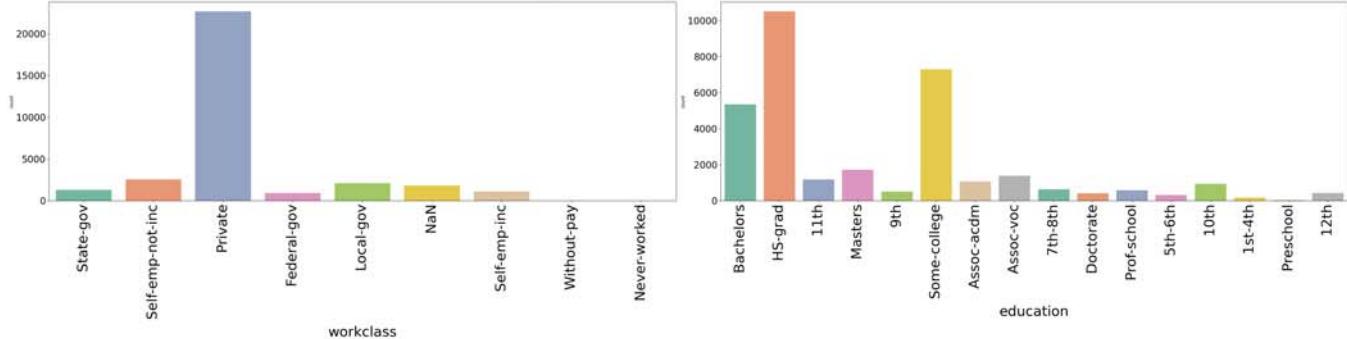
```
education
HS-grad      45.154747
Some-college 22.380674
Bachelors    16.452039
Masters      5.292436
Assoc-voc    4.247472
Assoc-acdm   3.279344
Prof-school  1.770292
Doctorate    1.269324
Preschool    0.153671
Name: education, dtype: float64
=====
```

```
marital_status
Married-civ-spouse 46.009159
Never-married     32.784215
Divorced          13.649076
Separated         3.150260
Widowed           3.051910
Married-spouse-absent 1.284691
Married-AF-spouse 0.070689
Name: marital_status, dtype: float64
=====
```

```
occupation
Prof-specialty 18.376003
Craft-repair    12.582598
Exec-managerial 12.493469
Adm-clerical   11.580662
Sales          11.217998
Other-service   10.114639
Machine-op-inspct 6.146848
Transport-moving 4.908258
Handlers-cleaners 4.207518
Farming-fishing 3.048837
Tech-support    2.849064
Protective-serv 1.994652
Priv-house-serv 0.451793
Armed-Forces    0.027661
Name: occupation, dtype: float64
```

```
=====
relationship
Husband          40.529244
Not-in-family    25.484833
Own-child        15.563820
Unmarried         10.587946
Wife              4.819129
Other-relative    3.015029
Name: relationship, dtype: float64
```

```
1 plt.figure(figsize=(40, 60))
2 plt.suptitle('Univariate Analysis of Categorical Features: Count Plots', fontsize=40, fontweight='bold', alpha=0.8, y=1.)
3 for i in range(0, len(categorical_cols)):
4     plt.subplot(5, 2, i+1)
5     sns.countplot(x=income[categorical_cols[i]], palette="Set2")
6     plt.xlabel(categorical_cols[i], fontsize = 30)
7     plt.xticks(rotation=90, fontsize = 30)
8     plt.yticks(fontsize = 20)
9     plt.tight_layout()
10
```

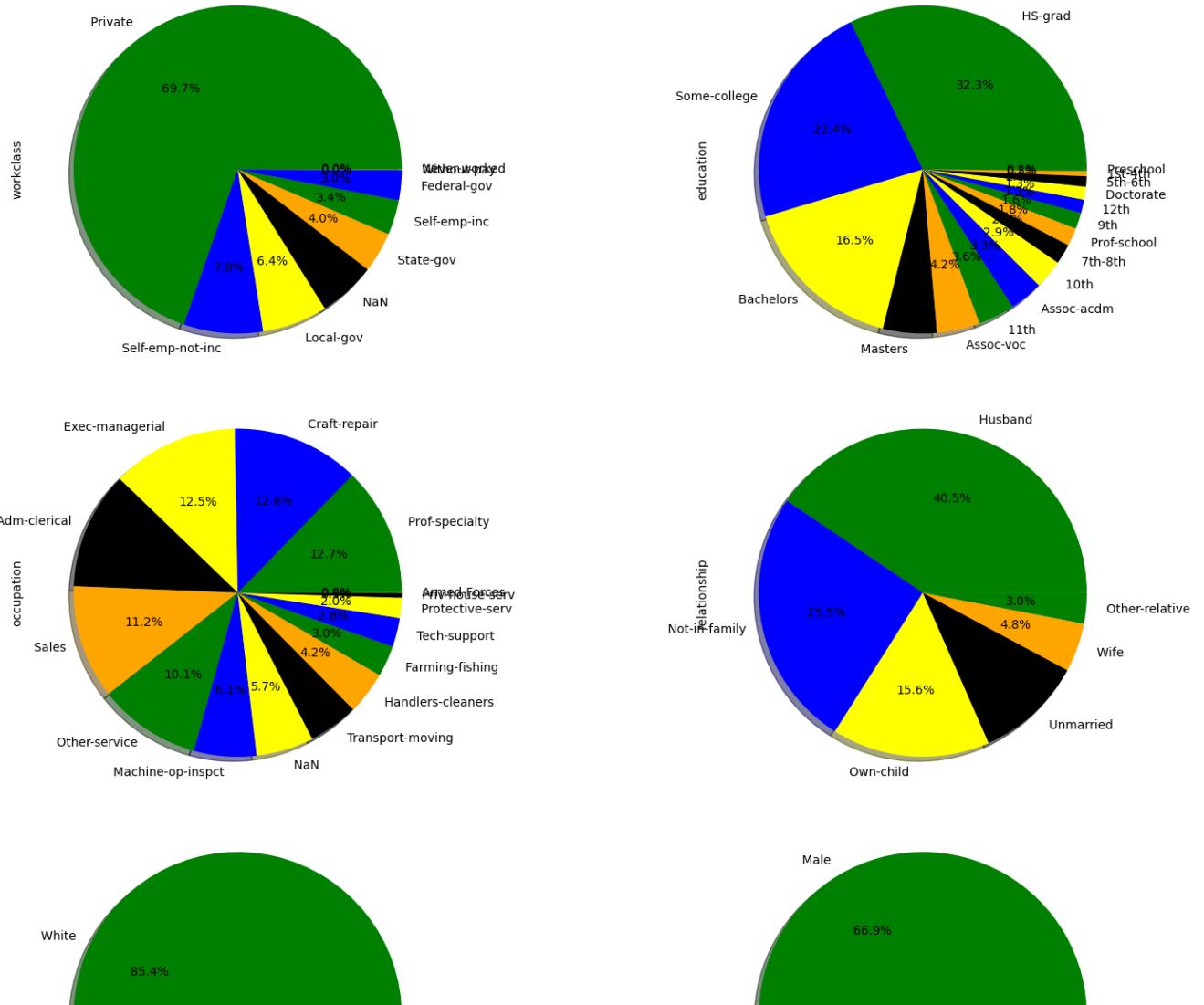
Univariate Analysis of Categorical Features: Count Plots

```

1 plt.figure(figsize=(15, 20))
2 plt.suptitle('Univariate Analysis of Categorical Features: ', fontsize=20, fontweight='bold', alpha=0.8, y=1.)
3 plt.subplot(421)
4 income['workclass'].value_counts().plot.pie(y=income['workclass'], autopct='%1.1f%%', shadow=True, colors=['green', 'blue', 'yellow', 'black'])
5 plt.subplot(422)
6 income['education'].value_counts().plot.pie(y=income['education'], autopct='%1.1f%%', shadow=True, colors=['green', 'blue', 'yellow', 'black'])
7 plt.subplot(423)
8 income['occupation'].value_counts().plot.pie(y=income['occupation'], autopct='%1.1f%%', shadow=True, colors=['green', 'blue', 'yellow', 'black'])
9 plt.subplot(424)
10 income['relationship'].value_counts().plot.pie(y=income['relationship'], autopct='%1.1f%%', shadow=True, colors=['green', 'blue', 'yellow'])
11 plt.subplot(425)
12 income['race'].value_counts().plot.pie(y=income['race'], autopct='%1.1f%%', shadow=True, colors=['green', 'blue', 'yellow', 'black', 'orange'])
13 plt.subplot(426)
14 income['sex'].value_counts().plot.pie(y=income['sex'], autopct='%1.1f%%', shadow=True, colors=['green', 'blue', 'yellow', 'black', 'orange'])
15 plt.subplot(427)
16 income['native_country'].value_counts().plot.pie(y=income['native_country'], autopct='%1.1f%%', shadow=True, colors=['green', 'blue', 'yellow'])
17 plt.tight_layout()
18 plt.show()
19

```

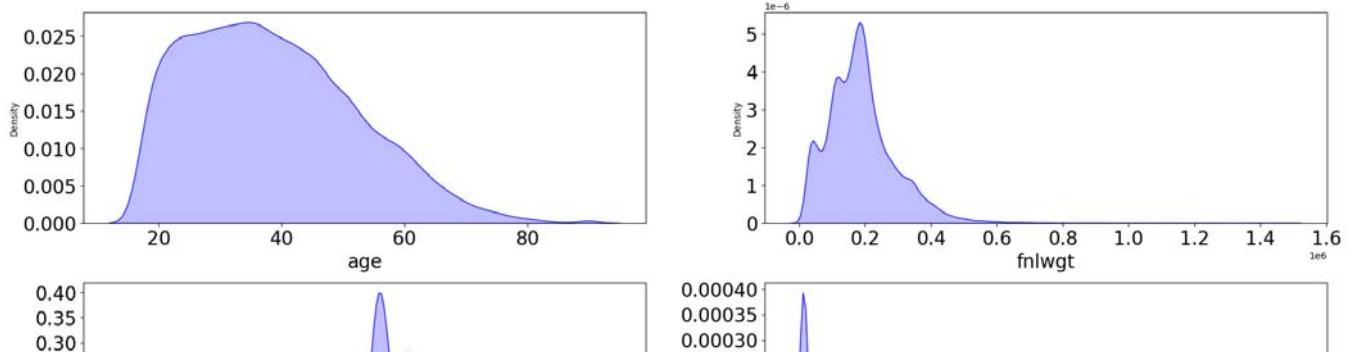
Univariate Analysis of Categorical Features:



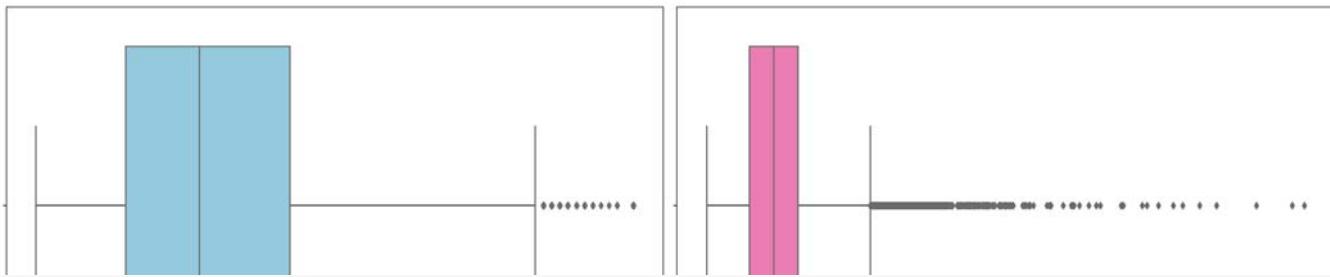
```

1 plt.figure(figsize=(20, 20))
2 plt.suptitle('Univariate Analysis of Numerical Features: Data Distribution', fontsize=20, fontweight='bold', alpha=0.8, y=1.)
3
4 for i in range(0, len(numeric_cols)):
5     plt.subplot(5, 2, i+1)
6     sns.kdeplot(x=income[numeric_cols[i]], shade=True, color='b')
7     plt.xlabel(numeric_cols[i], fontsize = 20)
8     plt.xticks(fontsize = 20)
9     plt.yticks(fontsize = 20)
10    plt.tight_layout()

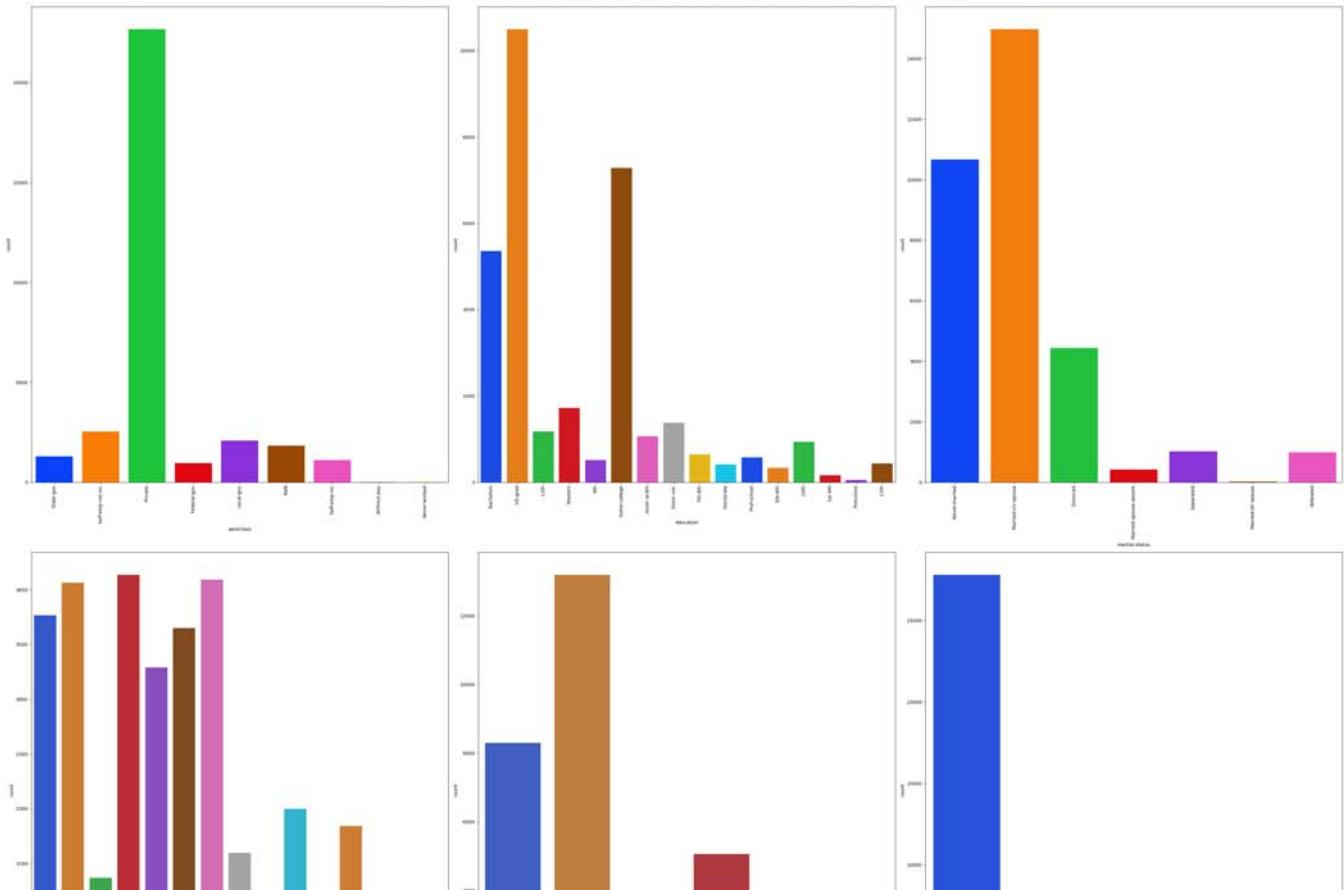
```

Univariate Analysis of Numerical Features: Data Distribution

```
1
2 plt.subplots(3,2,figsize=(20,20))
3 plt.suptitle('Univariate Analysis of Numerical Features: Data Distribution And Outlier Detection', fontsize=20, fontweight='bold', alpha=0.8)
4 plt.subplot(321)
5 sns.boxplot(income['age'],color='skyblue')
6 plt.subplot(322)
7 sns.boxplot(income['fnlwgt'],color='hotpink')
8 plt.subplot(323)
9 sns.boxplot(income['education-num'],color='yellow')
10 plt.subplot(324)
11 sns.boxplot(income['capital_gain'],color='lightgreen')
12 plt.subplot(325)
13 sns.boxplot(income['capital_loss'],color='lightblue')
14 plt.subplot(326)
15 sns.boxplot(income['hour_per_week'],color='orange')
16 plt.tight_layout()
17 plt.show()
18
```

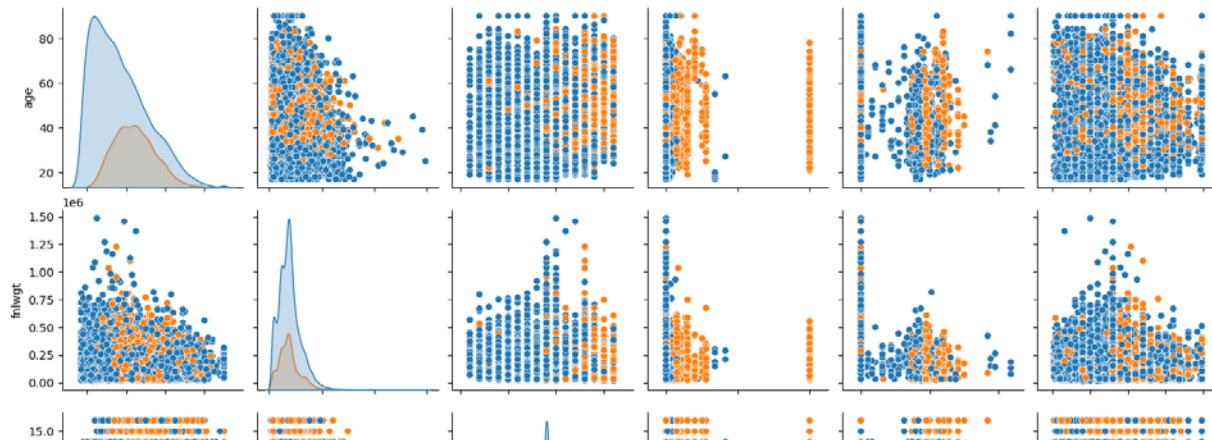
Univariate Analysis of Numerical Features: Data Distribution And Outlier Detection

```
1 #
2 plt.subplots(3,3,figsize=(40,50))
3 plt.suptitle('Univariate Analysis of Categorical Features: ', fontsize=50, fontweight='bold', alpha=0.8, y=1.)
4 plt.xticks(rotation=90)
5 plt.subplot(331)
6 sns.countplot(x=income['workclass'],data=income,palette = 'bright',saturation=0.95)
7 plt.xticks(rotation=90)
8 plt.subplot(332)
9 sns.countplot(x=income['education'],data=income,palette = 'bright',saturation=0.80)
10 plt.xticks(rotation=90)
11 plt.subplot(333)
12 sns.countplot(x=income['marital-status'],data=income,palette = 'bright',saturation=0.90)
13 plt.xticks(rotation=90)
14 plt.subplot(334)
15 sns.countplot(x=income['occupation'],data=income,palette = 'bright',saturation=0.60)
16 plt.xticks(rotation=90)
17 plt.subplot(335)
18 sns.countplot(x=income['relationship'],data=income,palette = 'bright',saturation=0.50)
19 plt.xticks(rotation=90)
20 plt.subplot(336)
21 sns.countplot(x=income['race'],data=income,palette = 'bright',saturation=0.70)
22 plt.xticks(rotation=90)
23 plt.subplot(337)
24 sns.countplot(x=income['sex'],data=income,palette = 'bright',saturation=0.85)
25 plt.xticks(rotation=90)
26 plt.subplot(338)
27 sns.countplot(x=income['native_country'],data=income,palette = 'bright',saturation=0.40)
28 plt.xticks(rotation=90)
29 plt.tight_layout()
30 plt.show()
31
```

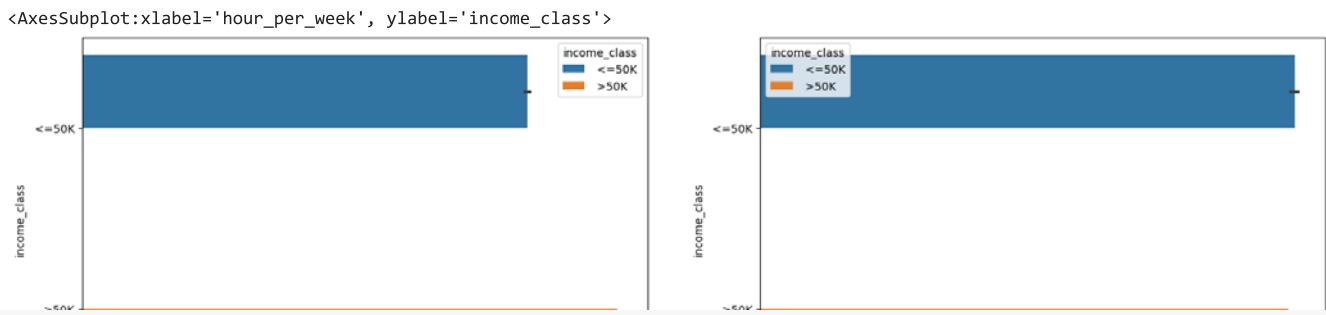
Univariate Analysis of Categorical Features:

```
1 # Correlation among numeric features
2 plt.figure(figsize=(40, 40))
3 plt.title('MultiVariate Analysis: Feature Correlation', fontsize=20, fontweight='bold', alpha=0.8, y=1.)
4 sns.pairplot(income, hue='income_class')
5 plt.tight_layout()
6 plt.show()
7
```

```
<seaborn.axisgrid.PairGrid at 0x2a110a01670>
<Figure size 4000x4000 with 0 Axes>
```



```
1 # Checkin the relationship between numeric columns and the label feature
2 plt.figure(figsize=(20,20))
3 plt.subplot(3,2,1)
4 sns.barplot(x=income['age'], y=income['income_class'], hue=income['income_class'])
5 plt.subplot(3,2,2)
6 sns.barplot(x=income['fnlwgt'], y=income['income_class'], hue=income['income_class'])
7 plt.subplot(3,2,3)
8 sns.barplot(x=income['education-num'], y=income['income_class'], hue=income['income_class'])
9 plt.subplot(3,2,4)
10 sns.barplot(x=income['capital_gain'], y=income['income_class'], hue=income['income_class'])
11 plt.subplot(3,2,5)
12 sns.barplot(x=income['capital_loss'], y=income['income_class'], hue=income['income_class'])
13 plt.subplot(3,2,6)
14 sns.barplot(x=income['hour_per_week'], y=income['income_class'], hue=income['income_class'])
15
16
```



```
1 # Checkin the relationship between numeric columns and the label feature
2 plt.figure(figsize=(20,20))
3 plt.subplot(3,2,1)
4 sns.boxplot (y=income['age'], x=income['income_class'], hue=income['income_class'])
5 plt.subplot(3,2,2)
6 sns.boxplot (y=income['fnlwgt'], x=income['income_class'], hue=income['income_class'])
7 plt.subplot(3,2,3)
8 sns.boxplot (y=income['education-num'], x=income['income_class'], hue=income['income_class'])
9 plt.subplot(3,2,4)
10 sns.boxplot (y=income['capital_gain'], x=income['income_class'], hue=income['income_class'])
11 plt.subplot(3,2,5)
12 sns.boxplot (y=income['capital_loss'], x=income['income_class'], hue=income['income_class'])
13 plt.subplot(3,2,6)
14 sns.boxplot (y=income['hour_per_week'], x=income['income_class'], hue=income['income_class'])
```

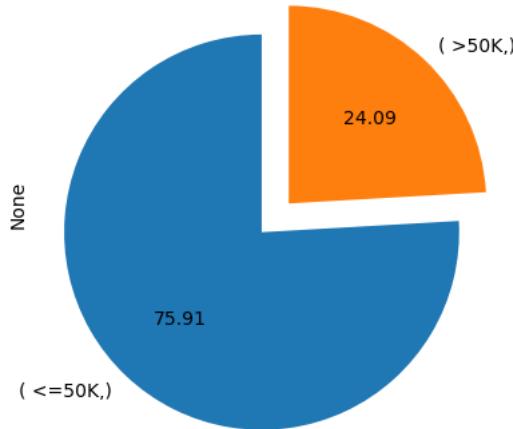
<AxesSubplot:xlabel='income_class', ylabel='hour_per_week'>

1 # Proportion of people that earn more than 50k

2 income.value_counts(['income_class']).plot.pie(y='income_class', startangle=90, explode=(0.2,0), title='Proportions of income class',)

<AxesSubplot:title={'center':'Proportions of income class'}, ylabel='None'>

Proportions of income class



```
1 # Top 5 workclass earning highest hour per week  
2 income.groupby(['workclass'])['hour_per_week'].sum().sort_values(ascending = False).reset_index()  
3
```

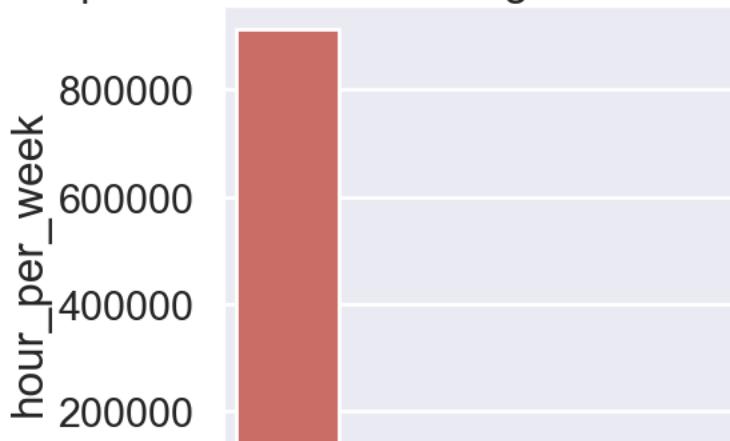
	workclass	hour_per_week
0	Private	913065
1	Self-emp-not-inc	112836
2	Local-gov	85777
3	NaN	58604
4	Self-emp-inc	54481
5	State-gov	50663
6	Federal-gov	39724
7	Without-pay	458
8	Never-worked	199

```
1 plt.figure(figsize=(5,5))
2 Top_5 = income.groupby(['workclass'])['hour_per_week'].sum().sort_values(ascending = False).reset_index()[:4]
3 #sns.set_context("poster")
4 sns.set_style("darkgrid")
5 plt.title('Top 5 Workclass with Highest Hour Per Week')
6 sns.barplot(data = Top_5, y='hour_per_week', x='workclass', palette='hls')
7 plt.xticks(rotation=90)
```



```
(array([0, 1, 2, 3]),
 [Text(0, 0, ' Private'),
  Text(1, 0, ' Self-emp-not-inc'),
  Text(2, 0, ' Local-gov'),
  Text(3, 0, ' NaN')])
```

Top 5 Workclass with Highest Hour Per Week

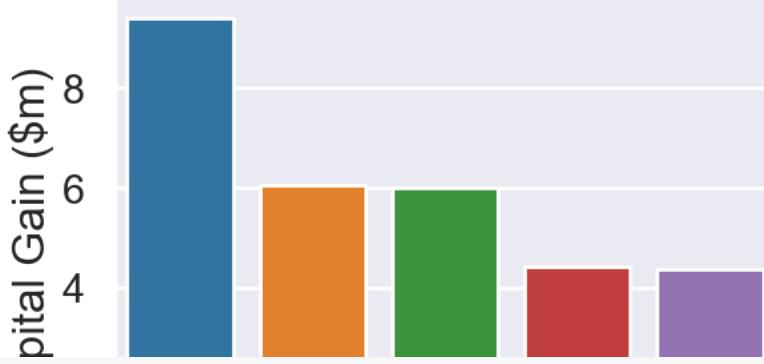


```
1 # Top 5 Education class that pay highest capital gain
2 income.groupby(['education'])['capital_gain'].sum().sort_values(ascending=False).reset_index()
```

	education	capital_gain
0	Bachelors	9404984
1	HS-grad	6056978
2	Prof-school	5998704
3	Masters	4415297
4	Some-college	4366027
5	Doctorate	1970070
6	Assoc-voc	988201
7	Assoc-acdm	683306
8	10th	377468
9	11th	252740
10	9th	175834
11	7th-8th	151125
12	12th	123010
13	5th-6th	58615
14	Preschool	45818
15	1st-4th	21147

```
1 Top_5_Cgain=income.groupby(['education'])['capital_gain'].sum().sort_values(ascending=False).reset_index()[:5]
2 plt.title('Top 5 Education class that pay highest capital gain')
3 sns.barplot(data=Top_5_Cgain, x = 'education', y='capital_gain')
4 plt.xticks(rotation=45)
5 plt.ylabel('Capital Gain ($m)')
6 plt.show()
```

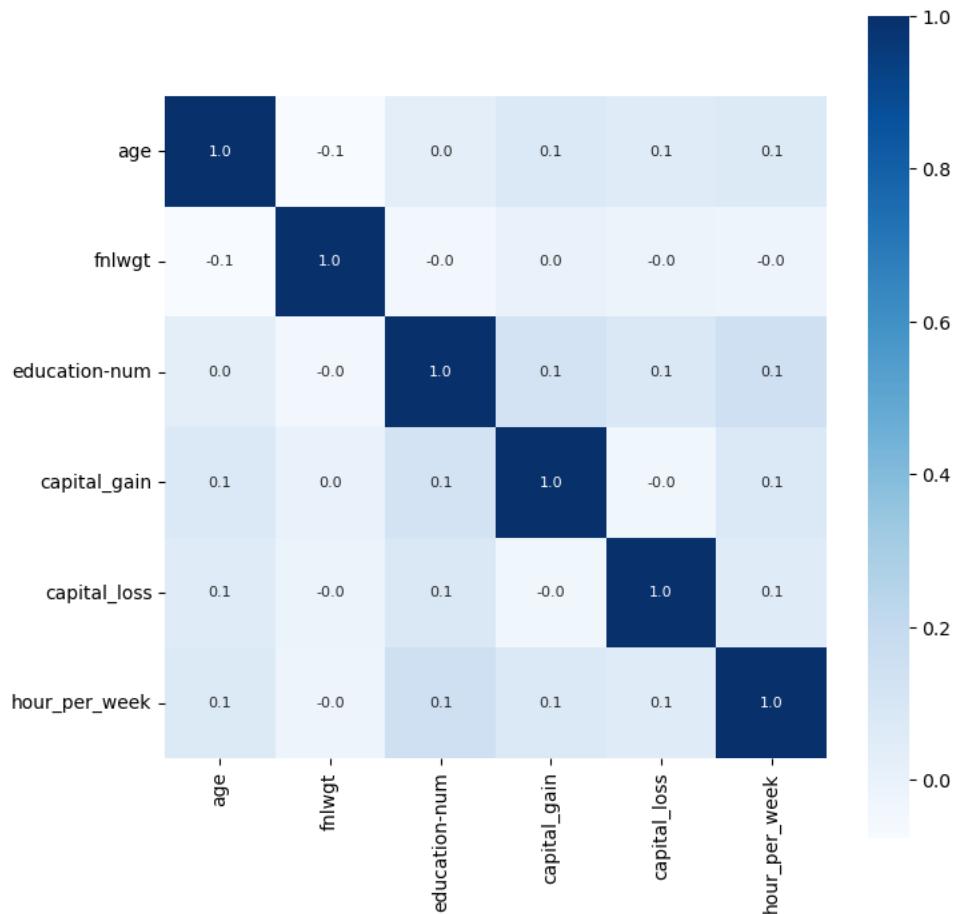
Top 5 Education class that pay highest capital gain



```

1 plt.figure(figsize=(8,8))
2 sns.heatmap(income.corr(),cbar=True, square=True, fmt='.1f', annot=True, annot_kws={'size':8}, cmap='Blues')
3 plt.show()

```



```

1 sns.pairplot(income, hue='income_class')
2 plt.show()

```



▼ Feature Engineering involves handling the following

- missing values
- handling outlier
- handling categorical encoding
- feature scaling
- feature selection

```

1 # Check the categorical columns propability of occurrence/frequency/count
2 for col in categorical_cols:
3     print(f"{col}")
4     print(f"{income[col].value_counts(normalize=True)*100}")
5     print("====")

```

```

Portugal          0.1131/
Nicaragua        0.104496
Peru             0.095276
France           0.089129
Greece           0.089129
Ecuador          0.086056
Ireland          0.073762
Hong             0.061468
Cambodia         0.058395
Trinidad&Tobago 0.058395
Laos             0.055322
Thailand          0.055322
Yugoslavia       0.049175
Outlying-US(Guam-USVI-etc) 0.043028
Honduras          0.039955
Hungary           0.039955
Scotland          0.036881
Holand-Netherlands 0.003073
Name: native_country, dtype: float64
=====
income_class
<=50K      75.907428
>50K      24.092572
Name: income_class, dtype: float64
=====
```

```

1 # Handling Columns Encoding
2 ordinal_cols = ["education", "relationship", "marital_status", "workclass"]
3 oho_cols = ["native_country", "sex", "race", "occupation"]
4
5 # Ranking of ordinal_cols
6 education = ['Preschool', 'HS-grad', 'Some-college', 'Assoc-voc', 'Assoc-acdm', 'Bachelors', 'Masters', 'Prof-school',
7   'Doctorate' ]
8 relationship = ['Not-in-family', 'Other-relative', 'Unmarried', 'Own-child', 'Wife', 'Husband']
9 marital_status = ['Never-married', 'Married-civ-spouse', 'Divorced', 'Married-spouse-absent'
10 'Separated', 'Married-AF-spouse', 'Widowed']
11 workclass = ['Never-worked', 'Without-pay', 'Self-emp-inc', 'Local-gov', 'Federal-gov', 'Private', 'Self-emp-not-inc', 'State-gov']
12

1 # Ordinal Encoding Using Replace Function
2 income['education'].replace({'Preschool':0, 'HS-grad':1, 'Some-college':2, 'Assoc-voc':3, 'Assoc-acdm':4, 'Bachelors':5, 'Masters':6, 'Prof-s
3 'Doctorate':8 }, inplace=True)
4 # For relationship column
5 income['relationship'].replace({'Not-in-family':0, 'Other-relative':1, 'Unmarried':2, 'Own-child':3, 'Wife':4, 'Husband':5 }, inplace=True)
6 # For marital_status
7 income['marital_status'].replace({'Never-married':0, 'Married-civ-spouse':1, 'Divorced':2, 'Married-spouse-absent':3,
8 'Separated':4, 'Married-AF-spouse':5, 'Widowed':6 }, inplace=True)
9 # For workclass
10 income['workclass'].replace({'Never-worked':0, 'Without-pay':1, 'Self-emp-inc':2, 'Local-gov':3, 'Federal-gov':4, 'Private':5, 'Self-emp-not
11

1 # Check The Encoding
2 for ord in ordinal_cols:
3     print(f"The ordinal encoding result for {ord} is: {income[ord].unique()}")
```

```

The ordinal encoding result for education is: [5 1 6 2 4 3 8 7 0]
The ordinal encoding result for relationship is: [0 5 4 3 2 1]
The ordinal encoding result for marital_status is: [0 1 2 3 4 5 6]
The ordinal encoding result for workclass is: [7 6 5 4 3 2 1 0]
```

```

1 # OneHot Encoding Using Python Function
2
3 # For Sex column
4 income['sex'] = np.where(income["sex"]=="Female", 0, 1)
5 print(income['sex'].unique())
6
7 # For Native_Country column
8 # American has almost 92% of the entire dataset. To avoid unnecessary multiplicity of columns using
9 # onehot encoding, I have basically encoded America with "1" while other nation with less than 8% with "0"
10 income['native_country'] = np.where(income["native_country"]=="United-States", 1, 0)
11 print(income['native_country'].unique())
12

[1 0]
[1 0]
```

```

1
2 income.head()
```

	age	workclass	fnlwgt	education	education-num	marital_status	occupation	relationship	race	sex	capital_gain	capital_loss	hours_per_week
0	39	7	77516	5	13	0	Adm-clerical		0	White	1	2174	0
1	50	6	83311	5	13	1	Exec-managerial		5	White	1	0	0
2	38	5	215646	1	9	2	Handlers-cleaners		0	White	1	0	0
3	53	5	234721	1	7	1	Handlers-cleaners		5	Black	1	0	0

```

1 # OneHotEncoding for other columns
2
3 # For Race and Occupation
4 ohe=pd.get_dummies(income[["race","occupation"]],drop_first=True)
5 ohe

```

	race_Asian-Pac-Islander	race_Black	race_Other	race_White	occupation_Armed-Forces	occupation_Craft-repair	occupation_Exec-managerial	occupation_Farming-fishing
0	0	0	0	1	0	0	0	0
1	0	0	0	1	0	0	1	0
2	0	0	0	1	0	0	0	0
3	0	1	0	0	0	0	0	0
4	0	1	0	0	0	0	0	0
...
32556	0	0	0	1	0	0	0	0
32557	0	0	0	1	0	0	0	0
32558	0	0	0	1	0	0	0	0
32559	0	0	0	1	0	0	0	0
32560	0	0	0	1	0	0	1	0

32537 rows × 17 columns

```

1 # Concatenate original and the dummy dataset
2 income= pd.concat([income,ohe],axis="columns")
3 income.head()

```

	age	workclass	fnlwgt	education	education-num	marital_status	occupation	relationship	race	sex	...	occupation_Farming-fishing	occ
0	39	7	77516	5	13	0	Adm-clerical		0	White	1	...	0
1	50	6	83311	5	13	1	Exec-managerial		5	White	1	...	0
2	38	5	215646	1	9	2	Handlers-cleaners		0	White	1	...	0
3	53	5	234721	1	7	1	Handlers-cleaners		5	Black	1	...	0
4	28	5	338409	5	13	1	Prof-specialty		4	Black	0	...	0

5 rows × 32 columns

```
1 income.columns
```

```

Index(['age', 'workclass', 'fnlwgt', 'education', 'education-num',
       'marital_status', 'occupation', 'relationship', 'race', 'sex',
       'capital_gain', 'capital_loss', 'hour_per_week', 'native_country',
       'income_class', 'race_Asian-Pac-Islander', 'race_Black', 'race_Other',
       'race_White', 'occupation_Armed-Forces', 'occupation_Craft-repair',
       'occupation_Exec-managerial', 'occupation_Farming-fishing',
       'occupation_Handlers-cleaners', 'occupation_Machine-op-inspct',
       'occupation_Other-service', 'occupation_Priv-house-serv',
       'occupation_Prof-specialty', 'occupation_Protective-serv',
       'occupation_Sales', 'occupation_Tech-support'],
      dtype='object')

```

```
'occupation_Transport-moving'],
dtype='object')
```

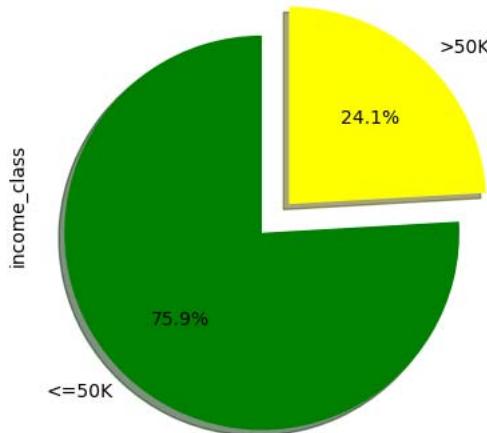
```
1 # Drop existing or original columns i.e Race and Occupation columns
2 income.drop(columns=["race","occupation"], axis = 1, inplace=True)
```

```
1 income.shape
```

```
(32537, 30)
```

```
1 # Calculating Lower and Upper Limit For columns And Capping Variables At Same
2 for col in df.columns[0:23]:
3     lower_limit, upper_limit = find_limits(df,col)
4     df[col]=np.where(df[col] > upper_limit, upper_limit,
5                         np.where(df[col] < lower_limit, lower_limit, df[col]))
6     print(f"For {col} variable: \n")
7     print(f"The Lower Limit: {lower_limit}")
8     print(f"The Upper Limit: {upper_limit}")
9     print("=====")
10
```

```
1 # Handling Imbalance Dataset
2 print(income["income_class"].value_counts())
3 income["income_class"].value_counts().plot.pie(y=income['income_class'], startangle=90, explode=(0.2,0), autopct='%1.1f%%', shadow=True, c
4 plt.show()
```



```
1 # Convert the target column categories into "0" and "1"
2 income['income_class']=income['income_class'].apply(lambda x:1 if x == ">50K" else 0)
3 income['income_class'].unique()
```

```
array([0, 1], dtype=int64)
```

```
1 # From the above, the dataset is not balanced and we need to oversample the imbalance class using SMOTE library
2 from imblearn.over_sampling import SMOTE
3
```

```
1 X = income.drop('income_class',axis=1)
2 y = income['income_class']
```

```
1 X_res,y_res = SMOTE(random_state=42).fit_resample(X,y)
```

```
1 X_res
```

	age	workclass	fnlwgt	education	education-num	marital_status	relationship	sex	capital_gain	capital_loss	...	occupation
0	39	7	77516	5	13	0	0	1	2174	0	0	...
1	50	6	83311	5	13	1	5	1	0	0	0	...
2	38	5	215646	1	9	2	0	1	0	0	0	...
3	53	5	234721	1	7	1	5	1	0	0	0	...
4	28	5	338409	5	13	1	4	0	0	0	0	...

```
1 y_res
```

```
0      0
1      0
2      0
3      0
4      0
..
49391  1
49392  1
49393  1
49394  1
49395  1
Name: income_class, Length: 49396, dtype: int64
```

```
1 # Let us combine the X_res and y_res back for further analysis after the oversampling
```

```
2 new_income=X_res
3 new_income["income_class"]=y_res
4 new_income
```

	age	workclass	fnlwgt	education	education-num	marital_status	relationship	sex	capital_gain	capital_loss	...	occupation
0	39	7	77516	5	13	0	0	1	2174	0	0	...
1	50	6	83311	5	13	1	5	1	0	0	0	...
2	38	5	215646	1	9	2	0	1	0	0	0	...
3	53	5	234721	1	7	1	5	1	0	0	0	...
4	28	5	338409	5	13	1	4	0	0	0	0	...
..
49391	43	5	280509	4	12	1	4	0	15024	0	0	...
49392	43	4	198096	4	12	1	5	1	7688	0	0	...
49393	48	4	242222	4	12	1	5	1	4386	0	0	...
49394	43	5	332633	4	12	0	4	1	6868	0	0	...
49395	60	5	238913	4	12	1	5	1	0	0	0	...

49396 rows × 30 columns

```
1 # Checking the dataset balance after ovemsampling
2 # Handling Imbalance Dataset
3 print(new_income["income_class"].value_counts())
4 new_income["income_class"].value_counts().plot.pie(y=new_income['income_class'], startangle=90, explode=(0.2,0), autopct='%1.1f%%', shadow=True)
5 plt.show()
```

```
0    24698
1    24698
Name: income_class, dtype: int64
```

```
1 # The education_num needs to be droped since we already have a more encoded and
2 # representative of same in education column
3 new_income.drop("education-num",axis=1, inplace=True)
4 new_income.shape
5
```

```
(49396, 29)
0 50.0% 1
```

```
1 # Check info about the new resampled dataset
2 new_income.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 49396 entries, 0 to 49395
Data columns (total 29 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   age              49396 non-null   int64  
 1   workclass        49396 non-null   int64  
 2   fnlwgt           49396 non-null   int64  
 3   education        49396 non-null   int64  
 4   marital_status   49396 non-null   int64  
 5   relationship     49396 non-null   int64  
 6   sex              49396 non-null   int32  
 7   capital_gain    49396 non-null   int64  
 8   capital_loss    49396 non-null   int64  
 9   hour_per_week   49396 non-null   int64  
 10  native_country  49396 non-null   int32  
 11  race_Aisan-Pac-Islander 49396 non-null   uint8  
 12  race_Black      49396 non-null   uint8  
 13  race_Other      49396 non-null   uint8  
 14  race_White      49396 non-null   uint8  
 15  occupation_Armed-Forces 49396 non-null   uint8  
 16  occupation_Craft-repair 49396 non-null   uint8  
 17  occupation_Exec-managerial 49396 non-null   uint8  
 18  occupation_Farming-fishing 49396 non-null   uint8  
 19  occupation_Handlers-cleaners 49396 non-null   uint8  
 20  occupation_Machine-op-inspct 49396 non-null   uint8  
 21  occupation_Other-service 49396 non-null   uint8  
 22  occupation_Priv-house-serv 49396 non-null   uint8  
 23  occupation_Prof-specialty 49396 non-null   uint8  
 24  occupation_Protective-serv 49396 non-null   uint8  
 25  occupation_Sales      49396 non-null   uint8  
 26  occupation_Tech-support 49396 non-null   uint8  
 27  occupation_Transport-moving 49396 non-null   uint8  
 28  income_class       49396 non-null   int64  
dtypes: int32(2), int64(10), uint8(17)
memory usage: 4.9 MB
```

```
1 # Save the cleaned dataset to a csv file named "final.csv"
2 new_income.to_csv("final.csv")
```

▼ Model Building, Hyper-parameter Tuning And Model Evaluation

```
1 # Import The Key Packages
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.metrics import accuracy_score,classification_report,roc_auc_score
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler
```

```
1 # Splitting Dataset Into Dependent And Independent Features
2 X =new_income.drop("income_class", axis =1)
3 y = new_income["income_class"]
```

```
1 # Split Dataset Into Train And Test Splits
2 X_train,X_test,y_train,y_test = train_test_split(X,y, test_size=0.2,random_state=42)
```

```
1 # Dataset Scaling And Standardisation
2 data_scaling = StandardScaler()
3 X_train =data_scaling.fit_transform(X_train)
4 X_test = data_scaling.transform(X_test)
```

```
1 print(f"X_train shape: {X_train.shape}")
2 print(f"X_test shape: {X_test.shape}")
```

```
3 print(f"X_train shape: {y_train.shape}")
4 print(f"X_test shape: {y_test.shape}")
  X_train shape: (39516, 28)
  X_test shape: (9880, 28)
  X_train shape: (39516,)
  X_test shape: (9880,)
```

▼ Logistic Regression Without Hyper-parameter

```
1 # Model Building
2 log_regression=LogisticRegression(verbose=False)
3 log_regression.fit(X_train, y_train)
4 y_pred = log_regression.predict(X_test)

1 # Model Evaluation
2 print("====")
3 print("Accuracy_Score")
4 print(accuracy_score(y_test, y_pred))
5 print("====")
6 print(classification_report(y_test, y_pred))

=====
Accuracy_Score
0.8545546558704453
=====

      precision    recall   f1-score   support

          0       0.85      0.86      0.86      4963
          1       0.86      0.85      0.85      4917

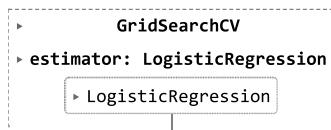
   accuracy                           0.85      9880
    macro avg       0.85      0.85      0.85      9880
weighted avg       0.85      0.85      0.85      9880
```

▼ Logistic Regression With Hyper-parameter Tuning Using GridSearchCV

```
1 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV

1 param_grid = {
2     "penalty": ["l1", "l2", "elasticnet"],
3     "C" : [0.001, 0.01, 0.1, 1, 10],
4     "solver" : ["lbfgs", "liblinear", "newton-cg", "sag", "saga"],
5     "max_iter" : [100, 300, 500, 1000, 1500]
6 }
```

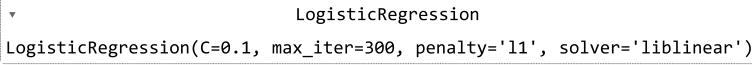
```
1 # Perform grid search with cross-validation
2 grid_search = GridSearchCV(log_regression, param_grid, cv=5, scoring='accuracy', n_jobs=-1)
3 grid_search.fit(X_train, y_train)
```



```
1 # Print the best hyperparameters
2 print("Best hyperparameters: ", grid_search.best_params_)
3

Best hyperparameters: {'C': 0.1, 'max_iter': 300, 'penalty': 'l1', 'solver': 'liblinear'}
```

```
1 # Train the model with the best hyperparameters
2 best_model = LogisticRegression(**grid_search.best_params_)
3 best_model.fit(X_train, y_train)
```

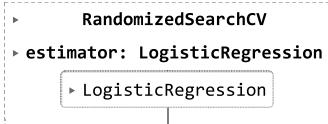


```
1 # Evaluate the performance on the validation(test) set
2 y_pred = best_model.predict(X_test)
3 accuracy = accuracy_score(y_test, y_pred)
4 print("Accuracy on validation set: {:.2f}".format(accuracy))

Accuracy on validation set: 0.85
```

▼ Logistic Regression With Hyper-parameter Tuning Using RandomizedSearchCV

```
1 # Perform random search with cross-validation
2 random_search = RandomizedSearchCV(log_regression, param_grid, cv=5, scoring='accuracy', n_jobs=-1, n_iter=20, random_state=42, verbose=0)
3 random_search.fit(X_train, y_train)
4
```



```
1 # Print the best hyperparameters
2 print("Best hyperparameters: ", random_search.best_params_)

Best hyperparameters: {'solver': 'liblinear', 'penalty': 'l1', 'max_iter': 100, 'C': 0.01}
```

```
1 # Train the model with the best hyperparameters
2 best_model = LogisticRegression(**random_search.best_params_)
3 best_model.fit(X_train, y_train)
```

```
graph TD; Model[LogisticRegression] --> Params[LogisticRegression(C=0.01, penalty='l1', solver='liblinear')]
```

```
1 # Evaluate the performance on the validation set
2 y_pred = best_model.predict(X_test)
3 accuracy = accuracy_score(y_test, y_pred)
4 print("Accuracy on validation set: {:.2f}".format(accuracy))
5
```

```
Accuracy on validation set: 0.85
```

```
1
```

