

```
In [1]: from pandas import read_csv
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.neighbors import KNeighborsClassifier
from sklearn.linear_model import LogisticRegression
import warnings
from sklearn.metrics import accuracy_score
warnings.filterwarnings('ignore')

In [2]: train = pd.read_csv('Users\SAURABH\Saurabh patil\DATA SCIENCE\Forecasting\SalaryData_Train(1).csv')
```

Out[2]:

| | age | workclass | education | educationno | maritalstatus | occupation | relationship | race | sex | capitalgain | capitalloss | hoursperweek | native | Salary | |
|-------|-----|------------------|------------|-------------|-------------------|---------------------|---------------|-------|--------|-------------|-------------|--------------|--------|---------------|-------|
| 0 | 39 | State-gov | Bachelors | 13 | Never-married | Adm-clerical | Not-in-family | White | Male | 2174 | 0 | 0 | 40 | United-States | <=50K |
| 1 | 50 | Self-emp-not-inc | Bachelors | 13 | Married-cv-spouse | Exec-managerial | Husband | White | Male | 0 | 0 | 0 | 13 | United-States | <=50K |
| 2 | 38 | Private | HS-grad | 9 | Divorced | Handiclers-cleaners | Not-in-family | White | Male | 0 | 0 | 0 | 40 | United-States | <=50K |
| 3 | 53 | Private | 11th | 7 | Married-cv-spouse | Handiclers-cleaners | Husband | Black | Male | 0 | 0 | 0 | 40 | United-States | <=50K |
| 4 | 28 | Private | Bachelors | 13 | Married-cv-spouse | Prof-specialty | Wife | Black | Female | 0 | 0 | 0 | 40 | Cuba | <=50K |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 30156 | 27 | Private | Assoc-acdm | 12 | Married-cv-spouse | Tech-support | Wife | White | Female | 0 | 0 | 0 | 38 | United-States | <=50K |
| 30157 | 40 | Private | HS-grad | 9 | Married-cv-spouse | Machine-op-inspct | Husband | White | Male | 0 | 0 | 0 | 40 | United-States | >50K |
| 30158 | 58 | Private | HS-grad | 9 | Widowed | Adm-clerical | Unmarried | White | Female | 0 | 0 | 0 | 40 | United-States | <=50K |
| 30159 | 22 | Private | HS-grad | 9 | Never-married | Adm-clerical | Own-child | White | Male | 0 | 0 | 0 | 20 | United-States | <=50K |
| 30160 | 52 | Self-emp-inc | HS-grad | 9 | Married-cv-spouse | Exec-managerial | Wife | White | Female | 15024 | 0 | 0 | 40 | United-States | >50K |

30161 rows x 14 columns

In [3]:

```
test = pd.read_csv('/Users/SAURABH/Saurabh patil/DATA SCIENCE/Forecasting/SalaryData_test(1).csv')
test
```

Out[3]:

| | age | workclass | education | educationno | maritalstatus | occupation | relationship | race | sex | capitalgain | capitalloss | hoursperweek | native | Salary | |
|---|-----|-----------|-----------|-------------|---------------|-------------------|--------------|-------|------|-------------|-------------|--------------|--------|---------------|-------|
| 0 | 25 | Private | 11th | 7 | Never-married | Machine-op-inspct | Own-child | Black | Male | 0 | 0 | 0 | 40 | United-States | <=50K |

```
In [3]: test = pd.read_csv('Users\SAURABH\Saurabh patil\DATA SCIENCE\Forecasting\SalaryData_Test(1).csv')
train.info()

Out[3]:
```

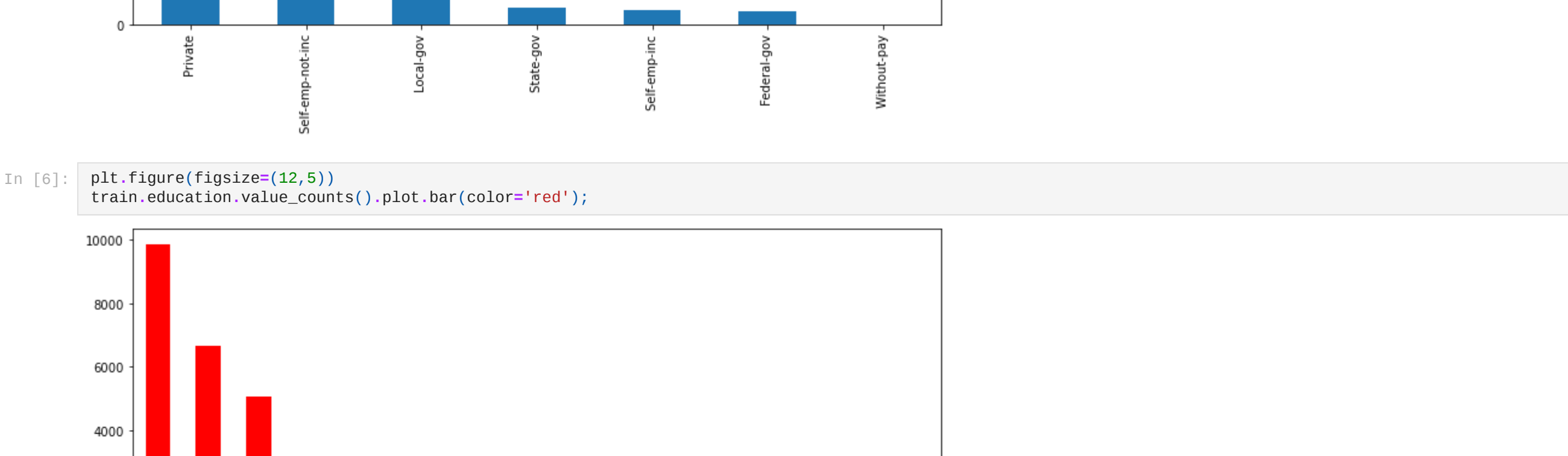
| | age | workclass | education | educationno | maritalstatus | occupation | relationship | race | sex | capitalgain | capitalloss | hoursperweek | native | Salary | |
|-------|-----|--------------|--------------|-------------|--------------------|-------------------|---------------|--------------------|--------|-------------|-------------|--------------|--------|---------------|-------|
| 0 | 25 | Private | 11th | 7 | Never-married | Machine-op-inspct | Own-child | Black | Male | 0 | 0 | 0 | 40 | United-States | <=50K |
| 1 | 38 | Private | HS-grad | 9 | Married-civ-spouse | Farming-fishing | Husband | White | Male | 0 | 0 | 0 | 50 | United-States | <=50K |
| 2 | 28 | Local-gov | Assoc-acdm | 12 | Married-civ-spouse | Protective-serv | Husband | White | Male | 0 | 0 | 0 | 40 | United-States | >50K |
| 3 | 44 | Private | Some-college | 10 | Married-civ-spouse | Machine-op-inspct | Husband | Black | Male | 7688 | 0 | 0 | 40 | United-States | >50K |
| 4 | 34 | Private | 10th | 6 | Never-married | Other-service | Not-in-family | White | Male | 0 | 0 | 0 | 30 | United-States | <=50K |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 15055 | 33 | Private | Bachelors | 13 | Never-married | Prof-specialty | Own-child | White | Male | 0 | 0 | 0 | 40 | United-States | <=50K |
| 15056 | 39 | Private | Bachelors | 13 | Divorced | Prof-specialty | Not-in-family | White | Female | 0 | 0 | 0 | 36 | United-States | <=50K |
| 15057 | 38 | Private | Bachelors | 13 | Married-civ-spouse | Prof-specialty | Husband | White | Male | 0 | 0 | 0 | 50 | United-States | <=50K |
| 15058 | 44 | Private | Bachelors | 13 | Divorced | Adm-clerical | Own-child | Asian-Pac-Islander | Male | 5455 | 0 | 0 | 40 | United-States | <=50K |
| 15059 | 35 | Self-emp-inc | Bachelors | 13 | Married-civ-spouse | Exec-managerial | Husband | White | Male | 0 | 0 | 0 | 60 | United-States | >50K |

15060 rows x 14 columns

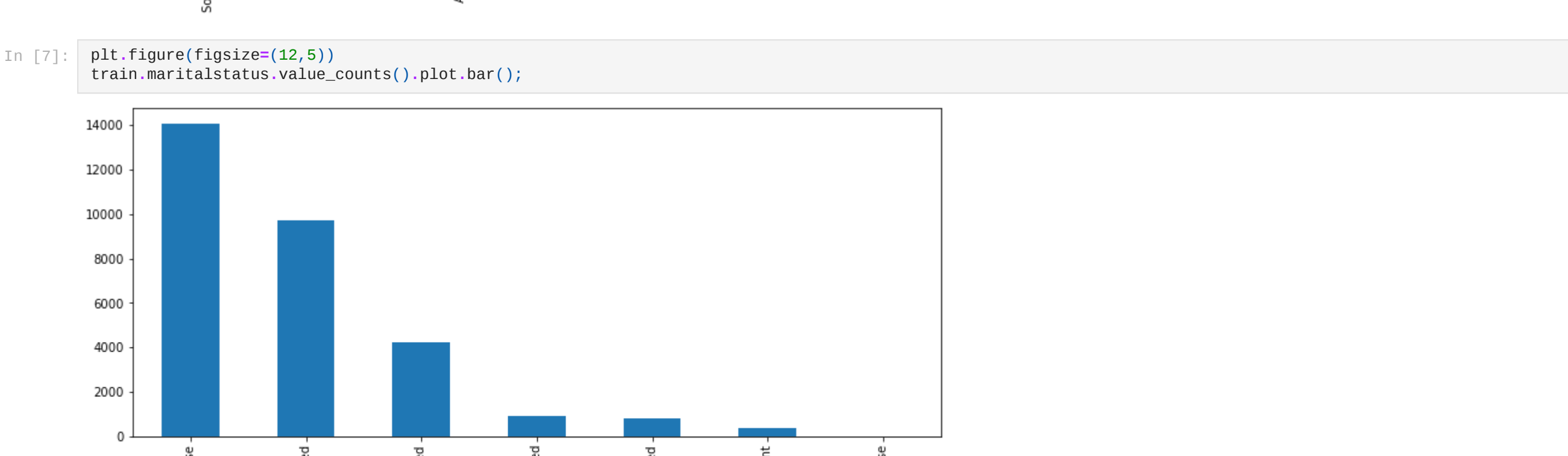
```
In [4]: #Checking for null values & data types
train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30161 entries, 0 to 30160
Data columns (total 14 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   age                   30161 non-null   int64
 1   workclass             30161 non-null   object
 2   education             30161 non-null   object
 3   educationno           30161 non-null   int64
 4   maritalstatus         30161 non-null   object
 5   occupation            30161 non-null   object
 6   relationship          30161 non-null   object
 7   race                 30161 non-null   object
 8   sex                  30161 non-null   object
 9   capitalgain           30161 non-null   int64
10   capitalloss           30161 non-null   int64
11   hoursperweek          30161 non-null   int64
12   native                30161 non-null   object
13   salary                30161 non-null   object
dtypes: int64(8), object(6)
memory usage: 2.2+ MB
```

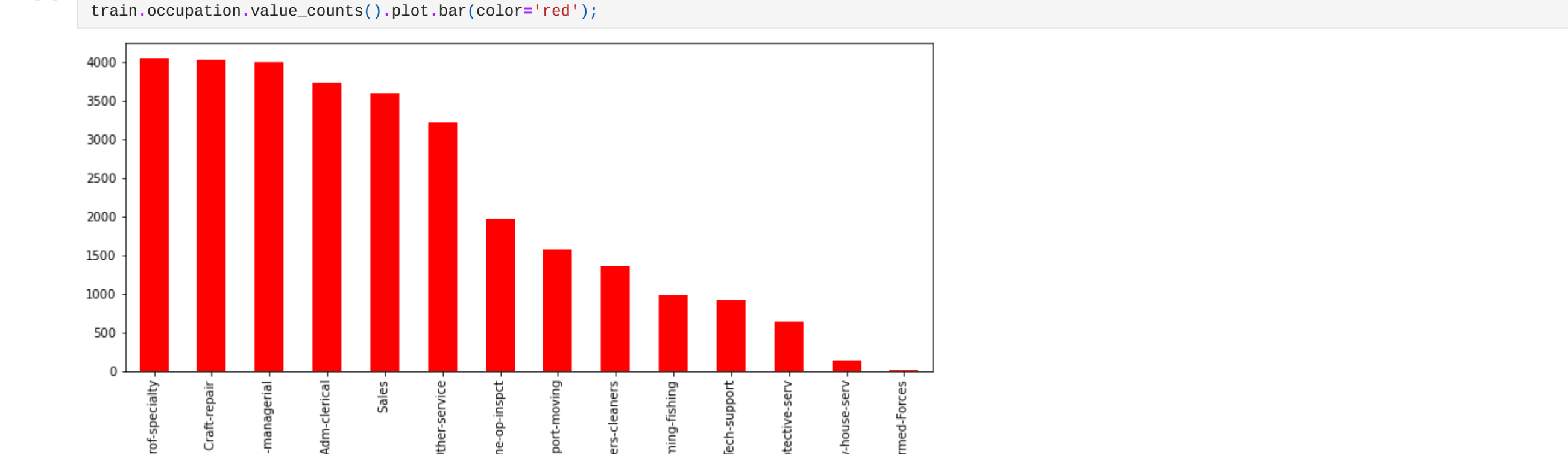
```
In [5]: plt.figure(figsize=(2,5))
train.workclass.value_counts().plot.bar();
```



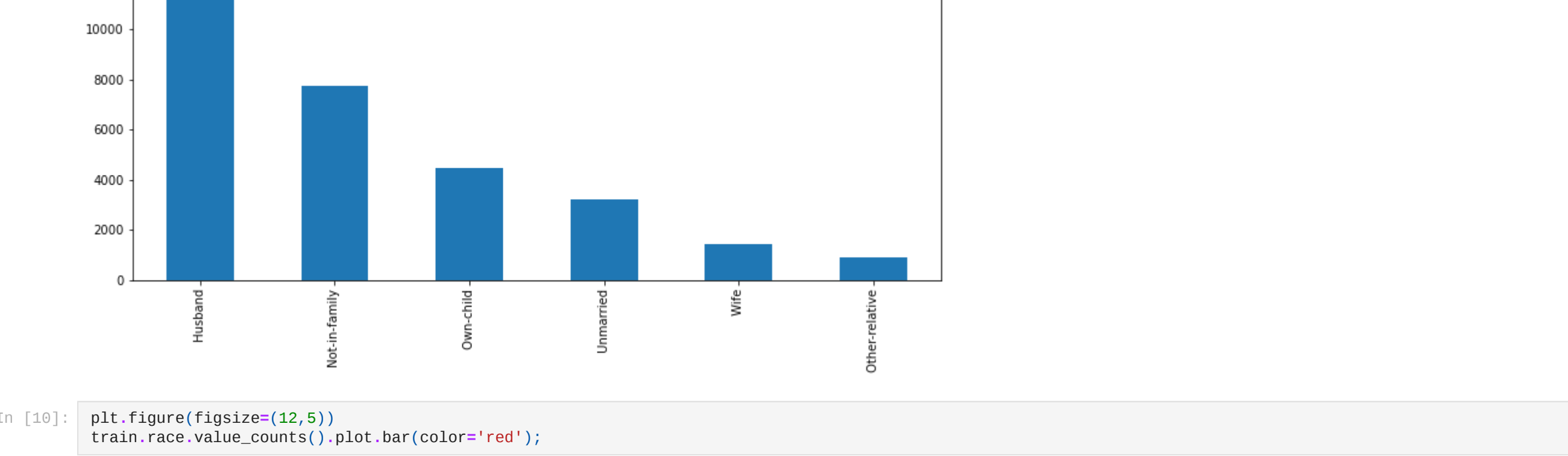
```
In [6]: plt.figure(figsize=(2,5))
train.education.value_counts().plot.bar(color='red');
```



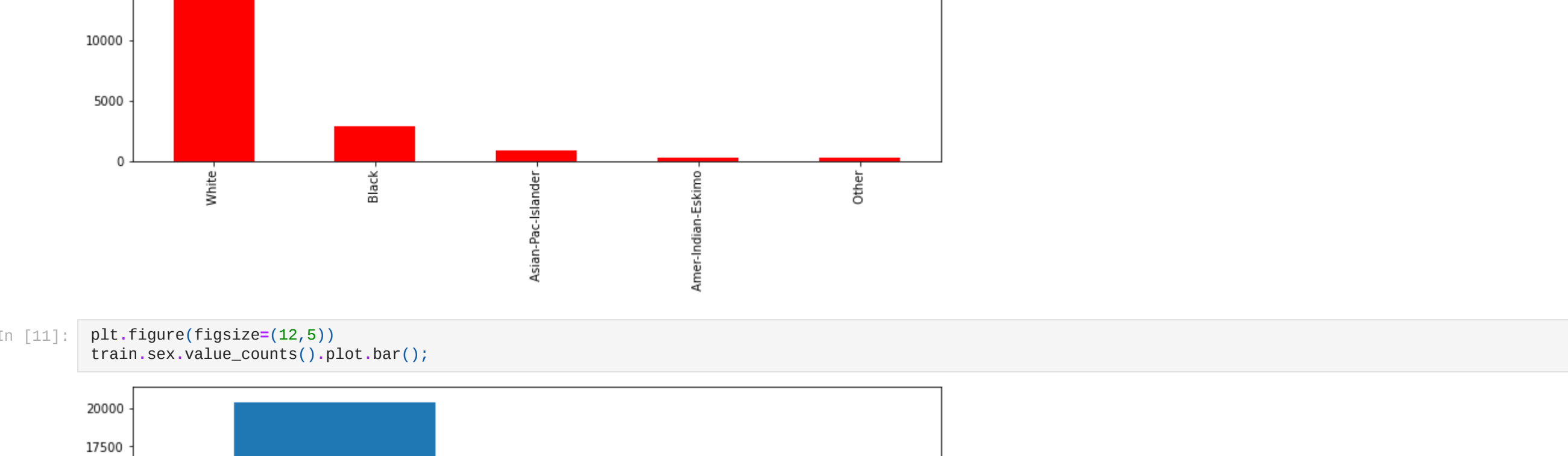
```
In [7]: plt.figure(figsize=(2,5))
train.maritalstatus.value_counts().plot.bar();
```



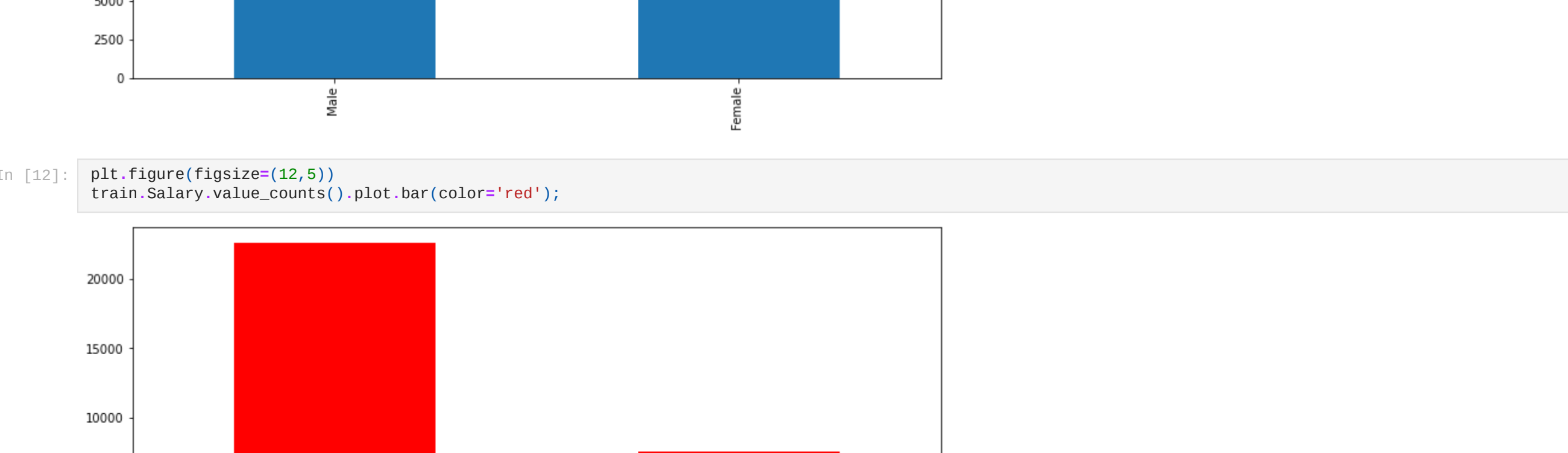
```
In [8]: plt.figure(figsize=(2,5))
train.occupation.value_counts().plot.bar(color='red');
```



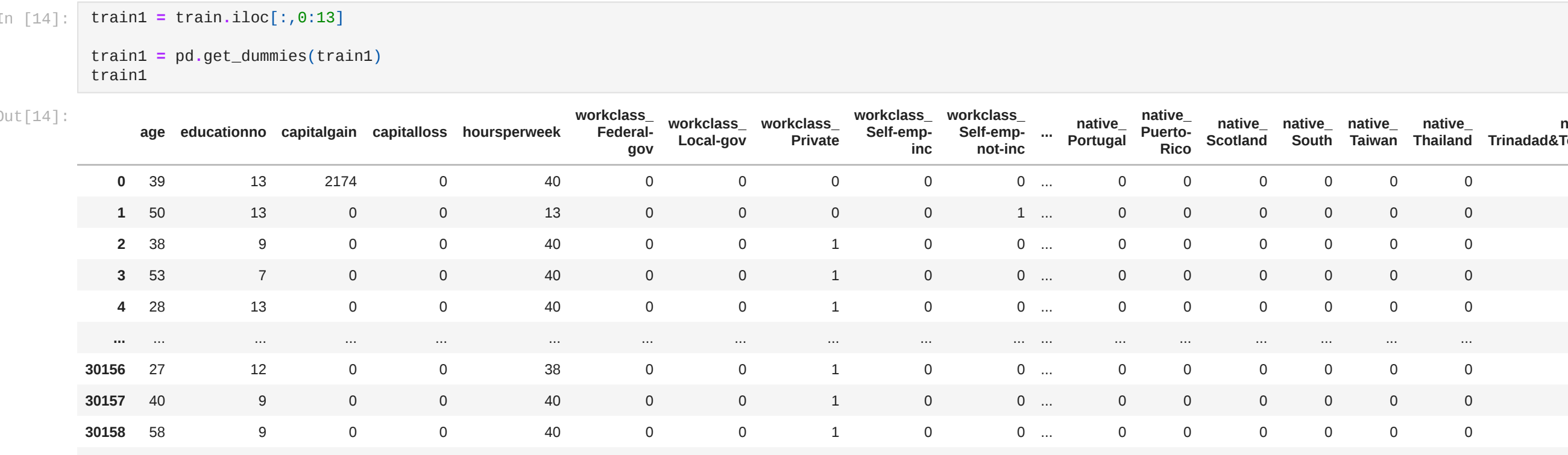
```
In [9]: plt.figure(figsize=(2,5))
train.relationship.value_counts().plot.bar();
```



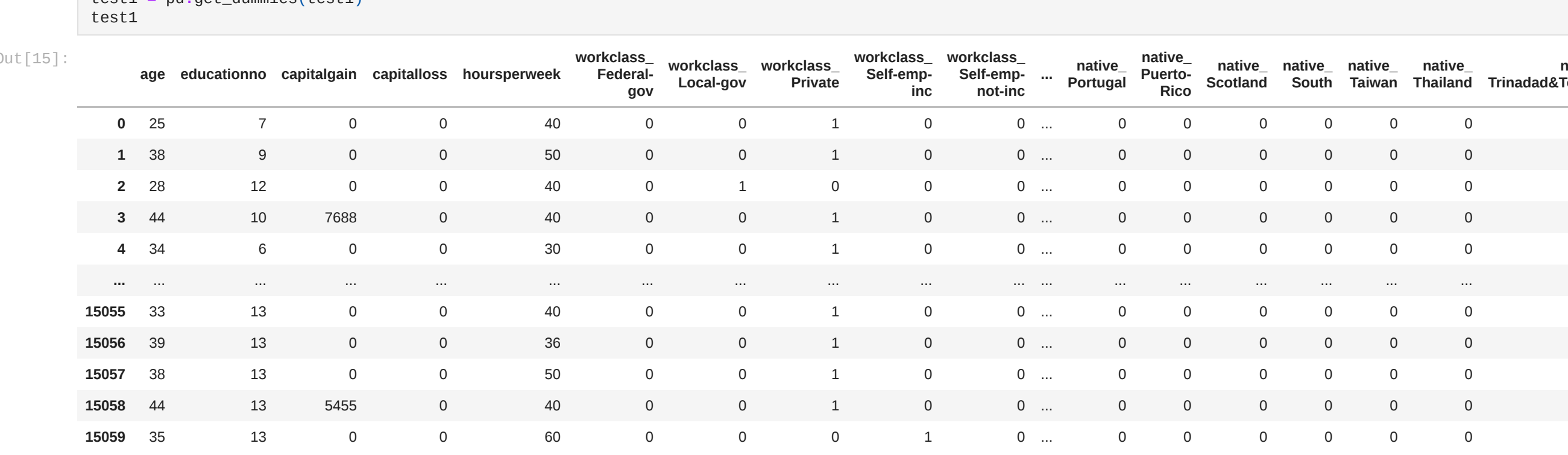
```
In [10]: plt.figure(figsize=(2,5))
train.race.value_counts().plot.bar(color='red');
```



```
In [11]: plt.figure(figsize=(2,5))
train.sex.value_counts().plot.bar();
```



```
In [12]: plt.figure(figsize=(2,5))
train.salary.value_counts().plot.bar(color='red');
```



```
In [14]: train1 = train.iloc[:,0:13]
train1 = pd.get_dummies(train1)
```

Out[14]:

| | age | educationno | capitalgain | capitalloss | hoursperweek | workclass_Federal_gov | workclass_Local_gov | workclass_Private | workclass_Self-emp-inc | workclass_Self-emp-not-inc | native_Puerto-Rico | native_Scotland | native_South | native_Taiwan | native_Thailand | native_Trinidad&Tobago | no |
|-------|-----|-------------|-------------|-------------|--------------|-----------------------|---------------------|-------------------|------------------------|----------------------------|--------------------|-----------------|--------------|---------------|-----------------|------------------------|-----|
| 0 | 39 | 13 | 2174 | 0 | 40 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 50 | 13 | 0 | 0 | 13 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 38 | 9 | 0 | 0 | 40 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 53 | 7 | 0 | 0 | 40 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 28 | 13 | 0 | 0 | 40 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 30156 | 27 | 12 | 0 | 0 | 38 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 30157 | 40 | 9 | 0 | 0 | 40 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 30158 | 58 | 9 | 0 | 0 | 40 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 30159 | 22 | 9 | 0 | 0 | 20 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 30160 | 52 | 9 | 15024 | 0 | 40 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

30161 rows x 102 columns

```
In [15]: test1 = test.iloc[:,0:13]
test1 = pd.get_dummies(test1)
```

Out[15]:

| | age | educationno | capitalgain | capitalloss | hoursperweek | workclass_Federal_gov | workclass_Local_gov | workclass_Private | workclass_Self-emp-inc | workclass_Self-emp-not-inc | native_Puerto-Rico | native_Scotland | native_South | native_Taiwan | native_Thailand | native_Trinidad&Tobago | no |
|-------|-----|-------------|-------------|-------------|--------------|-----------------------|---------------------|-------------------|------------------------|----------------------------|--------------------|-----------------|--------------|---------------|-----------------|------------------------|-----|
| 0 | 25 | 7 | 0 | 0 | 40 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 38 | 9 | 0 | 0 | 50 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 28 | 12 | 0 | 0 | 40 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 44 | 10 | 7688 | 0 | 40 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 34 | 6 | 0 | 0 | 30 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 15055 | 33 | 13 | 0 | 0 | 40 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 15056 | 39 | 13 | 0 | 0 | 36 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 15057 | 38 | 13 | 0 | 0 | 50 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 15058 | 44 | 13 | 5455 | 0 | 40 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 15059 | 35 | 13 | 0 | 0 | 60 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

15060 rows x 102 columns

Since number of features are more, let's use PCA

```
In [16]: #Scaling the data
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
train_norm = sc.transform(train1)
```

```
Out[16]: array([[ 0.84277892,  1.12898813,  0.14698988, ...,  0.31881205,
-0.84611353,  0.8239384 ],
[ 0.48820881,  1.12898813, -0.14744712, ...,  0.31881205,
-0.84611353,  0.8239384 ],
[ 0.6232558 ,  0.43972325, -0.14744712, ...,  0.31881205,
-0.84611353,  0.8239384 ],
[ 0.48933854,  0.43972325, -0.14744712, ...,  0.31881205,
-0.84611353,  0.8239384 ],
[ 1.55151256,  0.43972325, -0.14744712, ...,  0.31881205,
-0.84611353,  0.8239384 ],
[ 1.63253024,  0.43972325, -0.14744712, ...,  0.31881205,
-0.84611353,  0.8239384 ]])
```

```
In [17]: #Scaling the data
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
test_fit(test1)
test_norm = sc.transform(test1)
```

```
Out[17]: array([[ -1.02000513,  1.21850528, -0.14543845, ...,  0.30373366,
-0.85954172,  0.82156441],
[ 0.65742253,  0.43498824, -0.14543845, ...,  0.30373366,
-0.85954172,  0.82156441],
[ 0.80479376,  0.73195862, -0.14543845, ...,  0.30373366,
-0.85954172,  0.82156441],
[ 0.35742253,  1.12840499, -0.14543845, ...,  0.30373366,
-0.85954172,  0.82156441],
[ 0.30100821,  1.12840499, -0.14543845, ...,  0.30373366,
-0.85954172,  0.82156441],
[ 0.2816339 ,  1.12840499, -0.14543845, ...,  0.30373366,
-0.85954172,  0.82156441]])
```

```
In [18]: from sklearn.decomposition import PCA
train_pca = PCA(n_components = 382)
train_pca_values = train_pca.fit_transform(train_norm)
```

```
Out[18]: array([[ 5.58938898e-01,  2.38164988e+00,  5.91921169e-01, ...,
1.86102284e-15,  1.56196155e-15,  7.38126458e-16],
[ 2.81915815e-01,  1.53441835e+00,  1.46521931e-15],
[ 8.8236243e-02,  1.20633855e-02,  1.10277829e-02,  1.17761599e-02,
1.22845428e-02,  1.20633855e-02,  1.10277829e-02,  1.17761599e-02,
1.73932152e-01,  1.14590565e-01,  1.26215159e-01, ...,
3.53912850e-15,  1.64468924e-15,  6.83012194e-17],
[ 2.37835145e+00,  7.98690413e-01,  3.98105788e-01, ...,
-1.62601950e-17,  5.8465906e-17,  3.86829768e-17],
[ 0.87747329e+00,  1.37130326e-02,  9.82895489e-01,  1.82895489e+00,
1.02340338e-02,  1.02313312e-02,  1.0746044e-02,  1.88038858e-02,
1.08933098e-02,  1.08907488e-02,  9.37807138e-03,  9.83215141e-03,
9.9123258e-03,  9.87297873e-03,  9.80864172e-03,  9.85346688e-03,
9.83567041e-03,  9.82654350e-03,  9.82143339e-03,  9.81669030e-03,
9.81361940e-03,  9.89768489e-03,  9.80534229e-03,  9.80865163e-03,
9.78187138e-03,  9.77822386e-03,  9.71987789e-03,  9.78027058e-03,
9.72968212e-03,  9.71292926e-03,  9.68846212e-03,  9.68294717e-03,
9.66992730e-03,  9.67649451e-03,  9.62281979e-03,  9.58462673e-03,
9.56147128e-03,  9.52463771e-03,  9.42925368e-03,  9.43014688e-03,
9.36144053e-03,  9.38687731e-03,  9.26389704e-03,  9.11455862e-03,
9.07861657e-03,  9.07108230e-03,  9.04283193e-03,  8.93778300e-03,
8.72782590e-03,  8.59927945e-03,  8.45748487e-03,  8.40278797e-03,
9.31895498e-03,  9.17023620e-03,  9.16541797e-03,  7.17330330e-03,
5.13950548e-03,  4.76169851e-03,  4.27205376e-03,  2.37624813e-03,
1.95887190e-04,  4.46131301e-02,  2.23683176e-02,  1.45886206e-02,
3.35283508e-03,  8.27375302e-03,  7.32353682e-03,  4.67062116e-03,
2.48054012e-03,  1.43466435e-04]])
```

```
In [19]: from sklearn.decomposition import PCA
test_pca = PCA(n_components = 382)
test_pca_values = test_pca.fit_transform(test_norm)
```

```
Out[19]: array([[ 2.24293780e+00,  2.68318991e+00,  3.27616530e-01, ...,
2.23860506e-15,  1.53441835e+00,  1.46521931e-15],
[ 2.22690391e+00,  1.58471521e+00,  7.32087794e-01, ...,
1.57132732e-14,  5.28996417e-05,  6.51649996e-01],
[ 2.30704416e+00,  1.16883181e+00,  2.98521481e-01, ...,
2.30704416e+00,  1.16883181e+00,  2.98521481e-01],
[ 2.21130870e+00,  1.35451566e-15,  2.98521481e-01],
[ 2.35853218e+00,  1.46521931e+00,  3.27616530e-01, ...,
1.19801350e-16,  4.46131301e-02,  2.23683176e-02],
[ 1.84195534e-02,  1.03778312e-02,  9.82895489e-01,  1.82895489e+00,
1.02340338e-02,  1.02313312e-02,  1.0746044e-02,  1.88038858e-02,
1.08933098e-02,  1.08907488e-02,  9.37807138e-03,  9.83215141e-03,
9.9123258e-03
```