



Open
Data
Science




mlcourse.ai – Open Machine Learning Course

Author: [Yury Kashnitskiy](#)

Translated and edited by [Sergey Isaev](#), [Artem Trunov](#), [Anastasia Manokhina](#), and [Yuan Yuan Pao](#)

All content is distributed under the [Creative Commons CC BY-NC-SA 4.0](#) license.

Assignment #1 (demo). Solution

Exploratory data analysis with Pandas

In this task you should use Pandas to answer a few questions about the [Adult](#) dataset. (You don't have to download the data – it's already in the repository). Choose the answers in the [web-form](#).

Unique values of features (for more information please see the link above):

- `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, Trinidad&Tobago, Peru, Hong, Holand-Netherlands.
- `salary`: >50K, <=50K

In [1]:

```
import numpy as np
import pandas as pd
```

In [2]:

```
data = pd.read_csv('../input/adult.data.csv')
data.head()
```

Out[2]:

	age	workclass	fnlwgt	education	education-num	marital-status	occupation	relationship	race	sex	capital-gain	capital-loss	hours-per-week
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	0	40
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	13
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	40
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0	40
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	40

1. How many men and women (*sex* feature) are represented in this dataset?

In [3]:

```
data['sex'].value_counts()
```

Out[3]:

```
Male      21790
Female    10771
Name: sex, dtype: int64
```

2. What is the average age (*age* feature) of women?

In [4]:

```
data.loc[data['sex'] == 'Female', 'age'].mean()
```

Out[4]:

```
36.85823043357163
```

3. What is the proportion of German citizens (*native-country* feature)?

In [5]:

```
float((data['native-country'] == 'Germany').sum()) / data.shape[0]
```

Out[5]:

```
0.004207487485028101
```

4-5. What are mean value and standard deviation of the age of those who recieve more than 50K per year (*salary* feature) and those who receive less than 50K per year?

In [6]:

```
ages1 = data.loc[data['salary'] == '>50K', 'age']
ages2 = data.loc[data['salary'] == '<=50K', 'age']
print("The average age of the rich: {0} +- {1} years, poor - {2} +- {3} years.".format(
    round(ages1.mean()), round(ages1.std(), 1),
    round(ages2.mean()), round(ages2.std(), 1))
```

```
round(ages2.mean()), round(ages2.std(), 1))
```

The average age of the rich: 44 +- 10.5 years, poor - 37 +- 14.0 years.

6. Is it true that people who receive more than 50k have at least high school education? (*education - Bachelors, Prof-school, Assoc-acdm, Assoc-voc, Masters or Doctorate* feature)

In [7]:

```
data.loc[data['salary'] == '>50K', 'education'].unique() # No
```

Out[7]:

```
array(['HS-grad', 'Masters', 'Bachelors', 'Some-college', 'Assoc-voc',  
      'Doctorate', 'Prof-school', 'Assoc-acdm', '7th-8th', '12th',  
      '10th', '11th', '9th', '5th-6th', '1st-4th'], dtype=object)
```

7. Display statistics of age for each race (*race* feature) and each gender. Use *groupby()* and *describe()*. Find the maximum age of men of *Amer-Indian-Eskimo* race.

In [8]:

```
for (race, sex), sub_df in data.groupby(['race', 'sex']):  
    print("Race: {0}, sex: {1}".format(race, sex))  
    print(sub_df['age'].describe())
```

Race: Amer-Indian-Eskimo, sex: Female

```
count    119.000000  
mean      37.117647  
std       13.114991  
min       17.000000  
25%       27.000000  
50%       36.000000  
75%       46.000000  
max       80.000000
```

Name: age, dtype: float64

Race: Amer-Indian-Eskimo, sex: Male

```
count    192.000000  
mean      37.208333  
std       12.049563  
min       17.000000  
25%       28.000000  
50%       35.000000  
75%       45.000000  
max       82.000000
```

Name: age, dtype: float64

Race: Asian-Pac-Islander, sex: Female

```
count    346.000000  
mean      35.089595  
std       12.300845  
min       17.000000  
25%       25.000000  
50%       33.000000  
75%       43.750000  
max       75.000000
```

Name: age, dtype: float64

Race: Asian-Pac-Islander, sex: Male

```
count    693.000000  
mean      39.073593  
std       12.883944  
min       18.000000  
25%       29.000000  
50%       37.000000  
75%       46.000000  
max       90.000000
```

Name: age, dtype: float64

Race: Black, sex: Female

```
count    1555.000000  
mean      37.854019  
std       12.637197  
min       17.000000  
25%       28.000000
```

```

50%      37.000000
75%      46.000000
max       90.000000
Name: age, dtype: float64
Race: Black, sex: Male
count    1569.000000
mean      37.682600
std       12.882612
min       17.000000
25%       27.000000
50%       36.000000
75%       46.000000
max       90.000000
Name: age, dtype: float64
Race: Other, sex: Female
count     109.000000
mean      31.678899
std       11.631599
min       17.000000
25%       23.000000
50%       29.000000
75%       39.000000
max       74.000000
Name: age, dtype: float64
Race: Other, sex: Male
count     162.000000
mean      34.654321
std       11.355531
min       17.000000
25%       26.000000
50%       32.000000
75%       42.000000
max       77.000000
Name: age, dtype: float64
Race: White, sex: Female
count    8642.000000
mean      36.811618
std       14.329093
min       17.000000
25%       25.000000
50%       35.000000
75%       46.000000
max       90.000000
Name: age, dtype: float64
Race: White, sex: Male
count    19174.000000
mean      39.652498
std       13.436029
min       17.000000
25%       29.000000
50%       38.000000
75%       49.000000
max       90.000000
Name: age, dtype: float64

```

8. Among whom the proportion of those who earn a lot(>50K) is more: among married or single men (*marital-status* feature)? Consider married those who have a *marital-status* starting with *Married* (Married-civ-spouse, Married-spouse-absent or Married-AF-spouse), the rest are considered bachelors.

In [9]:

```

data.loc[(data['sex'] == 'Male') &
         (data['marital-status'].isin(['Never-married',
                                       'Separated',
                                       'Divorced',
                                       'Widowed']))], 'salary'].value_counts()

```

Out[9]:

```

<=50K      7552
>50K         697
Name: salary, dtype: int64

```

In [10]:

```
data.loc[(data['sex'] == 'Male') &
         (data['marital-status'].str.startswith('Married')), 'salary'].value_counts()
```

Out[10]:

```
<=50K      7576
>50K       5965
Name: salary, dtype: int64
```

In [11]:

```
data['marital-status'].value_counts()
```

Out[11]:

```
Married-civ-spouse      14976
Never-married           10683
Divorced                 4443
Separated                1025
Widowed                  993
Married-spouse-absent    418
Married-AF-spouse        23
Name: marital-status, dtype: int64
```

It's good to be married :)

9. What is the maximum number of hours a person works per week (*hours-per-week* feature)? How many people work such a number of hours and what is the percentage of those who earn a lot among them?

In [12]:

```
max_load = data['hours-per-week'].max()
print("Max time - {0} hours./week.".format(max_load))

num_workaholics = data[data['hours-per-week'] == max_load].shape[0]
print("Total number of such hard workers {0}".format(num_workaholics))

rich_share = float((data[(data['hours-per-week'] == max_load)
                        & (data['salary'] == '>50K')].shape[0]) / num_workaholics)
print("Percentage of rich among them {0}%".format(int(100 * rich_share)))
```

```
Max time - 99 hours./week.
Total number of such hard workers 85
Percentage of rich among them 29%
```

10. Count the average time of work (*hours-per-week*) those who earning a little and a lot (*salary*) for each country (*native-country*).

Simple method:

In [13]:

```
for (country, salary), sub_df in data.groupby(['native-country', 'salary']):
    print(country, salary, round(sub_df['hours-per-week'].mean(), 2))
```

```
? <=50K 40.16
? >50K 45.55
Cambodia <=50K 41.42
Cambodia >50K 40.0
Canada <=50K 37.91
Canada >50K 45.64
China <=50K 37.38
China >50K 38.9
Columbia <=50K 38.68
Columbia >50K 50.0
Cuba <=50K 37.99
Cuba >50K 42.44
Dominican-Republic <=50K 42.34
```

Dominican-Republic >50K 47.0
Ecuador <=50K 38.04
Ecuador >50K 48.75
El-Salvador <=50K 36.03
El-Salvador >50K 45.0
England <=50K 40.48
England >50K 44.53
France <=50K 41.06
France >50K 50.75
Germany <=50K 39.14
Germany >50K 44.98
Greece <=50K 41.81
Greece >50K 50.62
Guatemala <=50K 39.36
Guatemala >50K 36.67
Haiti <=50K 36.33
Haiti >50K 42.75
Holand-Netherlands <=50K 40.0
Honduras <=50K 34.33
Honduras >50K 60.0
Hong <=50K 39.14
Hong >50K 45.0
Hungary <=50K 31.3
Hungary >50K 50.0
India <=50K 38.23
India >50K 46.48
Iran <=50K 41.44
Iran >50K 47.5
Ireland <=50K 40.95
Ireland >50K 48.0
Italy <=50K 39.62
Italy >50K 45.4
Jamaica <=50K 38.24
Jamaica >50K 41.1
Japan <=50K 41.0
Japan >50K 47.96
Laos <=50K 40.38
Laos >50K 40.0
Mexico <=50K 40.0
Mexico >50K 46.58
Nicaragua <=50K 36.09
Nicaragua >50K 37.5
Outlying-US (Guam-USVI-etc) <=50K 41.86
Peru <=50K 35.07
Peru >50K 40.0
Philippines <=50K 38.07
Philippines >50K 43.03
Poland <=50K 38.17
Poland >50K 39.0
Portugal <=50K 41.94
Portugal >50K 41.5
Puerto-Rico <=50K 38.47
Puerto-Rico >50K 39.42
Scotland <=50K 39.44
Scotland >50K 46.67
South <=50K 40.16
South >50K 51.44
Taiwan <=50K 33.77
Taiwan >50K 46.8
Thailand <=50K 42.87
Thailand >50K 58.33
Trinidad&Tobago <=50K 37.06
Trinidad&Tobago >50K 40.0
United-States <=50K 38.8
United-States >50K 45.51
Vietnam <=50K 37.19
Vietnam >50K 39.2
Yugoslavia <=50K 41.6
Yugoslavia >50K 49.5

Elegant method:

In [14]:

```
pd.crosstab(data['native-country'], data['salary'],
            values=data['hours-per-week'], aggfunc=np.mean).T
```

Out[14]:

native-country	?	Cambodia	Canada	China	Columbia	Cuba	Dominican-Republic	Ecuador	El-Salvador	England	Fra
salary											
<=50K	40.164760	41.416667	37.914634	37.381818	38.684211	37.985714	42.338235	38.041667	36.030928	40.483333	41.058
>50K	45.547945	40.000000	45.641026	38.900000	50.000000	42.440000	47.000000	48.750000	45.000000	44.533333	50.750