

Double-click (or enter) to edit

Project: Machine Learning of Salary and Demographic Factors Name: Shaohua Feng Supervisor:

Double-click (or enter) to edit

```
1 from google.colab import drive
2
3 # Mount Google Drive
4 drive.mount('/content/drive')

Mounted at /content/drive

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt

1 # read in data loaded in google drive
2 file_path_1 = '/content/drive/My Drive/adult.data'
3 adult_1= pd.read_csv(file_path_1,header=None)
4 file_path_2 = '/content/drive/My Drive/adult.test.txt'
5 adult_2= pd.read_csv(file_path_2,header=None)
6 adult=pd.concat([adult_1, adult_2], ignore_index=True)

1 # add column names
2 cols=['age','workclass','fnlwtg','education','education-num','marital-status','occupation','relationship','race','sex','capital-gain','cap:
3 adult.columns=cols
4 adult.head(10)
```

	age	workclass	fnlwtg	education	education-num	marital-status	occupation	relationship	race
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	Wt
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	Wt
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	Wt
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Blk
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Blk
5	37	Private	284582	Masters	14	Married-civ-spouse	Exec-managerial	Wife	Wt
6	49	Private	160187	9th	5	Married-spouse-absent	Other-service	Not-in-family	Blk
7	52	Self-emp-not-inc	209642	HS-grad	9	Married-civ-spouse	Exec-managerial	Husband	Wt

```
1 # add y column to data frame, y=1 for label '>50k' and y=0 for label '<=50k'
2 adult['y']=np.where(adult['label']==' >50K',1,0)
3 adult.head(20)
```

	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
5	37	Private	284582	Masters	14	Married-civ-spouse	Exec-managerial	Wife
6	49	Private	160187	9th	5	Married-spouse-absent	Other-service	Not-in-family
7	52	Self-emp-not-inc	209642	HS-grad	9	Married-civ-spouse	Exec-managerial	Husband
8	31	Private	45781	Masters	14	Never-married	Prof-specialty	Not-in-family
9	42	Private	159449	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband
10	37	Private	280464	Some-college	10	Married-civ-spouse	Exec-managerial	Husband
11	30	State-gov	141297	Bachelors	13	Married-civ-spouse	Prof-specialty	Husband
12	23	Private	122272	Bachelors	13	Never-married	Adm-clerical	Own-child
13	32	Private	205019	Assoc-acdm	12	Never-married	Sales	Not-in-family
14	40	Private	121772	Assoc-voc	11	Married-civ-spouse	Craft-repair	Husband

```
1 print(adult.describe())
2 adult.dtypes
3 adult.info()
```

	age	fnlwgt	education-num	capital-gain	capital-loss	\
count	48842.000000	4.884200e+04	48842.000000	48842.000000	48842.000000	
mean	38.643585	1.896641e+05	10.078089	1079.067626	87.502314	
std	13.710510	1.056040e+05	2.570973	7452.019058	403.004552	
min	17.000000	1.228500e+04	1.000000	0.000000	0.000000	
25%	28.000000	1.175505e+05	9.000000	0.000000	0.000000	
50%	37.000000	1.781445e+05	10.000000	0.000000	0.000000	
75%	48.000000	2.376420e+05	12.000000	0.000000	0.000000	
max	90.000000	1.490400e+06	16.000000	99999.000000	4356.000000	

  

	hours-per-week	y
count	48842.000000	48842.000000
mean	40.422382	0.160538

```

std      12.391444      0.367108
min       1.000000      0.000000
25%      40.000000      0.000000
50%      40.000000      0.000000
75%      45.000000      0.000000
max      99.000000      1.000000
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 48842 entries, 0 to 48841
Data columns (total 16 columns):
#   Column                Non-Null Count  Dtype
---  -
0   age                   48842 non-null  int64
1   workclass             48842 non-null  object
2   fnlwgt               48842 non-null  int64
3   education             48842 non-null  object
4   education-num         48842 non-null  int64
5   marital-status        48842 non-null  object
6   occupation            48842 non-null  object
7   relationship          48842 non-null  object
8   race                  48842 non-null  object
9   sex                   48842 non-null  object
10  capital-gain           48842 non-null  int64
11  capital-loss           48842 non-null  int64
12  hours-per-week         48842 non-null  int64
13  native-country         48842 non-null  object
14  label                  48842 non-null  object
15  y                       48842 non-null  int64
dtypes: int64(7), object(9)
memory usage: 6.0+ MB

```

```

1 # explore: find factor levels
2 print(adult['workclass'].unique())
3 print(adult['occupation'].unique())
4 print(adult['native-country'].unique())
5 print(type(adult['occupation']))

[' State-gov' ' Self-emp-not-inc' ' Private' ' Federal-gov' ' Local-gov'
 ' ?' ' Self-emp-inc' ' Without-pay' ' Never-worked']
[' 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']
[' 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']
<class 'pandas.core.series.Series'>

```

```

1 adult['workclass']=adult['workclass'].replace(' ?',None)
2 adult['occupation']=adult['occupation'].replace(' ?',None)
3 adult['native-country']=adult['native-country'].replace(' ?',None)

```

```

1 print(adult['occupation'].unique())
2 print(adult['occupation'].unique())
3 print(adult['occupation'].unique())

[' Adm-clerical' ' Exec-managerial' ' Handlers-cleaners' ' Prof-specialty'
 ' Other-service' ' Sales' ' Craft-repair' ' Transport-moving'
 ' Farming-fishing' ' Machine-op-inspct' ' Tech-support' None
 ' Protective-serv' ' Armed-Forces' ' Priv-house-serv']
[' Adm-clerical' ' Exec-managerial' ' Handlers-cleaners' ' Prof-specialty'
 ' Other-service' ' Sales' ' Craft-repair' ' Transport-moving'
 ' Farming-fishing' ' Machine-op-inspct' ' Tech-support' None
 ' Protective-serv' ' Armed-Forces' ' Priv-house-serv']
[' Adm-clerical' ' Exec-managerial' ' Handlers-cleaners' ' Prof-specialty'
 ' Other-service' ' Sales' ' Craft-repair' ' Transport-moving'
 ' Farming-fishing' ' Machine-op-inspct' ' Tech-support' None
 ' Protective-serv' ' Armed-Forces' ' Priv-house-serv']

```

```

1 # check how many missing vales in columns workclass, occupation and native-country
2 print(adult['workclass'].isnull().sum())
3 print(adult['occupation'].isnull().sum())
4 print(adult['native-country'].isnull().sum())

```

```
2799
2809
857
```

```
1 # characterical columns
2 cols_cat=['workclass','fnlwgt','education','marital-status','occupation','relationship','race','sex','native-country']
3
4 for x in cols_cat:
5     adult[x] = adult[x].astype('category')
6     #print(x)
7
8 adult.dtypes
```

```
age                int64
workclass          category
fnlwgt             category
education          category
education-num      int64
marital-status     category
occupation         category
relationship       category
race              category
sex               category
capital-gain       int64
capital-loss       int64
hours-per-week     int64
native-country     category
label              object
y                 int64
dtype: object
```

```
1 # delete missing value
2 adult_cleaned=adult.dropna()
3 print(len(adult_cleaned))
```

```
45222
```

```
1 print("Check for NaN values:")
2 print(adult_cleaned.isna().any())
```

```
Check for NaN values:
age                False
workclass          False
fnlwgt             False
education          False
education-num      False
marital-status     False
occupation         False
relationship       False
race              False
sex               False
capital-gain       False
capital-loss       False
hours-per-week     False
native-country     False
label              False
y                 False
dtype: bool
```

```
1 # Grouped by factors
2 factor_cols=['workclass','education','marital-status','occupation','relationship','race','sex','native-country']
3
4 for factor_col in factor_cols:
5     # Group by the current factor column and calculate the mean
6     grouped_data = adult.groupby(factor_col)['y'].mean().reset_index()
7
8     # Sort the grouped data by percentage of salary>50,000
9     sorted_data = grouped_data.sort_values(by='y', ascending=False)
10    # Print the results
11    print(f"Grouped by {factor_col}:\n{sorted_data}\n")
12
13    # plot the sorted data
14    #plt.bar(sorted_data[factor_col], sorted_data['y'])
15    #plt.xlabel(f'factor_col')
16    #plt.ylabel('% salary>50,000')
17    #plt.title(f'Group Mean from Highest to Lowest')
18    #plt.show
19
```

```
20
21
22 #####
23 #grouped_df = df.groupby('Category')['Value'].mean().reset_index()
24
25 # Sort the DataFrame by mean values
26 #sorted_df = grouped_df.sort_values(by='Value', ascending=False)
27
28 # Plot the sorted data
29 #plt.bar(sorted_df['Category'], sorted_df['Value'])
30 #plt.xlabel('Category')
31 #plt.ylabel('Mean Value')
32 #plt.title('Group Mean from Highest to Lowest')
33 #plt.show()
34
35
```

```

Grouped by workclass:
  workclass      y
4    Self-emp-inc 0.366962
0    Federal-gov  0.259078
1    Local-gov   0.196747
5    Self-emp-not-inc 0.187468
6    State-gov   0.178193
3    Private     0.146375
2    Never-worked 0.000000
7    Without-pay 0.000000

```

```

Grouped by education:
  education      y
10  Doctorate    0.515152
14  Prof-school  0.507194
12  Masters      0.360933
9   Bachelors    0.276760
8   Assoc-voc    0.175158
7   Assoc-acdm   0.165522
15  Some-college 0.127505
11  HS-grad      0.106120
2   12th         0.050228
0   10th         0.044636
5   7th-8th      0.041885
6   9th          0.035714
1   11th         0.033113
4   5th-6th      0.031434
3   1st-4th      0.024291
13  Preschool    0.000000

```

```

Grouped by marital-status:
  marital-status      y
2    Married-civ-spouse 0.299030
1    Married-AF-spouse  0.270270
0          Divorced     0.069803
6          Widowed      0.055995
3    Married-spouse-absent 0.054140
5          Separated     0.043137
4          Never-married 0.030465

```

```

Grouped by occupation:
  occupation      y
3    Exec-managerial 0.323365
9    Prof-specialty  0.301199
10   Protective-serv 0.214649
12   Tech-support    0.195712
11   Sales           0.178597
2    Craft-repair    0.151996
13   Transport-moving 0.135881
0    Adm-clerical    0.090358
6    Machine-op-inspct 0.082727
4    Farming-fishing 0.077181
1    Armed-Forces    0.066667
5    Handlers-cleaners 0.041506
7    Other-service    0.027829
8    Priv-house-serv  0.004132

```

```

1 # I am very interested in the relationship between salary level and race and education combination
2 factor_columns = ['race', 'education']
3
4 for combination in adult.groupby(factor_columns):
5     group_name = combination[0]
6     group_data = combination[1]
7     mean_salary = group_data['y'].mean()
8
9     print(f"Factors: {'', '.join(f'{col}={val}' for col, val in zip(factor_columns, group_name))} | Salary>50,000: {mean_salary}")
10

```

```

Factors: race= Amer-Indian-Eskimo, education= 10th | Salary>50,000: 0.0
Factors: race= Amer-Indian-Eskimo, education= 11th | Salary>50,000: 0.07692307692307693
Factors: race= Amer-Indian-Eskimo, education= 12th | Salary>50,000: 0.0
Factors: race= Amer-Indian-Eskimo, education= 1st-4th | Salary>50,000: 0.0
Factors: race= Amer-Indian-Eskimo, education= 5th-6th | Salary>50,000: 0.0
Factors: race= Amer-Indian-Eskimo, education= 7th-8th | Salary>50,000: 0.0
Factors: race= Amer-Indian-Eskimo, education= 9th | Salary>50,000: 0.0
Factors: race= Amer-Indian-Eskimo, education= Assoc-acdm | Salary>50,000: 0.07692307692307693
Factors: race= Amer-Indian-Eskimo, education= Assoc-voc | Salary>50,000: 0.03225806451612903
Factors: race= Amer-Indian-Eskimo, education= Bachelors | Salary>50,000: 0.27586206896551724
Factors: race= Amer-Indian-Eskimo, education= Doctorate | Salary>50,000: 0.6666666666666666
Factors: race= Amer-Indian-Eskimo, education= HS-grad | Salary>50,000: 0.0625
Factors: race= Amer-Indian-Eskimo, education= Masters | Salary>50,000: 0.23076923076923078
Factors: race= Amer-Indian-Eskimo, education= Preschool | Salary>50,000: 0.0
Factors: race= Amer-Indian-Eskimo, education= Prof-school | Salary>50,000: 1.0
Factors: race= Amer-Indian-Eskimo, education= Some-college | Salary>50,000: 0.04838709677419355
Factors: race= Asian-Pac-Islander, education= 10th | Salary>50,000: 0.0625
Factors: race= Asian-Pac-Islander, education= 11th | Salary>50,000: 0.037037037037037035
Factors: race= Asian-Pac-Islander, education= 12th | Salary>50,000: 0.06666666666666667
Factors: race= Asian-Pac-Islander, education= 1st-4th | Salary>50,000: 0.0
Factors: race= Asian-Pac-Islander, education= 5th-6th | Salary>50,000: 0.10714285714285714
Factors: race= Asian-Pac-Islander, education= 7th-8th | Salary>50,000: 0.0
Factors: race= Asian-Pac-Islander, education= 9th | Salary>50,000: 0.1
Factors: race= Asian-Pac-Islander, education= Assoc-acdm | Salary>50,000: 0.16326530612244897
Factors: race= Asian-Pac-Islander, education= Assoc-voc | Salary>50,000: 0.16981132075471697
Factors: race= Asian-Pac-Islander, education= Bachelors | Salary>50,000: 0.23774509803921567
Factors: race= Asian-Pac-Islander, education= Doctorate | Salary>50,000: 0.391304347826087
Factors: race= Asian-Pac-Islander, education= HS-grad | Salary>50,000: 0.10119047619047619
Factors: race= Asian-Pac-Islander, education= Masters | Salary>50,000: 0.30714285714285716
Factors: race= Asian-Pac-Islander, education= Preschool | Salary>50,000: 0.0
Factors: race= Asian-Pac-Islander, education= Prof-school | Salary>50,000: 0.46551724137931033
Factors: race= Asian-Pac-Islander, education= Some-college | Salary>50,000: 0.10927152317880795
Factors: race= Black, education= 10th | Salary>50,000: 0.03296703296703297
Factors: race= Black, education= 11th | Salary>50,000: 0.027777777777777776
Factors: race= Black, education= 12th | Salary>50,000: 0.047619047619047616
Factors: race= Black, education= 1st-4th | Salary>50,000: 0.041666666666666664
Factors: race= Black, education= 5th-6th | Salary>50,000: 0.0
Factors: race= Black, education= 7th-8th | Salary>50,000: 0.022222222222222223
Factors: race= Black, education= 9th | Salary>50,000: 0.036036036036036036
Factors: race= Black, education= Assoc-acdm | Salary>50,000: 0.11801242236024845
Factors: race= Black, education= Assoc-voc | Salary>50,000: 0.10909090909090909
Factors: race= Black, education= Bachelors | Salary>50,000: 0.19047619047619047
Factors: race= Black, education= Doctorate | Salary>50,000: 0.5625
Factors: race= Black, education= HS-grad | Salary>50,000: 0.048314606741573035
Factors: race= Black, education= Masters | Salary>50,000: 0.27972027972027974
Factors: race= Black, education= Preschool | Salary>50,000: 0.0
Factors: race= Black, education= Prof-school | Salary>50,000: 0.38095238095238093
Factors: race= Black, education= Some-college | Salary>50,000: 0.07977736549165121
Factors: race= Other, education= 10th | Salary>50,000: 0.09090909090909091
Factors: race= Other, education= 11th | Salary>50,000: 0.0
Factors: race= Other, education= 12th | Salary>50,000: 0.0
Factors: race= Other, education= 1st-4th | Salary>50,000: 0.0
Factors: race= Other, education= 5th-6th | Salary>50,000: 0.043478260869565216
Factors: race= Other, education= 7th-8th | Salary>50,000: 0.0
Factors: race= Other, education= 9th | Salary>50,000: 0.0
Factors: race= Other, education= Assoc-acdm | Salary>50,000: 0.2
Factors: race= Other, education= Assoc-voc | Salary>50,000: 0.0
Factors: race= Other, education= Bachelors | Salary>50,000: 0.1

```

```

1 # correlation matrix
2 cor_matrix = adult_cleaned.corr()
3 print(cor_matrix)

```

```

          age    fnlwtg  education-num  capital-gain  capital-loss  \
age          1.000000 -0.075792      0.037623      0.079683      0.059351
fnlwtg      -0.075792  1.000000     -0.041993     -0.004110     -0.004349
education-num  0.037623 -0.041993  1.000000      0.126907      0.081711
capital-gain   0.079683 -0.004110   0.126907      1.000000     -0.032102
capital-loss   0.059351 -0.004349   0.081711     -0.032102      1.000000
hours-per-week 0.101992 -0.018679   0.146206      0.083880      0.054195

```