

Using KNN to do Predictions in the Adult Dataset

PMR3508: Machine Learning and Pattern Recognition

Student name: Rodolfo Lima (PMR3508-2021-6)

0- Libraries

```
In [31]: from IPython.core.display import display, HTML
display(HTML("<style>.container { width:100% !important; }</style>"))

import ssl
ssl._create_default_https_context = ssl._create_unverified_context

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sea

#sklearn functions
from sklearn.preprocessing import RobustScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import cross_val_score
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import accuracy_score
```

1- Data Cleaning

```
In [93]: #Training data
df_train = pd.read_csv('train_data.csv')
df_train = df_train.set_index('Id')

#testing Data
df_test = pd.read_csv("adult.test.txt",
                      names=[
                          'age', 'workclass', 'fnlwgt', 'education', 'education.num',
                          'marital.status', 'occupation', 'relationship', 'race', 'sex',
                          'capital.gain', 'capital.loss', 'hours.per.week', 'native.country',
                          'income'],
                      sep='|',\n",
                      engine='python')
df_test.drop(df_test.index[0], axis = 0, inplace = True)
df_train.head()
```

```
Out[93]:
```

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss
16280	34	Private	204991	Some-college	10	Divorced	Exec-managerial	Own-child	White	Male	0	
16281	58	Local-gov	310085	10th	6	Married-civ-spouse	Transport-moving	Husband	White	Male	0	
16282	25	Private	146117	Some-college	10	Never-married	Machine-op-inspct	Not-in-family	White	Male	0	
16283	24	Private	138938	Some-college	10	Divorced	Adm-clerical	Not-in-family	White	Female	0	
16284	57	Self-emp-inc	258883	HS-grad	9	Married-civ-spouse	Transport-moving	Husband	White	Male	5178	

1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 84

Feature education values: 'Assoc-voc', 'HS-grad', 'Bachelors', '12th', '10th', 'Masters', '7th-8th', 'Some-college', 'not-high', 'Preschool', '11th', 'Assoc-acdm', 'Doctorate', '9th', '5th-6th', 'Prof-school'

Feature education.num values: (1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16)

Feature marital.status values: {'Married-spouse-absent', 'Divorced', 'Separated', 'Widowed', 'Married-spouse', 'Never-married', 'Married-civ-spouse'}

Feature occupation values: {'Machine-op-inspct', 'Handlers-cleaners', 'Craft-repair', 'Priv-house-serv', 'Farming-fishing', '?', 'Sales', 'Exec-managerial', 'Tech-support', 'Armed-Forces', 'Protective-serv', 'Adm-clerical', 'Prof-specialty', 'Other-service', 'Transport-moving', 'Transport-moving'}

Feature relationship values: {'Husband', 'Own-child', 'Not-in-family', 'Unmarried', 'Wife', 'Other-relative'}

Feature race values: {'Asian-Pac-Islander', 'Black', 'Other', 'Amer-Indian-Eskimo', 'White'}

Feature sex values: {'Female', 'Male'}

Feature capital.gain values: (0, 2050, 4101, 7688, 14344, 2062, 2580, 22040, 3103, 1055, 25124, 2597, 4650, 18481, 2105, 5178, 1086, 3137, 6723, 2635, 4687, 594, 20051, 1111, 5721, 3674, 2653, 1639, 676, 114, 2174, 1151, 2176, 2296, 25236, 1173, 2202, 98989, 5084, 15020, 15024, 4787, 27828, 2228, 684, 9, 3781, 3273, 7896, 6360, 3818, 13550, 2290, 3325, 4865, 14084, 1797, 7430, 2829, 6418, 7443, 10520, 2329, 4386, 1831, 2346, 7978, 34095, 3883, 2354, 1848, 4416, 4931, 3908, 2885, 10566, 4934, 5455, 34, 11, 25, 36, 1307, 3418, 4110, 6486, 6227, 5754

Feature hours.per.week values: (1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 2, 0, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 4, 5, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 70, 7, 2, 73, 74, 75, 76, 77, 78, 80, 81, 82, 84, 85, 86, 87, 88, 89, 90, 91, 92, 94, 95, 96, 97, 98, 99)

Feature native.country values: {'Scotland', 'Japan', 'Hungary', 'Dominican-Republic', 'Ireland', 'Nicaragua', 'Taiwan', 'Puerto-Rico', 'Laos', 'India', 'England', 'United-States', 'Yugoslavia', 'Hong', 'Thailand', 'Peru', 'Cambodia', 'South', 'Guatemala', 'France', 'Ecuador', 'Canada', 'Italy', 'Germany', 'Honduras', 'Philippines', 'Greece', 'Iran', 'Outlying-US(Guam-USVI-etc)', 'Mexico', 'Jamaica', 'China', 'Cuba', 'Trinidad&Tobago', 'I', 'Holland-Netherlands', 'Portugal', 'Columbia', 'Poland', 'Vietnam', 'El-Salvador', 'Haiti'}

Feature income values: '>50K', '<=50K'

```
In [61]: # I have to deal with "?" strings in categorical variables and possible missing data
df_train = df_train.replace('?', np.nan)
df_test = df_test.replace('?', np.nan)
rows_nan = []
for index, row in df_train.isnull().iterrows():
    if (row.values).any():
        rows_nan.append(index)
len(rows_nan) #number of rows that have any missing data

Out [61]: 2399
```

```
In [62]: print('Rows that have any missing value: ',str(len(rows_nan)/df_train.shape[0]*100) + '%')

Rows that have any missing value: 7.3679361793618%
```

```
In [63]: #owing to the small % of rows having missing data, I'll get rid of them
df_train = df_train.dropna(axis = 0)
display(df_train)

#df_test = df_test.dropna(axis = 0)
```

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss
id												
16280	34	Private	204991	Some-college	10	Divorced	Exec-managerial	Own-child	White	Male		0
16281	58	Local-gov	310085	10th	6	Married-civ-spouse	Transport-moving	Husband	White	Male		0
16282	25	Private	146117	Some-college	10	Never-married	Machine-op-inspct	Not-in-family	White	Male		0
16283	24	Private	136938	Some-college	10	Divorced	Adm-clerical	Not-in-family	White	Female		0
16284	57	Self-emp-inc	258883	HS-grad	9	Married-civ-spouse	Transport-moving	Husband	White	Male	5178	
...
48835	42	Private	384236	Masters	14	Married-civ-spouse	Prof-specialty	Husband	White	Male		7688
48836	23	Private	129042	HS-grad	9	Never-married	Machine-op-inspct	Unmarried	Black	Female		0
48837	30	Private	195488	HS-grad	9	Never-married	Priv-house-serv	Own-child	White	Female		0
48838	18	Private	27620	HS-grad	9	Never-married	Adm-clerical	Not-in-family	White	Female		0
48839	47	Local-gov	203067	Bachelors	13	Divorced	Prof-specialty	Not-in-family	White	Male		0
30161 rows × 13 columns												

```
In [64]: df_train.columns

Out [64]: Index(['age', 'workclass', 'fnlwgt', 'education', 'education.num', 'marital.status', 'occupation', 'relationship', 'race', 'sex', 'capital.gain', 'capital.loss', 'hours.per.week', 'native.country', 'income'], dtype='object')
```

```
In [65]: #I will delete education columns once I already have 'education.num'
df_train.drop(labels = 'education', axis = 1, inplace = True)
df_test.drop(labels = 'education', axis = 1, inplace = True)
```

2- Exploratory Data Analysis (EDA) and Data Preparation

2.1 - Correlation Matrix

- Our target is the income, it's a categorical variable with two label
- I'm changing the target in order to have two quantitative value: 0 to <=50 and >50 to 1. Then, I'll be able to measure the correlations among the variables

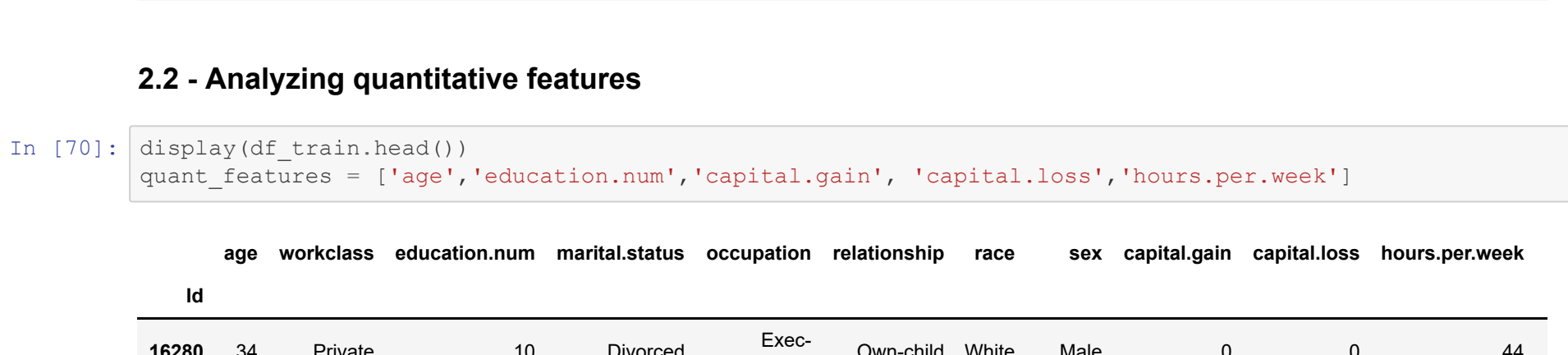
```
In [66]: df_train_corr = df_train.copy()
df_train_corr['income'] = df_train_corr['income'].apply(lambda x: 0 if x == '<=50K' else 1)
set(df_train_corr['income'])

Out [66]: {0, 1}
```

```
In [67]: sea.heatmap(round(df_train_corr.corr(),2), annot = True, cmap = 'viridis', vmin = -1, vmax = 1)
plt.title('Correlation (Spearman) Matrix', fontweight = 'bold', fontsize = 12)
```

```
plt.gcf().set_size_inches(8,8)
plt.gca().set_dpi(100)
plt.subplots_adjust(bottom = .25)
plt.xticks(size = 12)
plt.yticks(size = 12)
```

```
Out [67]: (array([0.5, 1.5, 2.5, 3.5, 4.5, 5.5, 6.5]),
<a list of 7 Text yticklabel objects>)
```



- Comparing the target with other features, we can see the pearson correlation between income and fnlwgt is close to 0. So we can consider take this feature out

```
In [68]: df_train.drop(labels = 'fnlwgt', axis = 1, inplace = True)
df_test.drop(labels = 'fnlwgt', axis = 1, inplace = True)
```

```
In [69]: df_train.head()

Out [69]:
```

	age	workclass	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss	hours.per.week
id											
16280	34	Private	10	Divorced	Exec-managerial	Own-child	White	Male	0	0	44
16281	58	Local-gov	6	Married-civ-spouse	Transport-moving	Husband	White	Male	0	0	40
16282	25	Private	10	Never-married	Machine-op-inspct	Not-in-family	White	Male	0	0	42
16283	24	Private	10	Divorced	Adm-clerical	Not-in-family	White	Female	0	0	40
16284	57	Self-emp-inc	9	Married-civ-spouse	Transport-moving	Husband	White	Male	5178	0	60

2.2 - Analyzing quantitative features

```
In [70]: display(df_train.head())
quant_features = ['age', 'education.num', 'capital.gain', 'capital.loss', 'hours.per.week']
```

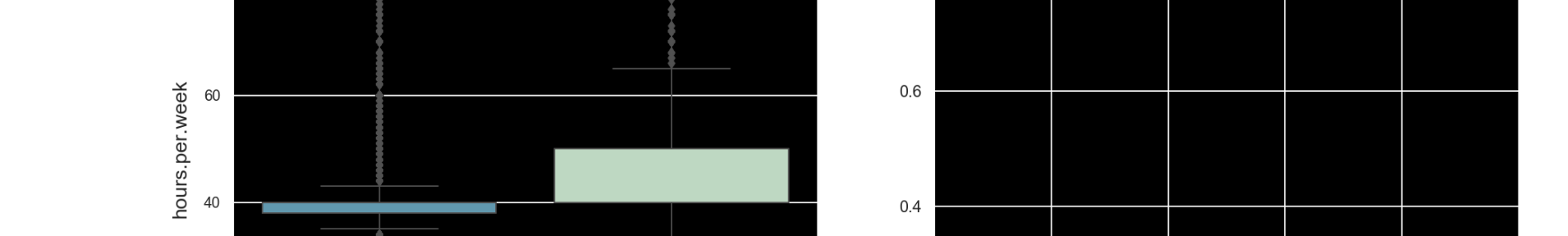
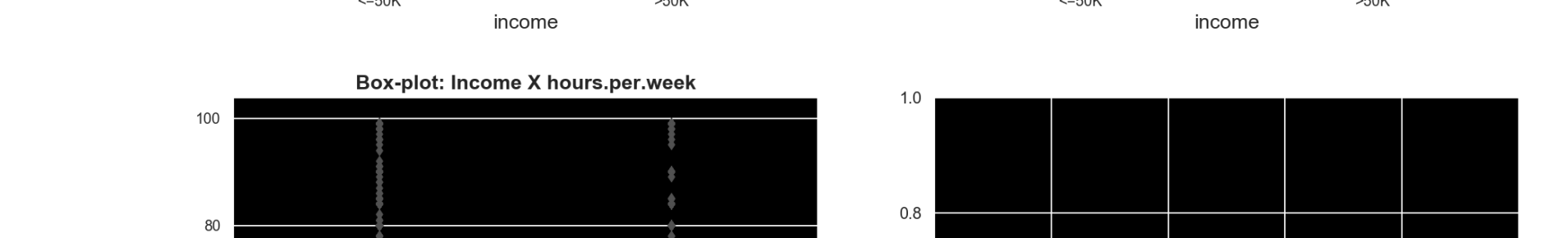
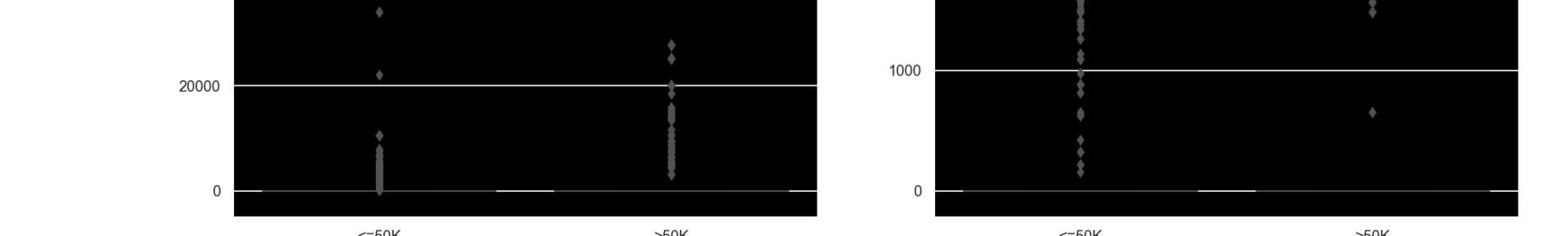
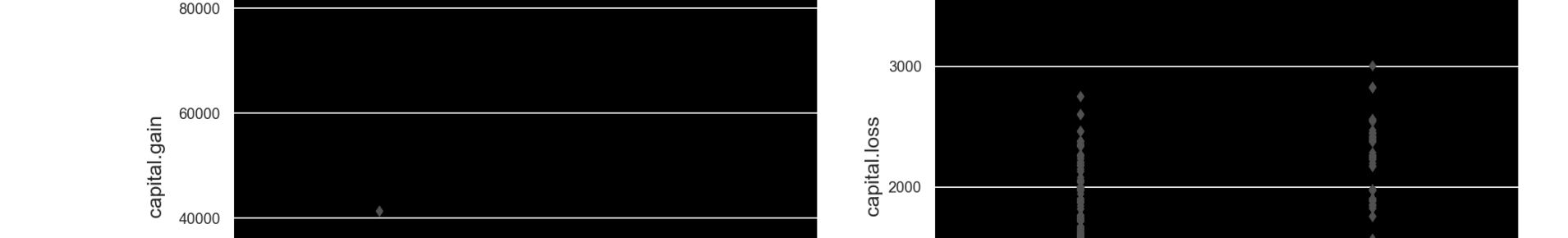
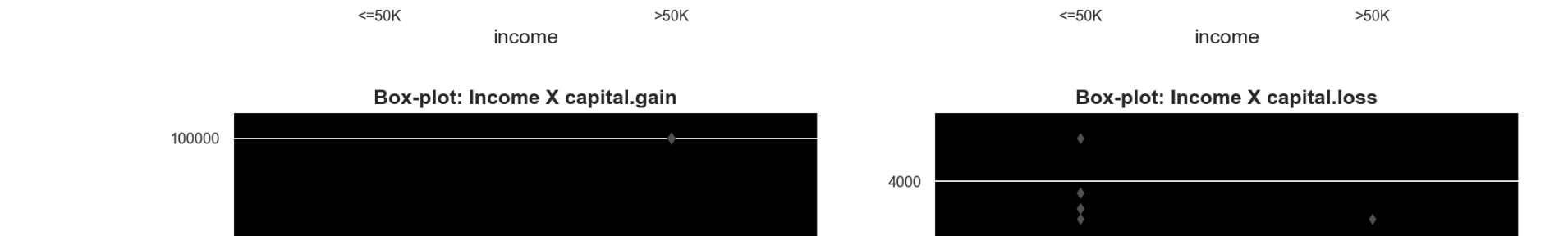
	age	workclass	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss	hours.per.week
id											
16280	34	Private	10	Divorced	Exec-managerial	Own-child	White	Male	0	0	44
16281	58	Local-gov	6	Married-civ-spouse	Transport-moving	Husband	White	Male	0	0	40
16282	25	Private	10	Never-married	Machine-op-inspct	Not-in-family	White	Male	0	0	42
16283	24	Private	10	Divorced	Adm-clerical	Not-in-family	White	Female	0	0	40
16284	57	Self-emp-inc	9	Married-civ-spouse	Transport-moving	Husband	White	Male	5178	0	60

```
In [71]: sea.set(rc = {'axes.facecolor':'black'})
fig, (ax1,ax2), (ax3,ax4), (ax5,ax6) = plt.subplots(3,2)
axes = [ax1,ax2,ax3,ax4,ax5,ax6]
```

```
for feature,ax in zip(quant_features,axes):
    #Plotting
    sea.boxplot(data = df_train, x = 'income', y = feature, palette = "GnBu_r", orient = "v", linewidth = 1, saturation=0.35, ax=ax)
    #sea.boxplot(data = df_train, x = 'income', y = feature, palette = "GnBu_r", orient = "v", size = 7,alpha=1, edgecolor='grey', ax = ax)
    ax.set_xlabel('income', fontsize = 15)
    ax.set_ylabel(feature, fontsize = 15)
    ax.set_title("Box-plot: Income X {}".format(feature), size = 15, fontweight = 'bold')
```

```
#Setting
plt.gcf().set_size_inches(17,12)
plt.gca().set_dpi(130)
plt.subplots_adjust(bottom = -0.8)
plt.xticks(fontsize = 12)
plt.yticks(fontsize = 12)
```

```
Out [71]: (array([0., 0.2, 0.4, 0.6, 0.8, 1.]), <a list of 6 Text yticklabel objects>)
```



- We can realize there are many outliers from the boxplots that need to be handled
- I'm using Robust Scaler method to deal with outliers

```
In [72]: df_train[['capital.gain', 'capital.loss', 'hours.per.week']] = RobustScaler().fit_transform(df_train[['capital.gain', 'capital.loss', 'hours.per.week']])
df_test[['capital.gain', 'capital.loss', 'hours.per.week']] = RobustScaler().fit_transform(df_test[['capital.gain', 'capital.loss', 'hours.per.week']])
```

2.3 - Analyzing qualitative features

```
In [73]: display(df_train.head())
qual_features = ['workclass', 'marital.status', 'occupation', 'relationship', 'race', 'sex', 'native.country']
```

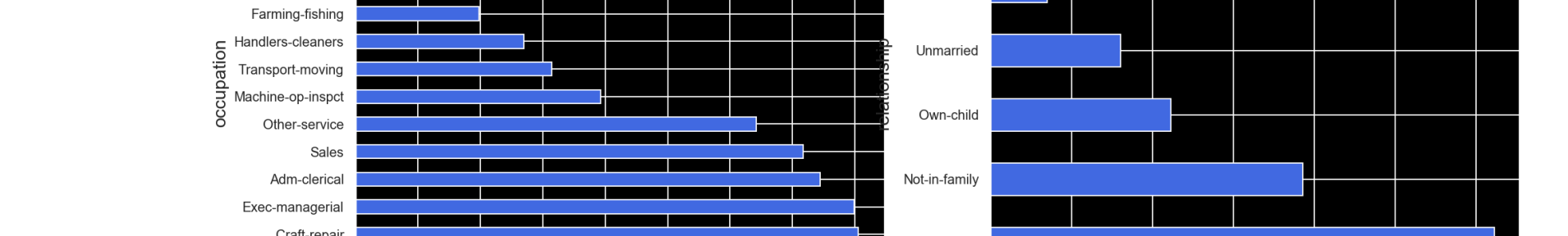
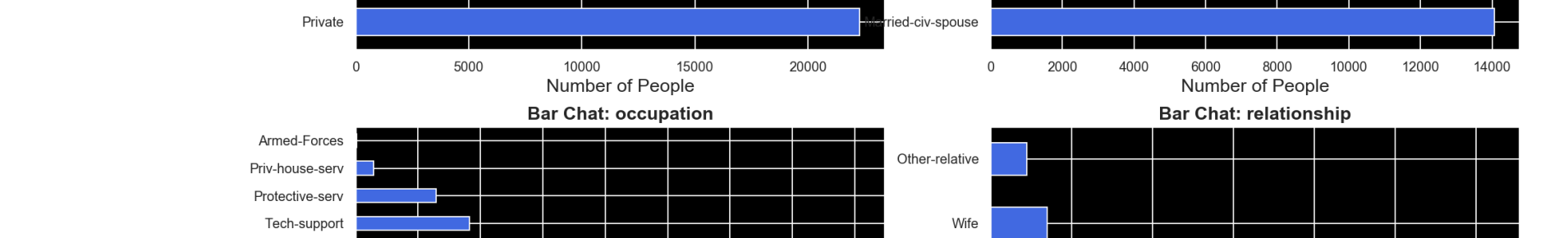
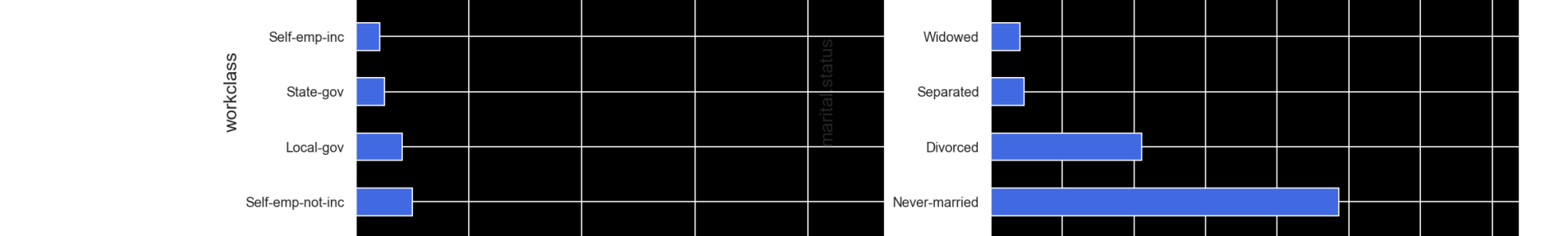
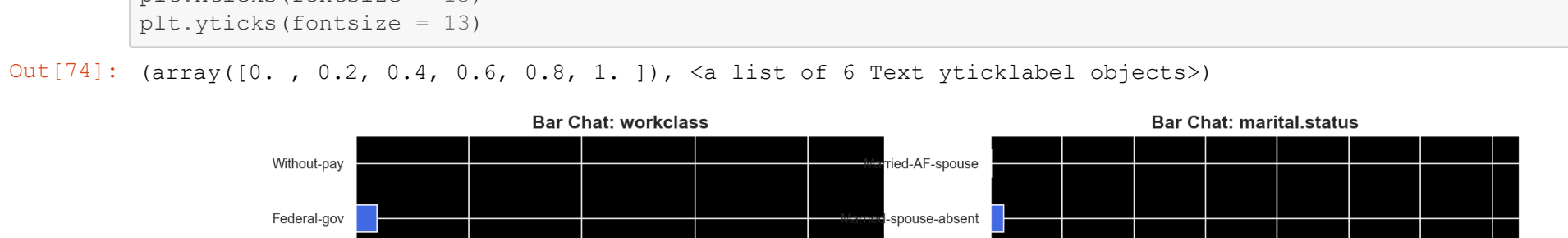
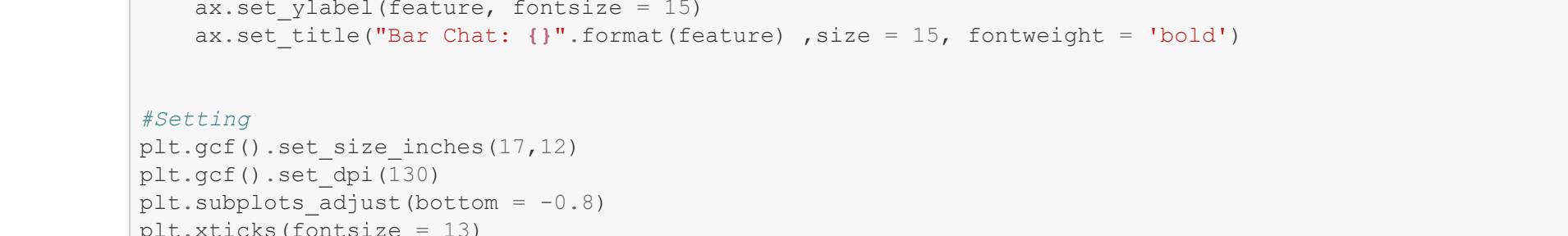
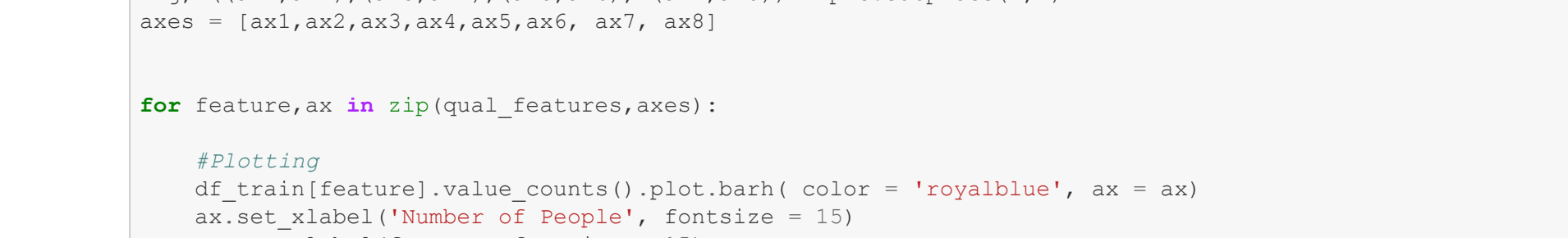
	age	workclass	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss	hours.per.week
id											
16280	34	Private	10	Divorced	Exec-managerial	Own-child	White	Male	0.0	0.0	0.8
16281	58	Local-gov	6	Married-civ-spouse	Transport-moving	Husband	White	Male	0.0	0.0	0.0
16282	25	Private	10	Never-married	Machine-op-inspct	Not-in-family	White	Male	0.0	0.0	0.4
16283	24	Private	10	Divorced	Adm-clerical	Not-in-family	White	Female	0.0	0.0	0.0
16284	57	Self-emp-inc	9	Married-civ-spouse	Transport-moving	Husband	White	Male	5178.0	0.0	4.0

```
In [74]: sea.set(rc = {'axes.facecolor':'black'})
fig, ((ax1,ax2), (ax3,ax4), (ax5,ax6), (ax7,ax8)) = plt.subplots(4,2)
axes = [ax1,ax2,ax3,ax4,ax5,ax6,ax7,ax8]
```

```
for feature,ax in zip(qual_features,axes):
    #Plotting
    df_train[feature].value_counts().plot.barh( color = 'royalblue', ax = ax)
    ax.set_xlabel('Number of People', fontsize = 15)
    ax.set_ylabel(feature, fontsize = 15)
    ax.set_title("Bar Chat: {}".format(feature), size = 15, fontweight = 'bold')
```

```
#Setting
plt.gcf().set_size_inches(17,12)
plt.gca().set_dpi(130)
plt.subplots_adjust(bottom = -0.8)
plt.xticks(fontsize = 13)
plt.yticks(fontsize = 13)
```

```
Out [74]: (array([0., 0.2, 0.4, 0.6, 0.8, 1.]), <a list of 6 Text yticklabel objects>)
```



- The fact that catch my attention the most is the too way greater number of records from the US

```
In [75]: df_train.groupby(by = 'native.country').agg({'native.country': lambda x: str(round(np.count_nonzero(x)/len(df_train) * 100,2)) + ' %'})
```

	native.country
1	Cambodia 0.06 %
2	Canada 0.35 %
3	China 0.23 %
4	Columbia 0.19 %
5	Cuba 0.31 %
6	Dominican-Republic 0.22 %
7	England 0.09 %
8	El-Salvador 0.33 %
9	France 0.29 %
10	Germany 0.42 %
11	Greece 0.1 %
12	Guatemala 0.21 %
13	Haiti 0.14 %
14	Holland-Netherlands 0.0 %
15	Honduras 0.04 %
16	Hong 0.06 %
17	Hungary 0.04 %
18	India 0.33 %
19	Indonesia 0.14 %
20	Ireland 0.08 %
21	Italy 0.23 %
22	Jamaica 0.27 %
23	Japan 0.2 %
24	Leao 0.06 %
25	Mexico 2.02 %
26	Nicaragua 0.05 %
27	Outlying-US(Guam-USVI-etc) 0.05 %
28	Peru 0.1 %
29	Philippines 0.62 %
30	Poland 0.19 %
31	Portugal 0.11 %
32	Puerto-Rico 0.36 %
33	Scotland 0.04 %
34	South 0.24 %
35	Taiwan 0.14 %
36	Thailand 0.06 %
37	Trinidad&Tobago 0.06 %
38	United-States 91.19 %
39	Vietnam 0.21 %
40	Yugoslavia 0.05 %

- We can notice that only the United-States accounts for more than 91% of the dataset records. We can conclude that the dataset is totally unbalanced
- Very well. I'll handle it placing 1 to people from the US and 0 to non people from the US

```
In [76]: df_train['USA'] = df_train['native.country'].apply(lambda x: 1 if x == 'United-States' else 0)
df_test['USA'] = df_test['native.country'].apply(lambda x: 1 if x == 'United-States' else 0)
```

- Now I have to convert the labels to numerical values in order to apply the KNN classifier

```
In [78]: df_test

Out [78]:
```

	age	workclass	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss	hours.per.week	native
1	25	Private	7.0	Never-married	Machine-op-inspct	Own-child	Black	Male	0.0	0.0	0.0	0.0
2	38	Private	9.0	Married-civ-spouse	Farming-fishing	Husband	White	Male	0.0	0.0	2.0	2.0
3	28	Local-gov	12.0	Married-civ-spouse	Protective-serv	Husband	Black	Male	0.0	0.0	0.0	0.0
4	44	Private	10.0	Married-civ-spouse	Machine-op-inspct	Husband	White	Female	7688.0	0.0	0.0	0.0
...
16277	39	Private	13.0	Divorced	Prof-specialty	Not-in-family	White	Female	0.0	0.0	0.0	-0.8
16278	64	NaN	8.0	Widowed	NaN	Other-relative	Black	Male	0.0	0.0	0.0	0.0
16279	38	Private	13.0	Married-civ-spouse	Prof-specialty	Husband	White	Male	0.0	0.0	2.0	2.0
16280	44	Private	13.0	Divorced	Adm-clerical	Own-child	Asian-Pac-Islander	Male	5455.0	0.0	0.0	0.0
16281	35	Self-emp-inc	13.0	Married-civ-spouse	Exec-managerial	Husband	White	Male	0.0	0.0	4.0	4.0

```
In [81]: df_test['workclass']

Out [81]:
```

1	Private
2	Private
3	Local-gov
4	Private
5	NaN
...	...
16277	Private
16278	NaN
16279	Private
16280	Private
16281	Self-emp-inc
Name: workclass, Length: 16281, dtype: object	

```
In [83]: #training dataset
for feature in qual_features:
    df_train[feature] = LabelEncoder().fit_transform(df_train[feature])

#test dataset
for feature in qual_features:
    print(feature)
    df_test[feature] = LabelEncoder().fit_transform(df_test[feature]).astype(str)
```

workclass
marital.status
occupation
relationship
race
sex
native.country

3- Using KNN (K-Nearest Nighbor) Classifier to Predict

3.1 - Feature Selection

```
In [84]: df_train.columns

Out [84]: Index(['age', 'workclass', 'education.num', 'marital.status', 'occupation', 'relationship', 'race', 'sex', 'capital.gain', 'capital.loss', 'hours.per.week', 'native.country', 'income', 'USA'], dtype='object')
```

```
In [85]: feature_select = ['age', 'workclass', 'education.num', 'marital.status', 'occupation', 'relationship', 'race', 'sex', 'capital.gain', 'capital.loss', 'hours.per.week', 'USA']
target = 'income'

x_train = df_train[feature_select]
y_train = df_train[target]
x_test = df_test[feature_select]
y_test = df_test[target]
```

3.2 Apply the Classifier to the training dataset

```
In [86]: #I'm using cross validation method with 10 folders to select K that presents more accuracy
from sklearn import cross_validation
score = cross_val_score(KNeighborsClassifier(n_neighbors = K), x_train, y_train, cv = 10, scoring = "accuracy").mean()
print('K-Nearest Neighbor: (K) | Accuracy: (score)')
```

K-Neighbor: 10 Accuracy: 0.85163284531966
K-Neighbor: 15 Accuracy: 0.85132768339377
K-Neighbor: 20 Accuracy: 0.850767659287716
K-Neighbor: 25 Accuracy: 0.849908320670976
K-Neighbor: 30 Accuracy: 0.850336866227574
K-Neighbor: 35 Accuracy: