

Names: Christa Mauldin Shofa
Numpy: 200537014

Kelas: D018
Tugas 1

File adult data berisikan data orang-orang dengan berbagai karakteristik. Dari karakteristik setiap orang, kemudian dianotasikan penghasilan mereka apakah di bawah 50 ribu dolar atau di atas 50 ribu dolar per tahunnya. Dari data ini kita bisa melatih mesin untuk memprediksi apakah seseorang memiliki penghasilan di bawah 50 ribu dolar atau di atas 50 ribu dolar dari karakteristik datanya. What you need to do are:

1. Understand the data. Tell us about the data. Visualize it. Describe it.
2. Lakukan penambangan data menggunakan SVM dan Neural Network. Bagaimana hasilnya. Mana kah yang memberikan performance yang lebih baik, apa sebabnya.

```
In [40]: # Import
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import sklearn.preprocessing as scp
from numpy.polynomial.polynomial import polyfit
from sklearn.model_selection import train_test_split
from sklearn import svm as svm
from sklearn.neural_network import MLPClassifier as MLPClassifier
from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score
from sklearn.preprocessing import StandardScaler as StandardScaler
import math
```

1. Understand the data. Tell us about the data. Visualize it. Describe it.

Data yang digunakan adalah adult.csv. File tersebut berisikan data mengenai informasi diri banyak orang yang akan dicari hubungannya dengan penghasilan mereka dalam satu tahun dalam mata uang dolar. Informasi diri terdiri dari umur, kelas pekerjaan, fnlwgt value, edukasi, status pernikahan, pekerjaan, hubungan, ras, jenis kelamin, capital-gain, capital-loss, jam kerja per minggu, asal negara, dan threshold pemasukan per tahun di bawah 50 ribu dolar atau di atasnya.

```
In [41]: df = pd.read_csv('adult.csv',
sep=',',
names=[ 'age', 'workclass', 'fnlwgt', 'education', 'education-num', 'marital-status', 'occupation',
'relationship', 'race', 'sex', 'capital-gain', 'capital-loss', 'hours-per-week', 'native-country',
'income'],
engine='python')
print('The shape of this data is ' + str(df.shape))
display(df.head())
```

The shape of this data is (32561, 15)

	age	workclass	fnlwgt	education	education-num	marital-status	occupation	relationship	race	sex	capital-gain	capital-loss	hours-per-week	native-country	income
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	40	United-States	<=50K
4	28	Private	138409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	40	Cuba	<=50K

Membersihkan data kategorikal

Pada tugas kali ini, kita memiliki acuan value data yang benar yaitu melalui dokumen `adult.names`. Value value yang didapat dari `adult.csv` masih belum sesuai dengan dokumen tersebut. Oleh karena itu kita perlu mencocokkan value yang ada dengan yang seharusnya ada.

Cara membersihkannya datanya adalah dengan mencari tahu berapa jumlah tiap value kategorikal pada setiap kolom. Dengan kode yang sederhana nupa dirangka agar menampilkan hasil yang diinginkan, kita bisa melihat apakah ada value yang tidak seharusnya berada di sana. Misalnya, value yang salah ketik, hanya berisi tanda tanya (?), dan lain-lain.

Value value tersebut tidak akan saya buang, melainkan saya akan mengganti valuenya dengan modus data. Di bawah ini adalah pengecekan data tidak sesuai dan pembersihannya.

Kolom 'workclass'

Value yang seharusnya: Private, Self-emp-not-inc, Self-emp-inc, Federal-gov, Local-gov, State-gov, Without-pay, Never-worked. Value tidak sesuai: "?" dan "Private".

Cara membersihkannya:

- "Private" akan diubah menjadi "Private" (asumsi salah ketik)
- "?" diubah menjadi modus dari data workclass, yaitu "Private"

Di bawah ini adalah before dan after hasil pembersihan datanya.

```
In [41]: # ===== BEFORE =====
display(df.groupby('workclass').size())

workclass
? 1836
Federal-gov 960
Local-gov 2093
Never-worked 7
Private 22695
Private 546
Self-emp-inc 1116
Self-emp-not-inc 2541
State-gov 1298
Without-pay 14
dtype: int64
```

```
In [42]: # ===== AFTER =====
df['workclass'] = df['workclass'].replace(["?", "Private"], "Private")
display(df.groupby('workclass').size())

workclass
Federal-gov 960
Local-gov 2093
Never-worked 7
Private 24532
Self-emp-inc 1116
Self-emp-not-inc 2541
State-gov 1298
Without-pay 14
dtype: int64
```

Kolom 'education'

Value yang seharusnya: Bachelors, Some-college, 11th, HS-grad, Prof-school, Assoc-acdm, Assoc-voc, 9th, 7th-8th, 12th, Masters, 1st-4th, 10th, Doctorate, 5th-6th, Preschool.

Value tidak sesuai: "?" dan "Some-colleges".

Cara membersihkannya: > "Some-colleges" akan diubah menjadi "Some-college" (asumsi salah ketik)

> "?" diubah menjadi modusnya yaitu "HS-grad"

Di bawah ini adalah before dan after hasil pembersihan datanya.

```
In [43]: # BEFORE
df.groupby('education').size()

Out[43]: education
10th 933
11th 1175
12th 433
1st-4th 168
5th-6th 333
7th-8th 646
9th 514
? 1
Assoc-acdm 1867
Assoc-voc 1382
Bachelors 5354
Doctorate 413
HS-grad 18581
Masters 1723
Preschool 51
Prof-school 576
Some-college 7290
Some-colleges 1
dtype: int64
```

```
In [44]: # ===== AFTER =====
df['education'] = df['education'].replace(["?", "HS-grad"], "HS-grad")
df['education'] = df['education'].replace(["Some-colleges"], "Some-college")
df.groupby('education').size()

Out[44]: education
10th 933
11th 1175
12th 433
1st-4th 168
5th-6th 333
7th-8th 646
9th 514
Assoc-acdm 1867
Assoc-voc 1382
Bachelors 5354
Doctorate 413
HS-grad 18582
Masters 1723
Prof-school 576
Some-college 7291
dtype: int64
```

Kolom 'marital-status'

Value yang seharusnya: Married-civ-spouse, Divorced, Never-married, Separated, Widowed, Married-spouse-absent, Married-AF-spouse

Value tidak sesuai: "?" Married-civ-spouse dan "Never-mind"

Cara membersihkannya:

- "Married-civ-spouse" akan diubah menjadi "Married-civ-spouse" (asumsi salah ketik)
- "Never-mind" diubah menjadi modus dari data workclass, yaitu "Private"

Di bawah ini adalah before dan after hasil pembersihan datanya.

```
In [45]: # ===== BEFORE =====
df.groupby('marital-status').size()

Out[45]: marital-status
Divorced 4443
Married-AF-spouse 23
Married-civ-spouse 14575
Married-civ-spouse 1
Married-spouse-absent 418
Never-married 18682
Never-mind 1
Separated 1025
Widowed 993
dtype: int64
```

```
In [46]: # ===== AFTER =====
df['marital-status'] = df['marital-status'].replace(["Married-civ-spouse", "Never-mind"], "Married-civ-spouse")
df.groupby('marital-status').size()

Out[46]: marital-status
Divorced 4443
Married-AF-spouse 23
Married-civ-spouse 14577
Married-spouse-absent 418
Never-married 18682
Separated 1025
Widowed 993
dtype: int64
```

Kolom 'occupation'

Value yang seharusnya: Tech-support, Craft-repair, Other-service, Sales, Exec-managerial, Prof-specialty, Handlers-cleaners, Machine-op-inspect, Adm-clerical, Farming-fishing, Transport-moving, Priv-house-serv, Protective-serv, Armed-Forces.

Value tidak sesuai: "?" dan "?"

Cara membersihkannya:

- "?" diubah menjadi modus dari data workclass, yaitu "Craft-repair"

Di bawah ini adalah before dan after hasil pembersihan datanya.

```
In [47]: # BEFORE
df.groupby('occupation').size()

Out[47]: occupation
? 1
Adm-clerical 3769
Armed-Forces 9
Craft-repair 4899
Exec-managerial 4866
Farming-fishing 994
Handlers-cleaners 1370
Machine-op-inspect 2082
Other-service 3295
Priv-house-serv 149
Prof-specialty 4140
Protective-serv 649
Sales 3650
Tech-support 928
Transport-moving 1597
dtype: int64
```

```
In [48]: # AFTER
df['occupation'] = df['occupation'].replace(["?", "?" ], "Craft-repair")
df.groupby('occupation').size()

Out[48]: occupation
Adm-clerical 3769
Armed-Forces 9
Craft-repair 5943
Exec-managerial 4866
Farming-fishing 994
Handlers-cleaners 1370
Machine-op-inspect 2082
Other-service 3295
Priv-house-serv 149
Prof-specialty 4140
Protective-serv 649
Sales 3650
Tech-support 928
Transport-moving 1597
dtype: int64
```

Kolom 'relationship'

Value yang seharusnya: Wife, Own-child, Husband, Not-in-family, Other-relative, Unmarried

Value tidak sesuai: "?" dan "?"

Cara membersihkannya:

- Tidak ada yang perlu dibersihkan

```
In [49]: df.groupby('relationship').size()

Out[49]: relationship
Husband 13193
Not-in-family 8385
Other-relative 981
Own-child 5868
Unmarried 3446
Wife 15187
dtype: int64
```

Kolom 'race'

Value yang seharusnya: White, Asian-Pac-Islander, Amer-Indian-Eskimo, Other, Black

Value tidak sesuai: "Yellow"

Cara membersihkannya:

- "Yellow" akan diubah menjadi "Other" (asumsi maksud dari yellow tidak diketahui dan bukan termasuk ras mudus)

Di bawah ini adalah before dan after hasil pembersihan datanya.

```
In [42]: # BEFORE
df.groupby('race').size()

Out[42]: race
Amer-Indian-Eskimo 311
Asian-Pac-Islander 1039
Black 3124
Other 272
White 27815
Yellow 1
dtype: int64
```

```
In [42]: # AFTER
df['race'] = df['race'].replace(["Yellow", "Other"])
df.groupby('race').size()

Out[42]: race
Amer-Indian-Eskimo 311
Asian-Pac-Islander 1039
Black 3124
Other 272
White 27815
dtype: int64
```

Kolom 'sex'

Value yang seharusnya: Female, Male

Value tidak sesuai: "?" dan "Male/Female"

Cara membersihkannya:

- "Male/Female" diubah menjadi modus dari data workclass, yaitu "Male"
- "?" diubah menjadi modus dari data workclass, yaitu "Male"

Di bawah ini adalah before dan after hasil pembersihan datanya.

```
In [42]: # BEFORE
df.groupby('sex').size()

Out[42]: sex
? 1
Female 10771
Male 21788
Male/Female 1
dtype: int64
```

```
In [42]: # AFTER
df['sex'] = df['sex'].replace(["?", "Male/Female"], "Male")
df.groupby('sex').size()

Out[42]: sex
Female 10771
Male 21790
dtype: int64
```

Kolom 'native-country'

Value yang seharusnya: United States, Canada, 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

Value tidak sesuai: "United-States-America" dan "?"

Cara membersihkannya:

- "United-States-America" akan diubah menjadi "United-States" (asumsi salah ketik)
- "?" diubah menjadi modus dari data workclass, yaitu "United-States"

Di bawah ini adalah before dan after hasil pembersihan datanya.

```
In [42]: # BEFORE
df.groupby('native-country').size()

Out[42]: native-country
? 583
Cambodia 19
Canada 121
China 75
Columbia 59
Cuba 95
Dominican-Republic 70
Ecuador 28
El-Salvador 106
England 90
France 29
Germany 137
Greece 29
Guatemala 64
Haiti 44
Holand-Netherlands 1
Honduras 13
Hong 20
Hungary 13
India 100
Iran 43
Ireland 24
Italy 73
Jamaica 81
Japan 62
Laos 18
Mexico 64
Nicaragua 34
Outlying-US(Guam-USVI-etc) 14
Peru 31
Philippines 198
Poland 60
Portugal 37
Puerto-Rico 114
Scotland 12
South 11
Taiwan 51
Thailand 18
Trinidad&Tobago 19
United-States 29753
United-States-America 1
Vietnam 67
Yugoslavia 16
dtype: int64
```

```
In [42]: # AFTER
df['native-country'] = df['native-country'].replace(["United-States-America", "?"], "United-States")
df.groupby('native-country').size()

Out[42]: native-country
Cambodia 19
Canada 121
China 75
Columbia 59
Cuba 95
Dominican-Republic 70
Ecuador 28
El-Salvador 106
England 90
France 29
Germany 137
Greece 29
Guatemala 64
Haiti 44
Holand-Netherlands 1
Honduras 13
Hong 20
Hungary 13
India 100
Iran 43
Ireland 24
Italy 73
Jamaica 81
Japan 62
Laos 18
Mexico 64
Nicaragua 34
Outlying-US(Guam-USVI-etc) 14
Peru 31
Philippines 198
Poland 60
Portugal 37
Puerto-Rico 114
Scotland 12
South 11
Taiwan 51
Thailand 18
Trinidad&Tobago 19
United-States 29753
Vietnam 67
Yugoslavia 16
dtype: int64
```

Kolom 'income'

Value yang seharusnya: <50K, >50K

Value tidak sesuai: "?" dan "income"

Cara membersihkannya:

```
In [42]: df.groupby('income').size()

Out[42]: income
<50K 24720
>50K 7841
dtype: int64
```

Membersihkan data continuous

Marshi merujuk kepada dokumen `data.names`, kita bisa mencari tahu kolom mana saja yang tipe datanya kontinu (berupa integer). Rencana saya dalam membersihkan kolom dengan tipe data ini adalah dengan menghapus kolom yang valuenya tidak relevan dengan tujuan dari tugas ini.

Sebelum data dibersihkan, cocokan terlebih dahulu tipe datanya dengan tipe data yang seharusnya (mengacu ke dokumen data.names). Ternyata kolom 'capital-gain' masih bertipe data object yang mana seharusnya bertipe int64 (lihat data di bawah). Setelah diubah, operasi pembersihan data kontinu baru bisa dimulai.

```
In [42]: display(df.dtypes)

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

Kolom 'capital-gain'

Ternyata yang membuat capital-gain menjadi tipe data object adalah adanya value "?". Oleh karena itu, kita ganti terlebih dahulu value "?" dengan modus dari data, akan tetapi, setelah dicheck size dari setiap value, "?" hanya terdapat di satu baris data yang mana efeknya tidak signifikan pada keseluruhan data.

```
In [42]: # Menunjukkan jumlah setiap value pada capital-gain
display(df.groupby('capital-gain').size())

capital-gain
0 29847
10520 43
105 25
10566 6
10605 12
dtype: int64

Kemudian, kita akan observasi saya, baris bernomor '0' yang berjumlah 29847 data memenuhi ~92% dari baris tabel. Karena angka '0' terlalu banyak dan dirasa tidak akan menghasilkan sesuatu yang bermanfaat yang berkaitan dengan soal tugas, saya memutuskan untuk menghapus kolom capital-gain. Di bawah ini merupakan bentuk tabel setelah capital-gain dihapus
```

```
In [42]: # Menghapus kolom capital-gain
df.drop(['capital-gain'], axis=1, inplace=True)

Setelah permasalahan teknis/data type selesai, saya akan mengecek kolom-kolom dengan data kontinu lainnya.
```

Kolom 'age'

Berdasarkan data di bawah, karena data terlihat merata, maka 'age' bisa langsung kita pakai.

```
In [43]: plt.figure(figsize=(15,2))
plt.xticks(size=8)
plt.yticks(size=10)
sns.countplot(x='age', data=df)

Out[43]: <AxesSubplot: xlabel='age', ylabel='count'>
```


Kolom 'fnlwgt'

Remove kolom 'fnlwgt' karena datanya sangat tersebar sehingga akan sulit menarik informasi dari data ini.

```
In [43]: display(df.groupby('fnlwgt').size())

# Remove 'fnlwgt'
df.drop(['fnlwgt'], axis=1, inplace=True)

fnlwgt
12285 1
13769 1
14878 1
18827 1
19214 1
1268339 1
1366128 1
1455435 1
1484705 1
22296450 1
Length: 21648, dtype: int64
```

Kolom 'education-num'

Berdasarkan data di bawah, karena data terlihat merata, maka data bisa langsung kita pakai.

```
In [43]: plt.figure(figsize=(15,2))
plt.xticks(size=8)
plt.yticks(size=10)
display(sns.countplot(x='education-num', data=df))

<AxesSubplot: xlabel='education-num', ylabel='count'>
```


Kolom 'capital-loss'

Berdasarkan observasi saya, baris bernomor '0' yang berjumlah 31041 data memenuhi ~95% dari baris tabel. Karena angka '0'-nya terlalu banyak dan dirasa tidak akan menghasilkan sesuatu yang bermanfaat yang berkaitan dengan soal tugas, saya memutuskan untuk menghapus kolom capital-loss.

```
In [43]: # Mendapat 31041 value '0' pada capital-loss
display(df.groupby('capital-loss').size())

capital-loss
0 31041
156 1
233 4
323 3
419 3
3004 2
3683 2
3770 1
3900 2
4356 3
Length: 93, dtype: int64
```

Di bawah ini merupakan bentuk tabel setelah capital-loss dihapus

```
In [43]: # Menghapus kolom capital-loss
df.drop(['capital-loss'], axis=1, inplace=True)
```

Kolom 'hours-per-week'

Meskipun pada kolom hours-per-week terdapat satu value yang kuantitasnya sangat banyak sehingga membuat data lain 'tergeser' statusnya menjadi outlier, saya tidak remove data tersebut. Asumsi saya, meskipun banyak dari mereka bekerja 40 jam seminggu, setiap orang memiliki budaya, peraturan, dan kepatuhan kerja yang berbeda-beda bahkan mungkin saja ada pekerjaan yang tidak memiliki jam kerja tetap sehingga perbedaan jam kerja adalah hal yang wajar.

```
In [43]: plt.figure(figsize=(20,5))
plt.xticks(size=8)
plt.yticks(size=10)
display(sns.countplot(x='hours-per-week', data=df))

<AxesSubplot: xlabel='hours-per-week', ylabel='count'>
```


Hasil pembersihan data continuous

Berikut adalah hasil akhir tabel setelah kolom fnlwgt, capital-gain, capital-loss, kolom dihapus.

```
In [43]: df.head()

Out[43]:
```

	age	workclass	education	education-num	marital-status	occupation	relationship	race	sex	hours-per-week	native-country	income
0	39	State-gov	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	40	United-States	<=50K
1	50	Self-emp-not-inc	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	13	United-States	<=50K
2	38	Private	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	40	United-States	<=50K
3	53	Private	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	40	United-States	<=50K
4	28	Private	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	40	Cuba	<=50K

2. Lakukan penambangan data menggunakan SVM dan Neural Network. Bagaimana hasilnya. Mana kah yang memberikan performance yang lebih baik, apa sebabnya.

Encode

Pada bagian ini, data kategorikal akan diencode. Karena kolom 'education' sudah direpresentasikan oleh 'education-num', kita bisa hapus kolom 'education'.

```
In [43]: # Drop kolom education
df.drop(['education'], axis=1, inplace=True)
df.head()
```

```
Out[43]:
```

	age	workclass	education-num	marital-status	occupation	relationship	race	sex	hours-per-week	native-country	income
0	39	State-gov	13	Never-married	Adm-clerical	Not-in-family	White	Male	40	United-States	<=50K
1	50	Self-emp-not-inc	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	13	United-States	<=50K
2	38	Private	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	40	United-States	<=50K
3	53	Private	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	40	United-States	<=50K
4	28	Private	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	40	Cuba	<=50K

Encode kolom income

Kita akan encode income dengan label encoding yaitu 0 untuk <=50K dan 1 untuk >50K

```
In [43]: df_encoded = df.copy()
df_encoded['income'] = df_encoded['income'].replace(['<=50K', '>50K'], 0, 1)

Encode kategorikal tidak ordinal
```

Encode data kategorikal yang tidak ordinal (tidak ada urutannya) menggunakan one-hot-encoding

```
In [43]: df_encoded = pd.get_dummies(df_encoded)
display(df_encoded.head())
```

	age	education-num	hours-per-week	income	workclass_Federal-gov	workclass_Local-gov	workclass_Never-worked	workclass_Private	workclass_Self-emp-inc	workclass_Self-emp-not-inc
0	39	13	40	0	0	0	0	0	0	0
1	50	13	13	0	0	0	0	0	0	1
2	38	9	40	0	0	0	0	1	0	0
3	53	7	40	0	0	0	0	1	0	0
4	28	13	40	0	0	0	0	1	0	0

Korelasinya dengan 'income'

Pada bagian ini, kita akan mencari tahu berapa nilai korelasi antara kolom-column di atas dengan income. Menurut

https://aisi.law.com/amesesia-analisis/ai/text/2018/05/00/correlation_coefficients_appropriate_use_and_50asp, data yang bisa diabaikan adalah ketika korelasinya dengan target bernilai < 0.1

```
In [44]: all_corr = df_encoded.corr()[['income']]
all_corr

Out[44]:
```

	age	education-num	hours-per-week	income	workclass_Federal-gov	workclass_Local-gov	workclass_Never-worked	workclass_Private	workclass_Self-emp-inc	workclass_Self-emp-not-inc
education-num	0.234037	</								


```
accuracy: 0.8495254048017867
precision score: 0.7617713804484304
recall: 0.5327321610580949
f1: 0.6269896193771626
```

Neural Network

```
In [474]: # Neural Network
mn = MLPClassifier(alpha=1e-5, hidden_layer_sizes=(12, 6, 3), random_state=1, max_iter=500)
mn.fit(x_test, y_test)
y_pred = mn.predict(x_test)

print("accuracy: ", accuracy_score(y_test, y_pred))
print("precision score: ", precision_score(y_test, y_pred))
print("recall: ", recall_score(y_test, y_pred))
print("f1: ", f1_score(y_test, y_pred))

accuracy: 0.832166387490801
precision score: 0.6085159214592275
recall: 0.6032928263426107
f1: 0.6306084818684695
```

Penutup

Yang memberikan performa lebih baik adalah Neural Network (NN). Neural network membutuhkan input data yang banyak jika dibandingkan dengan SVM. Semakin banyak data yang dimasukkan ke dalam network, maka akan semakin men-generate hasil yang lebih baik dan akurat sehingga membuat prediksi memiliki lebih sedikit kesalahan. Data yang kita miliki cukup banyak, oleh karena itu NN lebih cocok digunakan.

In []: