# 1) Pengumpulan Data

# 2) Menelaah Data

```python
import pandas as pd
import numpy as np
import re
import itertools
```

Load Data  masukkan dataset yang dibutuhkan

```python
#dir merupakan variabel yang digunakan untuk menyimpan data yang akan
diproses
dir = 'hungarian.data'

#membuka file dengan path atau nama file yang disimpan dalam variabel
dir.
with open(dir, encoding='Latin1') as file:
  #membaca setiap baris, menghapus whitespace di awal dan akhir setiap
baris, dan menyimpannya dalam list lines.
  lines = [line.strip() for line in file]
#Kemudian, 10 baris pertama dari list tersebut dicetak.
lines[0:10]
```

```
['1254 0 40 1 1 0 0',
 '-9 2 140 0 289 -9 -9 -9',
 '0 -9 -9 0 12 16 84 0',
 '0 0 0 0 150 18 -9 7',
 '172 86 200 110 140 86 0 0',
 '0 -9 26 20 -9 -9 -9 -9',
 '-9 -9 -9 -9 -9 -9 -9 12',
 '20 84 0 -9 -9 -9 -9 -9',
 '-9 -9 -9 -9 -9 1 1 1',
 '1 1 -9. -9. name']
```

```python
import itertools

#mengambil setiap kelompok 10 baris dari list lines, menggabungkannya
menjadi satu string, membaginya menjadi list hingga panjang 76
data = itertools.takewhile(
    lambda x: len(x) ==76,
    (' '.join(lines[i:(i+10)]).split() for i in range(0,
len(lines),10))
)
#Hasilnya dimasukkan ke dalam DataFrame.
df = pd.DataFrame.from_records(data)
```

```python
#lima baris terakhir dari DataFrame dicetak.
df.tail()
```

```
        0   1   2   3  4  5  6   7  8    9  ...  66  67 68 69 70 71 72
73  \
289  1053   0  48   0  0  0  0  -9  2   -9  ...  -9  -9  1  1  1  1  1
-9.
290  1054   0  36   1  1  0  0  -9  2  120  ...  -9  -9  1  1  1  1  1
-9.
291  5001   0  48   1  0  0  0  -9  3  110  ...  -9  -9  1  1  1  1  1
-9.
292  5000   0  47   0  0  0  0  -9  2  140  ...  -9  -9  1  1  1  1  1
-9.
293  5002   0  53   1  1  1  1  -9  4  130  ...   1   1  1  1  1  1  1
-9.

      74     75
289  -9.   name
290  -9.   name
291  -9.   name
292  -9.   name
293  -9.   name

[5 rows x 76 columns]
```

```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 294 entries, 0 to 293
Data columns (total 76 columns):
 #    Column   Non-Null Count   Dtype
---   ------   --------------   -----
 0    0        294 non-null     object
 1    1        294 non-null     object
 2    2        294 non-null     object
 3    3        294 non-null     object
 4    4        294 non-null     object
 5    5        294 non-null     object
 6    6        294 non-null     object
 7    7        294 non-null     object
 8    8        294 non-null     object
 9    9        294 non-null     object
 10   10       294 non-null     object
 11   11       294 non-null     object
 12   12       294 non-null     object
 13   13       294 non-null     object
 14   14       294 non-null     object
 15   15       294 non-null     object
 16   16       294 non-null     object
 17   17       294 non-null     object
```

```
18  18        294 non-null     object
19  19        294 non-null     object
20  20        294 non-null     object
21  21        294 non-null     object
22  22        294 non-null     object
23  23        294 non-null     object
24  24        294 non-null     object
25  25        294 non-null     object
26  26        294 non-null     object
27  27        294 non-null     object
28  28        294 non-null     object
29  29        294 non-null     object
30  30        294 non-null     object
31  31        294 non-null     object
32  32        294 non-null     object
33  33        294 non-null     object
34  34        294 non-null     object
35  35        294 non-null     object
36  36        294 non-null     object
37  37        294 non-null     object
38  38        294 non-null     object
39  39        294 non-null     object
40  40        294 non-null     object
41  41        294 non-null     object
42  42        294 non-null     object
43  43        294 non-null     object
44  44        294 non-null     object
45  45        294 non-null     object
46  46        294 non-null     object
47  47        294 non-null     object
48  48        294 non-null     object
49  49        294 non-null     object
50  50        294 non-null     object
51  51        294 non-null     object
52  52        294 non-null     object
53  53        294 non-null     object
54  54        294 non-null     object
55  55        294 non-null     object
56  56        294 non-null     object
57  57        294 non-null     object
58  58        294 non-null     object
59  59        294 non-null     object
60  60        294 non-null     object
61  61        294 non-null     object
62  62        294 non-null     object
63  63        294 non-null     object
64  64        294 non-null     object
65  65        294 non-null     object
66  66        294 non-null     object
```

```
67   67            294 non-null        object
68   68            294 non-null        object
69   69            294 non-null        object
70   70            294 non-null        object
71   71            294 non-null        object
72   72            294 non-null        object
73   73            294 non-null        object
74   74            294 non-null        object
75   75            294 non-null        object
dtypes: object(76)
memory usage: 174.7+ KB
```

```
#menghapus colomn pertama dan terakhir
df = df.iloc[:,:-1]
df = df.drop(df.columns[0], axis=1)
```

mengubah tipe data file daaset menjadi tipe data float sesuai dengan nilai null yaitu -9.0

```
df = df.astype(float)

df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 294 entries, 0 to 293
Data columns (total 74 columns):
 #    Column   Non-Null Count   Dtype
---   ------   --------------   -----
 0    1            294 non-null      float64
 1    2            294 non-null      float64
 2    3            294 non-null      float64
 3    4            294 non-null      float64
 4    5            294 non-null      float64
 5    6            294 non-null      float64
 6    7            294 non-null      float64
 7    8            294 non-null      float64
 8    9            294 non-null      float64
 9    10           294 non-null      float64
 10   11           294 non-null      float64
 11   12           294 non-null      float64
 12   13           294 non-null      float64
 13   14           294 non-null      float64
 14   15           294 non-null      float64
 15   16           294 non-null      float64
 16   17           294 non-null      float64
 17   18           294 non-null      float64
 18   19           294 non-null      float64
 19   20           294 non-null      float64
 20   21           294 non-null      float64
 21   22           294 non-null      float64
 22   23           294 non-null      float64
```

```
23   24        294 non-null      float64
24   25        294 non-null      float64
25   26        294 non-null      float64
26   27        294 non-null      float64
27   28        294 non-null      float64
28   29        294 non-null      float64
29   30        294 non-null      float64
30   31        294 non-null      float64
31   32        294 non-null      float64
32   33        294 non-null      float64
33   34        294 non-null      float64
34   35        294 non-null      float64
35   36        294 non-null      float64
36   37        294 non-null      float64
37   38        294 non-null      float64
38   39        294 non-null      float64
39   40        294 non-null      float64
40   41        294 non-null      float64
41   42        294 non-null      float64
42   43        294 non-null      float64
43   44        294 non-null      float64
44   45        294 non-null      float64
45   46        294 non-null      float64
46   47        294 non-null      float64
47   48        294 non-null      float64
48   49        294 non-null      float64
49   50        294 non-null      float64
50   51        294 non-null      float64
51   52        294 non-null      float64
52   53        294 non-null      float64
53   54        294 non-null      float64
54   55        294 non-null      float64
55   56        294 non-null      float64
56   57        294 non-null      float64
57   58        294 non-null      float64
58   59        294 non-null      float64
59   60        294 non-null      float64
60   61        294 non-null      float64
61   62        294 non-null      float64
62   63        294 non-null      float64
63   64        294 non-null      float64
64   65        294 non-null      float64
65   66        294 non-null      float64
66   67        294 non-null      float64
67   68        294 non-null      float64
68   69        294 non-null      float64
69   70        294 non-null      float64
70   71        294 non-null      float64
71   72        294 non-null      float64
```

```
 72   73        294 non-null      float64
 73   74        294 non-null      float64
dtypes: float64(74)
memory usage: 170.1 KB
```

```python
##menganti nilai -9.0 menjadi NaN
df.replace(-9.0, np.nan, inplace=True)

##informasi jumlah NaN dalam setiap DataFrame
df.isnull().sum()
```

```
1        0
2        0
3        0
4        0
5        0
     ...
70       0
71       0
72       0
73     266
74     294
Length: 74, dtype: int64
```

```python
df.head()
```

```
    1     2    3    4    5    6   7    8      9    10  ...  65   66
67   68  \
0  0.0  40.0  1.0  1.0  0.0  0.0 NaN  2.0  140.0  0.0  ... NaN  NaN
NaN  1.0
1  0.0  49.0  0.0  1.0  0.0  0.0 NaN  3.0  160.0  1.0  ... NaN  NaN
NaN  1.0
2  0.0  37.0  1.0  1.0  0.0  0.0 NaN  2.0  130.0  0.0  ... NaN  NaN
NaN  1.0
3  0.0  48.0  0.0  1.0  1.0  1.0 NaN  4.0  138.0  0.0  ... NaN  2.0
NaN  1.0
4  0.0  54.0  1.0  1.0  0.0  1.0 NaN  3.0  150.0  0.0  ... NaN  1.0
NaN  1.0

    69    70    71    72   73   74
0  1.0   1.0   1.0   1.0 NaN NaN
1  1.0   1.0   1.0   1.0 NaN NaN
2  1.0   1.0   1.0   1.0 NaN NaN
3  1.0   1.0   1.0   1.0 NaN NaN
4  1.0   1.0   1.0   1.0 NaN NaN

[5 rows x 74 columns]
```

```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 294 entries, 0 to 293
Data columns (total 74 columns):
 #    Column   Non-Null Count   Dtype
---   ------   --------------   -----
 0    1        294 non-null     float64
 1    2        294 non-null     float64
 2    3        294 non-null     float64
 3    4        294 non-null     float64
 4    5        294 non-null     float64
 5    6        294 non-null     float64
 6    7        0 non-null       float64
 7    8        294 non-null     float64
 8    9        293 non-null     float64
 9    10       293 non-null     float64
 10   11       271 non-null     float64
 11   12       12 non-null      float64
 12   13       1 non-null       float64
 13   14       0 non-null       float64
 14   15       286 non-null     float64
 15   16       21 non-null      float64
 16   17       1 non-null       float64
 17   18       293 non-null     float64
 18   19       294 non-null     float64
 19   20       294 non-null     float64
 20   21       294 non-null     float64
 21   22       293 non-null     float64
 22   23       292 non-null     float64
 23   24       293 non-null     float64
 24   25       293 non-null     float64
 25   26       293 non-null     float64
 26   27       285 non-null     float64
 27   28       292 non-null     float64
 28   29       104 non-null     float64
 29   30       292 non-null     float64
 30   31       293 non-null     float64
 31   32       293 non-null     float64
 32   33       293 non-null     float64
 33   34       293 non-null     float64
 34   35       293 non-null     float64
 35   36       293 non-null     float64
 36   37       293 non-null     float64
 37   38       292 non-null     float64
 38   39       294 non-null     float64
 39   40       104 non-null     float64
 40   41       293 non-null     float64
 41   42       294 non-null     float64
 42   43       4 non-null       float64
 43   44       0 non-null       float64
 44   45       0 non-null       float64
```

```
45   46        0 non-null      float64
46   47        3 non-null      float64
47   48        0 non-null      float64
48   49        2 non-null      float64
49   50       28 non-null      float64
50   51       27 non-null      float64
51   52       17 non-null      float64
52   53        0 non-null      float64
53   54      294 non-null      float64
54   55      294 non-null      float64
55   56      294 non-null      float64
56   57      294 non-null      float64
57   58       19 non-null      float64
58   59       58 non-null      float64
59   60       48 non-null      float64
60   61       18 non-null      float64
61   62       59 non-null      float64
62   63        9 non-null      float64
63   64       23 non-null      float64
64   65        5 non-null      float64
65   66       50 non-null      float64
66   67       25 non-null      float64
67   68      294 non-null      float64
68   69      294 non-null      float64
69   70      294 non-null      float64
70   71      294 non-null      float64
71   72      294 non-null      float64
72   73       28 non-null      float64
73   74        0 non-null      float64
dtypes: float64(74)
memory usage: 170.1 KB
```

# 4) Menentukan Object Data

memilih 14 fitur yang akan digunakan sesuai deskripsi dataset

```
#mementukan object data yang sudah diketahui pada dataset
df_selected = df.iloc[0:, [1, 2, 7, 8, 10, 14, 17, 30, 36, 38, 39, 42,
49, 56]]

df_selected.head()

        2     3     8      9      11    15    18      31     37    39    40    43
50    57
0   40.0   1.0   2.0  140.0  289.0   0.0   0.0  172.0   0.0   0.0  NaN  NaN
NaN   0.0
1   49.0   0.0   3.0  160.0  180.0   0.0   0.0  156.0   0.0   1.0  2.0  NaN
NaN   1.0
```

```
2   37.0   1.0   2.0   130.0   283.0   0.0   1.0    98.0   0.0   0.0   NaN NaN
NaN   0.0
3   48.0   0.0   4.0   138.0   214.0   0.0   0.0   108.0   1.0   1.5   2.0 NaN
NaN   3.0
4   54.0   1.0   3.0   150.0     NaN   0.0   0.0   122.0   0.0   0.0   NaN NaN
NaN   0.0
```

```
df_selected.tail()
```

```
          2     3     8       9      11     15     18      31     37     39     40    43
50   \
289   48.0   0.0   2.0     NaN   308.0   0.0   1.0     NaN    NaN    2.0    1.0   NaN
NaN
290   36.0   1.0   2.0   120.0   166.0   0.0   0.0   180.0    0.0    0.0    NaN   NaN
NaN
291   48.0   1.0   3.0   110.0   211.0   0.0   0.0   138.0    0.0    0.0    NaN   NaN
6.0
292   47.0   0.0   2.0   140.0   257.0   0.0   0.0   135.0    0.0    1.0    1.0   NaN
NaN
293   53.0   1.0   4.0   130.0   182.0   0.0   0.0   148.0    0.0    0.0    NaN   NaN
NaN

        57
289   0.0
290   0.0
291   0.0
292   0.0
293   0.0
```

```
df_selected.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 294 entries, 0 to 293
Data columns (total 14 columns):
 #    Column   Non-Null Count   Dtype
---   ------   --------------   -----
 0    2        294 non-null     float64
 1    3        294 non-null     float64
 2    8        294 non-null     float64
 3    9        293 non-null     float64
 4    11       271 non-null     float64
 5    15       286 non-null     float64
 6    18       293 non-null     float64
 7    31       293 non-null     float64
 8    37       293 non-null     float64
 9    39       294 non-null     float64
 10   40       104 non-null     float64
 11   43       4 non-null       float64
 12   50       28 non-null      float64
 13   57       294 non-null     float64
```

```
dtypes: float64(14)
memory usage: 32.3 KB
```

mengganti nama kolom sesuai dengan 14 nama kolom yang ada pada deskripsi dataset

```python
import re

columns_mapping = {
    2: 'age',
    3: 'sex',
    8: 'cp',
    9: 'trestbps',
    11: 'chol',
    15: 'fbs',
    18: 'restecg',
    31: 'thalach',
    37: 'exang',
    39: 'oldpeak',
    40: 'slope',
    43: 'ca',
    50: 'thal',
    57: 'target'
}

df_selected.rename(columns = columns_mapping, inplace=True)

df_selected.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 294 entries, 0 to 293
Data columns (total 14 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   age       294 non-null    float64
 1   sex       294 non-null    float64
 2   cp        294 non-null    float64
 3   trestbps  293 non-null    float64
 4   chol      271 non-null    float64
 5   fbs       286 non-null    float64
 6   restecg   293 non-null    float64
 7   thalach   293 non-null    float64
 8   exang     293 non-null    float64
 9   oldpeak   294 non-null    float64
 10  slope     104 non-null    float64
 11  ca        4 non-null      float64
 12  thal      28 non-null     float64
 13  target    294 non-null    float64
dtypes: float64(14)
memory usage: 32.3 KB
```

```
#melihat nilai mana yang paling sering muncul dan paling jarang
df_selected.value_counts()
```

```
age    sex   cp    trestbps   chol    fbs   restecg   thalach   exang   oldpeak
slope   ca    thal   target
47.0   1.0   4.0   150.0      226.0   0.0   0.0       98.0      1.0     1.5
2.0     0.0   7.0    1.0        1
dtype: int64
```

# 5) Membersihkan data

```
#melihat banyak nilai null pada setiap object
df_selected.isnull().sum()
```

```
age            0
sex            0
cp             0
trestbps       1
chol          23
fbs            8
restecg        1
thalach        1
exang          1
oldpeak        0
slope        190
ca           290
thal         266
target         0
dtype: int64
```

```
#menghapus object yang memiliki nilai null yang sangat banyak
columns_to_drop = ['ca','slope','thal']
df_selected = df_selected.drop(columns_to_drop, axis =1)

df_selected.isnull().sum()
```

```
age            0
sex            0
cp             0
trestbps       1
chol          23
fbs            8
restecg        1
thalach        1
exang          1
oldpeak        0
target         0
dtype: int64
```

```python
#menghapus nilai NaN sehingga nanti bisa diisi nilai baru
meanTBPS = df_selected['trestbps'].dropna()
meanChol = df_selected['chol'].dropna()
meanFbs = df_selected['fbs'].dropna()
meanRestCG = df_selected['restecg'].dropna()
meanThalach = df_selected['thalach'].dropna()
meanExang = df_selected['exang'].dropna()
#konversi tipe data object menjadi float
meanTBPS = meanTBPS.astype(float)
meanChol = meanChol.astype(float)
meanFbs = meanFbs.astype(float)
meanRestCG = meanRestCG.astype(float)
meanThalach = meanThalach.astype(float)
meanExang = meanExang.astype(float)
#membuat nilai baru dari perhitungan rata2 dari nilai yang kemudian
dibulatkan
meanTBPS = round(meanTBPS.mean())
meanChol = round(meanChol.mean())
meanFbs = round(meanFbs.mean())
meanRestCG = round(meanRestCG.mean())
meanThalach = round(meanThalach.mean())
meanExang = round(meanExang.mean())

#mengisi nilai NaN tadi dengan nilai rata2 yang sebelumnya sudah
dihitung
fill_values = {'trestbps' : meanTBPS,'chol' : meanChol,
               'fbs' : meanFbs,'thalach' : meanThalach,
               'exang' : meanExang,'restecg' : meanRestCG}
dfClean = df_selected.fillna(value=fill_values)

dfClean.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 294 entries, 0 to 293
Data columns (total 11 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   age       294 non-null    float64
 1   sex       294 non-null    float64
 2   cp        294 non-null    float64
 3   trestbps  294 non-null    float64
 4   chol      294 non-null    float64
 5   fbs       294 non-null    float64
 6   restecg   294 non-null    float64
 7   thalach   294 non-null    float64
 8   exang     294 non-null    float64
 9   oldpeak   294 non-null    float64
 10  target    294 non-null    float64
dtypes: float64(11)
memory usage: 25.4 KB
```

```
dfClean.isnull().sum()

age          0
sex          0
cp           0
trestbps     0
chol         0
fbs          0
restecg      0
thalach      0
exang        0
oldpeak      0
target       0
dtype: int64
```

```
#mengidentifikasi baris-baris yang merupakan duplikat
duplicate_rows = dfClean.duplicated()
dfClean[duplicate_rows]
```

```
        age  sex   cp  trestbps   chol  fbs  restecg  thalach  exang
oldpeak  \
163    49.0  0.0  2.0     110.0  251.0  0.0      0.0    160.0    0.0
0.0

       target
163       0.0
```

```
#melihat baris duplikat
print("All Duplicate_rows")
dfClean[dfClean.duplicated(keep=False)]
```

```
All Duplicate_rows

        age  sex   cp  trestbps   chol  fbs  restecg  thalach  exang
oldpeak  \
90     49.0  0.0  2.0     110.0  251.0  0.0      0.0    160.0    0.0
0.0
163    49.0  0.0  2.0     110.0  251.0  0.0      0.0    160.0    0.0
0.0

       target
90        0.0
163       0.0
```

```
#menghapus semua baris duplikat
dfClean = dfClean.drop_duplicates()
#mengecek kembali apakah ada baris yang masih duplikat
print("All Duplicate Rows : ")
dfClean[dfClean.duplicated(keep=False)]
```

```
All Duplicate Rows :
```

```
Empty DataFrame
Columns: [age, sex, cp, trestbps, chol, fbs, restecg, thalach, exang,
oldpeak, target]
Index: []
```

```
dfClean.head()
```

```
      age   sex    cp   trestbps    chol   fbs   restecg   thalach   exang
oldpeak   \
0    40.0   1.0   2.0      140.0   289.0   0.0       0.0     172.0     0.0
0.0
1    49.0   0.0   3.0      160.0   180.0   0.0       0.0     156.0     0.0
1.0
2    37.0   1.0   2.0      130.0   283.0   0.0       1.0      98.0     0.0
0.0
3    48.0   0.0   4.0      138.0   214.0   0.0       0.0     108.0     1.0
1.5
4    54.0   1.0   3.0      150.0   251.0   0.0       0.0     122.0     0.0
0.0

    target
0      0.0
1      1.0
2      0.0
3      3.0
4      0.0
```

```python
#melihat jumlah kemunculan setiap nilai unik dalam kolom target
dfClean['target'].value_counts()
```

```
0.0     187
1.0      37
3.0      28
2.0      26
4.0      15
Name: target, dtype: int64
```

```python
import seaborn as sns
import matplotlib.pyplot as plt

#meghitung matriks korelasi
dfClean.corr()
```

```
             age        sex         cp    trestbps       chol
fbs    \
age     1.000000   0.014516   0.146616   0.246571   0.087101   0.181130

sex     0.014516   1.000000   0.245769   0.082064   0.027695   0.044372

cp      0.146616   0.245769   1.000000   0.081293   0.134697   0.031930
```

```
trestbps   0.246571   0.082064   0.081293   1.000000   0.080818   0.096222

chol       0.087101   0.027695   0.134697   0.080818   1.000000   0.107686

fbs        0.181130   0.044372   0.031930   0.096222   0.107686   1.000000

restecg    0.050672  -0.108656  -0.016372   0.011256   0.048081   0.047988

thalach   -0.460514  -0.106959  -0.367819  -0.181824  -0.122038  -0.069722

exang      0.239223   0.154925   0.494674   0.211507   0.161055   0.115503

oldpeak    0.178172   0.115959   0.351735   0.204000   0.106743   0.063179

target     0.210429   0.220732   0.427536   0.214898   0.256027   0.154319


            restecg    thalach      exang    oldpeak     target
age         0.050672  -0.460514   0.239223   0.178172   0.210429
sex        -0.108656  -0.106959   0.154925   0.115959   0.220732
cp         -0.016372  -0.367819   0.494674   0.351735   0.427536
trestbps    0.011256  -0.181824   0.211507   0.204000   0.214898
chol        0.048081  -0.122038   0.161055   0.106743   0.256027
fbs         0.047988  -0.069722   0.115503   0.063179   0.154319
restecg     1.000000   0.006084   0.041290   0.042193   0.042643
thalach     0.006084   1.000000  -0.400508  -0.300458  -0.367525
exang       0.041290  -0.400508   1.000000   0.624965   0.571710
oldpeak     0.042193  -0.300458   0.624965   1.000000   0.580732
target      0.042643  -0.367525   0.571710   0.580732   1.000000
```

```python
#visualiasi matriks korelasi
cor_mat = dfClean.corr()
# Membuat objek gambar (fig) dan sumbu (ax) menggunakan subplots dari
matplotlib
fig,ax = plt.subplots(figsize=(15,10))
#membuat heatmap dari matriks korelasi. Parameter annot=True
menambahkan label nilai korelasi pada sel heatmap. linewidth=0.5
menentukan lebar garis pembatas antar sel. fmt=".3f" menentukan format
nilai desimal untuk label (tiga desimal).
sns.heatmap(cor_mat, annot=True, linewidth=0.5, fmt=" .3f")
```

```
<Axes: >
```