

# Tugas Praktikum

Wisconsin Breast Cancer

## Deskripsi Tugas

Pada tugas pratikum ini Anda akan menggunakan data "Wisconsin Breast Cancer". Data tersebut terdiri dari 569 data yang digunakan untuk mendiagnosis jenis kanker Malignant (M) dan Benign (B). Tugas Anda adalah,

1. Pisahkan antara variabel yang dapat digunakan dan variabel yang tidak dapat digunakan.

```
import pandas as pd
```

```
dpath = 'Folder Data/wbc.csv'
```

```
df = pd.read_csv(dpath)
```

```
df.head()
```

|             | id       | diagnosis | radius_mean | texture_mean | perimeter_mean |
|-------------|----------|-----------|-------------|--------------|----------------|
| area_mean \ |          |           |             |              |                |
| 0           | 842302   | M         | 17.99       | 10.38        | 122.80         |
| 1001.0      |          |           |             |              |                |
| 1           | 842517   | M         | 20.57       | 17.77        | 132.90         |
| 1326.0      |          |           |             |              |                |
| 2           | 84300903 | M         | 19.69       | 21.25        | 130.00         |
| 1203.0      |          |           |             |              |                |
| 3           | 84348301 | M         | 11.42       | 20.38        | 77.58          |
| 386.1       |          |           |             |              |                |
| 4           | 84358402 | M         | 20.29       | 14.34        | 135.10         |
| 1297.0      |          |           |             |              |                |

|               | smoothness_mean | compactness_mean | concavity_mean | concave |
|---------------|-----------------|------------------|----------------|---------|
| points_mean \ |                 |                  |                |         |
| 0             | 0.11840         | 0.27760          | 0.3001         |         |
| 0.14710       |                 |                  |                |         |
| 1             | 0.08474         | 0.07864          | 0.0869         |         |
| 0.07017       |                 |                  |                |         |
| 2             | 0.10960         | 0.15990          | 0.1974         |         |
| 0.12790       |                 |                  |                |         |
| 3             | 0.14250         | 0.28390          | 0.2414         |         |
| 0.10520       |                 |                  |                |         |
| 4             | 0.10030         | 0.13280          | 0.1980         |         |
| 0.10430       |                 |                  |                |         |

| ...                | texture_worst | perimeter_worst | area_worst |
|--------------------|---------------|-----------------|------------|
| smoothness_worst \ |               |                 |            |
| 0 ...              | 17.33         | 184.60          | 2019.0     |
|                    |               |                 | 0.1622     |

|   |     |       |        |        |        |
|---|-----|-------|--------|--------|--------|
| 1 | ... | 23.41 | 158.80 | 1956.0 | 0.1238 |
| 2 | ... | 25.53 | 152.50 | 1709.0 | 0.1444 |
| 3 | ... | 26.50 | 98.87  | 567.7  | 0.2098 |
| 4 | ... | 16.67 | 152.20 | 1575.0 | 0.1374 |

|                  |                   |                 |                      |
|------------------|-------------------|-----------------|----------------------|
|                  | compactness_worst | concavity_worst | concave points_worst |
| symmetry_worst \ |                   |                 |                      |
| 0                | 0.6656            | 0.7119          | 0.2654               |
| 0.4601           |                   |                 |                      |
| 1                | 0.1866            | 0.2416          | 0.1860               |
| 0.2750           |                   |                 |                      |
| 2                | 0.4245            | 0.4504          | 0.2430               |
| 0.3613           |                   |                 |                      |
| 3                | 0.8663            | 0.6869          | 0.2575               |
| 0.6638           |                   |                 |                      |
| 4                | 0.2050            | 0.4000          | 0.1625               |
| 0.2364           |                   |                 |                      |

|   |                         |             |
|---|-------------------------|-------------|
|   | fractal_dimension_worst | Unnamed: 32 |
| 0 | 0.11890                 | NaN         |
| 1 | 0.08902                 | NaN         |
| 2 | 0.08758                 | NaN         |
| 3 | 0.17300                 | NaN         |
| 4 | 0.07678                 | NaN         |

[5 rows x 33 columns]

df.info()

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 569 entries, 0 to 568

Data columns (total 33 columns):

| #   | Column                 | Non-Null Count | Dtype   |
|-----|------------------------|----------------|---------|
| --- | -----                  | -----          | -----   |
| 0   | id                     | 569 non-null   | int64   |
| 1   | diagnosis              | 569 non-null   | object  |
| 2   | radius_mean            | 569 non-null   | float64 |
| 3   | texture_mean           | 569 non-null   | float64 |
| 4   | perimeter_mean         | 569 non-null   | float64 |
| 5   | area_mean              | 569 non-null   | float64 |
| 6   | smoothness_mean        | 569 non-null   | float64 |
| 7   | compactness_mean       | 569 non-null   | float64 |
| 8   | concavity_mean         | 569 non-null   | float64 |
| 9   | concave points_mean    | 569 non-null   | float64 |
| 10  | symmetry_mean          | 569 non-null   | float64 |
| 11  | fractal_dimension_mean | 569 non-null   | float64 |

|    |                         |     |          |         |
|----|-------------------------|-----|----------|---------|
| 12 | radius_se               | 569 | non-null | float64 |
| 13 | texture_se              | 569 | non-null | float64 |
| 14 | perimeter_se            | 569 | non-null | float64 |
| 15 | area_se                 | 569 | non-null | float64 |
| 16 | smoothness_se           | 569 | non-null | float64 |
| 17 | compactness_se          | 569 | non-null | float64 |
| 18 | concavity_se            | 569 | non-null | float64 |
| 19 | concave points_se       | 569 | non-null | float64 |
| 20 | symmetry_se             | 569 | non-null | float64 |
| 21 | fractal_dimension_se    | 569 | non-null | float64 |
| 22 | radius_worst            | 569 | non-null | float64 |
| 23 | texture_worst           | 569 | non-null | float64 |
| 24 | perimeter_worst         | 569 | non-null | float64 |
| 25 | area_worst              | 569 | non-null | float64 |
| 26 | smoothness_worst        | 569 | non-null | float64 |
| 27 | compactness_worst       | 569 | non-null | float64 |
| 28 | concavity_worst         | 569 | non-null | float64 |
| 29 | concave points_worst    | 569 | non-null | float64 |
| 30 | symmetry_worst          | 569 | non-null | float64 |
| 31 | fractal_dimension_worst | 569 | non-null | float64 |
| 32 | Unnamed: 32             | 0   | non-null | float64 |

dtypes: float64(31), int64(1), object(1)

memory usage: 146.8+ KB

df.isnull().sum()

|                        |   |
|------------------------|---|
| id                     | 0 |
| diagnosis              | 0 |
| radius_mean            | 0 |
| texture_mean           | 0 |
| perimeter_mean         | 0 |
| area_mean              | 0 |
| smoothness_mean        | 0 |
| compactness_mean       | 0 |
| concavity_mean         | 0 |
| concave points_mean    | 0 |
| symmetry_mean          | 0 |
| fractal_dimension_mean | 0 |
| radius_se              | 0 |
| texture_se             | 0 |
| perimeter_se           | 0 |
| area_se                | 0 |
| smoothness_se          | 0 |
| compactness_se         | 0 |
| concavity_se           | 0 |
| concave points_se      | 0 |
| symmetry_se            | 0 |
| fractal_dimension_se   | 0 |
| radius_worst           | 0 |
| texture_worst          | 0 |

```

perimeter_worst      0
area_worst           0
smoothness_worst     0
compactness_worst    0
concavity_worst      0
concave points_worst 0
symmetry_worst       0
fractal_dimension_worst 0
Unnamed: 32          569
dtype: int64

```

```

df = df.drop(columns=['id', 'Unnamed: 32'])
df.head()

```

|   | diagnosis | radius_mean | texture_mean | perimeter_mean | area_mean | \ |
|---|-----------|-------------|--------------|----------------|-----------|---|
| 0 | M         | 17.99       | 10.38        | 122.80         | 1001.0    |   |
| 1 | M         | 20.57       | 17.77        | 132.90         | 1326.0    |   |
| 2 | M         | 19.69       | 21.25        | 130.00         | 1203.0    |   |
| 3 | M         | 11.42       | 20.38        | 77.58          | 386.1     |   |
| 4 | M         | 20.29       | 14.34        | 135.10         | 1297.0    |   |

|   | smoothness_mean | compactness_mean | concavity_mean | concave points_mean | \ |
|---|-----------------|------------------|----------------|---------------------|---|
| 0 | 0.11840         | 0.27760          | 0.3001         | 0.14710             |   |
| 1 | 0.08474         | 0.07864          | 0.0869         | 0.07017             |   |
| 2 | 0.10960         | 0.15990          | 0.1974         | 0.12790             |   |
| 3 | 0.14250         | 0.28390          | 0.2414         | 0.10520             |   |
| 4 | 0.10030         | 0.13280          | 0.1980         | 0.10430             |   |

|   | symmetry_mean | ... | radius_worst | texture_worst | perimeter_worst | \ |
|---|---------------|-----|--------------|---------------|-----------------|---|
| 0 | 0.2419        | ... | 25.38        | 17.33         | 184.60          |   |
| 1 | 0.1812        | ... | 24.99        | 23.41         | 158.80          |   |
| 2 | 0.2069        | ... | 23.57        | 25.53         | 152.50          |   |
| 3 | 0.2597        | ... | 14.91        | 26.50         | 98.87           |   |
| 4 | 0.1809        | ... | 22.54        | 16.67         | 152.20          |   |

|   | area_worst | smoothness_worst | compactness_worst | concavity_worst | \ |
|---|------------|------------------|-------------------|-----------------|---|
| 0 | 2019.0     | 0.1622           | 0.6656            | 0.7119          |   |
| 1 | 1956.0     | 0.1238           | 0.1866            | 0.2416          |   |
| 2 | 1709.0     | 0.1444           | 0.4245            | 0.4504          |   |
| 3 | 567.7      | 0.2098           | 0.8663            | 0.6869          |   |
| 4 | 1575.0     | 0.1374           | 0.2050            | 0.4000          |   |

|   | concave points_worst | symmetry_worst | fractal_dimension_worst |
|---|----------------------|----------------|-------------------------|
| 0 | 0.2654               | 0.4601         | 0.11890                 |

|   |        |        |         |
|---|--------|--------|---------|
| 1 | 0.1860 | 0.2750 | 0.08902 |
| 2 | 0.2430 | 0.3613 | 0.08758 |
| 3 | 0.2575 | 0.6638 | 0.17300 |
| 4 | 0.1625 | 0.2364 | 0.07678 |

[5 rows x 31 columns]

1. Lakukan proses encoding pada kolom "diagnosis".

```
from sklearn.preprocessing import LabelEncoder, StandardScaler

le = LabelEncoder() # membuat objek dari LabelEncoder
df['diagnosis'] = le.fit_transform(df['diagnosis']) # proses encoding

df.head()
```

|   | diagnosis | radius_mean | texture_mean | perimeter_mean | area_mean | \ |
|---|-----------|-------------|--------------|----------------|-----------|---|
| 0 | 1         | 17.99       | 10.38        | 122.80         | 1001.0    |   |
| 1 | 1         | 20.57       | 17.77        | 132.90         | 1326.0    |   |
| 2 | 1         | 19.69       | 21.25        | 130.00         | 1203.0    |   |
| 3 | 1         | 11.42       | 20.38        | 77.58          | 386.1     |   |
| 4 | 1         | 20.29       | 14.34        | 135.10         | 1297.0    |   |

|             | smoothness_mean | compactness_mean | concavity_mean | concave |
|-------------|-----------------|------------------|----------------|---------|
| points_mean |                 |                  |                |         |
| 0           | 0.11840         | 0.27760          | 0.3001         |         |
| 0.14710     |                 |                  |                |         |
| 1           | 0.08474         | 0.07864          | 0.0869         |         |
| 0.07017     |                 |                  |                |         |
| 2           | 0.10960         | 0.15990          | 0.1974         |         |
| 0.12790     |                 |                  |                |         |
| 3           | 0.14250         | 0.28390          | 0.2414         |         |
| 0.10520     |                 |                  |                |         |
| 4           | 0.10030         | 0.13280          | 0.1980         |         |
| 0.10430     |                 |                  |                |         |

|   | symmetry_mean | ... | radius_worst | texture_worst | perimeter_worst | \ |
|---|---------------|-----|--------------|---------------|-----------------|---|
| 0 | 0.2419        | ... | 25.38        | 17.33         | 184.60          |   |
| 1 | 0.1812        | ... | 24.99        | 23.41         | 158.80          |   |
| 2 | 0.2069        | ... | 23.57        | 25.53         | 152.50          |   |
| 3 | 0.2597        | ... | 14.91        | 26.50         | 98.87           |   |
| 4 | 0.1809        | ... | 22.54        | 16.67         | 152.20          |   |

|   | area_worst | smoothness_worst | compactness_worst | concavity_worst | \ |
|---|------------|------------------|-------------------|-----------------|---|
| 0 | 2019.0     | 0.1622           | 0.6656            | 0.7119          |   |
| 1 | 1956.0     | 0.1238           | 0.1866            | 0.2416          |   |
| 2 | 1709.0     | 0.1444           | 0.4245            | 0.4504          |   |
| 3 | 567.7      | 0.2098           | 0.8663            | 0.6869          |   |
| 4 | 1575.0     | 0.1374           | 0.2050            | 0.4000          |   |

| concave | points_worst | symmetry_worst | fractal_dimension_worst |
|---------|--------------|----------------|-------------------------|
|---------|--------------|----------------|-------------------------|

|   |        |        |         |
|---|--------|--------|---------|
| 0 | 0.2654 | 0.4601 | 0.11890 |
| 1 | 0.1860 | 0.2750 | 0.08902 |
| 2 | 0.2430 | 0.3613 | 0.08758 |
| 3 | 0.2575 | 0.6638 | 0.17300 |
| 4 | 0.1625 | 0.2364 | 0.07678 |

[5 rows x 31 columns]

1. Lakukan proses standarisasi pada semua kolom yang memiliki nilai numerik.

```
num_colm = df.drop(columns=['diagnosis']).columns
std = StandardScaler()
df[num_colm] = std.fit_transform(df[num_colm])
df.head()
```

|   | diagnosis | radius_mean | texture_mean | perimeter_mean | area_mean \ |
|---|-----------|-------------|--------------|----------------|-------------|
| 0 | 1.297676  | 1.097064    | -2.073335    | 1.269934       | 0.984375    |
| 1 | 1.297676  | 1.829821    | -0.353632    | 1.685955       | 1.908708    |
| 2 | 1.297676  | 1.579888    | 0.456187     | 1.566503       | 1.558884    |
| 3 | 1.297676  | -0.768909   | 0.253732     | -0.592687      | -0.764464   |
| 4 | 1.297676  | 1.750297    | -1.151816    | 1.776573       | 1.826229    |

|   | smoothness_mean | compactness_mean | concavity_mean | concave points_mean \ |
|---|-----------------|------------------|----------------|-----------------------|
| 0 | 1.568466        | 3.283515         | 2.652874       | 2.532475              |
| 1 | -0.826962       | -0.487072        | -0.023846      | 0.548144              |
| 2 | 0.942210        | 1.052926         | 1.363478       | 2.037231              |
| 3 | 3.283553        | 3.402909         | 1.915897       | 1.451707              |
| 4 | 0.280372        | 0.539340         | 1.371011       | 1.428493              |

|   | symmetry_mean | ... | radius_worst | texture_worst | perimeter_worst \ |
|---|---------------|-----|--------------|---------------|-------------------|
| 0 | 2.217515      | ... | 1.886690     | -1.359293     | 2.303601          |
| 1 | 0.001392      | ... | 1.805927     | -0.369203     | 1.535126          |
| 2 | 0.939685      | ... | 1.511870     | -0.023974     | 1.347475          |
| 3 | 2.867383      | ... | -0.281464    | 0.133984      | -0.249939         |
| 4 | -0.009560     | ... | 1.298575     | -1.466770     | 1.338539          |

|   | area_worst | smoothness_worst | compactness_worst | concavity_worst \ |
|---|------------|------------------|-------------------|-------------------|
| 0 | 2.001237   | 1.307686         | 2.616665          | 2.109526          |
| 1 | 1.890489   | -0.375612        | -0.430444         | -0.146749         |
| 2 | 1.456285   | 0.527407         | 1.082932          | 0.854974          |
| 3 | -0.550021  | 3.394275         | 3.893397          | 1.989588          |
| 4 | 1.220724   | 0.220556         | -0.313395         | 0.613179          |

|  | concave points_worst | symmetry_worst | fractal_dimension_worst |
|--|----------------------|----------------|-------------------------|
|--|----------------------|----------------|-------------------------|

|   |          |           |           |
|---|----------|-----------|-----------|
| 0 | 2.296076 | 2.750622  | 1.937015  |
| 1 | 1.087084 | -0.243890 | 0.281190  |
| 2 | 1.955000 | 1.152255  | 0.201391  |
| 3 | 2.175786 | 6.046041  | 4.935010  |
| 4 | 0.729259 | -0.868353 | -0.397100 |

[5 rows x 31 columns]

1. Lakukan proses stratified split data untuk membuat data latih dan data uji dengan rasio 80:20.

```
#Split data
from sklearn.model_selection import train_test_split

#Split data training dan dan lainnya
#data lainnya, akan kita split lagi menjadi validasi dan testing.
#Rasio yang akan kita gunakan adalah 8:1:1
df_train, df_unseen = train_test_split(df, test_size=0.2,
random_state=0, stratify=df['diagnosis'])

#Split lagi antara validasi dan testing
df_val, df_test = train_test_split(df_unseen, test_size=0.5,
random_state=0, stratify=df_unseen['diagnosis'])

#Cek masing-masing ukuran data

print(f'Jumlah label data asli:\n{df.diagnosis.value_counts()}')
print(f'Jumlah label data train:\n{df_train.diagnosis.value_counts()}')
print(f'Jumlah label data val:\n{df_val.diagnosis.value_counts()}')
print(f'Jumlah label data test:\n{df_test.diagnosis.value_counts()}')
```

Jumlah label data asli:

|           |     |
|-----------|-----|
| diagnosis |     |
| -0.770609 | 357 |
| 1.297676  | 212 |

Name: count, dtype: int64

Jumlah label data train:

|           |     |
|-----------|-----|
| diagnosis |     |
| -0.770609 | 285 |
| 1.297676  | 170 |

Name: count, dtype: int64

Jumlah label data val:

|           |    |
|-----------|----|
| diagnosis |    |
| -0.770609 | 36 |
| 1.297676  | 21 |

Name: count, dtype: int64

Jumlah label data test:

|           |    |
|-----------|----|
| diagnosis |    |
| -0.770609 | 36 |

```
1.297676    21  
Name: count, dtype: int64
```