

Exploratory Data Analysis - Data Preparation

6604067636348

ముబ్రీస్రవార్ష టైటిల్ గాలుకా



[8] # Install and Import ydata_profiling (Auto EDA tools)
!pip install ydata_profiling
from ydata_profiling import ProfileReport

Show hidden output
▶ ติดตั้ง library "ydata_profiling" ใน Google Colab
▶ นำเข้า class "ProfileReport" ลง library "ydata_profiling"

[6] import numpy as np
import tensorflow as tf
from tensorflow import keras
import pandas as pd
from matplotlib import pyplot as plt
%matplotlib inline

▶ นำเข้า library "Numpy" และตัวแปร np ใช้สำหรับตรี噪และจัดการ Array
▶ นำเข้า library "TensorFlow" และตัวแปร tf ใช้ Machine Learning และ Deep Learning
▶ นำเข้า Module "keras" ลง TensorFlow เพื่อใช้สร้าง Model "Neural Network"
▶ นำเข้า library "Pandas" และตัวแปร pd สำหรับจัดการข้อมูลเชิงstructured (Data Flow)
▶ นำเข้า Module "pyplot" ลง library "Matplotlib" และตัวชี้บอร์ด plt ผ่านช่องทาง Matplotlib ของ Python

// คลิปบันทึก Jupyter / Colab เมื่อ 9 เดือนที่แล้ว บน "Matplotlib" ประสบความ notebooke กันดี

▶ นำเข้า module "files" ลง module google.colab เป็นตัวรับไฟล์ข้อมูลไฟล์ใน google.colab

Show hidden output

1. เรียกใช้เมธอด upload() ต้องติดตั้งลงเครื่องก่อนใน google.colab

2. นำร่องตัวต่อไฟล์ CSV ที่ได้รับโดย upload() ลงเป็น dictionary ให้เข้าใจว่ามีค่าอะไรบ้าง

[11] df = pd.read_csv("/content/xxx_dataset.csv")

▶ ใช้ฟังก์ชัน "read_csv" ลง library pandas เพื่อโหลดไฟล์ CSV ที่อยู่ใน path /content/xxx_dataset.csv

* ไฟล์ CSV ถูกแปลงเป็น DataFrame (โครงสร้าง) และเก็บไว้ในตัวแปร df ซึ่งใช้สำหรับการทำสถิติและจัดการข้อมูลใน Python

[12] # Generate the Profiling Report
profile = ProfileReport(
 df, title = "xxx_dataset", html={"style": {"full_width": True}}, sort=None)

▶ ตั้งแต่เริ่มต้นต้องมีไฟล์ "xxx_dataset"

▶ ไฟล์ DataFrame (df) ห้ามลบเด็ดขาด

1. รันคำสั่ง generate_profile() หน่วย "ProfileReport" ลง library "ydata_profiling"

2. รับผลลัพธ์ของตัวชี้วัด ภาระ computational, และ การตรวจสอบข้อมูล (Missing values)

[13] # The HTML report in an iframe
profile.to_notebook_iframe()

▶ เนื่องจากเราต้องติดตั้ง profile ลงใน Jupyter Notebook บน Google Colab

```
[14] df.shape  
→ เป็นค่าที่ใช้ตรวจสอบ DataFrame df  
⇒ (400, 11) // คือมี tuple 400 บรรทัด, 11 คอลัมน์
```

```
[15] df.columns  
→ ชื่อหัวข้อ ของคอลัมน์ใน DataFrame df  
⇒ Index(['gender', 'age', 'hypertension', 'heart_disease', 'Marriage',  
          'work_type', 'Living_type', 'avg_glucose', 'bmi', 'smoking_status',  
          'illness'],  
          dtype='object') // คือเป็น index object ที่มีค่าเป็นชื่อคอลัมน์เช่น gender, age, hypertension, heart_disease, Marriage, work_type, Living_type, avg_glucose, bmi, smoking_status, illness
```

```
[16] df.isnull().any()  
→ ตรวจสอบว่า ค่าข้อมูล (missing values) ในแต่ละ column ของ DataFrame df ยังไง  
⇒  


|                | θ     |
|----------------|-------|
| gender         | False |
| age            | False |
| hypertension   | False |
| heart_disease  | False |
| Marriage       | False |
| work_type      | False |
| Living_type    | False |
| avg_glucose    | False |
| bmi            | False |
| smoking_status | False |
| illness        | False |



dtype: bool



→ คือค่าเป็น Series ให้เห็นว่าค่าอะไรบ้างที่ True หรือ False



- True คือ missing ใน column นั้น
- False ไม่มีค่า missing ใน column นั้น

```

```
[17] df.unique()  
→ ตรวจสอบ จำนวนค่าต่างๆ (unique value) ในแต่ละ column ของ DataFrame df  
⇒  


|                | θ   |
|----------------|-----|
| gender         | 2   |
| age            | 79  |
| hypertension   | 2   |
| heart_disease  | 2   |
| Marriage       | 2   |
| work_type      | 4   |
| Living_type    | 2   |
| avg_glucose    | 383 |
| bmi            | 199 |
| smoking_status | 4   |
| illness        | 2   |



dtype: int64



→ คือค่าเป็น series ที่แสดงจำนวน ค่าต่างๆ ที่มีอยู่ในแต่ละ คอลัมน์


```

```
[18] df.describe(include = 'all')
```

	gender	age	hypertension	heart_disease	Marriage	work_type	Living_type	avg_glucose	bmi	smoking_status	illness
count	400	400.00000	400.000000	400.000000	400	400	400	400.000000	400.000000	400	400.000000
unique	2	NaN	NaN	NaN	2	4	2	NaN	NaN	4	NaN
top	Female	NaN	NaN	NaN	Yes	Private	Urban	NaN	NaN	never smoked	NaN
freq	214	NaN	NaN	NaN	306	231	201	NaN	NaN	164	NaN
mean	NaN	55.26780	0.180000	0.132500	NaN	NaN	NaN	119.391950	29.481750	NaN	0.500000
std	NaN	22.51279	0.384669	0.339458	NaN	NaN	NaN	54.377459	6.488354	NaN	0.500626
min	NaN	0.80000	0.000000	0.000000	NaN	NaN	NaN	56.070000	15.600000	NaN	0.000000
25%	NaN	44.00000	0.000000	0.000000	NaN	NaN	NaN	80.460000	25.575000	NaN	0.000000
50%	NaN	59.00000	0.000000	0.000000	NaN	NaN	NaN	97.665000	28.600000	NaN	0.500000
75%	NaN	74.25000	0.000000	0.000000	NaN	NaN	NaN	144.345000	33.025000	NaN	1.000000
max	NaN	82.00000	1.000000	1.000000	NaN	NaN	NaN	271.740000	48.900000	NaN	1.000000

1. វិភាគនៃព័ត៌មានរបស់ក្រុមហ៊ុនបញ្ជូនទិន្នន័យ Data Frame នៃ df

2. សារព័ត៌មានចាប់អាមេរិក នៅក្នុងការបង្ហាញ នៅក្នុងការបង្ហាញ ជា Max , ជា Min , និងសំណងការការពារ ដែលបានសរុបដោយការបង្ហាញ

3. រាយការ include = 'all' នៅលើ colwise ដើម្បី គ្រប់គ្រង នៅក្នុងការបង្ហាញ (ខ្លះ colwise ការបង្ហាញរបស់ក្រុមហ៊ុនបានលើក)

នៅក្នុងការបង្ហាញ ទាំងឡាយ គ្រប់គ្រង នៅក្នុងការបង្ហាញ NULL , និងសំណងការបង្ហាញរបស់ក្រុមហ៊ុន

```

[50]: df['gender'] = df['gender'].map({'Female': 1, 'Male': 0})
line 2 : df['Marriage'] = df['Marriage'].map({'Yes': 1, 'No': 0})
line 3 : df['work_type'] = df['work_type'].map({'Private': 0, 'Self-employed': 1, 'Govt_job': 2, 'children': 3})
line 4 : df['Living_type'] = df['Living_type'].map({'Urban': 1, 'Rural': 0})
line 5 : df['smoking_status'] = df['smoking_status'].map({'formerly smoked': 0, 'never smoked': 1, 'smoked': 2, 'Unknown': 3})

```

line 1 • ផ្លូវការ gender តាមលក្ខណៈ Female នឹង 1 និង Male នឹង 0

- នីមួយៗ map() ដែលត្រូវបានរាយជាតិមុខ ទាំងអស់

line 2 • ផ្លូវការ Marriage តាមលក្ខណៈ Yes នឹង 1 និង No នឹង 0

- នីមួយៗ map() ដែលត្រូវបានរាយជាតិមុខ ទាំងអស់

line 3 • ផ្លូវការលើ column work_type តាមលក្ខណៈប្រភេទការងារដែលត្រូវបានរាយជាតិមុខ Private នឹង 0, Self-employed នឹង 1 ពីរប៉ុណ្ណោះ

line 4 • ផ្លូវការលើ column Living_type តាមលក្ខណៈ Urban នឹង 1 និង Rural នឹង 0

- នីមួយៗ map() ដែលត្រូវបានរាយជាតិមុខ ទាំងអស់

line 5 • ផ្លូវការលើ column smoking_status តាមលក្ខណៈ formerly smoked នឹង 0, never smoked នឹង 1, smoked នឹង 2, និង Unknown នឹង 3

[55] df

	gender	age	hypertension	heart_disease	Marriage	work_type	Living_type	avg_glucose	bmi	smoking_status	illness	
0	0	67.0	0	1	1	0	1	228.69	36.6	0	1	
1	0	80.0	0	1	1	0	0	105.92	32.5	1	1	
2	1	49.0	0	0	1	0	1	171.23	34.4	2	1	
3	1	79.0	1	0	1	1	0	174.12	24.0	1	1	
4	0	81.0	0	0	1	0	1	186.21	29.0	0	1	
...	
395	1	54.0	0	1	1	0	1	140.28	37.1	0	0	
396	0	67.0	0	0	1	2	0	244.28	29.4	0	0	
397	0	53.0	0	0	1	0	1	124.16	31.7	1	0	
398	0	47.0	0	0	1	0	0	93.55	31.4	1	0	
399	0	44.0	0	0	1	0	1	99.34	33.1	1	0	

→ Column នេះ នឹងត្រូវបានរាយជាតិមុខ ទាំងអស់ ដើម្បីធ្វើការសម្រេចការគ្រែកសាន្តនៃការសរស់របស់ Model Machine Learning ក្នុងការអនុវត្តការណ៍

```
[23] y = df['illness']  
x = df.drop(columns=['illness'])  
→ ก็จะมี column illness หายไป เป็น column X ของ DataFrame df  
  
• ถ้าแยกตัวไป x เป็น DataFrame ที่ไม่รวม column illness
```

- หรือ `drop(columns=['illness'])` จะนำ column illness ออก dari df

```
[97] x.head()  
+---+-----+-----+-----+-----+-----+-----+-----+-----+-----+  
| 0 | gender | age | hypertension | heart_disease | Marriage | work_type | Living_type | avg_glucose | bmi | smoking_status |  
+---+-----+-----+-----+-----+-----+-----+-----+-----+-----+  
| 0 | 0 | 67.0 | 0 | 1 | 1 | 0 | 1 | 228.69 | 36.6 | 0 |  
| 1 | 0 | 80.0 | 0 | 0 | 1 | 1 | 0 | 105.92 | 32.5 | 1 |  
| 2 | 1 | 49.0 | 0 | 0 | 1 | 0 | 0 | 171.23 | 34.4 | 2 |  
| 3 | 1 | 79.0 | 1 | 0 | 1 | 1 | 0 | 174.12 | 24.0 | 1 |  
| 4 | 0 | 81.0 | 0 | 0 | 1 | 0 | 1 | 188.21 | 38.0 | 0 |  
+---+-----+-----+-----+-----+-----+-----+-----+-----+-----+
```

```
[99] y.value_counts()  
+---+  
| count |  
| illness |  
+---+  
| 1 | 200 |  
| 0 | 200 |  
+---+  
dtype: int64  
→ คุณสามารถเห็นว่า column y ซึ่งเป็น illness นั้น column illness ที่เราต้องการจะใช้เป็น target variable (ตัวที่ต้องการคำนวณ)
```

```

01 from sklearn.preprocessing import MinMaxScaler // line 1
02
03 scale=MinMaxScaler() // line 2
04 x_scaled=pd.DataFrame(scale.fit_transform(x),columns=x.columns) // line 3
05 x_scaled // line 4

```

	gender	age	hypertension	heart_disease	Marriage	work_type	Living_type	avg_glucose	bmi	smoking_status
0	0.0	0.815271	0.0	1.0	1.0	0.000000	1.0	0.800389	0.630631	0.000000
1	0.0	0.975369	0.0	1.0	1.0	0.000000	0.0	0.231140	0.507508	0.333333
2	1.0	0.593596	0.0	0.0	1.0	0.000000	1.0	0.533964	0.564565	0.666667
3	1.0	0.963054	1.0	0.0	1.0	0.333333	0.0	0.547364	0.252252	0.333333
4	0.0	0.987685	0.0	0.0	1.0	0.000000	1.0	0.603422	0.402402	0.000000
...
395	1.0	0.655172	0.0	1.0	1.0	0.000000	1.0	0.390458	0.645646	0.000000
396	0.0	0.815271	0.0	0.0	1.0	0.666667	0.0	0.872676	0.414414	0.000000
397	0.0	0.642857	0.0	0.0	1.0	0.000000	1.0	0.315714	0.483483	0.333333
398	0.0	0.568966	0.0	0.0	1.0	0.000000	0.0	0.173784	0.474474	0.333333
399	0.0	0.532020	0.0	0.0	1.0	0.000000	1.0	0.200631	0.525526	0.333333

400 rows × 10 columns

line 1 ผู้เข้า MinMaxScaler ณ sklearn.preprocessing ชี้ให้เรียกข้อมูลไปยังตัวที่เราต้องการใช้ใน Class หลักต่อไปนี้ 0 ถึง 17 ที่เราได้ระบุไว้ตอนนี้จะเป็นตัวแปรที่ใช้ใน model Machine Learning

จุดนี้ ภารกิจจะยกเว้นค่า 0 หรือ 1 ที่ไม่สามารถคำนวณได้

line 2 สร้างตัวแปร scale เป็น object นั้น object จะเป็น MinMax Scaler ซึ่งจะใช้มาแทนตัวของค่า

line 3 • ใช้ fit_transform(x) เนื่องด้วยตัวนี้คือการนำข้อมูลตัวอย่างใน DataFrame x และเปลี่ยนรูปแบบ x ให้มีต้องบูนิ่ง 0 ถึง 1 โดยใช้ MinMax Scaler

- แล้วเราจะเปลี่ยนรูปแบบข้อมูลนี้เป็น DataFrame ใหม่ x_scaled โดยให้ชื่อ columns ตามเดิม x.columns

line 4 • ต้องขออนุญาติใน x_scaled จะถูกยกเว้นตัว 0,1 ซึ่งจะถูกยกเว้นโดยทางด้านนี้เองที่ต้องการกันมาก

- รูปแบบนี้จะถูกนำไปใช้ใน model Machine Learning หรือใน model ที่สามารถรับข้อมูลได้ลักษณะ

```
[103] from sklearn.model_selection import train_test_split // line 1  
x_train, x_test, y_train, y_test = train_test_split(x_scaled, y, test_size=0.30, random_state=42) // line 2
```

line 1 ការប្រើប្រាស់ `train-test-split` នៃ `sklearn.model_selection` ដើម្បីបង្កើតឡាយចំណែកលើការបង្ហាញ (training set), និងសមាស្រាវ (test set) ដើម្បីរាយការណ៍នៃការសរស់សរស់ `Machine Learning`

line 2 • `train-test-split()` បង្កើតឡាយចំណែក `x-scaled` នាម `y` ទាំងអស់ ដោយអាចរាយការណ៍តាមលានសាលាដែលបានរាយការណ៍ឡើង។ `test_size` គឺជាទឹកចំណែក 30% រាយការណ៍ឡើង។

- ក្នុងការប្រើប្រាស់ `train-test-split()` មានរាយការណ៍ខាងក្រោម

- `x_train` ចំណែកដែលមានសាលាដែលបានរាយការណ៍ឡើង

- `x-test` ចំណែកដែលមិនមែនសាលាដែលបានរាយការណ៍ឡើង

- `y-train` តារាងតម្រាករបស់ `y` ដែលបានរាយការណ៍ឡើង

- `y-test` តារាងតម្រាករបស់ `y` ដែលមិនមែនសាលាដែលបានរាយការណ៍ឡើង

- `test_size = 0.30` ការអនុវត្តនភាពចូលរួមទៅក្នុងការបង្ហាញ 30% ទូទៅនូវការបង្ហាញ (តាមព័ត៌មាធទី 0 ទៅមិនមែនទី 1)

- `random_state = 42` ការអនុវត្តនភាព `random_state` ដើម្បីការបង្ហាញបានលាស់ឡើងឡើង ដែលមិនមែនតាមលទ្ធផលដែលក្នុងការបង្ហាញ។ `random_state` តាមព័ត៌មាធទី 0 ទៅ 999 ដើម្បីការបង្ហាញបានលាស់ឡើងឡើង។

Example

សម្រាប់រាយការណ៍ដែលបានបង្ហាញ `x-scaled` នាមឈ្មោះ `target` `y`

- នៅលើ `run code` ចំណែកលើការបង្ហាញ `x` ទូទៅ

1. 70% នៃចំណែកនេះនឹងរាយការណ៍ (train)

2. 30% នៃចំណែកនេះនឹងរាយការណ៍ឡើង (test)

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text RAM Disk Gemini

```
# Install and Import ydata_profiling (Auto EDA tools)
!pip install ydata_profiling
from ydata_profiling import ProfileReport
```

Requirement already satisfied: ydata_profiling in /usr/local/lib/python3.10/dist-packages (4.12.1)

Requirement already satisfied: scipy<1.14,>=1.4.1 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (1.13.1)

Requirement already satisfied: pandas!=1.4.0,<3.1.1 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (2.2.2)

Requirement already satisfied: matplotlib<3.10,>=3.5 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (3.9.4)

Requirement already satisfied: pydantic>=2 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (2.10.4)

Requirement already satisfied: PyYAML<6.1,>=5.0.0 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (6.0.2)

Requirement already satisfied: jinja2<3.2,>=2.11.1 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (3.1.5)

Requirement already satisfied: visions<0.7.7,>=0.7.5 in /usr/local/lib/python3.10/dist-packages (from visions[`type_image_path`]<0.7.7,>=0.7.5->ydata_profiling) (0.7.6)

Requirement already satisfied: numpy<2.2,>=1.16.0 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (1.26.4)

Requirement already satisfied: htmlmin==0.1.12 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (0.1.12)

Requirement already satisfied: phik<0.13,>=0.11.1 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (0.12.4)

Requirement already satisfied: requests<3,>=2.24.0 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (2.32.3)

Requirement already satisfied: tqdm<5,>=4.48.2 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (4.67.1)

Requirement already satisfied: seaborn<0.14,>=0.10.1 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (0.13.2)

Requirement already satisfied: multimethod<2,>=1.4 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (1.12)

Requirement already satisfied: statsmodels<1,>=0.13.2 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (0.14.4)

Requirement already satisfied: typeguard<5,>=3 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (4.4.1)

Requirement already satisfied: imagehash==4.3.1 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (4.3.1)

Requirement already satisfied: wordcloud>=1.9.3 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (1.9.4)

Requirement already satisfied: dacite>=1.8 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (1.8.1)

Requirement already satisfied: numba<1,>=0.56.0 in /usr/local/lib/python3.10/dist-packages (from ydata_profiling) (0.60.0)

Requirement already satisfied: PyLavender in /usr/local/lib/python3.10/dist-packages (from imagehash==4.3.1->ydata_profiling) (1.8.0)

Requirement already satisfied: pillow in /usr/local/lib/python3.10/dist-packages (from imagehash==4.3.1->ydata_profiling) (11.1.0)

Requirement already satisfied: MarkupSafe==2.0 in /usr/local/lib/python3.10/dist-packages (from jinja2<3.2,>=2.11.1->ydata_profiling) (3.0.2)

Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata_profiling) (1.3.1)

Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata_profiling) (0.12.1)

Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata_profiling) (4.55.3)

Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata_profiling) (1.4.8)

Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata_profiling) (24.2)

Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata_profiling) (3.2.1)

Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib<3.10,>=3.5->ydata_profiling) (2.8.2)

Requirement already satisfied: llvmlite<0.44,>=0.43.0dev in /usr/local/lib/python3.10/dist-packages (from numba<1,>=0.56.0->ydata_profiling) (0.43.0)

Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas!=1.4.0,<3,>1.1->ydata_profiling) (2024.2)

Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.10/dist-packages (from pandas!=1.4.0,<3,>1.1->ydata_profiling) (2024.2)

Requirement already satisfied: joblib>=0.14.1 in /usr/local/lib/python3.10/dist-packages (from phik<0.13,>=0.11.1->ydata_profiling) (1.4.2)

Requirement already satisfied: annotated-types>=0.6.0 in /usr/local/lib/python3.10/dist-packages (from pydantic>=2->ydata_profiling) (0.7.0)

Requirement already satisfied: pydantic-core==2.27.2 in /usr/local/lib/python3.10/dist-packages (from pydantic>=2->ydata_profiling) (2.27.2)

Requirement already satisfied: typing_extensions>=4.12.2 in /usr/local/lib/python3.10/dist-packages (from pydantic>=2->ydata_profiling) (4.12.2)

Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.24.0->ydata_profiling) (3.4.1)

Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.24.0->ydata_profiling) (3.10)

Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.24.0->ydata_profiling) (2.3.0)

Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.24.0->ydata_profiling) (2024.12.14)

Requirement already satisfied: patry>=0.5.6 in /usr/local/lib/python3.10/dist-packages (from statsmodels<1,>=0.13.2->ydata_profiling) (1.0.1)

Requirement already satisfied: attrs>=19.3.0 in /usr/local/lib/python3.10/dist-packages (from visions[`type_image_path`]<0.7.7,>=0.7.5->ydata_profiling) (24.3.0)

Requirement already satisfied: networkx>=2.4 in /usr/local/lib/python3.10/dist-packages (from visions[`type_image_path`]<0.7.7,>=0.7.5->ydata_profiling) (3.4.2)

Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.7->matplotlib<3.10,>=3.5->ydata_profiling) (1.17.0)

```
[84] import numpy as np
import tensorflow as tf
from tensorflow import keras
import pandas as pd
from matplotlib import pyplot as plt
%matplotlib inline
```

```
[85] from google.colab import files
uploaded = files.upload() # download CSV file and upload to your colab
```

Choose File xxx_dataset.csv

- xxx_dataset.csv (ext/csv) - 23316 bytes, last modified: 1/9/2025 - 100% done

Saving xxx_dataset.csv to xxx_dataset.csv

```
[86] df = pd.read_csv("/content/xxx_dataset.csv")
```

```
[87] # Generate the Profiling Report
profile = ProfileReport(
    df, title = "xxx_dataset" , html={"style": {"full_width": True}}, sort=None
)
```

```
[88] # The HTML report in an iframe
profile.to_notebook_iframe()
```

Summarize dataset: 100% [29/29] [00:03<00:00, 4.04it/s, Completed]

Generate report structure: 100% [1/1] [00:03<00:00, 3.91s/it]

Render HTML: 100% [1/1] [00:00<00:00, 1.38it/s]

xxx_dataset

Overview Variables Interactions Correlations Missing values Sample

Overview

Brought to you by YData

Overview Alerts 5 Reproduction

Dataset statistics

Number of variables	11
Number of observations	400
Missing cells	0
Missing cells (%)	0.0%
Duplicate rows	0
Duplicate rows (%)	0.0%
Total size in memory	34.5 KiB
Average record size in memory	88.3 B

Variable types

Categorical	7
Numeric	3
Boolean	1

Variables

Select Columns ▾

✓ [89] df.shape

→ (400, 11)

[90] df.columns

```
Index(['gender', 'age', 'hypertension', 'heart_disease', 'Marriage',  
       'work_type', 'Living_type', 'avg_glucose', 'bmi', 'smoking_status',  
       'illness'],  
      dtype='object')
```

[91] df.isnull().any()

gender	False
age	False
hypertension	False
heart_disease	False
Marriage	False
work_type	False
Living_type	False
avg_glucose	False
bmi	False
smoking_status	False
illness	False

dtype: bool

✓ [92] df.nunique()

gender	2
age	79
hypertension	2
heart_disease	2
Marriage	2
work_type	4
Living_type	2
avg_glucose	393
bmi	199
smoking_status	4
illness	2

dtype: int64

✓ [93]: df.describe(include = 'all')

	gender	age	hypertension	heart_disease	Marriage	work_type	Living_type	avg_glucose	bmi	smoking_status	illness
count	400	400.000000	400.000000	400.000000	400	400	400	400.000000	400.000000	400	400.000000
unique	2	NaN	NaN	NaN	2	4	2	NaN	NaN	4	NaN
top	Female	NaN	NaN	NaN	Yes	Private	Urban	NaN	NaN	never smoked	NaN
freq	214	NaN	NaN	NaN	306	231	201	NaN	NaN	164	NaN
mean	NaN	55.26780	0.180000	0.132500	NaN	NaN	NaN	119.391950	29.481750	NaN	0.500000
std	NaN	22.51279	0.384669	0.339458	NaN	NaN	NaN	54.377459	6.488354	NaN	0.500626
min	NaN	0.80000	0.000000	0.000000	NaN	NaN	NaN	56.070000	15.600000	NaN	0.000000
25%	NaN	44.00000	0.000000	0.000000	NaN	NaN	NaN	80.460000	25.575000	NaN	0.000000
50%	NaN	59.00000	0.000000	0.000000	NaN	NaN	NaN	97.665000	28.600000	NaN	0.500000
75%	NaN	74.25000	0.000000	0.000000	NaN	NaN	NaN	144.345000	33.025000	NaN	1.000000
max	NaN	82.00000	1.000000	1.000000	NaN	NaN	NaN	271.740000	48.900000	NaN	1.000000

```
[94]: df['gender'] = df['gender'].map({'Female': 1, 'Male': 0})
df['Marriage'] = df['Marriage'].map({'Yes': 1, 'No': 0})
df['work_type'] = df['work type'].map({'Private': 0, 'Self-employed': 1, 'Govt_job': 2, 'children': 3})
df['Living_type'] = df['Living type'].map({'Urban': 1, 'Rural': 0})
df['smoking_status'] = df['smoking status'].map({'formerly smoked': 0, 'never smoked': 1, 'smokes': 2, 'Unknown': 3})
```

	gender	age	hypertension	heart_disease	Marriage	work_type	Living_type	avg_glucose	bmi	smoking_status	illness
0	0	67.0	0	1	1	0	1	228.69	36.6	0	1
1	0	80.0	0	1	1	0	0	105.92	32.5	1	1
2	1	49.0	0	0	1	0	1	171.23	34.4	2	1
3	1	79.0	1	0	1	1	0	174.12	24.0	1	1
4	0	81.0	0	0	1	0	1	186.21	29.0	0	1

```
395    1  54.0      0      1      1      0      1    140.28  37.1      0      0
396    0  67.0      0      0      1      2      0    244.28  29.4      0      0
397    0  53.0      0      0      1      0      1    124.16  31.7      1      0
398    0  47.0      0      0      1      0      0    93.55  31.4      1      0
399    0  44.0      0      0      1      0      1    99.34  33.1      1      0
```

400 rows × 11 columns

Next steps: [Generate code with df](#) [View recommended plots](#) [New interactive sheet](#)

```
[96] y = df['illness']
x = df.drop(columns=['illness'])
```

```
[97] x.head()
```

	gender	age	hypertension	heart_disease	Marriage	work_type	Living_type	avg_glucose	bmi	smoking_status
0	0	67.0	0	1	1	0	1	228.69	36.6	0
1	0	80.0	0	1	1	0	0	105.92	32.5	1
2	1	49.0	0	0	1	0	1	171.23	34.4	2
3	1	79.0	1	0	1	1	0	174.12	24.0	1
4	0	81.0	0	0	1	0	1	186.21	29.0	0

Next steps: [Generate code with x](#) [View recommended plots](#) [New interactive sheet](#)

```
[99] y.value_counts()
```

```
count
```

illness	count
1	200
0	200

dtype: int64

```
from sklearn.preprocessing import MinMaxScaler
scale=MinMaxScaler()
x_scaled=pd.DataFrame(scale.fit_transform(x),columns=x.columns)
x_scaled
```

	gender	age	hypertension	heart_disease	Marriage	work_type	Living_type	avg_glucose	bmi	smoking_status
0	0.0	0.815271	0.0	1.0	1.0	0.000000	1.0	0.800389	0.630631	0.000000
1	0.0	0.975369	0.0	1.0	1.0	0.000000	0.0	0.231140	0.507508	0.333333
2	1.0	0.593596	0.0	0.0	1.0	0.000000	1.0	0.533964	0.564565	0.666667
3	1.0	0.963054	1.0	0.0	1.0	0.333333	0.0	0.547364	0.252252	0.333333
4	0.0	0.987685	0.0	0.0	1.0	0.000000	1.0	0.603422	0.402402	0.000000
...
395	1.0	0.655172	0.0	1.0	1.0	0.000000	1.0	0.390458	0.645646	0.000000
396	0.0	0.815271	0.0	0.0	1.0	0.666667	0.0	0.872676	0.414414	0.000000
397	0.0	0.642857	0.0	0.0	1.0	0.000000	1.0	0.315714	0.483483	0.333333
398	0.0	0.568966	0.0	0.0	1.0	0.000000	0.0	0.173784	0.474474	0.333333
399	0.0	0.532020	0.0	0.0	1.0	0.000000	1.0	0.200631	0.525526	0.333333

400 rows × 10 columns

Next steps: [Generate code with x_scaled](#) [View recommended plots](#) [New interactive sheet](#)

```
[103] from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x_scaled, y, test_size=0.30, random_state=42)
```

Colab paid products - Cancel contracts here

✓ 0s completed at 4:46AM