

Mục lục

1	Biến cố ngẫu nhiên và xác suất	1
1.1	Khái niệm	1
1.2	Mô hình xác suất cổ điển	3
1.3	Mô hình xác suất hình học	7
1.4	Công thức cộng và nhân xác suất	9
1.5	Công thức xác suất đầy đủ và công thức Bayes	15
1.6	Dãy thử Bernoulli	17
2	Đại lượng ngẫu nhiên	22
2.1	Khái niệm	23
2.2	Hàm phân bố xác suất	26
2.3	Hàm phụ thuộc đại lượng ngẫu nhiên	28
2.4	Các đặc trưng số của đại lượng ngẫu nhiên	31
2.5	Các phân bố xác suất thường gặp	37
3	Véc tơ ngẫu nhiên	45
3.1	Khái niệm	46
3.2	Hàm phân bố xác suất đồng thời	50
3.3	Xác định luật phân bố thành phần	52
3.4	Các đại lượng ngẫu nhiên độc lập	53
3.5	Phân bố có điều kiện	55
3.6	Tổng các đại lượng ngẫu nhiên	59
3.7	Momen tương quan và Hệ số tương quan	63
4	Các định lý giới hạn	70
5	Mẫu và phân bố mẫu	74
5.1	Mẫu ngẫu nhiên đơn giản	74
5.2	Các đặc trưng mẫu	75
5.3	Các phân bố thường gặp trong thống kê	80

5.4 Phân bố mẫu	82
6 Ước lượng tham số	84
7 Kiểm định giả thuyết thống kê	86
7.1 Khái niệm	86
7.2 Kiểm định giả thuyết về giá trị trung bình và xác suất	87
7.3 Tiêu chuẩn phù hợp χ^2	94
8 Tương quan và hồi quy	107
8.1 Hồi quy	107
8.2 Hồi quy tuyến tính	108
8.3 Dữ liệu lớn và học máy	112
1 Biến cố ngẫu nhiên và xác suất	122
2 Đại lượng ngẫu nhiên	125
3 Vectơ ngẫu nhiên	129
4 Các định lý giới hạn	136
5 Mẫu và phân bố mẫu	138
7 Kiểm định giả thuyết thống kê	140
8 Tương quan và hồi quy	144
Phụ lục	145
A Python	146
A.1 Thư viện, mô đun, phương thức	146

Chương 8

Tương quan và hồi quy

8.1	Hồi quy	107
8.2	Hồi quy tuyến tính	108
8.2.1	Hồi quy tuyến tính lý thuyết	108
8.2.2	Hệ số tương quan mẫu	109
8.2.3	Hồi quy tuyến tính thực nghiệm	110
8.3	Dữ liệu lớn và học máy	112

8.1 Hồi quy

Xét bài toán

$$E \left\{ [Y - \varphi(X)]^2 \right\} \rightarrow \min_{\varphi}.$$

Nghiệm của bài toán ký hiệu là $Y^* = \varphi^*(X)$ và giá trị min gọi là sai số bình phương trung bình.

Định nghĩa 8.1. Hàm hồi quy của Y theo X , là đại lượng ngẫu nhiên $Z = E(Y | X)$ tức là khi $X = x$ thì $Z = E(Y | X = x)$.

a) Xét (X, Y) rời rạc. Nếu $X = x_i$ thì $W = (Y | X = x_i)$ có luật phân bố

$$P(W = y_j) = \frac{p_{ij}}{p_{i*}} \quad \forall j.$$

Khi đó

$$Z = EW = \sum_j y_j \frac{p_{ij}}{p_{i*}} = \frac{1}{p_{i*}} \sum y_j p_{ij}.$$

b) Xét (X, Y) liên tục. Nếu $X = x$ thì $W = (Y | X = x)$ có hàm mật độ (công thức 3.15)

$$\psi(y | x) = \frac{f(x, y)}{f_1(x)} \quad \forall y.$$

Kỳ vọng của W , tức là kỳ vọng có điều kiện của Y khi biến cố $X = x$ đã xảy ra là

$$\begin{aligned} EW &= E(Y | X = x) = \int_{-\infty}^{\infty} y \psi(y | x) dy \\ &= \int_{-\infty}^{\infty} y \frac{f(x, y)}{f_1(x)} dy = \frac{1}{f_1(x)} \int_{-\infty}^{\infty} y f(x, y) dy. \end{aligned} \quad (8.1)$$

Định lý 8.1. $\varphi^*(X) = E(Y | X)$.

Ví dụ 8.1. (B)

Giải.

X	$Z = E(Y X)$
0	$E(Y X = 0) = \frac{(-1) \cdot 0.08 + 0 \cdot 0.16 + 2 \cdot 0.12 + 3 \cdot 0.04}{0.4} = 0.7$
1	$E(Y X = 1) = \frac{(-1) \cdot 0.05 + 0 \cdot 0.1 + 2 \cdot 0.08 + 3 \cdot 0.02}{0.25} = 0.68$
4	$E(Y X = 4) = \frac{(-1) \cdot 0.07 + 0 \cdot 0.14 + 2 \cdot 0.1 + 3 \cdot 0.04}{0.35} = 0.7143$

□

Ví dụ 8.2. (C)

Giải. Khi $X = x$ thì

$$Z = E(Y | X = x) = \frac{1}{f_1(x)} \int_{-\infty}^{\infty} y f(x, y) dy = \begin{cases} \frac{10x}{9}, & 0 \leq x < 1 \\ 0, & x \notin [0, 1). \end{cases}$$

□

8.2 Hồi quy tuyến tính

8.2.1 Hồi quy tuyến tính lý thuyết

Bài toán

$$E\{[Y - aX - b]^2\} \rightarrow \min_{a, b}$$

có nghiệm

$$\begin{aligned} a &= \rho \sqrt{\frac{DY}{DX}} \\ b &= EY - aEX. \end{aligned} \quad (8.2)$$

và sai số $\sigma_{Y/X}^2 = (1 - \rho^2) DY$.

Hàm $\varphi(X) = aX + b$ gọi là hàm hồi quy tuyến tính lý thuyết của Y theo X . Đường thẳng $y = ax + b$ gọi là đường hồi quy tuyến tính lý thuyết.

Ví dụ 8.3. (B) $EX = 1.65, EY = 0.7, DX = 3.1275, DY = 1.81, \rho = 0.006305$.

Giải.

$$\begin{aligned} a &= 0.006305 \sqrt{\frac{1.81}{3.1275}} = 0.04797, b = 0.7 - 0.04797 \cdot 1.65 = 0.6208 \\ \sigma_{Y/X}^2 &= (1 - 0.006305^2) 1.81 = 1.8099. \end{aligned}$$

□

Ví dụ 8.4. (C) $EX = \frac{3}{4}, EY = \frac{5}{6}, DX = \frac{3}{80}, DY = \frac{43}{180}, \rho = \frac{5}{\sqrt{129}}$.

Giải.

$$\begin{aligned} a &= \frac{5}{\sqrt{129}} \sqrt{\frac{43/180}{3/80}} = \frac{10}{9}, b = \frac{5}{6} - \frac{10}{9} \cdot \frac{3}{4} = 0 \\ \sigma_{Y/X}^2 &= \left[1 - \left(\frac{5}{\sqrt{129}} \right)^2 \right] \frac{3}{80} = \frac{13}{430} = 0.03023. \end{aligned}$$

□

8.2.2 Hệ số tương quan mẫu

Nhắc lại công thức (3.21) về hệ số tương quan giữa hai đại lượng ngẫu nhiên X và Y

$$\rho = \rho_{XY} = \frac{E(XY) - EX \cdot EY}{\sqrt{DX} \sqrt{DY}}.$$

Giả sử $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ là mẫu ngẫu nhiên từ véc tơ ngẫu nhiên (X, Y) . Hệ số tương quan mẫu giữa X và Y là đại lượng ước lượng cho hệ số tương quan, xác định bởi

$$r = \frac{\overline{XY} - \bar{X} \cdot \bar{Y}}{S_X S_Y} \quad (8.3)$$

trong đó $\overline{XY} = \frac{1}{n} \sum_{i=1}^n X_i Y_i$.

Nếu mẫu cụ thể được rút gọn dưới dạng bảng với (x_i, y_j) lặp lại n_{ij} lần thì $n = \sum_{i,j} n_{ij}$ và

$$\overline{xy} = \frac{1}{n} \sum_{i,j} n_{ij} x_i y_j, \quad \bar{x} = \sum_i n_{i*} x_i, \quad \bar{y} = \sum_j n_{*j} y_j.$$

8.2.3 Hồi quy tuyến tính thực nghiệm

Cho mẫu ngẫu nhiên $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ từ vectơ ngẫu nhiên (X, Y) . Bài toán

$$\frac{1}{n} \sum_{i=1}^n (Y_i - aX_i - b)^2 \rightarrow \min_{a,b}$$

có nghiệm

$$\begin{aligned} a &= r \frac{S_Y}{S_X} \\ b &= \bar{Y} - a\bar{X}. \end{aligned} \quad (8.4)$$

và sai số $S_{Y/X}^2 = (1 - r^2) S_Y^2$.

Hàm $\tilde{\varphi}(X) = aX + b$ gọi là hàm hồi quy tuyến tính thực nghiệm của Y theo X , đường thẳng $y = ax + b$ gọi là đường hồi quy tuyến tính thực nghiệm.

Ví dụ 8.5. Cho mẫu từ vectơ ngẫu nhiên (X, Y)

X	1	1.5	3	4.5	5
Y	1.25	1.4	1.5	1.75	2.25

- Mô tả dữ liệu bởi các điểm trên mặt phẳng.
- Tính hệ số tương quan mẫu.
- Tìm hàm hồi quy tuyến tính thực nghiệm của Y theo X . Vẽ đường hồi quy tuyến tính thực nghiệm cùng với các điểm ở ý (a).

Giải. b) $\overline{xy} = \frac{1 \cdot 1.25 + 1.5 \cdot 1.4 + 3 \cdot 1.5 + 4.5 \cdot 1.75 + 5 \cdot 2.25}{5} = 5.395$.

$$\bar{x} = 3, \bar{y} = 1.63, s_X^2 = 2.5, s_Y^2 = 0.1226 \Rightarrow r = \frac{5.395 - 3 \cdot 1.63}{\sqrt{2.5 \cdot 0.1226}} = 0.9122.$$

Cách 1:

$$\begin{aligned} 1 \quad X &= [1, 1.5, 3, 4.5, 5] \\ 2 \quad Y &= [1.25, 1.4, 1.5, 1.75, 2.25] \end{aligned}$$

```

4 from scipy.stats import pearsonr
5 r, _ = pearsonr(X, Y) # r là thành phần thứ nhất trong hai
                        thành phần của kết quả

```

Cách 2:

```

1 import numpy as np
2 np.corrcoef([X, Y]) # ma trận hệ số tương quan mẫu  $(r_{ij})_n$ , trong
                      đó  $r_{ij} = r_{X_j X_i}$ 
3 r = _[0, 1]

```

Cách 3:

```

1 import pandas as pd
2 df = pd.DataFrame({'X': X, 'Y': Y})

4 r = df['X'].corr(df['Y']) # hoặc gọi từ ma trận hệ số tương
                           quan bằng lệnh df.corr()['X']['Y'], df.corr().at['X', 'Y'] hoặc
                           df.corr().iat[0, 1]

```

$$c) a = 0.9122 \sqrt{\frac{0.1226}{2.5}} = 0.2020, b = 1.63 - 0.202 \cdot 3 = 1.0240.$$

Hàm hồi quy tuyến tính thực nghiệm của Y theo X là $\tilde{\varphi}(X) = 0.2020X + 1.0240$.

Khi $X = 2$, ta dự báo $Y = 0.2020 \cdot 2 + 1.0240 = 1.428$.

Sai số $s_{Y/X}^2 = (1 - 0.9122^2) 0.1226 = 0.02059$.

```

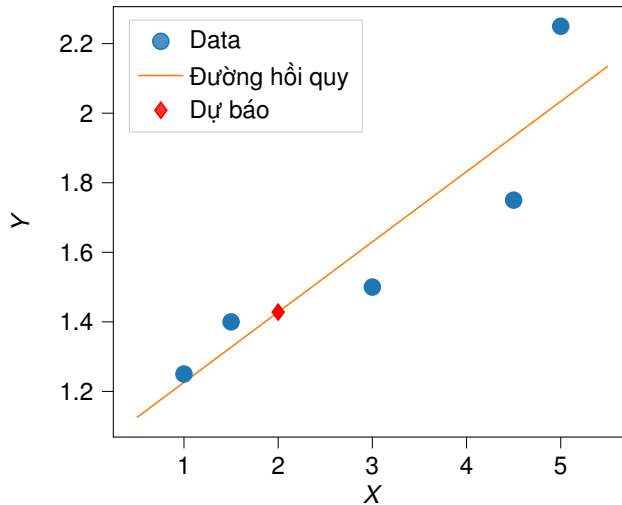
1 from scipy.stats import linregress
2 a, b, r, _, _ = linregress(X, Y) # kết quả có 5 thành phần

4 import matplotlib.pyplot as plt

6 plt.plot(X, Y, 'o', label='Data')
7 plt.plot([0.5, 5.5], [a*0.5 + b, a*5.5 + b], label='Đường hồi
    quy')

8 plt.plot([2], [a*2 + b], 'rD', label='Dự báo')
9 plt.legend()
10 plt.xlabel('X')
11 plt.ylabel('Y')

```



□

8.3 Dữ liệu lớn và học máy

Ví dụ 8.6. Giải bài toán Ví dụ 8.5(c) bằng

- d) Mô hình hồi quy trong học máy, với thư viện `scikit-learn`.
- e) Mạng thần kinh (học sâu), với thư viện `TensorFlow`

Giải. d) Mô hình dự báo

```
1 X = [[1], [1.5], [3], [4.5], [5]]
2 y = [1.25, 1.4, 1.5, 1.75, 2.25]
3 from sklearn.linear_model import LinearRegression

5 model = LinearRegression()
6 model.fit(X, y)

8 model.coef_          # a
9 model.intercept_     # b

11 model.predict([[2]]) # dự báo
```

và đánh giá mô hình

```
1 y_pred = model.predict(X)

3 from sklearn.metrics import mean_squared_error, r2_score
```



```

5 mean_squared_error(y, y_pred) #  $s_{y/x}^2$ 
6 r2_score(y, y_pred)           #  $0 \leq R^2 \leq 1$ , càng gần 1 nghĩa là mô
    hình càng hiệu quả

```

e)

```

1 import tensorflow as tf

3 X = tf.constant([[1], [1.5], [3], [4.5], [5]])
4 y = tf.constant([1.25, 1.4, 1.5, 1.75, 2.25])

6 # Bước 1: xây dựng mô hình
7 model = tf.keras.Sequential([
8     tf.keras.layers.InputLayer(input_shape=(1,)),
9     tf.keras.layers.Dense(100, activation='relu'),
10    tf.keras.layers.Dense(100, activation='relu'),
11    tf.keras.layers.Dense(100, activation='relu'),
12    tf.keras.layers.Dense(1, activation=None),
13 ])

15 # Bước 2: chọn phương pháp tối ưu mô hình
16 model.compile(
17     loss=tf.keras.losses.mse,
18     optimizer=tf.keras.optimizers.SGD()
19 )

21 # Bước 3: tối ưu mô hình với số liệu
22 model.fit(X, y, epochs=100) # epochs là số bước lặp

24 model.predict([2]) # khi X=2, dự báo Y=1.3260

```

Sai số của mô hình tại các bước lặp 1, 10, 20,...

Bước	Sai số	Bước	Sai số
1	2.9204	60	0.0522
10	0.1415	70	0.0445
20	0.1121	80	0.0387
30	0.0909	90	0.0341
40	0.0745	100	0.0306
50	0.0619		

□

Khi xây dựng mô hình dự báo với tập dữ liệu lớn về (X, y) , gồm $(X^1, y_1), (X^2, y_2), \dots, (X^N, y_N)$.

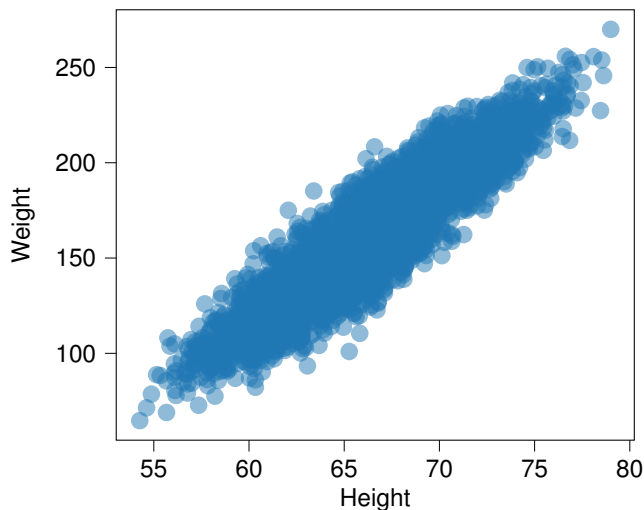
- 1) Chia dữ liệu thành hai phần $X = X_{train} \cup X_{test}$, $y = y_{train} \cup y_{test}$, tương ứng theo các cặp. Thông thường, tỷ lệ dữ liệu học $75\% \leq \frac{|X_{train}|}{|X|} = \frac{|y_{train}|}{|y|} \leq 80\%$.
- 2) Tối ưu mô hình với số liệu X_{train} và y_{train} .
- 3) Dự báo các giá trị của y khi X nhận giá trị trong X_{test} , lưu vào dãy y_{pred} .
- 4) Đánh giá độ hiệu quả của mô hình bằng cách so sánh độ lệch giữa giá trị thực y_{test} và giá trị dự báo y_{pred} .

Ví dụ 8.7. Tiếp tục với số liệu về giới tính, chiều cao, cân nặng trong **Ví dụ 5.2**.

- a) Mô tả dữ liệu về chiều cao – cân nặng bằng biểu đồ điểm.
- b) Dự đoán chiều cao theo cân nặng, áp dụng mô hình hồi quy tuyến tính của học máy. Mô tả mô hình bằng hình vẽ.
- c) Thực hiện yêu cầu của ý (b) theo mô hình học sâu.

Giải. a)

```
1 import pandas as pd
2 df = pd.read_csv('https://raw.githubusercontent.com/
    Dataweekends/zero_to_deep_learning_video/master/data/
    weight-height.csv')
4 df.plot(kind='scatter', x='Height', y='Weight', alpha=0.5)
```



```

1 X, y = df[['Height']], df['Weight']
2 # Chia dữ liệu thành hai phần: dữ liệu học và dữ liệu thử
3 from sklearn.model_selection import train_test_split
4 X_train, X_test, y_train, y_test = train_test_split(X, y)

```

b)

```

1 from sklearn.linear_model import LinearRegression
2 model = LinearRegression()
3 model.fit(X_train, y_train)

5 y_pred = model.predict(X_test)

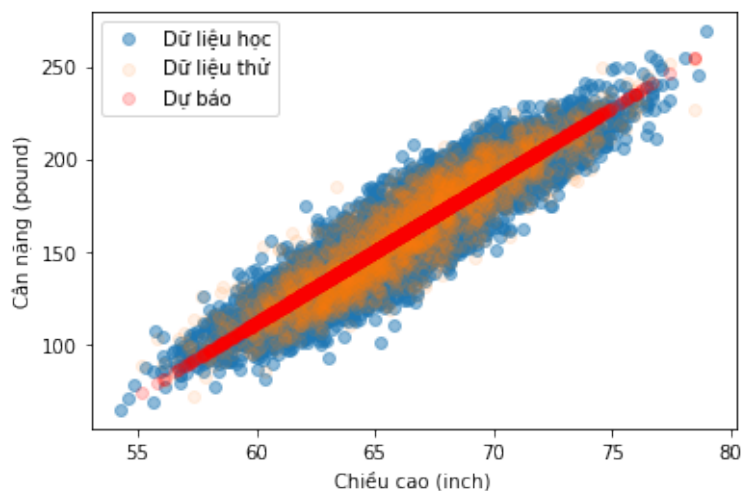
7 from sklearn.metrics import mean_squared_error, r2_score
8 mean_squared_error(y_test, y_pred) # 149.70
9 r2_score(y_test, y_pred) # 0.8515

11 mean_squared_error(y_train, model.predict(X_train)) # 149.17
12 r2_score(y_train, model.predict(X_train)) # 0.8563

1 import matplotlib.pyplot as plt
2 plt.plot(X_train, y_train, 'o', alpha=0.5, label='Dữ liệu
   học')

3 plt.plot(X_test, y_test, 'o', alpha=0.1, label='Dữ liệu thử')
4 plt.plot(X_test, y_pred, 'or', alpha=0.2, label='Dự báo')
5 plt.legend()
6 plt.xlabel('Chiều cao (inch)')
7 plt.ylabel('Cân nặng (pound)')

```



Hình 8.1: Mô hình hồi quy tuyến tính của cân nặng theo chiều cao

c)

```

1 import tensorflow as tf

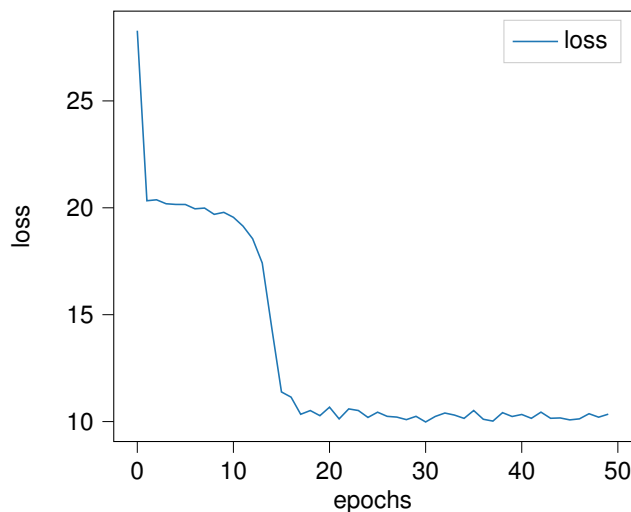
3 model = tf.keras.Sequential([
4     tf.keras.layers.InputLayer(input_shape=(1,)),
5     tf.keras.layers.Dense(100, activation='relu'),
6     tf.keras.layers.Dense(100, activation='relu'),
7     tf.keras.layers.Dense(100, activation='relu'),
8     tf.keras.layers.Dense(1, activation=None),
9 ])

11 model.compile(
12     loss=tf.keras.losses.mae,
13     optimizer=tf.keras.optimizers.Adam(),
14 )

16 history = model.fit(X_train, y_train, epochs=50)

18 pd.DataFrame(history.history).plot(xlabel='epochs', ylabel='
    loss') # biểu đồ sai số sau mỗi bước
    lặp

```



Sai số của mô hình học sâu cũng xấp xỉ các sai số ở ý (b)

```

1 mean_squared_error(y_test, y_pred) # 149.81
2 r2_score(y_test, y_pred)           # 0.8545

```

và đồng thời, đồ thị mô tả giống hệt [Hình 8.1](#).

□

Một ưu điểm của thuật toán học sâu với thư viện TensorFlow, chẳng hạn, bài toán phân loại được cấu hình khá tương tự bài toán hồi quy.

Ví dụ 8.8. Nối tiếp ví dụ [Ví dụ 8.7](#)

- d) Mô tả số liệu chiều cao – cân nặng theo giới tính trên cùng biểu đồ điểm.
- e) Dự đoán những người có chiều cao 65, 66, 67, 68, 69, và 70 inch là nam hay nữ.
- f) Cùng có chiều cao 70 inch, dự đoán người có cân nặng 160, 170, và 180 pound là nam hay nữ.

Giải. d)

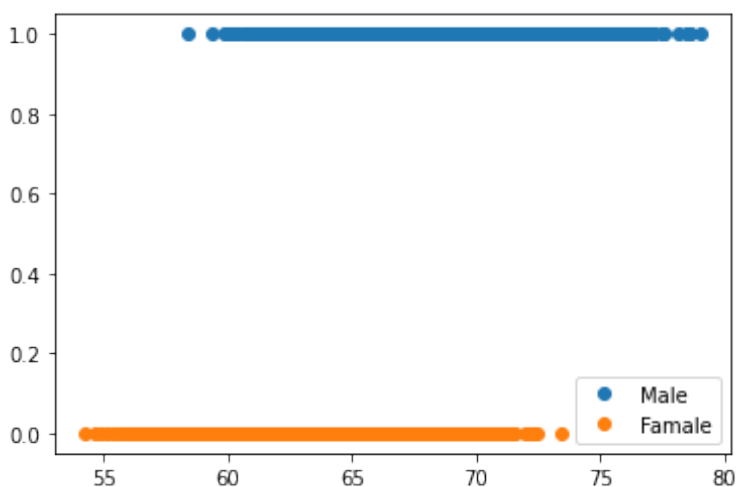
```

1 X = df[['Height']]
3 y = [0]*len(X)
4 for i in range(len(X)):
5     if df['Gender'][i] == 'Male':
6         y[i] = 1
8 y = tf.constant(y)

10 X_Male = df[ df['Gender'] == 'Male' ]['Height']
11 plt.plot(X_Male, [1]*len(X_Male), 'o', label='Male')

13 X_Female = df[ df['Gender'] == 'Female' ]['Height']
14 plt.plot(X_Female, [0]*len(X_Female), 'o', label='Famale')
15 plt.legend();

```



e)

```

1 model = tf.keras.Sequential([
2     tf.keras.layers.Dense(100, activation='relu'),
3     tf.keras.layers.Dense(100, activation='relu'),
4     tf.keras.layers.Dense(100, activation='relu'),
5     tf.keras.layers.Dense(1, activation='sigmoid'),
6 ])

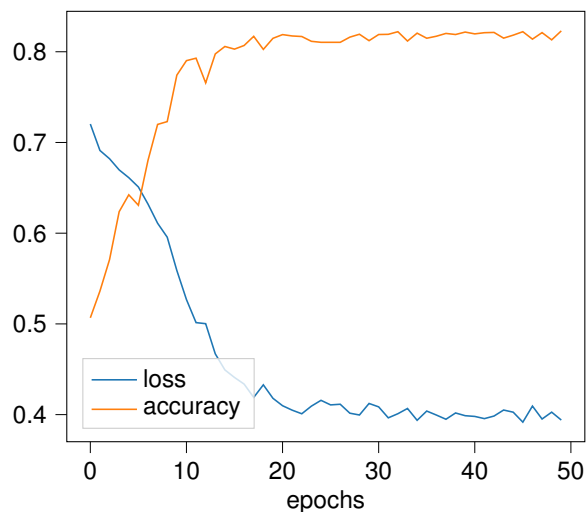
8 model.compile(
9     loss=tf.keras.losses.BinaryCrossentropy(),
10    optimizer=tf.keras.optimizers.Adam(),
11    metrics=['accuracy']
12 )

14 history = model.fit(X, y, epochs=50, verbose=0)

16 pd.DataFrame(history.history).plot()

18 model.predict([62, 64, 66, 68, 70])

```



Chiều cao	Xác suất dự đoán	Kết luận
62	0.08283	Nữ
64	0.2277	Nữ
66	0.5038	Nam
68	0.7776	Nam
70	0.9233	Nam

f)

```

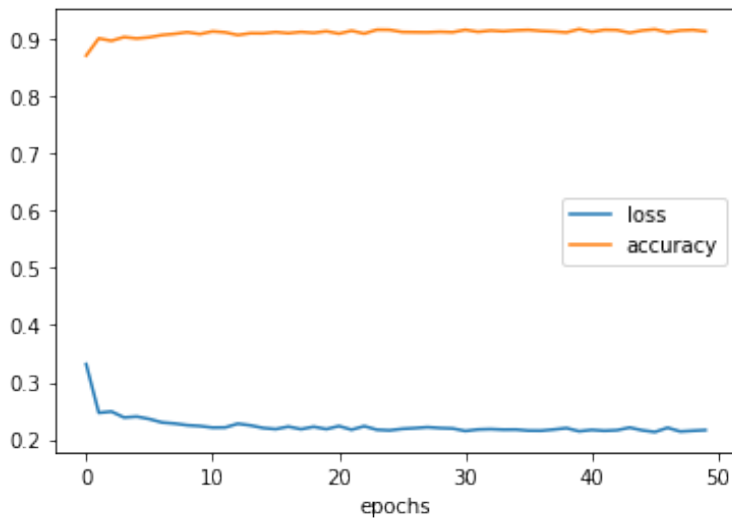
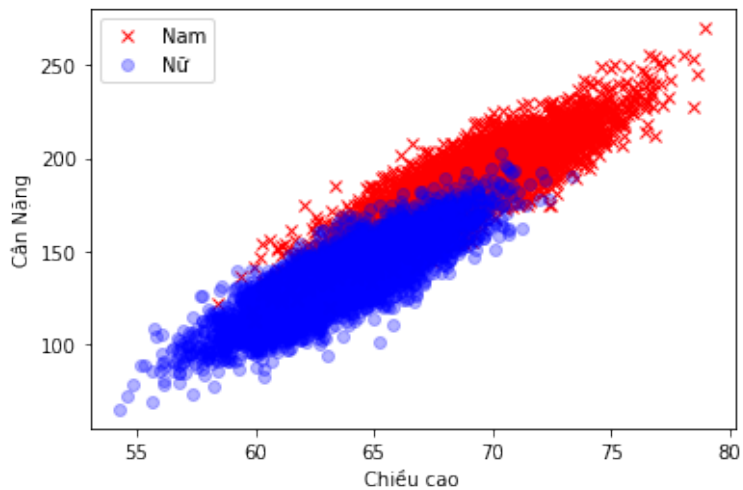
1 X = np.array( pd.DataFrame({'Height': df['Height'], 'Weight'
2     : df['Weight']}) )
2 plt.plot(X[y==1, 0], X[y==1, 1], 'xr', label='Nam')

```

```

3 plt.plot(X[y==0, 0], X[y==0, 1], 'ob', alpha=0.3, label='Nữ'
4 )
5 plt.legend()
6 plt.xlabel('Chiều cao')
7 plt.ylabel('Cân Nặng')
8 # → Giữ nguyên phần cài đặt, đánh giá mô hình như ý (e)
9
10 model.predict([[70, 160], [70, 170], [70, 180]])

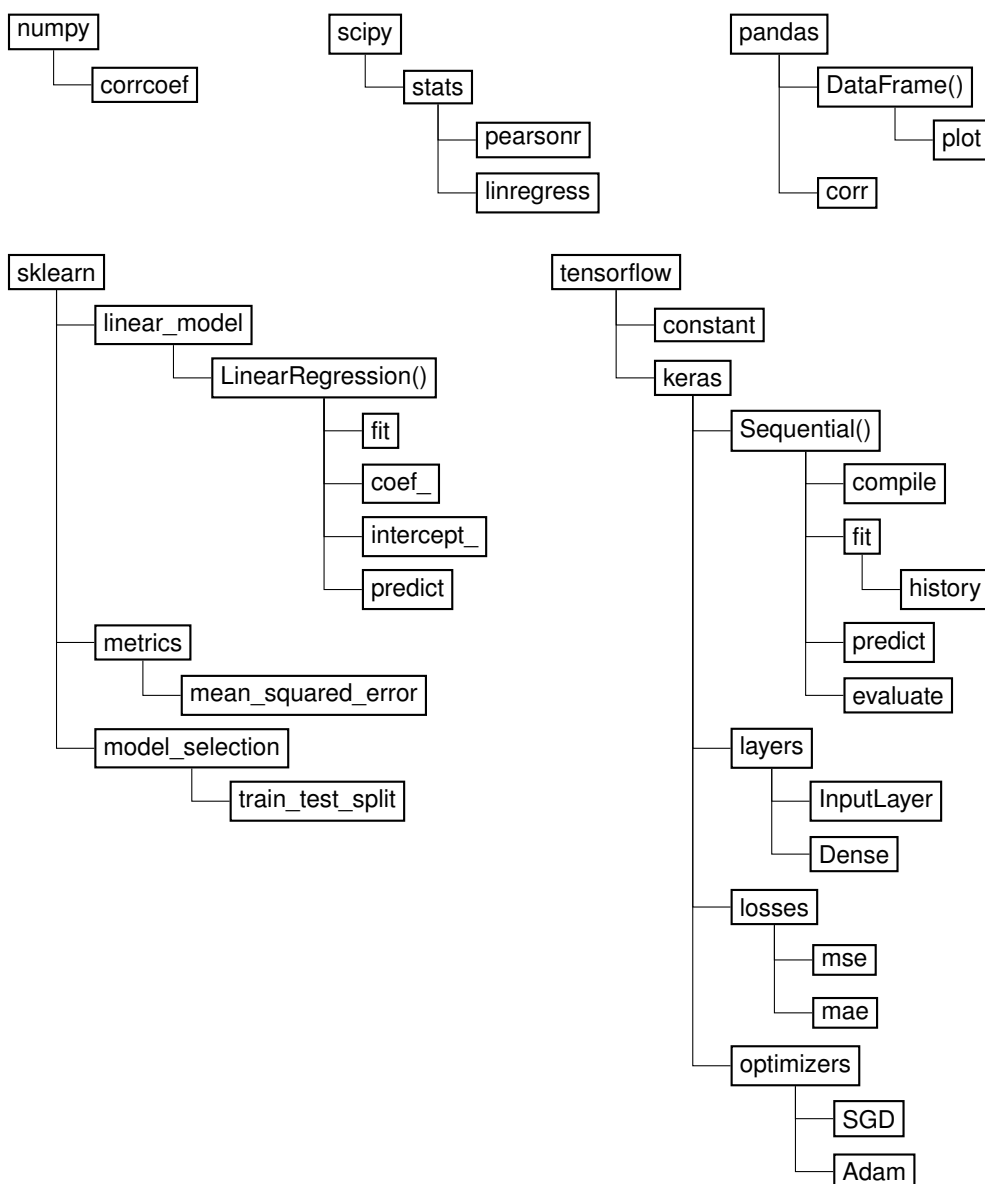
```



Chiều cao	Cân nặng	Xác suất dự báo	Kết luận
70	160	0.08519	Nữ
70	170	0.4667	Nữ
70	180	0.8859	Nam

□

Tóm tắt về Python



Bài tập Chương 8

Với mẫu từ (X, Y)

- Tính hệ số tương quan mẫu giữa X và Y .
- Tìm hồi quy bình phương trung bình tuyến tính thực nghiệm của Y theo X . Ước lượng sai số.
- Viết phương trình đường thẳng hồi quy bình phương trung bình thực nghiệm của Y

đối với X và vẽ nó.

d) Dự báo giá trị của Y khi X nhận giá trị cho trước.

8.1.

x_i	5	3	4	6	5	5.5	6	7
y_i	2.5	1.5	3	7	6.5	5.5	5	11

Dự báo Y khi $X = 3.5$.

8.2.

$y \backslash X$	7	8	9
200	41	7	
300	1	51	1
400		8	40

Dự báo Y khi $X = 8.5$.

Chương 8

Tương quan và hồi quy

8.1 a) $r = 0.86$

b) $\tilde{\varphi} = 2.0855X - 5.5684, s_Y^2(1 - r^2) = 2.0994.$

c) $y = 2.0855x - 5.5684$

d) Khi $X = 3.5$, ta dự báo $Y = 1.7308.$

8.2 a) $r = 0.9075$

b) $\tilde{\varphi} = 97.5982X - 480.1350, s_Y^2(1 - r^2) = 1129.6972.$

c) $y = 97.5982x - 480.135.$

d) Khi $X = 8.5$, ta dự báo $Y = 349.4498.$

