# Ứng dụng Trí tuệ nhân tạo trong Nuôi trồng thủy sản

#### NGUYỄN HẢI TRIỀU<sup>1</sup>

 $^1$  Bộ môn Kỹ thuật phần mềm, Khoa Công nghệ thông tin, Trường ĐH Nha Trang

NhaTrang, September 2024

- Using Logistic Regression for Classification
  - Single Layer Neural Networks
  - Activation Function
  - Logistic regression loss function
  - Model Training with Stochastic Gradient Descent
  - Automatic Differentiation in PyTorch
  - The PyTorch API
  - Training a Logistic Regression Model in PyTorch
  - Feature Normalization

## Training a Logistic Regression Model I

Trong phần này, chúng ta sẽ tìm hiểu cách huấn luyện mô hình hồi quy logistic trong PyTorch. Chúng ta sẽ sử dụng lại code và bộ dữ liệu từ chương 2 (chi tiết về code xem file Chuong\_3.3\_code)

```
import pandas as pd
    import numpy as np
    import matplotlib.pyplot as plt
    df = pd.read_csv("perceptron_toydata-truncated.txt", sep="\t")
    X_{\text{train}} = df[["x1", "x2"]].values
    y_train = df["label"].values
    np.bincount(y_train) # kiểm tra phân bố số lượng nhãn của lớp 0, 1.
    # vẽ phân bố các điểm dữ liệu đầu vào với 2 đặc trưng
10
    plt.plot(
     X train[v train == 0, 0].
11
     X_train[y_train == 0, 1],
12
     marker="D",
13
```

Trieu Hai Nguyen Image Processing BM. KTPM-Khoa CNTT 3/35

## Training a Logistic Regression Model II

```
markersize=10,
14
      linestyle="",
15
      label="Class 0",
16
17
18
    plt.plot(
19
      X_train[y_train == 1, 0],
20
      X_{train}[y_{train} == 1, 1],
      marker="^",
      markersize=13.
      linestyle="",
24
      label="Class 1".
25
26
27
28
    plt.legend(loc=2)
29
    plt.xlim([-5, 5])
30
    plt.ylim([-5, 5])
31
32
    plt.xlabel("Feature $x_1$", fontsize=12)
33
```

## Training a Logistic Regression Model III

```
34 plt.ylabel("Feature $x_2$", fontsize=12)
35
36 plt.grid()
37 plt.show()
38
39 # chuẩn hoá dữ liệu đầu vào
40 X train = (X train - X train.mean(axis=0)) / X train.std(axis=0)
```

### Chuẩn hoá dữ liệu

We can now see that the data set features are clustered around zero or they are centered at zero. And also the spread is a little bit smaller here. This is an important step to make the logistic regression training smoother.

## Cài đặt mô hình hồi quy logistic trong PyTorch

Đầu tiên chúng ta tạo lớp LogisticRegression kế thừa từ torch.nn.Module, sử dụng lớp Linear thay cho cho việc khởi tạo trọng số như mạng Perceptron. Tiếp theo định nghĩa phương thước forward để tính giá trị đầu ra của lớp Linear và hàm kích hoạt sigmoid.

```
import torch

class LogisticRegression(torch.nn.Module):

def __init__(self, num_features):
    super().__init__()
    self.linear = torch.nn.Linear(num_features, 1)

def forward(self, x):
    logits = self.linear(x)
    probas = torch.sigmoid(logits)
    return probas
```

### Khởi tạo một ví dụ và xem kết quả dự đoán

Khởi tạo mô hình Logistic regression với một mẫu huấn luyện có 2 đặc trưng x=torch.tensor([1.1,2.1])

```
torch.manual_seed(1)
model = LogisticRegression(num_features=2)
x = torch.tensor([1.1, 2.1])
with torch.no_grad():
proba = model(x)
print(proba)
```

Trong đó, with torch.no\_grad() có nghĩa là quá trình dự đoán sẽ không liên quan đến sơ đồ tính toán gradient mà PyTorch thực hiện. Điều này giúp tiết kiệm bộ nhớ và giúp cho mô hình chạy dự đoán nhanh hơn.

#### Data loader

Bỏ qua ví dụ chỉ có 1 mẫu huấn luyện ở trước, ta đi đến phần chuẩn bị *Data loader* là dữ liệu dùng để huấn luyện mô hình.

```
from torch.utils.data import Dataset, DataLoader
    class MyDataset(Dataset):
     def __init__(self, X, y):
       self.features = torch.tensor(X, dtype=torch.float32)
       self.labels = torch.tensor(y, dtype=torch.float32)
     def __getitem__(self, index):
       x = self.features[index]
       y = self.labels[index]
       return x, y
     def len (self):
       return self.labels.shape[0]
15
    train_ds = MyDataset(X_train, y_train)
16
17
```

```
train_loader = DataLoader(
dataset=train_ds,
batch_size=10,
shuffle=True,

22 )
```

#### MyDataset

Chúng ta sẽ định nghĩa lớp dữ liệu huấn của mình MyDataset dựa trên 2 lớp Dataset, DataLoader từ thư viện torch.utils.data. Với đầu vào là tensor đặc trưng  $X\_train$  và tensor label  $y\_train$  tương ứng, ta định nghĩa phương thức  $\_\_getitem\_\_$  để lấy từng mẫu dữ liệu huấn luyện và nhãn thông qua index. Phương thức  $\_\_len\_\_$  để lấy số lượng mẫu có trong bộ dữ liệu huấn luyên.

#### DataLoader

Để sử dung thuật toán huấn luyên minibatch stochastic gradient descent, bô dữ liêu huấn luyên cần chia thành các batches, ở đây chúng ta sử dụng lớp *Dataloader* với thuộc tính *batch* size của PyTorch. Bên cạnh đó ta cũng xáo trộn bộ dữ liệu shuffle=True trước mỗi epoch để đảm bảo quá trình học được nhanh hơn.

## Huấn luyện mô hình I

Quá trình huấn luyện mô hình Logistic regression, ta khởi tạo mô hình với lớp  $LogisticRegression(num\_features=2)$  và thuật toán huấn luyện cập nhật tham số mô hình là torch.optim.SGD(model.parameters(), lr=0.05). Chúng ta sẽ huấn luyện trong vòng 20 epochs và learning rate = 0.05. Các tham số này có thể thay đổi tuỳ thực nghiệm.

```
import torch.nn.functional as F

torch.manual_seed(1)

model = LogisticRegression(num_features=2)

optimizer = torch.optim.SGD(model.parameters(), lr=0.05)

num_epochs = 20

for epoch in range(num_epochs):

model = model.train()
```

# Huấn luyện mô hình II

```
for batch_idx, (features, class_labels) in enumerate(train_loader):
11
       probas = model(features)
12
       loss = F.binary_cross_entropy(probas, class_labels.view(probas.shape))
14
       optimizer.zero_grad()
16
       loss.backward()
17
       optimizer.step()
18
19
       ### LOGGING
20
       print(f'Epoch: {epoch+1:03d}/{num_epochs:03d}')
21
           f' | Batch {batch_idx:03d}/{len(train_loader):03d}'
           f' | Loss: {loss:.2f}')
23
```

## Huấn luyện mô hình III

Khi lặp qua mỗi epoch, nên thiết lập chế độ model.train().

Tiếp theo, for batch\_idx, (features, class\_labels) in enumerate(train\_loader): sẽ giúp chúng ta duyệt qua các mini batch trong train\_loader, ứng với mỗi mini batch ta có các số lượng mẫu huấn luyện tương ứng (features, class\_labels). Trong phần DataLoader trước, ta chọn  $batch\_size = 10$ , tức là mỗi mini batch sẽ có 10 dòng dữ liệu mẫu huấn luyện. Với bộ dữ liệu có 20 dữ liệu huấn luyện, PyTorch sẽ chia ra được 2 mini batches,  $batch\_idx = \{0, 1\}$ . Dĩ nhiên,  $batch\_size$  là siêu tham số, chúng ta có thể điều chỉnh dựa vào learning experience.

## Huấn luyện mô hình IV

Ö vòng lặp bên trong,

```
probas = model(features) loss = F.binary_cross_entropy(probas, class_labels.view(probas.shape)) chúng ta thu được class membership probabilities probas cho các training example trong mini batch. Khi tính loss ta thu được một giá trị tương ứng tương ứng cho một mini batch. Lưu \acute{y}: để tính được loss thì true label class (hiện tại torch.Size([10])) phải cùng kích thước với probas (torch.Size([10,1])), do đó ta gọi class\_labels.view(probas.shape) để quy về cùng 1 kích thước tensor.
```

## Huấn luyện mô hình V

Tiếp theo, chúng ta tiến hành cập nhật trọng số bằng thuật toán minibatch stochastic gradient descent:

```
optimizer.zero_grad() # tránh cộng dồn gradient ởbatch trước loss.backward() # tính đạo hàm riêng theo các các tham số mô hình optimizer.step() # cập nhật trọng số đểcực tiểu loss.
```

Trieu Hai Nguyen Image Processing BM. KTPM-Khoa CNTT

## Huấn luyện mô hình VI

## Kết quả huấn luyện

Dựa vào kết quả, ta thấy rằng hàm Loss giảm nhỏ dần so với ban đầu và đủ nhỏ để nhận thấy rằng quá trình huấn luyện mô hình Logistic regression có hiệu quả.

```
Epoch: 001/020 | Batch 000/002 | Loss: 0.67
Epoch: 001/020 | Batch 001/002 | Loss: 0.73
Epoch: 002/020 | Batch 000/002 | Loss: 0.67
Epoch: 002/020 | Batch 001/002 | Loss: 0.67
Epoch: 003/020 | Batch 001/002 | Loss: 0.60
...
Epoch: 019/020 | Batch 000/002 | Loss: 0.34
Epoch: 019/020 | Batch 001/002 | Loss: 0.39
Epoch: 020/020 | Batch 000/002 | Loss: 0.33
Epoch: 020/020 | Batch 001/002 | Loss: 0.38
```

Trieu Hai Nguyen Image Processing BM. KTPM-Khoa CNTT 16/35

# Ước lượng hiệu quả của mô hình I

Để đánh giá mô hình sau khi huấn luyện có dự đoán tốt không, chúng ta có thể sử hàm " *compute\_accuracy*"

```
def compute_accuracy(model, dataloader):
       model = model.eval()
       correct = 0.0
       total examples = 0
       for idx, (features, class_labels) in enumerate(dataloader):
           with torch.no_grad():
               probas = model(features)
           pred = torch.where(probas > 0.5, 1, 0)
11
           lab = class_labels.view(pred.shape).to(pred.dtype)
14
           compare = lab == pred
           correct += torch.sum(compare)
```

## Ước lượng hiệu quả của mô hình II

```
total_examples += len(compare)
return correct / total_examples
train_acc = compute_accuracy(model, train_loader)
print(f"Accuracy: {train_acc*100}%")
```

#### evaluation mode

Trước khi tính đánh giá hiệu quả huấn luyện mô hình, chúng ta cần đặt model vào evaluation mode: model.eval(). Điều này giúp mô hình tách rời với sơ đồ tính toán gradient và sẽ thực hiện tính toán khác so với huấn luyện ở các lớp BatchNorm, BatchNorm . . . được nhắc ở chương sau.

## Ước lượng hiệu quả của mô hình III

## torch.where(probas > 0.5, 1, 0)

Nhắc lại rằng probas hiện tại đang ở dạng xác suất thành phần của class, chúng ta cần hàm threshold để phân lớp dữ liệu mẫu thuộc 0 hoặc 1 dựa vào ngưỡng 0.5. Vậy nên chúng ta sử dụng torch.where(probas > 0.5,1,0).

#### compare = lab == pred

Sử dụng compare chúng ta có thể tính được số lượng nhãn dữ đoán đúng với nhãn thật sự của các điểm dữ liệu trong mini batch bằng cách correct += torch.sum(compare).

Ở phần huấn luyện mạng Logistic regression, chúng ta đã thực hiện chuẩn hoá dữ liệu đâu vào. *Câu hỏi đặt ra tại sao cần phải chuẩn hoá?* Hãy xem bộ dữ liệu thực tế Wine dataset với mong muốn dự đoán chất lượng của rượu từ 13 thông số đặc trưng khác nhau.

	Alcohol	Malic acid	Ash	Alcalinity of ash	Magnesium	Total phenols	Flavanoids	Nonflavanoid phenols	Proanthocyanins	Color intensity	Hue	OD280/OD315 of diluted wines	Proline
0	14.23	1.71	2.43	15.6	127	2.80	3.06	0.28	2.29	5.64	1.04	3.92	1065
1	13.20	1.78	2.14	11.2	100	2.65	2.76	0.26	1.28	4.38	1.05	3.40	1050
2	13.16	2.36	2.67	18.6	101	2.80	3.24	0.30	2.81	5.68	1.03	3.17	1185
3	14.37	1.95	2.50	16.8	113	3.85	3.49	0.24	2.18	7.80	0.86	3.45	1480
4	13.24	2.59	2.87	21.0	118	2.80	2.69	0.39	1.82	4.32	1.04	2.93	735
173	13.71	5.65	2.45	20.5	95	1.68	0.61	0.52	1.06	7.70	0.64	1.74	740
174	13.40	3.91	2.48	23.0	102	1.80	0.75	0.43	1.41	7.30	0.70	1.56	750
175	13.27	4.28	2.26	20.0	120	1.59	0.69	0.43	1.35	10.20	0.59	1.56	835
176	13.17	2.59	2.37	20.0	120	1.65	0.68	0.53	1.46	9.30	0.60	1.62	840
177	14.13	4.10	2.74	24.5	96	2.05	0.76	0.56	1.35	9.20	0.61	1.60	560
178 r	ows × 13	column	s										

#### Đễ dàng quan sát được rằng các đặc trung chênh lệch rất lớn.

Diều này ảnh hương tiêu cực đến thuật toán học SGD

	Alcohol	Malic acid	Ash	Alcalinity of ash	Magnesium	Total phenols	Flavanoids	Nonflavanoid phenols	Proanthocyanins	Color intensity	Hue	of diluted wines	Proline
0	14.23	1.71	2.43	15.6	127	2.80	3.06	0.28	2.29	5.64	1.04	3.92	1065
1	13.20	1.78	2.14	11.2	100	2.65	2.76	0.26	1.28	4.38	1.05	3.40	1050
2	13.16	2.36	2.67	18.6	101	2.80	3.24	0.30	2.81	5.68	1.03	3.17	1185
3	14.37	1.95	2.50	16.8	113	3.85	3.49	0.24	2.18	7.80	0.86	3.45	1480
4	13.24	2.59	2.87	21.0	118	2.80	2.69	0.39	1.82	4.32	1.04	2.93	735
173	13.71	5.65	2.45	20.5	95	1.68	0.61	0.52	1.06	7.70	0.64	1.74	740
174	13.40	3.91	2.48	23.0	102	1.80	0.75	0.43	1.41	7.30	0.70	1.56	750
175	13.27	4.28	2.26	20.0	120	1.59	0.69	0.43	1.35	10.20	0.59	1.56	835
176	13.17	2.59	2.37	20.0	120	1.65	0.68	0.53	1.46	9.30	0.60	1.62	840
177	14.13	4.10	2.74	24.5	96	2.05	0.76	0.56	1.35	9.20	0.61	1.60	560
178 r	ows × 13	column	s										

Trieu Hai Nguyen Image Processing BM. KTPM-Khoa CNTT 21/35

Để hiểu lý do tại sao, ta xem lại công thức cập nhật trọng số của SGD

$$\mathbf{w} = \mathbf{w} + \alpha \nabla L_{\mathbf{w}}.$$

Trong đó,  $\nabla L_{\mathbf{w}} = [\partial L/\partial w_1, \partial L/\partial w_2, \dots, \partial L/\partial w_m]^T$ . Rõ ràng,  $\partial L/\partial w_1, \dots, \partial L/\partial w_m$  phụ thuộc vào các đặc trưng đầu vào. Vậy nên, khi các đặc trưng có tỉ lệ chênh lệch lớn thì các hệ số  $\partial L/\partial w_i$  cũng chênh lệch lớn.

Do đó, về mặt tổng thể chúng ta sẽ không thể tìm được một  $learning rate <math>\alpha$  nào có thể thoả mãn được hết  $tất \ cả \ các \ hệ \ số$   $chênh \ lệch \ dó$ .

#### Main Advantages of Normalized Features

Nếu đặc trung đầu vào được chuẩn hoá, chúng ta sẽ thu được:

- Easier to find a good learning rate
- ② Getting numerically more stable gradients
- The training will be faster due to faster convergence, which essentially means that we need fewer training epochs.

## Common Feature Normalization Techniques

Hiện nay, chúng ta thường sử dụng 2 kỹ thuật chuẩn hoá dữ liệu trong ML/DL cho dữ liệu có cấu trúc:

- 0-1 Normalization/Min-max normalization
- Z-Score Standardization

## 0-1 Normalization/Min-max normalization

Now one common normalization technique is 0-1 normalization or sometimes also called min-max normalization:

$$norm(x_j^{[i]}) = \frac{x_j^{[i]} - min(x_j)}{max(x_j) - min(x_j)}$$

Để hiểu hơn công thức, ta xem ví dụ về Wine dataset trước đó:

		Alcohol	Malic acid	Ash	Alcalinity of ash	Magnesium	Total phenols	Flavanoids	Nonflavanoid phenols	Proanthocyanins	Color intensity	Hue	of diluted wines	Proline
	0	14.23	1.71	2.43	15.6	127	2.80	3.06	0.28	2.29	5.64	1.04	3.92	1065
	1	13.20	1.78	2.14	11.2	100	2.65	2.76	0.26	1.28	4.38	1.05	3.40	1050
	2	13.16	2.36	2.67	18.6	101	2.80	3.24	0.30	2.81	5.68	1.03	3.17	1185
	3	14.37	1.95	2.50	16.8	113	3.85	3.49	0.24	2.18	7.80	0.86	3.45	1480
	4	13.24	2.59	2.87	21.0	118	2.80	2.69	0.39	1.82	4.32	1.04	2.93	735
1	73	13.71	5.65	2.45	20.5	95	1.68	0.61	0.52	1.06	7.70	0.64	1.74	740
1	74	13.40	3.91	2.48	23.0	102	1.80	0.75	0.43	1.41	7.30	0.70	1.56	750
1	75	13.27	4.28	2.26	20.0	120	1.59	0.69	0.43	1.35	10.20	0.59	1.56	835

# 0-1 Normalization/Min-max normalization I

#### Load và hiển thị bộ dữ liệu wine dataset

Trieu Hai Nguyen Image Processing BM. KTPM-Khoa CNTT

## 0-1 Normalization/Min-max normalization II

	Alcohol	Malic Xacid	Ash	Alcalinity X of ash	Magnesium	Total phenols	Flavanoids	Nonflay ph	anoid enois	Proanthocyanins	Color intensity	Hue	OD280/OD315 of diluted wines	Proline
0	14.23	1.]	2.43	15.6	127	2.80	3.06	• • •	0.28	2.29	5.64	1.04	3.92	1065
1	13.20	1.78	2.14	[i] <sub>11.2</sub>	100	2.65	2.76		0.26	1.28	4.38	1.05	3.40	1050
2	13.16	2.36	2.67	18.6	101	2.80	3.24		0.30	2.81	5.68	1.03	3.17	1185
3	14.37	1.95	2.50	16.8	113	3.85	3.49		0.24	2.18	7.80	0.86	3.45	1480
4	13.24	2.59	2.87	21.0	118	2.80	2.69		0.39	1.82	4.32	1.04	2.93	735
								•••						
173	13.71	5.65	2.45	20.5	95	1.68	0.61		0.52	1.06	7.70	0.64	1.74	740
174	13.40	3.91	2.48	23.0	102	1.80	0.75		0.43	1.41	7.30	0.70	1.56	750
175	13.27	4.28	2.26	200	120	1.59	0.69		0.43	1.35	10.20	0.59	1.56	835
176	13.17	2.59	2.37	20.0	120	1.65	0.68		0.53	1.46	9.30	0.60	1.62	840
177	14.13	4.10	2.74	24.5	96	2.05	0.76		0.56	1.35	9.20	0.61	1.60	560
178 rd	ows × 13	column	m	1n/r	nax(2	<b>(</b> j )								

Trieu Hai Nguyen Image Processing BM. KTPM-Khoa CNTT

## 0-1 Normalization/Min-max normalization III

### Tính $\min(x_j)/\max(x_j)$

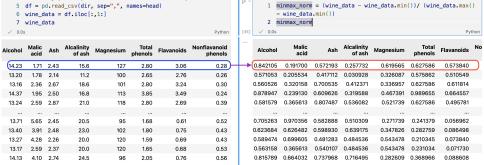
The minimum and maximum values of each feature column are wildly different.

<pre>1 wine_data.min()</pre>		1 wine_data.max()	
✓ 0.0s	Python	[40]	
Alcohol	11.03	··· Alcohol 14	.83
Malic acid	0.74	Malic acid 5	80
Ash	1.36	Ash 3	3.23
Alcalinity of ash	10.60	Alcalinity of ash 30	.00
Magnesium	70.00	Magnesium 162	2.00
Total phenols	0.98	Total phenols	8.88
Flavanoids	0.34	Flavanoids 5	.08
Nonflavanoid phenols	0.13	Nonflavanoid phenols	.66
Proanthocyanins	0.41	Proanthocyanins 3	.58
Color intensity	1.28	Color intensity 13	.00
Hue	0.48	Hue 1	.71
OD280/OD315 of diluted wines	1.27	OD280/OD315 of diluted wines 4	.00
Proline	278.00	Proline 1680	.00
dtype: float64		dtype: float64	

Trieu Hai Nguyen Image Processing BM. KTPM-Khoa CNTT 28 / 35

# 0-1 Normalization/Min-max normalization

#### Sau khi áp dụng 0-1 Normalization, ta thu được



Trieu Hai Nguyen Image Processing BM. KTPM-Khoa CNTT 29/35

## 0-1 Normalization/Min-max normalization

Rõ ràng như tên gọi 0-1 Normalization, sau khi chuẩn hoá các đặc trưng mới có giá trị nhỏ nhất là 0 và lớn nhất là 1.

<pre>1 minmax_norm.min()</pre>		1 minmax_norm.max()
✓ 0.0s	Python	[47]
Alcohol	0.0	··· Alcohol 1.0
Malic acid	0.0	Malic acid 1.0
Ash	0.0	Ash 1.0
Alcalinity of ash	0.0	Alcalinity of ash 1.0
Magnesium	0.0	Magnesium 1.0
Total phenols	0.0	Total phenols 1.0
Flavanoids	0.0	Flavanoids 1.0
Nonflavanoid phenols	0.0	Nonflavanoid phenols 1.0
Proanthocyanins	0.0	Proanthocyanins 1.0
Color intensity	0.0	Color intensity 1.0
Hue	0.0	Hue 1.0
OD280/OD315 of diluted wines	0.0	OD280/OD315 of diluted wines 1.0
Proline	0.0	Proline 1.0
dtype: float64	_	dtype: float64

Trieu Hai Nguyen Image Processing BM. KTPM-Khoa CNTT

#### **Z-Score Standardization**

Một chuẩn hoá tiếp theo thường hay được sử dụng hơn 0-1 Normalization là **Z-Score Standardization** 

$$standar(x_j^{[i]}) = \frac{x_j^{[i]} - mean(x_j)}{std(x_j)}.$$

Trong đó,  $mean(x_j)$ ,  $std(x_j)$  lần lượt là là giá trị trung bình và độ lệch chuẩn của cột đặc trưng thứ j.

After this z-score standardization, the mean of these features will be centered at *zero* and the standard deviation of the features will be *one*.

#### **Z-Score Standardization**

Sau khi áp dụng chuẩn hoá Z-Score Standardization, ta thu được các đặc trưng có giá trị như hình dưới.

```
1 znorm = (wine_data - wine_data.mean())/wine_data.std()
```

2 znorm

v/ 0.0s

V 0.08											Python
Alcohol	Malic acid	Ash	Alcalinity of ash	Magnesium	Total phenols	Flavanoids	Nonflavanoid phenols	Proanthocyanins	Color intensity	Hue	OD280/OD315 of diluted wines
1.514341	-0.560668	0.231400	-1.166303	1.908522	0.806722	1.031908	-0.657708	1.221438	0.251009	0.361158	1.842721
0.245597	-0.498009	-0.825667	-2.483841	0.018094	0.567048	0.731565	-0.818411	-0.543189	-0.292496	0.404908	1.110317
0.196325	0.021172	1.106214	-0.267982	0.088110	0.806722	1.212114	-0.497005	2.129959	0.268263	0.317409	0.786369
1.686791	-0.345835	0.486554	-0.806975	0.928300	2.484437	1.462399	-0.979113	1.029251	1.182732	-0.426341	1.180741
0.294868	0.227053	1.835226	0.450674	1.278379	0.806722	0.661485	0.226158	0.400275	-0.318377	0.361158	0.448336
0.873810	2.966176	0.304301	0.300954	-0.331985	-0.982841	-1.420891	1.270726	-0.927563	1.139596	-1.388840	-1.227742
0.491955	1.408636	0.413653	1.049555	0.158126	-0.791103	-1.280731	0.547563	-0.316058	0.967055	-1.126341	-1.481267
0.331822	1.739837	-0.388260	0.151234	1.418411	-1.126646	-1.340800	0.547563	-0.420888	2.217979	-1.607590	-1.481267
0.208643	0.227053	0.012696	0.151234	1.418411	-1.030776	-1.350811	1.351077	-0.228701	1.829761	-1.563840	-1.396759
1.391162	1.578712	1.361368	1.498716	-0.261969	-0.391646	-1.270720	1.592131	-0.420888	1.786626	-1.520090	-1.424928

Trieu Hai Nguyen Image Processing BM. KTPM-Khoa CNTT 32 / 35

#### **Z-Score Standardization**

#### Zero Mean & Standard Deviation 1

Sau khi chuẩn hoá Z-Score Standardization, kiểm tra xem trung bình của các cột đặc trưng bằng  $\theta$  và độ lệch chuẩn bằng  $\theta$ .

<pre>1 znorm.mean()</pre>			1 znorm.std()	
✓ 0.0s	Python	[56]	✓ 0.0s	
Alcohol	-9.181170e-16		Alcohol	1.0
Malic acid	0.000000e+00		Malic acid	1.0
Ash	-8.070947e-16		Ash	1.0
Alcalinity of ash	-7.983626e-17		Alcalinity of ash	1.0
Magnesium	-1.995907e-17		Magnesium	1.0
Total phenols	3.991813e-17		Total phenols	1.0
Flavanoids	-3.592632e-16		Flavanoids	1.0
Nonflavanoid phenols	3.592632e-16		Nonflavanoid phenols	1.0
Proanthocyanins	-1.596725e-16		Proanthocyanins	1.0
Color intensity	1.995907e-17		Color intensity	1.0
Hue	1.995907e-16		Hue	1.0
OD280/OD315 of diluted wines	3.193450e-16		OD280/OD315 of diluted wines	1.0
Proline	-7.983626e-17		Proline	1.0
dtype: float64			dtype: fleat64	

Lưu ý, khi chuẩn hoá dữ liệu đầu vào thì phải chuẩn hoá cho cả 3 tập: training, validation và test set.

#### Chuẩn hoá cho 3 tập sử dụng 0-1 Normalization

```
return (df - train_min)/ (train_max - train_min)

train_min, train_max = df_train.min(), df_train.max()

df_train_norm = normalize(df_train,train_min,train_max)

df_val_norm = normalize(df_val,train_min,train_max)

df_test_norm = normalize(df_test,train_min,train_max)
```

def normalize (df, train\_min, train\_max):

Thực hiện tương tự cho Z-score standardization trên 3 tập.

## Tài liệu tham khảo

- Sebastian Raschka, Yuxi (Hayden) Liu, Vahid Mirjalili Machine Learning with PyTorch and Scikit-Learn: Develop machine learning and deep learning models with Python (2022). Published by Packt Publishing Ltd, ISBN 978-1-80181-931-2.
- Sebastian Raschka
  MACHINE LEARNING Q AND AI: 30 Essential Questions
  and Answers on Machine Learning and AI (2024). ISBN-13:
  978-1-7185-0377-9 (ebook).
- LightningAI LightningAI: PyTorch Lightning (2024) .