

- How to split Dataset

[illegible]

- ▼ make dataset

```
# input data = {fish_length, weight}

import numpy as np

fish_data = [[l,w] for l, w in zip(fish_length, fish_weight)]
# [[25.4, 242.0, [26.3, 290.0], ...]

print(len(fish_data))

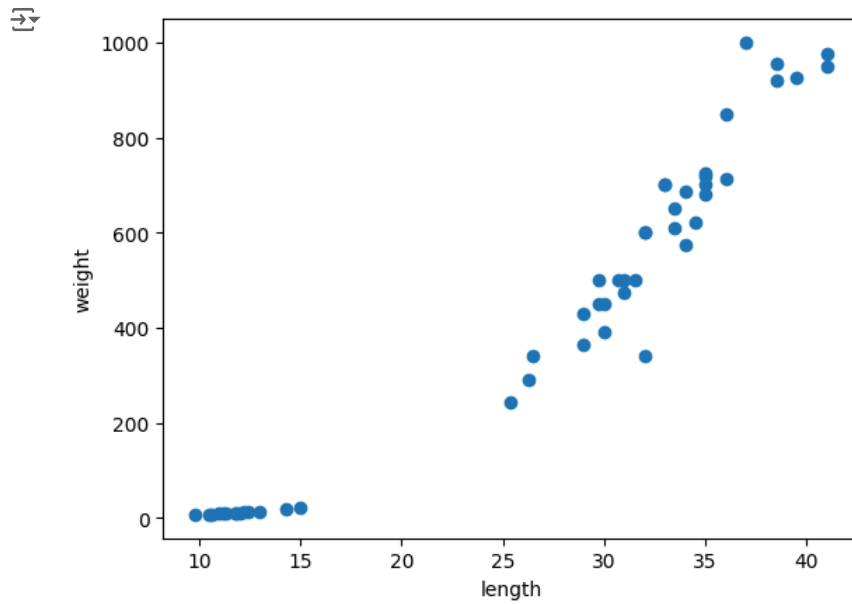
input_arr = np.array(fish_data) # R(49x2)
target_arr = np.array(fish_target) #R(49x1)
```

49

- analysis of dataset

```
import matplotlib.pyplot as plt

plt.scatter(input_arr[:, 0], input_arr[:, 1])
plt.xlabel('length')
plt.ylabel('weight')
plt.show()
```

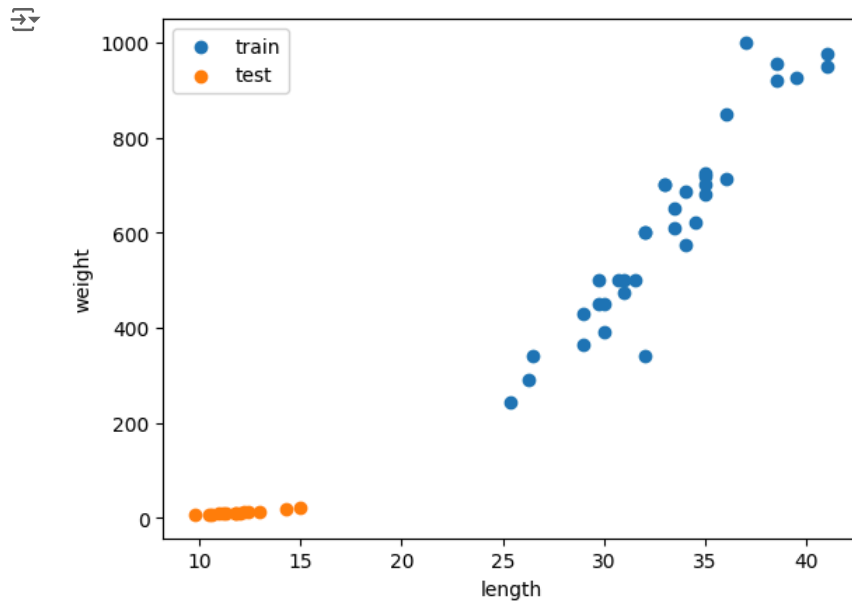


✓ spilit dataset (try #1)

```
train_input = input_arr[:35]
train_target = target_arr[:35]
```

```
test_input = input_arr[35:]
test_target = target_arr[35:]
```

```
plt.scatter(train_input[:, 0], train_input[:, 1], label='train')
plt.scatter(test_input[:, 0], test_input[:, 1], label='test')
plt.xlabel('length')
plt.ylabel('weight')
plt.legend()
plt.show()
```



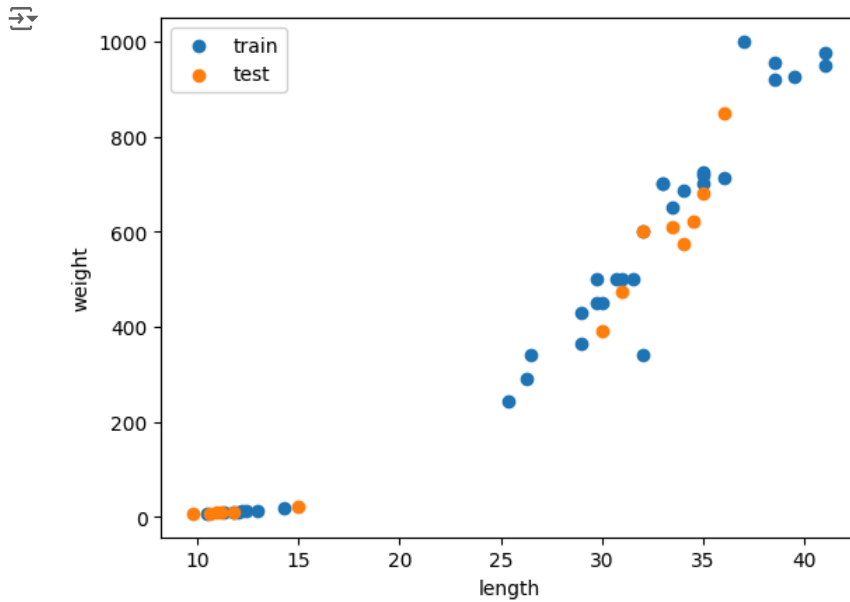
✓ spilit dataset (try #2: random shuffle)

```
np.random.seed(42)
index = np.arange(len(input_arr))
np.random.shuffle(index)

train_input = input_arr[index[:35]]
train_target = target_arr[index[:35]]

test_input = input_arr[index[35:]]
test_target = target_arr[index[35:]]

plt.scatter(train_input[:, 0], train_input[:, 1], label='train')
plt.scatter(test_input[:, 0], test_input[:, 1], label='test')
plt.xlabel('length')
plt.ylabel('weight')
plt.legend()
plt.show()
```



✓ spilit dataset (try #3: using sklearn)

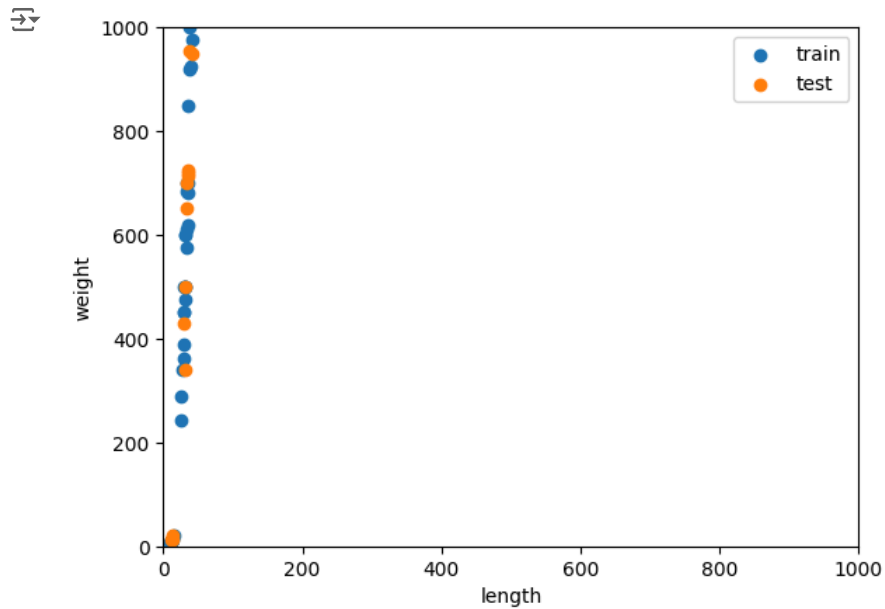
```
from sklearn.model_selection import train_test_split

train_input, test_input, train_target, test_target = train_test_split(
    fish_data, fish_target, random_state=42)

#print(train_target)

train_input = np.array(train_input)
train_target = np.array(train_target)
test_input = np.array(test_input)
test_target = np.array(test_target)

plt.scatter(train_input[:, 0], train_input[:, 1], label='train')
plt.scatter(test_input[:, 0], test_input[:, 1], label='test')
plt.xlim(0,1000)
plt.ylim(0,1000)
plt.xlabel('length')
plt.ylabel('weight')
plt.legend()
plt.show()
```



▼ split dataset (try #4: standardization)

```
#print(train_input.shape)
mean = np.mean(train_input, axis=0)

std = np.std(train_input, axis=0)

#print(mean, std)

train_scaled = (train_input - mean) / std
test_scaled = (test_input - mean) / std

#print(train_scaled, test_scaled)

plt.scatter(train_scaled[:, 0], train_scaled[:, 1], label='train')
plt.scatter(test_scaled[:, 0], test_scaled[:, 1], label='test')
plt.xlim(0,2)
plt.ylim(0,2)
plt.xlabel('length')
plt.ylabel('weight')
plt.legend()
plt.show()
```

