# Linear Regression using Normal Equation

Jisen Fang

Original code: [Chapter 4 - Training Linear Models - Colab](#)

What I did:

1.  Modify the code to read data from abalone_train.csv in the local drive instead of reading random data.
2.  Modify the input array to fit the size of abalone_train.csv

# Read data from local

```python
import numpy as np
import pandas as pd
%matplotlib inline
import matplotlib.pyplot as plt

from google.colab import files
uploaded = files.upload()  # upload the data from local

import io
abalone = pd.read_csv(
    io.BytesIO(uploaded['abalone_train.csv']),
    names=["Length", "Diameter", "Height", "Whole weight", "Shucked weight",
           "Viscera weight", "Shell weight", "Age"]) # read the data into a pandas DataFrame

X1 = abalone["Length"]

X2 = np.array(X1)
X = X2.reshape(-1, 1)

y1 = abalone["Height"]
y2 = np.array(y1)
y = y2.reshape(-1, 1)
```

# Modify the input array

```python
X_b = np.c_[np.ones((3320, 1)), X]   # 100 to 3320 since the size of abalone_train.csv is 3320
theta_best = np.linalg.inv(X_b.T.dot(X_b)).dot(X_b.T).dot(y)
```

# All output

```python
1  import numpy as np
2  import pandas as pd
3  %matplotlib inline
4  import matplotlib.pyplot as plt
5
6  from google.colab import files
7  uploaded  = files.upload()
8
9  import io
10 abalone = pd.read_csv(
11     io.BytesIO(uploaded['abalone_train.csv']),
12     names=["Length", "Diameter", "Height", "Whole weight", "Shucked weight",
13            "Viscera weight", "Shell weight", "Age"])
14
15 X1 = abalone["Length"]
16
17 X2 = np.array(X1)
18
19 X = X2.reshape(-1, 1)
20
21 y1 = abalone["Height"]
22 y2 = np.array(y1)
23 y = y2.reshape(-1, 1)
```
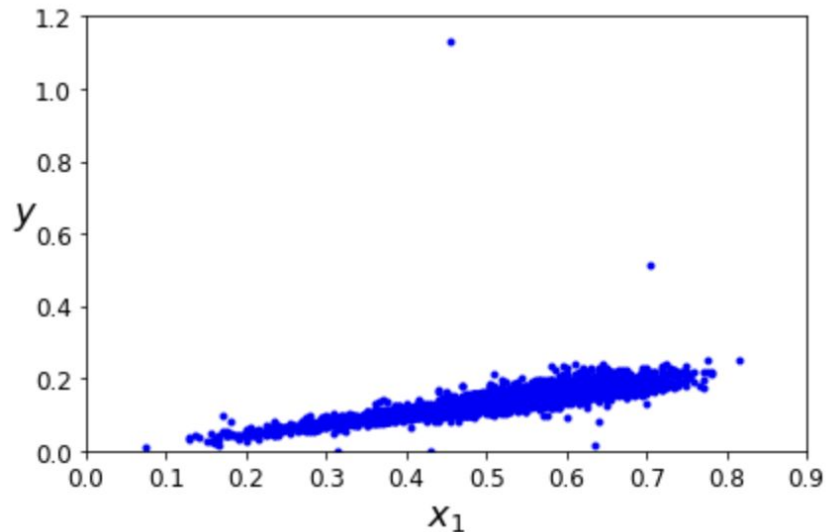
Browse... abalone_train.csv
**abalone_train.csv**(text/csv) - 145915 bytes, last modified: n/a - 100% done
Saving abalone_train.csv to abalone_train (9).csv

```
1 plt.plot(X, y, "b.")
2 plt.xlabel("$x_1$", fontsize=18)
3 plt.ylabel("$y$", rotation=0, fontsize=18)
4 plt.axis([0, 0.9, 0, 1.2])
5 save_fig("generated_data_plot")
6 plt.show()
```

Saving figure generated_data_plot

```python
1 X_b = np.c_[np.ones((3320, 1)), X]  # add x0 = 1 to each instance
2 theta_best = np.linalg.inv(X_b.T.dot(X_b)).dot(X_b.T).dot(y)
```

[25]
```python
1 theta_best
```

```
array([[-0.0108267 ],
       [ 0.28716253]])
```

[26]
```python
1 X_new = np.array([[0], [2]])
2 X_new_b = np.c_[np.ones((2, 1)), X_new]  # add x0 = 1 to each instance
3 y_predict = X_new_b.dot(theta_best)
4 y_predict
```
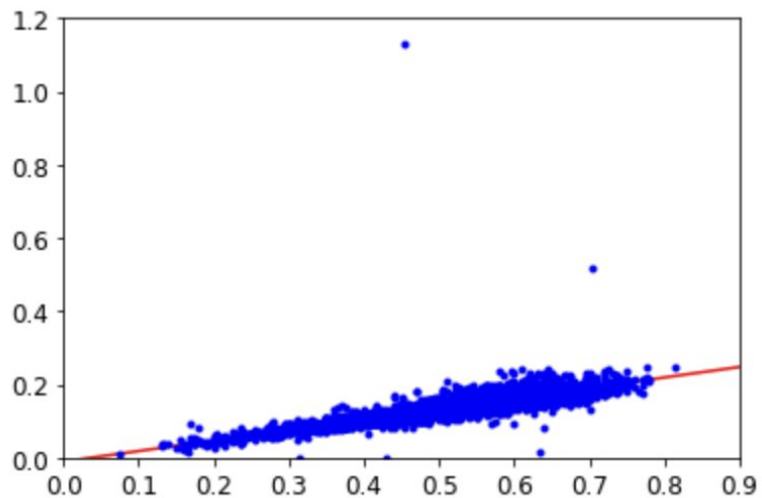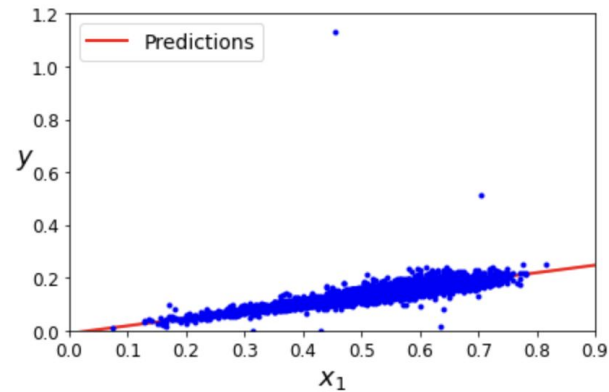
```
array([[-0.0108267 ],
       [ 0.56349837]])
```

```
1 plt.plot(X_new, y_predict, "r-")
2 plt.plot(X, y, "b.")
3 plt.axis([0, 0.9, 0, 1.2])
4 plt.show()
```

```
[37]  1 plt.plot(X_new, y_predict, "r-", linewidth=2, label="Predictions")
      2 plt.plot(X, y, "b.")
      3 plt.xlabel("$x_1$", fontsize=18)
      4 plt.ylabel("$y$", rotation=0, fontsize=18)
      5 plt.legend(loc="upper left", fontsize=14)
      6 plt.axis([0, 0.9, 0, 1.2])
      7 save_fig("linear_model_predictions_plot")
      8 plt.show()
```

Saving figure linear_model_predictions_plot

```
[31]  1 from sklearn.linear_model import LinearRegression
      2
      3 lin_reg = LinearRegression()
      4 lin_reg.fit(X, y)
      5 lin_reg.intercept_, lin_reg.coef_
```

(array([-0.0108267]), array([[0.28716253]]))

```
[32]  1 lin_reg.predict(X_new)
```

array([[-0.0108267 ],
       [ 0.56349837]])

```
[33]  1 theta_best_svd, residuals, rank, s = np.linalg.lstsq(X_b, y, rcond=1e-6)
      2 theta_best_svd
```

array([[-0.0108267 ],
       [ 0.28716253]])

```
[34]  1 np.linalg.pinv(X_b).dot(y)
```

array([[-0.0108267 ],
       [ 0.28716253]])