

CS550 - Machine Learning and Business Intelligence

Jupyter: Training Linear Models

By: Maranata Muluneh

Instructor: Dr. Chang, Henry

2023

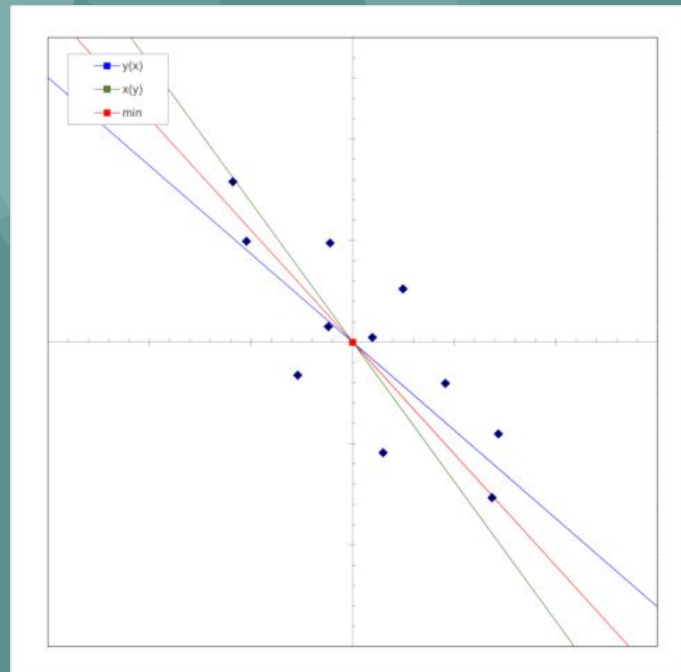




Table of Content

- Introduction
- Theory
- Implementation
- References



Introduction

The normal equation is a method to find the optimal solution for linear regression problems with a single variable or multiple variables. It is an analytical approach to find the values of the coefficients (weights) that minimize the mean squared error between the predicted output and the true output. This method is efficient for small to medium sized datasets, but becomes computationally expensive for large datasets as the computational complexity is $O(n^3)$.



Theory

Machine Learning Process

During Machine Learning process we have three steps:

- Collecting data
- Creating model
- Prediction



Theory

Collect Data

- After we collected the data we needed we have to classify the collected data in to **training, validation and test** categories .
- Training (50% of the data)
- Evaluation(25% of the data)
- Test (25% of the data)



Theory

Create Model

- To create a model we have different Machine Learning (ML) types to evaluate the models, for this slide **we are focusing on supervised learning specifically on Linear Regression.**



Theory

Supervised learning

- Regression is one of the types if supervised learning includes Linear and nonlinear regression.
- We use overfitting model to evaluate the model either MSE or RMSE, while MSE is for numbers evaluation RMSE is for strings.

Mean Squared Error (MSE) , Root Mean Squared Error (RMSE)



Theory

Linear Regression

- A linear regression model follows a very particular form.
- Regression Equation(y) = $a + bx$
- Slope(b) = $(N\sum XY - (\sum X)(\sum Y)) / (N\sum X^2 - (\sum X)^2)$
- Intercept(a) = $(\sum Y - b(\sum X)) / N$



Theory

Linear Regression

Where:

- x and y are the variables.
- b = The slope of the regression line
- a = The intercept point of the regression line and the y axis.
- N = Number of values or elements
- X = First Score
- Y = Second Score
- ΣXY = Sum of the product of first and Second Scores
- ΣX = Sum of First Scores
- ΣY = Sum of Second Scores
- ΣX^2 = Sum of square First Scores



Theory

Non Linear Regression

- A Non linear regression model is flexible and will not follow a specific rule or form.
- Regression Equation(y) = $a + b^2$
- Slope(b) = $(N\sum PY - (\sum P)(\sum Y)) / (N\sum P^2 - (\sum P)^2)$
- Intercept(a) = $(\sum Y - b(\sum P)) / N$
- Where $P = X * X$

For this slide we are focusing on Linear Regression



Implementation

Environment : Colab, Tensorflow 2

Programming Language: Python

Libraries 

```
import matplotlib.pyplot as plt
```

```
import numpy as np
```

```
import sys
```

```
import tensorflow as tf
```

```
from sklearn import datasets
```

```
from tensorflow.python.framework import ops
```

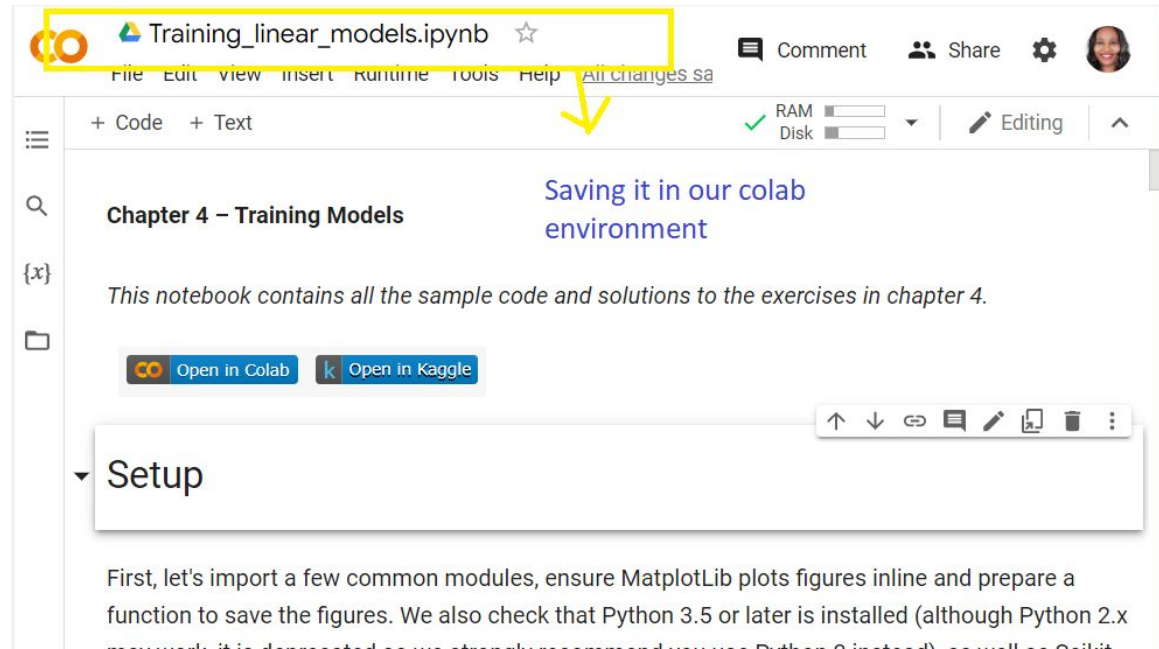
```
ops.reset_default_graph()
```

```
sess = tf.Session()
```

```
iris = datasets.load_iris()
```

Implementation

Opening the sample code in colab, saving it and running it .



The screenshot shows a Google Colab notebook interface. At the top, the title bar displays the Colab logo, the filename 'Training_linear_models.ipynb', and a star icon. Below the title bar is a menu bar with options: File, Edit, View, Insert, Runtime, Tools, Help, and All changes saved. A yellow box highlights the title bar area, and a yellow arrow points from the text 'Saving it in our colab environment' to the 'All changes saved' link. The main content area has a tab labeled '+ Code' and '+ Text'. The notebook content includes the title 'Chapter 4 – Training Models', a paragraph stating 'This notebook contains all the sample code and solutions to the exercises in chapter 4.', and two buttons: 'Open in Colab' and 'Open in Kaggle'. Below this is a section titled 'Setup' with a dropdown arrow. The text under 'Setup' reads: 'First, let's import a few common modules, ensure Matplotlib plots figures inline and prepare a function to save the figures. We also check that Python 3.5 or later is installed (although Python 2.x is supported as well, although deprecated, we recommend using Python 3 instead), as well as Colab'.

Training_linear_models.ipynb ☆

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

Chapter 4 – Training Models

This notebook contains all the sample code and solutions to the exercises in chapter 4.

Open in Colab Open in Kaggle

Setup

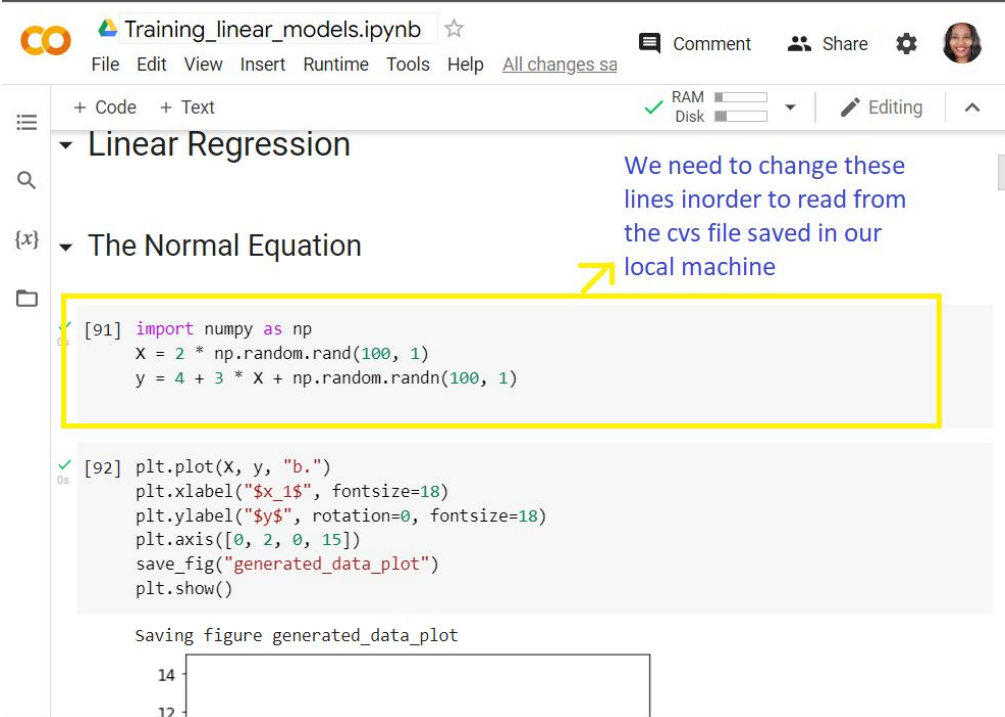
First, let's import a few common modules, ensure Matplotlib plots figures inline and prepare a function to save the figures. We also check that Python 3.5 or later is installed (although Python 2.x is supported as well, although deprecated, we recommend using Python 3 instead), as well as Colab

Implementation

Instead of reading random data

We are going to use the csv file.

We are going to use the csv file.



The screenshot shows a Jupyter Notebook titled "Training_linear_models.ipynb". The notebook has two sections: "Linear Regression" and "The Normal Equation". A yellow box highlights the code in the "Linear Regression" section, which generates random data for a linear regression model. A yellow arrow points from the text "We need to change these lines in order to read from the csv file saved in our local machine" to the highlighted code block.

```
[91] import numpy as np
X = 2 * np.random.rand(100, 1)
y = 4 + 3 * X + np.random.randn(100, 1)
```

```
[92] plt.plot(X, y, "b.")
plt.xlabel("$x_1$", fontsize=18)
plt.ylabel("$y$", rotation=0, fontsize=18)
plt.axis([0, 2, 0, 15])
save_fig("generated_data_plot")
plt.show()
```

Saving figure generated_data_plot

14

12

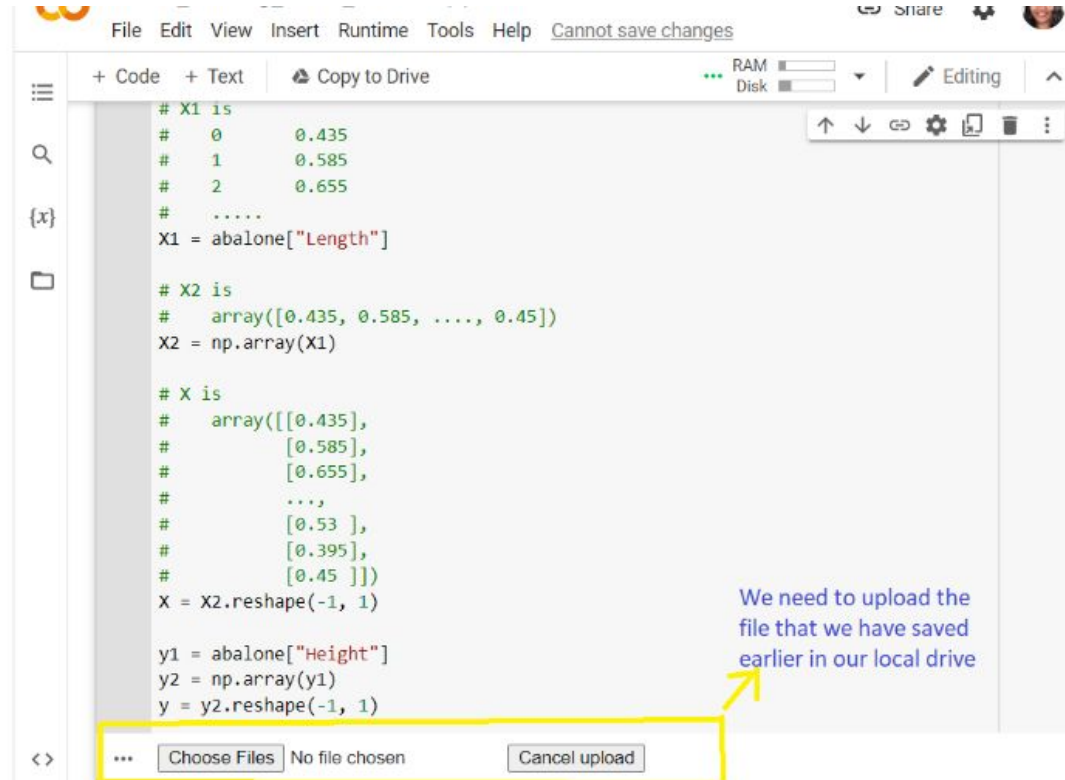


Implementation

```
import numpy as np
import pandas as pd
# X = 2 * np.random.rand(100, 1)
# y = 4 + 3 * X + np.random.randn(100, 1)
from google.colab import files
uploaded = files.upload()
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"])
# X1 is
# 0    0.435
# 1    0.585
# 2    0.655
#
```

```
X1 = abalone["Length"]
# X2 is
# array([0.435, 0.585, ..., 0.45])
X2 = np.array(X1)
# X is
# array([[0.435],
#        [0.585],
#        [0.655],
#        ...,
#        [0.53 ],
#        [0.395],
#        [0.45 ]])
X = X2.reshape(-1, 1)
y1 = abalone["Height"]
y2 = np.array(y1)
y = y2.reshape(-1, 1)
```

Implementation



```
File Edit View Insert Runtime Tools Help Cannot save changes
+ Code + Text Copy to Drive RAM Disk Editing
# X1 is
# 0 0.435
# 1 0.585
# 2 0.655
# .....
X1 = abalone["Length"]

# X2 is
# array([0.435, 0.585, ..., 0.45])
X2 = np.array(X1)

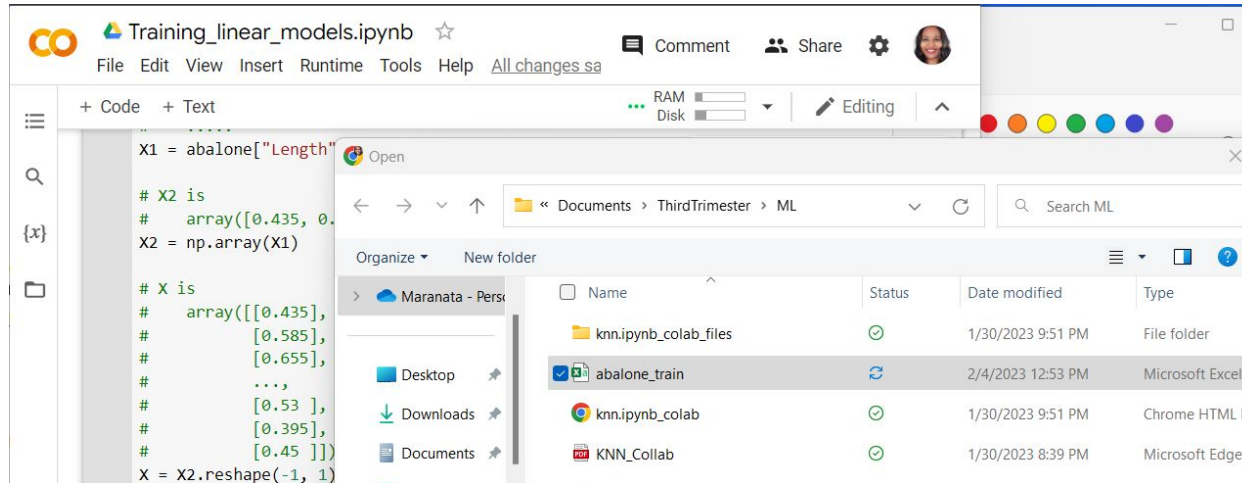
# X is
# array([[0.435],
#        [0.585],
#        [0.655],
#        ...,
#        [0.53 ],
#        [0.395],
#        [0.45 ]])
X = X2.reshape(-1, 1)

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

... Choose Files No file chosen Cancel upload

We need to upload the file that we have saved earlier in our local drive

Implementation



The screenshot displays a Jupyter Notebook environment with a file explorer window open. The notebook, titled "Training_linear_models.ipynb", shows the following code in the "Code" cell:

```
X1 = abalone["Length"]

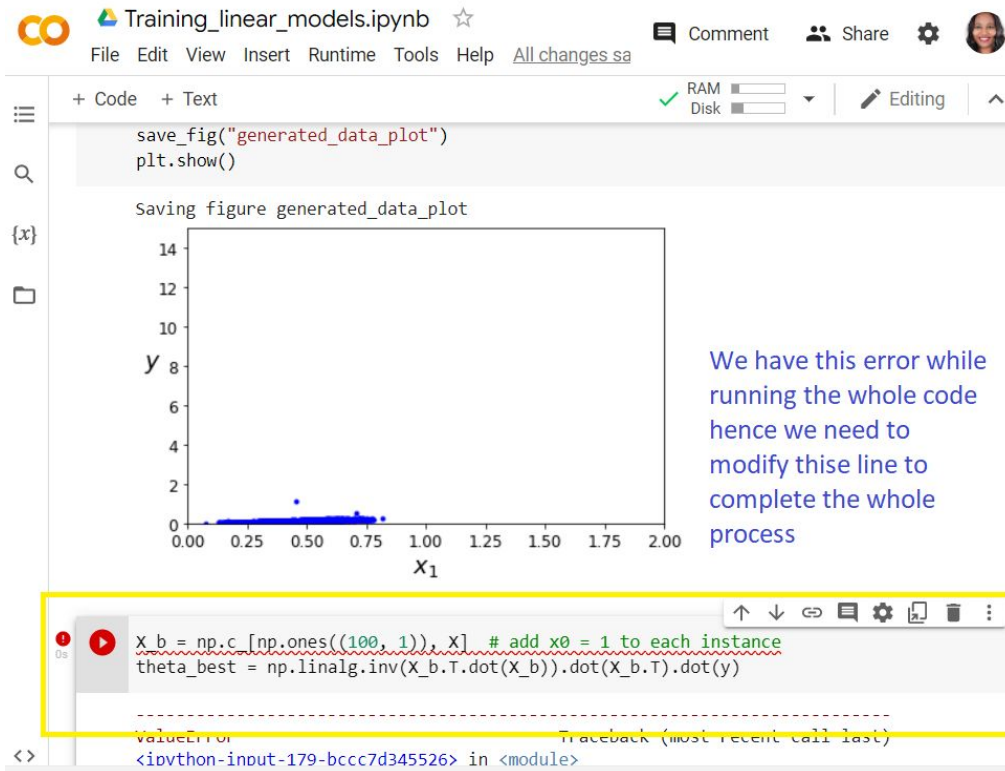
# X2 is
# array([[0.435, 0.585],
#        [0.655, 0.53 ],
#        ...,
#        [0.395, 0.45 ]])
X2 = np.array(X1)

# X is
# array([[0.435],
#        [0.585],
#        [0.655],
#        ...,
#        [0.395],
#        [0.45 ]])
X = X2.reshape(-1, 1)
```

The file explorer window, titled "Open", shows the directory path "Documents > ThirdTrimester > ML". It lists several files and folders:

Name	Status	Date modified	Type
knn.ipynb_colab_files	✓	1/30/2023 9:51 PM	File folder
abalone_train	🔄	2/4/2023 12:53 PM	Microsoft Excel C
knn.ipynb_colab	✓	1/30/2023 9:51 PM	Chrome HTML Do
KNN_Collab	✓	1/30/2023 8:39 PM	Microsoft Edge P

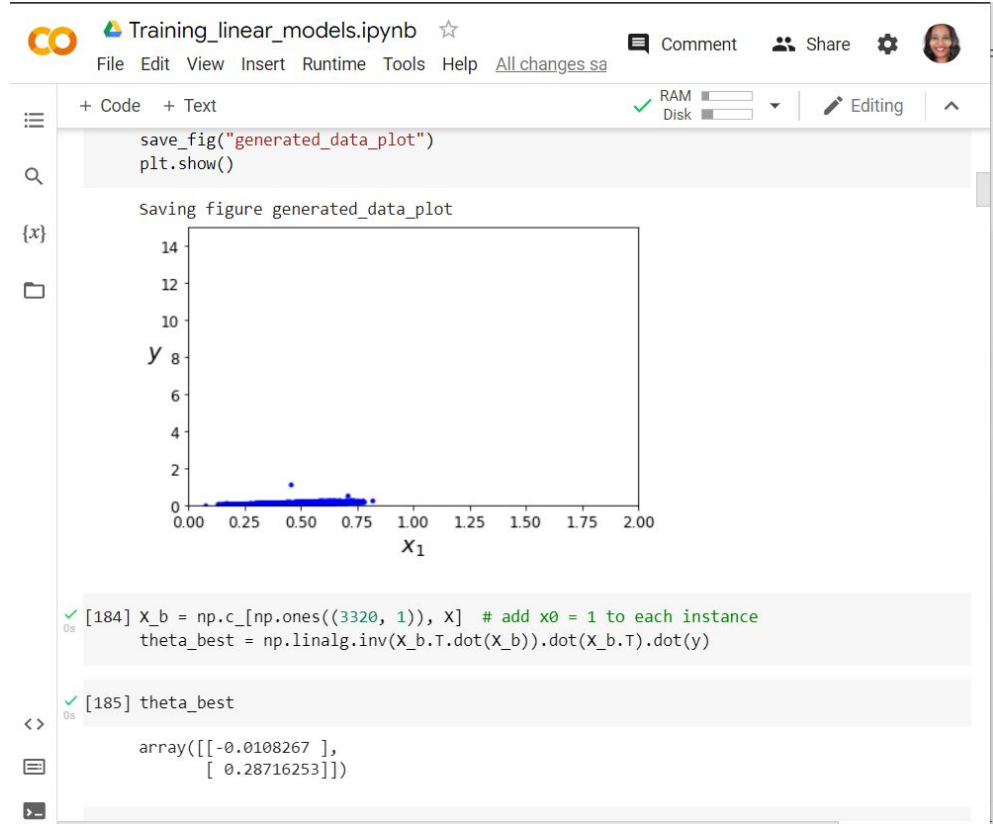
Implementation



The total number of data is 3320 hence we need to change the value.

Implementation

After we changed the value of the size we have the complete process running.





References

- Google. (n.d.). Google colaboratory. Google Colab. Retrieved February 4, 2023, from https://colab.research.google.com/github/ageron/handson-ml2/blob/master/04_training_linear_models.ipynb
- McClure, N. (n.d.). Tensorflow machine learning cookbook. O'Reilly Online Learning. Retrieved February 4, 2023, from <https://learning.oreilly.com/library/view/tensorflow-machine-learning/9781786462169/ch03s04.html>