**ABERDEEN 2040**

# GitHub, Python 3 and an Introduction to Machine Learning with Linear Regression

Exploratory Data Analysis

California Housing CSV Dataset

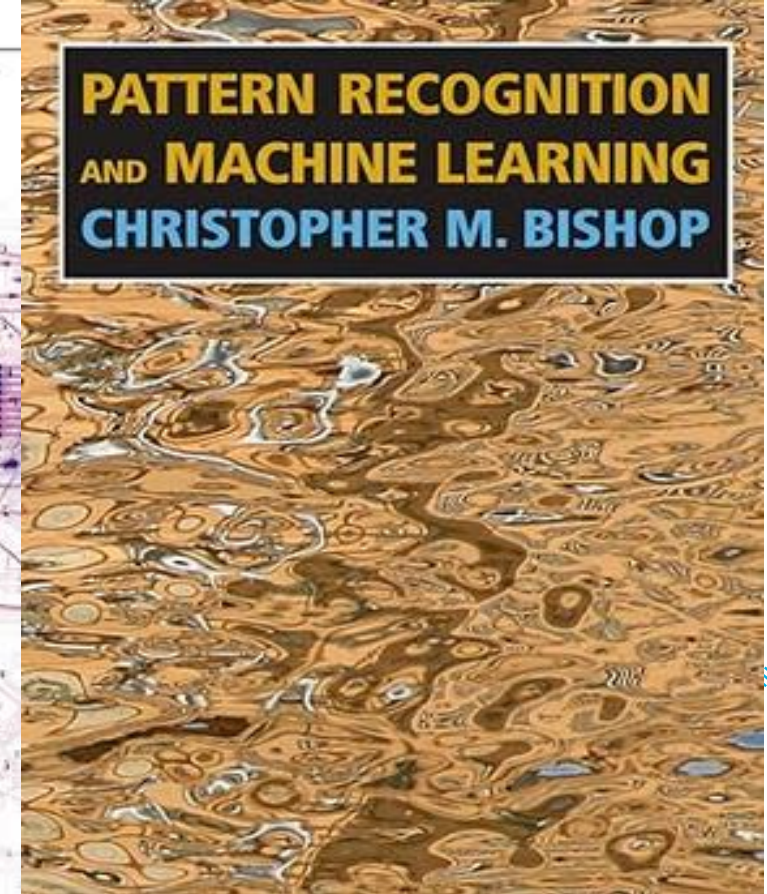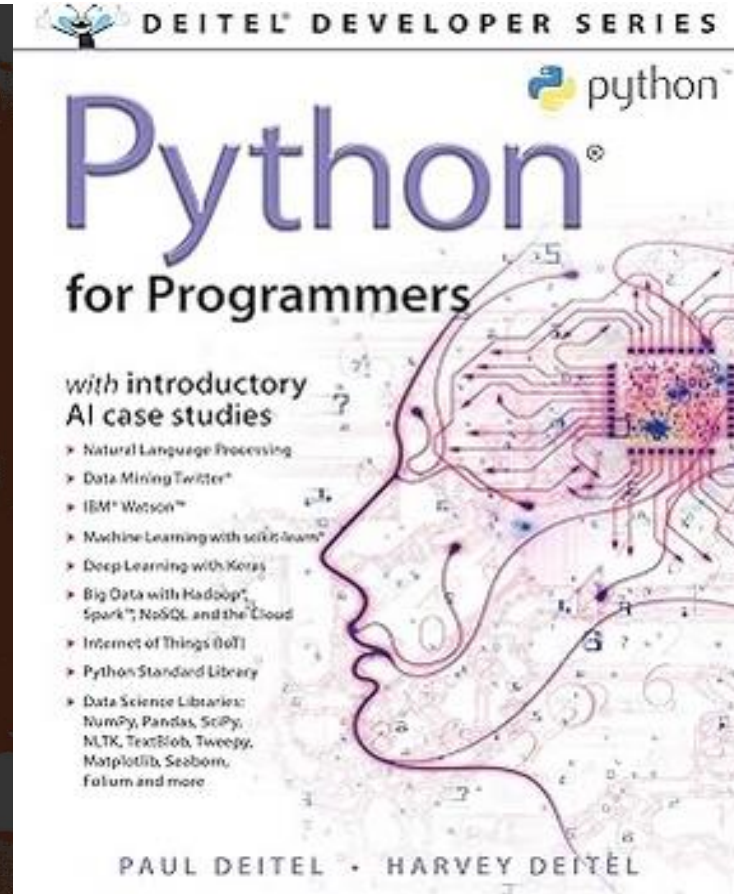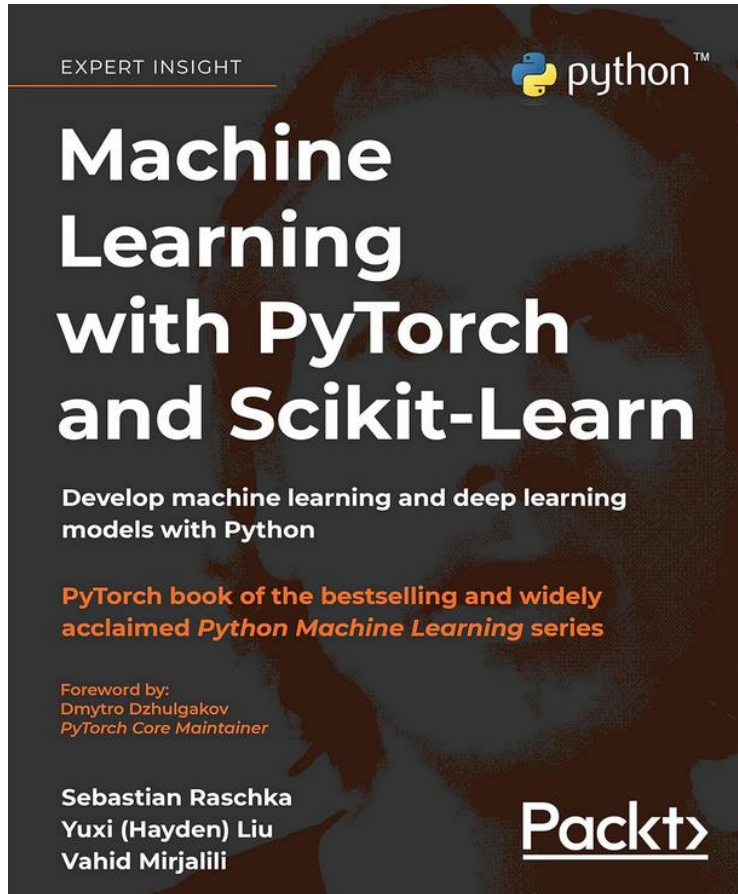**Coding, GIS and Remote Sensing Data for Glaciology**

27 February 2025

St Mary's G15

# Overview

1. GitHub

2. Python 3

3. Machine Learning Cross Validation Methodology

4. Simple Linear Regression

5. Multiple Linear Regression

# Machine Learning Reading

# Anaconda Cheat Sheet



https://docs.conda.io/projects/conda/en/4.6.0/_downloads/52a95608c49671267e40c689e0bc00ca/conda-cheatsheet.pdf

# Python Code for the Session

The slides and the code for the session can be found implemented at:

GitHub ID: Stevieee83

https://github.com/Stevieee83

Machine Learning Library: ScikitLearn

# GitHub Web Browser Demonstration

# Python 3

# Python 3 ipython Interpreter

# Machine Learning with Linear Regression

# Machine Learning Libraries

- ScikitLearn (well recognised machine learning library)

- Pandas (data processing Python library)

- NumPy (numerical Python library)

- Matplotlib (data visualisation Python library)

- Seaborn (data visualisation Python library)

- SciPy (statistics Python library)

# Gradient Descent



$$X = X - learning\_rate * gradient$$

# Exploratory Data Analysis Code Demonstration

# Data Pre-Processing

# Data Normalisation



original data     zero-centered data     normalized data

https://cs231n.stanford.edu/

# **Hold Out Cross Validation**

- Split the dataset to 80% training and 20% test data

- Split the training dataset again to generate a 20% validation dataset (20% over all of the examples is recommended, not just the training dataset)

- Only use the test dataset after training and validating the model with the training and validation datasets.

# Hold Out Cross Validation



100 Total Overall Dataset Examples

| Training Dataset | Validation Dataset | Test Dataset |
|---|---|---|
| 60% to 70% of the Dataset | 15% to 20% Dataset | 15% to 20% Dataset |
| 60 to 70 Dataset Examples | 15 to 20 Dataset Examples | 15 to 20 Dataset Examples |

# Simple Linear Regression

# Simple Linear Regression



Regression Plot MedInc vs MedHouseValue

$$y(X, W) = W_o + (W1X1)$$

# Training Linear Regression

- Set the Learning Rate

- Adapt the training iterations if required

- Tune any regularisation parameters

# Loss in Linear Regression Models

- Mean Squared Error (MSE)

- Root Mean Squared Error (RMSE)

- Mean Absolute Error (MAE)

MAE does not penalise Linear Regression models as much when there are outliers in the data compared to MSE and RMSE

# Evaluating the Linear Regression

- R2 Score

- Pearson Linear Correlation Coefficient (PLCC)

- Spearman Rank Correlation Coefficient (SRCC)

- Kendal Rank Correlation Coefficient (KRCC)

# Simple Linear Regression Evaluation (ScikitLearn)

| Dataset Split | MSE Loss |
|---|---|
| Training | 0.52 |
| Validation | 0.55 |
| Test | 0.53 |

# Simple Linear Regression Evaluation (Python)

| Dataset Split | MSE Loss |
|---------------|----------|
| Training | 0.52 |
| Validation | 0.52 |
| Test | 0.53 |

# Simple Linear Regression Evaluation (ScikitLearn)

| Dataset Split | R2 Score | PLCC | SRCC | KRCC |
|---|---|---|---|---|
| Test | 0.46 | 0.68 | 0.67 | 0.49 |

# Simple Linear Regression Evaluation (Python)

| Dataset Split | R2 Score | PLCC | SRCC | KRCC |
|---|---|---|---|---|
| Test | 0.46 | 0.68 | 0.67 | 0.49 |

# Multiple Linear Regression

# Multiple Linear Regression

- Exactly the same concept as Simple Linear Regression except there is more than one feature for the model to use for prediction.

- 8 features are in the example notebooks for the session.

$$y(X, W) = W_o + (W_1 X_1) + \ldots\ldots + (W_D X_D)$$

$$y(X, W) = W_o + (W_1 X_1) + \ldots\ldots + (W_8 X_8)$$

# Multiple Linear Regression Evaluation (ScikitLearn)

| Dataset Split | MSE Loss |
|---|---|
| Training | 0.39 |
| Validation | 0.40 |
| Test | 0.42 |

# Multiple Linear Regression Evaluation (Python)

| Dataset Split | MSE Loss |
|---|---|
| Training | 0.39 |
| Validation | 0.39 |
| Test | 0.42 |

# Multiple Linear Regression Evaluation (ScikitLearn)

| Dataset Split | R2 Score | PLCC | SRCC | KRCC |
|---|---|---|---|---|
| Test | 0.58 | 0.76 | 0.81 | 0.62 |

# Multiple Linear Regression Evaluation (Python)

| Dataset Split | R2 Score | PLCC | SRCC | KRCC |
|---|---|---|---|---|
| Test | 0.58 | 0.76 | 0.81 | 0.62 |

# Summary

- Exploratory data analysis and data pre-processing are essential before training any machine learning model with any dataset.

- Understanding gradient descent and how the learning rate effects the weight parameters are essential for training any machine learning model that uses gradient descent .

- Linear Regression models can only be used on data where the features have linear relationships to the output predictions made by the model over the complete range.

# References

[1] Christopher M. Bishopl. Pattern Recognition and Machine Learning. Basingstoke, UK: Springer, 2007. isbn: 978-0387310732.

[2] Paul Deitel and Harvey Deitel. Python For Programmers. London, UK: Pearson, 2019. isbn: 978-0-13-522433-5.

[3] Li, F.-F. CS231N: Deep Learning for Computer Vision, Stanford University CS231n: Deep Learning for Computer Vision. Available at: https://cs231n.stanford.edu/ (Accessed: 24 April 2024).

[4] Sebastian Raschka, Yuxi (Hayden) Liu, and Vahid Mirjalili. Machine Learning with PyTorch and Scikit-Learn. Birmingham, UK: Packt Publishing, 2022. isbn: 978-1801819312.