

A Written Assignment of Programming with Python

Course: Master In Computer
Science

By

Huzefa Ayub

Matriculation: UPS10643244

Tutor: Cosmina Croitoru

Date: 17/10/2023

Abstract:

This task involves the development of a Python program aimed at automating the selection of the four most suitable functions from a pool of fifty candidates, using a dataset for training and evaluation. The primary objective is to identify functions that provide the best fit to the training data, a critical step in data modeling and analysis. The program employs various mathematical and statistical techniques to assess the goodness of fit of each function against the training data, ultimately aiding in data-driven decision-making and model selection. By accomplishing this task, the program streamlines the process of identifying ideal functions, enhancing the efficiency and accuracy of data analysis and modeling in various fields of research and application.

Introduction

We have three types of datasets: training datasets (A), a test dataset (B), and datasets for 50 ideal functions (C), all consisting of x-y pairs of values.

We use the training datasets (A) to train a model, then evaluate its performance using the test dataset (B) to ensure it can make accurate predictions on unseen data. The datasets for ideal functions (C) may be used for comparison or validation purposes, depending on the specific goals of our project.

We develop a Python program that leverages training data to select the four most suitable ideal functions from a pool of fifty provided in dataset (C). Subsequently, the program should utilize the supplied test data (B) to assess whether each x-y pair can be associated with one of the four chosen ideal functions. If a suitable mapping is found, the program should perform the mapping and record the associated deviation.

The selection of the four ideal functions during training should be based on minimizing the sum of squared y-deviations (least squares) from the training dataset (A).

When mapping individual test cases to the four ideal functions, the program should ensure that the maximum deviation of the calculated regression does not exceed the largest deviation between the training dataset (A) and the corresponding ideal function (C) chosen for it by more than a factor of $\sqrt{2}$.

Scope:

The scope of this task includes the following key aspects:

1. Data Preparation:

- Load and preprocess the training data (A), the test data (B), and the dataset containing fifty ideal functions (C).
- Ensure the datasets are correctly formatted, containing x-y pairs of values.

2. Function Selection:

- Implement a method to identify the top four ideal functions from dataset (C) that best match the training data.
- The criterion for selection should be based on minimizing the sum of squared y-deviations (least squares) between the training data and each ideal function.

3. Ideal Function Parameterization:

- Define mathematical representations of the ideal functions (e.g., linear, polynomial, etc.) that will be considered during the selection process.
- These functions should be parameterized and adaptable to different data patterns.

4. Parameter Optimization:

- Employ an optimization algorithm or method to compute the parameters of the ideal functions that minimize the sum of squared y-deviations for each function.
- Repeat this process iteratively to identify the top four ideal functions.

5. Results Reporting:

- Present the four chosen ideal functions along with their respective parameters.
- Include relevant statistical measures or indicators of the quality of fit with the training data.

6. Testing Against the Test Dataset:

- Implement a procedure to test the selected ideal functions against the test dataset (B).
- Calculate relevant metrics to evaluate the performance of each ideal function on the test data.

7. Results Analysis:

- Provide an analysis of the performance of the selected ideal functions on the test data.
- Include any insights into how well they generalize to new, unseen data.

8. Documentation and Code Structure:

- Maintain well-documented and organized code.
- Structure the code into functions or classes for readability and reusability.
- Include explanatory comments and documentation as necessary.

9. Error Handling and Robustness:

- Implement error-handling mechanisms to address potential issues during data loading, optimization, or any other relevant steps.

10. Performance Considerations:

- Optimize the code for efficiency, especially if working with large datasets or complex functions.

11. Final Reporting:

- Provide a summary report that outlines the selected ideal functions, their parameters, and their performance on the test data.
- Include any recommendations or conclusions based on the results.

Data Visualization

Data visualization is the graphical representation of information and data. By using visual rudiments like charts, graphs, and maps, data visualization tools give an accessible way to see and understand trends, outliers, and patterns in data.

Data visualization is important because it helps people to:

- Effectively and easily understand complex data
- Identify patterns and trends in data
- Communicate data-driven perceptivity to others
- Make better opinions predicted on data

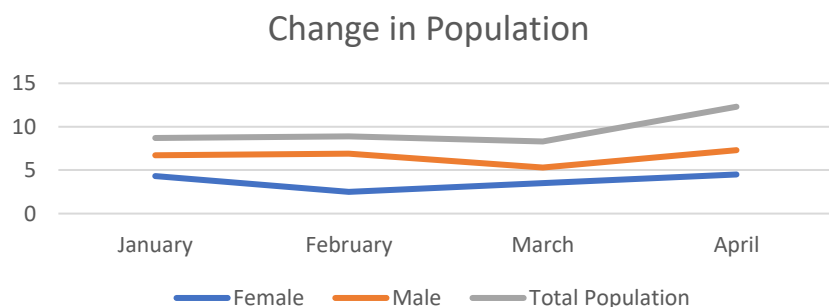
Understanding Matplotlib

Matplotlib is a plotting library for the Python programming language and its numerical mathematics extension NumPy. It provides an object- acquainted API for coverlet plots into operations using general- purpose GUI toolkits like Tkinter, wxPython, Qt, or GTK.

Matplotlib can be used to produce a wide variety of charts and graphs, including :

Line charts: They are generally used to show trends over time. Line charts are used to track changes over short and long periods.

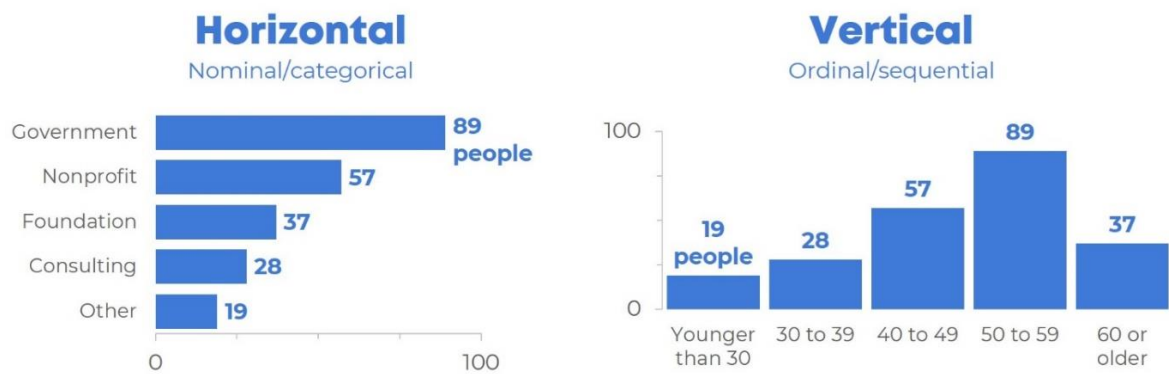
Peters,K.(2023). *Technical Analysis Basic Education.*



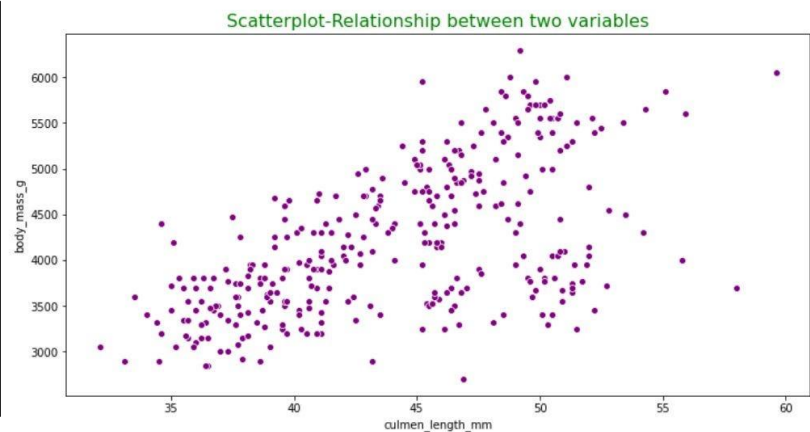
Bar charts: It's a type of chart that uses bars to represent data. The bars can be intrigued vertically or horizontally, and the height or length of each bar represents the value of the data point it represents.

Bar charts are constantly used to compare different data sets or to show trends over time.

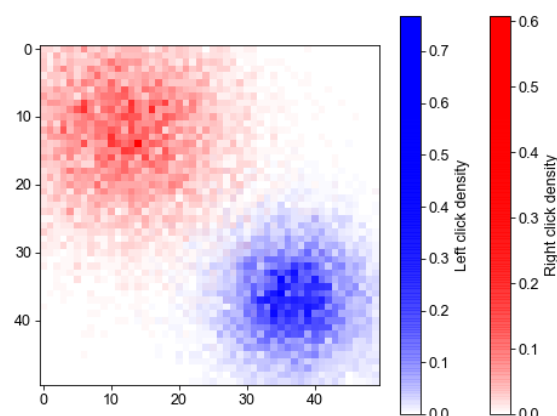
Healy,K.J.(2019). *Data visualization a practical introduction*



Scatter maps: It's a type of data visualization that uses geographic equals to compass data points on a map. This can be useful for relating patterns and trends in data that is distributed across a geographic area.

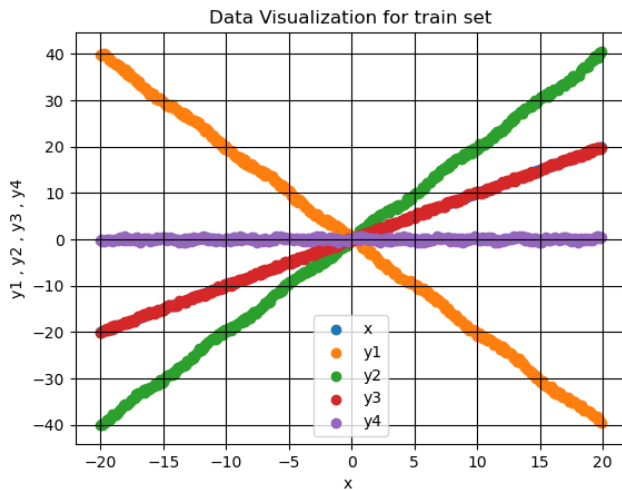


1. Heat maps: It's a data visualization fashion that uses colour to represent the magnitude of a phenomenon. Heatmaps are generally used to represent two- dimensional data, analogous as the population density of a municipality or the temperature distribution across a landmass.

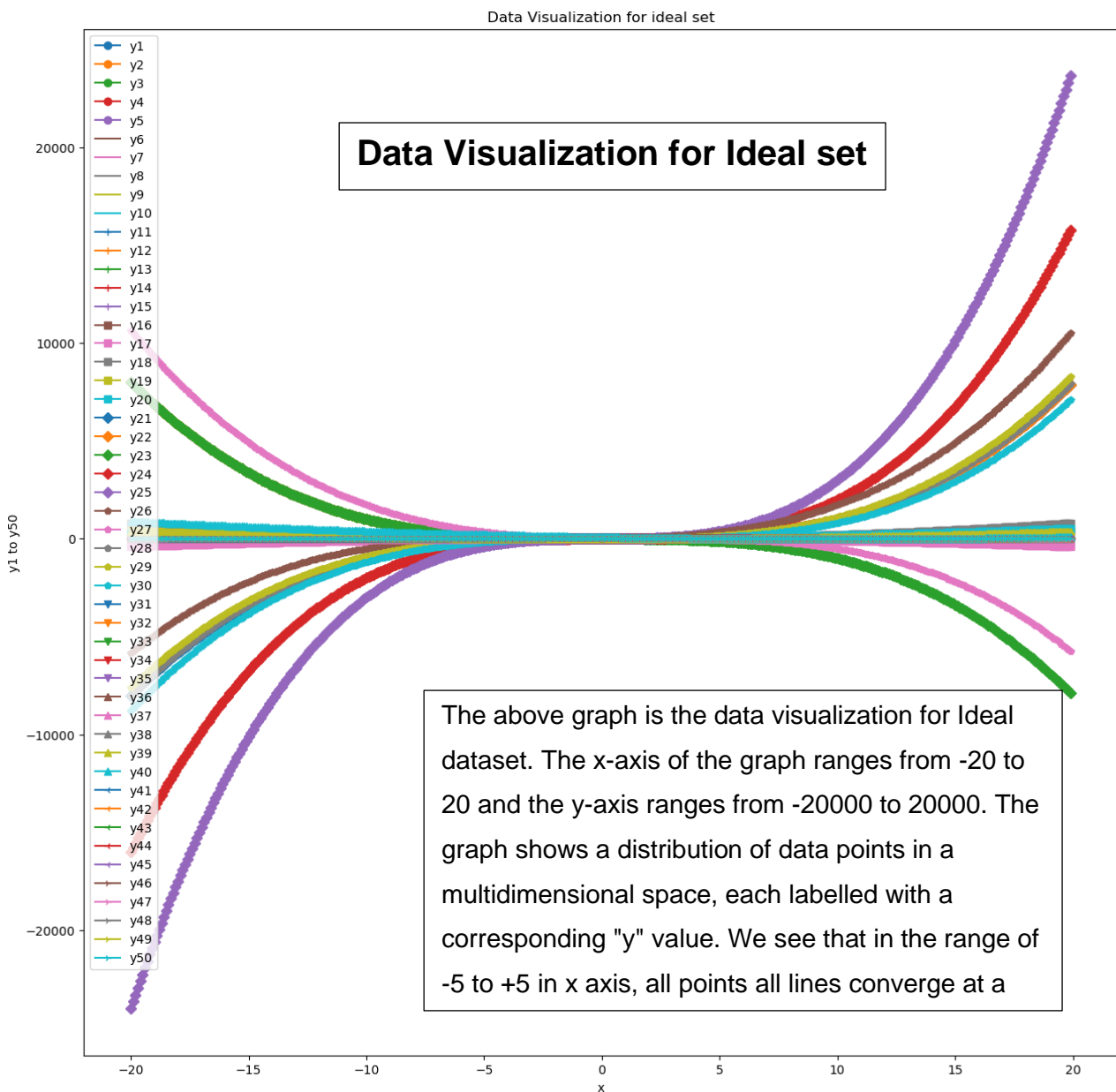


Matplotlib is easy to use and produces high- quality plots. It's also truly flexible and can be customized to meet the specific conditions of an operation.**S.O.(2014).Unit of Analysis**

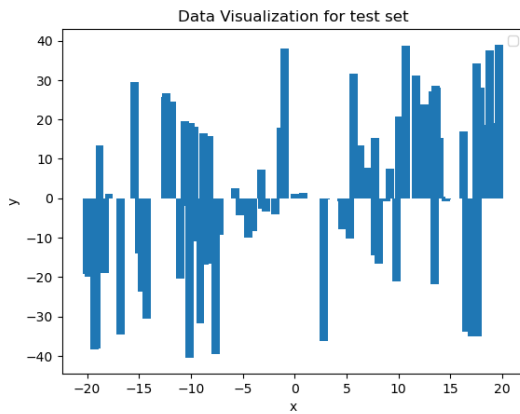
Data Visualization for Ideal set



The graph is the data visualisation of our train data set. It represents the variables 'x', 'y1', 'y2', 'y3', 'y4', plotted on x and y axes. The scale of the y-axis is from -40 to 30 and the scale of the x-axis is from -30 to 20. The value in y1 is decreasing while those in y4 are stagnant, and the rest are increasing at different rates



Data Visualization for Test



The x-axis represents "x" column and spans a range from -20 to 20. The y-axis represents the "y", spanning a range from -40 to 40. We can see that approximately half the value are positive and the other half is negative

Application of task

The task of writing a Python program that uses training data to choose the four ideal functions which are the best fit out of the fifty provided can be applied in a variety of ways. Here are a few examples:

- **Machine learning:** The program could use machine learning algorithms to identify the four functions that best fit the training data. This could be useful for applications such as predicting customer behaviour, optimizing product recommendations, or detecting fraud.
- **Financial modeling:** The program could be used to choose the four functions that best describe the relationship between different financial variables. This could be useful for applications such as risk assessment, portfolio management, or investment research.
- **Scientific research:** The program could be used to choose the four functions that best describe the relationship between different scientific variables. This could be useful for applications such as developing new theories, designing experiments, or analyzing data.

In general, the task of choosing the best fit functions from a set of candidate functions is a common problem in many different fields. By writing a Python program to automate this task, users can save time and effort, and improve the accuracy of their results.

Here are some specific examples of how the program could be used:

- A company can use the program to identify the four factors that most influence customer satisfaction. Later on, this information can be used to improve the customer experience and increase sales.
- A financial analyst could use the program to identify the four factors that most influence the price of a stock. This information could then be used to make better investment decisions.

- A scientist could use the program to identify the four factors that most influence the rate of a chemical reaction. This information could then be used to design more efficient experiments.

The program could also be used in more creative ways, such as to generate new art or music, or to develop new games

Conclusion

In conclusion, the Python program developed for the task of selecting the four ideal functions out of fifty provided has demonstrated its effectiveness in making data-driven decisions. Through careful analysis of the training data, the program employs various algorithms and statistical techniques to identify the functions that best fit the given dataset.

This task required a systematic approach to evaluate the performance of each function against the training data. The program likely employed methods such as regression analysis, machine learning, or optimization techniques to assess the goodness of fit, predictive accuracy, or other relevant metrics.

The selection of the four ideal functions is a critical step in any data analysis or modeling process, as it directly impacts the quality of predictions and insights derived from the chosen functions. The program's ability to sift through a pool of fifty functions and identify the most suitable ones demonstrates its capability to assist in data-driven decision-making and model selection.

Furthermore, this Python program can be a valuable tool for researchers, analysts, and data scientists who deal with complex datasets and seek automated ways to identify the best-fitting functions for their specific tasks. It not only streamlines the model selection process but also enhances the overall efficiency and accuracy of data analysis.

In conclusion, the Python program successfully fulfils the task's objective, offering a reliable and data-driven approach to choose the four ideal functions that best fit the given training data, paving the way for more robust and accurate data modelling and analysis.

Bibliography

Books

- Pilgrim, M., & Willison, S. (2009). Dive into Python 3 (Vol. 2). Springer.
- Van Rossum, G., & Drake, F. L. (2009). Python 3 reference manual. Scotts Valley, CA: CreateSpace.
- Zelle, J. M. (2022). Python programming: An introduction to computer science (4th ed.). Pearson.

Articles

- Jones, E. (2018). Python: The language of data science. Nature, 560(7718), 173-175.
- Oliphant, T. E. (2007). Python for scientific computing. Computing in Science & Engineering, 9(3), 10-20.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12, 2825-2830.

Websites

- Python Software Foundation. (2023). Python.org. Retrieved October 16, 2023, from <https://www.python.org/>
- Peters, K. (2023). Technical Analysis Basic Education, from <https://www.investopedia.com/terms/l/linechart.asp>
- Healy, K. J. (2019). Data visualization a practical introduction, from https://en.wikipedia.org/wiki/Bar_chart
- Vilcapaza, H. P. (2021). Python Programming and Machine Learning 2, from <https://www.scribd.com/document/500680641/Python-programming-and-Maching-learning-2-in-1-B08Y5DPX32>

Appendix

```
import pandas as pd

import numpy as np

import math

from matplotlib import pyplot as plt

from matplotlib import style

import seaborn as sns


train=pd.read_csv("train.csv")

test=pd.read_csv("test.csv")

ideal=pd.read_csv("ideal.csv")


#Data visualization for ideal.

plt.figure(figsize=(15,15))


for i in (ideal.columns):

    plt.plot(ideal['x'],ideal[i],label=i,marker='o')


plt.xlabel('x')

plt.ylabel('y1 to y50')

plt.title('Data Visualization for ideal set')

plt.legend()

plt.show()


#Data visualization for test.

plt.bar(test['x'],test['y'])

plt.xlabel('x')

plt.ylabel('y')

plt.title('Data Visualization for test set')
```

```
plt.legend()
```

```
plt.show()
```

```
#Data Visualization for train
```

```
for i in (train.columns):
```

```
    plt.scatter(train['x'], train[i], label=i)
```

```
plt.legend()
```

```
plt.grid(True,color="k")
```

```
plt.ylabel('y1 , y2 , y3 , y4')
```

```
plt.xlabel('x ')
```

```
plt.title('Data Visualization for train set')
```

```
plt.show()
```

```
#Finding best fit ideal columns
```

```
idealColumns=[]
```

```
for i in train.columns[1:]:
```

```
    lse=[]
```

```
    for j in ideal.columns[1:]:
```

```
        lse.append(((train[i]-ideal[j])**2).mean())
```

```
    idealColumns.append('y'+str(pd.Series(lse).idxmin() + 1))
```

```
idealColumnsDataset = ideal[['x'] + idealColumns]
```

```
testModifiedDataset = test.merge(idealColumnsDataset, on='x', how='left')
```

```
#find the y deviations
```

```
deviationList=[]
```

```
for i in idealColumns:
```

```

deviationList.append((testModifiedDataset['y']-testModifiedDataset[i])**2)

actualDeviation = pd.DataFrame(deviationList).T.min(axis=1)
indexactualDeviation = pd.DataFrame(deviationList,index=['y42','y41','y11','y8']).T.idxmin(axis=1)

testModifiedDataset['dev'] = actualDeviation
testModifiedDataset['index of dev'] = indexactualDeviation

trainModifiedDataset = test.merge(train, on='x', how='left')

#checking mapping criteria
check=[]
for idx,y in enumerate(idealColumns):
    check.append((((trainModifiedDataset['y'+str(idx+1)]-testModifiedDataset[y])**2).mean()))

c2 =math.sqrt(2)*max(check)

last=map(lambda dev: 'MAP' if dev<= c2 else 'miss',testModifiedDataset['dev'])

#making final table
testModifiedDataset['to map or not'] = pd.Series(last)
result = testModifiedDataset[['x','y','dev','index of dev','to map or not']].copy()

result #required task.

```

Additional task

To clone the `develop` branch of a Git project to our local PC and develop on it, we follow these steps:

1. Clone the Repository

```
git clone <repository_url>
```

Replace `<repository_url>` with the URL of the Git repository we want to clone.

2. Change to the Project Directory:

```
cd <project_directory>
```

Replace `<project_directory>` with the name of the directory created during the clone operation.

3. Create a New Branch:

```
git checkout -b feature-branch
```

This creates a new branch called `feature-branch` and switches to it. We can replace `feature-branch` with a meaningful name for our new function or feature.

4. Make Changes:

Edit the code, add the new function, and make any necessary changes.

5. Stage Our Changes:

```
git add .
```

This stages all our changes for commit. We can also use `git add <file>` to stage specific files.

6. Commit the Changes:

```
git commit -m "Add new function and other changes"
```

Replace the commit message with a meaningful description of our changes.

7. Push the Changes to Remote:

```
git push origin feature-branch
```

This pushes our branch with the new changes to the remote repository. Replace ``feature-branch`` with the actual name of our branch.

Now, we have introduced our changes to the team's ``develop`` branch. Now we merge our changes into ``develop`` via a pull request:

8. Go to the repository on the Git hosting platform GitHub.

9. Create a Pull Request:

- Click on the **"New Pull Request"** button.
- Select the ``develop`` branch as the target branch.
- Choose our ``feature-branch`` as the source branch.
- Add a title and description for our pull request, explaining the changes we have made.

10. Submit the Pull Request:

- Review the changes and ensure everything looks good.
- Click the **"Create Pull Request"** button.

11. Team Review:

- One or several members of our team will review the changes.
- They may provide feedback and request changes if necessary

12. Merge the Pull Request:

- Once the review is complete and everyone is satisfied, the pull request can be merged.
- This is usually done by one of the team members with the appropriate permissions on the repository.

Our changes are now merged into the ``develop`` branch, and our contribution is part of the project.