

# **Introduction to Data Science**

Ethan Meyers

# Table of contents

<b>Welcome</b>	<b>3</b>
<b>1 Introduction</b>	<b>4</b>
1.1 What is Data Science? . . . . .	4
1.2 A brief history of Data Science . . . . .	5
1.2.1 A brief history of data . . . . .	5
1.2.2 A brief history of Statistics . . . . .	5
1.2.3 A brief history of computation . . . . .	5
1.2.4 The creation of the field of Data Science . . . . .	5
<b>2 Python basics</b>	<b>6</b>
2.1 Expressions and Statements . . . . .	6
2.1.1 Expressions . . . . .	6
2.1.2 Syntax . . . . .	7
2.1.3 Statements . . . . .	7
2.1.4 The difference between expressions and statements . . . . .	9
2.2 Data types . . . . .	9
2.2.1 Numbers . . . . .	9
2.2.2 Character strings . . . . .	9
2.2.3 Booleans . . . . .	9
2.3 Functions . . . . .	9
<b>3 Array computations</b>	<b>10</b>
3.1 Data & Methods . . . . .	11
3.2 Conclusion . . . . .	11
References . . . . .	11
<b>4 Data tables</b>	<b>13</b>
4.1 Data & Methods . . . . .	14
4.2 Conclusion . . . . .	14
References . . . . .	14
<b>5 Data visualization</b>	<b>16</b>
5.1 Data & Methods . . . . .	17
5.2 Conclusion . . . . .	17
References . . . . .	17

# Welcome



This book gives an introduction to Data Science using the Python programming language.

# 1 Introduction

In this chapter we will discuss what the field of Data Science is, and give a brief history of how the field developed.

This book is your guide to understanding the exciting and increasingly influential field of data science. Whether you're curious about how data shapes our world or are looking to explore the possibilities of data-driven insights, this book will provide you with a foundational understanding of what data science is and why it matters.

## 1.1 What is Data Science?

Data science is a dynamic and interdisciplinary field that combines techniques and theories from statistics, computer science, and specialized knowledge in various areas to extract valuable knowledge and insights from data [Chapter 1]. This data can come in many forms, whether neatly organized in databases or existing as unstructured information like text or images.

At its core, data science follows a systematic process for analyzing data. This includes a range of crucial steps, starting with data collection and ensuring the data is in a usable state through data cleaning. Once prepared, the data is explored to uncover initial patterns and relationships (data exploration). Data scientists then apply various modeling techniques to identify deeper insights, which need to be carefully interpreted to draw meaningful conclusions. Finally, the findings are communicated effectively to inform decisions and understanding [Chapter 1].

The field of data science has experienced remarkable growth in recent years. This surge in prominence can be attributed to several key factors: - The explosion in the amount of data being generated across all sectors, from social media to scientific research. - Significant advancements in computing power, enabling the processing and analysis of these vast datasets. - The development of increasingly sophisticated analytical tools and techniques that allow for more complex and insightful data exploration.

By delving into data science, you can gain practical analytical skills that are applicable across a wide array of fields [Chapter 1, 62]. You'll learn how to approach real-world data, identify key questions, and use data-driven methods to find answers and understand the world around us [Chapter 1, 62]. As a lighthearted starting point, you might hear the quip that "A Data Scientist is a Statistician who lives in San Francisco" [Chapter 1, 11]. While humorous, this

simple definition hints at the combination of statistical thinking with the technological innovation often associated with data science. Throughout this book, we will move beyond simplistic definitions to explore the rich and multifaceted nature of this vital field.

Key points - Despite the fact that humans have been collecting data for millenia, and doing sophisticated analyses of data for centuries, the field of data science” (or at least the name) is relatively new. -

## **1.2 A brief history of Data Science**

### **1.2.1 A brief history of data**

### **1.2.2 A brief history of Statistics**

### **1.2.3 A brief history of computation**

Computational devices also have a long history.

### **1.2.4 The creation of the field of Data Science**

## 2 Python basics

Now that we have discussed what data science is, let's learn some of the basics of the Python programming language so we can begin to learn how to analyze data.

In this chapter we will discuss: - The basic syntax of programming in Python that is needed to write Python code that can run. - Common type of data, such as numbers and character strings, that we frequently use. - Key data structures that we can use to store data.

### 2.1 Expressions and Statements

#### 2.1.1 Expressions

A **Python expression** is **any piece of code that produces a value..** For example, the following is an expression that simply creates the number 21.

```
21
```

21

Similarly, an expression could be a series of mathematical operations that evaluate to number. For example, if want want to add 5 plus 2 and then multiple the result by 6 we can write:

```
6 * (5 + 2)
```

42

As mentioned above, the defining features of a *python expression* is that it produces a value. Expressions are one of the fundamental building blocks of data analysis and they will appear frequently throughout this book.

### 2.1.2 Syntax

**Syntax** is the set of rules that defines how Python code **must** be written. One that think of syntax as the grammar of the Python programming language. In order for Python to be able to run your code, it **must** use the correct syntax. If incorrect syntax is used, then one will get a “syntax error”, and the code will not run.

To illustrate this, let’s calculate the value of 3 squared ( $3^2$ ) which hopefully you remember is equal to the value of 9. In Python, if we want to take a value  $x$  to the power  $y$  (i.e., to calculate  $x^y$ ) we use the syntax `x**y`. So, if we wanted to calculate  $3^2$  we would write the following Python code:

```
3**2
```

```
9
```

Since we have written the correct syntax, the code runs and the result of 9 is calculated as expected.

However, if we accidentally put an extra space between the two `*` symbols, Python will not know how to interpret the expression and we will get a syntax error as shown below:

```
3* *2
```

```
SyntaxError: invalid syntax (3783800731.py, line 1)
```

```
Cell In[4], line 1
```

```
3* *2
```

```
^
```

```
SyntaxError: invalid syntax
```

We there is a syntax error, Python will print out `SyntaxError` and give you an indication where the syntax error has occurred using a `^` symbol. As we can see here, Python is trying to show that the syntax error has occurred due to the extra space between the `*` symbols.

The ability to be able to spot and fix syntax errors is a fundamental skill you will develop as become proficient in analyzing data in Python.

### 2.1.3 Statements

A Python statement is a line of code that performs an action; i.e., it tells Python that it should do something other than just calculating a value.

One of the most fundamental statement types is an *assignment statement*, where we have Python store a value in a named **variable**. For example, the following code assigns the value 10 to the variable `a`:

```
a = 10
```

We can then refer back to the variable `a` later in our code to retrieve the stored value. For example, if we just write `a` on a line by itself, it will print out the value stored in `a`.

```
a
```

```
10
```

As we can see, the value printed out is 10 which is the value we had previously stored in the name `a`.

If we then assign the name `a` to another value, it will overwrite the previously stored value and `a` will store the new value.

```
a = 21  
a
```

```
21
```

We can also do mathematical operations on values stored in variables, such as adding and multiplying variables together. For example, we can assign the variable `h` to store the value 24, and the variable `d` to store the value 7, and then we can multiply these together and store the result in the variable `t`.

```
h = 24  
d = 7  
t = h * d  
t
```

```
168
```

### 2.1.3.1 Variable names

Must start with a character Should be meaningful name (`a` bad variable name).

The first is a law (syntax error). The second is good style to make code readable (like wearing sweatpants to class).



```
hours = 24
days = 7
total_hours_in_a_week = hours * days
total_hours_in_a_week
```

168

**Exercise:** Calculate CT yearly wage if making minimum wage...

So far we have only described “assignment statements”. However, later in this book will see many other types of statements, such as “for loops” and “conditional statements” which us to do more complex processing in Python.

### 2.1.4 The difference between expressions and statements

If Python can evaluate it to a value, it’s an expression. If it’s doing something (like a loop, a function definition, or an assignment), it’s a statement.

So,  $2 + 2$  is an expression, but  $x = 2 + 2$  is a statement that contains an expression.

```
x = 5
```

Term	Description	Example
<b>Expression</b>	Returns a value	$2 + 2$
<b>Statement</b>	Performs an action (e.g., assignment, loop, etc.)	$x = 2 + 2$

## 2.2 Data types

### 2.2.1 Numbers

Ints and floats

### 2.2.2 Character strings

### 2.2.3 Booleans

## 2.3 Functions

### 3 Array computations

This is a book created from markdown and executable code.

```
import matplotlib.pyplot as plt
import numpy as np
eruptions = [1492, 1585, 1646, 1677, 1712, 1949, 1971, 2021]

plt.figure(figsize=(6, 1))
plt.eventplot(eruptions, lineoffsets=0, linelengths=0.1, color='black')
plt.gca().axes.get_yaxis().set_visible(False)
plt.ylabel('')
plt.show()
```

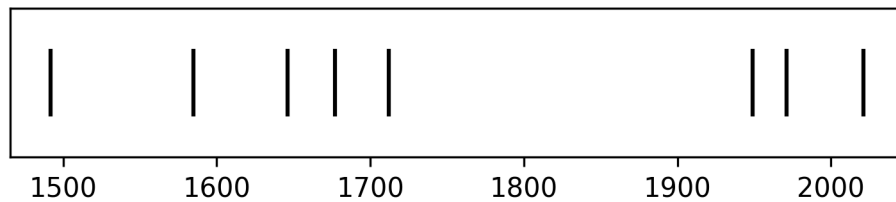


Figure 3.1: Timeline of recent earthquakes on La Palma

```
avg_years_between_eruptions = np.mean(np.diff(eruptions[:-1]))
avg_years_between_eruptions
```

Based on data up to and including 1971, eruptions on La Palma happen every 79.8 years on average.

Studies of the magma systems feeding the volcano, such as Marrero et al. (2019), have proposed that there are two main magma reservoirs feeding the Cumbre Vieja volcano; one in the mantle (30-40km depth) which charges and in turn feeds a shallower crustal reservoir (10-20km depth).

Eight eruptions have been recorded since the late 1400s (Figure 5.1).

Data and methods are discussed in Section 5.1.

Let  $x$  denote the number of eruptions in a year. Then,  $x$  can be modeled by a Poisson distribution

$$p(x) = \frac{e^{-\lambda} \lambda^x}{x!} \quad (3.1)$$

where  $\lambda$  is the rate of eruptions per year. Using Equation 5.1, the probability of an eruption in the next  $t$  years can be calculated.

Table 3.1: Recent historic eruptions on La Palma

Name	Year
Current	2021
Teneguía	1971
Nambroque	1949
El Charco	1712
Volcán San Antonio	1677
Volcán San Martin	1646
Tajuya near El Paso	1585
Montaña Quemada	1492

Table 5.1 summarises the eruptions recorded since the colonization of the islands by Europeans in the late 1400s.

Figure 5.2 shows the location of recent Earthquakes on La Palma.

## 3.1 Data & Methods

## 3.2 Conclusion

## References

Marrero, José, Alicia García, Manuel Berrocoso, Ángeles Llinares, Antonio Rodríguez-Losada, and R. Ortiz. 2019. “Strategies for the Development of Volcanic Hazard Maps in Monogenetic Volcanic Fields: The Example of La Palma (Canary Islands).” *Journal of Applied Volcanology* 8 (July). <https://doi.org/10.1186/s13617-019-0085-5>.

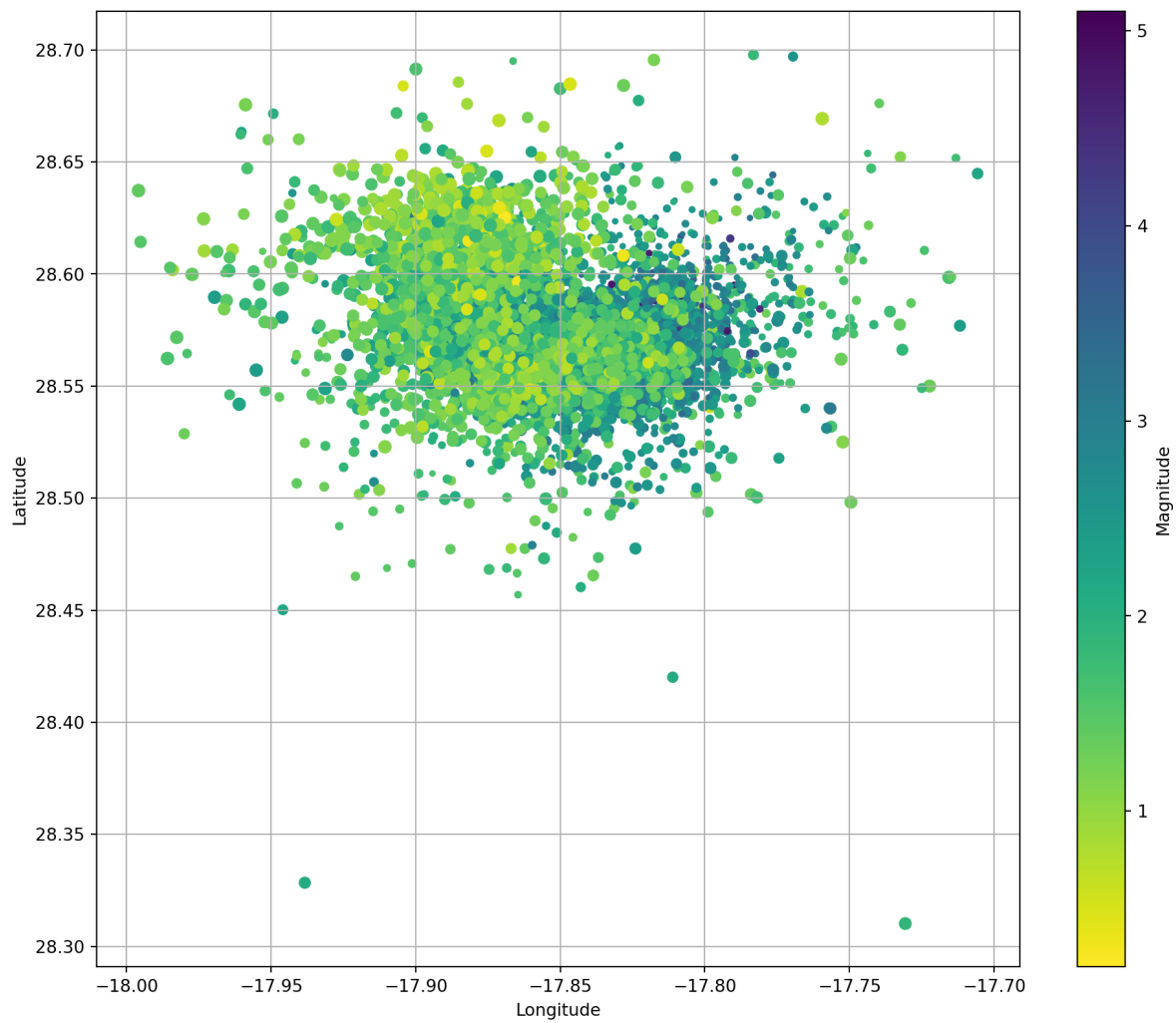


Figure 3.2: Locations of earthquakes on La Palma since 2017.

## 4 Data tables

This is a book created from markdown and executable code.

```
import matplotlib.pyplot as plt
import numpy as np
eruptions = [1492, 1585, 1646, 1677, 1712, 1949, 1971, 2021]

plt.figure(figsize=(6, 1))
plt.eventplot(eruptions, lineoffsets=0, linelengths=0.1, color='black')
plt.gca().axes.get_yaxis().set_visible(False)
plt.ylabel('')
plt.show()
```

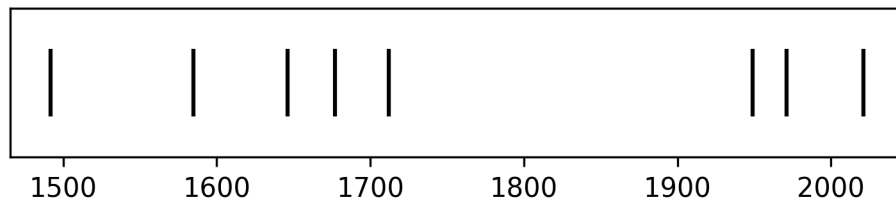


Figure 4.1: Timeline of recent earthquakes on La Palma

```
avg_years_between_eruptions = np.mean(np.diff(eruptions[:-1]))
avg_years_between_eruptions
```

Based on data up to and including 1971, eruptions on La Palma happen every 79.8 years on average.

Studies of the magma systems feeding the volcano, such as Marrero et al. (2019), have proposed that there are two main magma reservoirs feeding the Cumbre Vieja volcano; one in the mantle (30-40km depth) which charges and in turn feeds a shallower crustal reservoir (10-20km depth).

Eight eruptions have been recorded since the late 1400s (Figure 5.1).

Data and methods are discussed in Section 5.1.

Let  $x$  denote the number of eruptions in a year. Then,  $x$  can be modeled by a Poisson distribution

$$p(x) = \frac{e^{-\lambda} \lambda^x}{x!} \quad (4.1)$$

where  $\lambda$  is the rate of eruptions per year. Using Equation 5.1, the probability of an eruption in the next  $t$  years can be calculated.

Table 4.1: Recent historic eruptions on La Palma

Name	Year
Current	2021
Teneguía	1971
Nambroque	1949
El Charco	1712
Volcán San Antonio	1677
Volcán San Martin	1646
Tajuya near El Paso	1585
Montaña Quemada	1492

Table 5.1 summarises the eruptions recorded since the colonization of the islands by Europeans in the late 1400s.

Figure 5.2 shows the location of recent Earthquakes on La Palma.

## 4.1 Data & Methods

## 4.2 Conclusion

## References

Marrero, José, Alicia García, Manuel Berrocoso, Ángeles Llinares, Antonio Rodríguez-Losada, and R. Ortiz. 2019. “Strategies for the Development of Volcanic Hazard Maps in Monogenetic Volcanic Fields: The Example of La Palma (Canary Islands).” *Journal of Applied Volcanology* 8 (July). <https://doi.org/10.1186/s13617-019-0085-5>.

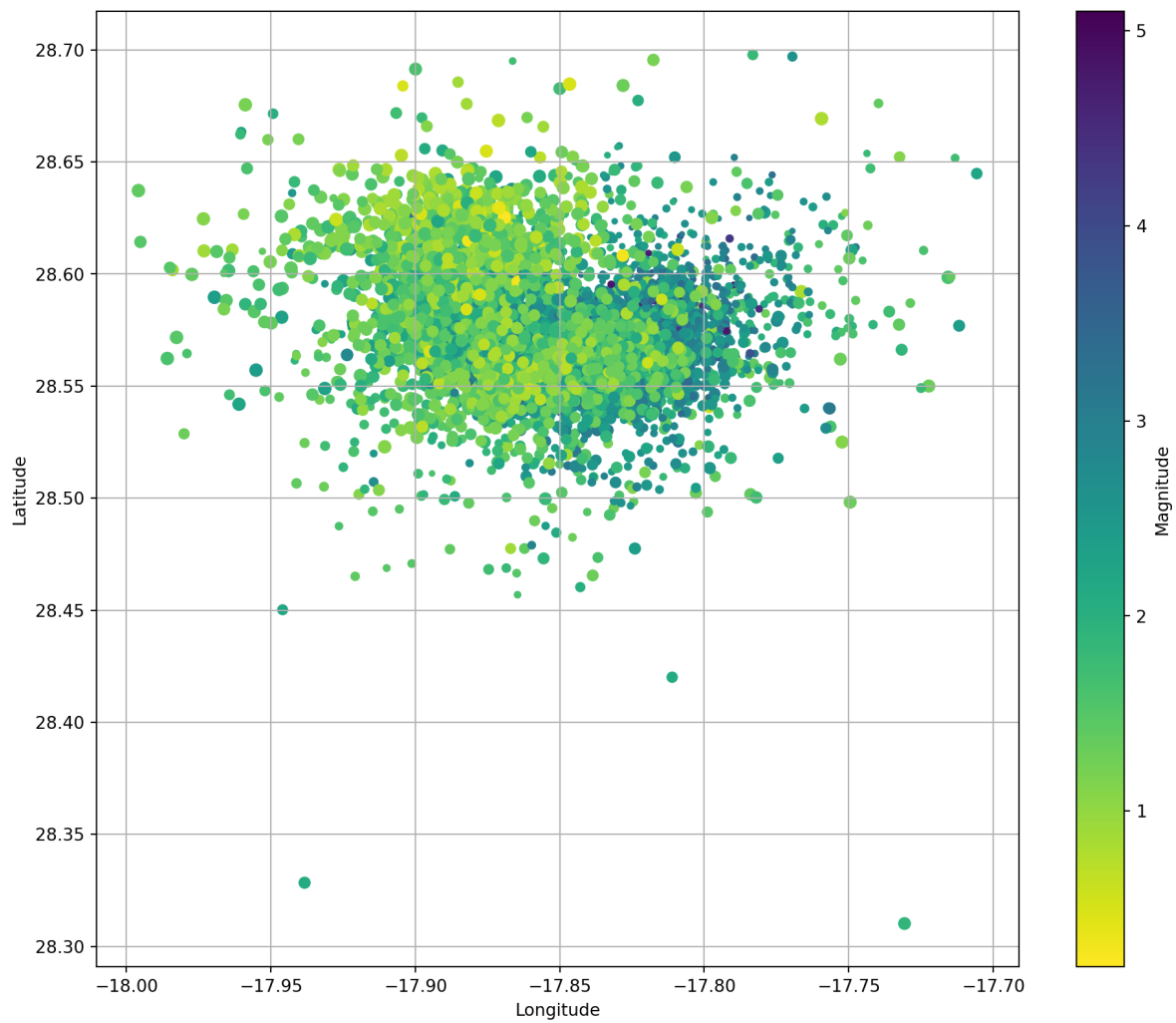


Figure 4.2: Locations of earthquakes on La Palma since 2017.

## 5 Data visualization

This is a book created from markdown and executable code.

```
import matplotlib.pyplot as plt
import numpy as np
eruptions = [1492, 1585, 1646, 1677, 1712, 1949, 1971, 2021]

plt.figure(figsize=(6, 1))
plt.eventplot(eruptions, lineoffsets=0, linelengths=0.1, color='black')
plt.gca().axes.get_yaxis().set_visible(False)
plt.ylabel('')
plt.show()
```

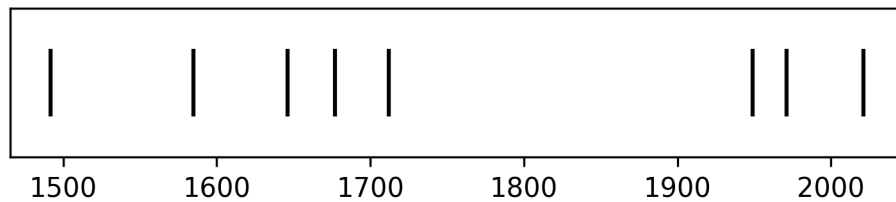


Figure 5.1: Timeline of recent earthquakes on La Palma

```
avg_years_between_eruptions = np.mean(np.diff(eruptions[:-1]))
avg_years_between_eruptions
```

Based on data up to and including 1971, eruptions on La Palma happen every 79.8 years on average.

Studies of the magma systems feeding the volcano, such as Marrero et al. (2019), have proposed that there are two main magma reservoirs feeding the Cumbre Vieja volcano; one in the mantle (30-40km depth) which charges and in turn feeds a shallower crustal reservoir (10-20km depth).

Eight eruptions have been recorded since the late 1400s (Figure 5.1).



Data and methods are discussed in Section 5.1.

Let  $x$  denote the number of eruptions in a year. Then,  $x$  can be modeled by a Poisson distribution

$$p(x) = \frac{e^{-\lambda} \lambda^x}{x!} \quad (5.1)$$

where  $\lambda$  is the rate of eruptions per year. Using Equation 5.1, the probability of an eruption in the next  $t$  years can be calculated.

Table 5.1: Recent historic eruptions on La Palma

Name	Year
Current	2021
Teneguía	1971
Nambroque	1949
El Charco	1712
Volcán San Antonio	1677
Volcán San Martin	1646
Tajuya near El Paso	1585
Montaña Quemada	1492

Table 5.1 summarises the eruptions recorded since the colonization of the islands by Europeans in the late 1400s.

Figure 5.2 shows the location of recent Earthquakes on La Palma.

## 5.1 Data & Methods

## 5.2 Conclusion

## References

Marrero, José, Alicia García, Manuel Berrocoso, Ángeles Llinares, Antonio Rodríguez-Losada, and R. Ortiz. 2019. “Strategies for the Development of Volcanic Hazard Maps in Monogenetic Volcanic Fields: The Example of La Palma (Canary Islands).” *Journal of Applied Volcanology* 8 (July). <https://doi.org/10.1186/s13617-019-0085-5>.

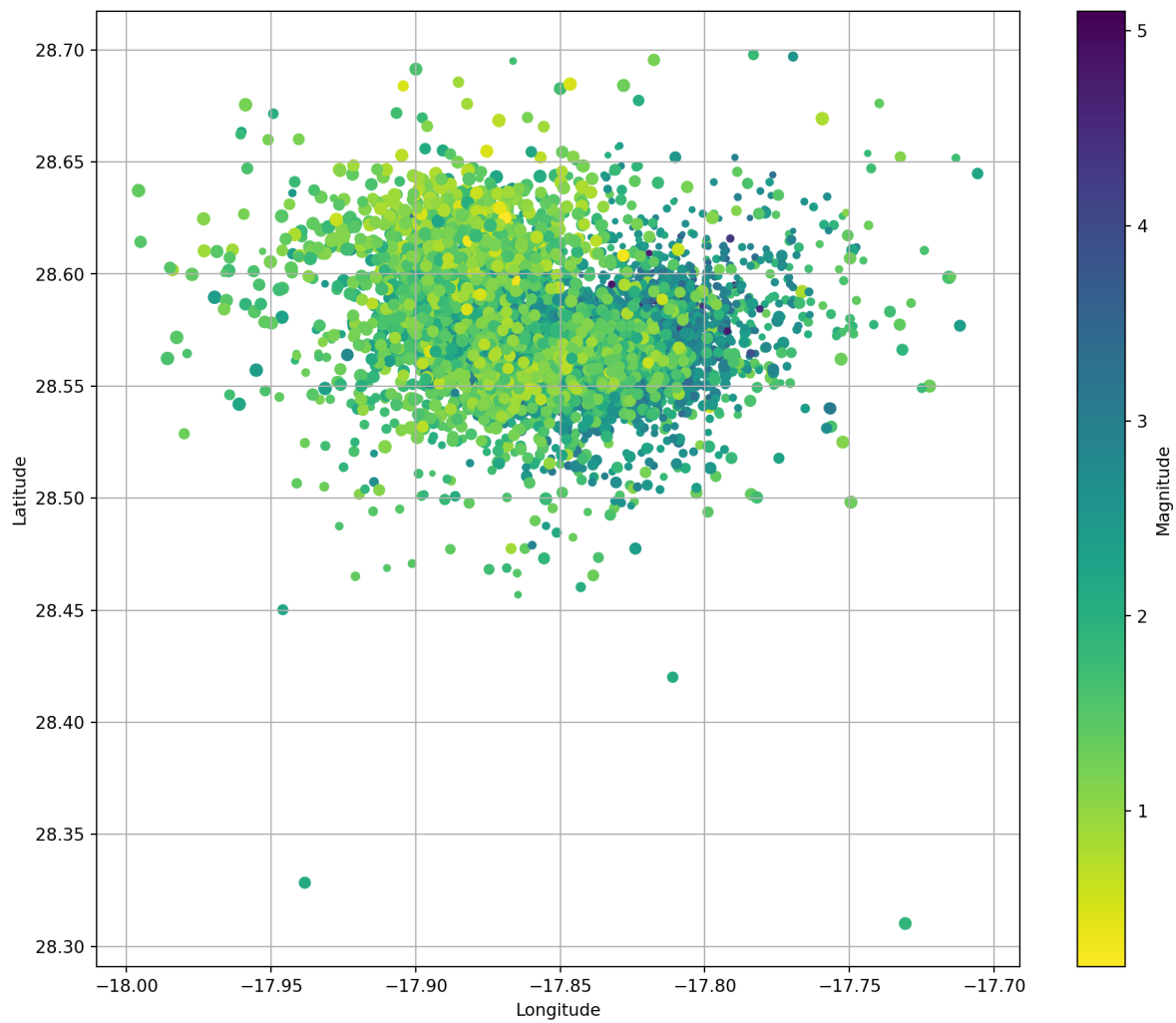


Figure 5.2: Locations of earthquakes on La Palma since 2017.