

Jupyter Data\_visualisation1 Last Checkpoint: 09/06/2025 (autosaved)

In [2]:

```
import numpy as np
import pandas as pd
import seaborn as sns

pd.set_option("display.max_rows", 20)
```

## Load Data

In [3]:

```
penguin = pd.read_csv(r"C:\Users\Srilakshmi N Murthy\Desktop\MS 2nd sem\Data visualization\penguins.csv")
```

In [4]:

```
penguin.head(20)
```

Out[4]:

	studyName	Sample Number	Species	Region	Island	Stage	Individual ID	Clutch Completion	Date Egg	Culmen Length (mm)	Culmen Depth (mm)	Flipper Length (mm)	Body Mass (g)	Sex	Delta 15 N (o/oo)	Delta 13 C (o/oo)
0	PAL0708	1	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A1	Yes	11/11/07	39.1	18.7	181.0	3750.0	MALE	NaN	NaN
1	PAL0708	2	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A2	Yes	11/11/07	39.5	17.4	186.0	3800.0	FEMALE	8.94956	-24.69454

23:47 23-09-2025

Jupyter Data\_visualisation1 Last Checkpoint: 09/06/2025 (autosaved)

In [5]:

```
print("Rows, columns", penguin.shape)
```

Rows, columns (344, 17)

In [6]:

```
penguin.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 344 entries, 0 to 343
Data columns (total 17 columns):
 #   Column           Non-Null Count  Dtype  
 --- 
 0   studyName        344 non-null    object  
 1   Sample Number    344 non-null    int64  
 2   Species          344 non-null    object  
 3   Region           344 non-null    object  
 4   Island            344 non-null    object  
 5   Stage             344 non-null    object  
 6   Individual ID    344 non-null    object  
 7   Clutch Completion 344 non-null    object  
 8   Date Egg          344 non-null    object  
 9   Culmen Length (mm) 342 non-null    float64 
 10  Culmen Depth (mm) 342 non-null    float64 
 11  Flipper Length (mm) 342 non-null    float64 
 12  Body Mass (g)     342 non-null    float64 
 13  Sex               344 non-null    object  
 14  Delta 15 N (o/oo) 334 non-null    float64 
 15  Delta 13 C (o/oo) 331 non-null    float64
```

23:48 23-09-2025

```
In [7]: penguin.isna().mean().sort_values(ascending=False).head(10).to_frame("missing_column")
```

```
Out[7]:
```

	missing_column
Comments	0.924419
Delta 15 N (o/oo)	0.040698
Delta 13 C (o/oo)	0.037791
Sex	0.029070
Culmen Length (mm)	0.005814
Body Mass (g)	0.005814
Flipper Length (mm)	0.005814
Culmen Depth (mm)	0.005814
studyName	0.000000
Sample Number	0.000000

## FIX DTYPES

```
In [8]: penguin.dtypes.head(18)
```

```
Out[8]:
```

	dtype
studyName	object
Sample Number	int64
Species	object
Region	object
Island	object
Stage	object
Individual ID	object
Clutch Completion	object
Date Egg	object
Culmen Length (mm)	float64
Culmen Depth (mm)	float64
Flipper Length (mm)	float64
Body Mass (g)	float64
Sex	object
Delta 15 N (o/oo)	float64
Delta 13 C (o/oo)	float64
Comments	object
dtype	object

```
In [9]: #converting strings to categories
```

```
cat_cols = ("studyName", "Species", "Region", "Island", "Stage", "Individual ID", "Clutch Completion", "Sex", "Comments")
for x in cat_cols:
    if x in penguin:
        penguin[x] = penguin[x].astype("category")
```

localhost:8889/notebooks/Data\_visualisation1.ipynb

jupyter Data\_visualisation1 Last Checkpoint: 09/06/2025 (autosaved)

In [10]: penguin.dtypes.head(17)

```
Out[10]: studyName      category
Sample Number    int64
Species         category
Region          category
Island          category
Stage           category
Individual ID   category
Clutch Completion category
Date Egg        object
Culmen Length (mm) float64
Culmen Depth (mm) float64
Flipper Length (mm) float64
Body Mass (g)   float64
Sex             category
Delta 15 N (o/o) float64
Delta 13 C (o/o) float64
Comments        category
dtype: object
```

## handling duplicates

In [11]: penguin.shape

```
Out[11]: (344, 17)
```

In [12]: penguin= penguin.drop\_duplicates()
print("After dropping duplicates", penguin.shape)

Type here to search 23°C 23:50 23-09-2025

localhost:8889/notebooks/Data\_visualisation1.ipynb

jupyter Data\_visualisation1 Last Checkpoint: 09/06/2025 (autosaved)

In [11]: penguin.shape

```
Out[11]: (344, 17)
```

In [12]: penguin= penguin.drop\_duplicates()
print("After dropping duplicates", penguin.shape)

After dropping duplicates (344, 17)

In [13]: penguin.duplicated().sum() #no duplicates

```
Out[13]: 0
```

## Missing values

In [14]: missing\_column = penguin.isna().mean()
mostly\_missing = missing\_column[missing\_column > 0.6].index.tolist()
penguin\_clean = penguin.drop(columns = mostly\_missing)
print("Dropped mostly missing values",mostly\_missing)

```
Dropped mostly missing values ['Comments']
```

In [15]: num\_cols = penguin\_clean.select\_dtypes(include = "number").columns
cat\_cols = penguin\_clean.select\_dtypes(include = "category").columns

penguin\_clean[num\_cols]=penguin\_clean[num\_cols].fillna(penguin\_clean[num\_cols].median())

for c in cat\_cols:
 penguin\_clean[c]=penguin\_clean[c].bfill()

Type here to search 23°C 23:50 23-09-2025

A screenshot of a Jupyter Notebook interface running on a Windows desktop. The notebook is titled "Data\_visualisation1" and shows a single cell of Python code. The code prints a message about missing values and then performs data cleaning on a dataset named "penguin\_clean". The output shows that all columns have zero missing values.

```
print("Dropped mostly missing values",mostly_missing)
Dropped mostly missing values ['Comments']

In [15]: num_cols = penguin_clean.select_dtypes(include = "number").columns
cat_cols = penguin_clean.select_dtypes(include = "category").columns

for c in cat_cols:
    penguin_clean=penguin_clean.fillna(penguin_clean[num_cols].median())

penguin_clean.isna().sum()

Out[15]:
```

	studyName	Sample Number	Species	Region	Island	Stage	Individual ID	Clutch Completion	Date Egg	Culmen Length (mm)	Culmen Depth (mm)	Flipper Length (mm)	Body Mass (g)	Sex	Delta 15 N (o/oo)	Delta 13 C (o/oo)	dtype: int64
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

A screenshot of a Jupyter Notebook interface running on a Windows desktop. The notebook is titled "Data\_visualisation1" and shows a cell of Python code for text cleaning. The code iterates through categorical columns and applies normalization steps like stripping whitespace and replacing non-breaking spaces with regular spaces. The output shows the result of the cleaning process.

```
In [16]: penguin_clean.isna().sum().sum()
Out[16]: 0

Every listed column has 0 missing values.
#So your penguin dataset is already perfectly clean in terms of NaNs.

Standardize text by lowercasing, removing unwanted characters, trimming whitespace, and applying consistent casing. This avoids problems like treating "Southampton", "southampton", and " Southampton" as different categories.



### Text cleaning



In [17]: for col in cat_cols:
    if col in penguin_clean.columns:
        penguin_clean[col + ".norm"] = (
            penguin_clean[col]
            .astype(str)
            .str.strip()
            .str.replace(r"\s+", " ", regex=True)
            .str.title()
        )
```

localhost:8889/notebooks/Data\_visualisation1.ipynb

jupyter Data\_visualisation1 Last Checkpoint: 09/06/2025 (autosaved)

In [18]:

```
penguin_clean = penguin_clean.rename(columns={  
    "Culmen Length (mm)": "Culmen_Length_mm",  
    "Culmen Depth (mm)": "Culmen_Depth_mm",  
    "Flipper Length (mm)": "Flipper_Length_mm",  
    "Body Mass (g)": "Body_Mass_g",  
    "Delta 15 N (o/oo)": "Delta15N_o_per_oo",  
    "Delta 13 C (o/oo)": "Delta13C_o_per_oo"  
})
```

In [19]:

```
cols_to_show = [c for c in penguin_clean.columns]  
penguin_clean[cols_to_show].head(8)
```

Out[19]:

	studyName	Sample Number	Species	Region	Island	Stage	Individual ID	Clutch Completion	Date Egg	Culmen_Length_mm	...	Delta15N_o_per_oo	Delta13C_o_per_oo
0	PAL0708	1	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A1	Yes	11/11/07	39.10	...	8.652405	-25.83352
1	PAL0708	2	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A2	Yes	11/11/07	39.50	...	8.949560	-24.69454
2	PAL0708	3	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A1	Yes	11/16/07	40.30	...	8.368210	-25.33302
3	PAL0708	4	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A2	Yes	11/16/07	44.45	...	8.652405	-25.83352

localhost:8889/notebooks/Data\_visualisation1.ipynb

jupyter Data\_visualisation1 Last Checkpoint: 09/06/2025 (autosaved)

In [20]:

```
p99 = penguin_clean["Body_Mass_g"].quantile(0.99)  
penguin_clean["Body_Mass_g_capped"] = penguin_clean["Body_Mass_g"].clip(lower=p99)  
penguin_clean[["Body_Mass_g", "Body_Mass_g_capped"]].describe()
```

Out[20]:

	Body_Mass_g	Body_Mass_g_capped
count	344.000000	344.000000
mean	4200.872093	4199.604651
std	799.696532	796.664777
min	2700.000000	2700.000000
25%	3550.000000	3550.000000
50%	4050.000000	4050.000000
75%	4750.000000	4750.000000
max	6300.000000	5978.500000

Feature Engineering

In [21]:

```
penguin_clean["Body_Mass_kg"] = penguin_clean["Body_Mass_g"] / 1000
```

localhost:8889/notebooks/Data\_visualisation1.ipynb

jupyter Data\_visualisation1 Last Checkpoint: 09/06/2025 (autosaved)

File Edit View Insert Cell Kernel Help

Kernel error Trusted Python 3.8 (py38)

## Feature Engineering

In [21]:

```
penguin_clean["Body_Mass(kg)"] = penguin_clean["Body_Mass_g"]/1000  
penguin_clean[["Body_Mass(kg)"]]
```

Out[21]:

	Body Mass(kg)
0	3.75
1	3.80
2	3.25
3	4.05
4	3.45
...	...
339	4.05
340	4.85
341	5.75
342	5.20
343	5.40

344 rows × 1 columns

In [22]:

```
penguin_clean["Flipper_per_Culmen"] = penguin_clean["Flipper_Length_mm"]/penguin_clean["Culmen_Depth_mm"]  
penguin_clean[["Flipper_per_Culmen"]]
```

Windows Taskbar:

- Type here to search
- Start button
- File Explorer
- File History
- OneDrive
- Recycle Bin
- Task View
- Google Chrome
- Microsoft Edge
- File Manager
- Visual Studio Code
- PowerShell
- Windows Terminal
- File
- 23°C
- 23:51
- 23-09-2025

localhost:8889/notebooks/Data\_visualisation1.ipynb

jupyter Data\_visualisation1 Last Checkpoint: 09/06/2025 (autosaved)

File Edit View Insert Cell Kernel Help

Kernel error Trusted Python 3.8 (py38)

344 rows × 1 columns

In [22]:

```
penguin_clean["Flipper_per_Culmen"] = penguin_clean["Flipper_Length_mm"]/penguin_clean["Culmen_Depth_mm"]  
penguin_clean[["Flipper_per_Culmen"]]
```

Out[22]:

	Flipper_per_Culmen
0	9.679144
1	10.689655
2	10.833333
3	11.387283
4	10.000000
...	...
339	11.387283
340	15.034965
341	14.140127
342	14.324324
343	13.229814

344 rows × 1 columns

## Groupby and aggregate

Windows Taskbar:

- Type here to search
- Start button
- File Explorer
- File History
- OneDrive
- Recycle Bin
- Task View
- Google Chrome
- Microsoft Edge
- File Manager
- Visual Studio Code
- PowerShell
- Windows Terminal
- File
- 23°C
- 23:53
- 23-09-2025

localhost:8889/notebooks/Data\_visualisation1.ipynb

jupyter Data\_visualisation1 Last Checkpoint: 09/06/2025 (autosaved)

File Edit View Insert Cell Kernel Help

Kernel error Trusted Python 3.8 (py38)

## Groupby and aggregate

Comput average survival rates grouped by class and sex. GroupBy operations are central to data prep, letting us summarize and check trends quickly.

penguin\_clean.groupby("Species") → creates a grouped object (by species).

["Body\_Mass\_g"] → tells pandas which numeric column you want to aggregate (take the mean of).

```
In [23]: #Average body mass per species
Average_body_mass_per_species=(penguin_clean
    .groupby("Species")["Body_Mass_g"]
    .mean()
    .rename("body_mass")
    .reset_index()
    .sort_values(["Species"]))
Average_body_mass_per_species
```

```
Out[23]:
Species      body_mass
0   Adelie Penguin (Pygoscelis adeliae) 3702.960526
1   Chinstrap penguin (Pygoscelis antarctica) 3733.088235
2   Gentoo penguin (Pygoscelis papua) 5067.741935
```

```
In [24]: avg_mass_species_sex = (
    penguin_clean
    .groupby(["Species", "Sex"])["Body_Mass_g"]
    .mean())
avg_mass_species_sex
```

localhost:8889/notebooks/Data\_visualisation1.ipynb

jupyter Data\_visualisation1 Last Checkpoint: 09/06/2025 (autosaved)

File Edit View Insert Cell Kernel Help

Kernel error Trusted Python 3.8 (py38)

```
In [24]: avg_mass_species_sex = (
    penguin_clean
    .groupby(["Species", "Sex"])["Body_Mass_g"]
    .mean()
    .reset_index(name="mean_mass_g"))
avg_mass_species_sex
```

```
Out[24]:
Species      Sex      mean_mass_g
0   Adelie Penguin (Pygoscelis adeliae)  FEMALE  3378.040541
1   Adelie Penguin (Pygoscelis adeliae)  MALE    4011.217949
2   Adelie Penguin (Pygoscelis adeliae)  .        NaN
3   Chinstrap penguin (Pygoscelis antarctica)  FEMALE  3527.205882
4   Chinstrap penguin (Pygoscelis antarctica)  MALE    3938.970588
5   Chinstrap penguin (Pygoscelis antarctica)  .        NaN
6   Gentoo penguin (Pygoscelis papua)  FEMALE  4875.000000
7   Gentoo penguin (Pygoscelis papua)  MALE    4669.067797
8   Gentoo penguin (Pygoscelis papua)  .        5438.281250
```

```
#for my reference
#A pivot table is a way to summarize data in a grid by grouping rows and columns.
#Key difference
```

**Pivot tables**

```
In [ ]:
```

```
In [25]: pivot_table = avg_mass_species_sex.pivot(
    index="Species",
    columns="Sex",
    values="mean_mass_B"
).round(3)

pivot_table
```

```
Out[25]:
```

Species	SEX	FEMALE	MALE
Adelie Penguin (Pygoscelis adeliae)	NaN	3378.041	4011.218
Chinstrap penguin (Pygoscelis antarctica)	NaN	3527.206	3938.971
Gentoo penguin (Pygoscelis papua)	4875.0	4669.068	5438.281

**joins**

```
In [26]: avg_mass = penguin_clean.groupby("Species")["Body_Mass_g"].mean().reset_index(name ="avg_mass")
avg_flipper=penguin_clean.groupby("Species")["flipper_length_mm"].mean().reset_index(name="avg_flipper")
merged = pd.merge(avg_mass, avg_flipper, on = "Species", how = "left")
merged
```

```
Out[26]:
```

Species	avg_mass	avg_flipper
Adelie Penguin (Pygoscelis adeliae)	3702.960526	190.000000
Chinstrap penguin (Pygoscelis antarctica)	3733.088235	195.823529
Gentoo penguin (Pygoscelis papua)	5067.741935	217.024194

**Time series reshape**

```
In [27]: wide = penguin_clean.pivot_table(index = "Date_Egg", columns = "Species", values= "Body_Mass_g")
long_again = wide.reset_index().melt(id_vars="Date_Egg", var_name = "Species", value_name= "Body_Mass_g")
long_again.head()
```

```
Out[27]:
```

Date Egg	Species	Body_Mass_g
11/10/07	Adelie Penguin (Pygoscelis adeliae)	3675.0
11/10/08	Adelie Penguin (Pygoscelis adeliae)	3500.0
11/10/09	Adelie Penguin (Pygoscelis adeliae)	3987.5
11/11/07	Adelie Penguin (Pygoscelis adeliae)	3775.0

Blind 75 | (13) Search | vanshb0 | Inbox (2) | Reminders | Home Page | Data\_vis | Files | React bc | +

localhost:8889/notebooks/Data\_visualisation1.ipynb

jupyter Data\_visualisation1 Last Checkpoint: 09/06/2025 (autosaved)

File Edit View Insert Cell Kernel Help

Kernel error Trusted Python 3.8 (py38)

## Validation checks

```
In [28]: assert penguin_clean["Body_Mass_g"].ge(0).all()
assert penguin_clean["Flipper Length_mm"].ge(0).all()
print("Basic checks passed ✅")
```

Basic checks passed ✅

## Plotting

### Line Chart

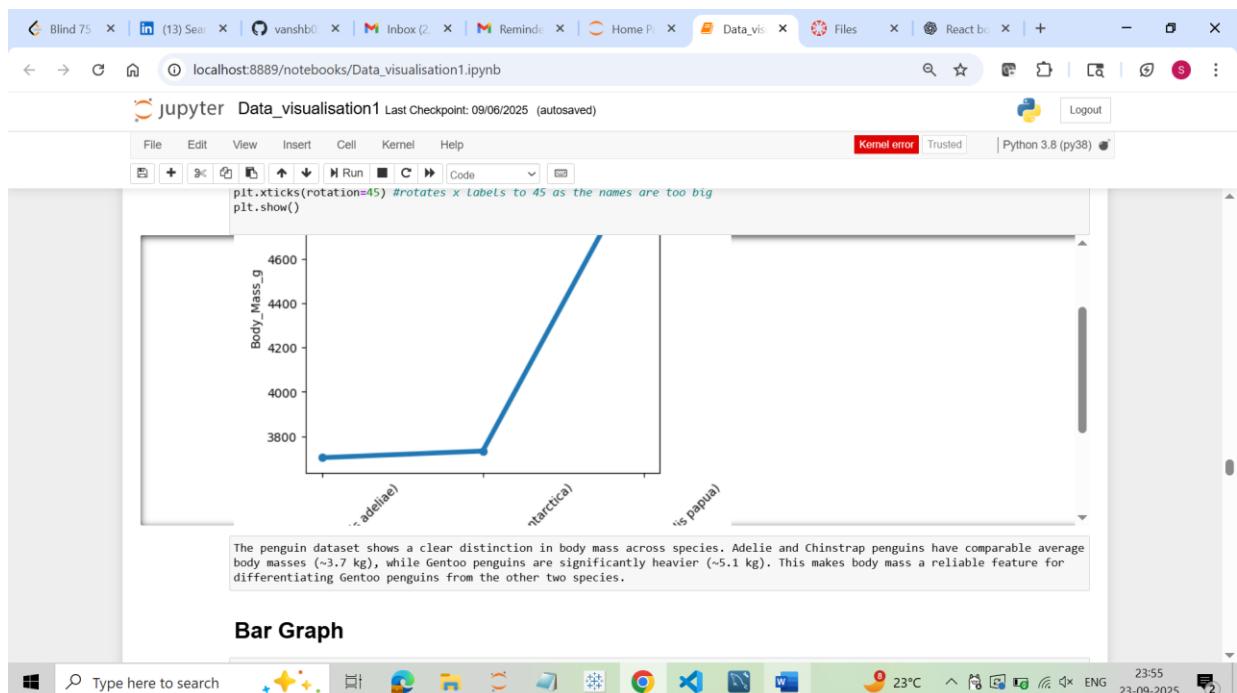
```
In [29]: import matplotlib.pyplot as plt
avg_mass = penguin_clean.groupby("Species")["Body_Mass_g"].mean().reset_index()

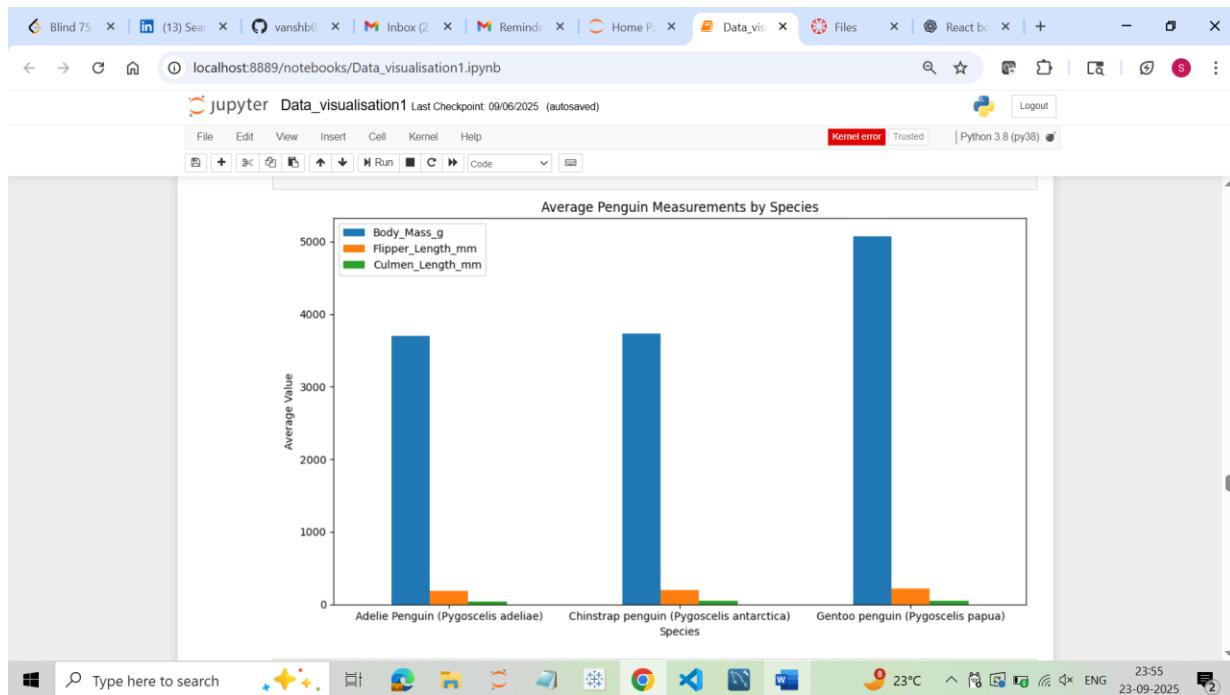
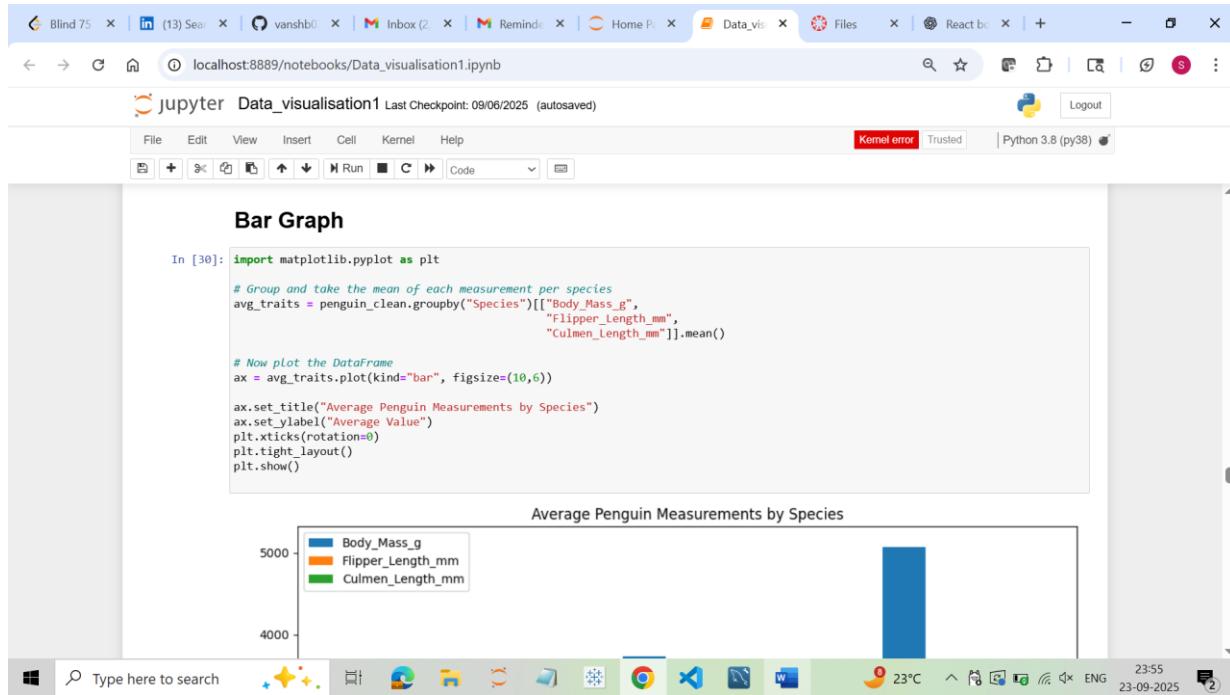
fig, ax = plt.subplots()
ax.plot(avg_mass["Species"], avg_mass["Body_Mass_g"], marker="o", linewidth=4)

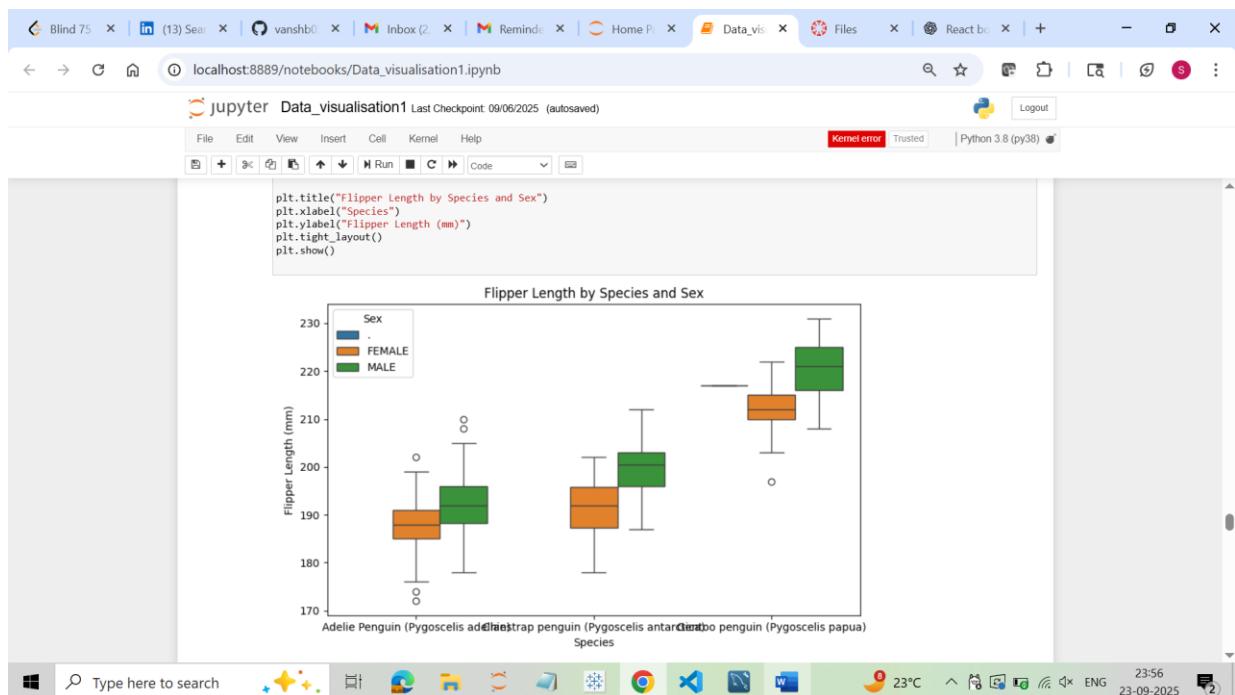
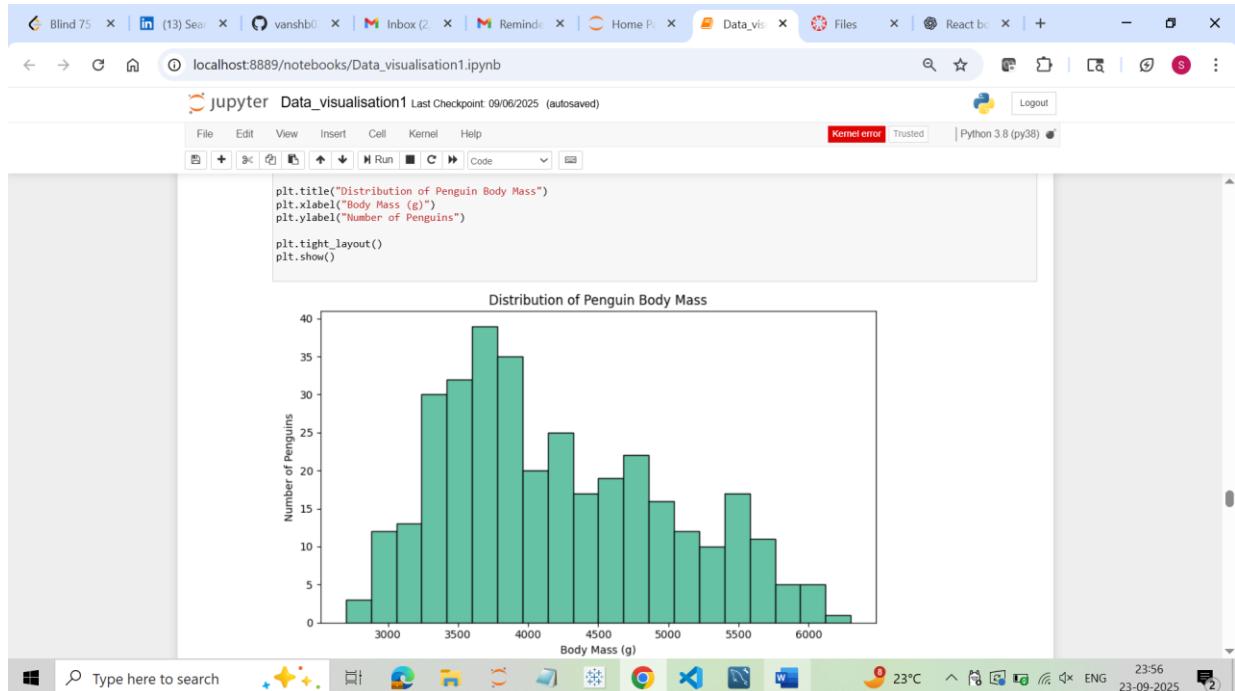
ax.set_title("Average Penguin Body Mass by Species")
ax.set_xlabel("Species")
ax.set_ylabel("Body_Mass_g")

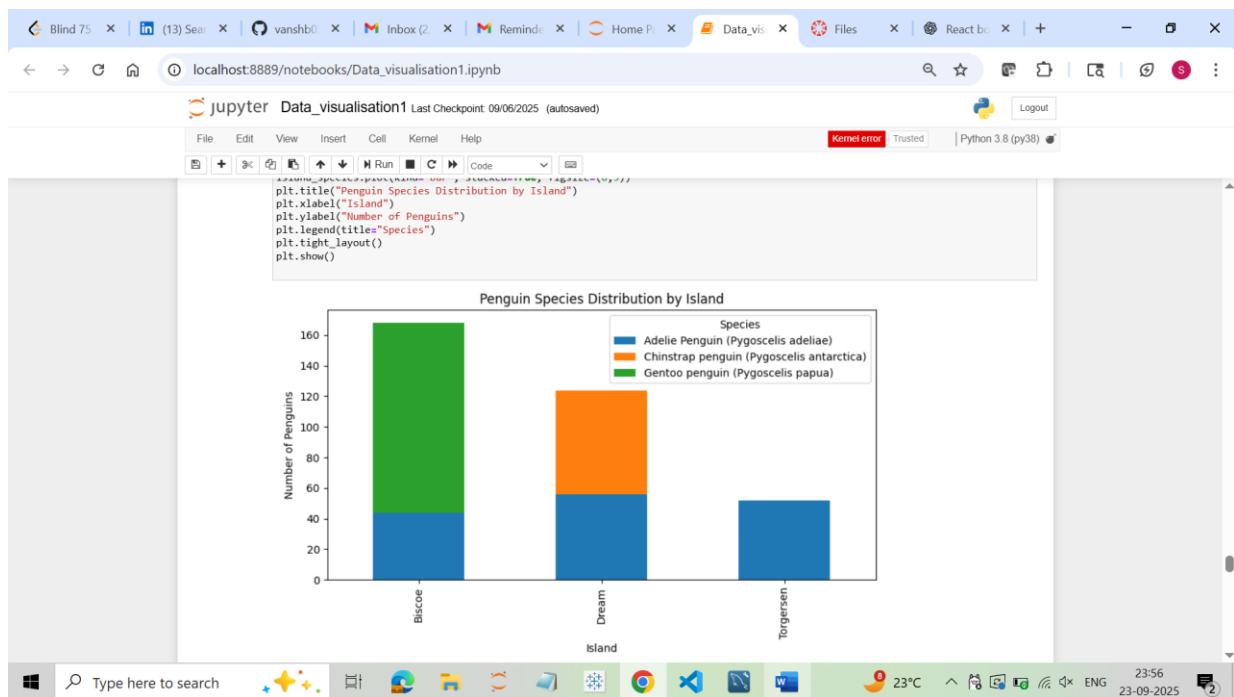
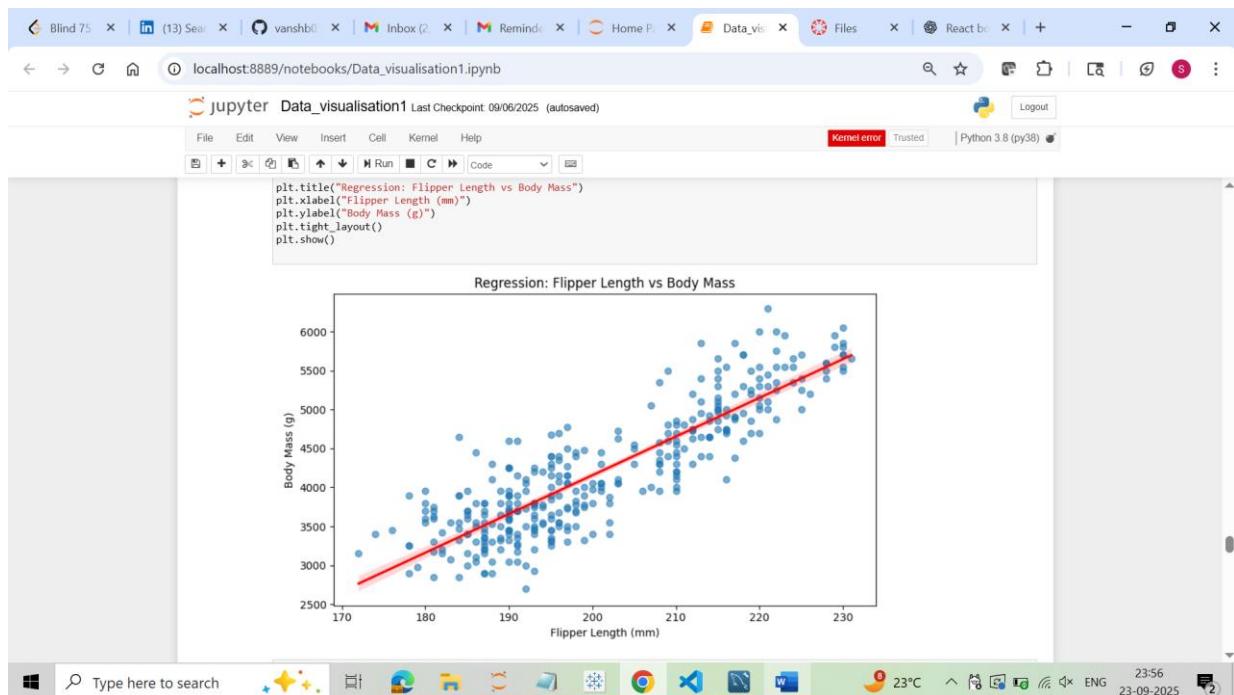
plt.tight_layout() #automatically adjusts the spacing between subplots and around the edges so everything fits nicely.
plt.xticks(rotation=45) #rotates x labels to 45 as the names are too big
plt.show()
```

Type here to search 23:54 23-09-2025









localhost:8889/notebooks/Data\_visualisation1.ipynb

jupyter Data\_visualisation1 Last Checkpoint 09/06/2025 (autosaved)

In [36]:

```
import seaborn as sns
import matplotlib.pyplot as plt

# Compute correlation matrix
num_cols = ["Body_Mass_g", "Flipper_Length_mm", "Culmen_Length_mm", "Culmen_Depth_mm"]
corr = penguin_clean[num_cols].corr()

plt.figure(figsize=(6,5))
sns.heatmap(corr, annot=True, cmap="YlGnBu", fmt=".2f")
plt.title("Correlation Heatmap of Penguin Measurements")
plt.tight_layout()
plt.show()
```

Correlation Heatmap of Penguin Measurements

	Body_Mass_g	Flipper_Length_mm	Culmen_Length_mm	Culmen_Depth_mm
Body_Mass_g	1.00	0.87	0.59	-0.47
Flipper_Length_mm	0.87	1.00	0.66	-0.58
Culmen_Length_mm	0.59	0.66	1.00	-0.24
Culmen_Depth_mm	-0.47	-0.58	-0.24	1.00

