**Create a new google collab file that only contains the code for the exercises.**

Please add a text cell that includes:
- title "problem set 4"
- team lab 1 or 2
- name of the team members.

The screenshot shows a Google Colab notebook with the following content:

Problem Set 4

Team Lab 1/2:

Slava Jankin

Paulina Garcia Corral

```python
# Use this data

penguins = pd.read_csv("penguins_classification.csv")
culmen_columns = ["Culmen Length (mm)", "Culmen Depth (mm)"]
target_column = "Species"

from sklearn.model_selection import train_test_split

data, target = penguins[culmen_columns], penguins[target_column]
data_train, data_test, target_train, target_test = train_test_split(
    data, target, random_state=0
)
```

Create a decision tree classifier with a maximum depth of 2 levels and fit the training data. Once this classifier is trained, plot the data and the decision boundary to see the benefit of increasing the depth. To plot the decision tree.

```python
[14] # Write your code here
```

Did we make use of the feature "Culmen Length"? Plot the tree using the function sklearn.tree.plot_tree to find out!
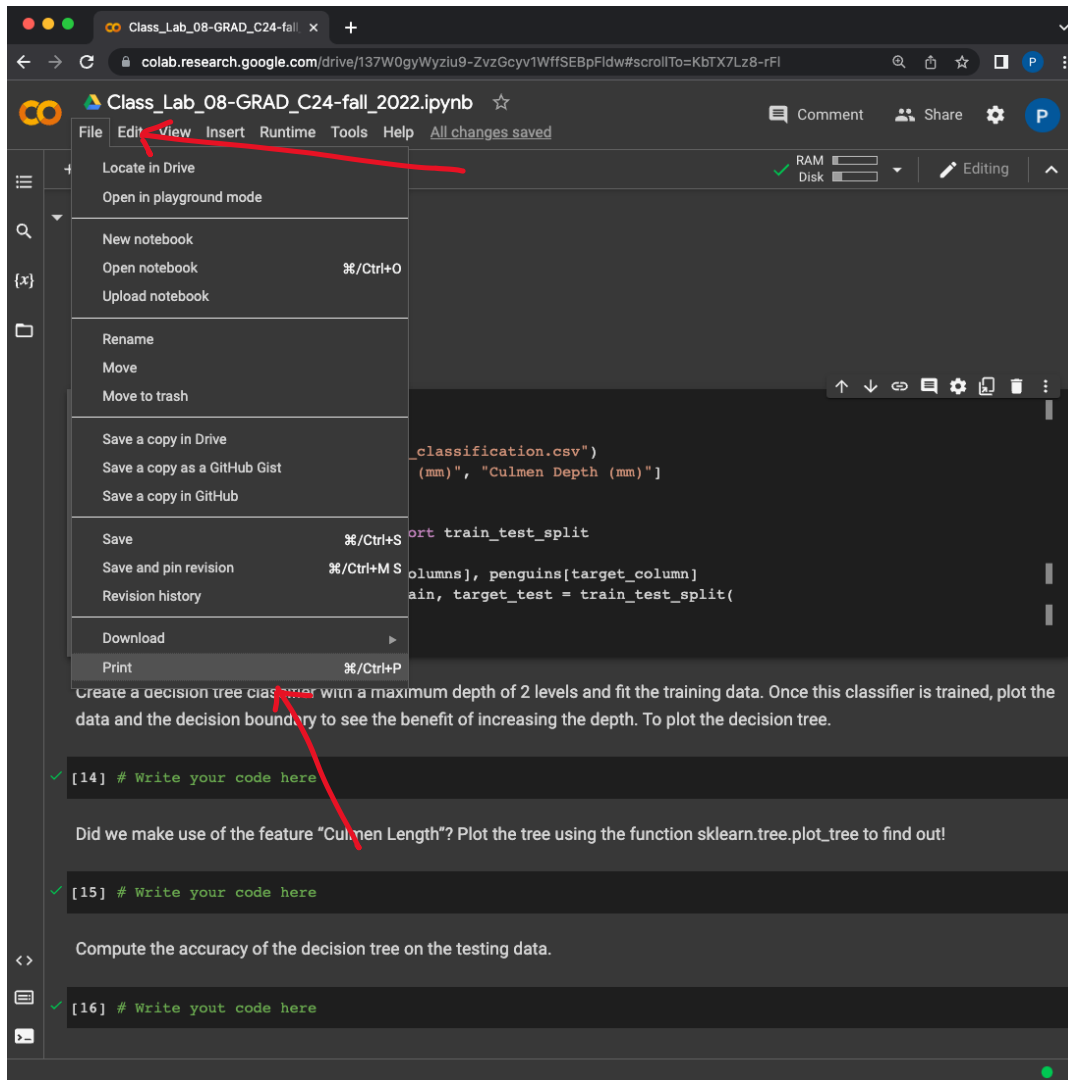
```python
[15] # Write your code here
```

Compute the accuracy of the decision tree on the testing data.

```python
[16] # Write yout code here
```

Submit the code as a PDF with all the plots visible.

Go to "file" and "print"

Select "Save as a PDF" for destination.

Then save.

Submit this file to Moodle, if and only if:
- only exercises code
- has Lab #
- all the team member names
- PDF file