# OCMP5310: Principles of Data Science

## Week 1 Live Session

**Presented by**

Daniela Rivas

THE UNIVERSITY OF
SYDNEY

# WELCOME!

# Daniela Rivas

- Bachelor in Engineering.

- PhD in Computational Biophysics.

- Postdoctoral Research Fellow at the Daffodil Centre.
  - Cervical Cancer and HPV modelling.

- Casual Academic at the University of Sydney.

# LIVE SESSIONS

# Live Sessions Details

- 6 Weeks.

- Every **Monday** at **6:30 PM**.

- First Live Session: Monday 17th of April.

- Last Live Session: Monday 22nd of May.

- 90 minutes.

# Live Sessions Expectations

**Before the session:**

– Watch the lecture videos.

– Work on the practical exercises.

**During the session:**

– Reflect on the main topics of the week.

– Run some additional group exercises.

– Give feedback on the exercises and the assessment tasks.

# ASSESSMENTS

# Assessments

- The official syllabus is the authoritative source of assessment information.
    - https://www.sydney.edu.au/units/OCMP5310/2023-S1CRB-OL-OP
- 10%: Weekly Review Quizzes.
- 10%: Project Stage 1 (Week 3, 7 May 2023, 23:59*).
- 10%: Project Stage 2 (Week 5, 16 May 2023, 23:59*).
- 5%: Project Oral Presentation (Week 6, 22 May 2023, during Live Session).
- 15%: Project Stage 3 (Week 6, 28 May 2023, 23:59*).
- 50%: Final exam (Week 8, 5 June 2023, 18:30*).

*Sydney time

# Project

This is a small data science project that will be developed incrementally through the course, and which focuses on **understanding the data science process** from data ingestion and cleaning, over understanding data and gathering descriptive statistics, to building a (simple) predictive model. There are three main submissions, as well as a short oral presentation that is designed to train presentation skills.

# Project

- **Stage 0:** Project Proposal.
- **Stage 1:** Data Acquisition and Cleaning, and Problem definition.
- **Project Stage 2:** Data Summarisation and Analysis.
- **Project Stage 3:** Predictive Model and Evaluation.

- More details [here](#).

# WEEK 1: SESSION ACTIVITY

# Week 1: Session Activity

- Installing Python and PostgreSQL.

- Exploratory Data Analysis with Spreadsheets/Python (Calculating Descriptive Statistics, plotting).

- Descriptive Statistics and Data visualization with Python.

# Installing Python and PostgreSQL
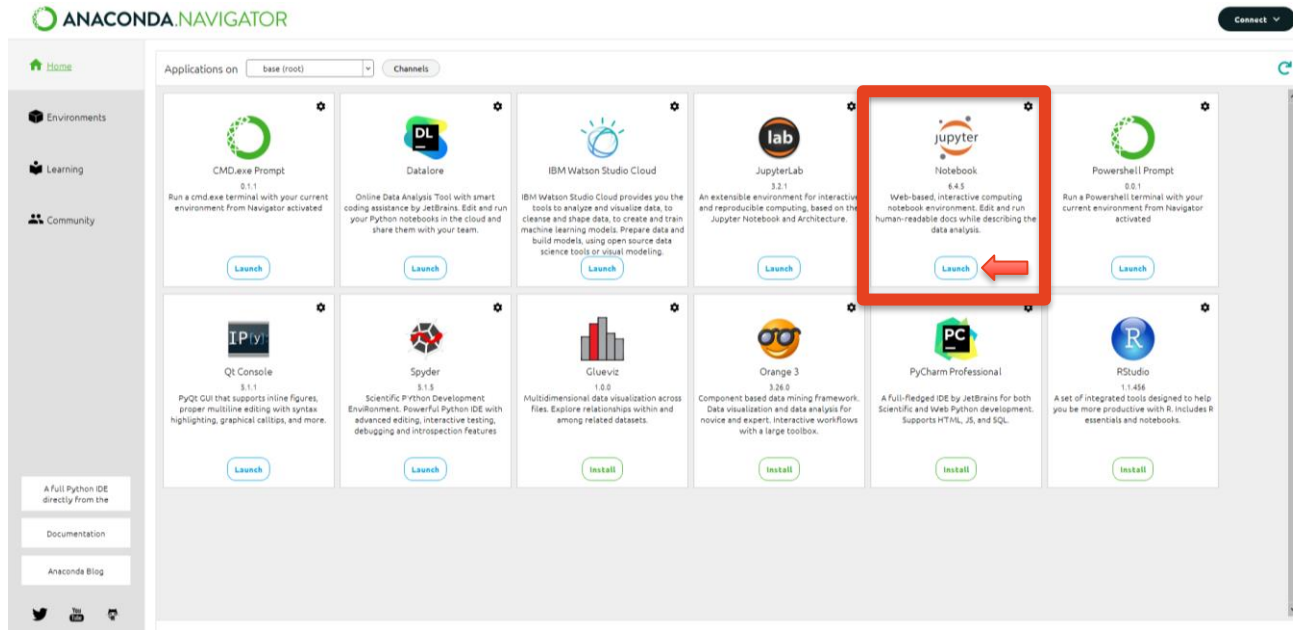
# Installing Anaconda

– Go to https://www.anaconda.com/products/distribution

– Choose the right installer for your Operating System and download it.

## Anaconda Installers

| Windows ⊞ | MacOS  | Linux 🐧 |
|---|---|---|
| Python 3.9 | Python 3.9 | Python 3.9 |
| 64-Bit Graphical Installer (621 MB) | 64-Bit Graphical Installer (688 MB) | 64-Bit (x86) Installer (737 MB) |
| | 64-Bit Command Line Installer (681 MB) | 64-Bit (Power8 and Power9) Installer (360 MB) |
| | 64-Bit (M1) Graphical Installer (484 MB) | 64-Bit (AWS Graviton2 / ARM64) Installer (534 MB) |
| | 64-Bit (M1) Command Line Installer (472 MB) | 64-bit (Linux on IBM Z & LinuxONE) Installer (282 MB) |

– Open the installer and follow the instructions to install.

# Opening Jupyter Notebook

– Open the Anaconda Navigator and click on "Launch" under Jupyter Notebook from the options available.
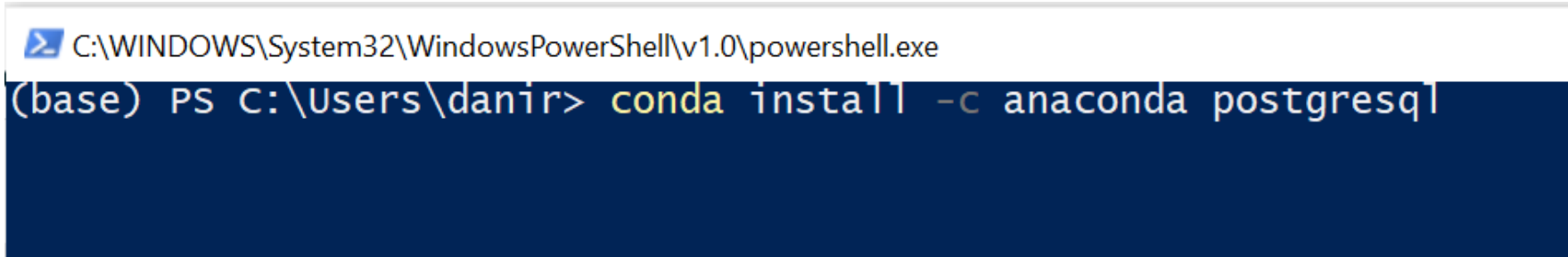
# Opening a New Jupyter Notebook

- In the upper right select [New] –> [Python 3] to open a new Python 3 notebook.

# Installing PostgreSQL on Anaconda

- Open your computer terminal or the Anaconda Powershell Prompt and type:
  - **`conda install –c anaconda postgresql`**



```
C:\WINDOWS\System32\WindowsPowerShell\v1.0\powershell.exe
(base) PS C:\Users\danir> conda install -c anaconda postgresql
```

# Installing PostgreSQL on your computer

- Go to https://www.postgresql.org/download/
- Choose the right option for your Operating System and click on **download the installer**.

## Interactive installer by EDB

Download the installer certified by EDB for all supported PostgreSQL versions.

**Note!** This installer is hosted by EDB and not on the PostgreSQL community servers. If you have issues with the website it's hosted on, please contact webmaster@enterprisedb.com.

This installer includes the PostgreSQL server, pgAdmin; a graphical tool for managing and developing your databases, and StackBuilder; a package manager that can be used to download and install additional PostgreSQL tools and drivers. Stackbuilder includes management, integration, migration, replication, geospatial, connectors and other tools.
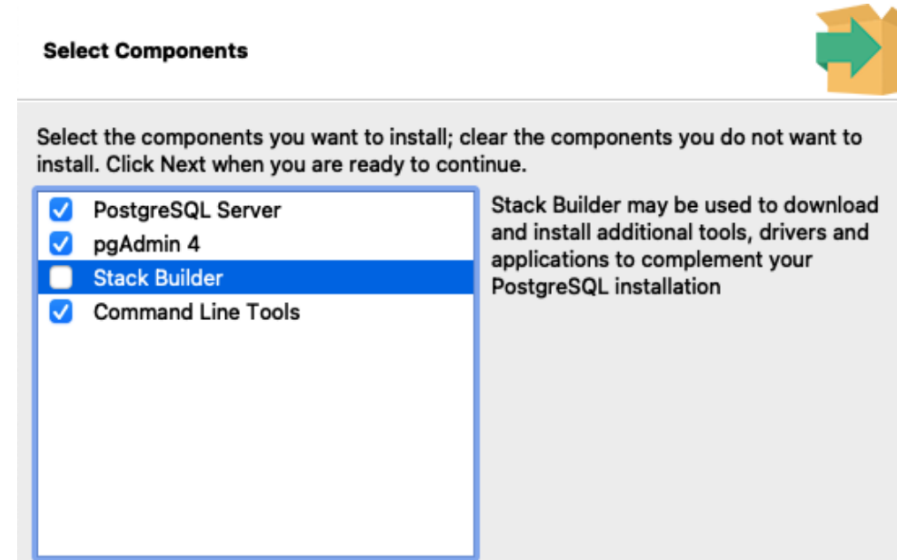
This installer can run in graphical or silent install modes.

The installer is designed to be a straightforward, fast way to get up and running with PostgreSQL on Windows.

*Advanced users* can also download a zip archive of the binaries, without the installer. This download is intended for users who wish to include PostgreSQL as part of another application installer.

# Installing PostgreSQL on your computer

- Select the latest PostgreSQL version available for your OS and download it.
- Open the installer and follow the instructions to install.
  - In "Select Components", select:
    - The PostgreSQL server.
    - pgAdmin.
    - Command Line Tools.
  - In "Password"
    - Type a password.
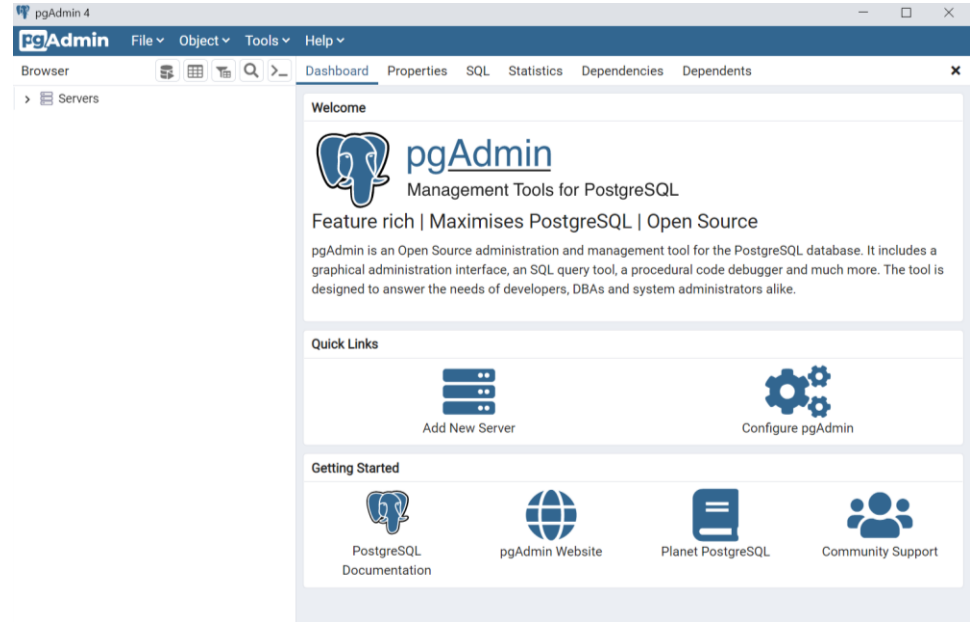    - **Remember this password!**
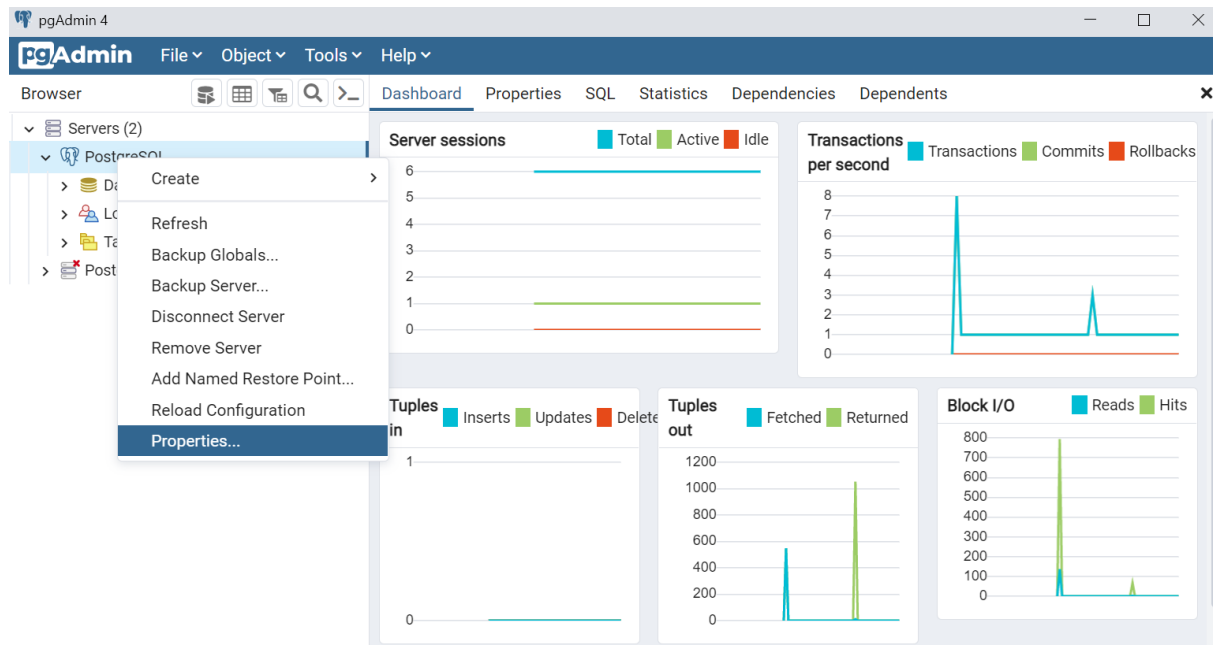
# Launching PostgreSQL



– Open pgAdmin4.

pgAdmin 4

– Type in the **password** you used during the installation.

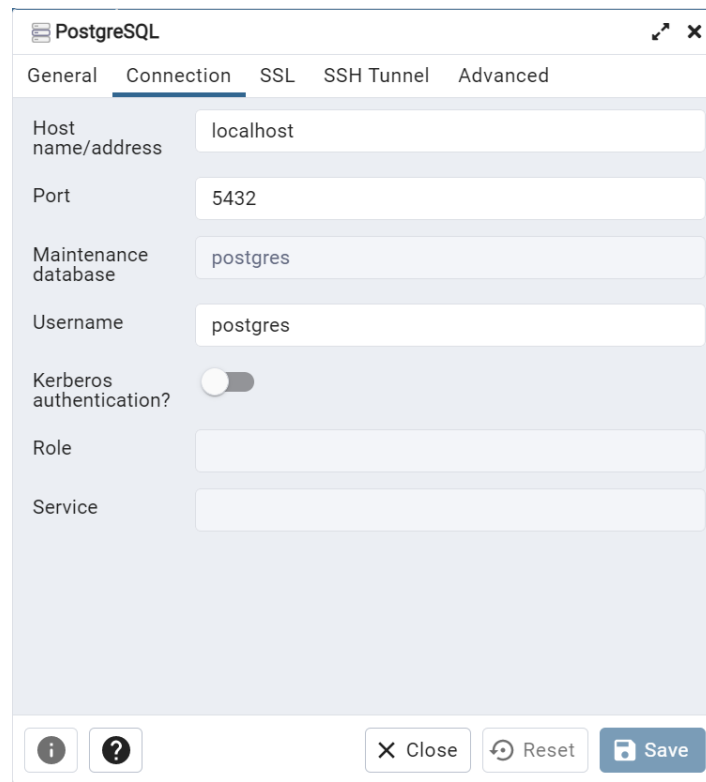– After entering your password, you will see the pgAdmin homepage.

# Launching PostgreSQL

– Find PostgreSQL Server in the left tree, right-click 'PostgresSQL', and select 'Properties'.

# Launching PostgreSQL

– Go to 'Connection' and select:

  – host name: localhost

  – port: 5432

  – username: postgres

# Connecting to the server

## Option 1

- Open Jupyter Notebook, and select New > Terminal.
- Once a new terminal is opened, run the following command:
  - `psql -h localhost -U <username>`

# Connecting to the server

## Option 2

– Only for the **first time you connect**, in your Jupyter Notebook, type:

– `!pip install psycopg2`

– To connect to the database server, copy the code on the right to your Jupyter Notebook (Do not forget to **change the password** to your own password)

```python
import psycopg2
def pgconnect():
    """ Connect to the PostgreSQL database server """
    conn = None
    try:

        # connect to the PostgreSQL server
        print('Connecting to the PostgreSQL database...')
        conn = psycopg2.connect(host = 'localhost',
                                database = 'postgres',
                                user = 'postgres',
                                password = 'abcd1234')
        print("connected")

    except Exception as e:
        print("unable to connect to the database")
        print(e)
    return conn
conn = pgconnect()
```

# Data exploration with Python

# Activity

- In Canvas, go to:
  - Exercise: Data Acquisition and Data Cleaning with Python.
- Download Jupyter Notebook:
  - data_exploration_with_Python.ipynb
- Download WFH survey responses:
  - WFH-Survey-Responses-NSW.csv

# Project Stage 0

# Project Stage 0

- Identify and explore possible problems and datasets.
- Select project dataset, define problem, submit a project proposal before April 23rd.
- For those without own proposal, a project will be chosen and assigned.

# Where to find a dataset/problem

- WHO - World health organisation
- The World Bank
- International Labour Organization
- United Nations
- OECD
- Earth Data (NASA)
- UNICEF
- Australian Government

- Australian Bureau of Statistics
- Research Data Australia
- UCI Machine Learning Repository
- Kaggle Datasets
- AIHW Data

# QUESTIONS?