# Software for the Course

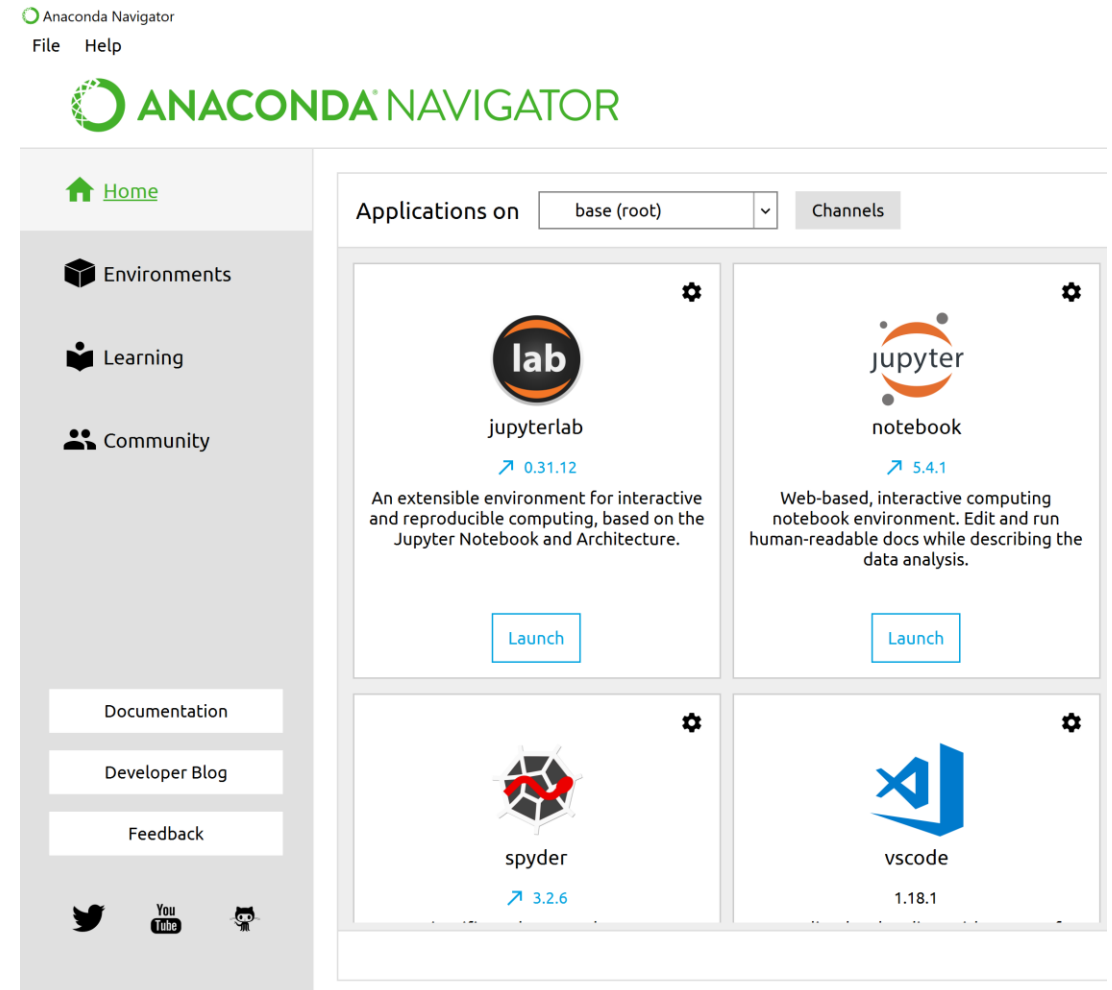Joel Klein

# Setting up Python

# Python Distributions

- Using a distribution simplifies the process of setting up your python environment, includes necessary data packages, and integrate useful tools (IDE's, notebooks, etc.)

- In class we will be using the **Anaconda** Distribution which is based on the conda package manager

- It provides integration with Jupyter, Virtual Envs, etc.
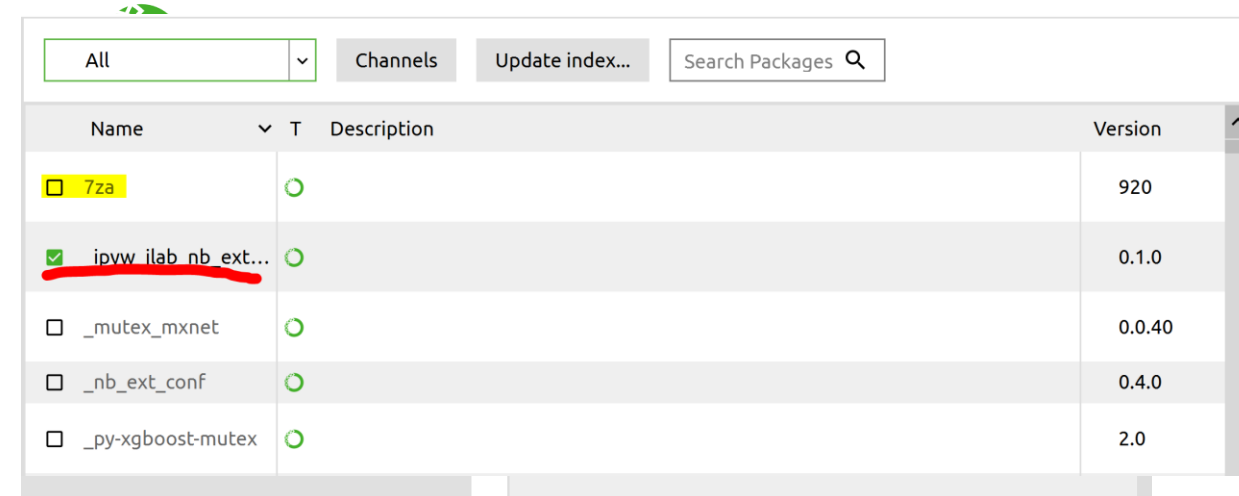
# Anaconda

# Anaconda Navigator

- The Navigator is a main landing page for working with your python environment.
- Here we can launch editors (spyder, jupyter notebook, etc.) to write and develop python code
- In addition we can manage (install packages, etc.) our python environment

# Anaconda Environments

- Clicking on the "Environment" tab will show us what environments are available in Anaconda
  - In the simplest terms, an anaconda "environment" is a self-contained collection of python packages.
- From the "Environment" tab we can see which packages are installed and which packages are available for installation.
  - If you click on a package for installation, you'll be prompted to Apply your changes

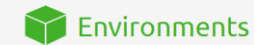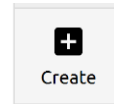# Setting Up Class Environment

- For this class I've provided an environment file (and a package list) on blackboard. This environment should include all of the packages necessary for the class and can be installed as follows:
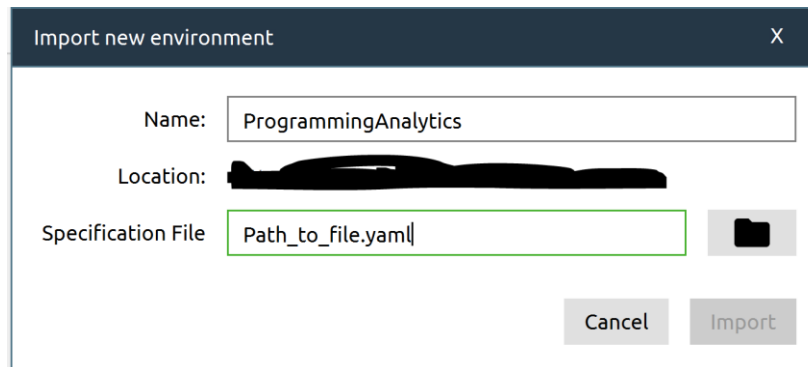  1. Navigate to the "Environment" tab in Anaconda.
  2. Click on the "Create" button
  3. Next fill in the specification file by navigating to the provided .yaml file
  4. Import

# Anaconda Applications

- On the home page we can choose which environment (base(root) in the image) we want to launch applications from.
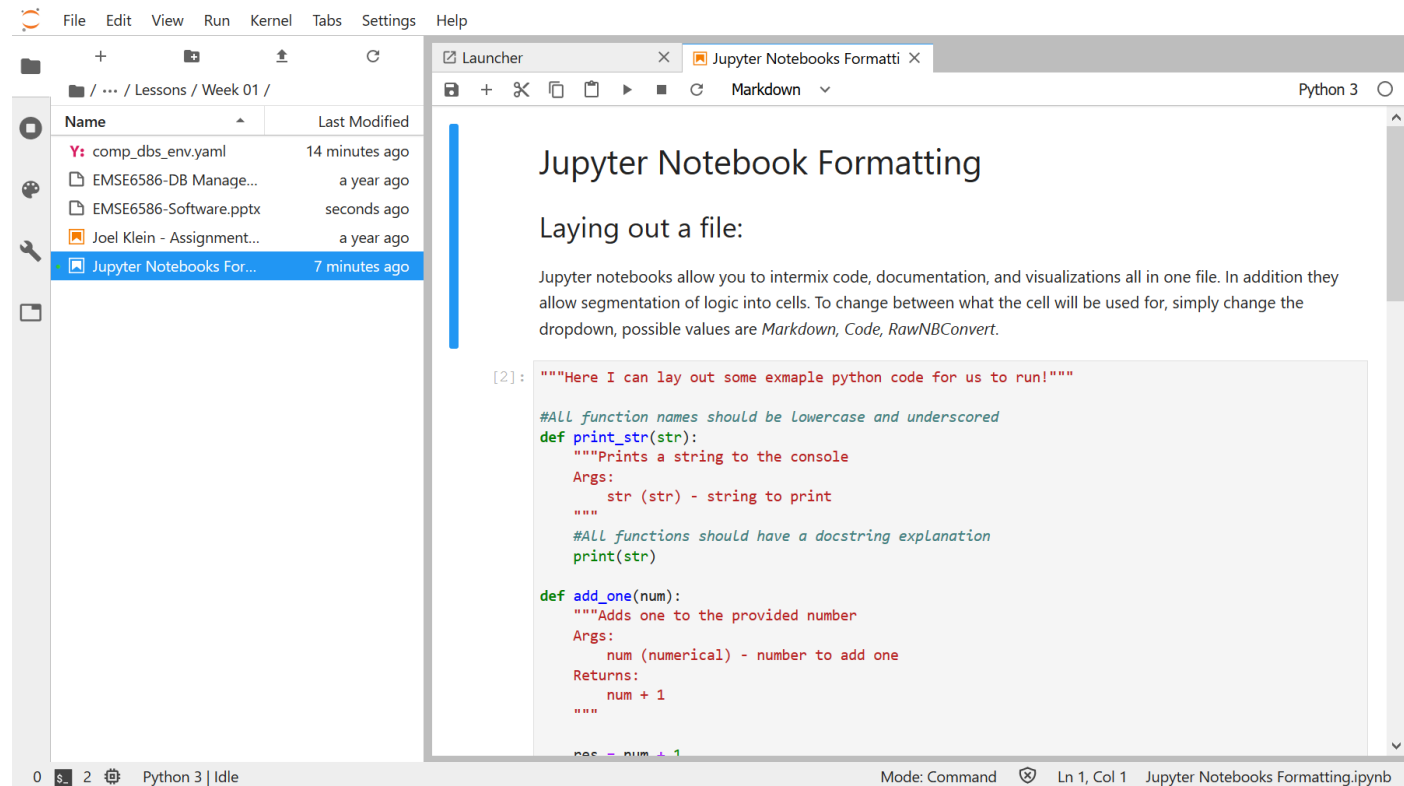- Clicking the "Launch" button on any of these applications will launch a separate window.
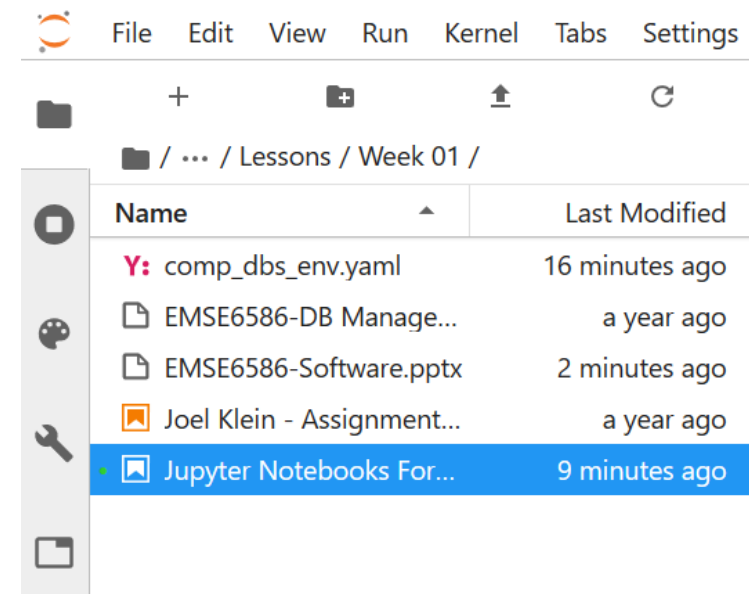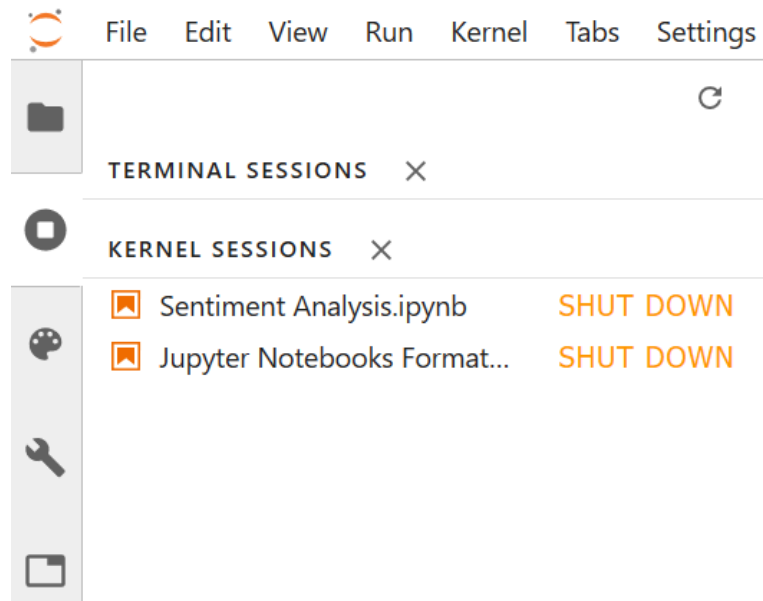
# Jupyterlab

# JupyterLab

- Jupyter Notebooks and Jupyterlab are the IDEs we will use for this course.

- Jupyter Notebooks are used heavily throughout this course, as they provide a means to intermix text, code, and graphics.
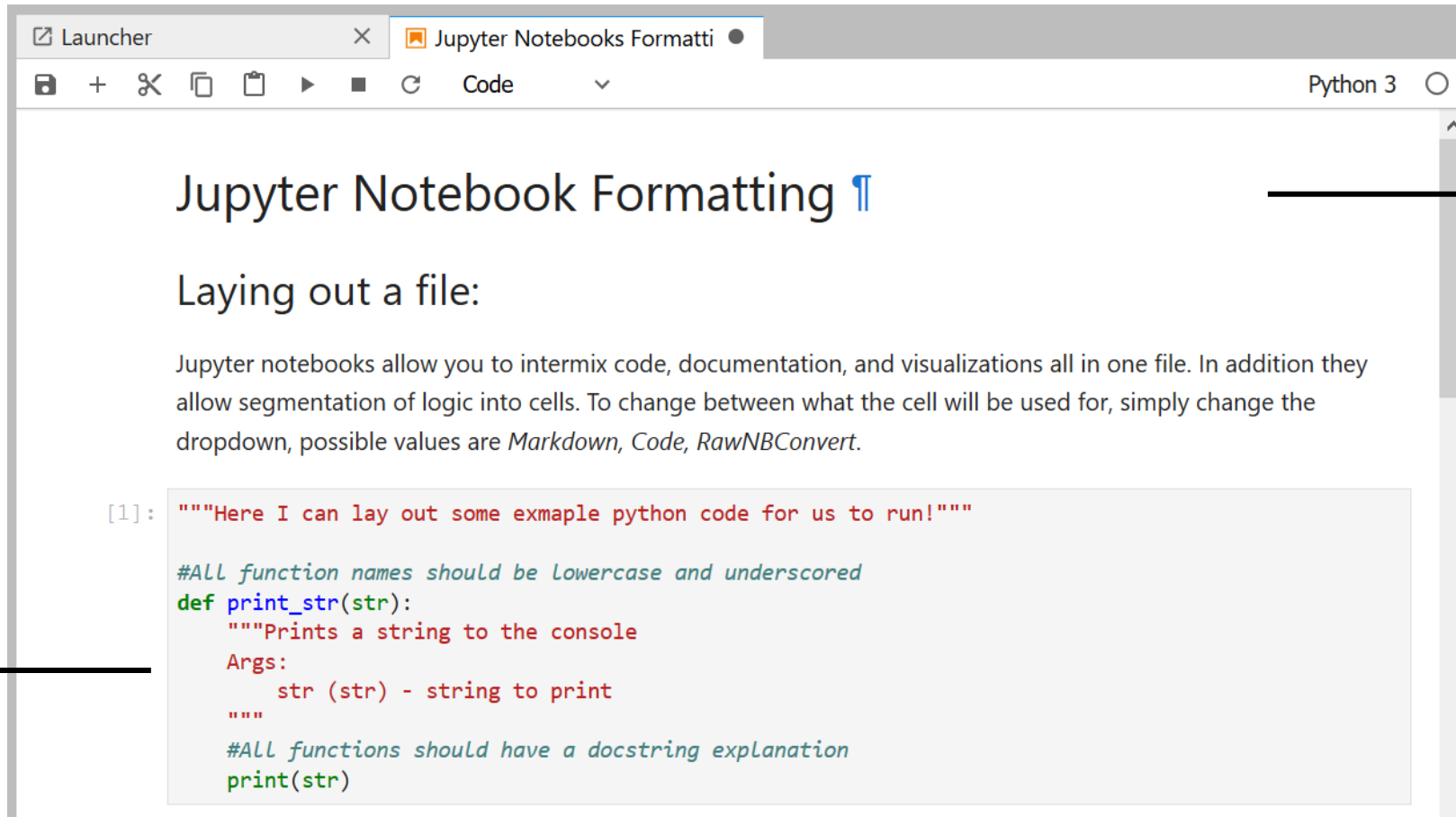
# File Navigation and Kernels

We can navigate the file system using the left hand panel

In addition we can monitor the running kernels within that same panel

# Working with Notebooks



Notebooks let us mark up (using markdown) our code to provide context and information.

The primary function is to provide cells of execution for our code/logic
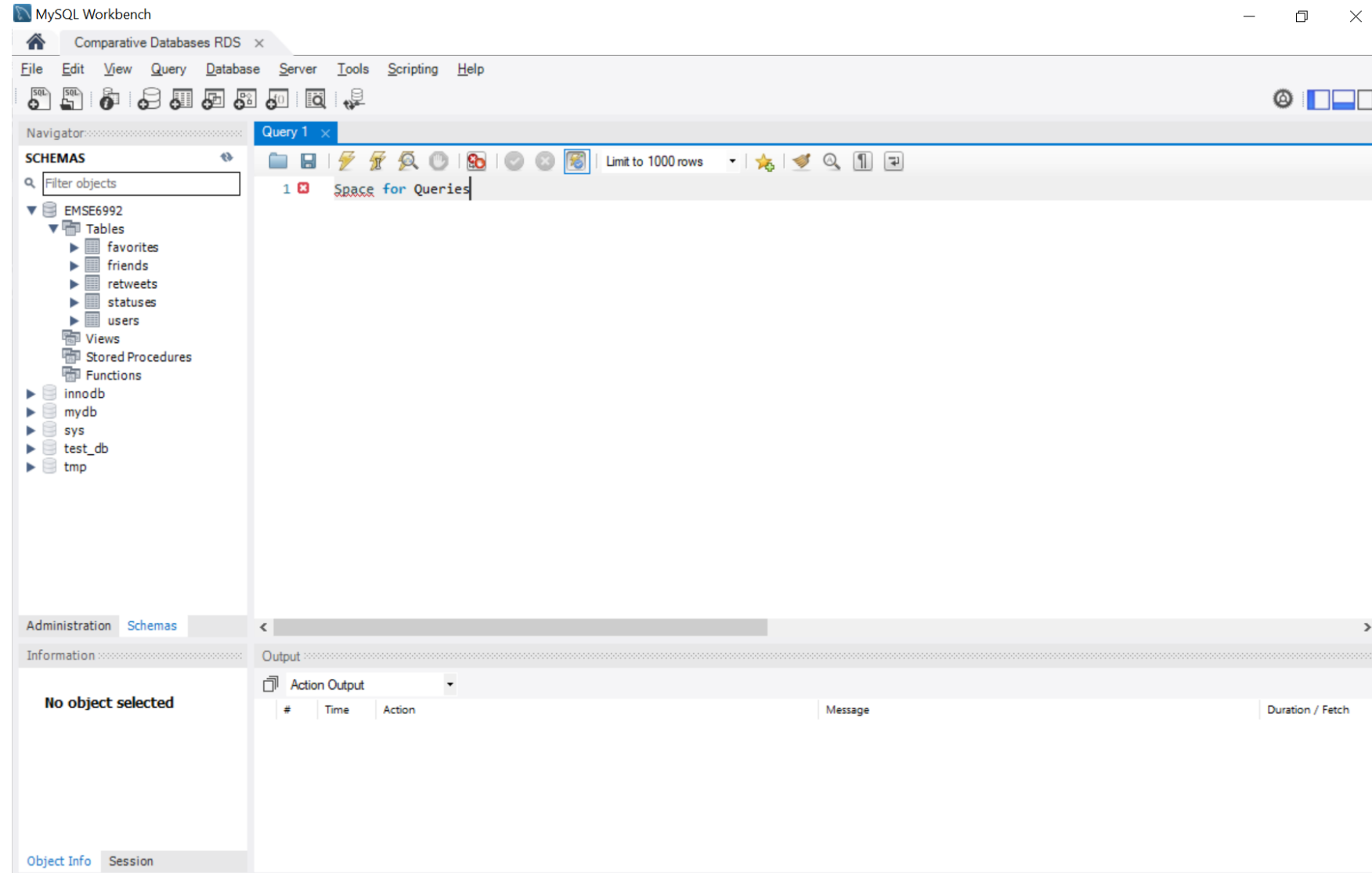
# The Databases

# MySQL
## Type – SQL | Environment – AWS RDS

- We will be interacting with an AWS RDS instance running MySQL to learn about SQL query languages
  - RDS supports several different SQL types

- For non-programmatic access we will be utilizing MySQL Workbench (https://www.mysql.com/products/workbench/)

- For programmatic access we will be using the **pymysql** package for python
  - This can be downloaded through anaconda or through pip

# MySQL Workbench

This is a free tool for managing MySQL instances and running queries and working with database schemas

You can download it for free at:
https://www.mysql.com/products/workbench/

# Mongo DB:
# Type – Document | Environment – EC2 Instance

- For this course we will be using an instance of Mongo DB running on an AWS EC2 instance
  - Alternatively, there is Azure Cosmos and AWS DocumentDB which both support the Mongo API in managed environments

- For non-programmatic access we will be using Robo3T

- For programmatic access we will be using python and the pymongo package

# Robo3T

This is a free tool for querying databases that support the mongo API and provides some minimal management options

You can download it for free at:
https://robomongo.org/download

# Arango:
## Type – Graph | Environment – EC2 Instance

- For this course we will be using an instance of Arango running on an AWS EC2 instance
  - Alternatives are Azure Cosmos and AWS Neptune for managed solutions

- For non-programmatic access we will be using a web browser (http://18.219.151.47:8529)

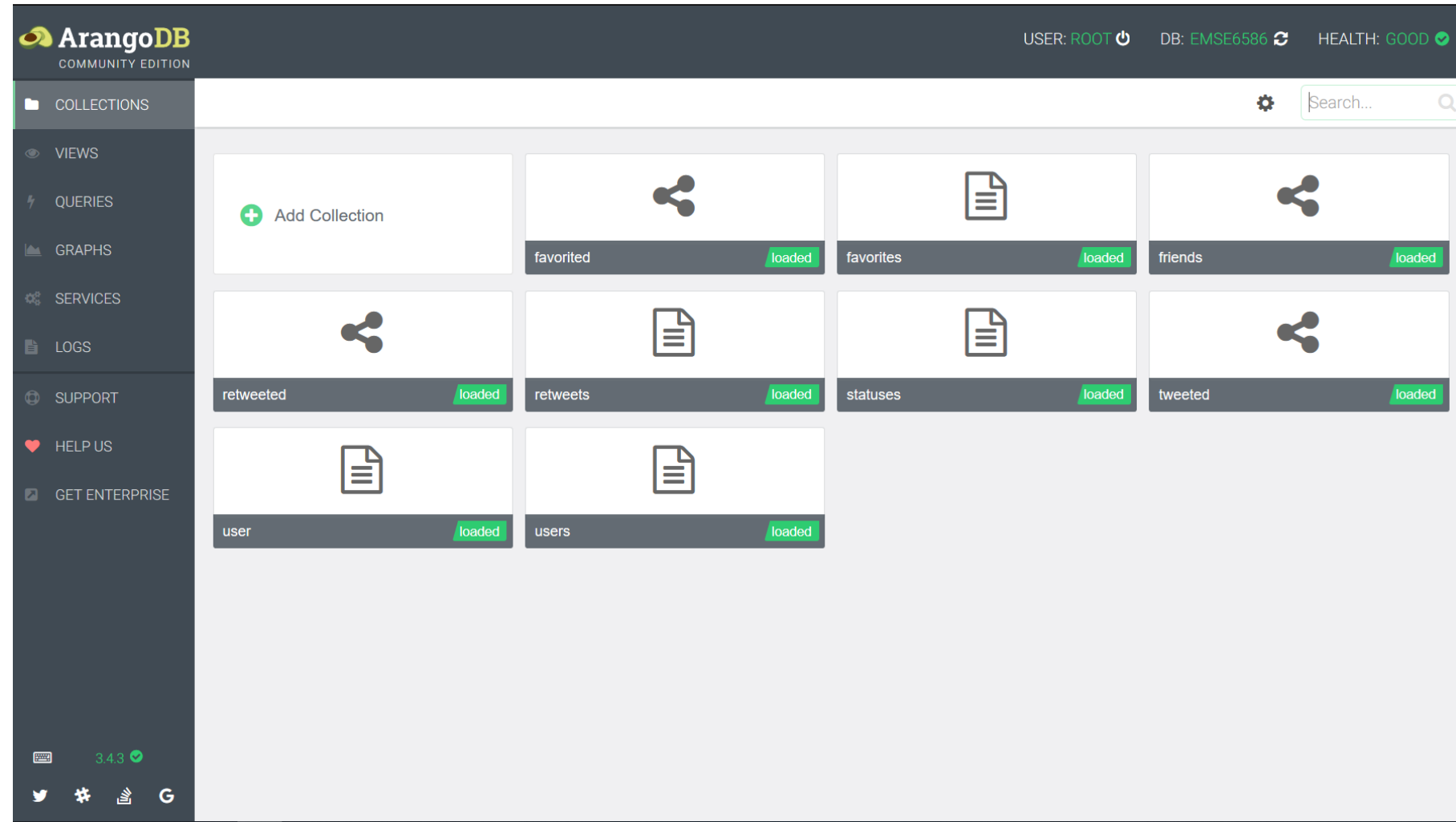- For programmatic access we will be using python and python-arango package

# Arango Web-GUI

Arango natively supports a web GUI for managing the database, running queries, and visualizing graphs

You can download it for free at:
http://18.219.151.47:8529

# Cloud Database Solutions

# Cloud Databases

- While the means in which I host our database resources leverage cloud infrastructure, they are all open-source solutions not dependent on cloud infrastructure

- The following database solutions are approaches to handling data that are more unconventional and rely more heavily on services.

# Azure Cosmos

- Azure Cosmos is a multi-model database provided by Microsoft on their Azure cloud.
  - Multi-Model means that the database supports multiple models of interacting with the data

- The benefits of Cosmos is that it provides the following:
  - Managed Solution
  - "Infinitely" scalable
  - Multi-Model
  - Etc.

# AWS Athena

- AWS Athena is Amazon's approach to SQL-like access to static files.
  - Essentially Athena is focused on providing queryable access to files stored in S3

- The benefits of Athena are:
  - Quick and Flexible access to existing data
  - A Serverless approach to SQL-like interfaces

# Hadoop and Spark

# Hadoop/Spark Infrastructure

- Time permitting we will also dive into the Hadoop and Spark architecture

- Hadoop is a pseudo database architecture that has evolved into a ecosystem of tools and services for processing large volumes of data
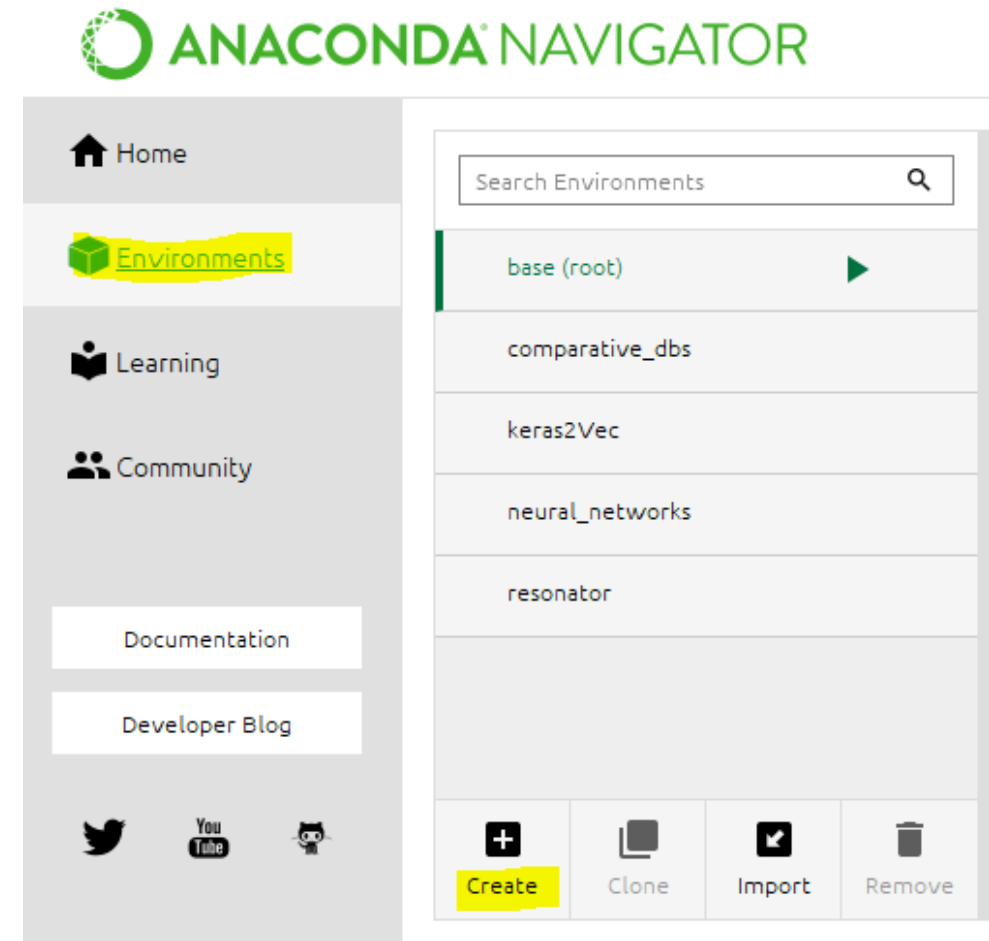
# End Slide

## EMSE 6586 – DBMS for Data Analytics
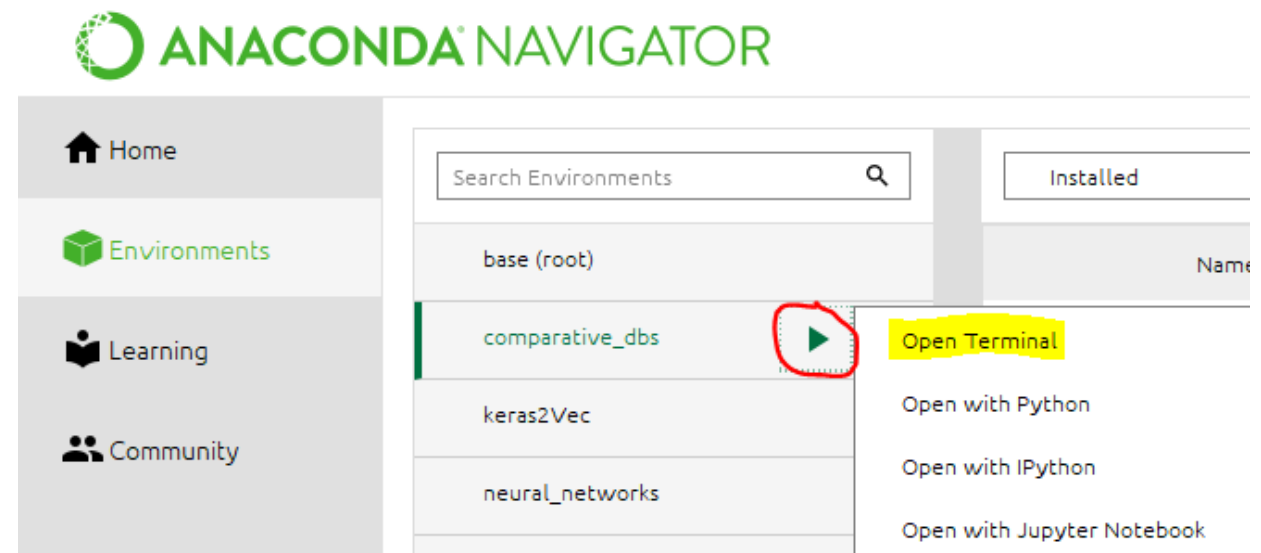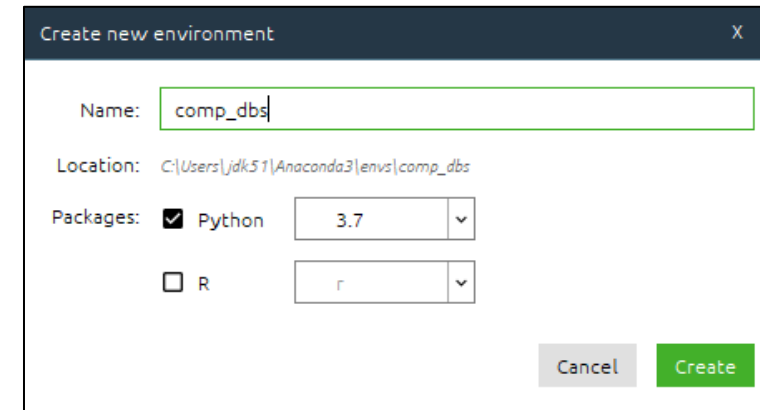
# Environment Problems?

# Pip Install to Environment

- If Anaconda gives you some problems with environment yaml file, then pip is your next best option.

- To install via pip, you'll want to create a new environment manually.

  1. Click on  Environments on the left-hand side
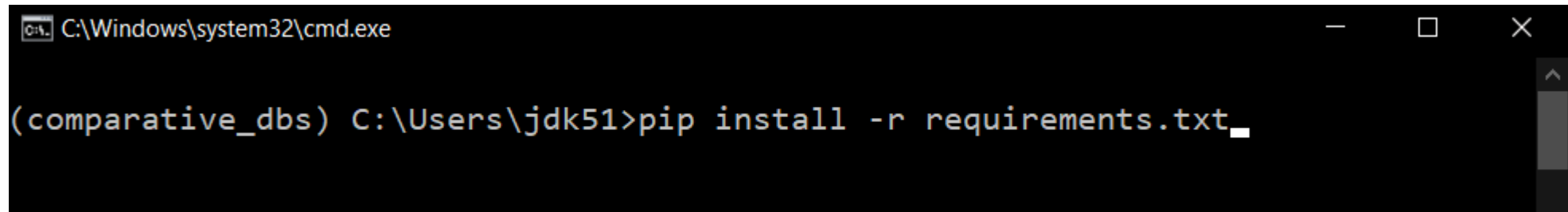  2. Click on  at the bottom

# Creating a New Environment

- In the resulting pop-up
    1. Provide a name for the environment
    2. Select "Python 3.7"

- Once the environment is created, you'll need to open a terminal.
    1. Click on your environment
    2. Hit ▶ that appears
    3. Click on "Open Terminal"

# Finally Pip Install

- In the resulting window you'll be able to install the requirements for the environment.

- Simply run "pip install –r requirements.txt"
  - Note: Make sure that you provide the correct pathing to the requirements.txt file