

EE3801

Data Engineering Principles

A/Prof. Bharadwaj Veeravalli
Dept of ECE, NUS

&

Dr. Yen Shih-Cheng
EDIC, CDE, NUS

elebv@nus.edu.sg

shihcheng@nus.edu.sg

Lesson Plan – *What you will learn in this module?*

Part 1 - A/Prof Bharadwaj V

- Introduction to Data Engineering (DE) and computing platforms, DE concepts data wrangling and ETL/ELT, Data warehouses, Data pipelining concepts, current day big data architectures, design example of a BDP, introduction to Cluster and Cloud/DC platform architectures, resource sharing in cluster platforms

Part 2 - Dr. Yen Shih-Cheng

- Amazon Web Services Elastic Compute Cloud, Parallel Cluster, Elastic Block System file system, SLURM Workload Manager, data pipeline construction and operation, high-throughput data visualization and inspection, vectorization techniques in data analysis.

Assessment – Marks Distribution

- Labs 1 to 3: 35%
- Quiz 1: 15%
- Labs 4 to 8: 35%
- Quiz 2: 15%

Part 1: All Lectures, Lab sessions, and Quiz are held in this venue LT3(Mon)/LT7(Thurs).

Part 2: All Lectures and Quiz will be held in LT7 on Thursdays, while Lab sessions will be held in E2-03-08/09 on Mondays.

Part 1 – Contents, Assessments & Marks Distribution

Part 1: 50%

- Lab sessions 35% (3 Lab sessions)

Quiz 1 (Chapters 1, 2, and 4): 15%

Date: Oct 2, 2023

Time: 12pm-1pm (60 mins)

Venue: LT3

Part 1 Lab Schedule

LAB 1

Aug 30 - Lab 1 assignment release
Aug 31 - Lab 1 briefing by GAs (45 mins)
Sept 7 - Deadline Lab 1

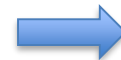
LAB 2

Sept 6 - Lab 2 assignment release
Sept 7 - Lab 2 briefing by GAs (45 mins)
Sept 14 - Deadline Lab 2

Work on
your
individual
PCs/Laptops

Lab 3:

Sept 13 Lab 3 assignment release
Sept 14 - Lab 3 briefing (45 mins)
Sept 21 - Deadline Lab 3



(Your first ride on AWS! 😊)

Part 2 - Assessments & Marks Distribution

Part 2: 50%

- Lab sessions: 35%
5 Lab sessions all equally weighted
- Quiz 2 (hands on DE): 15%
Date: Nov 16, 2023 (Thursday)
Time: 3 pm
Venue: LT7

Lesson Plan – Part 2

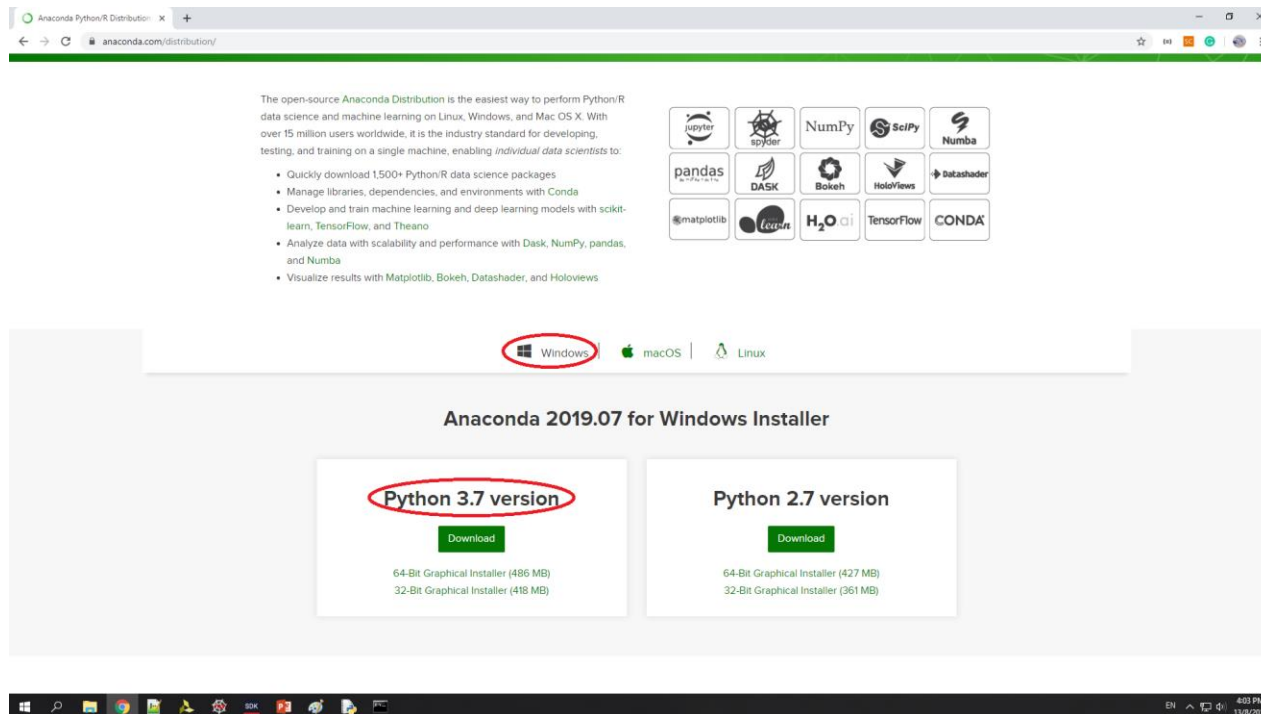
Week	Monday	Thursday
7	Quiz 1 (Part 1)	Lecture
8	Lab 4	Lecture
9	Lab 5	Lecture
10	Lab 6	Lecture
11	Lab 7	Lecture
12	Lab 8	Lecture
13	Public Holiday	Quiz 2 (Part 2)

Labs will be held in E2-03-08/09

Lectures and Quiz 2 will be held in LT7

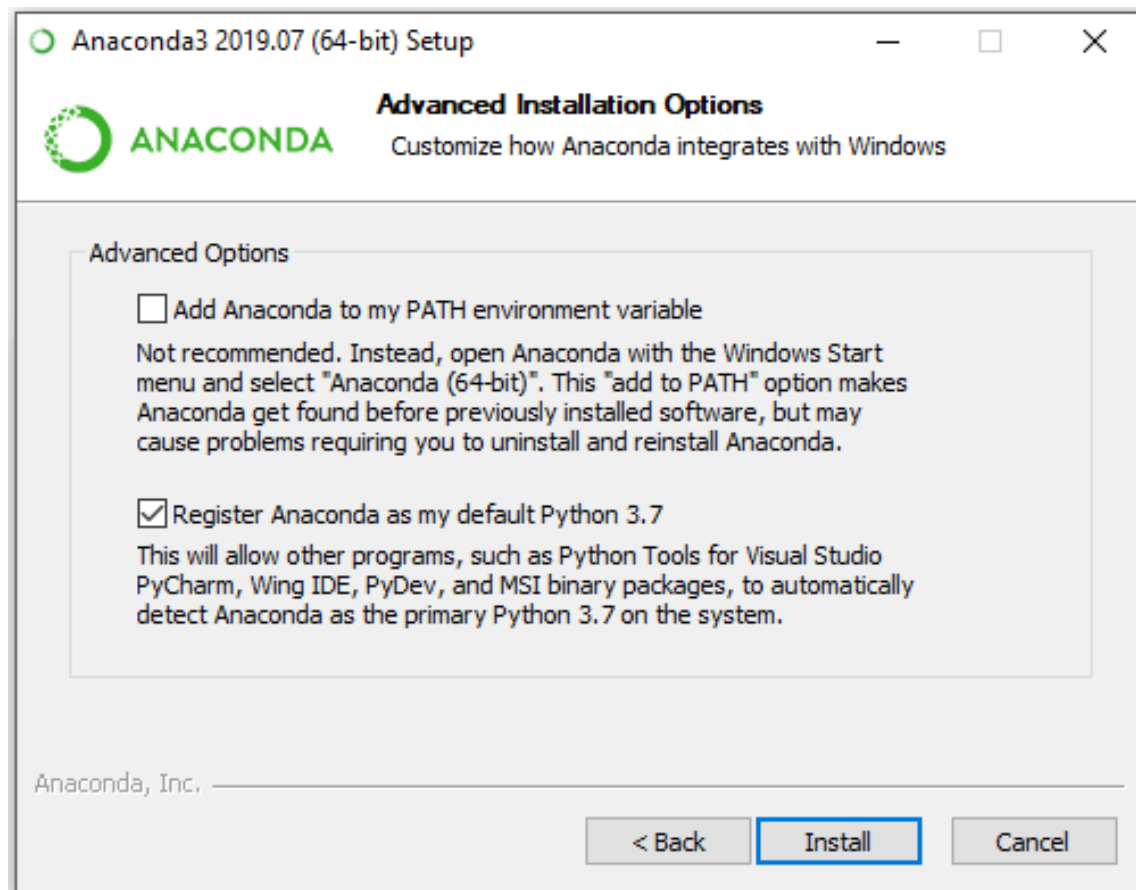
Anaconda Installation

- Download Anaconda for Python 3.x from here:
<https://www.anaconda.com/distribution/#download-section>
- <https://docs.conda.io/en/latest/miniconda.html> (**Miniconda**) [~250 MB]
- <https://katiecodes.com/setup-python-windows-miniconda/> (Useful)



Anaconda Installation

- Click Next until the screen below. Select the options as shown in the screenshot. Click Install.



Anaconda Installation

- Visit: <https://www.datacamp.com/community/tutorials/installing-anaconda-windows>

The above site takes you through a step-by-step installation process;

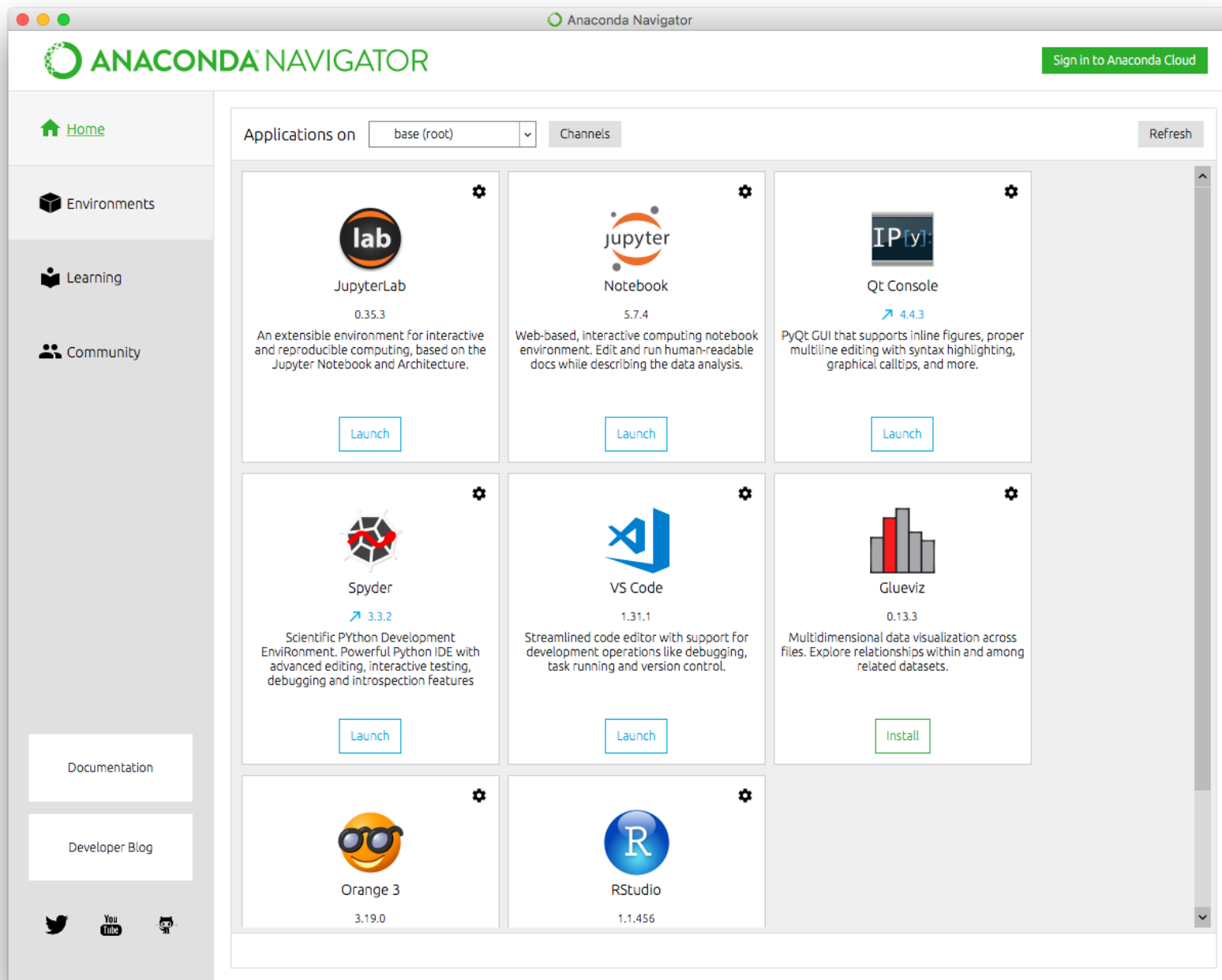
Note that with Anaconda lots of packages (~400+) will be installed; So, if you need other packages, you need to install on your own:

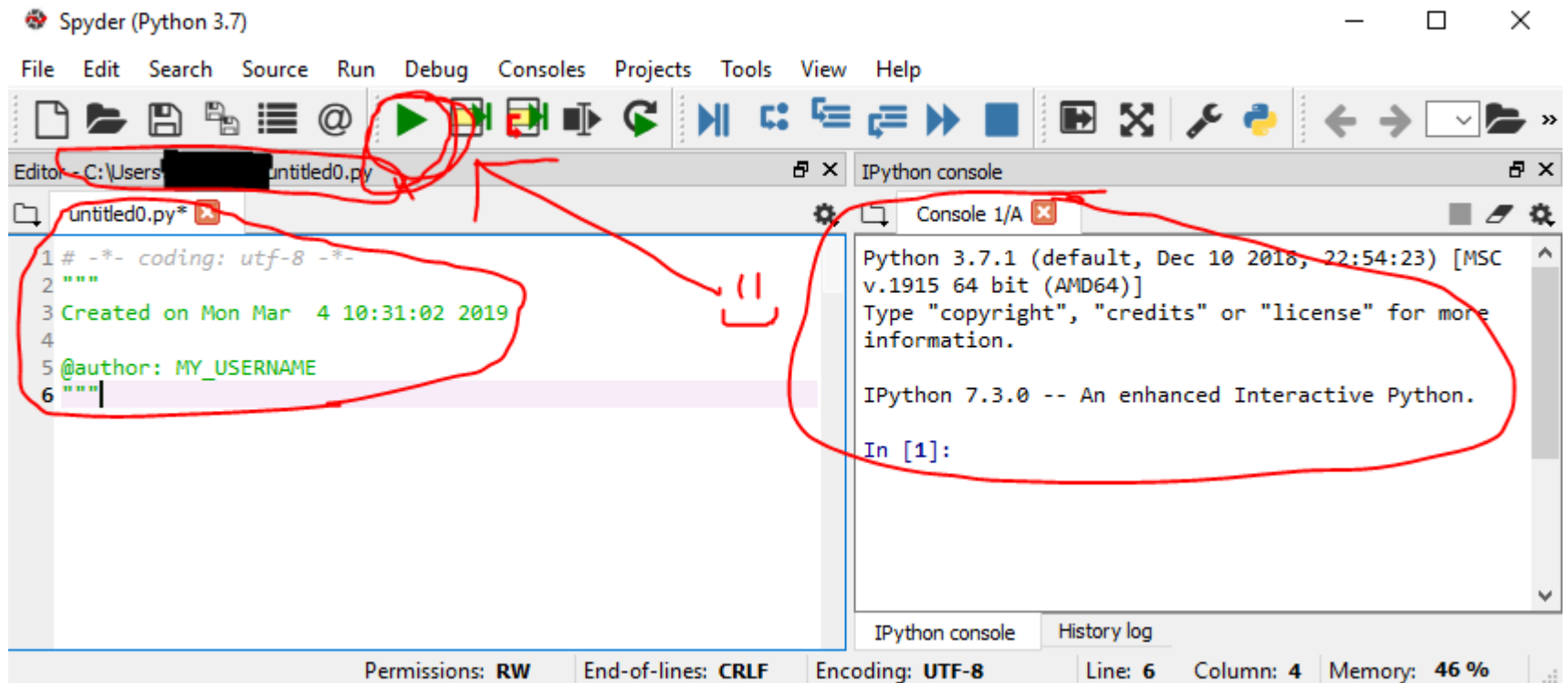
Use:

- In Miniconda – IDLE is available; Jupyter notebook is NOT available; So, just like other pkgs, you can install using ***conda install jupyter*** should install it;
- With Anaconda *Jupyter, Spyder, and IDLE* are already built-in.
- Second part of this module exclusively uses Spyder editor; For the first part you can use either IDLE or Spyder;

Installing packages required:

- ***pip install pkg-name*** (or) ***> conda install pkg-name***





What next?

Before Aug 30, 2023:

- Install Python (*miniconda / Anaconda*)
- Editors – IDLE / Spyder
- Python packages that will/may be used:
 - Numpy;
 - Scipy;
 - Pandas;
 - Matplotlib* / Seaborn; (* Recommended)

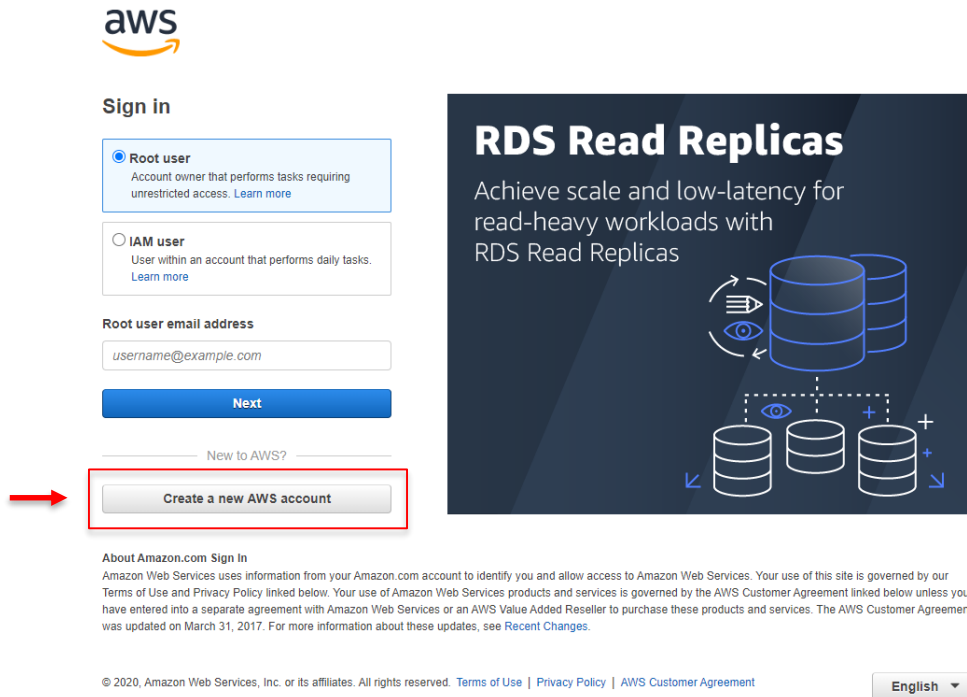
What next?

Before Aug 30, 2023:

- Get ready by revisiting your *vanilla* Python fundamentals!
- Python
 - All data structures: Lists, tuples, sets, dictionaries;
 - Handling Python arrays – 1D & 2D array;
 - Basic python methods used in the above data structures;
 - Slicing concepts;
 - Converting one data structure to other types
- Go over the definitions of *certain basic statistical quantities* – mean, median, mode, std. deviation, normal distribution, etc.

Sign up for AWS

- https://console.aws.amazon.com/console/home?nc2=h_ct&src=header-signin



aws

Sign in

☒ **Root user**
Account owner that performs tasks requiring unrestricted access. [Learn more](#)

☐ **IAM user**
User within an account that performs daily tasks. [Learn more](#)

Root user email address

Next

New to AWS?

Create a new AWS account

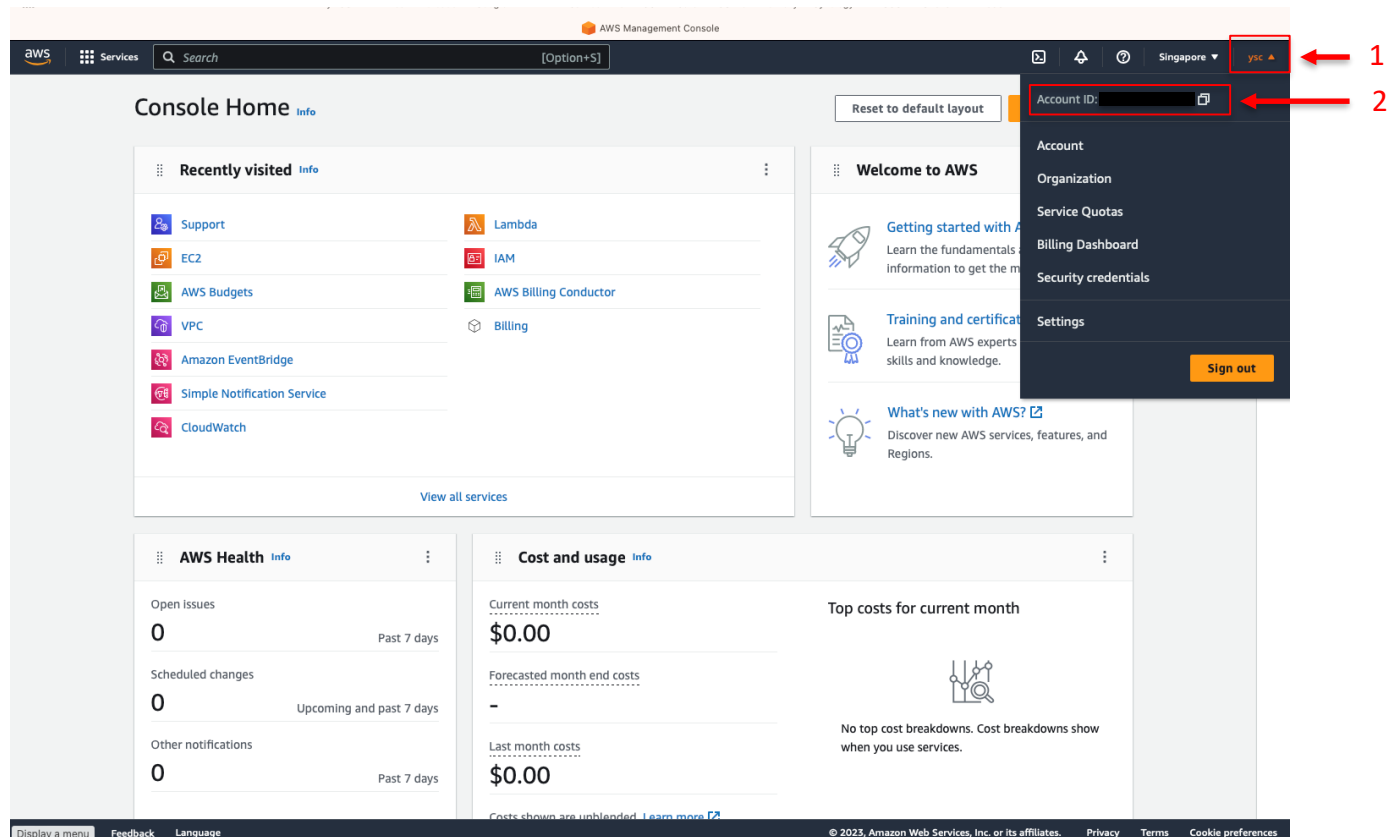
RDS Read Replicas
Achieve scale and low-latency for read-heavy workloads with RDS Read Replicas

About Amazon.com Sign In
Amazon Web Services uses information from your Amazon.com account to identify you and allow access to Amazon Web Services. Your use of this site is governed by our [Terms of Use](#) and [Privacy Policy](#) linked below. Your use of Amazon Web Services products and services is governed by the [AWS Customer Agreement](#) linked below unless you have entered into a separate agreement with Amazon Web Services or an AWS Value Added Reseller to purchase these products and services. The AWS Customer Agreement was updated on March 31, 2017. For more information about these updates, see [Recent Changes](#).

© 2020, Amazon Web Services, Inc. or its affiliates. All rights reserved. [Terms of Use](#) | [Privacy Policy](#) | [AWS Customer Agreement](#)

English

Get Account Number from AWS Management Console



Submit AWS and Zoom Information

- Canvas Ungraded Survey (under Quizzes)
 - Information for Labs
 - Submit 12-digit AWS account number
 - Indicate what kind of computer you will be using to complete the labs (e.g. Windows, Mac, or Linux)
 - Submit by Aug 27 (Sun), 11:59 pm

Annex – Quick review of basic Python methods

Python - List methods

(Note: This is not an exhaustive list)

FUNCTION	DESCRIPTION
<code>append()</code>	Add an element to the end of the list
<code>extend()</code>	Add all elements of a list to the another list
<code>insert()</code>	Insert an item at the defined index
<code>remove()</code>	Removes an item from the list
<code>pop()</code>	Removes and returns an element at the given index
<code>clear()</code>	Removes all items from the list
<code>index()</code>	Returns the index of the first matched item
<code>count()</code>	Returns the count of number of items passed as an argument
<code>sort()</code>	Sort items in a list in ascending order
<code>reverse()</code>	Reverse the order of items in the list
<code>copy()</code>	Returns a copy of the list
<code>len()</code>	Returns the length of the list (# of elements in the list)

DIY!

Create an empty list and demonstrate some of the above methods; Observe how elements are accessed!

Python – Tuples

(Note: This is not an exhaustive list)

Method	Description
count()	returns occurrences of element in a tuple
index()	returns smallest index of element in tuple
len()	returns Length of an Object
max()	returns largest element
min()	returns smallest element
reversed()	returns reversed iterator of a sequence
slice()	creates a slice object specified by range()
sum()	Add items of an iterable
tuple()	Creates a Tuple

DIY!

Create an empty tuple and demonstrate some of the above methods

Python – Sets

(Note: This is not an exhaustive list)

<i>Method</i>	<i>Description</i>
<code>add()</code>	Adds an element to the set
<code>clear()</code>	Removes all the elements from the set
<code>copy()</code>	Returns a copy of the set
<code>difference()</code>	Returns a set containing the difference between two or more sets
<code>discard()</code>	Remove the specified item
<code>intersection()</code>	Returns a set, that is the intersection of two other sets
<code>isdisjoint()</code>	Returns whether two sets have a intersection or not
<code>issubset()</code>	Returns whether another set contains this set or not
<code>issuperset()</code>	Returns whether this set contains another set or not
<code>pop()</code>	Removes an element from the set
<code>remove()</code>	Removes the specified element
<code>union()</code>	Return a set containing the union of sets
<code>update()</code>	Update the set with the union of this set and others

DIY!

Python – Dictionary

(Note: This is not an exhaustive list)

Method	Description
<code>clear()</code>	Removes all the elements from the dictionary
<code>copy()</code>	Returns a copy of the dictionary
<code>fromkeys()</code>	Returns a dictionary with the specified keys and values
<code>get()</code>	Returns the value of the specified key
<code>items()</code>	Returns a list containing a tuple for each key value pair
<code>keys()</code>	Returns a list containing the dictionary's keys
<code>pop()</code>	Removes the element with the specified key
<code>popitem()</code>	Removes the last inserted key-value pair
<code>update()</code>	Updates the dictionary with the specified key-value pairs
<code>values()</code>	Returns a list of all the values in the dictionary

DIY!

Python – Arrays

- Some useful array methods include:

append() insert() pop() reverse()
remove() index() count() extend() ... *and more*

- Numpy – Matrix operations

DIY!

Create an array of integers and demonstrate some of the above methods; Observe how elements are accessed!

Ethics in Computer Programming

Following **Code of Ethics in Computer Programming** is based on defunct International Programmer's Guild.

A programmer must...

- ...never create or distribute malware.

- ...never write code that is **obfuscated or intentionally difficult to follow**.

- ...never write **documentation that is intentionally confusing or inaccurate**.

- ...never reuse **copyrighted code** unless the proper license is purchased or permission is obtained.

- ...**acknowledge** (verbally and in source code comments) the work of other programmers on which the code is based, even if substantial changes are made.

- ...never write code that is **deliberately inefficient** with the intent of later claiming credit for making efficiency improvements.

- ...**never intentionally introduce bugs** with the intent of later claiming credit for fixing the bugs, or to stimulate the uptake of later versions.

- ...never write code that **intentionally breaks another programmer's code for the purpose of elevating one's status**.

- ...never hide known obstacles to a project's completion during any phase of development, especially the design phase.

- ...never **dishonestly downplay the difficulty of completing a project**.

Ethics in Computer Programming (Cont'd)...

- ...report any **illegal activities** of the employer.
- ...never **defame the profession**.
- ...**never falsely deny the presence of bugs**.
- ...never reveal the secret corporate knowledge of an employer.
- ...**never accept compensation from multiple parties for the same work** unless permission is given.
- ...never **perform competitive work without the employer's knowledge**.
- ...never conceal pertinent information from other members of the development team.
- ...never conceal from the employer their financial interest in development resources.
- ...never **conceal any conflict of interest that may affect the project**.
- ...never seek external profit from a project that was funded by a second party without permission. If permission is given to resell a product, the work should be discounted.
- ...never maliciously injure the reputation of an employer or members of the development team.
- ...never **misrepresent their knowledge, experience, or abilities**.
- ...never take credit for another's work.
- ...**never steal software, especially development tools**.
- ...never conceal the deficiencies of other programmers by writing code for them and allowing them to pass it off as their own work.
- ...**never install third-party applications without the user's permission**. Preferably not at all.
- ...stay current on the advancement of the field of Computer Science.
- ...never **force updates** on a user without their knowledge and approval.