

# Numpy and Pandas

Python for Economics

Morning, Aug 17, 2023

Zhenhua He

# The need for better number-crunching

Lists don't have many built-in methods for interacting with data.

The base Python data types also take up a lot of space compared to other languages.

The **Numpy** and **Pandas** modules offer powerful tools for improving performance when you're using lots of data and doing lots of operations on them.

# Arrays and Series: Arrays

Numpy Arrays supports common operations, such as arithmetic, on an element-by-element (or “vectorized”) basis.

Example:

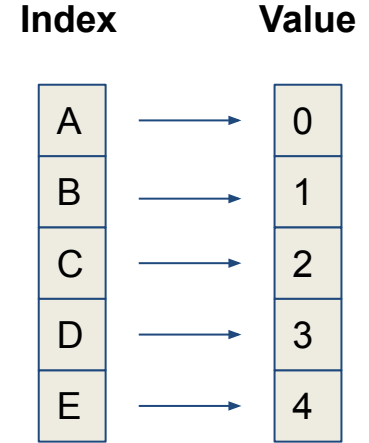
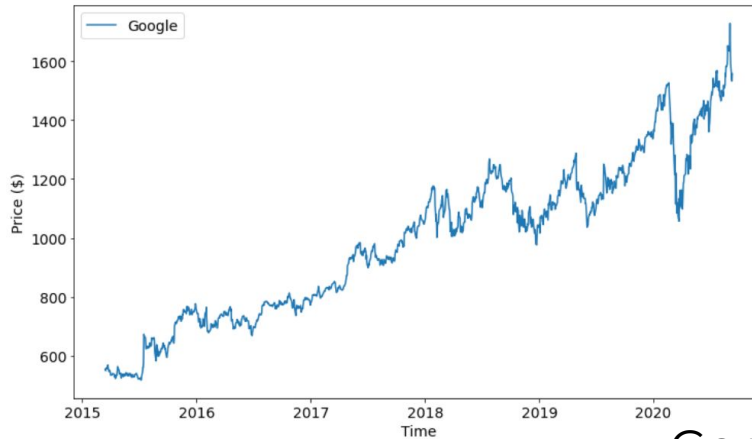
```
array C = array A + array B
```

This adds the elements of A and B pair-wise (Instead of concatenating the elements as would happen with lists).

Pandas Series and DataFrames further expand on this.

# Arrays and Series: Series

- One-dimensional labeled array
- Capable of holding any data type (integers, strings, floating point numbers, etc.)
- Example:



→ Go to notebook to practice

# Break Time Reminder Slide

10 minutes break



# Matplotlib

Python for Economics

Morning, Aug 17, 2023

Zhenhua He

# This Module

1. Line Plots
2. Scatter plots

# Matplotlib Setup

Matplotlib does lots of plotting for us—specifically using the “pyplot” submodule. Import a module with a nickname using `as`:

```
import matplotlib.pyplot as plt
import numpy as np
```

(The following slides show some plot style reference, but otherwise...)

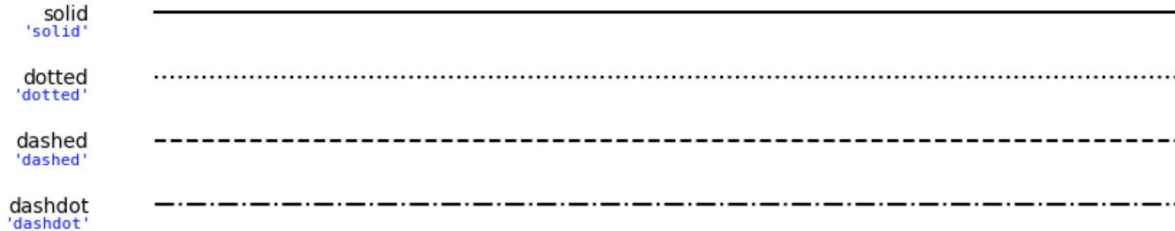
→ Go to notebook to practice



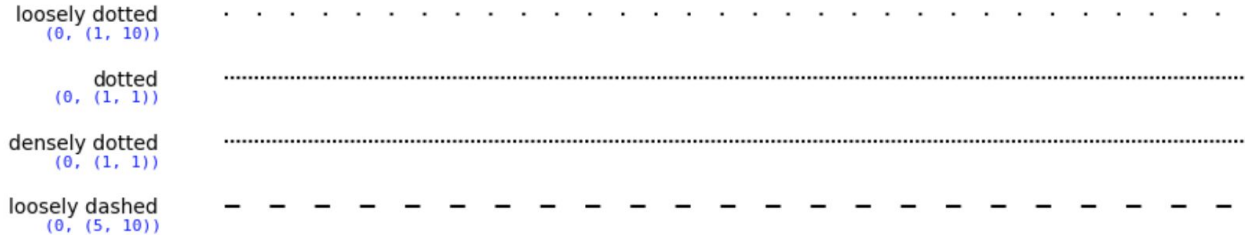
# Line Plots Styles

Simple line styles can be defined using the strings "solid", "dotted", "dashed" or "dashdot".

## Named linestyles



## Parametrized linestyles



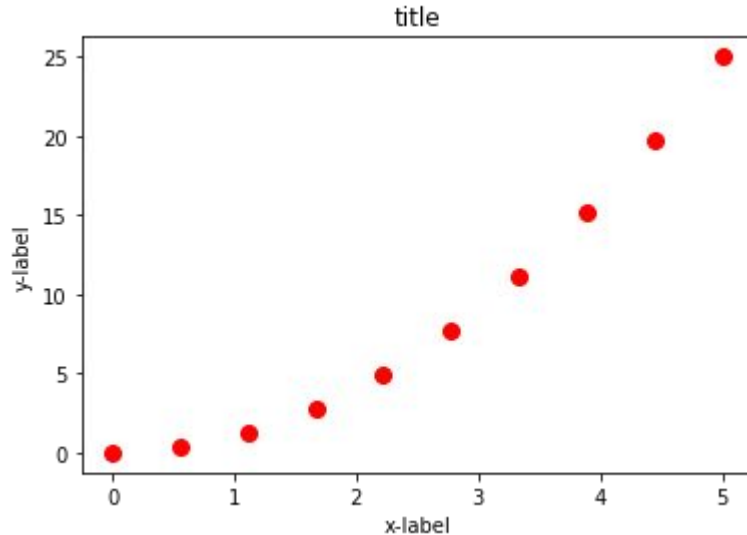
# Anatomy of a Plot

## Marker

- style
- size
- color

## Figure

- title
- xlabel
- ylabel



# Plot Marker symbols

marker	symbol	description
"."	•	point
","	,	pixel
"o"	●	circle
"v"	▼	triangle_down
"^"	▲	triangle_up
"<"	◀	triangle_left
">"	▶	triangle_right
"1"	⋮	tri_down
"2"	⋊	tri_up
"3"	↵	tri_left
"4"	↻	tri_right
"8"	⬢	octagon
"s"	■	square
"p"	⬠	pentagon
"P"	⊕	plus (filled)
"*"	★	star

# Break Time Reminder Slide

10 minutes break



# Data Manipulation

Python for Economics

Morning and Afternoon, Aug 17, 2023

Wesley Brashear, Josh Winchell



# This Module

1. Array/Series data selection
2. DataFrames
3. Columns and Filtering
4. DataFrame Methods

# Intro: Pandas VS NumPy

NumPy	Pandas
Faster mathematical operations ✓	Slower mathematical operations
Only supports integer index	Customized index ✓
must use structured arrays	Easily handles different data types ✓
better performance when number of rows is 500K or less	better performance when number of rows is 500K or more ✓
more complicated to read and write files	simpler to read and write more file formats ✓

# Array/Series Data Selection

Say we have a lot of data—and now that we have matplotlib we want to plot it... but only some of it.

Arrays provide us with ways to select data that are more nuanced than the options provided by plain lists.

→ Go to notebook to practice



# DataFrames

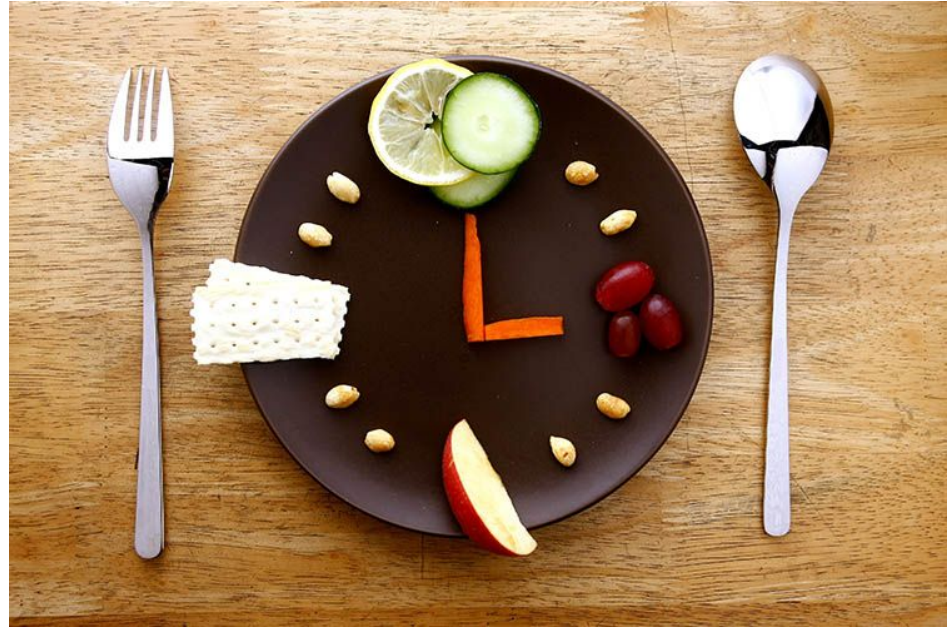
- Primary Pandas data structure
- A dict-like container for Series objects
- Two-dimensional size-mutable
- Heterogeneous tabular data structure

	C1	C2	C3	C4
A	0	x	0.1	True
B	1	y	2.4	False
C	2	z	1.9	True
D	NA	w	8.3	False
E	9	a	6.8	False

→ Go to notebook to practice

# Lunch Break Reminder Slide

1 hour break



# Columns and Filtering

Like arrays and series,  
DataFrames can be indexed,  
sliced, and filtered.

You can select specific rows  
and/or columns by name or  
based on some criteria.

*Say we only want columns 1  
and 4 when column 4 is "True"...*

		C1	C2	C3	C4
A	→	0	x	0.1	True
B	→	1	y	2.4	False
C	→	2	z	1.9	True
D	→	NA	w	8.3	False
E	→	9	a	6.8	False

→ Go to notebook to practice

# DataFrame Methods

There's a lot we can do via DataFrame Methods:

- Selecting/slicing and filtering
- Sorting or grouping by index or values
- Reading or writing to files
- Plotting
- Data cleanup
- Data merging

→ Go to notebook to practice

# Break Time Reminder Slide

10 minutes break



# Pandas Cheat Sheet (continued learning)

Data Wrangling  
with pandas  
Cheat Sheet  
<http://pandas.pydata.org>

## Syntax – Creating DataFrames

```
df = pd.DataFrame(
    {"a": [4, 5, 6],
     "b": [7, 8, 9],
     "c": [10, 11, 12]},
    index = [1, 2, 3])
Specify values for each column.

df = pd.DataFrame(
    [[4, 7, 10],
     [5, 8, 11],
     [6, 9, 12]],
    index=[1, 2, 3],
    columns=['a', 'b', 'c'])
Specify values for each row.
```

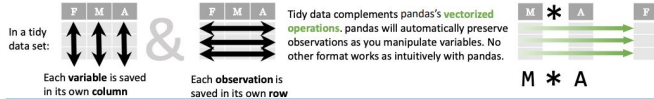
```
df = pd.DataFrame(
    {"a": [4, 5, 6],
     "b": [7, 8, 9],
     "c": [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [(1, 'd'), (2, 'e'), (3, 'f')],
        names=['n', 'v'])
Create DataFrame with a MultiIndex
```

## Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```
df = (pd.melt(df)
     .rename(columns={
         'variable': 'var',
         'value': 'val'})
     .query('val > 200'))
```

## Tidy Data – A foundation for wrangling in pandas



## Reshaping Data – Change the layout of a data set

```
df.sort_values('mpg')
Order rows by values of a column (low to high).

df.sort_values('mpg', ascending=False)
Order rows by values of a column (high to low).

df.rename(columns = {'y': 'year'})
Rename the columns of a DataFrame

df.sort_index()
Sort the index of a DataFrame

df.reset_index()
Reset index of DataFrame to row numbers, moving index to columns.

df.drop(columns=['Length', 'Height'])
Drop columns from DataFrame
```

```
pd.melt(df)
pd.pivot(columns='var', values='val')
pd.concat([df1, df2])
pd.concat([df1, df2], axis=1)
```

## Subset Observations (Rows)

```
df[df.Length > 7]
Extract rows that meet logical criteria.

df.drop_duplicates()
Remove duplicate rows (only considers columns).

df.head(n)
Select first n rows.

df.tail(n)
Select last n rows.

df.sample(frac=0.5)
Randomly select fraction of rows.

df.sample(n=10)
Randomly select n rows.

df.iloc[10:20]
Select rows by position.

df.nlargest(n, 'value')
Select and order top n entries.

df.nsmallest(n, 'value')
Select and order bottom n entries.
```

## Subset Variables (Columns)

```
df[['width', 'length', 'species']]
Select multiple columns with specific names.

df['width'] or df.width
Select single column with specific name.

df.filter(regex='regex')
Select columns whose name matches regular expression regex.

regex (Regular Expressions) Examples
```

Regex	Matches
`.`	Matches strings containing a period `.`
`.Length\$`	Matches strings ending with word `Length`
`.Sepal`	Matches strings beginning with the word `Sepal`
`.^x[1-3]\$`	Matches strings beginning with `x` and ending with 1,2,3,4,5
`.^(?!species)\$`	Matches strings except the string `Species`

Logic in Python (and pandas)		
<	Less than	!=
>	Greater than	df.column.isin(values)
=	Equals	pd.isnull(obj)
<=	Less than or equals	pd.notnull(obj)
>=	Greater than or equals	df.any(), df.all()

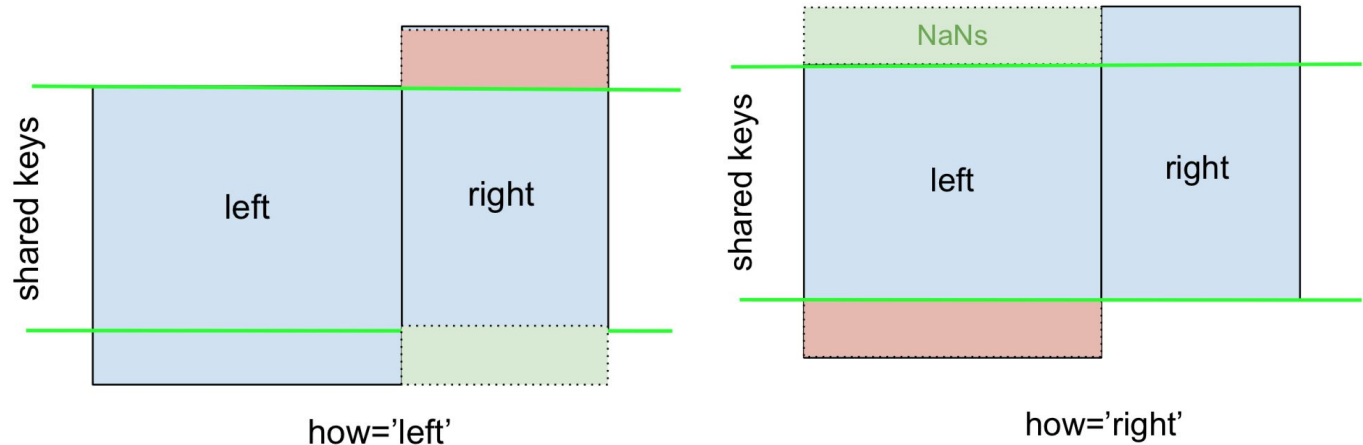
<http://pandas.pydata.org>. This cheat sheet inspired by RStudio Data Wrangling Cheatsheet (<https://www.rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheat-sheet.pdf>). Written by Iv Louie, @ibouie, @conaldata

[https://pandas.pydata.org/Pandas\\_Cheat\\_Sheet.pdf](https://pandas.pydata.org/Pandas_Cheat_Sheet.pdf)

# DataFrame - Merging Data

Merge DataFrame with a database-style join.

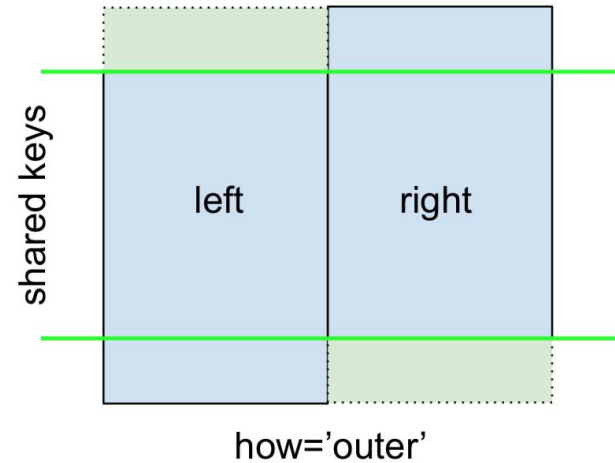
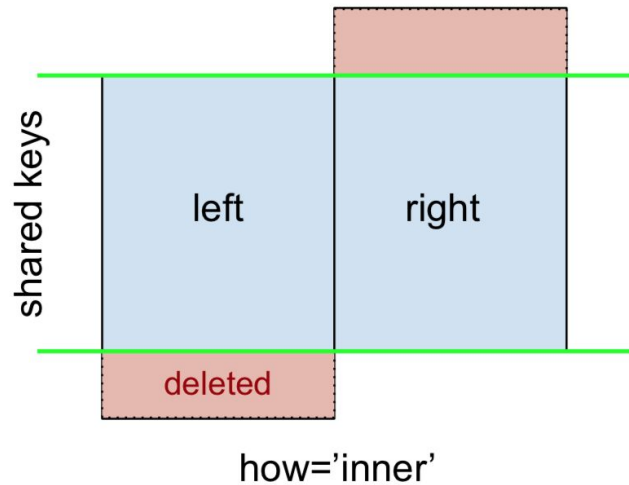
- left join
- right join



# DataFrame - Merging Data

Merge DataFrame with a database-style join.

- inner join
- outer join





# APIs

Python for Economics  
Afternoon, Aug 17, 2023  
Richard Lawrence

# This Module

1. JSON Format
2. Requests
3. FRED API
  
4. Capstone Project
5. ACES

# JSON - JavaScript Object Notation

- A text format for storing data
- language-independent (why they should know json)

JSON string examples:

```
'{"name":"Jack", "age":20, "major":"computer science"}'
```

```
'{ "args": {}, "data": "", "files": {}, "form": {  
"soup": "hot soup" }, .... }'
```

loads() function -> Python dictionary

# JSON module

- Python built-in module **json**
- `json.loads()` : converts JSON string to Python dictionary
- Example:

```
import json
text = "{ keys : values, ... }"
dict = json.loads(text)
```

# Requests

**Requests** library for HTTP activities.

Replicate the experience of visiting a web page,  
but in a Notebook instead of a Browser.

# FRED API

Accessing Federal Reserve Economic Data

# Web Scraping API Exercise

## **Fred API**

Retrieve economic data from the FRED® and ALFRED® websites hosted by the Economic Research Division of the Federal Reserve Bank of St. Louis

Reference: <https://fred.stlouisfed.org/docs/api/fred/>

# Get an API Key

[Register](#) and log into your fredaccount.stlouisfed.org user account and request your API Key.

Most web services require an API key to identify who owns a request.



# Break Time Reminder Slide

10 minutes break



# Capstone

Putting It All Together

# ACES

Supercomputing Cluster at Texas A&M

# ACES

## Accelerating Computing for Emerging Sciences

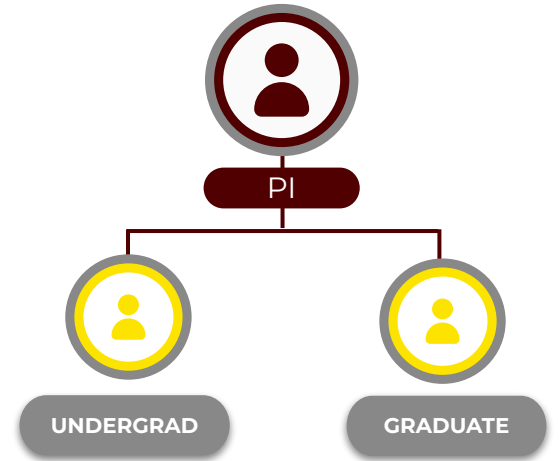
### Our Mission:

- Offer an accelerator testbed for numerical simulations and **AI/ML workloads**
- Provide consulting, technical guidance, and training to researchers
- Collaborate on computational and data-enabled research.



# Getting on ACES

- You must have an [ACCESS](#) account (we did this yesterday)
- PI's can apply for allocations directly
- Students will use our training allocation today
- Email us at [help@hprc.tamu.edu](mailto:help@hprc.tamu.edu) for questions, comments, and concerns.



**PIs** can apply for an account and sponsor accounts for their researchers.

# HPRC Portal



TEXAS A&M HIGH PERFORMANCE RESEARCH COMPUTING



Home User Services Resources Research Policies Events Training About Portal

Terra Portal

Grace Portal

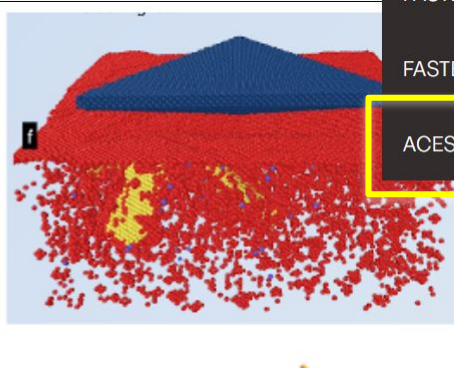
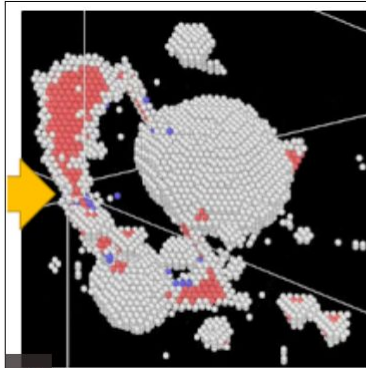
FASTER Portal

FASTER Portal (ACCESS)

ACES Portal (ACCESS)

## Quick Links

- New User Information
- Accounts
  - Apply for Accounts
  - Manage Accounts
- User Consulting
- Training
- Knowledge Base
- Software
- FAQ



# Accessing ACES via the HPRC Portal (ACCESS)

Log-in using your ACCESS credentials.

The screenshot shows the ACCESS portal interface. At the top left is the ACCESS logo, and at the top right is the 'Powered By CILogon' logo. A teal header bar contains a dropdown menu labeled 'Consent to Attribute Release'. Below this is a white box with the text: 'TAMU FASTER ACCESS OOD requests access to the following information. If you do not approve this request, do not proceed.' followed by a bulleted list: 'Your CILogon user identifier', 'Your name', 'Your email address', and 'Your username and affiliation from your identity provider'. Below the consent box is a teal header bar labeled 'Select an Identity Provider'. Underneath, there is a dropdown menu showing 'ACCESS CI (XSEDE)' with a question mark icon. Below the dropdown is a checkbox labeled 'Remember this selection' and a teal 'Log On' button. At the bottom of this section, it says 'By selecting "Log On", you agree to the [privacy policy](#).' At the very bottom of the page, there is a footer with links for 'For questions about this site, please see [FAQs](#) or send email to [help@cilogon.org](mailto:help@cilogon.org)' and 'Know your responsibilities using the CILogon Service. See [acknowledgments](#) for support for this site.'

The screenshot shows the ACCESS portal login page. At the top left is the ACCESS logo, and at the top right is the CILogon logo. The main heading is 'Login to CILogon'. Below this are two input fields: 'ACCESS Username' and 'ACCESS Password'. There is a checkbox labeled 'Don't Remember Login' and a teal 'Login' button. To the right of the login form, there is a teal box with the CILogon logo and the text 'CILogon facilitates secure access to CyberInfrastructure (CI)'. Below this are several links: 'If you had an XSEDE account, please enter your XSEDE username and password for ACCESS login', 'Register for an ACCESS Account', 'Forgot your password?', and 'Need Help?'. At the bottom of the page, there is a link that says 'Click Here for Assistance'.

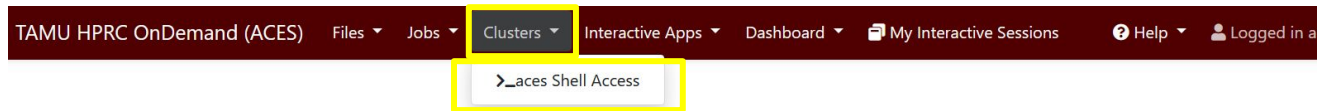
This is a close-up of the 'Select an Identity Provider' dropdown menu. The dropdown is highlighted with a yellow border and shows the option 'ACCESS CI (XSEDE)' with a question mark icon to its right.

Select the Identity Provider appropriate for your account. You need an ACCESS account, but can choose to log in with your TAMU NetID here.

# Shell access via the HPRC Portal

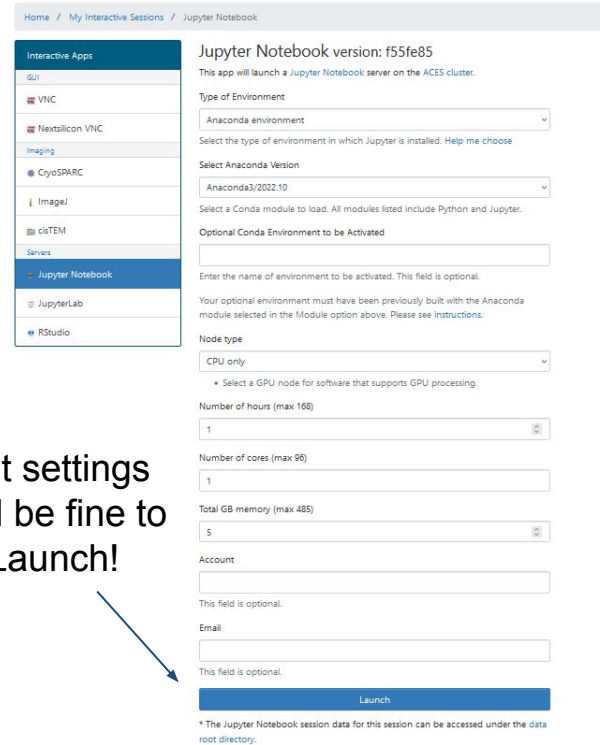
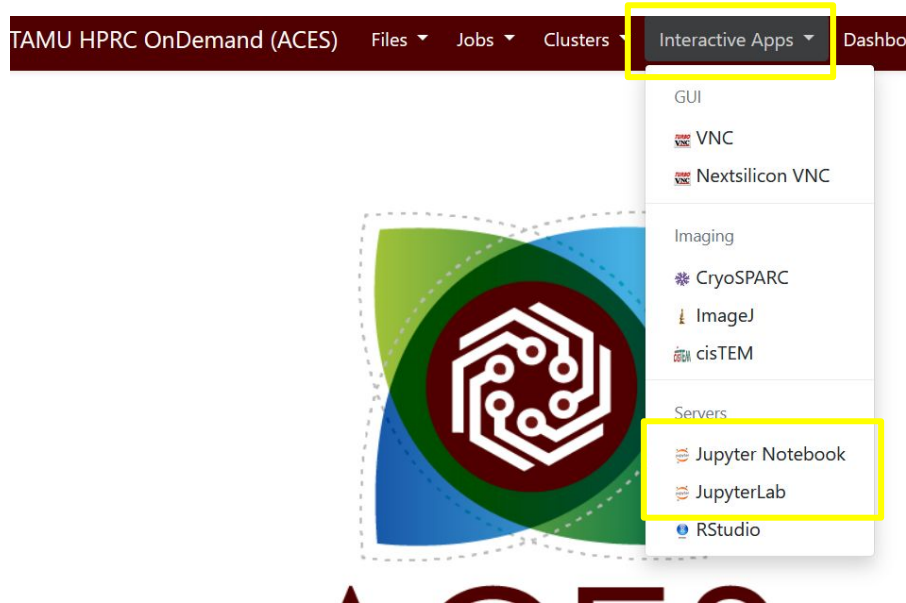
Access through (most) web browsers

-Top Banner Menu “Clusters” -> “Shell Access”





# Accessing Jupyter Notebooks on ACES



Default settings should be fine to start. Launch!

Session was successfully created. [X]

Home / My Interactive Sessions

Interactive Apps
GUI
VNC
Nextsilicon VNC
Imaging
CryoSPARC
ImageJ

### Jupyter Notebook (3773)

1 node | 1 core | Starting

**Created at:** 2023-08-10 16:16:36 CDT Delete

**Time Remaining:** 58 minutes

**Session ID:** 3c563bd9-302d-4827-a450-b6183d84a50d

Your session is currently starting... Please be patient as this process can take a few minutes.

Session was successfully created. [X]

Home / My Interactive Sessions

Interactive Apps
GUI
VNC
Nextsilicon VNC
Imaging
CryoSPARC
ImageJ
ImageJ-TEAM

### Jupyter Notebook (3773)

1 node | 1 core | Running

**Host:** >\_ac956 Delete

**Created at:** 2023-08-10 16:16:36 CDT

**Time Remaining:** 58 minutes

**Session ID:** 3c563bd9-302d-4827-a450-b6183d84a50d

[Connect to Jupyter](#)

