# DATA SCIENCE TOOLS

*Katherine Bueno*

*Data Analyst/ Data Scientist at AKQA*

# LET'S DISCUSS THE CURRENT LESSON OBEJCTIVES

▸ Identify the data science toolkit we'll be using in class

▸ Learn how to navigate Git and the Command Line

▸ Download the course Git repository and practice some git commands

▸ Make a probability and odds IPython notebook (time permitting)

▸ Slides will be available on the course repository for your reference.

# INTRO

‣ Name

‣ Familiarity with programming, languages used

‣ Familiarity with UNIX/command line tools

‣ What you do for work/why you're taking data science

Please post your github username in the slack channel so I can add you to our demo repository!

# TOOLS OF THE TRADE

# TOOLS OF THE TRADE

‣ Today we are going to review some of the tools we use in data science.

‣ We'll see how they fit into the wider programming environment.

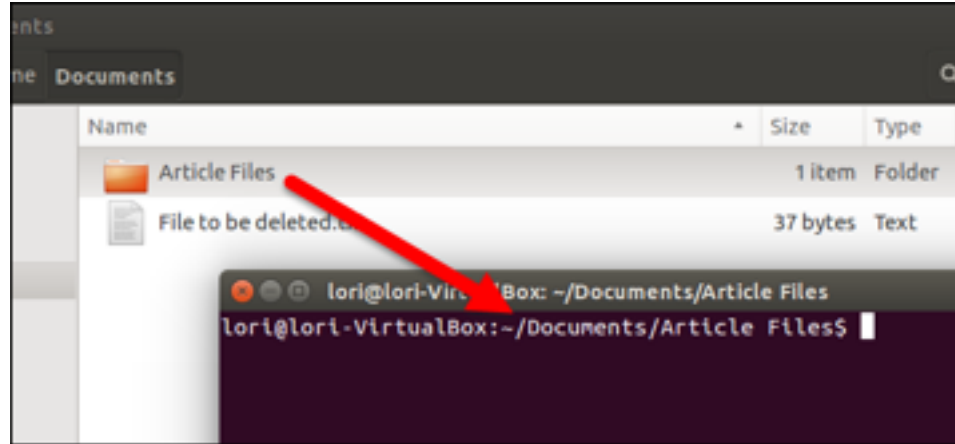‣ We'll start with the command line.  This is your portal to your computer and the outside world.

# LOCAL MACHINE

‣On your local computer, you have a variety of tools at your disposal.

    ‣Text editor

    ‣Programs/tools

    ‣Your files



‣All of these can be accessed through the terminal or through a GUI (Graphical User Interface).

‣You can navigate your files through the terminal or through Finder.

# DATA SCIENCE TOOLS

Outside World

Local Machine

Terminal/
Command Line

# COMMAND LINE

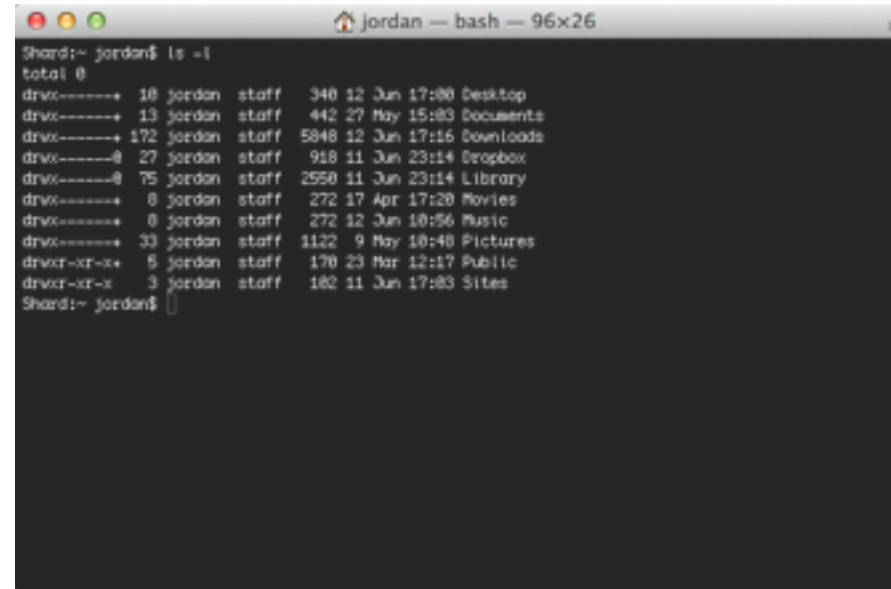# COMMAND LINE

‣ Let's walk through a few very basic commands.

    ‣ `cd`

    ‣ `pwd`

    ‣ `$home`

    ‣ `mkdir`

    ‣ `open`



‣ We can access many tools with the terminal.  Let's walk through a few.

# DATA SCIENCE TOOLS

# TEXT EDITORS

# TEXT EDITORS

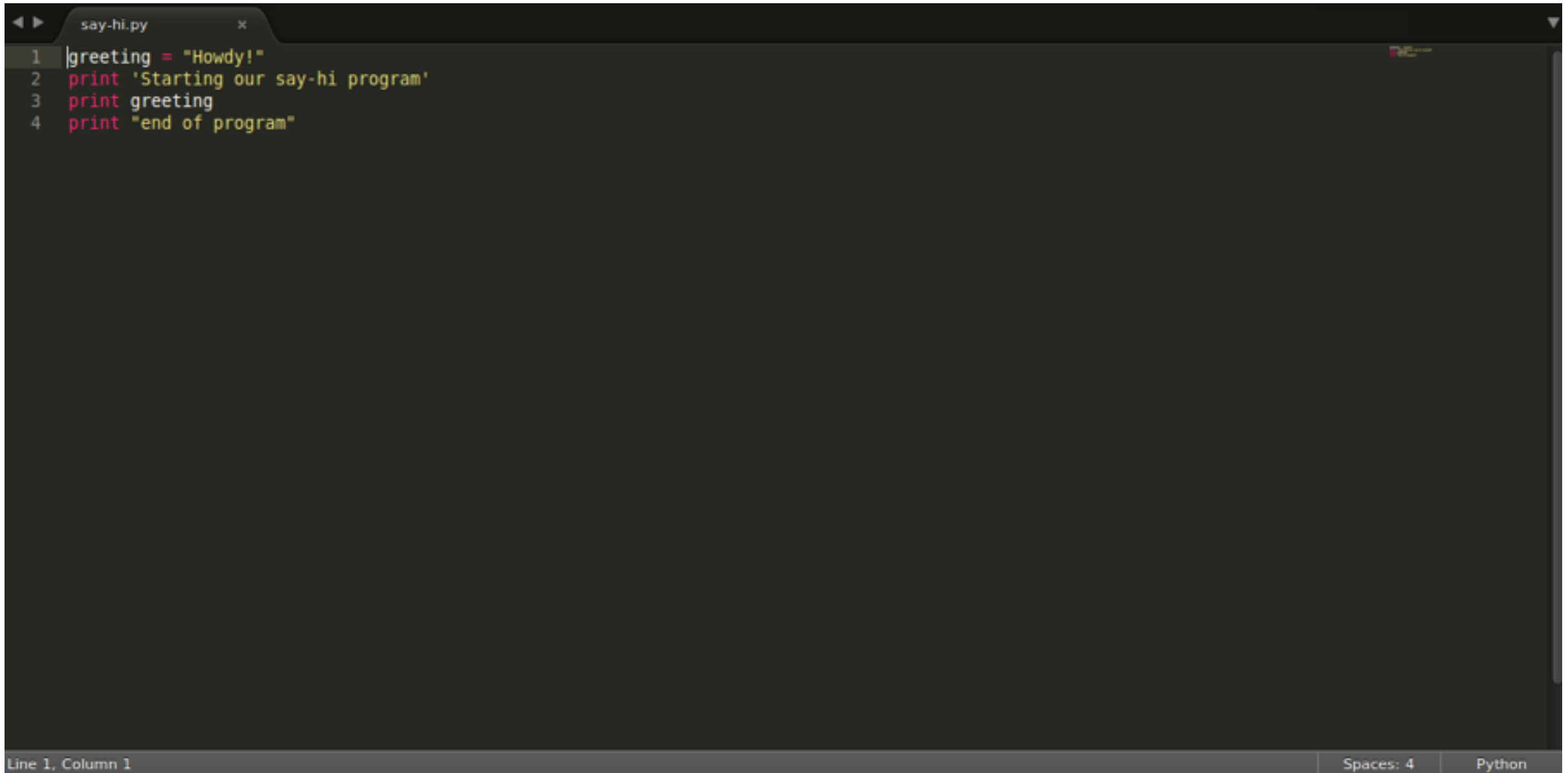‣So far, we've used iPython Notebooks in place of a text editor.
‣However, there are many options available
    ‣eMacs
    ‣Vim
    ‣Sublime Text
    ‣IDE (Integrated Development Environments) such as PyCharm or Eclipse
    ‣Many programmers love to [argue at length about which is the 'best' one.](#)
‣Let's see what Sublime Text looks like with Python!

# TEXT EDITORS

# TEXT EDITORS

‣ Open "say-hi.py", found in the lesson-02 folder of the class repo, in Sublime Text to see it for yourself.

‣ Check out the difference between saving a text file in MS Word and using sublime text

# DATA SCIENCE TOOLS

Outside World

Local Machine

Terminal/
Command Line

open, mkdir,
cd, rm

Your Files

Edit

Text Editor,
e.g. Sublime Text

# ACTIVITY: KNOWLEDGE CHECK

**EXERCISE**

## ANSWER THE FOLLOWING QUESTIONS

1. Who has a text editor they like that I didn't talk about? Why do you like it?
2. What are some good reasons to use a programming text editor instead of Microsoft Word or Google Docs?

## DELIVERABLE

Answers to the above questions

# IPYTHON NOTEBOOK

# IPYTHON NOTEBOOK

‣Where does iPython Notebook fit in?

‣We can refer to the iPython Notebook docs to get a better idea: the notebook combines the console, web apps, and markdown to capture the whole computation process.

‣IPython notebooks combine three components in sequential 'cells' that can be edited and run in any order:
   ‣Text (Markdown/HTML/Plain text)
   ‣Code
   ‣Output (plots, command line output, etc)

# IPYTHON NOTEBOOK

▸ IPython notebook demo:
  https://nature.tmpnb.org

# ACTIVITY: KNOWLEDGE CHECK

**ANSWER THE FOLLOWING QUESTIONS**

EXERCISE

1. What are the three components of IPython notebooks?

**DELIVERABLE**

Answers to the above questions

# PYTHON PACKAGES

# PYTHON PACKAGES
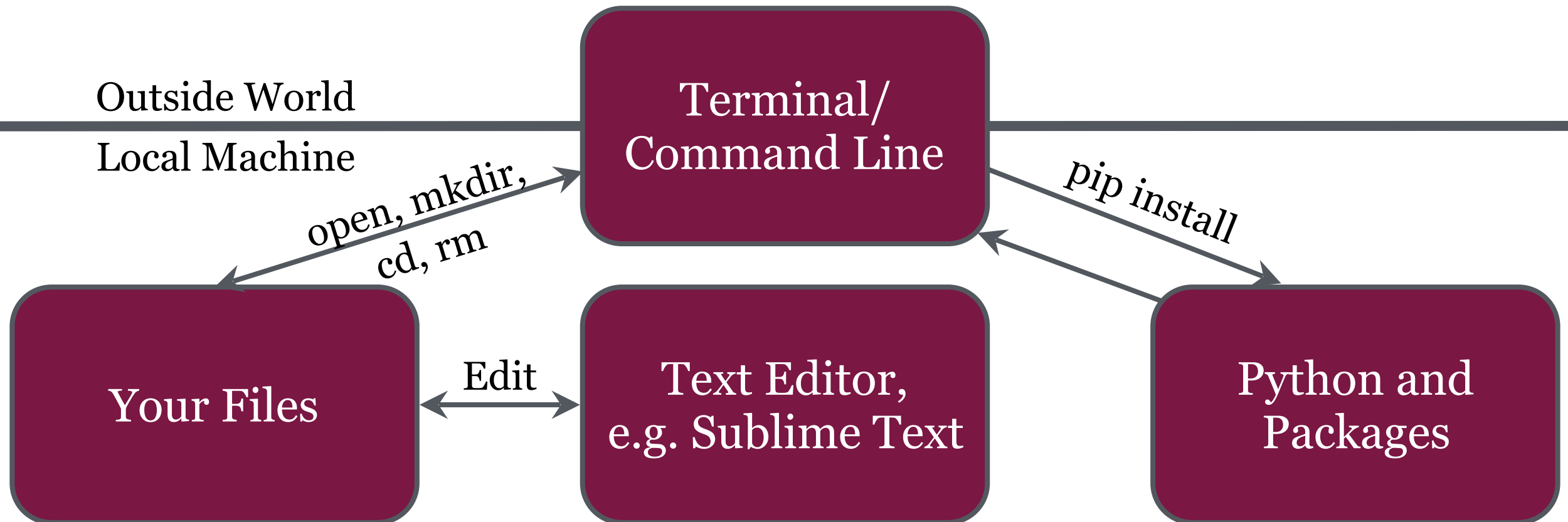
‣ The terminal allows us to run programs and reach out to the outside world.

‣ We can add programs and packages as needed.

‣ To add Python packages, we use a tool called *pip*.

‣ Let's `pip install` a package with the command line. We'll install Beautiful Soup, a HTML/XML parsing package.

```
pip install beautifulsoup4
```
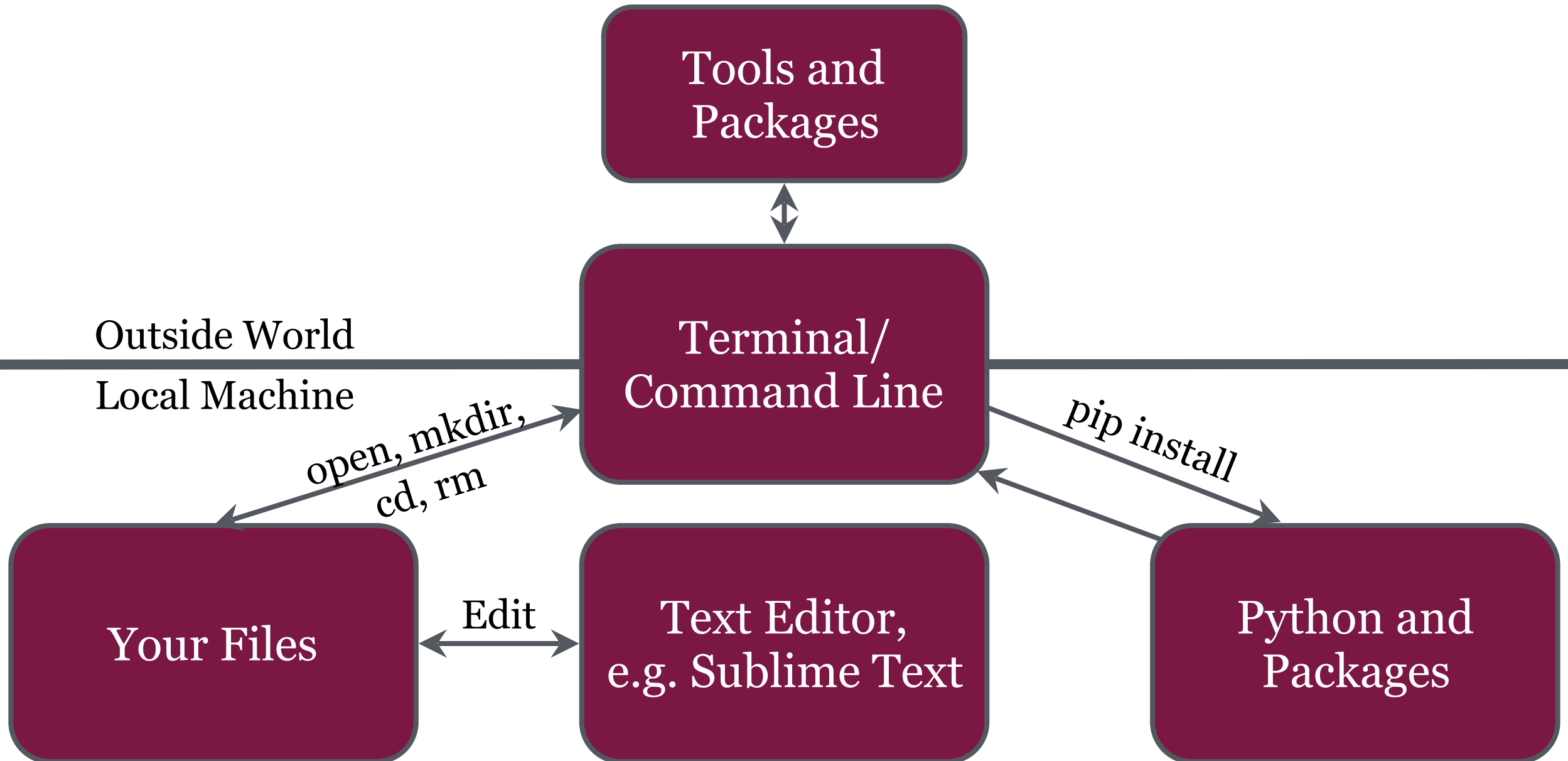
# DATA SCIENCE TOOLS

# THE OUTSIDE WORLD

# THE OUTSIDE WORLD

‣ The command line also allows you to download and use other tools and packages.

‣ There are many tools for different purposes available in the outside world.

# DATA SCIENCE TOOLS

# THE OUTSIDE WORLD

‣ As we saw with pip, the command line can connect us to the outside world. This becomes more important for data.

‣ We may have HIPAA protected data. This means we can't leave this sensitive data on our *local* machine (i.e. laptop).

‣ We need to communicate with a *remote* machine (i.e. server) to access the data via command line.

‣ Let's see a demonstration of this.

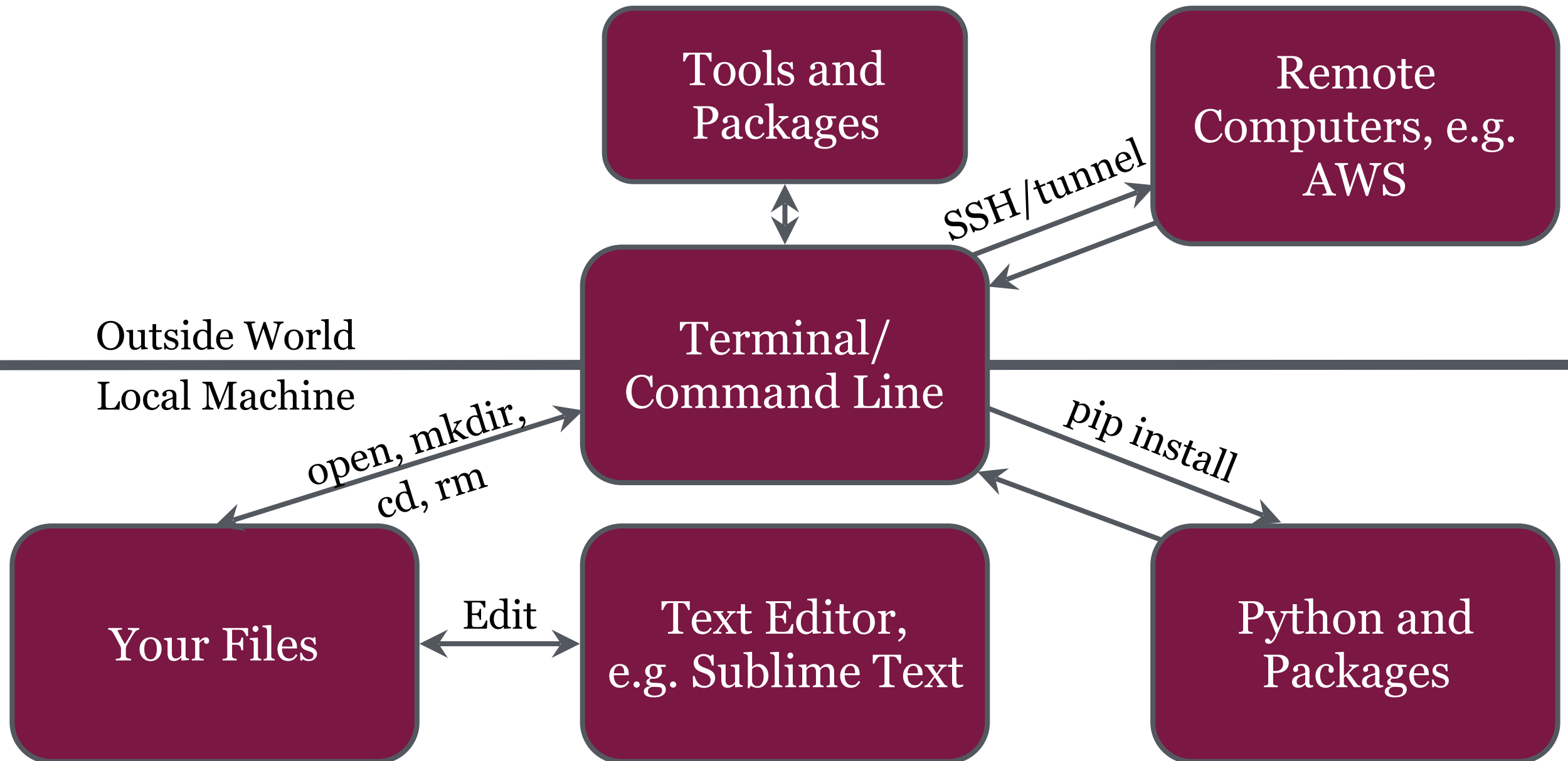# THE OUTSIDE WORLD

‣ AWS EC2 demo
   ‣ SSH
   ‣ Local database
   ‣ Text files
‣ S3 Demo

# DATA SCIENCE TOOLS

# BREAK

# 10 mins

# GIT

# GIT

‣ Version control is necessary when working on complex projects.

‣ Git is a way of tracking changes we've made to our programs that allows us to go back in time to fix errors.

‣ Combined with Github, Git is a powerful tool for collaborating with colleagues. You can work on different aspects of projects simultaneously and merge the changes together seamlessly.

‣ There are many different ways to use these tools.
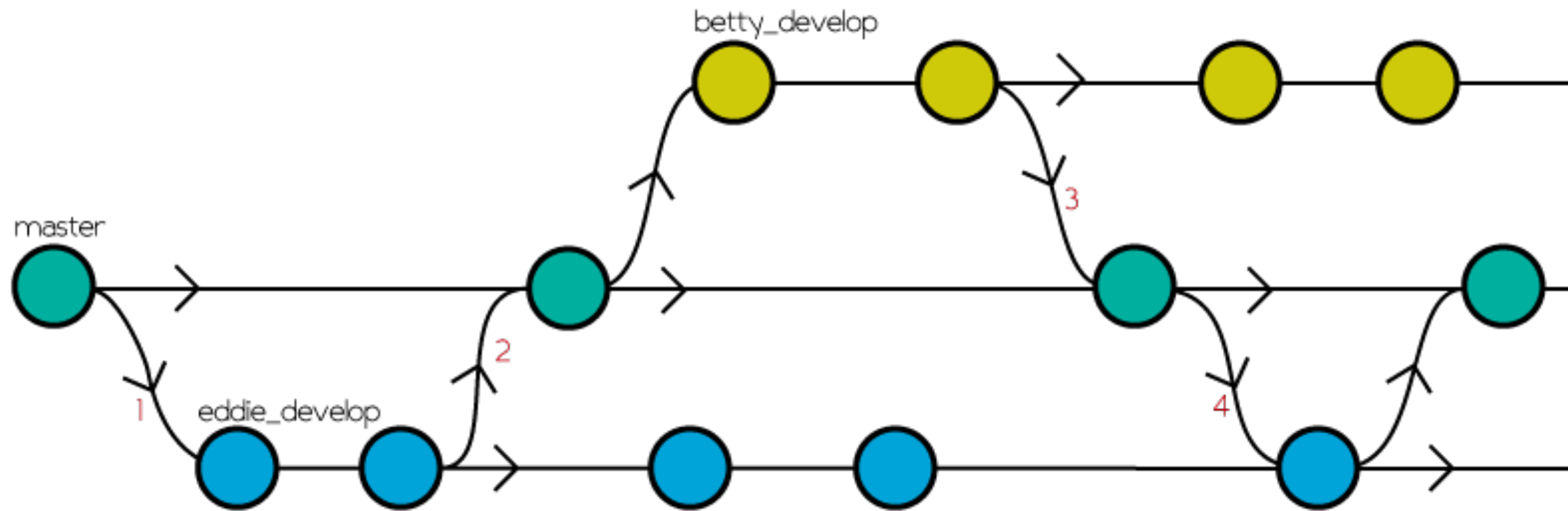
# GIT

‣ Let's see an example of using Git and Github.

‣ There are three primary commands we'll use.

    ‣ `git add`

    ‣ `git commit`

    ‣ `git push`

‣ When a colleague wants to implement our change, we may use the command `git pull`.

# GIT

# DATA SCIENCE TOOLS

# ACTIVITY:  KNOWLEDGE CHECK

**ANSWER THE FOLLOWING QUESTIONS**


EXERCISE

1. How are Python software modules shared?
2. What tools do you use to find and get new Python software?

**DELIVERABLE**

Answers to the above questions

# GIT AND COMMAND LINE

# ACTIVITY: GIT AND COMMAND LINE

**DIRECTIONS (30 minutes)**

**EXERCISE**

1. https://guides.github.com/activities/hello-world/
2. https://try.github.io/
3. Don't forget to paste your github username in the slack channel if you haven't done so already

**DELIVERABLE**

Github username

# Q & A

# BREAK

# 10 mins

# GIT EXERCISE

# GIT

▸ **Goal:** to fork and clone our own repositories on github, make sure we can edit a trivial file, and push it back up.

▸ **Secondary goal:** become comfortable wrestling the Git monster.

*When the office git expert has to come fix everything*

source: http://wheningit.tumblr.com/post...

# BRINGING IT ALL TOGETHER

## Activity (20 mins)

1. Break into groups of 4
2. Fork and clone the class repository at https://github.com/DAT-44-...
3. Create a text file in students/firstnamelastname.txt with your first name, last name, favorite food, and first place you would go if someone gave you a free plane ticket anywhere in the world.
4. 'Git add' the text file to your local repository
5. 'Git commit -m "type a commit message here"'
6. 'Git push origin'

## DELIVERABLE

A link to your text file on github in the slack channel

**EXERCISE**

# ODDS AND PROBABILITY

# ACTIVITY: ODDS & PROBABILITY

**EXERCISE**

## DIRECTIONS (20 minutes)

Some of you may already be familiar with odds and probability.

1. We will use the starter code in lesson-02 of the class repo to review the concepts of odds and probability.

## DELIVERABLE

Answer the questions in the notebook

# TOPIC REVIEW

# REVIEW

▸ What are some common data science tools?

▸ Why are these tools useful?

▸ Any other questions?

# BEFORE NEXT CLASS:

Practice Git, review today's lesson, and schedule office hours with me, if needed.

# EXIT TICKET

**DON'T FORGET TO FILL OUT YOUR EXIT TICKET**

=Will send link in slack