# DATA SCIENCE TOOLS

**Angelo Klin**

Katra Analytics

# LEARNING OBJECTIVES

◉ Identify the Data Science toolkit

◉ Navigate Git and the Command Line

◉ Describe Probability vs Odds

# PRE-WORK

# PRE-WORK REVIEW

- Use descriptive statistics to understand your data

# DATA SCIENCE TOOLS

# LET'S DISCUSS THE CURRENT LESSON OBJECTIVES

◉ Identify the Data Science toolkit

◉ Navigate Git and the Command Line

◉ Describe Probability vs Odds

# TOOLS OF THE TRADE

## LOCAL MACHINE

- On your local computer, you have a variety of tools at your disposal.
  - Text Editor
  - Programs, packages and tools
  - Your files

- All of these can be accessed through a **Terminal** or through a **GUI** (Graphical User Interface)

- You can navigate your files through the Terminal or through Finder/ Explorer

# TOOLS OF THE TRADE

- Today we are going to review **some** of the tools we use in Data Science

- We will see how they fit into the wider environment

- We will start with the Command Line
  - This is your portal to your computer and the outside world

# DATA SCIENCE TOOLS

Outside World

Local Machine

Terminal /
Command Line

# COMMAND LINE
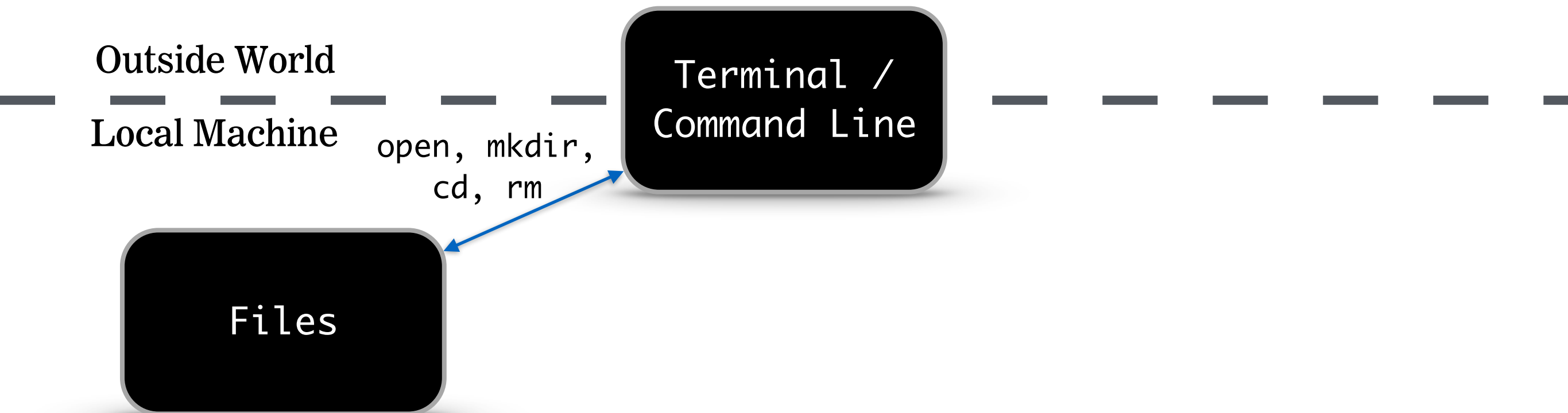
# COMMAND LINE

- We can access many tools with the Terminal

- Let's walk through a few commands
  - cd
  - pwd
  - mkdir
  - open

# DATA SCIENCE TOOLS

Outside World

Local Machine

Terminal /
Command Line

open, mkdir,
cd, rm

Files

# TEXT EDITORS

# TEXT EDITORS

- So far, we have used Jupyter Notebooks in place of a text editor

- However, there are many options available
  - Vim
  - Sublime Text
  - Atom

- Let's see how Python looks with Syntax Highlight

# TEXT EDITORS

```python
BOARD_SIZE = 8
class BailOut(Exception):
    pass

def validate(queens):
    left = right = col = queens[-1]
    for r in reversed(queens[:-1]):
        left, right = left - 1, right + 1
        if r in (left, col, right):
            raise BailOut

def add_queen(queens):
    for i in range(BOARD_SIZE):
        test_queens = queens + [i]
        try:
            validate(test_queens)
            if len(test_queens) == BOARD_SIZE:
                return test_queens
            else:
                return add_queen(test_queens)
        except BailOut:
            pass
    raise BailOut

queens = add_queen([])
print queens
print "\n".join(". " * q + "Q " + ". " * (BOARD_SIZE - q - 1) for q in queens)
```

# TEXT EDITORS

- Open the lesson 05 folder of the class repository and open the files
  - ~/lessons/lesson-05/code/say-hi.py
  - ~/lessons/lesson-05/code/eight-queens.py

- **NOTE**: These are Python source code, NOT Jupyter Notebooks!

# DATA SCIENCE TOOLS

Outside World

Local Machine

Terminal /
Command Line

open, mkdir,
cd, rm

Files

Text Editor

edit

# ACTIVITY: KNOWLEDGE CHECK

**EXERCISE**

## DIRECTIONS: ANSWER THE FOLLOWING QUESTIONS

1. What is a text editor?

2. Can you name any other examples?

# INTRODUCTION

# JUPYTER NOTEBOOK

# JUPITER NOTEBOOK

- Where does Jupyter Notebook fit in?
  - "The Jupyter Notebook is a web application that allows you to create and share documents that contain live code, equations, visualisations and explanatory text."

- Jupyter notebooks combine
  - The console
  - Web application
  - Markdown to capture the whole computation process

# PYTHON PACKAGES
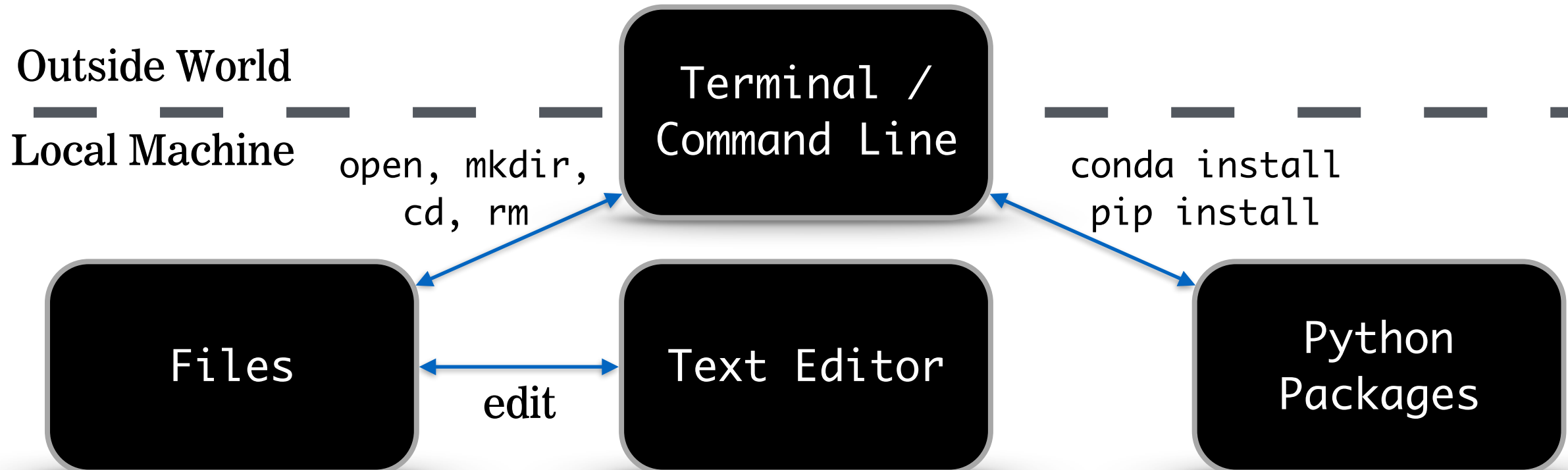
# PYTHON PACKAGES

◉ We can add programs and packages as needed

◉ To add Python packages, we use tools like conda and pip

◉ To install Beautiful Soup, a HTML/XML parsing package

◉ `conda install beautifulsoup4`

◉ `pip install beautifulsoup4`

# DATA SCIENCE TOOLS

Outside World

Local Machine

**Terminal / Command Line**

open, mkdir, cd, rm

conda install pip install

**Files**
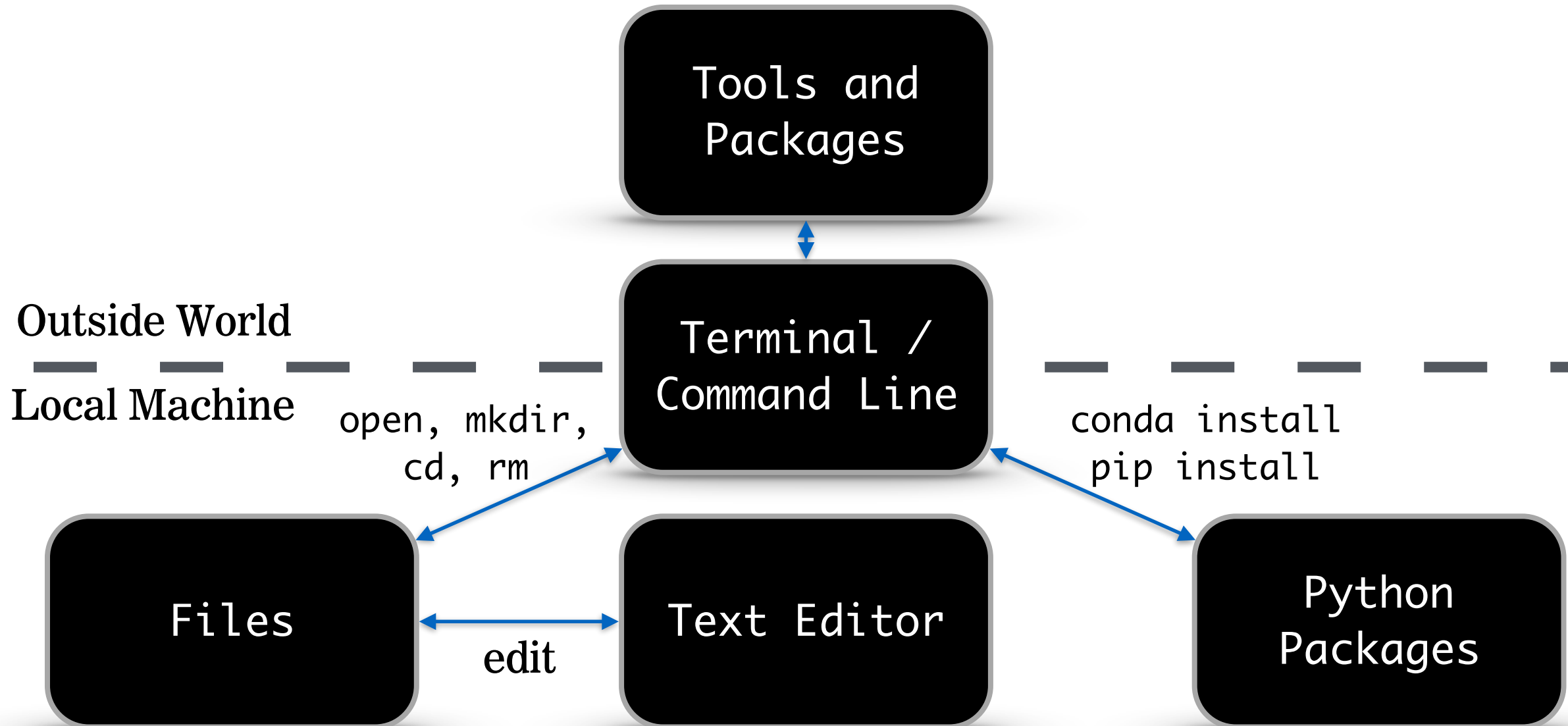
edit

**Text Editor**

**Python Packages**

# 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
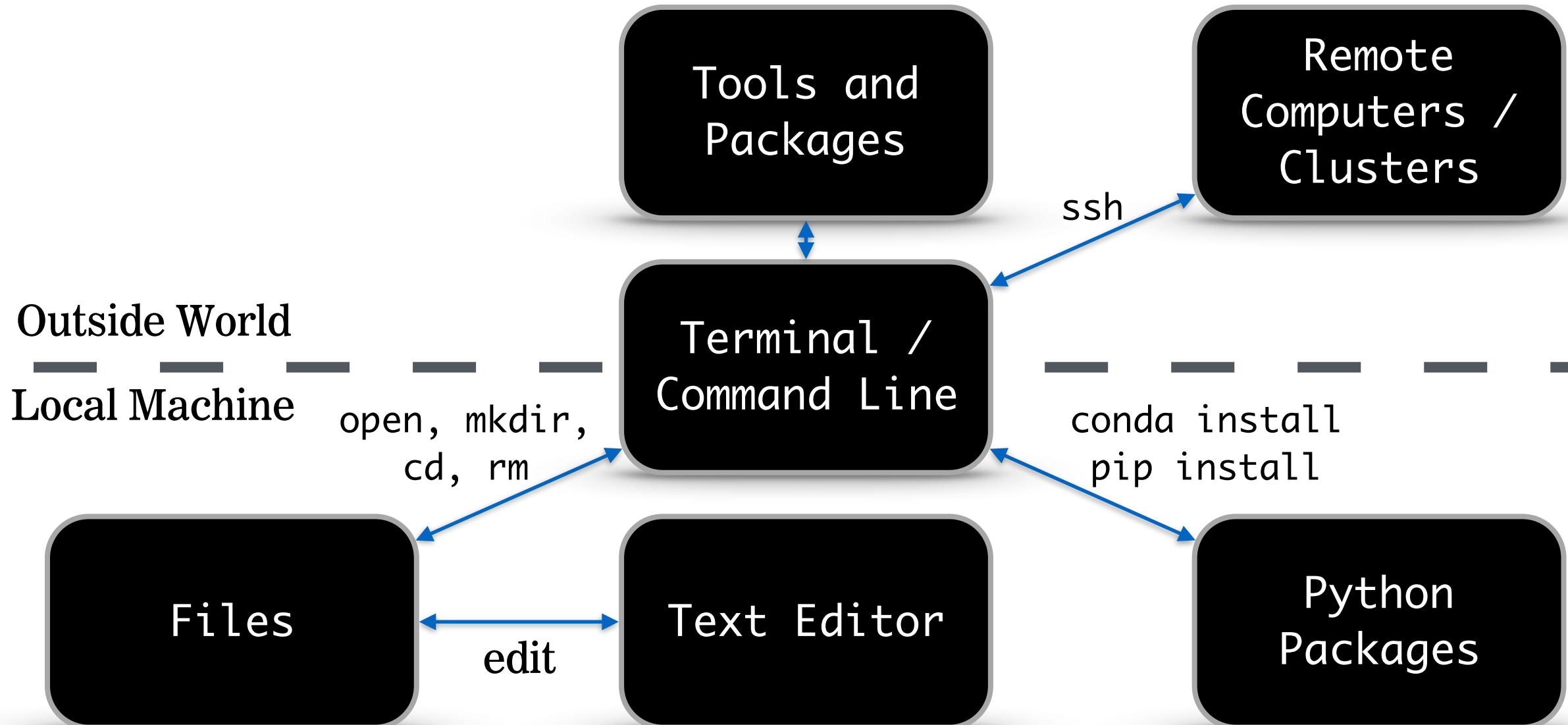
# DATA SCIENCE TOOLS

# THE OUTSIDE WORLD

- As we saw with `conda`/`pip`/`git`, 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

- 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

# DATA SCIENCE TOOLS

# GIT

# GIT

- Version control is necessary when working on complex projects

- Git is a way of tracking changes we have 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
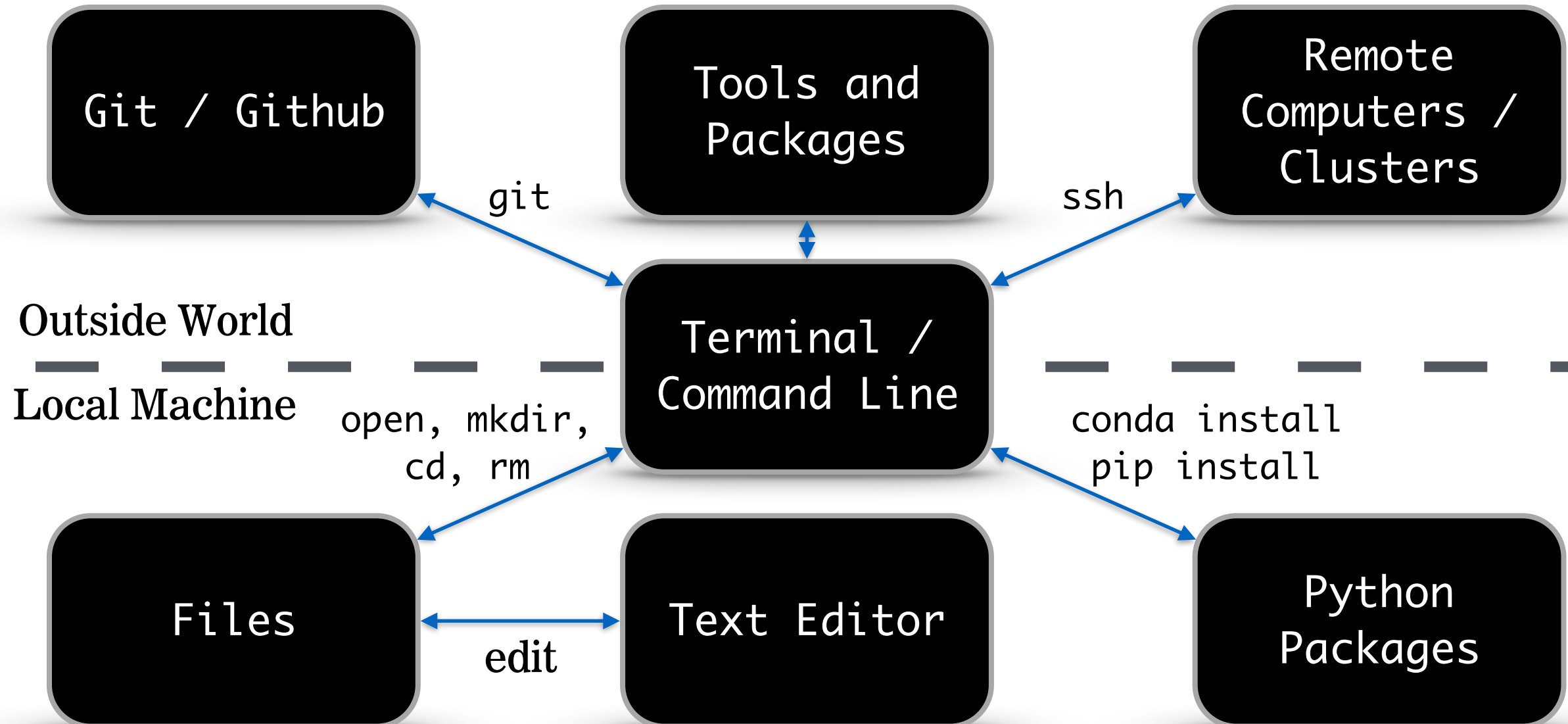
- There are many different ways to use these tools

# GIT

- Let's see an example of using Git and Github

- There are three primary commands we will use
  - `git add`
  - `git commit`
  - `git push`

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

# DATA SCIENCE TOOLS

# ACTIVITY: KNOWLEDGE CHECK

**DIRECTIONS: ANSWER THE FOLLOWING QUESTIONS**

1. What is a GUI?

2. What is the Command Line?

3. What are the big advantages of using the Command Line over a GUI?

**EXERCISE**

# GIT AND THE COMMAND LINE

# ACTIVITY: GIT AND THE COMMAND LINE

**EXERCISE**

## DIRECTIONS (35 MINUTES)

1. Let's review the exercises from Codecademy Python

2. Let's review the exercises from the GA's Command Line Tutorial

3. Are there any questions?

# ODDS AND PROBABILITY

# ACTIVITY: ODDS AND PROBABILITY

**DIRECTIONS (20 MINUTES)**

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

2. We will use the starter code in lesson 05 of the class repository to review the concepts of odds and probability.

    a.  ~/lessons/lesson-05/code/starter/starter-5.ipynb

**EXERCISE**

# TOPIC REVIEW

## TOPIC REVIEW

◉ What are some common Data Science tools?

◉ Why are these tools useful?

◉ Any other questions?

DATA SCIENCE

# BEFORE NEXT CLASS

# DUE DATE

- ◉ Project
  - ◉ Unit Project 2

# Q & A