# WELCOME TO GA!

## ORIENTATION



### **ABOUT ME**

- PhD in EE (learning and control in systems with uncertainty), MS in Statistics
- Researcher at IBM Watson
- Chief Technology Officer at a Data Science Healthcare startup
- Love to travel, hike and take photos

Hi!

## STUDENT-FACING SUPPORT

## **STUDENT SERVICES**

StudentServicesLA@generalassemb.ly











# GENERAL ASSEMBLY'S MISSION IS TO BUILD OUR COMMUNITY BY TRANSFORMING MILLIONS OF THINKERS INTO CREATORS.

# COMMUNITY OF INDIVIDUALS EMPOWERED TO PURSUE THE WORK WE LOVE.

## **ALL YOU CAN LEARN BUFFET**

## **CLASSES & WORKSHOPS**

Supplement and focus your learning during and after the course with a \$150 credit for 1 year from your start date

## STUDENT DISCOUNTS

Our students and alumni get free or discounted access to software and services from many vendors

## **GA COMMUNITY**

Expand your network at GA events, Meetups, and continue to collaborate and network with your classmates



### FEEDBACK / SUPPORT

- Access to Instructional Team: office hours, in class support
- Exit Tickets
- Mid-Course Feedback
- End of Course Feedback



### **Exit Tickets**

### Welcome to Today's Exit Ticket

This surveyin designed in halpymer instructor as self as the CLE vacer to understandflow you felfatiout the lesson and how offering the instruction was stoley. Year howest feedback will help you improve your howevery reprehense to red time.

\*Required



Please, write your full name, "
Counse/Colort # 1
What's the 'excents and of ' If you're not sure, despectable with your instructor or TA.   •
What was the took of the leasen? * If you're not sure, desearcheckwith your instructor or T.A.
My instructional team use effective in helping me while the barring objectives for this lesson ${}^{\bullet}\!$
1 1 0 4 0 5 7 1 9 10
Not iffective or a contract or a contractive
Interesonagents was well-organized and surricient time was given for each activity.
12245671912
De Net Agree + + + + + + + + + + + + Strengly Agree
I feel prepared to continue practicing this skill outside of class.
1224147100
Do Not Agree o o o o o o o o o o Strongly Agree

### GitHub



### You've been added to the generalassembly-studio organization!

Here are some quick tips for a first-time organization member.

 Use the switch context button in the upper left corner of this page to switch between your personal context (malloryjacobs) and organizations you are a member of.



×

After you switch contexts you'll see an organization-focused dashboard that lists out
organization repositories and activities.

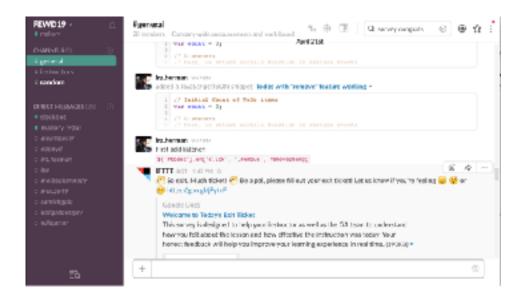


Welcome to GitHub! What's next? (15 days ago)

Create a repository
Tell us about yourself
Browse interesting repositories
Follow @github on Twitter

Slack

## #ds-sm-15

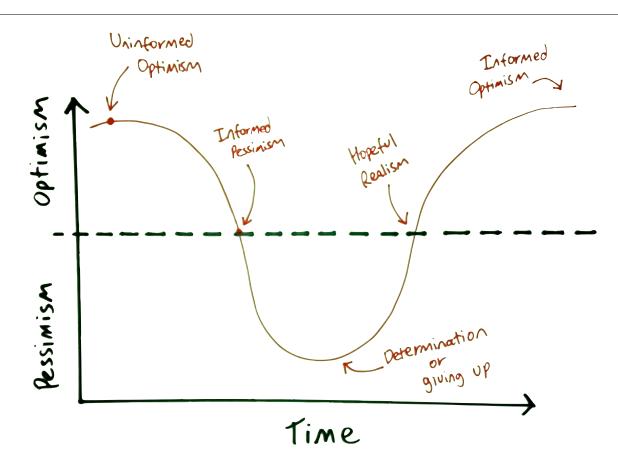


Slack preferred over email for quick questions





### **EMOTIONAL CYCLE OF CHANGE**



### Don't Forget

## Have you signed your enrollment agreement?

### **GA** Catalog:

- ✓ Attendance Policy
- ✓ Cancellation/Withdrawal
- ✓ Tuition Payments
- √Homework + Projects (80%)
- ✓ Letter of Completion

### 0

#### RELIA - WELCOME

#### Dear Student

Congratulations or your acceptance, and/welcome to General Acceptable!
You should have already received at official acceptance-small, but there are after zero maps to take to existly your confinent and reserve-your rest in the class.

as you may have bound, Gu is in the particle of getting hormority Californian regulatory agency for protocondary education. State have requires no to share some new documents with you as apart of the admissions process. We have credened them here for your redunding states, and described them in more detaillations.

Those documents are intension to provide you with important information about corolling in a LL, program, which is an investment that we take very seriously. When presently my a fam of logalous but those forms are required to contain specific language. Where possible, we've tried our best inequire democlars and may to make cond.

#### MECT STEPS

#### REVIEW THE GA CATALIS

Review the 6th Catalog Nazadanial love remined accepted fifth Catalog is your acceptance exact this top acceptance in appealanty in fur palanes and course information. We empourage you in read this carefully before you sign the Euroliseand Agreement.

#### REMIEW AND THE EMPOLLMENT ARRESTST

Review and sign the Euroliment Agreement. Whit forement is a matrix between you and GA. It indicates important information about your reuses, including all applicable loss and sur-posicion for relands, numerication, and slithdrawal.

#### REVIEW AND WORTH THE STANCET BLOCKINGS

Review and sign the School Performance Eact Street. This obscrame this introduction give you cars no exadent entreases, including course completion and job piaconsort suctes, where a opplication. Call was not offering all everses in California in 1811 or 1813, but we are reporting on this time period because the State of California requires us to publish data by, august as every your from the pelos to wy years. Syou'd like more every data as this processor, were formandations to state on.

Phonos short hardrate to much out to ourse Administrate Produce of you have use

### **GA Graduation Requirements**

ASSIGNMENTS
(COMPLETE 80% OF WORK/LABS)

ATTENDANCE
(MISS NO MORE THAN 2 CLASSES)

FINAL PROJECT



### **Our Grads**

# BUILD YOUR NETWORK

It's not just about altruism, your network is your most valuable asset. Alumni have started companies together and recruited other alumni to join their teams

# ACCESS TO OUR ALUMNI COMMUNITY

You're part of the alumni community forever

# OFFICIAL LETTER OF COMPLETION

Sent to your inbox via email after the course is completed. \*important if you're being reimbursed by your business PERKS!
DISCOUNTS OFF CLASSES
AND WORKSHOPS, TUITION
CREDIT

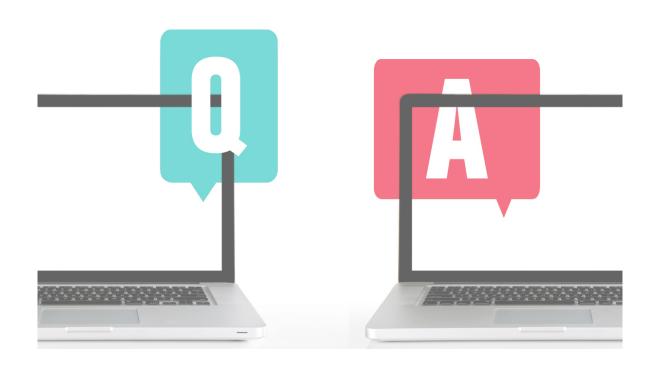
We can't wait to have you back on campus



### **HOW TO SUCCEED**

- Ask
- Share
- D0
- Do
- Do

## **ASK AWAY!**





# WELCOME TO DATA SCIENCE

Naumaan Nayyar

Chief Technology Officer, Vivace Systems

### **INTRODUCTION**

## INTRODUCTIONS

### **WELCOME TO DATA SCIENCE**

## **INTRODUCTIONS**

- ▶ Who are you?
- ▶ Tell us about yourself
- ▶ Share your professional experiences
- ▶ What do you expect to get out of this course?

### **WELCOME TO DATA SCIENCE**

## **LEARNING OBJECTIVES**

- Describe the roles and components of a successful learning environment
- Define data science and the data science workflow
- Setup your development environment and review python basics

### **DATA SCIENCE**

# PRE-WORK

### **PRE-WORK REVIEW**

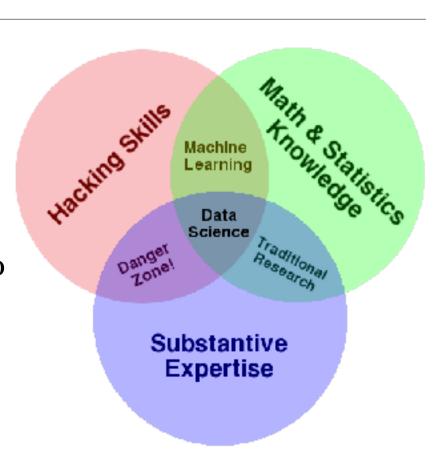
- ▶ Define basic data types used in object-oriented programming
- ▶ Recall the Python syntax for lists, dictionaries, and functions
- ▶ Create files and navigate directories using the command line interface

### INTRODUCTION

## WHAT IS DATA SCIENCE?

### WHAT IS DATA SCIENCE?

- Knowledge vs data
- A set of tools and techniques for data
- ▶ Interdisciplinary problem-solving
- ▶ Application of scientific techniques to practical problems



### WHAT IS DATA SCIENCE?

- Data Science
- **▶**Software Engineering
- ▶ Machine Learning (or Statistics...?)

"...better statisticians than your average programmer and they're better programmers than your average statistician"

### **WHO USES DATA SCIENCE?**

# NETFLIX









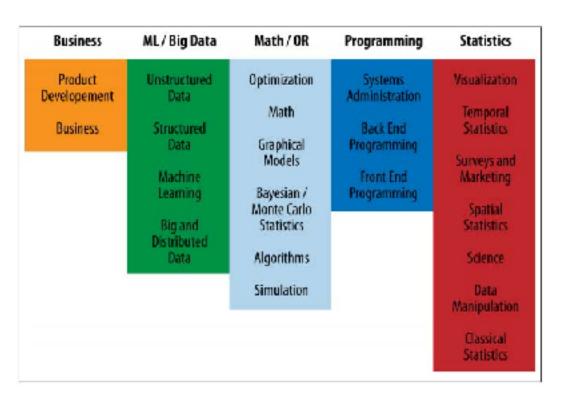


### WHO USES DATA SCIENCE?

How does your company or industry use data science?

### WHAT ARE THE ROLES IN DATA SCIENCE?

Data Science involves a variety of skill sets, not just one.



## DATA SCIENCE BASELINE

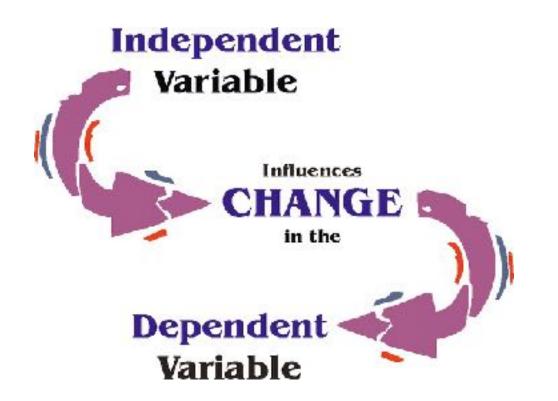
### **ACTIVITY: DATA SCIENCE BASELINE QUIZ**



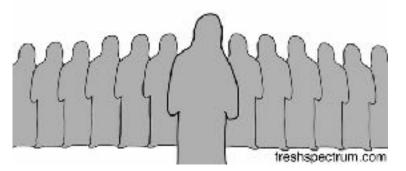
### DIRECTIONS (10 minutes)

- 1. Form groups of three.
- 2. Answer the following questions.
  - a. True or False: Gender (coded male=0, female=1) is a continuous variable.
  - b. Define the following:
    - i. Independent variable (Predictor)
    - ii. Dependent variable (Outcome)
    - iii. Covariate
  - c. Draw a normal distribution
  - d. True or False: Linear regression is an unsupervised learning algorithm.
  - e. What is a hypothesis test?

### **ACTIVITY: DATA SCIENCE BASELINE QUIZ**



# The default, the status quo I am already accepted, can only be rejected The burden of proof is on the alternative I am the null hypothesis

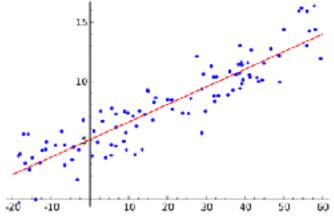


# THINGS YOU WILL BE ABLE TO DO

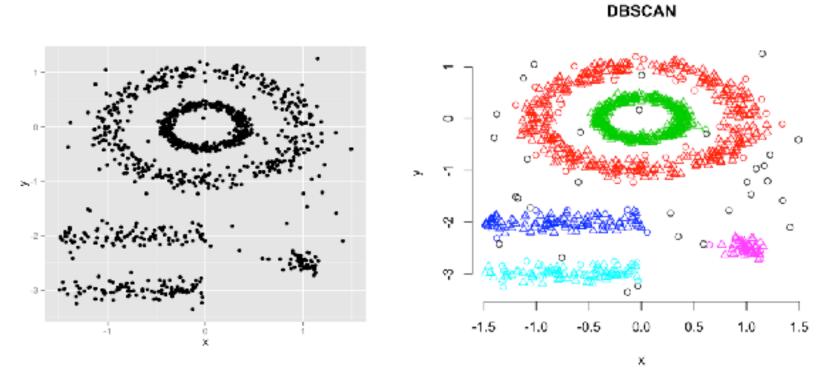
#### **SIMPLE LINEAR REGRESSION**

Def: Explanation of a continuous variable given a series of independent variables

- The simplest version is just a line of best fit: y = mx + b
- Explain the relationship between **x** and **y** using the starting point **b** and the power in explanation **m**.



#### **CLUSTERING: Density-Based**

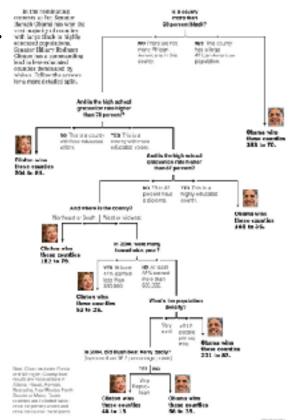


Source: http://www.sthda.com/english/wiki/dbscan-density-based-clustering-for-discovering-clusters-in-large-datasets-with-noise-unsupervised-machine-learning

#### INTUITION BEHIND DECISION TREES

- Decision trees are like the game "20 questions". They make decision by answering a series of questions, most often binary questions (yes or no).
- We want the smallest set of questions to get to the right answer.
- ▶ Each questions should reduce the search space as much as possible.

#### Decision Tree: The Obama-Clinton Divide



# THE DATA SCIENCE WORKFLOW

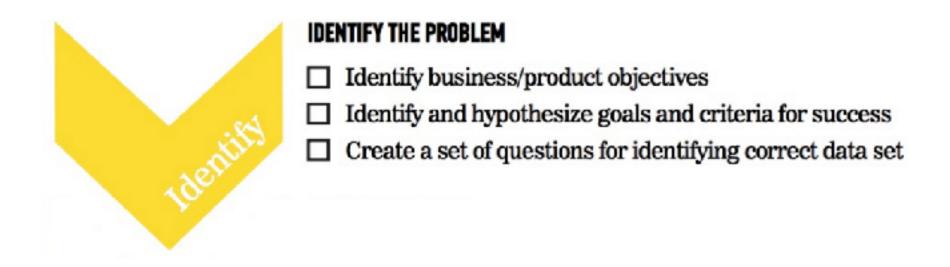
- ▶ A methodology for doing Data Science
- Similar to the scientific method
- ▶ Helps produce ? and ? results

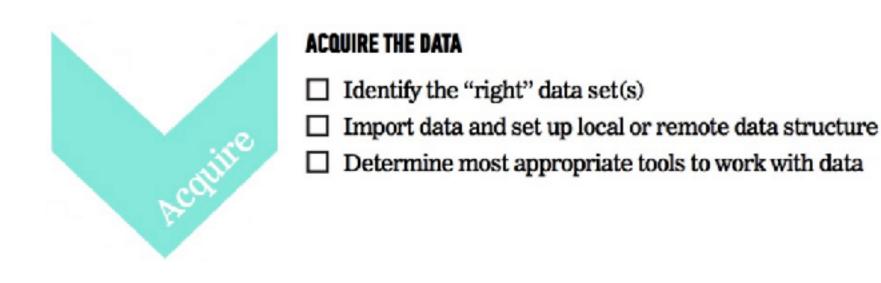
- ▶ A methodology for doing Data Science
- Similar to the scientific method
- ▶ Helps produce *reliable* and *reproducible* results
  - Reliable: Accurate findings
  - *Reproducible*: Others can follow your steps and get the same results

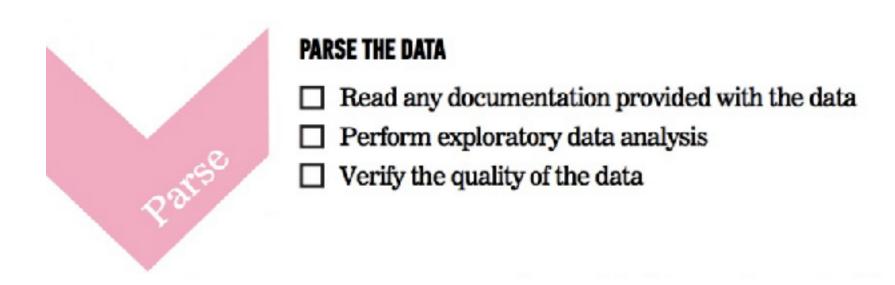
#### The steps:

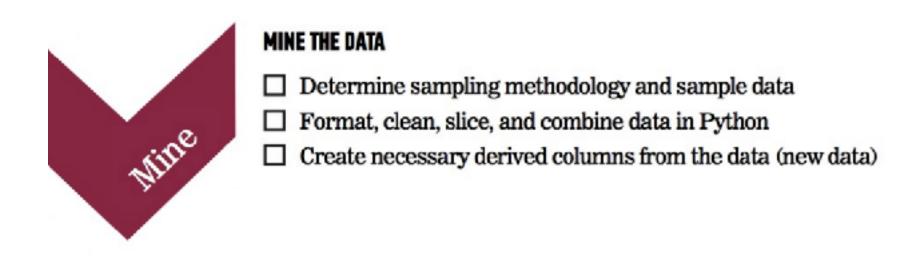
- 1. Identify the problem
- 2. Acquire the data
- 3. Parse the data
- 4. Mine the data
- 5. Refine the data
- 6. Build a data model
- 7. Present the results



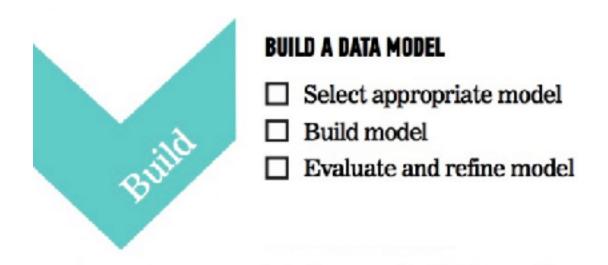














#### PRESENT THE RESULTS

- Summarize findings with narrative, storytelling techniques
- ☐ Present limitations and assumptions of your analysis
- ☐ Identify follow up problems and questions for future analysis

#### **FUTURAMA EXAMPLE**

Problem Statement: "Using Planet Express customer data from January 3001-3005, determine how likely previous customers are to request a repeat delivery using demographic information (profession, company size, location) and previous delivery data (days since last delivery, number of total deliveries)."



• We can use the Data Science workflow to work through this problem.

#### **FUTURAMA EXAMPLE: IDENTIFY THE PROBLEM**

- ▶ Identify the business/product objectives.
- Identify and hypothesize goals and criteria for success.
- ▶ Create a set of questions to help you identify the correct data set.

#### **FUTURAMA EXAMPLE: ACQUIRE THE DATA**

- Ideal data vs. data that is available
- Learn about limitations of the data.
- ▶ What data is available for this example?
- ▶ What kind of questions might we want to ask about the data?

#### **FUTURAMA EXAMPLE: ACQUIRE THE DATA**

- Questions to ask about the data
  - ▶ Is there enough data?
  - ▶Does it appropriately align with the question/problem statement?
  - •Can the dataset be trusted? How was it collected?
  - Is this dataset aggregated? Can we use the aggregation or do we need to get it preaggregated?

#### **FUTURAMA EXAMPLE: PARSE THE DATA**

▶ Secondary data = we didn't directly collect it ourselves

#### ▶ Example data dictionary

Variable	Description	Type of Variable
Profession	Title of the account owner	Categorical
Company Size	1- small, 2- medium, 3- large	Categorical
Location	Planet of the company	Categorical
Days Since Last Delivery	Integer	Continuous
Number of Deliveries	Integer	Continuous

#### **FUTURAMA EXAMPLE: PARSE THE DATA**

- Questions to ask while parsing
  - ▶ Is there documentation for the data? Is there a data dictionary?
  - ▶ What kind of filtering, sorting, or simple visualizations can help understand the data?
  - ▶ What information is contained in the data?
  - ▶ What data types are the variables?
  - Are there outliers? Are there trends?

#### **FUTURAMA EXAMPLE: MINE THE DATA**

- ▶ Think about sampling
- Get to know the data
- **▶**Explore outliers
- ▶ Address missing values
- Derive new variables (i.e. columns)

#### **FUTURAMA EXAMPLE: MINE THE DATA**

- ▶ Common steps while mining the data
  - Sample the data with appropriate methodology
  - ▶ Explore outliers and null values
  - Format and clean the data
  - ▶Determine how to address missing values
  - Format and combine data; aggregate and derive new columns

#### **FUTURAMA EXAMPLE: REFINE THE DATA**

▶ Use statistics and visualization to identify trends

► Example of basic statistics

Variable	Mean (STD) or Frequency (%)
Number of Deliveries	50.0 (10)
Earth	50 (10%)
Amphibios 9	100 (20%)
Bogad	100 (20%)
Colgate 8	100 (20%)
Other	150 (30%)

#### **FUTURAMA EXAMPLE: REFINE THE DATA**

- Descriptive stats help refine by
  - ▶Identifying trends and outliers
  - Deciding how to deal with outliers
  - Applying descriptive and inferential statistics
  - ▶ Determining visualization techniques for different data types
  - **▶**Transforming data

#### **FUTURAMA EXAMPLE: CREATE A DATA MODEL**

- Select a model based upon the outcome
- ▶ Example model statement: "We completed a logistic regression using Statsmodels v. XX. We calculated the probability of a customer placing another order with Planet Express."
- ▶ Steps for model building

#### **FUTURAMA EXAMPLE: CREATE A DATA MODEL**

- ▶ The steps for model building are
  - Select the appropriate model
  - ▶Build the model
  - Evaluate and refine the model
  - Predict outcomes and action items

#### **FUTURAMA EXAMPLE: PRESENT THE RESULTS**

- ▶ You have to effectively communicate your results for them to matter!
- Ranges from a simple email to a complex web graphic.
- ▶ Make sure to consider your audience.
- A presentation for fellow data scientists will be drastically different from a presentation for an executive.

#### **FUTURAMA EXAMPLE: PRESENT THE RESULTS**

- ▶ Key factors of a good presentation include
  - Summarize findings with narrative and storytelling techniques
  - ▶ Refine your visualizations for broader comprehension
  - Present both limitations and assumptions
  - Determine the integrity of your analyses
  - ▶ Consider the degree of disclosure for various stakeholders
  - ▶Test and evaluate the effectiveness of your presentation beforehand

#### **FUTURAMA EXAMPLE: PRESENT THE RESULTS**

- ▶ Example presentations and infographics
  - ▶512 Paths to the White House
  - ▶ Who Old Are You?
  - ▶2015 NFL Predictions

# DATA SCIENCE WORK FLOW

#### **ACTIVITY: DATA SCIENCE WORKFLOW (TIME PERMITTING)**



#### DIRECTIONS (25 minutes)

- 1. Divide into 4 groups, each located at a whiteboard.
- 2. **IDENTIFY**: Each group should develop 1 research question they would like to know about their classmates. Create a hypothesis to your question. Don't share your question yet! (5 minutes)
- **3. ACQUIRE**: Rotate from group to group to collect data for your hypothesis. Have other students write or tally their answers on the whiteboard. (10 minutes)
- **4. PRESENT**: Communicate the results of your analysis to the class. (10 minutes)
  - a. Create a narrative to summarize your findings.
  - b. Provide a basic visualization for easy comprehension.
  - c. Choose one student to present for the group.

#### **DELIVERABLE**

Presentation of the results

# ENVIRONMENT SETUP

#### **DEV ENVIRONMENT SETUP**

- Brief intro of tools
- **▶**Environment setup
  - ▶Create a Github account
  - Install Python 2.7 and Anaconda
  - ▶ Practice Python syntax, Terminal commands, and Pandas
- ▶iPython Notebook test and Python review

#### **DEV ENVIRONMENT SETUP**

- Test your new setup using the lesson 1 starter code available at / lessons/lesson-01/code/starter-code/ in the Github repo
- Ask your classmates and instructor for help if you have problems!

#### **CONCLUSION**

# REVIEW

#### **CONCLUSION**

- You should now be able to answer the following questions:
  - ▶ What is Data Science?
  - ▶ What is the Data Science workflow?
  - ▶ How can you have a successful learning experience at GA?

#### **DATA SCIENCE**

## BEFORE NEXT CLASS

#### **BEFORE NEXT CLASS**

- ▶ Complete <u>Learn Python the Hard Way</u> through Exercise 35, or as far as you have time
  - ▶ Remember the Python syntax for lists, dictionaries, and functions
- ▶ Create files and navigate through the CLI. GA Tutorial
- ▶ Go through <u>Learn Pandas</u> (up to Lesson 3)

#### **RESOURCES**

- ▶ Python quick reference
- ▶ Probability and Statistics <u>refresher</u>

#### **WELCOME TO DATA SCIENCE**

0 & A

#### **WELCOME TO DATA SCIENCE**

### **EXIT TICKET**

DON'T FORGET TO FILL OUT YOUR EXIT TICKET