# CSC380: Principles of Data Science

## Introduction and Course Overview

**Prof. Jason Pacheco**

**TA: Enfa Rose George**          **TA: Saiful Islam Salim**

# Outline

- COVID-19 Precautions

- Data Science Introduction

- Course Overview

- **Mask up in class**

- The vaccines are very safe and very effective

- Notify me if you fall ill and think it will impact coursework

**If we are forced to go remote**
  - I will schedule Zoom lectures
  - They will be accessible via D2L
  - I will notify everyone by email

Stop the Spread.

ACCESS

# Data Science Job Market

*A search of "data scientist" jobs in the US (on 8/20/2021) shows…*

## Many job options available
- Indeed: 42,000+ jobs
- Glassdoor: 24,000+ jobs
- LinkedIn: 63,000+ jobs

2021's #2 best job in America, according to Glassdoor.com (after Java Developer)

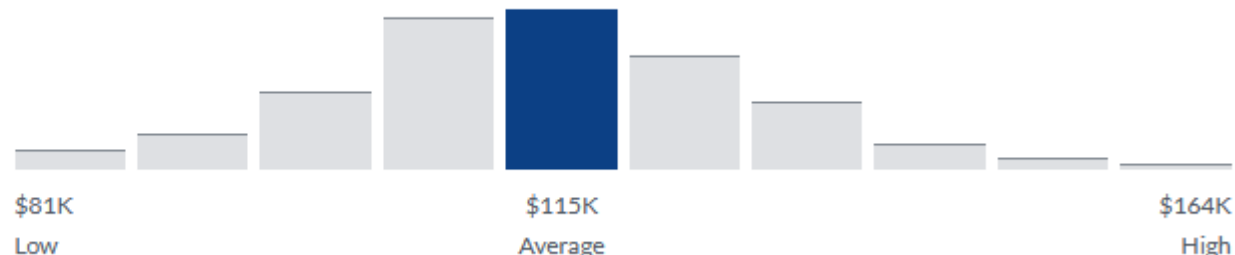## Lucrative pay (Glassdoor)

**Very High** Confidence
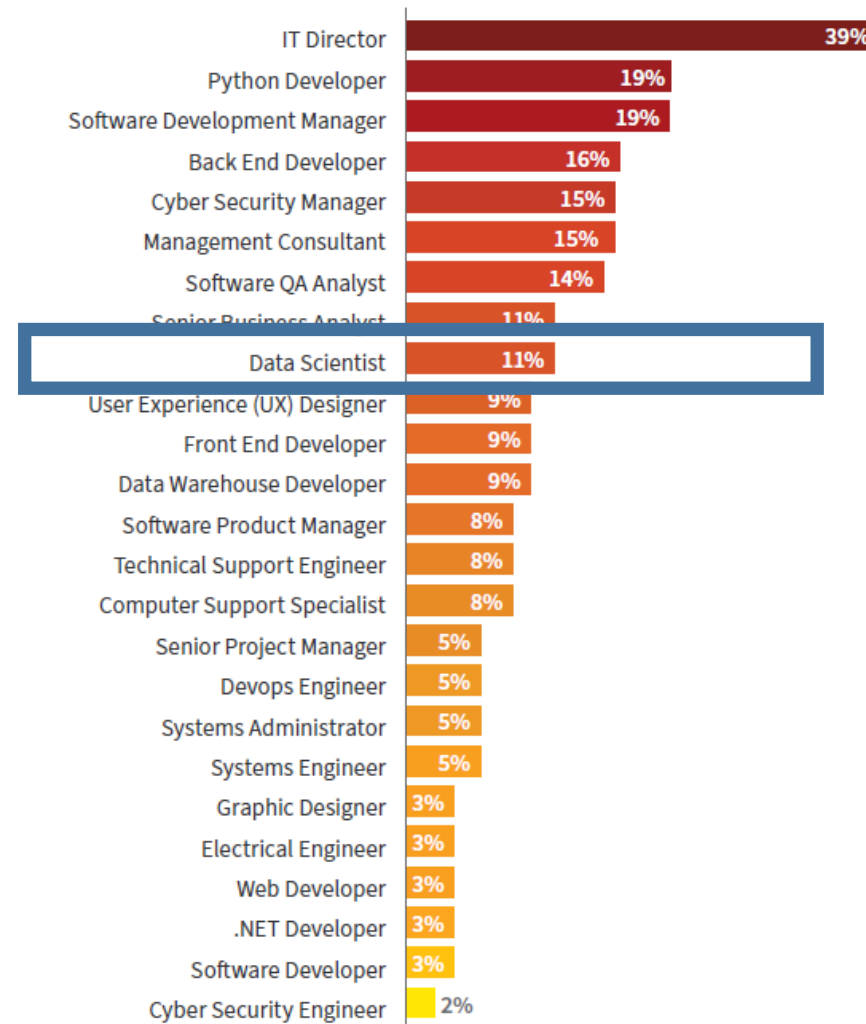
**$115,394** /yr

**Average Base Pay**

17,903 salaries

$81K
Low

$115K
Average

$164K
High

**Seniority Levels**

L2   Data Scientist
         $115,394 /yr

L3   Senior Data Scientist
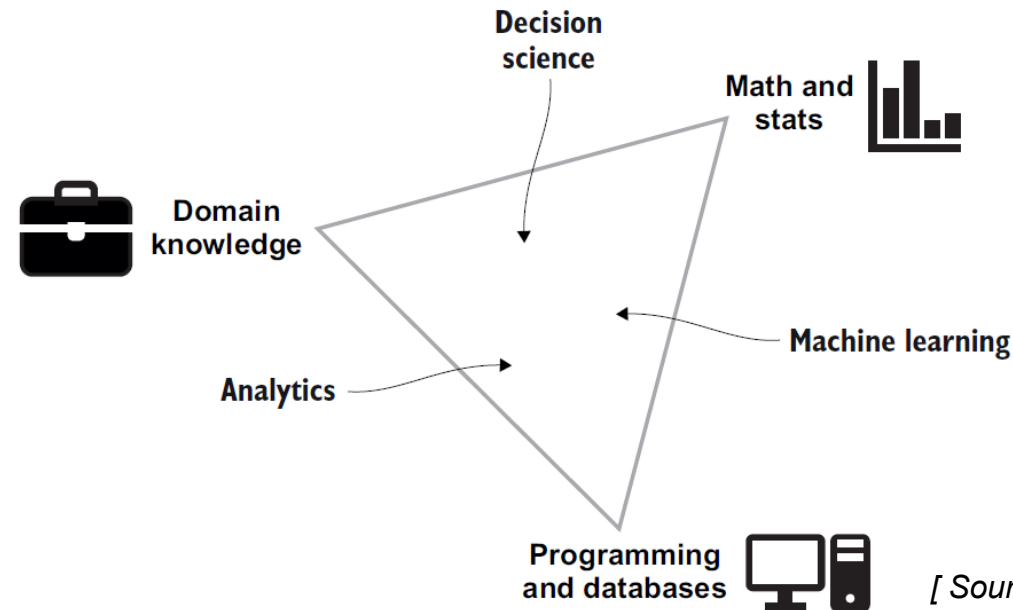         $139,919 /yr

L4   Data Scientist IV
         $135,775 /yr

# Data Science Job Market

*Among the top 10 fastest growing jobs in 2020*



| Job | Growth |
|-----|--------|
| IT Director | 39% |
| Python Developer | 19% |
| Software Development Manager | 19% |
| Back End Developer | 16% |
| Cyber Security Manager | 15% |
| Management Consultant | 15% |
| Software QA Analyst | 14% |
| Senior Business Analyst | 11% |
| Data Scientist | 11% |
| User Experience (UX) Designer | 9% |
| Front End Developer | 9% |
| Data Warehouse Developer | 9% |
| Software Product Manager | 8% |
| Technical Support Engineer | 8% |
| Computer Support Specialist | 8% |
| Senior Project Manager | 5% |
| Devops Engineer | 5% |
| Systems Administrator | 5% |
| Systems Engineer | 5% |
| Graphic Designer | 3% |
| Electrical Engineer | 3% |
| Web Developer | 3% |
| .NET Developer | 3% |
| Software Developer | 3% |
| Cyber Security Engineer | 2% |

**Source: Top Jobs in Dice Tech Q3 Report**

# What is "Data Science"?

**My Definition:** *The process of using data to answer questions, extract knowledge, and predict future outcomes.*



Caveat: I don't *love* this figure since it leaves out **visualization.**

[ *Source: Robinson, E. and Nolis, J.* ]

## Data Science Is:

- **Interdisciplinary**: Combines tools and techniques from Math / Statistics / CS
- **Exploratory**: Understanding data requires creative exploration and visualization
- **Applied Statistics & Probability** + extra stuff to handle, process, and visualize data

Data Science Applications

# Moneyball

**Problem** *How to assemble the best baseball team with a small budget?*

- 2002 Major League Baseball (MLB) draft
- Traditional team building relies on *scouts*
- Assumption: The collective wisdom of insiders is biased / flawed
- **SABRmetrics:** Data-driven and evidence-based approach to player quality evaluation
- *On-base %* and *Slugging %* are good indicators of offensive success
- Players with these "features" are cheaper compared to traditional statistics (stolen bases, runs batted in, batting average)

# *Moneyball:* Impact

- In 2002 Oakland A's ($44M budget) were competitive to the New York Yankees ($125M budget)

- Toronto Blue Jays hired full-time sabermetric analysts

- 2020 season "masters of Moneyball" Tampa Bay Rays reached world series with the 3$^{rd}$ lowest salary of all MLB

- In 2019 Liverpool Football/Soccer adopted this approach to nearly win the title (they lost to Manchester)

- Brad Pitt got a paycheck out of it for the movie (7.6/10 IMDB)…

# Data Science Workflow



[ Adapted from: Grolemund and Wickham, 2018 ]

This is a class about <u>data science</u> it is **not** a class about politics.  We will discuss election forecasting **only** in the context of <u>data science</u> and we will **ignore politics**.

**Problem** *Who will win the 2020 US presidential election?*

**Details**

- There are 2 primary candidates Donald Trump & Joe Biden*
- The *incumbent* (Trump) is the sitting president
- There are 50 states, each has a number of *electors*
- Each elector has a vote in the *electoral college*
- Electors for each state vote for the majority vote in that state
- Maine and Nebraska use a district method
- The winner has the majority of 538 electors (typically 270 or more votes)

*\* Secondary candidates do not have a realistic chance of winning, but cannot be ignored since they affect votes for primary candidates*

# Election Forecasting: The Model

*[FiveThirtyEight](#) uses a proprietary statistical model based on…*

## Poll aggregation model

Weight accounts for poll sample size, timeliness, historical accuracy

$$\text{prediction} = \sum_i \text{weight}_i \times \text{poll}_i + \text{random noise}$$

## Additional model inputs

- States grouped by demographic subcategories
- Per capita income
- Age distribution of residents
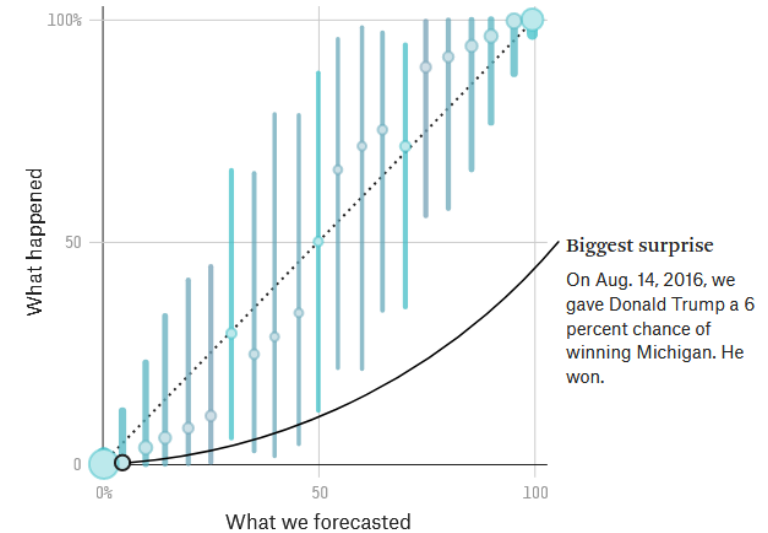- All features are *significant* to 85% level

## Important properties of the model

- Predictive statements are *probabilistic*
- Assigns higher probability to extreme outliers
- Accounts for correlation among states / polls

**Calibration plot**

Calibration plots show us whether events happened as often as we predicted they would.

Key   1,000 ○ ○ 10,000 observations
95% confidence



What happened

100%

50

0

0%        50        100

What we forecasted

Biggest surprise
On Aug. 14, 2016, we gave Donald Trump a 6 percent chance of winning Michigan. He won.

**FiveThirtyEight**

*Generative (Bayesian\*) model allows simulation of random realizations…*

*…visualizations targeted at communicating <u>uncertainty</u> about prediction.*

*Model also allows "what if" (e.g. counterfactual) analysis…*



*…this is a feature of model interpretability.*

# Bad Data Science & Statistics

# Types of Data

*Data come in many forms, each requiring different approaches & models*

**Natural Language**





**Timeseries**

**Image / Video**

*The number of types is endless, these are just some examples*

# Programming Languages for Data Science

*Python and R are both standard for data science these days*



We will use Python for this course since you should already know it

## Python Packages Covered

## Other Useful Python Packages

**CYVERSE** ®

Transforming Science Through Data-Driven Discovery

- Created and primarily house here at **University of Arizona!**
- Now a full-fledged infrastructure for shared data science across many fields
- Cloud storage and computing computing resources
- Makes sharing data, code, etc. easy across large distributed teams
- Virtual machine environments for high performance computing
- Supports most standard languages / tools / etc.
- Can package code into containers and easily share with others

# Course Overview: Resources

## Probabilistic Graphical Models

CSC 380 | Fall 2021 | TuTh 5:00-6:15pm | Modern Languages Rm 311



[ Image Source: Robinson, E. and Nolis, J. ]

### Description of Course

This course introduces students to the principles and tools of data science. This course will provide a foundation for properly collecting and analyzing data to draw insights and to answer data-driven questions. The course has three main components: applied probability and statistics, data analysis and visualization, and machine learning. In the first component students will be introduced to the fundamentals of applied probability and statistics, le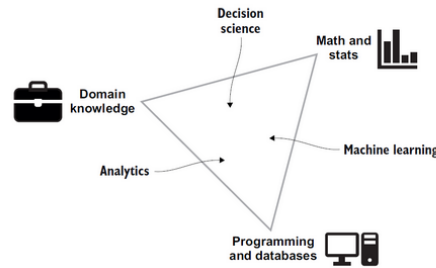arn how to interpret randomness, and how to assess predictive uncertainty. Students will then learn how to handle, clean, process, and visualize data of varying types using Python. Finally, the students will be introduced to the basics of machine learning to build predictive models. Students will further learn how to assess model validity and how to interpret the quality of model predictions.

### Primary Resources

All reading material will be made available through presentation slides or the course webpage. Students will find the following optional textbooks useful throughout this course:

**WL** : Wasserman, L. "All of Statistics: A Concise Course in Statistical Inference." Springer, 2004

**MK** : Murphy, K. "Machine Learning: A Probabilistic Perspective." MIT press, 2012

### Instructor and Contact Information:

**Instructor:** Jason Pacheco, GS 724, Email: pachecoj@cs.arizona.edu
**Office Hours:** TBD
**D2L:** https://d2l.arizona.edu/d2l/home/1072117
**Piazza:** https://piazza.com/arizona/fall2021/csc380
**Instructor Homepage:** http://www.pachecoj.com

---

Resources accessible on course website
pachecoj.com/courses/csc380_fall21/

## Specific resources

- D2L for assignment submission

- Piazza for **all** communication

- Readings and electronic textbooks

- Lecture slides (posted after class)

## Every lecture accompanied by reading

- These will not be graded but are required

- Homeworks will incorporate material from readings

- Reading for today's lecture:

   *Robinson and Nolis, "What's Data Science?"*

# Textbooks



Murphy, K. "Machine Learning: A Probabilistic Perspective." MIT press, 2012

( UA Library )



Wasserman, L. "All of Statistics." Springer, 2004

( Springer )

*Additional readings on the course webpage*

# Course TAs

*Your friendly course TAs…*

*Saiful Islam Salim*

*saifulislam@email.arizona.edu*

*Enfa Rose George*

*enfageorge@email.arizona.edu*

# Expected Skills

- This class will require a fair amount of math
  - Probability and Statistics (first few lectures)
  - Single-variable Calculus
  - Linear Algebra (two-lecture overview)

- This class will require a fair amount of **coding**
  - Reading in / cleaning / visualizing data
  - Simulating random processes
  - Training and evaluating machine learning models

- Early assignments will be mostly math, later will be coding

# Course Overview

**Course Objective** *Introduction to basic concepts in data science and machine learning.*

| Probability and Statistics | Data Handling and Visualization | Machine Learning | Ethics and Fairness |
|---|---|---|---|
| Random events / variables, distributions / densities, moments, descriptive stats, estimation | Reading & cleaning, transformation & preprocessing, visualization | Predictive models, supervised learning, unsupervised learning, model checking | Data privacy, ethics, fairness |

# Probability and Statistics

***Suppose we roll <u>two fair dice</u>…***

- ➢ What are the possible outcomes?
- ➢ What is the *probability* of rolling **even** numbers?

*… this is an **experiment** or **random process**.*

**We will learn how to…**

- ➢ Mathematically formulate outcomes and their probabilities?
- ➢ Describe characteristics of random processes
- ➢ Estimate unknown quantities (e.g. are the dice actually fair?)
- ➢ Characterize the uncertainty in random outcomes
- ➢ Identify and measure dependence among random quantities

# Data Handling and Visualization

## *In Data Handling we will learn to…*

➢ Collect data through population sampling

➢ Identify and avoid biased population samples

➢ Clean data and correct errors

➢ Transform and preprocess data (*wrangling*)

[ Image Source: Code A Star ]

## *In Data Visualization we will learn…*

➢ Why visualization is important

➢ Exploratory data analysis

➢ Common forms of visualization

➢ Pitfalls and gotchas

# Machine Learning

*How do use data to learn underlying patterns and predict unknowns?*



**Unlabelled Data** → K-means → **Labelled Clusters**

X = Centroid

## *In Machine Learning we will learn…*

➢ Principles of prediction
➢ Proper partitioning of training / validation / test data
➢ Unsupervised vs. supervised learning
➢ Linear and nonlinear models
➢

**We will preface this section with a Linear Algebra primer**

[ Image Source: Towards Data Science ]

**RESPONSIBLE DATA SCIENCE**

FAIRNESS · ACCURACY · CONFIDENTIALITY · TRANSPARENCY

## *In Ethics and Fairness we will learn…*

➢ The principles of data privacy

➢ Identifying, measuring, and ensuring fairness

➢ Measuring accuracy through model validation and checking

➢ Transparent communication of findings and predictions

[ Image Source: Responsible Data Science Consortium ]

*9 Homeworks + Midterm + Final Exam*

## Homeworks

- Generally, you will have 1 week per assignment
- There will be an assignment nearly every week
- Assignment *typically* out and due on Thursdays
- Grades by one week after due date
- Some irregularity around holidays
- No assignment over Thanksgiving break

## Grading Breakdown

- Homework: 60%
- Midterm: 20%
- Final: 20%

**First assignment out one week from today**

*Attendance is **strongly** encouraged, but not explicitly graded*

# Late Policy

*Late submissions impact other students, delay grading, and delay solutions*

**But sometimes we need a little extra time…**
- **No more than 1** assignment **no more than 1** day late without penalty
- All subsequent late assignments will receive a zero score
- D2L will accept late assignments but they will be flagged

**If you are struggling with time…**
- Notify me (Piazza) at least 24hrs before the deadline
- Submit the best version of what you have by the deadline
- In general I **will not** grant extra time, and will grade what has been submitted

*If you submit **all** assignments on time, it may benefit your final grade*

# Academic Integrity

*Assignments are to be done independently…*

**If I or the TA suspects you of having cheated**
- You will be notified immediately
- We will have a conference where you can plead your case
- If I am not swayed then you receive a zero for the assignment
- There is an appeals process if you are confident in your case

**Bottom line don't cheat**

# "Office" Hours

- Office hours will be held via Zoom, accessible via D2L

- I will hold two 1.5hr sessions each week

- The final office hour schedule will be announced next week

- If you have a conflict with the schedule, let me know (Piazza)

# Mental Wellbeing

*Some occasional stress / depression / anxiety is normal, but sometimes you may need extra help*

- Non-emergency UA resources at Counseling & Psych Services Mon-Fri
  - Phone: 520-621-3334
  - Web: https://health.arizona.edu/counseling-psych-services

- Emergency resources in Tucson in this Google Doc

# Inclusivity

*I want to foster a comfortable and inclusive classroom experience*

Please let me know if you feel excluded in any way, e.g.
- "Alice-and-Bob" style examples of material
- Improper use of pronouns
- Microagressions
- Miscellaneous statements / interactions

**You can message me on Piazza or discuss in person**