# Welcome to CSC 276
# Data Science

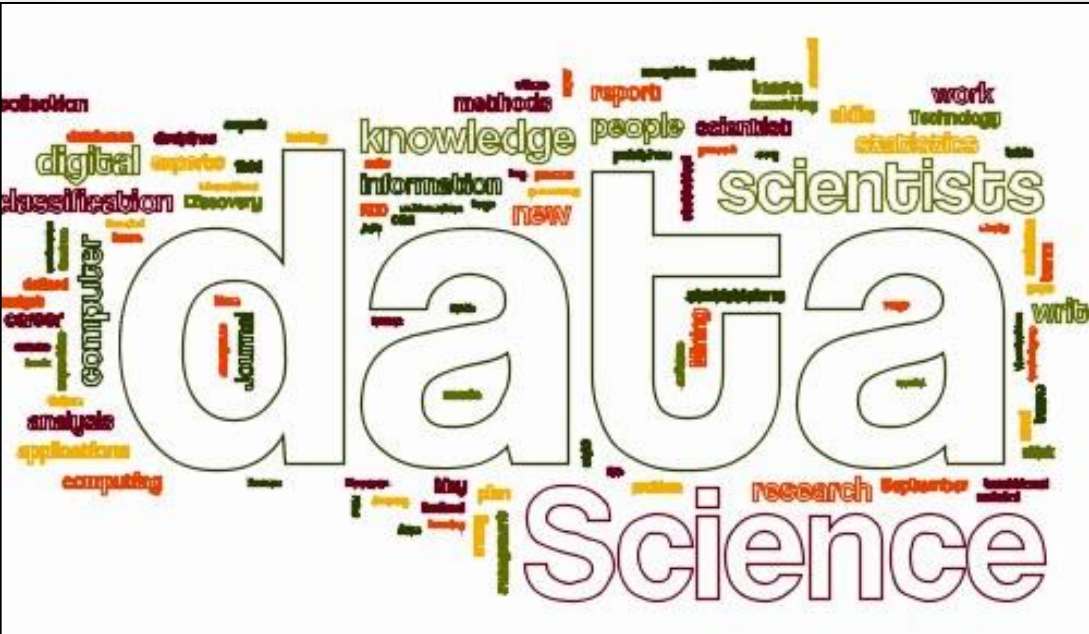# CSC 276: Data Science
# Lecture #2
# Introduction

Dr.Fatema Nafa

Fall 2022

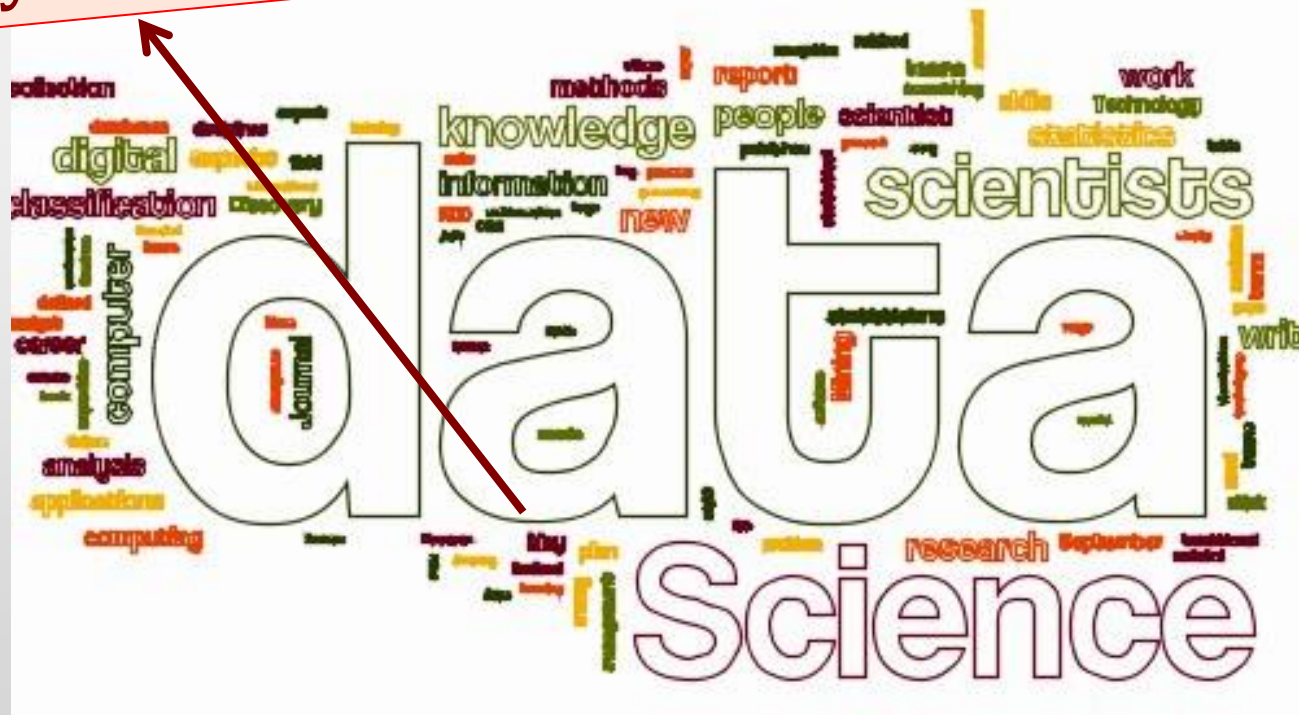# Welcome to CSC 276!



Data Science

# Welcome to CSC 276!

Data Science

This class is truly seminar-style: I'm here, as you are, in order to gain insights into this very new field… .

# Lecture Outline

- **The Art of Data Science**
- Volume, Velocity, Variety
- The Logic of Data Science
- How to Be Agile
- Treating Data as Evidence
- Python
    - Fundamentals of Data Manipulation
    - Basic Data Processing with Pandas
    - Answering Questions with Messy Data

# Data Science – A Definition

a data scientist is someone who asks unique, interesting questions of data based on formal or informal theory, to generate rigorous and useful insights.

It is likely to be an individual with multi-disciplinary training in computer science, business, economics, statistics, and armed with the **necessary quantity of domain knowledge relevant to the question at hand**.

The potential of the field is enormous for just a few well-trained data scientists armed with big data have the potential to transform organizations and societies.

In the narrower domain of business life, the role of the data scientist is to generate applicable business intelligence.
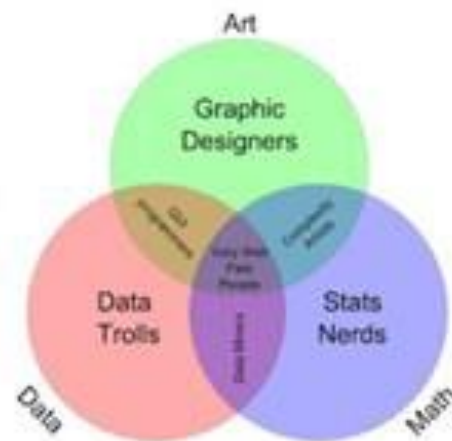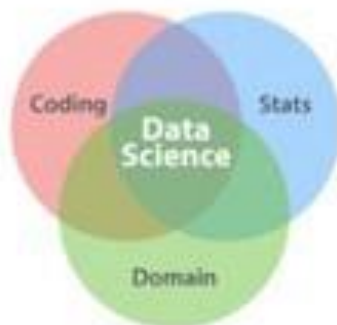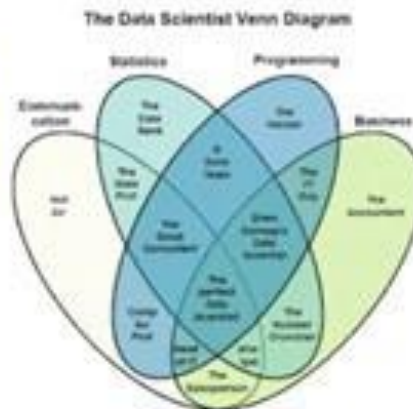
Among all the new buzzwords in business – and there are many –
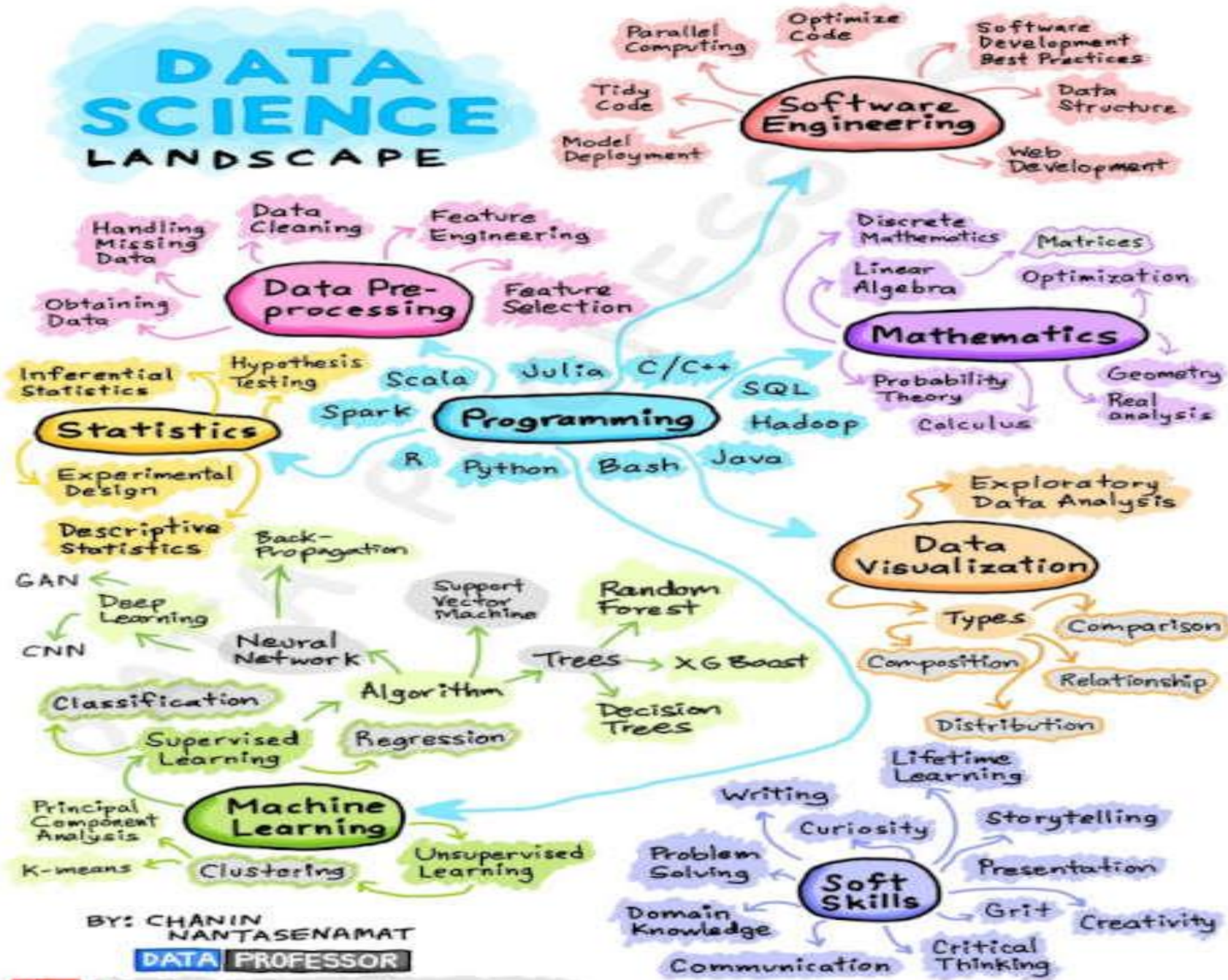
# Data Science – A Definition

Being a data scientist is inherently interdisciplinary.

Good questions come from many disciplines, and the best answers are likely to come from people who are interested in multiple fields, or at least from teams that co-mingle varied skill sets.
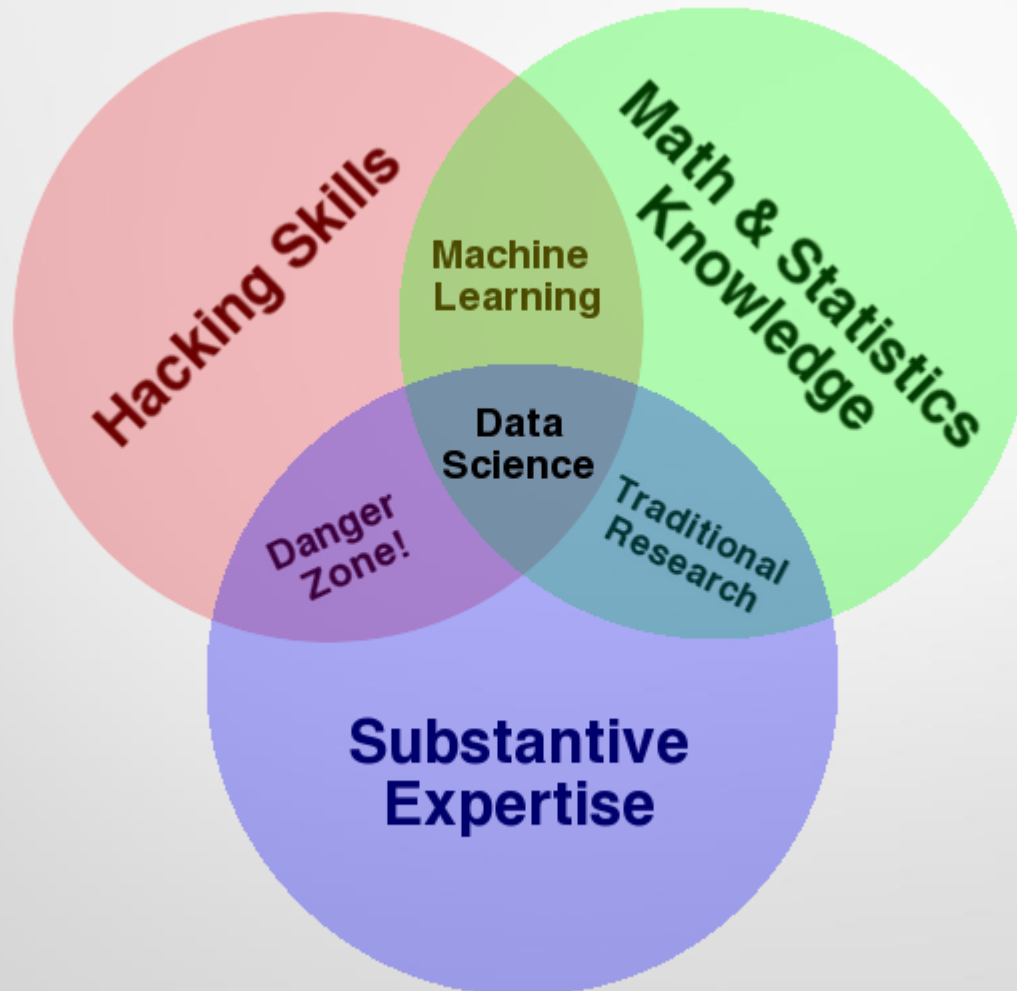
# Data Science – One Definition

# Data Scientists are in high demand

# Also, in academia

# Pays Well



**Big Data, Big Paycheck**

Median salary for analytics professionals and those specifically within data science, by level of experience.

| Experience | | |
|---|---|---|
| Up to 3 years | Analytics professionals | $65,000 |
| | Data scientists | $80,000 |
| 4 to 8 years | | $85,000 |
| | | $120,000 |
| 9+ years | | $115,000 |
| | | $150,000 |

Note: Data do not include managers    Source: Burtch Works    The Wall Street Journal

# Lecture Outline

- **The Art of Data Science**
- **Volume, Velocity, Variety**
- The Logic of Data Science
- How to Be Agile
- Treating Data as Evidence
- Fundamentals of Data Manipulation
- Basic Data Processing with Pandas
- Answering Questions with Messy Data

There are several "V"s of big data: three of these are volume, velocity, variety.[8] Big data exceeds the storage capacity of conventional databases. This is it's *volume* aspect. The scale of data generation is mind-boggling. Google's Eric Schmidt pointed out that until 2003, all of human kind had generated just 5 exabytes of data (an exabyte is $1000^6$ bytes or a billion-billion bytes). Today we generate 5 exabytes of data every two days. The main reason for this is the explosion of "interaction" data, a new phenomenon in contrast to mere "transaction" data. Interaction data comes from recording activities in our day-to-day ever more digital lives, such as browser activity, geo-location data, RFID data, sensors, personal digital recorders such as the fitbit and phones, satellites, etc. We now live in the "internet of things" (or iOT), and it's producing a wild quantity of data, all of which we seem to have an endless need to analyze. In some quarters it is better to speak of 4 Vs of big data, as shown in Figure 1.1.

# Goal of Data Science

Turn data into data products.

# Analysis

**What kinds of data will you use?**

- Almost anything is OK, except other predictions.

- **History**: individual or pair-wise?

- Team or players?

- Numerical or text?

- What kind of model will you build?

- What assumptions are safe to make?

# **Where does data come from?**

# "Big Data" Sources

## It's All Happening On-line

Every:
Click
Ad impression
Billing event
Fast Forward, pause,…
Server request
Transaction
Network message
Fault

…

## User Generated (Web & Mobile)

…

..

## Internet of Things / M2M

## Health/Scientific Computing

**Baseline information**

Cost of genome sequencing compared with Moore's law for computers

Log scale

Cost of computing (Moore's law)

$ per million DNA bases

100,000
10,000
1,000
100
10
1.0
0.1

1999  2002  04  06  08  10

Source: Broad Institute

# Graph Data

Lots of interesting data
has a graph structure:
- Social networks
- Communication networks
- Computer Networks
- Road networks
- Citations
- Collaborations/Relationships
- …

Some of these graphs can get
quite large (e.g., Facebook[*]
user graph)

# What can you do with the data?



Crowdsourcing      +   physical modeling      +   sensing   +   data assimilation

to produce:



From Alex Bayen, UCB

# Contrast: Databases

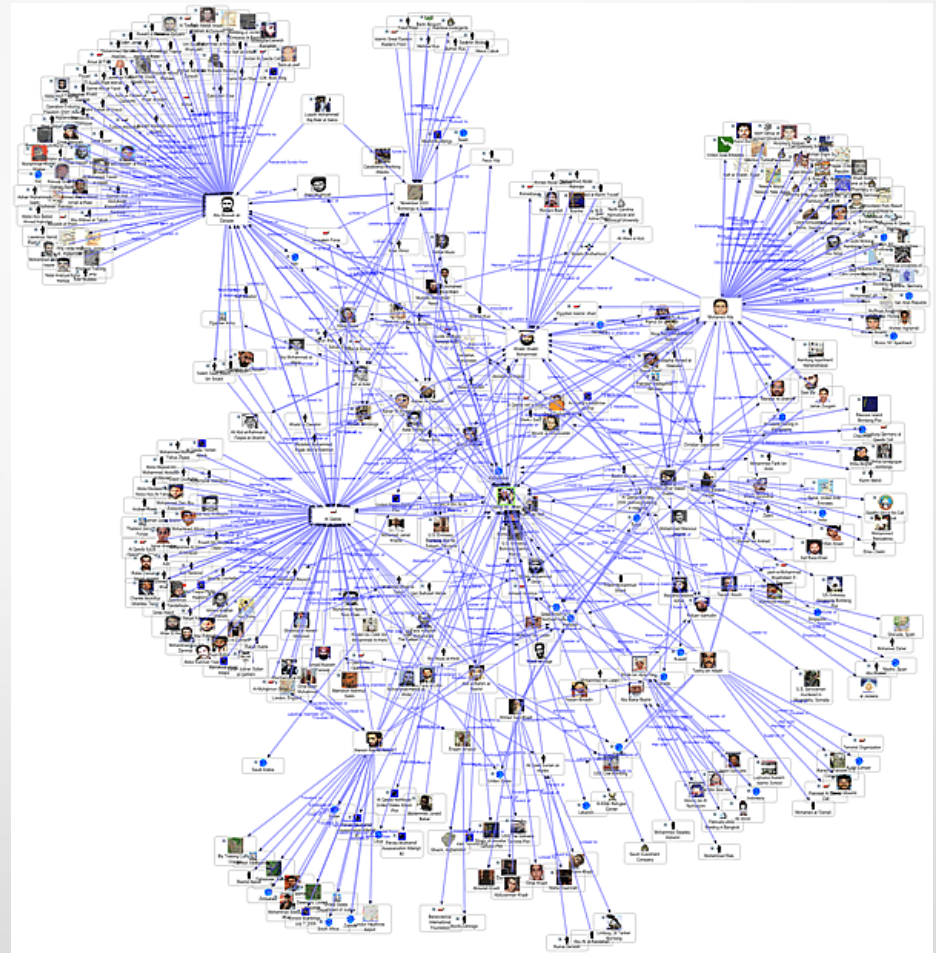|  | **Databases** | **Data Science** |
| --- | --- | --- |
| Data Value | "Precious" | "Cheap" |
| Data Volume | Modest | Massive |
| Examples | Bank records, Personnel records, Census, Medical records | Online clicks, GPS logs, Tweets, Building sensor readings |
| Priorities | Consistency, Error recovery, Auditability | Speed, Availability, Query richness |
| Structured | Strongly (Schema) | Weakly or none (Text) |
| Realizations | SQL | NoSQL: Riak, Memcached, Apache River, MongoDB, CouchDB, Hbase, Cassandra,… |

# Contrast: Databases

| Databases | Data Science |
|---|---|
| Querying the past | Querying the future |



**Business intelligence** (**BI**) is the transformation of raw data into meaningful and useful information for business analysis purposes. BI can handle enormous amounts of unstructured data to help identify, develop and otherwise create new strategic business opportunities - Wikipedia

# Contrast: Scientific Computing


NOAA GFDL CM2.1 Climate Model
Surface Air Temperature Change [°F]
(2050s average minus 1971-2000 average)    SRES A1B scenario

Image    General purpose classifier
Supernova
Not
Nugent group / C3 LBL

| Scientific Modeling | Data-Driven Approach |
|---|---|
| Physics-based models | General inference engine replaces model |
| Problem-Structured | Structure not related to problem |
| Mostly deterministic, precise | Statistical models handle true randomness, and **unmodeled complexity**. |
| Run on Supercomputer or High-end Computing Cluster | Run on cheaper computer Clusters (EC2) |

# Contrast: Computational Science



**CASP: A Worldwide, Biannual Protein Folding Contest**

**Brain Mapping: Allen Institute, White House, Berkeley**

| Quark | Raptor-X |
|---|---|
| Rich, Complex Energy Models | Data-intensive, general ML models |
| Faithful, Physical Simulation | Feature-based inference |
| | Conditional Neural Fields |

| Techniques (Massive ML) |
|---|
| Principal Component Analysis |
| Independent Component Analysis |
| Sparse Coding |
| Spatial (Image) Filtering |

# Contrast: Machine Learning

| Machine Learning | Data Science |
|---|---|
| Develop new (individual) models | Explore many models, build and tune hybrids |
| Prove mathematical properties of models | Understand empirical properties of models |
| Improve/validate on a few, relatively clean, small datasets | Develop/use tools that can handle massive datasets |
| Publish a paper | Take action! |

## How to Be Agile

The nature of data science is experimental. You don't know the answer to the question asked of you—or even if an answer exists. You don't know how long it will take to produce a result or how much data you need. The easiest approach is to just come up with an idea and work on it until you have something. But for those of us with deadlines and expectations, that approach doesn't fly. Companies that issue you regular paychecks usually want insight into your progress.

This is where being agile matters. An agile data scientist works in small iterations, pivots based on results, and learns along the way. Being agile doesn't guarantee that an idea will succeed, but it does decrease the amount of time it takes to spot a dead end. Agile data science lets you deliver results on a regular basis and it keeps stakeholders engaged.

The key to agile data science is delivering data products in defined time boxes—say, two- to three-week sprints. Short delivery cycles force us to be creative and break our research into small chunks that can be tested using minimum viable experiments (Figure 6-1). We deliver something tangible after almost every sprint for our stakeholders to review and give us feedback. Our stakeholders get better visibility into our work, and we learn early on if we are on track.

## Treating Data as Evidence

The logic of data science tells us what it means to treat data as evidence. But following the evidence does not necessarily lead to a smooth increase or decrease in confidence in a model. Models in real-world data science change, and sometimes these changes can be dramatic. New observations can change the models you should consider. New evidence can change confidence in a model. As we collected new employee satisfaction responses, factors like specific job titles became less important, while factors like advancement opportunities became crucial. We stuck with the methods described in this chapter, and as we collected more observations, our models became more stable and more reliable.

I believe that data science is the best technology we have for discovering business insights. At its best, data science is a competition of hypotheses about how a business really works. The logic of data science are the rules of the contest. For the practicing data scientist, simple rules like Ockham's Razor and Bayesian reasoning are all you need to make high-quality, real-world decisions.

# Lecture Outline

- **The Art of Data Science**
- Volume, Velocity, Variety
- The Logic of Data Science
- How to Be Agile
- Treating Data as Evidence
- **Python**
  - Fundamentals of Data Manipulation
  - Basic Data Processing with Pandas
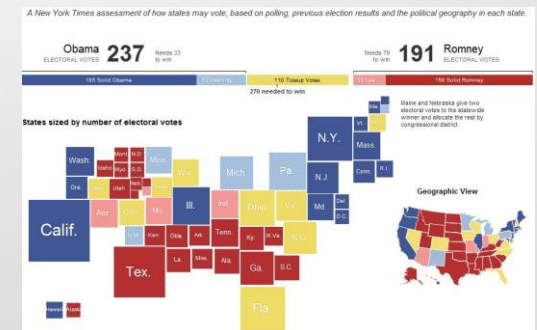  - Answering Questions with Messy Data

1. Data sources

2. Collect data(**download**)

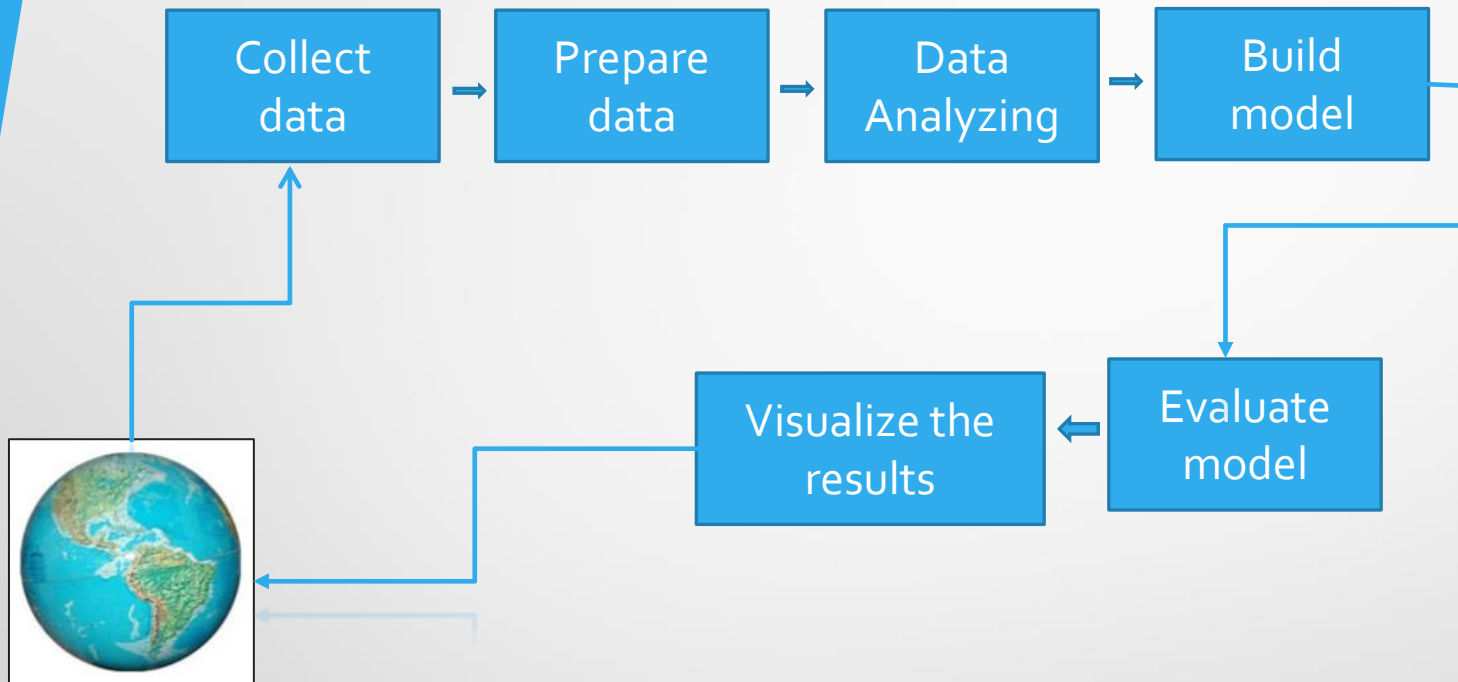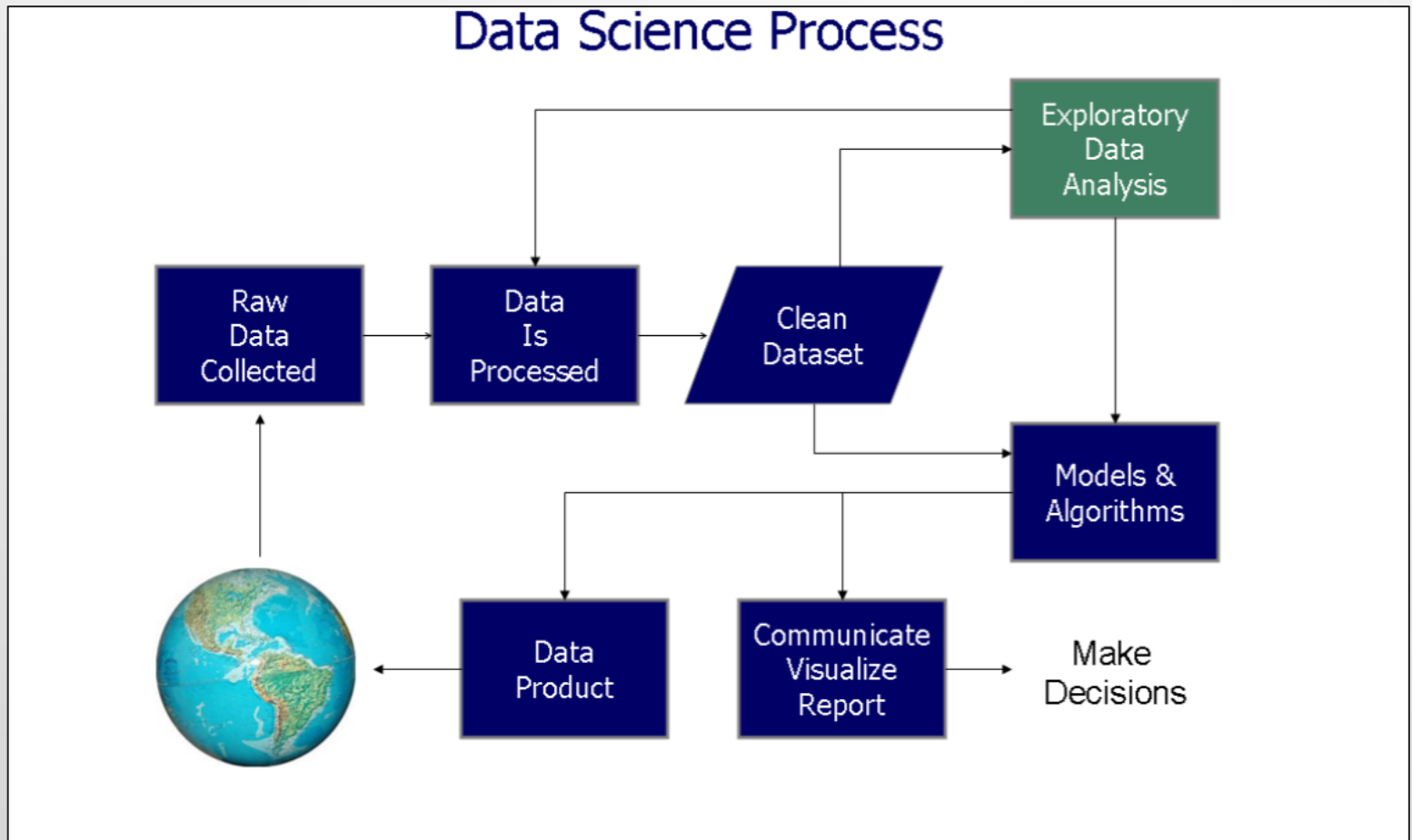3. Prepare data (integrate, transform, clean, filter, aggregate)

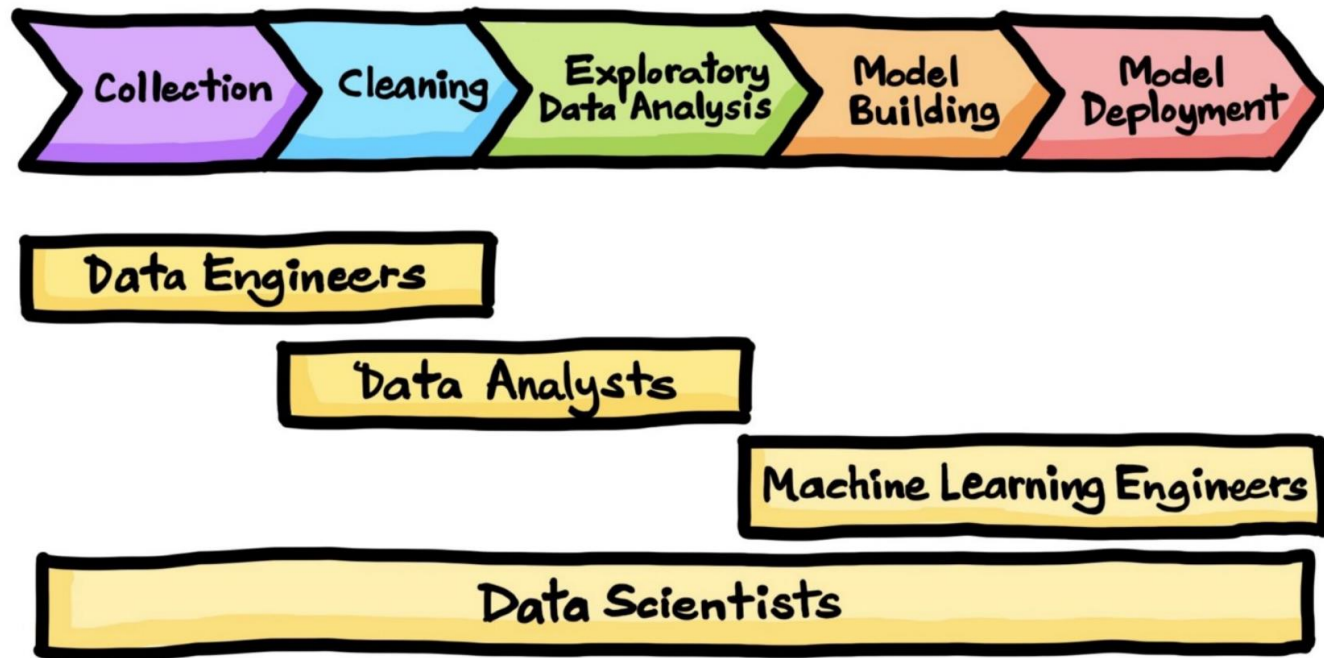4. Build model

5. Evaluate model

6. Visualize the results

Data Science Process

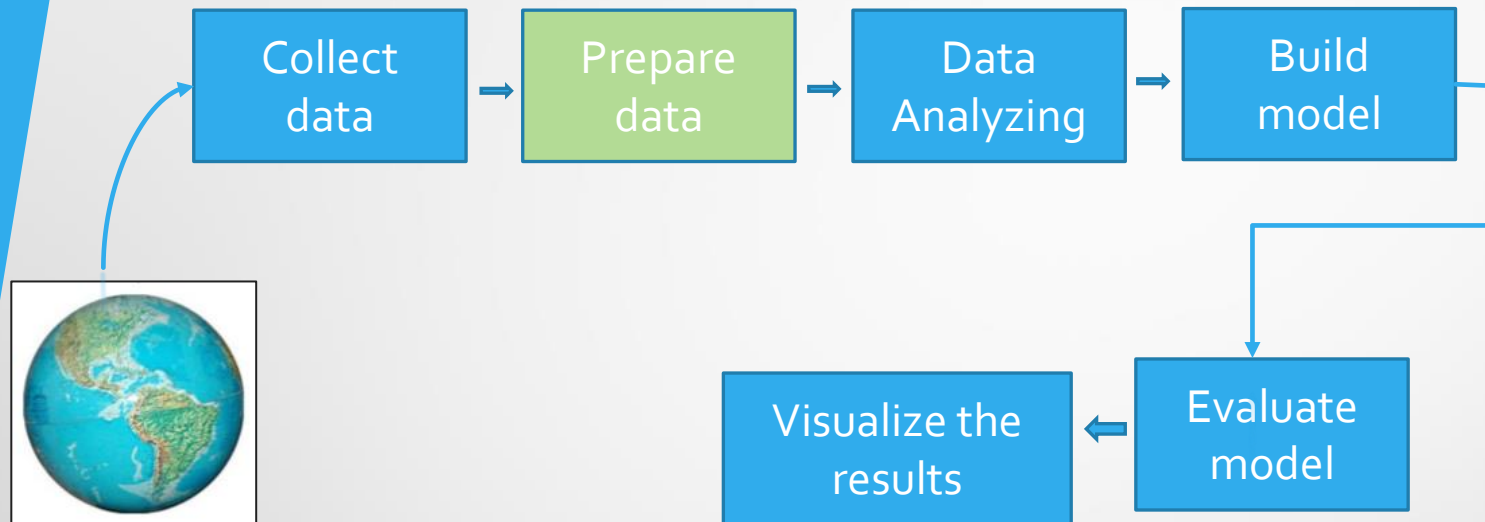# Asking Interesting Questions from Data

Good data scientists develop an inherent curiosity about the world around them, particularly in the associated domains and applications they are working on. They enjoy talking shop with the people whose data they work with. They ask them questions: What is the coolest thing you have learned about this field? Why did you get interested in it? What do you hope to learn by analyzing your data set? Data scientists always ask questions.

Good data scientists have wide-ranging interests. They read the newspaper every day to get a broader perspective on what is exciting. They understand that the world is an interesting place. Knowing a little something about everything equips them to play in other people's backyards. They are brave enough to get out of their comfort zones a bit, and driven to learn more once they get there.

Software developers are not really encouraged to ask questions, but data scientists are. We ask questions like:

- What things might you be able to learn from a given data set?

- What do you/your people really want to know about the world?

- What will it mean to you once you find out?

- After you understand what kind of information is available, try to come up with, say,

- 10 interesting questions you might explore/answer with access to the data set.

Collect data → Prepare data → Data Analyzing → Build model → Evaluate model → Visualize the results

# My DATA

Statistical information about my data

# What's Hard about Data Science

- Overcoming assumptions

- Making ad-hoc explanations of data patterns

- Overgeneralizing

- Communication

- Not checking enough (validate models, data pipeline integrity, etc.)

- Using statistical tests correctly

- Prototype → Production transitions

- Data pipeline complexity (who do you ask?)

# **Readings**

Read next week readings and complete it before next class.

- Chapter Two: Python Language Basics, IPython, and Jupyter Notebooks

- Chapter Three: Built-In Data Structures, Functions, and Files

python for data analysis 2nd edition pdf