**LECTURE 1**

# Course Overview

An overview of data science and the data science lifecycle

Data Science, Spring 2024 @ Knowledge Stream

Sana Jabbar

## Intro – Sana Jabbar

- BS Telecommunication Engineering @ FAST, Islamabad
- MS Electrical Engineering @ FAST, Lahore
- Since 2018 @ Lums, Research Associate in Computer Vision and Graphics Lab, formerly in Clinical and Translational Lab
- Background:  PhD Remote Sensing,
- Research interests
    - Interactive computational tools for earth observation and medicine.
    - Applications in Remote Sensing for Taxation Automation and in medical Diagnosis.

# What is Data Science?

Lecture 01

- Intros
- **What is data science?**
- The objective of this course?
- Course Overview
- Data Science Lifecycle

PRINCIPLES AND TECHNIQUES OF DATA SCIENCE

# What is Data Science?

**Definition:**

- Data science is an interdisciplinary field that uses scientific methods, processes, algorithms, and systems to extract insights and knowledge from structured and unstructured data.

**Joey Gonzalez**

**Data Science** is the application of data centric, computational, and inferential thinking to:

- Understand the world (science).
- Solve problems (engineering).

# What is Data Science?

**Definition:**
- Data science is an interdisciplinary field that uses scientific methods, processes, algorithms, and systems to extract insights and knowledge from structured and unstructured data.

**Interdisciplinary Nature:**
- Data science combines elements of computer science, mathematics, and domain expertise to solve complex problems.

**Structured Data:**
- Refers to data that is organized into a specific format, typically consisting of rows and columns, where each piece of information is stored in a well-defined manner.

**Unstructured Data:**
- Refers to data that lacks a specific, predefined structure, making it more challenging to organize and analyse compared to structured data.

# Objective

Lecture 01

- Intros
- What is data science?
- **The objective of this course?**
- Course Overview
- Data Science Lifecycle

PRINCIPLES AND TECHNIQUES OF DATA SCIENCE

# Course Objective

**Course Objectives:**

- Equip students with the **skills** needed to extract valuable insights from data.

**Course Objectives:**

- Equip students with the **skills** needed to extract valuable insights from data.
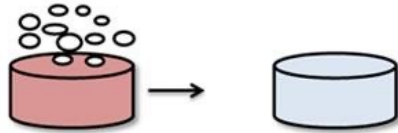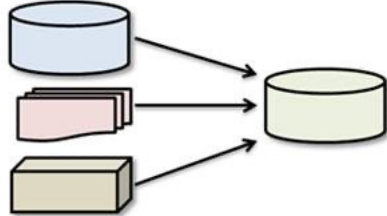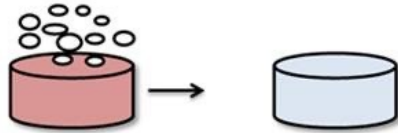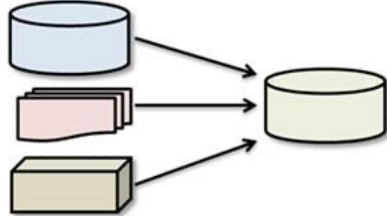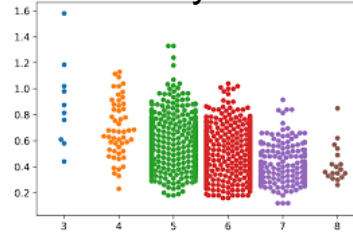
Data Collection

**Course Objectives:**

- Equip students with the **skills** needed to extract valuable insights from data.
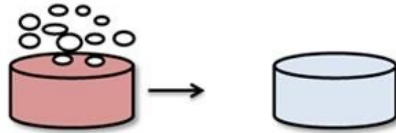
Data Collection

Data Cleaning

Data Integration

## Course Objectives:

- Equip students with the **skills** needed to extract valuable insights from data.
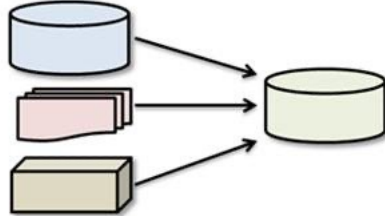
Data Collection

Data Cleaning

Data Integration

Exploratory Data Analysis

## Course Objectives:

- Equip students with the **skills** needed to extract valuable insights from data.

Data Cleaning

Data Collection

Data Integration

Exploratory Data Analysis

Machine Learning

# Course Objective

## Course Objectives:

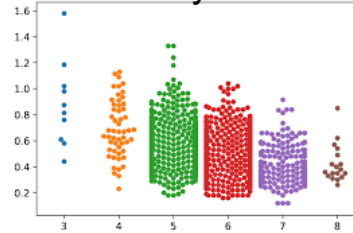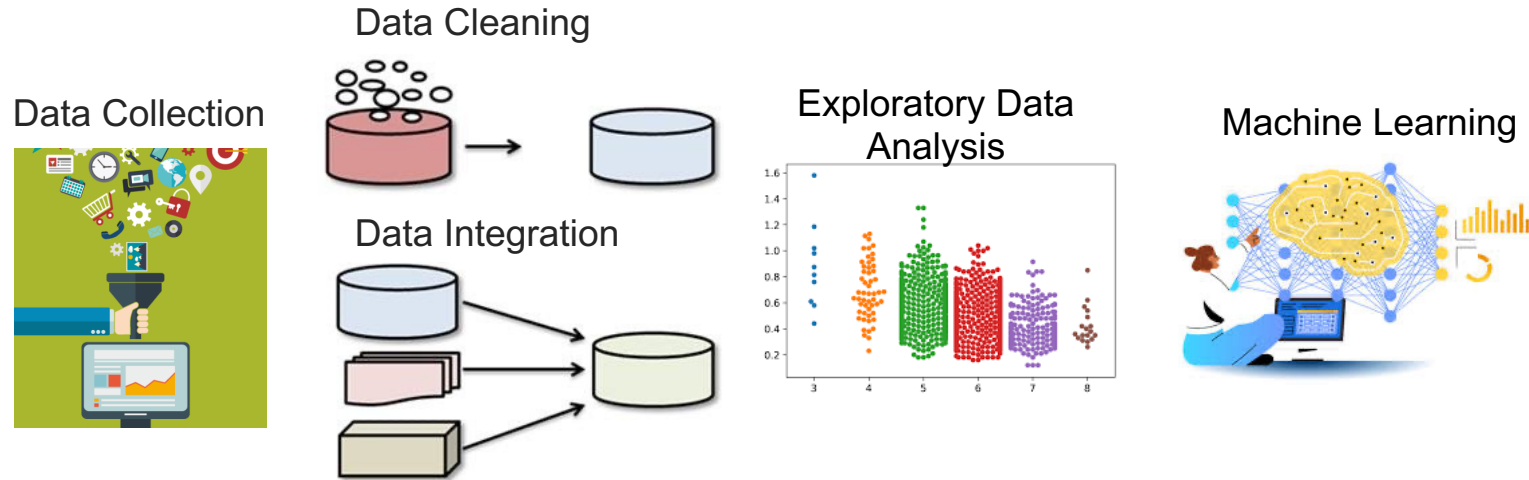- Equip students with the **skills** needed to extract valuable insights from data.



Data Cleaning

Data Collection

Data Integration

Exploratory Data Analysis

Machine Learning

- Prepare students for **real-world** data science challenges.
- Effectively communicate their findings to non-technical stakeholders

| Prepare | Prepare students for **data management**, **machine learning**, and **statistics**, by providing the necessary foundation and context. |
|---|---|
| Enable | Enable students to start careers as data scientists by providing experience working with **real-world data, tools, and techniques**. |
| Empower | Empower students to apply computational and inferential thinking to address **real-world problems**. |

## The world is complicated! Decisions are hard.

- Data science drives decision-making across various industries.

- There is a high demand for data scientists in today's job market.

- Data is used everywhere to answer hard questions and make tough decisions:

  - Science
  - Medicine
  - Engineering
  - Sports

Claims about data come up in discussing almost any important issue:
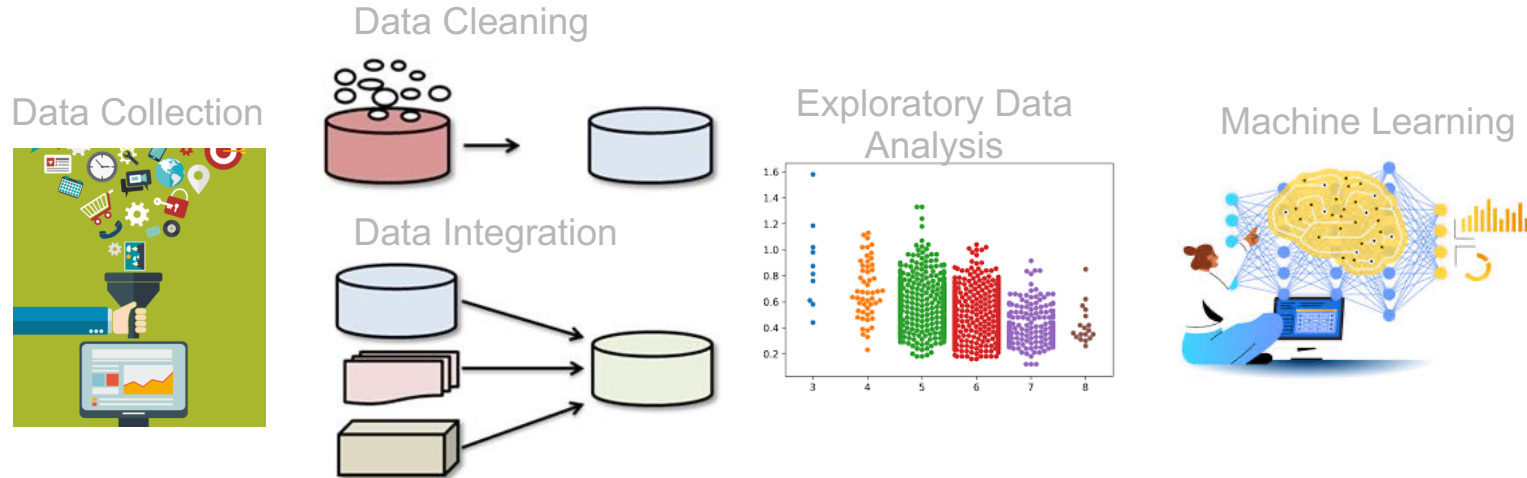
- Instead of "Alex says," now it's "the data says."
- It is usually not easy to tell what the data "says"
- **Empower yourself** to participate in the arguments that shape your life and your society

# Importance of Data Science

**Course Objectives:**

- Equip students with the **skills** needed to extract valuable insights from data.

Data Cleaning

Data Collection

Data Integration
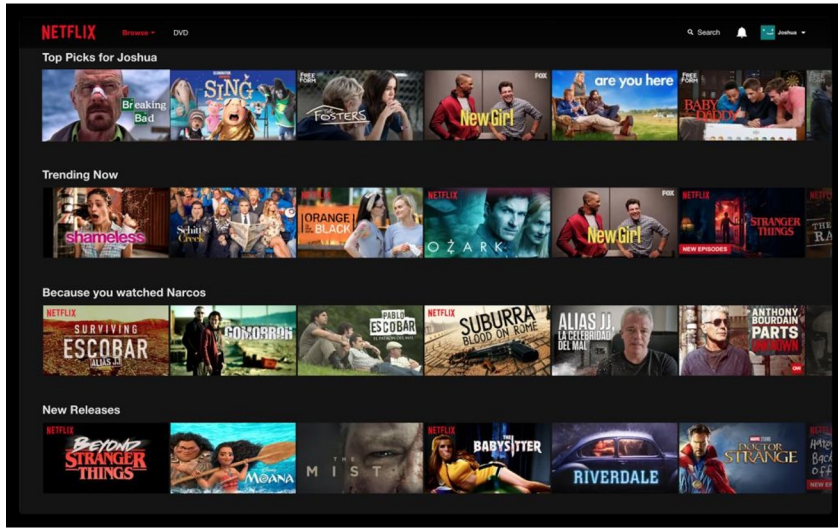
Exploratory Data Analysis

Machine Learning

- Prepare them for real-world data science challenges.
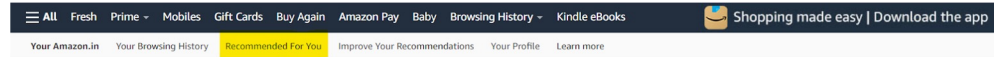
- Effectively communicate their findings to non-technical stakeholders

**Example:**

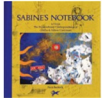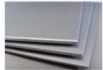- Imagine you work for an e-commerce

# Recommendation Systems

# First Image of a Black Hole

# Technology Trends

2020s ● **AI...?**

2010s ● Data Industry
➤ Collect and sell information

2000s ● Internet Industry
➤ Online retailers and services

1990s ● Software Industry
➤ Sold computer software

1980s ● Hardware Industry
➤ Sold computers

From Joey Gonzalez.

https://en.wikipedia.org/wiki/AlphaGo_versus_Lee_Sedol

21

# Data!

Knowledge is empowering.

Data science offers **immense potential** to address challenging problems facing society.

The future is in your hands, and I believe:

### You will use your knowledge for good.

…I am thrilled to teach Data Science :-)

**The world is complicated! Decisions are hard.**

Data science is a fundamentally human-centered field that facilitates decision-making by quantitatively balancing tradeoffs.

- To quantify things reliably we must:
  - Find relevant data;
  - Recognize its limitations;
  - Ask the right questions;
  - Make reasonable assumptions;
  - Conduct an appropriate analysis; and
  - Synthesize and explain our insights.

- Apply critical thinking and skepticism at every step

- Consider how our decisions affect others.

## The world is complicated! Decisions are hard.

Data science is a fundamentally human-centered field that facilitates decision-making by quantitatively balancing tradeoffs.

- To quantify things reliably we must:
  - Find relevant data;
  - Recognize its limitations;
  - Ask the right questions;
  - Make reasonable assumptions;
  - Conduct an appropriate analysis; and
  - Synthesize and explain our insights.

- Apply critical thinking and skepticism at every step

- Consider how our decisions affect others.

After this course, you should be able to take data and produce useful insights on the world's most challenging and ambiguous problems.

# Importance of Data Science



TASKS
(major involvement only)

**39%** ORGANIZING AND GUIDING TEAM PROJECTS

**36%** IMPLEMENTING MODELS/ ALGORITHMS INTO PRODUCTION

**43%** DEVELOPING PROTOTYPE MODELS

**43%** FEATURE EXTRACTION

**32%** COLLABORATING ON CODE PROJECTS (READING/EDITING OTHERS' CODE, USING GIT)

**47%** IDENTIFYING BUSINESS PROBLEMS TO BE SOLVED WITH ANALYTICS

**31%** TEACHING/TRAINING OTHERS

**30%** PLANNING LARGE SOFTWARE PROJECTS OR DATA SYSTEMS

**49%** CREATING VISUALIZATIONS

**30%** DEVELOPING DASHBOARDS

**53%** DATA CLEANING

**28%** COMMUNICATING WITH PEOPLE OUTSIDE YOUR COMPANY

**29%** ETL

**58%** COMMUNICATING FINDINGS TO BUSINESS DECISION-MAKERS

**20%** DEVELOPING DATA ANALYTICS SOFTWARE

**24%** SETTING UP / MAINTAINING DATA PLATFORMS

**61%** CONDUCTING DATA ANALYSIS TO ANSWER RESEARCH QUESTIONS

**19%** DEVELOPING PRODUCTS THAT DEPEND ON REAL-TIME DATA ANALYTICS

**19%** USING DASHBOARDS AND SPREADSHEETS (MADE BY OTHERS) TO MAKE DECISIONS

**69%** BASIC EXPLORATORY DATA ANALYSIS

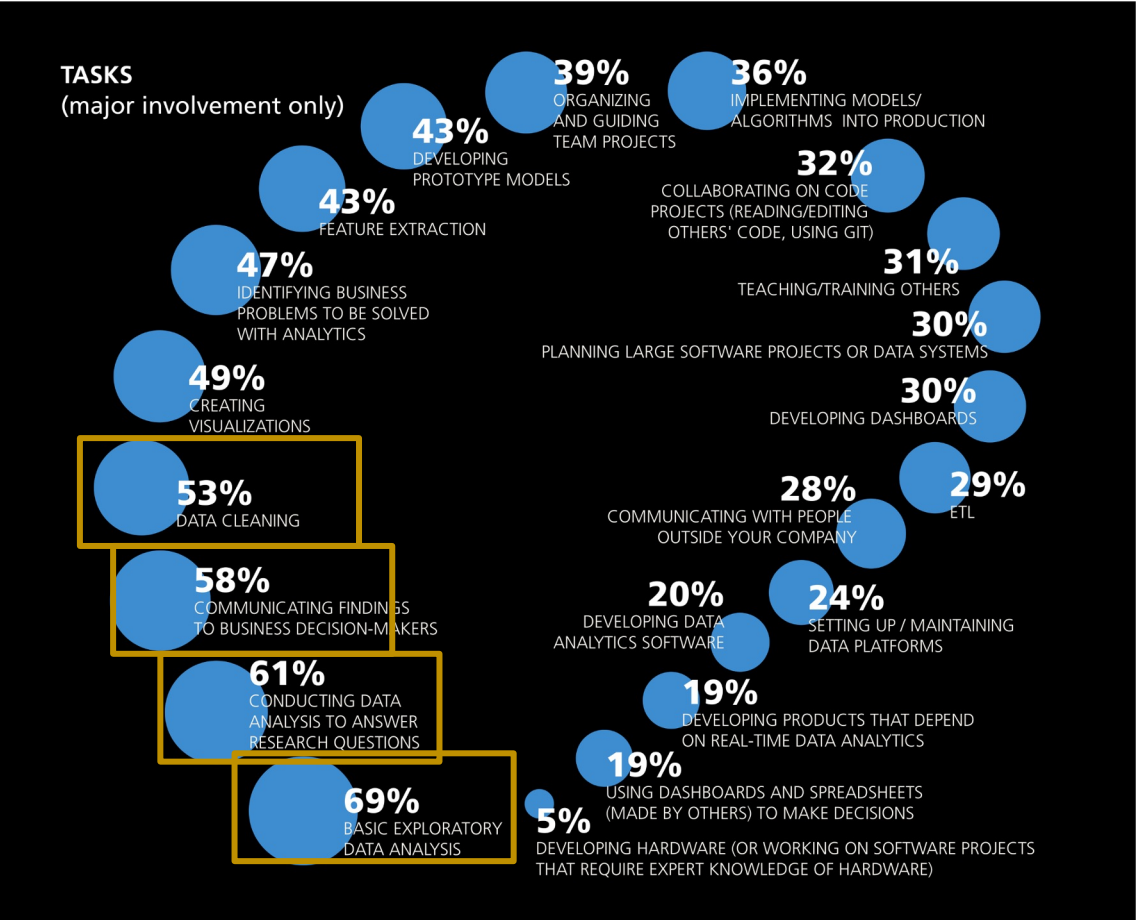**5%** DEVELOPING HARDWARE (OR WORKING ON SOFTWARE PROJECTS THAT REQUIRE EXPERT KNOWLEDGE OF HARDWARE)

The major tasks that data scientists say they work on regularly.

Based on the results of the 2016 Data Science Salary Survey.

25

# Importance of Data Science



TASKS
(major involvement only)

**43%** DEVELOPING PROTOTYPE MODELS

**43%** FEATURE EXTRACTION

**39%** ORGANIZING AND GUIDING TEAM PROJECTS

**36%** IMPLEMENTING MODELS/ ALGORITHMS INTO PRODUCTION

**32%** COLLABORATING ON CODE PROJECTS (READING/EDITING OTHERS' CODE, USING GIT)

**47%** IDENTIFYING BUSINESS PROBLEMS TO BE SOLVED WITH ANALYTICS

**31%** TEACHING/TRAINING OTHERS

**30%** PLANNING LARGE SOFTWARE PROJECTS OR DATA SYSTEMS

**49%** CREATING VISUALIZATIONS

**30%** DEVELOPING DASHBOARDS

**53%** DATA CLEANING

**28%** COMMUNICATING WITH PEOPLE OUTSIDE YOUR COMPANY

**29%** ETL

**58%** COMMUNICATING FINDINGS TO BUSINESS DECISION-MAKERS

**61%** CONDUCTING DATA ANALYSIS TO ANSWER RESEARCH QUESTIONS

**20%** DEVELOPING DATA ANALYTICS SOFTWARE

**24%** SETTING UP / MAINTAINING DATA PLATFORMS

**19%** DEVELOPING PRODUCTS THAT DEPEND ON REAL-TIME DATA ANALYTICS

**69%** BASIC EXPLORATORY DATA ANALYSIS

**19%** USING DASHBOARDS AND SPREADSHEETS (MADE BY OTHERS) TO MAKE DECISIONS

**5%** DEVELOPING HARDWARE (OR WORKING ON SOFTWARE PROJECTS THAT REQUIRE EXPERT KNOWLEDGE OF HARDWARE)

The major tasks that data scientists say they work on regularly.

Based on the results of the 2016 Data Science Salary Survey.

26

# Data Science Requires Engineering and Scientific Insight

**Good data analysis is not:**

- The simple application of a statistics recipe.
- Simple application of statistical software.

There are many **tools** out there for data science, but they are merely tools.

- **They don't do any of the important thinking!**

"The purpose of computing is insight, not numbers."

R. Hamming. *Numerical Methods for Scientists and Engineers (1962).*

27

# Example Questions in Data Science

Some (broad) questions we might try to answer with data science:

- What show should we recommend to our users to watch?
- In which markets should we focus our advertising campaign?
- Should I send my kids to daycare?
- Is the world getting better or worse?
- What areas of the world are at higher risk for climate change impact in 10 years? 20?
- What should we eat to avoid dying early of heart disease?
- Do immigrants from poor countries have a positive or negative impact on the economy?
- Which university will be the most appropriate for Data science engineering?

# Course Overview

Lecture 01

- Intros
- What is data science?
- The objective of this course?
- **Course Overview**
- Data Science Lifecycle

# Tentative List of Topics to be Covered in Data Science

- Pandas and NumPy
- Relational Databases & SQL
- Exploratory Data Analysis
- Regular Expressions
- Visualization
  - matplotlib
  - Seaborn
  - plotly
- Sampling

- Model design and loss formulation
- Linear Regression
- Feature Engineering
- Regularization, Bias-Variance Tradeoff, Cross-Validation
- Gradient Descent
- Data science in the physical world
- Logistic Regression
- Clustering
- PCA

# Programming Environment for our Course: Jupyter Notebook

# Data Science Lifecycle

Lecture 01

The data science lifecycle is a **high-level description** of the data science workflow.

Note the two distinct entry points!



Ask a Question

Obtain Data

Understand the World

Understand the Data

Reports, Decisions, and Solutions
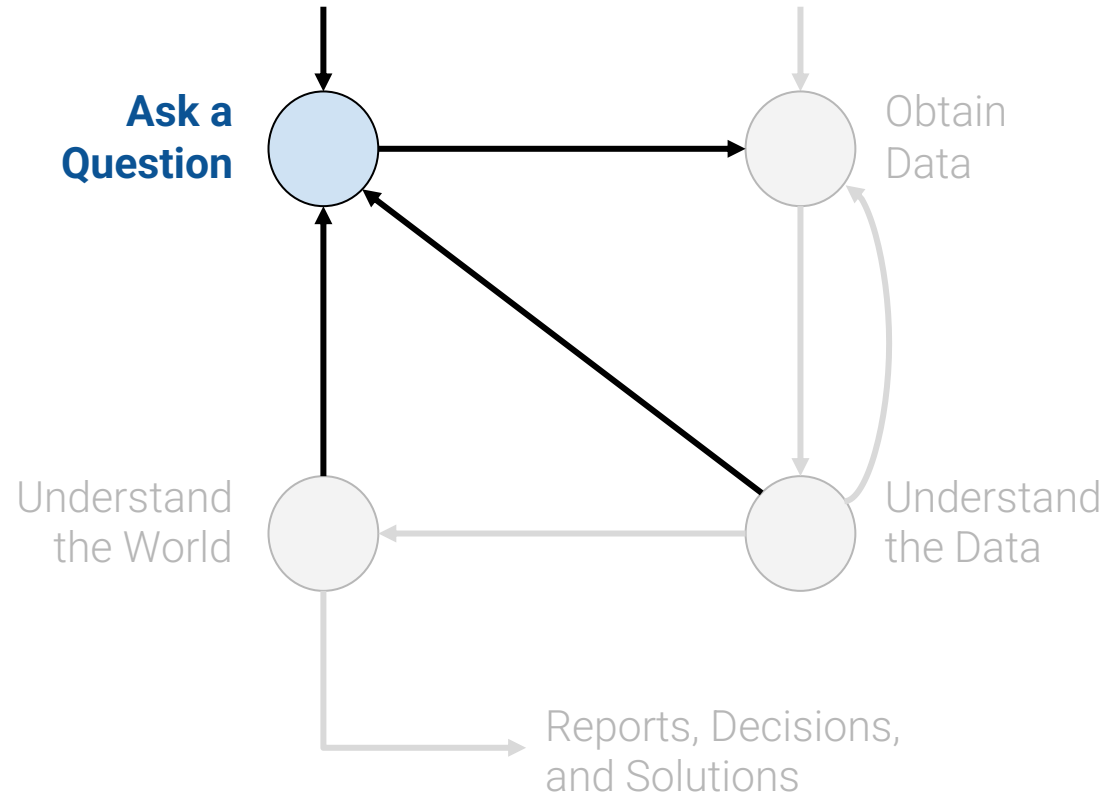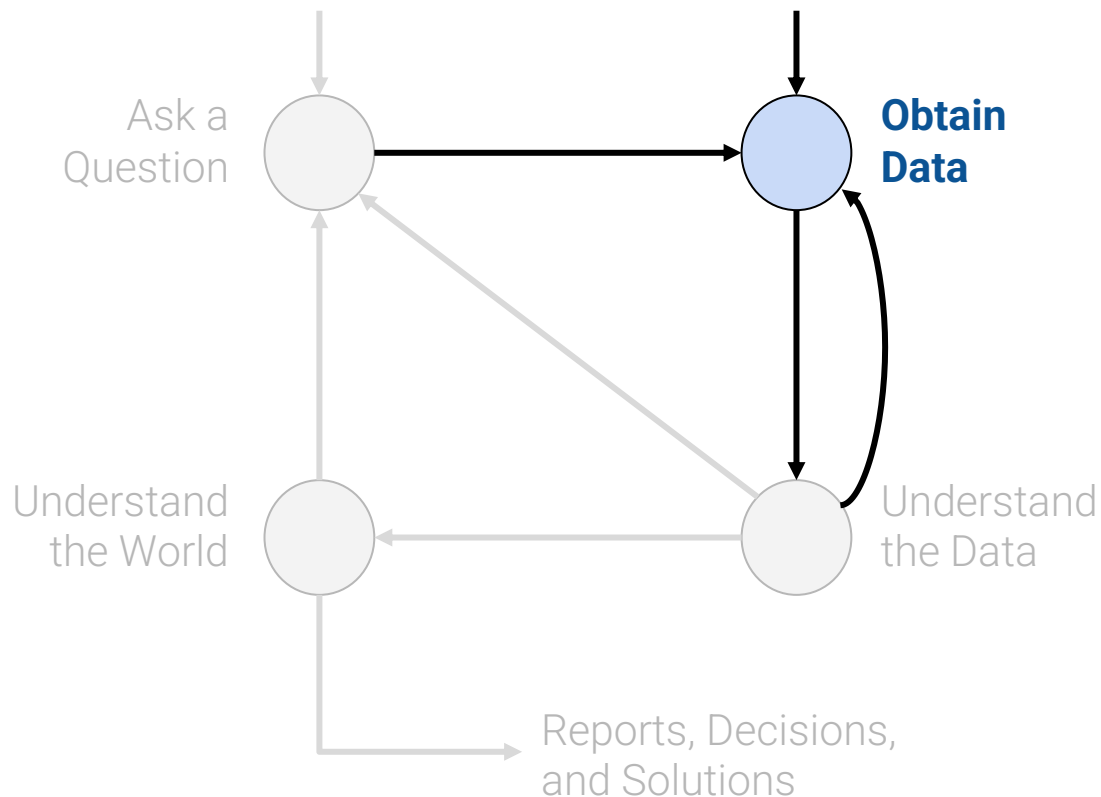
34

# 1. Question/Problem Formulation

- What do we want to know?
- What problems are we trying to solve?
- What hypotheses do we want to test?
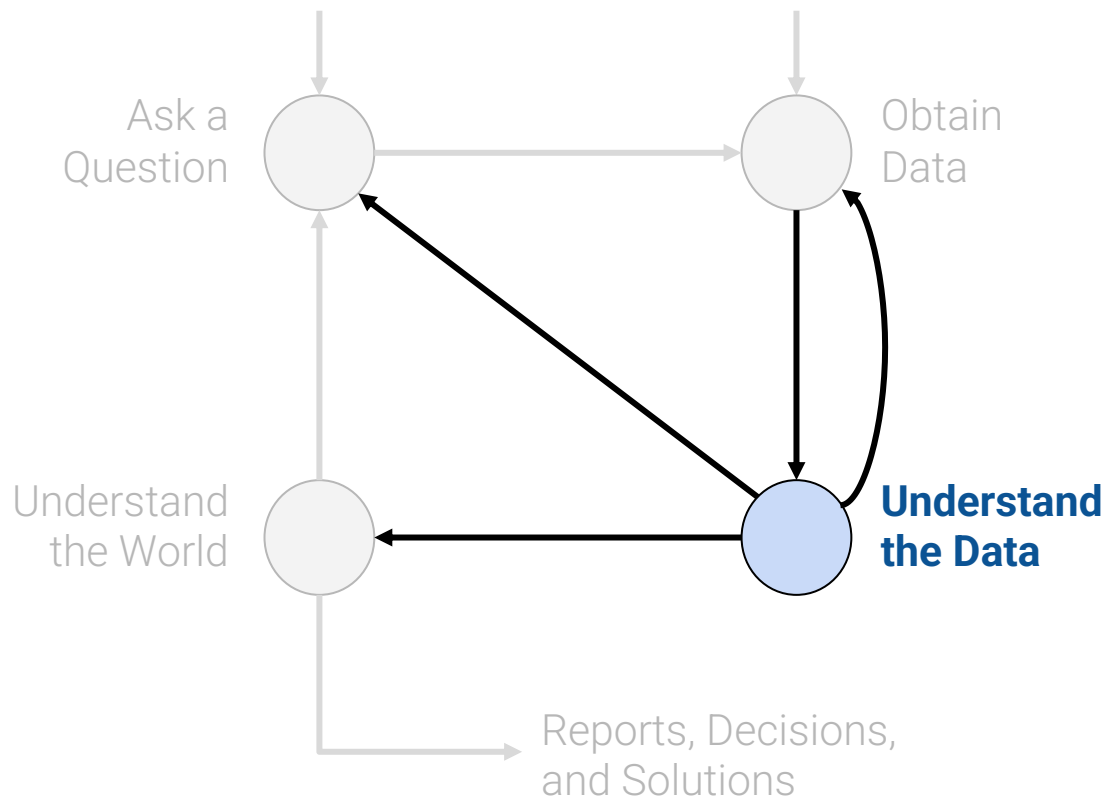- What are our metrics for success?

## 2. Data Acquisition and Cleaning

- What data do we have and what data do we need?
- How will we sample more data?
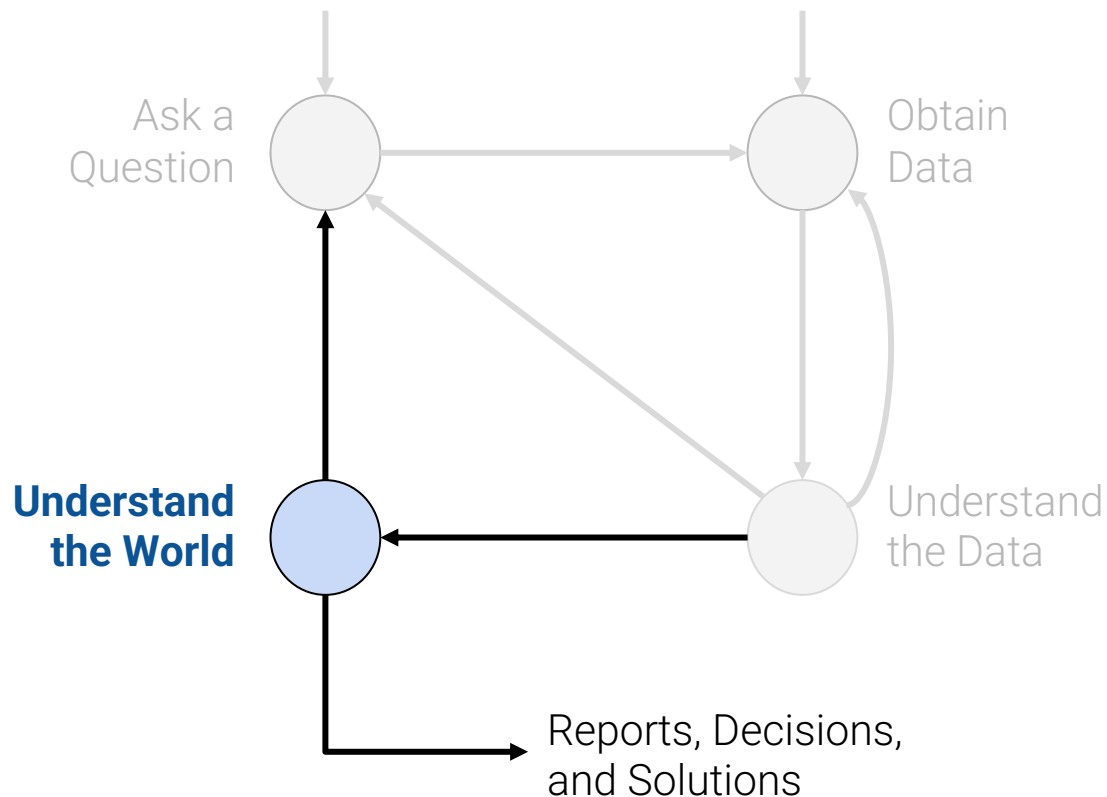- Is our data representative of the population we want to study?

# 3. Exploratory Data Analysis & Visualization

- How is our data organized and what does it contain?
- Do we already have relevant data?
- What are the biases, anomalies, or other issues with the data?
- How do we transform the data to enable effective analysis?



Ask a Question

Obtain Data

Understand the World

**Understand the Data**

Reports, Decisions, and Solutions

# 4. Prediction and Inference

- What does the data say about the world?
- Does it answer our questions or accurately solve the problem?
- How robust are our conclusions and can we trust the predictions?



Ask a Question

Obtain Data

**Understand the World**

Understand the Data

Reports, Decisions, and Solutions

# Setup Framework!