# Santa Monica Data Academy

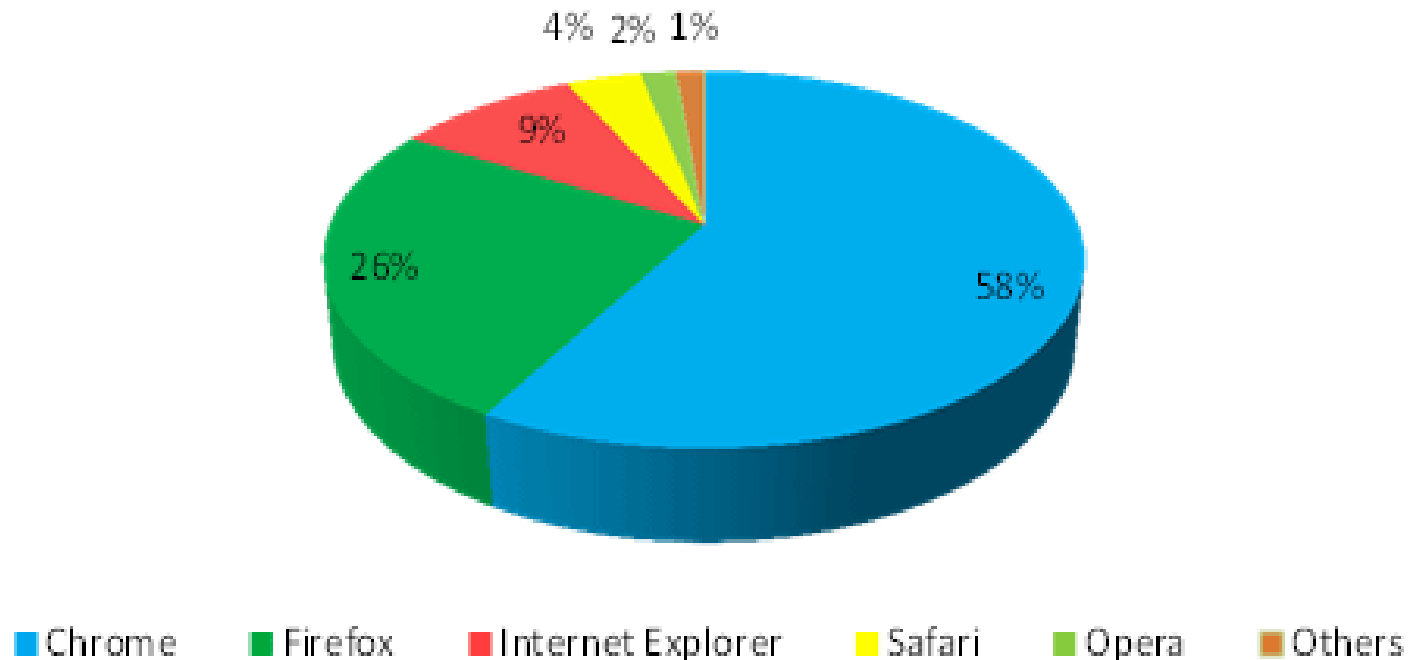**DA101**

Data Basics

# Welcome!

- Your name

- Your role (i.e. what you *actually do*)

- Why are you taking this class?

# Is it data?

- 4 examples

- In your group, decide if the answer is **Yes** or **No**

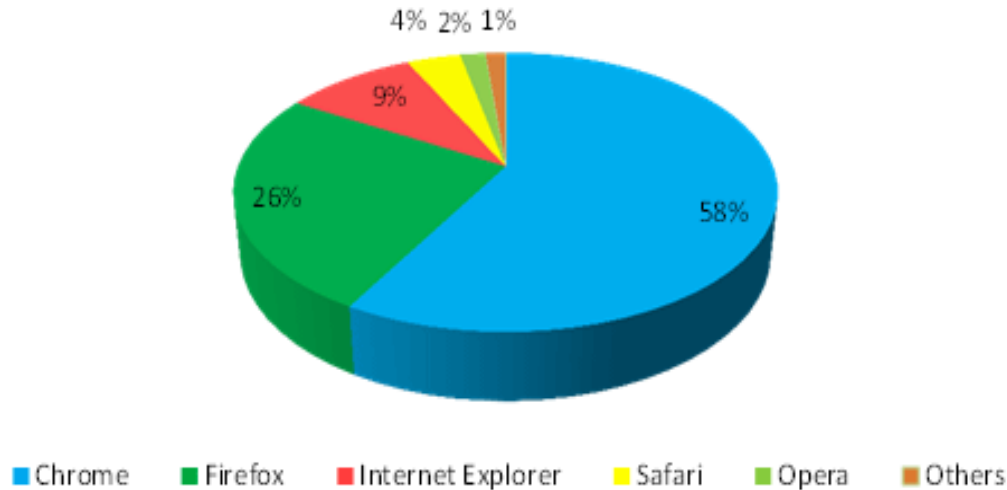- Come up with one or two reasons **why** you chose your answer

Browser Usage Statistics: 2014

**NO:** **This is a data visualization** (*pie chart*)



Browser Usage Statistics: 2014

- **What is the total that each percentage is derived from?**

- **How many "Others" are there, and what are they?**

# Is it data?

|  | Monthly Cash Flow | | |
|---|---|---|---|
|  | Actual | Budget | Variance |
| **Cash received** |  |  |  |
| Fees | $21,571 | $20,000 | $1,571 |
| Salary grants | 10,005 | 11,000 | (995) |
| Other | 76 |  | 76 |
|  | 31,652 | 31,000 | 652 |
| **Cash paid out** |  |  |  |
| Salaries and benefits | 21,575 | 20,000 | (1,575) |
| Food | 2,350 | 2,000 | (350) |
| Play supplies | 335 | 500 | 165 |
| Other | 3,270 | 1,500 | (1,770) |
|  | 27,530 | 24,000 | (3,530) |
| Excess of cash received over cash paid out | $4,122 | $7,000 | $(2,878) |

# Is it data?

**NO:** This is a **report**

| Monthly Cash Flow | Actual | Budget | Variance |
|---|---|---|---|
| **Cash received** | | | |
| Fees | $21,571 | $20,000 | $1,571 |
| Salary grants | 10,005 | 11,000 | (995) |
| Other | 76 | | 76 |
| | 31,652 | 31,000 | 652 |
| **Cash paid out** | | | |
| Salaries and benefits | 21,575 | 20,000 | (1,575) |
| Food | 2,350 | 2,000 | (350) |
| Play supplies | 335 | 500 | 165 |
| Other | 3,270 | 1,500 | (1,770) |
| | 27,530 | 24,000 | (3,530) |
| Excess of cash received over cash paid out | $4,122 | $7,000 | $(2,878) |

- **How does this month compare to last month?**

- **Which category is the most over-budget?**

# Is it data?

# Is it data?

**NO:** This is a **map/visualization**



- **Which regions/states are seeing the most change?**

- **What does a "report" mean? Are there numbers to back it up?**

# Is it data?

University of Dhaka

Faculty of Arts

## KHA - Unit

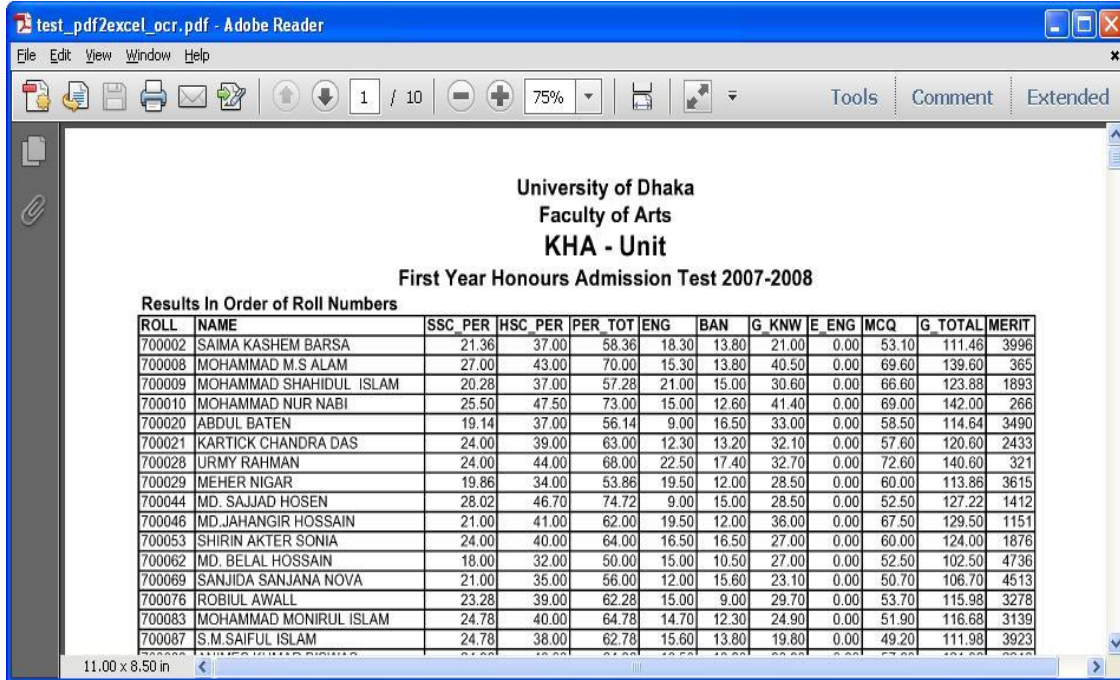First Year Honours Admission Test 2007-2008

**Results In Order of Roll Numbers**

| ROLL | NAME | SSC_PER | HSC_PER | PER_TOT | ENG | BAN | G_KNW | E_ENG | MCQ | G_TOTAL | MERIT |
|------|------|---------|---------|---------|-----|-----|-------|-------|-----|---------|-------|
| 700002 | SAIMA KASHEM BARSA | 21.36 | 37.00 | 58.36 | 18.30 | 13.80 | 21.00 | 0.00 | 53.10 | 111.46 | 3996 |
| 700008 | MOHAMMAD M.S ALAM | 27.00 | 43.00 | 70.00 | 15.30 | 13.80 | 40.50 | 0.00 | 69.60 | 139.60 | 365 |
| 700009 | MOHAMMAD SHAHIDUL  ISLAM | 20.28 | 37.00 | 57.28 | 21.00 | 15.00 | 30.60 | 0.00 | 66.60 | 123.88 | 1893 |
| 700010 | MOHAMMAD NUR NABI | 25.50 | 47.50 | 73.00 | 15.00 | 12.60 | 41.40 | 0.00 | 69.00 | 142.00 | 266 |
| 700020 | ABDUL BATEN | 19.14 | 37.00 | 56.14 | 9.00 | 16.50 | 33.00 | 0.00 | 58.50 | 114.64 | 3490 |
| 700021 | KARTICK CHANDRA DAS | 24.00 | 39.00 | 63.00 | 12.30 | 13.20 | 32.10 | 0.00 | 57.60 | 120.60 | 2433 |
| 700028 | URMY RAHMAN | 24.00 | 44.00 | 68.00 | 22.50 | 17.40 | 32.70 | 0.00 | 72.60 | 140.60 | 321 |
| 700029 | MEHER NIGAR | 19.86 | 34.00 | 53.86 | 19.50 | 12.00 | 28.50 | 0.00 | 60.00 | 113.86 | 3615 |
| 700044 | MD. SAJJAD HOSEN | 28.02 | 46.70 | 74.72 | 9.00 | 15.00 | 28.50 | 0.00 | 52.50 | 127.22 | 1412 |
| 700046 | MD.JAHANGIR HOSSAIN | 21.00 | 41.00 | 62.00 | 19.50 | 12.00 | 36.00 | 0.00 | 67.50 | 129.50 | 1151 |
| 700053 | SHIRIN AKTER SONIA | 24.00 | 40.00 | 64.00 | 16.50 | 16.50 | 27.00 | 0.00 | 60.00 | 124.00 | 1876 |
| 700062 | MD. BELAL HOSSAIN | 18.00 | 32.00 | 50.00 | 15.00 | 10.50 | 27.00 | 0.00 | 52.50 | 102.50 | 4736 |
| 700069 | SANJIDA SANJANA NOVA | 21.00 | 35.00 | 56.00 | 12.00 | 15.60 | 23.10 | 0.00 | 50.70 | 106.70 | 4513 |
| 700076 | ROBIUL AWALL | 23.28 | 39.00 | 62.28 | 15.00 | 9.00 | 29.70 | 0.00 | 53.70 | 115.98 | 3278 |
| 700083 | MOHAMMAD MONIRUL ISLAM | 24.78 | 40.00 | 64.78 | 14.70 | 12.30 | 24.90 | 0.00 | 51.90 | 116.68 | 3139 |
| 700087 | S.M.SAIFUL ISLAM | 24.78 | 38.00 | 62.78 | 15.60 | 13.80 | 19.80 | 0.00 | 49.20 | 111.98 | 3923 |

11.00 x 8.50 in

# Is it data?

## NO: This is a PDF

test_pdf2excel_ocr.pdf - Adobe Reader

File  Edit  View  Window  Help

1 / 10    75%    Tools    Comment    Extended

**University of Dhaka**
**Faculty of Arts**
**KHA - Unit**
**First Year Honours Admission Test 2007-2008**

Results In Order of Roll Numbers

| ROLL | NAME | SSC_PER | HSC_PER | PER_TOT | ENG | BAN | G_KNW | E_ENG | MCQ | G_TOTAL | MERIT |
|------|------|---------|---------|---------|-----|-----|-------|-------|-----|---------|-------|
| 700002 | SAIMA KASHEM BARSA | 21.36 | 37.00 | 58.36 | 18.30 | 13.80 | 21.00 | 0.00 | 53.10 | 111.46 | 3996 |
| 700008 | MOHAMMAD M.S ALAM | 27.00 | 43.00 | 70.00 | 15.30 | 13.80 | 40.50 | 0.00 | 69.60 | 139.60 | 365 |
| 700009 | MOHAMMAD SHAHIDUL  ISLAM | 20.28 | 37.00 | 57.28 | 21.00 | 15.00 | 30.60 | 0.00 | 66.60 | 123.88 | 1893 |
| 700010 | MOHAMMAD NUR NABI | 25.50 | 47.50 | 73.00 | 15.00 | 12.60 | 41.40 | 0.00 | 69.00 | 142.00 | 266 |
| 700020 | ABDUL BATEN | 19.14 | 37.00 | 56.14 | 9.00 | 16.50 | 33.00 | 0.00 | 58.50 | 114.64 | 3490 |
| 700021 | KARTICK CHANDRA DAS | 24.00 | 39.00 | 63.00 | 12.30 | 13.20 | 32.10 | 0.00 | 57.60 | 120.60 | 2433 |
| 700028 | URMY RAHMAN | 24.00 | 44.00 | 68.00 | 22.50 | 17.40 | 32.70 | 0.00 | 72.60 | 140.60 | 321 |
| 700029 | MEHER NIGAR | 19.86 | 34.00 | 53.86 | 19.50 | 12.00 | 28.50 | 0.00 | 60.00 | 113.86 | 3615 |
| 700044 | MD. SAJJAD HOSEN | 28.02 | 46.70 | 74.72 | 9.00 | 15.00 | 28.50 | 0.00 | 52.50 | 127.22 | 1412 |
| 700046 | MD.JAHANGIR HOSSAIN | 21.00 | 41.00 | 62.00 | 19.50 | 12.00 | 36.00 | 0.00 | 67.50 | 129.50 | 1151 |
| 700053 | SHIRIN AKTER SONIA | 24.00 | 40.00 | 64.00 | 16.50 | 16.50 | 27.00 | 0.00 | 60.00 | 124.00 | 1876 |
| 700062 | MD. BELAL HOSSAIN | 18.00 | 32.00 | 50.00 | 15.00 | 10.50 | 27.00 | 0.00 | 52.50 | 102.50 | 4736 |
| 700069 | SANJIDA SANJANA NOVA | 21.00 | 35.00 | 56.00 | 12.00 | 15.60 | 23.10 | 0.00 | 50.70 | 106.70 | 4513 |
| 700076 | ROBIUL AWALL | 23.28 | 39.00 | 62.28 | 15.00 | 9.00 | 29.70 | 0.00 | 53.70 | 115.98 | 3278 |
| 700083 | MOHAMMAD MONIRUL ISLAM | 24.78 | 40.00 | 64.78 | 14.70 | 12.30 | 24.90 | 0.00 | 51.90 | 116.68 | 3139 |
| 700087 | S.M.SAIFUL ISLAM | 24.78 | 38.00 | 62.78 | 15.60 | 13.80 | 19.80 | 0.00 | 49.20 | 111.98 | 3923 |

11.00 x 8.50 in

- **What is the average MERIT score?**

- **Who had the highest overall performance?**

# Learning Objectives

- Understand what **data** is*

- Learn basic data **vocabulary**

- Perform basic **operations** on data

What do we mean when we say **data**?

# What do we mean when we say data?

Wikipedia: *Data (computing)*, October 2018

**Data** is any sequence of one or more symbols given meaning by specific act(s) of interpretation.

Data [...] requires interpretation to become information.

# What do we mean when we say data?

Wikipedia: *Data (computing)*, October 2018

**Data** is any sequence of one or more symbols given meaning by specific act(s) of interpretation.

**Data [...] requires interpretation to become information.**

# What do we mean when we say data?

Wikipedia: *Data*, October 2018

**Data** is a set of values of qualitative or quantitative variables.

Data is measured, collected and reported, and analyzed, whereupon it can be visualized using graphs, images or other analysis tools.

# What do we mean when we say data?

**Data** is a set of values of **qualitative** or **quantitative** variables.

Data is measured, collected and reported, and analyzed, whereupon it can be visualized using graphs, images or other analysis tools.

# What do we mean when we say data?

Wikipedia: *Data*, October 2018 https://en.wikipedia.org/wiki/Data

**Data** is a set of values of qualitative or quantitative variables.

Data is **measured**, **collected** and **reported**, and **analyzed**, whereupon it can be **visualized** using graphs, images or other analysis tools.

What do we *really* mean when we say **data**?

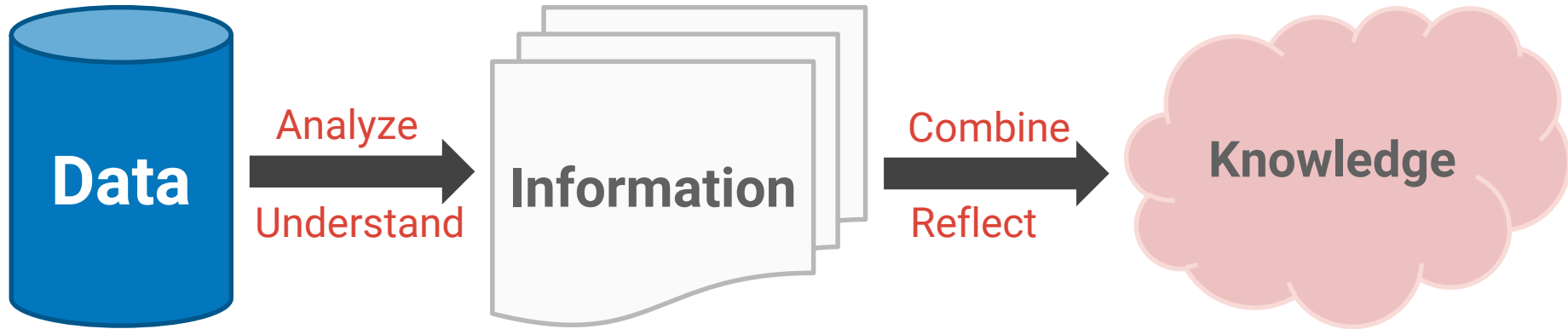# What do we *really* mean when we say data?

- Digital

- Raw

- Structured

# What do we *really* mean when we say data?

- **Digital** (so we can use software tools)

- As **Raw** As Possible

- As **Structured** As Possible

Let's **zoom out** a little…

# The Bigger Picture

# The Bigger Picture

**CCS Fees**

**CCS Activity Registrations**

Analyze

Understand

Avg. fee per Age Group

Ratio of Resident to Non-Resident payments

Activity level by Age Group

Combine

Reflect

Identify under-served population

This process is called

**data driven decision making**

# Organizing data with **tables**

# Also known as *tabular data*

| Employee ID | Name | Date of Birth | Zip Code |
|:---:|:---:|:---:|:---:|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

# The header

| Employee ID | Name | Date of Birth | Zip Code |
|:---:|:---:|:---:|:---:|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

Describes the data by naming individual *attributes*

Also known as the *fields* or *properties* of the data

| Employee ID | Name | Date of Birth | Zip Code |
|---|---|---|---|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

# The rows

| Employee ID | Name | Date of Birth | Zip Code |
|---|---|---|---|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

# The rows

Individual items or observations in the data

Also known as the *records* of data

A record contains *values* for all of the attributes

| Employee ID | Name | Date of Birth | Zip Code |
|---|---|---|---|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

# The columns

| Employee ID | Name | Date of Birth | Zip Code |
|:---:|:---:|:---:|:---:|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

# The columns

Contain values of the same *type* for all records

Ideally a column represents a *single attribute* of the data

| Employee ID | Name | Date of Birth | Zip Code |
|---|---|---|---|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

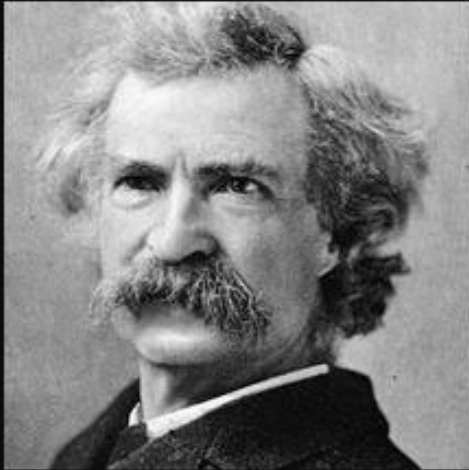Why does **structure** matter?

# Why does structure matter?

A **data structure** is a collection of data values, the relationships among them, and the functions or operations that can be applied to the data.

# Why does structure matter?

A **data structure** is a collection of data values, the relationships among them, and the **functions or operations that can be applied** to the data.
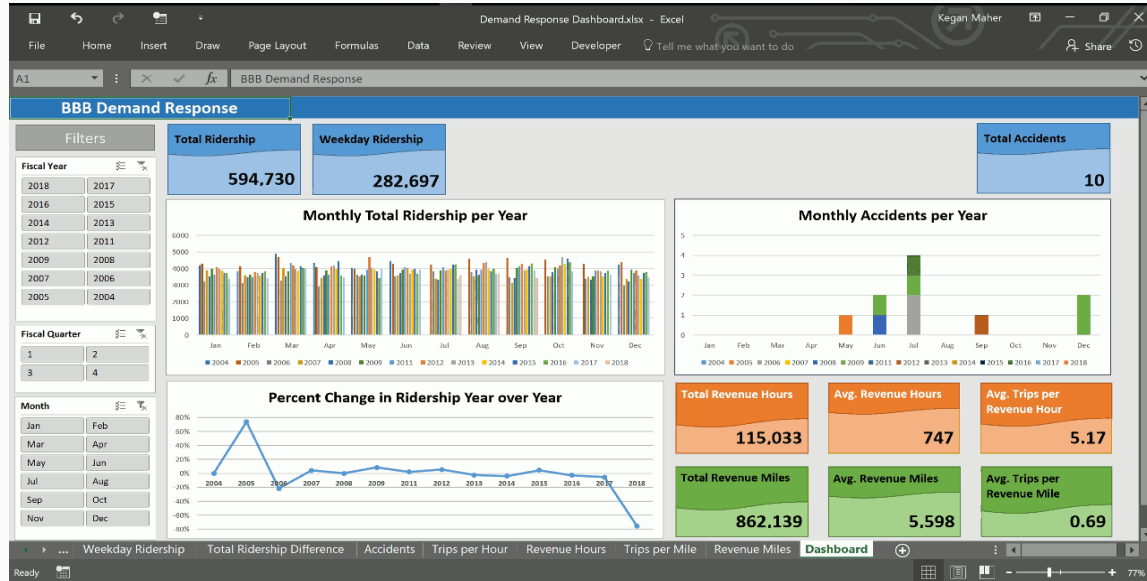
# Why does structure matter?



Data is like garbage. You'd better know what you are going to do with it before you collect it.

~ Mark Twain

# Data Collection and Storage

# Data Sources and Collection

# Spectrum of Data Sources and Destinations

# Spectrum of Data Sources and Destinations

## Initial investment/data collection





**Almost no time/$**

**How much time/$ do you have?**

## Flexibility for analysis





**Concrete**

**Water**

# Spectrum of Data Sources and Destinations



Pen & Paper

# Spectrum of Data Sources and Destinations



PP

Flat File
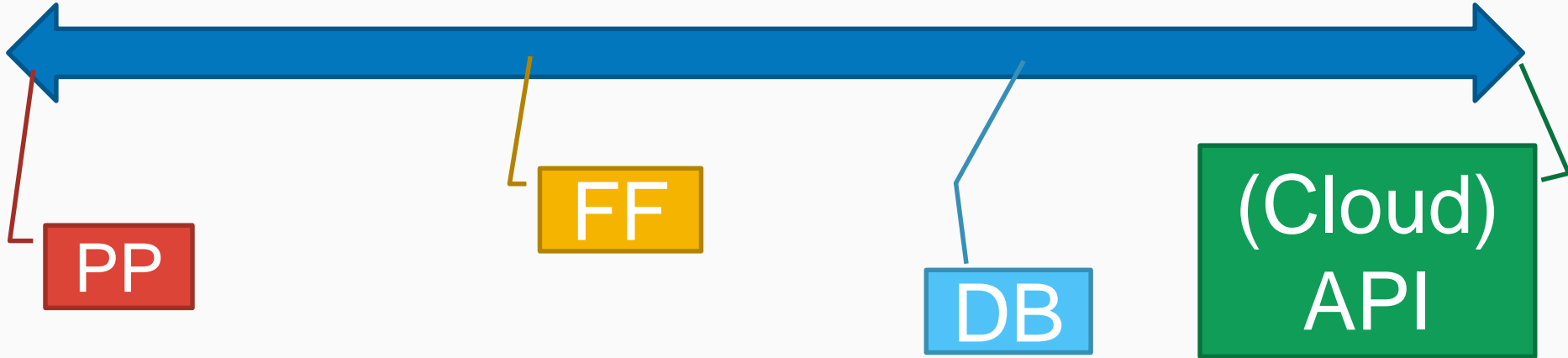
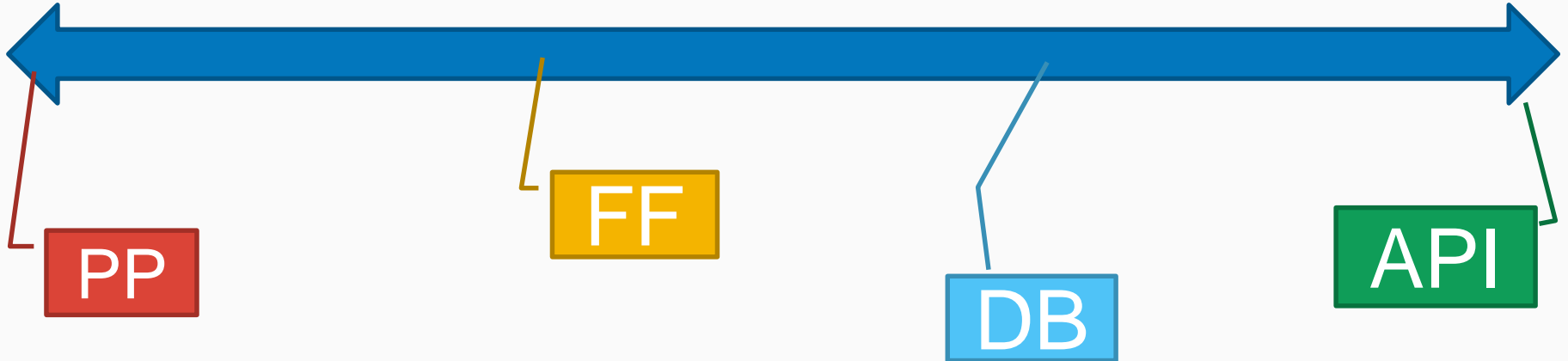# Spectrum of Data Sources and Destinations



PP

FF

(Local) Database

# Spectrum of Data Sources and Destinations

# Spectrum of Data Sources and Destinations



PP

FF

DB

API

Collecting data: a case study

**Homeless Demographic Survey**

# Homeless Demographic Survey

- 2016 and 2017

- Questions on age range, gender, race, medical conditions, etc.

## 2016 Homeless Demographic Survey

1. Have you ever been diagnosed with a serious medical condition?

2. Have you ever been diagnosed with a substance abuse issue?

3. Have you ever been diagnosed with a mental health issue?

# Collecting data: a case study

| ID | Year | Serious Medical Condition | Substance Abuse Issue | Mental Health Issue |
|----|------|---------------------------|-----------------------|---------------------|
| 0  | 2016 | Yes                       | Yes                   | No                  |
| 1  | 2016 | No                        | Yes                   | Yes                 |
| 2  | 2016 | No                        | No                    | Yes                 |

## 2017 Homeless Demographic Survey

1. Have you ever been diagnosed with any of the following?

a) Serious medical condition
b) Substance abuse issue
c) Mental health issue

# Collecting data: a case study

| ID | Year | Diagnosed With |
|----|------|----------------|
| 3 | 2017 | Serious medical condition |
| 4 | 2017 | Substance abuse issue, Mental health issue |
| 5 | 2017 | Serious medical condition, Mental health issue |

# Collecting data: a case study

| ID | Year | Serious Medical Condition | Substance Abuse Issue | Mental Health Issue |
|----|------|---------------------------|------------------------|---------------------|
| 0  | 2016 | Yes                       | Yes                    | No                  |
| 1  | 2016 | No                        | Yes                    | Yes                 |
| 2  | 2016 | No                        | No                     | Yes                 |

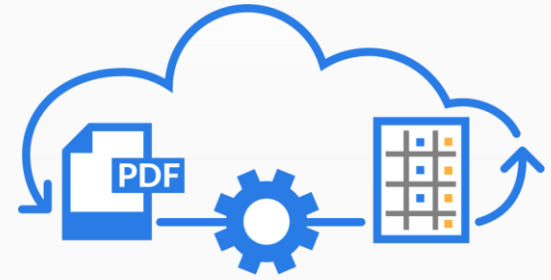| ID | Year | Diagnosed With | | |
|----|------|----------------|--|--|
| 3  | 2017 | Serious medical condition | | |
| 4  | 2017 | Substance abuse issue, Mental health issue | | |
| 5  | 2017 | Serious medical condition, Mental health issue | | |

# Santa Monica Data Academy

## 10 Minute Break

## BREAK IS OVER

# Design a data table (~15 minutes)

- Study the PDF form

- Come up with a list of columns
  Decide: *Qualitative* or *Quantitative?*

- Choose a *spokesperson* to share your group's design

## TIME TO PRESENT

# Organizing data with tables

| Employee ID | July 2018 Hours | December 2018 Hours |
|:---:|:---:|:---:|
| 0 | 44 | 0 |
| 1 | 20 | 28 |
| 2 | 48 | 48 |

# Organizing data with tables

| Employee ID | Date | Hours |
| --- | --- | --- |
| 0 | July 2018 | 44 |
| 0 | December 2018 | 0 |
| 1 | July 2018 | 20 |
| 1 | December 2018 | 28 |
| 2 | July 2018 | 48 |
| 2 | December 2018 | 48 |

# Does it matter?

| Employee ID | July 2018 Hours | December 2018 Hours |
|:---:|:---:|:---:|
| 0 | 44 | 0 |
| 1 | 20 | 28 |
| 2 | 48 | 48 |

| Employee ID | Date | Hours |
|:---:|:---:|:---:|
| 0 | July 2018 | 44 |
| 0 | December 2018 | 0 |
| 1 | July 2018 | 20 |
| 1 | December 2018 | 28 |
| 2 | July 2018 | 48 |
| 2 | December 2018 | 48 |

# Does it matter?

| Employee ID | Date | Hours | Hourly Rate |
|:---:|:---:|:---:|:---:|
| 0 | July 2018 | 44 | $18.00 |
| 0 | December 2018 | 0 | $19.00 |
| 1 | July 2018 | 20 | $15.00 |
| 1 | December 2018 | 28 | $15.00 |
| 2 | July 2018 | 48 | $25.00 |
| 2 | December 2018 | 48 | $30.00 |

This format of tabular data is called

**Tidy Data**

# Operations on tabular data

# Operations on tabular data

*What is the average age of the people in our data?*

| Employee ID | Name | DOB | Zip Code |
|:-----------:|:----:|:---:|:--------:|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

# Operations on tabular data

*What is the average age of the people in our data?*

| Employee ID | Name | DOB | Zip Code |
|:---:|:---:|:---:|:---:|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

# Operations on tabular data: Calculated/Derived Field

*What is the average age of the people in our data?*

| Employee ID | Name | DOB | Age | Zip Code |
|---|---|---|---|---|
| 20100 | Harry | 1980-11-12 | 38 | 90401 |
| 10401 | Tom | 1966-03-22 | 53 | 90405 |
| 52200 | Rick | 1986-09-28 | 32 | 90401 |

# Operations on tabular data: Column Aggregation (Sum)

*What is the average age of the people in our data?*

| Employee ID | Name | DOB | Age | Zip Code |
|:---:|:---:|:---:|:---:|:---:|
| 20100 | Harry | 1980-11-12 | 38 | 90401 |
| 10401 | Tom | 1966-03-22 | 53 | 90405 |
| 52200 | Rick | 1986-09-28 | 32 | 90401 |

# Operations on tabular data: Column Aggregation (Sum)

*What is the average age of the people in our data?*

**38 + 53 + 32 = 123**

| Employee ID | Name | DOB | Age | Zip Code |
|---|---|---|---|---|
| 20100 | Harry | 1980-11-12 | 38 | 90401 |
| 10401 | Tom | 1966-03-22 | 53 | 90405 |
| 52200 | Rick | 1986-09-28 | 32 | 90401 |

*What is the average age of the people in our data?*

$$(38 + 53 + 32 = \textbf{123}) / \textbf{3} = \textbf{41}$$

| Employee ID | Name | DOB | Age | Zip Code |
|---|---|---|---|---|
| 20100 | Harry | 1980-11-12 | 38 | 90401 |
| 10401 | Tom | 1966-03-22 | 53 | 90405 |
| 52200 | Rick | 1986-09-28 | 32 | 90401 |

# Operations Recap

- Calculate a new column *(DOB -> Age)*

- Aggregate values in a column *(average all the Ages)*

# Operations on tabular data

*What is the geographic distribution of people in our data?*

| Employee ID | Name | DOB | Zip Code |
|---|---|---|---|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

# Operations on tabular data: Group by Zip Code

*What is the geographic distribution of people in our data?*

| Employee ID | Name | DOB | Zip Code |
|---|---|---|---|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

# Operations on tabular data: Count within Groups

*What is the geographic distribution of people in our data?*

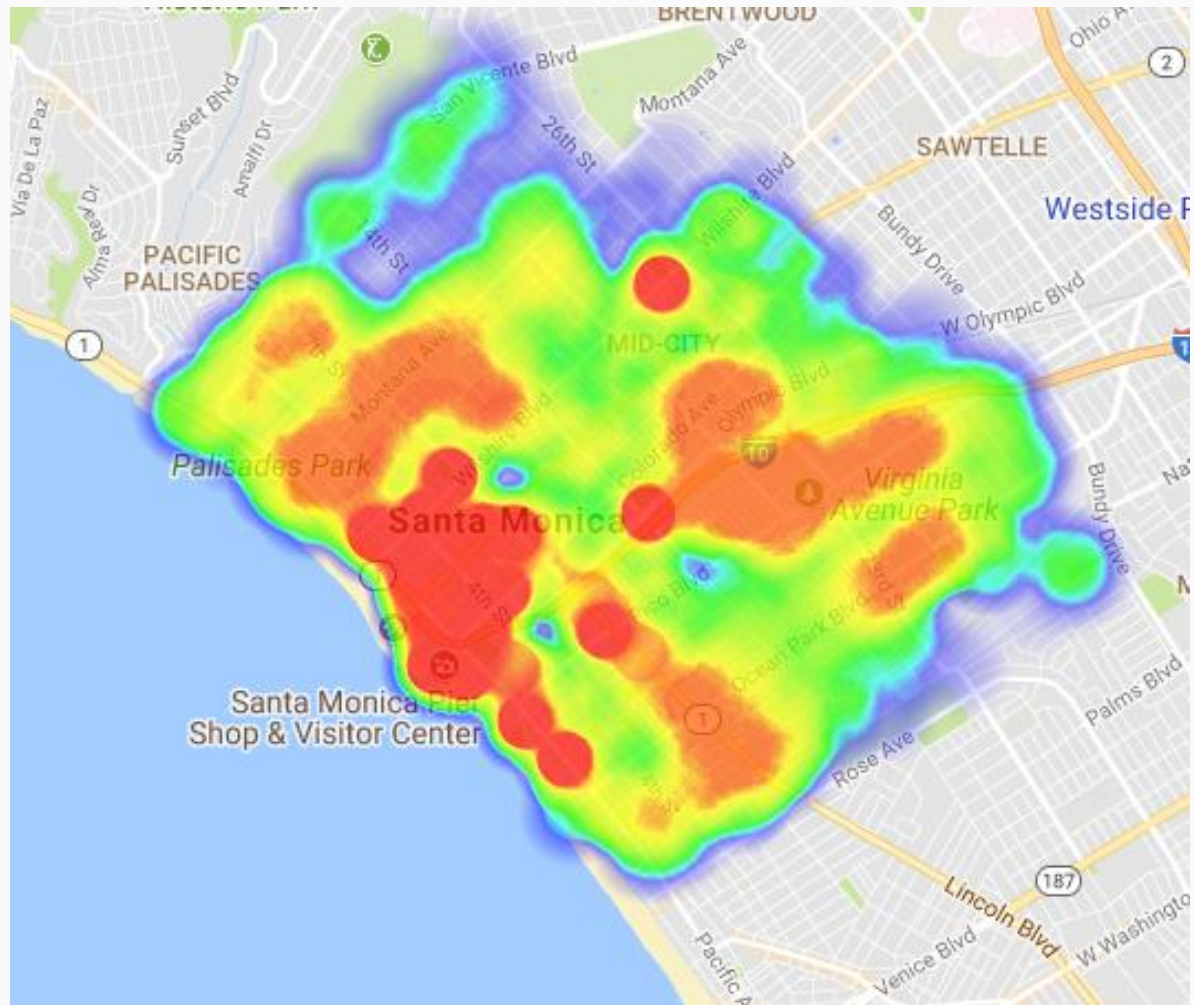| Employee ID | Name | DOB | Zip Code |
|---|---|---|---|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

# Operations on tabular data: Count within Groups

*What is the geographic distribution of people in our data?*

| Zip Code | Count |
|----------|-------|
| 90401 | 2 |
| 90405 | 1 |

# Operations on tabular data

## Visualize using a heatmap

# Operations Recap

- Group *(by Zip Code*)

- Count *(how many rows in each group*)

- Visualize *(a heatmap shows the distribution graphically*)

# Operations on tabular data

*Who is the youngest person that lives in 90401?*

| Employee ID | Name | DOB | Zip Code |
|---|---|---|---|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

*Who is the youngest person that lives in 90401?*

Zip Code = **90401**

| Employee ID | Name | DOB | Zip Code |
|---|---|---|---|
| 20100 | Harry | 1980-11-12 | 90401 |
| 10401 | Tom | 1966-03-22 | 90405 |
| 52200 | Rick | 1986-09-28 | 90401 |

# Operations on tabular data: **Filter**

*Who is the youngest person that lives in 90401?*

Zip Code = **90401**

| Employee ID | Name | DOB | Zip Code |
|:---:|:---:|:---:|:---:|
| 20100 | Harry | 1980-11-12 | 90401 |
| 52200 | Rick | 1986-09-28 | 90401 |

# Operations on tabular data: Column Aggregation (Max)

*Who is the youngest person that lives in 90401?*

| Employee ID | Name | DOB | Zip Code |
|:---:|:---:|:---:|:---:|
| 20100 | Harry | 1980-11-12 | 90401 |
| 52200 | Rick | 1986-09-28 | 90401 |

*Who is the youngest person that lives in 90401?*

Max(**1980-11-12, 1986-09-28**) = **1986-09-28**

| Employee ID | Name | DOB | Zip Code |
|:-----------:|:----:|:----------:|:--------:|
| 20100 | Harry | 1980-11-12 | 90401 |
| 52200 | Rick | 1986-09-28 | 90401 |

# Operations on tabular data: Filter

*Who is the youngest person that lives in 90401?*

DOB = **1986-09-28**

| Employee ID | Name | DOB | Zip Code |
|---|---|---|---|
| 20100 | Harry | 1980-11-12 | 90401 |
| 52200 | Rick | 1986-09-28 | 90401 |

*Who is the youngest person that lives in 90401?*

```
DOB = 1986-09-28
```

| Employee ID | Name | DOB | Zip Code |
|:---:|:---:|:---:|:---:|
| 52200 | Rick | 1986-09-28 | 90401 |

# Operations on tabular data: Select a column

*Who is the youngest person that lives in 90401?*

| Employee ID | Name | DOB | Zip Code |
|---|---|---|---|
| 52200 | Rick | 1986-09-28 | 90401 |

# Operations Recap

- Filter *(keep only rows with Zip Code = 90401)*

- Aggregate values in a column *(get the maximum DOB)*

- Filter *(keep only rows with a matching DOB)*

- Select *(the name field from the remaining row)*

We talked about quite a few **data operations**

# Data Operations

- **Calculate new columns**

- **Aggregate columns**

- **Select columns**

- **Count rows**

- **Group rows**

- **Filter rows**

- **Visualize**

# POP QUIZ

What day had the most **Requests Closed**?

07-Aug (14 *Full Release* closed)

What is the proportion of **Activities** marked for *Full Release* vs. *No Records Exist*?

**524 (Full Release)**
**74 (No Records Exist)**

What is the **Avg. Days to Close** per month?

Mar   3.9        Jun   6.1
Apr   6.9        Jul   7.8
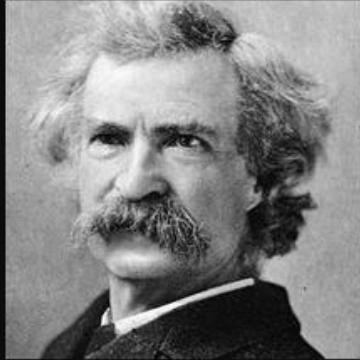May   6.6        Aug   7.3

# Wrapping up

# What do we *really* mean when we say data

- **Digital** (so we can use software tools)

- As **Raw** As Possible

- As **Structured** As Possible
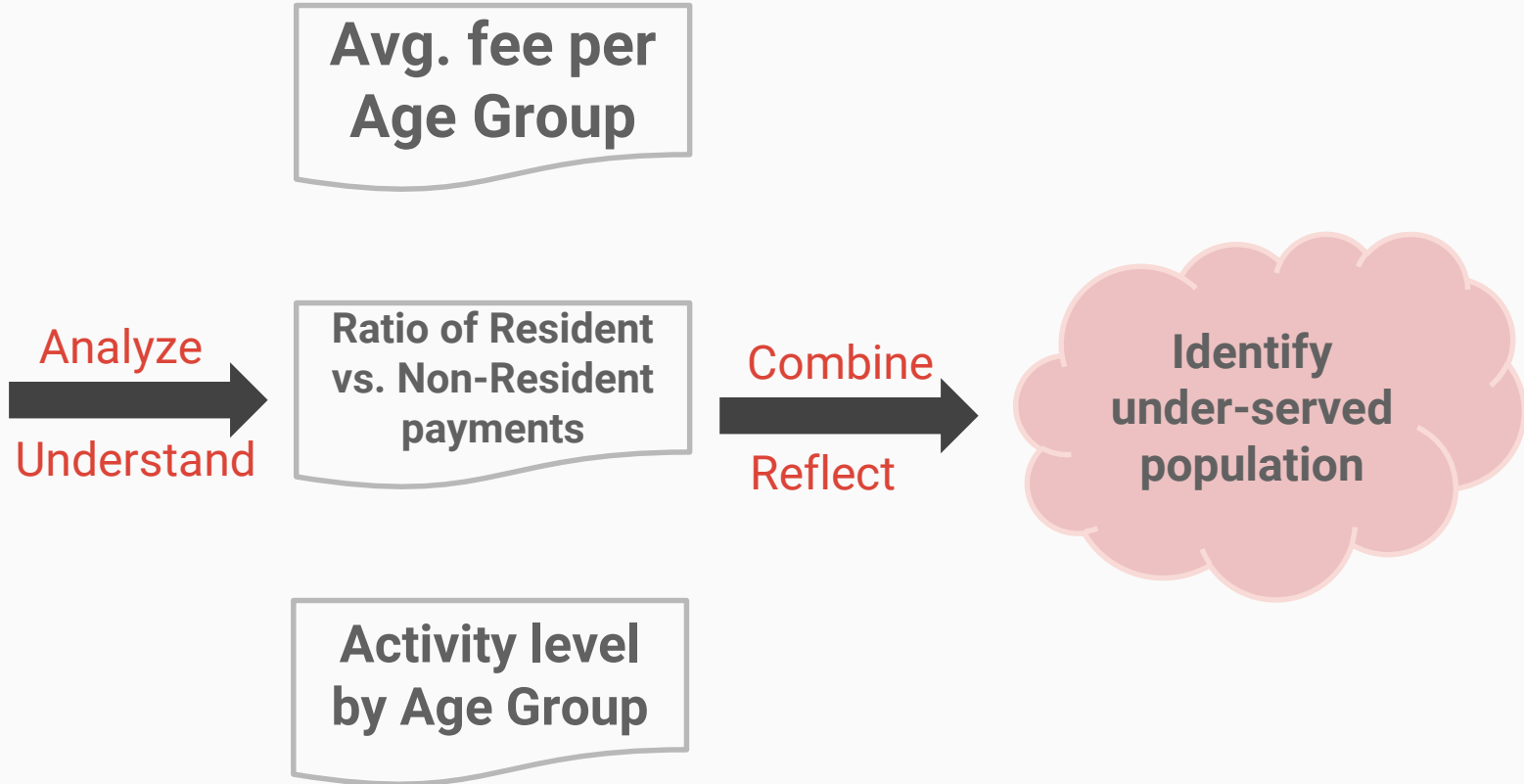
# Data source and structure matter!



Data is like garbage. You'd better know what you are going to do with it before you collect it.

~ Mark Twain

AZ QUOTES

# Operations on tabular data

- **Calculate new columns**

- **Aggregate values in a column**

- **Select columns**

- **Count rows**

- **Filter rows**

- **Group rows**

The Bigger Picture: Data Driven Decision Making

# Thank You For Joining Us!

**Please** fill out the feedback form before leaving ☺

**Materials** for today's course:
**santamonica.gov/DA101**

Questions, feedback anytime: **data@smgov.net**