# Descriptive Analytics
## Exploring and Visualizing Data

Veterans Analytics Course

September 16-17, 2020

Provided by: CANA Advisors

# Case Study
## National Parks

► National park visitors by year

► Data fields:
  ➢ Year
  ➢ National park name
  ➢ Region
  ➢ State
  ➢ Visitors
  ➢ Visit Type

# Data: structured vs. unstructured

► **Structured**
  ➢ Lists, data frames, spreadsheets, databases, 'big' data
  ➢ May contain -
    ▪ Numeric values
    ▪ Logicals (True/False)
    ▪ Factors
    ▪ Strings with set format

► **Unstructured**
  ➢ Multiple formats (no rigid structure)
  ➢ May contain -
    ▪ Images
    ▪ Free Text
    ▪ Speech
    ▪ Others

# Data Attributes

► **Accuracy** – How correct is the data?

► **Confidence** – what are the 'error bars' around the data provided?

► **Authority** – how authoritative is the source of the data?

*Question - Is it better to have Accurate or Authoritative data?*

# Tidy Data

Optimally organized data

# **Tidy** data

▶ Each variable must have its own column

▶ Each observation must have its own row

▶ Each value must have its own cell

## *For Excel, this means the first row is the header / list of field names and each row underneath is a record / observation*

▶ No blank rows

▶ No merged cells

# This data is tidy

► Each variable has its own column

► Each observation has its own row

► Each value has its own cell

| | A | B | C | D |
|---|---|---|---|---|
| 1 | ParkName | Rank | Value | PercentOfTotal |
| 2 | Grand Canyon NP | 1 | 320,032 | 16.42% |
| 3 | Lake Mead NRA | 2 | 219,510 | 11.26% |
| 4 | Yosemite NP | 3 | 155,578 | 7.98% |
| 5 | Olympic NP | 4 | 96,864 | 4.97% |
| 6 | Great Smoky Mour | 5 | 94,569 | 4.85% |
| 7 | Glen Canyon NRA | 6 | 89,782 | 4.61% |
| 8 | Canyonlands NP | 7 | 77,695 | 3.99% |
| 9 | Saint Croix NSR | 8 | 68,164 | 3.50% |
| 10 | Mount Rainier NP | 9 | 47,703 | 2.45% |
| 11 | Voyageurs NP | 10 | 45,262 | 2.32% |
| 12 | Rocky Mountain N | 11 | 41,213 | 2.11% |
| 13 | Yellowstone NP | 12 | 38,236 | 1.96% |

# Tidy or untidy?

| religion | <$10k | $10–20k | $20–30k | $30–40k | $40–50k | $50–75k |
|---|---|---|---|---|---|---|
| Agnostic | 27 | 34 | 60 | 81 | 76 | 137 |
| Atheist | 12 | 27 | 37 | 52 | 35 | 70 |
| Buddhist | 27 | 21 | 30 | 34 | 33 | 58 |
| Catholic | 418 | 617 | 732 | 670 | 638 | 1116 |
| Don't know/refused | 15 | 14 | 15 | 11 | 10 | 35 |
| Evangelical Prot | 575 | 869 | 1064 | 982 | 881 | 1486 |
| Hindu | 1 | 9 | 7 | 9 | 11 | 34 |
| Historically Black Prot | 228 | 244 | 236 | 238 | 197 | 223 |
| Jehovah's Witness | 20 | 27 | 24 | 24 | 21 | 30 |
| Jewish | 19 | 19 | 25 | 25 | 30 | 95 |

Table 4: The first ten rows of data on income and religion from the Pew Forum. Three columns, $75–100k, $100–150k and >150k, have been omitted.

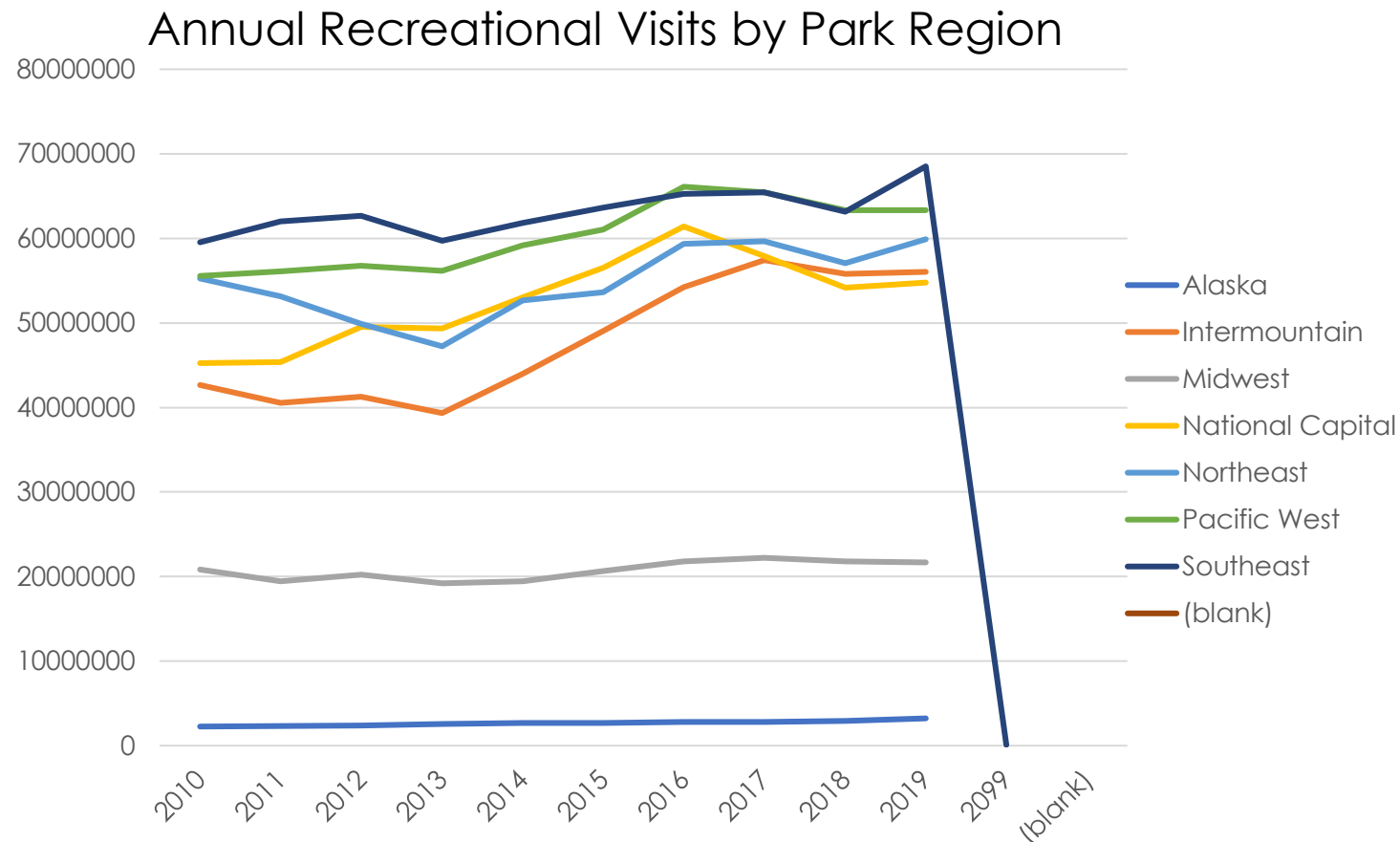Wickham, H. (2014). Tidy Data. Journal of Statistical Software, 59(10), 1 - 23. doi:http://dx.doi.org/10.18637/jss.v059.i10

# Common data issues

▶ Improperly stored data

▶ Duplicate records

▶ Outliers

▶ Missing data

▶ Invalid entries

# Garbage data, garbage visual
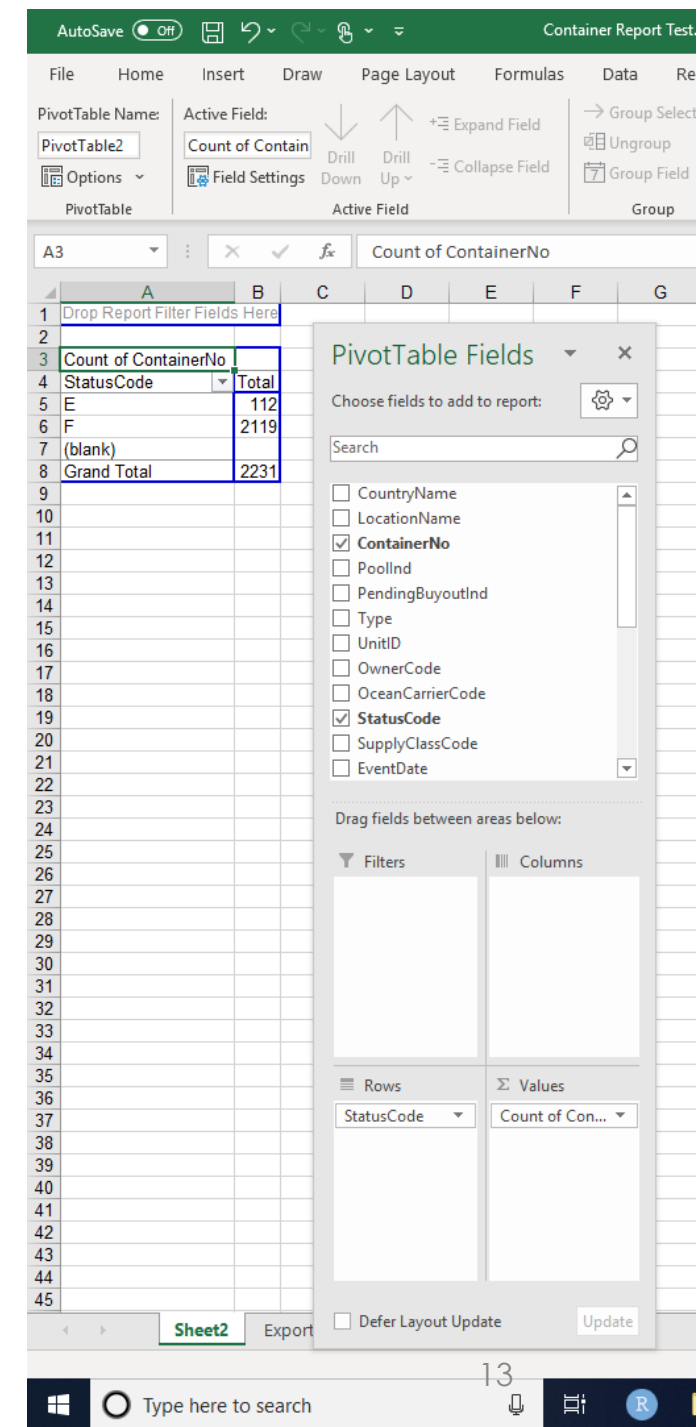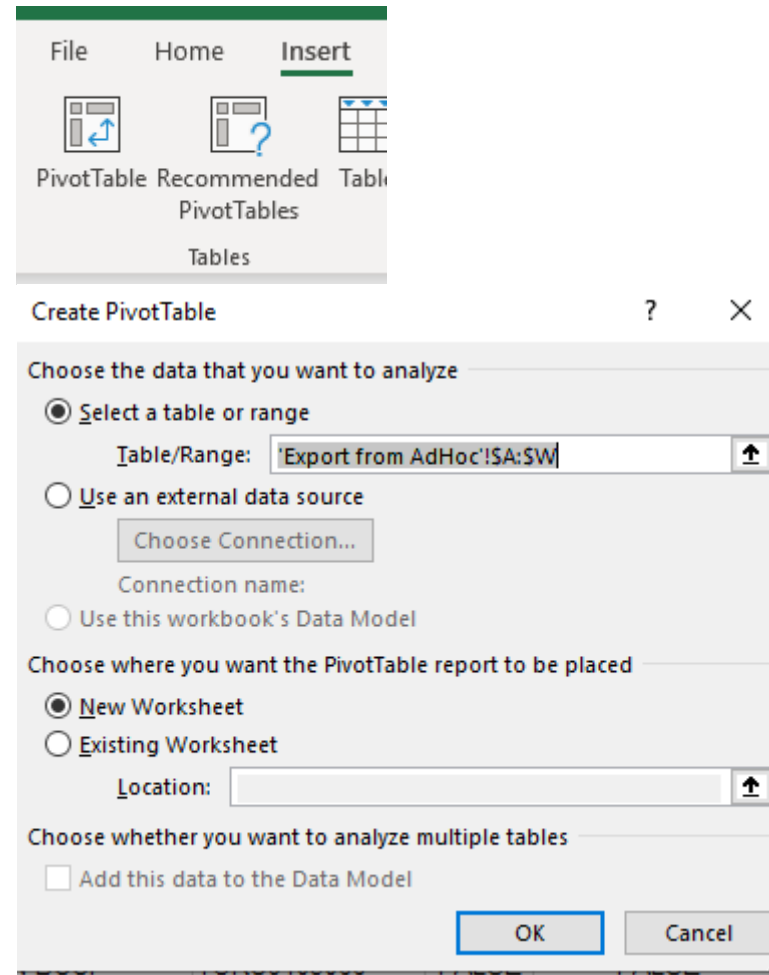
▶ What underlying issues impact this graph?

**Annual Recreational Visits by Park Region**



Legend:
- Alaska
- Intermountain
- Midwest
- National Capital
- Northeast
- Pacific West
- Southeast
- (blank)

# Summarizing Data

Slice and dice data with Pivot Tables

# **Functions** compute summary statistics

► SUM()

► MIN()

► MAX()

► MEDIAN()

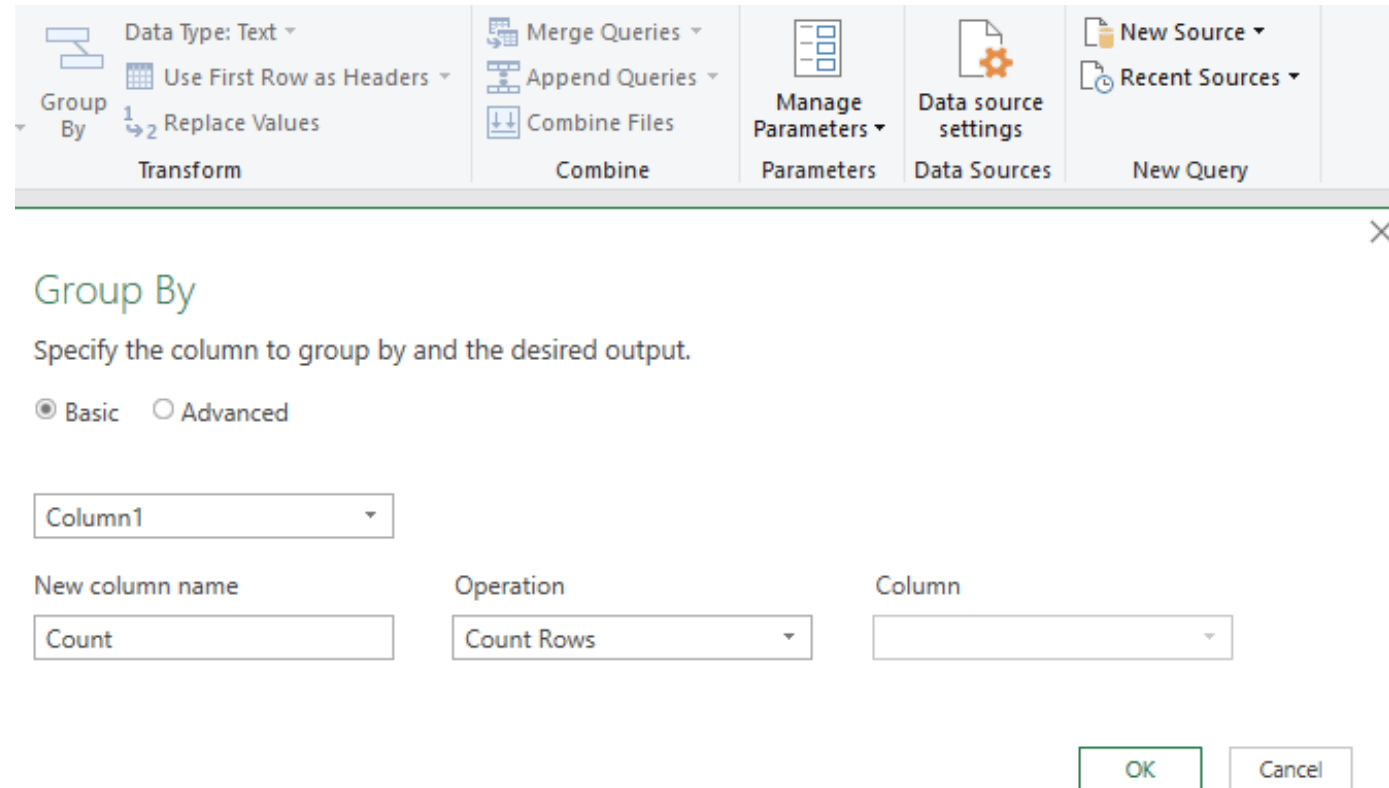► AVERAGE()

► COUNTA()

# PivotTables

- ► Allow you to quickly summarize data by groups

- ► Select Insert → PivotTable

- ► You typically want the default options:

  - ➤ Entire sheet as your range

  - ➤ PivotTable in a new window

# Power Query + PivotTables

▶ Power Query offers an intuitive way to construct PivotTables
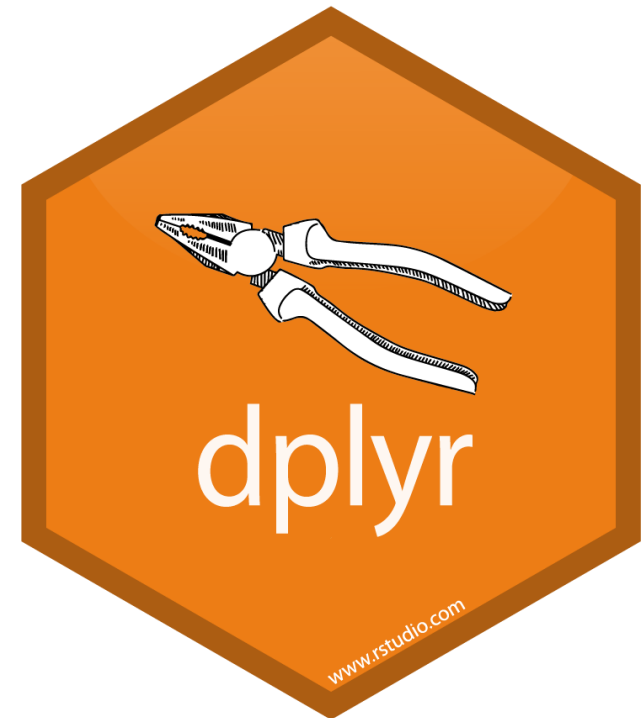
# **Practical** Examples

Let's see it in practice!

# Visualizing Data

Uncovering relationships in complex data
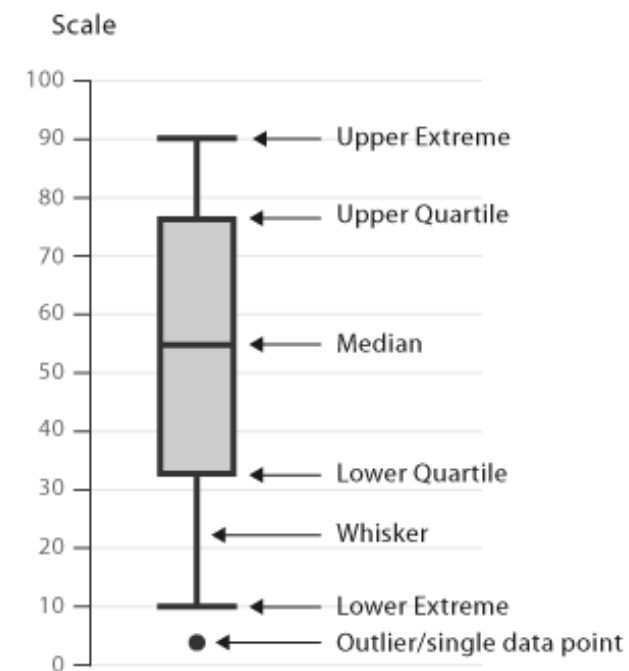
# Discovering relationships in data

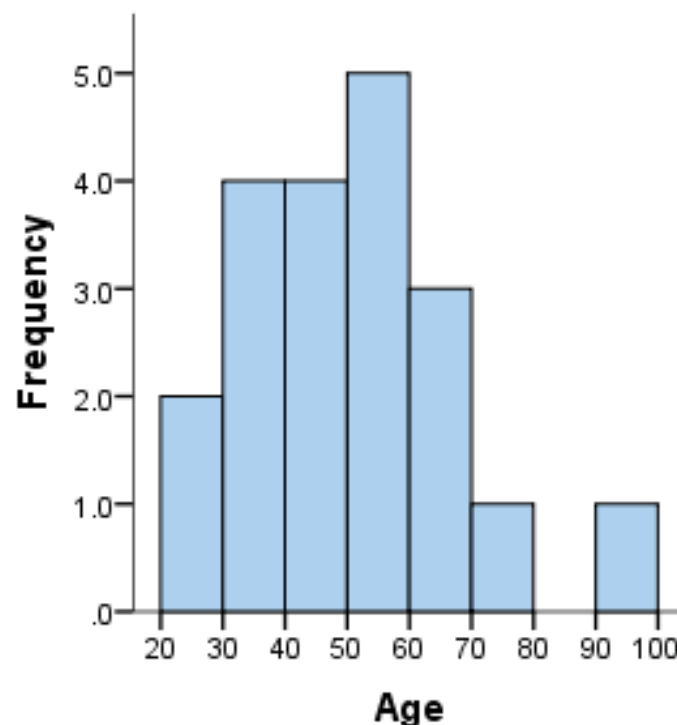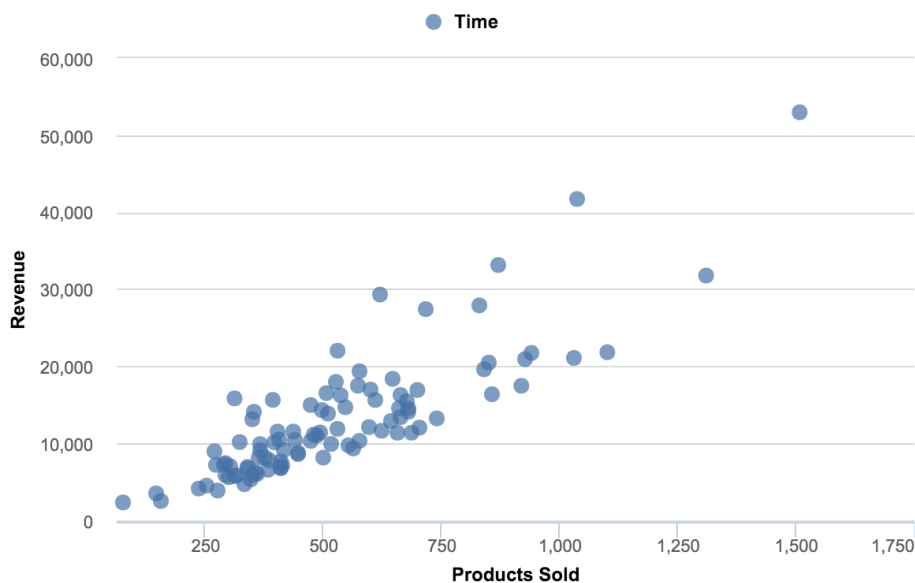► **Collecting and summarizing data**

➢ Basic statistics are helpful in summarizing data

➢ Box plots, scatter plots, box and whisker plots provide compact representations of how data is distributed

➢ Useful for exploring data but may not always be the best choice for communicating the data to your audience.

# Discovering relationships in data

▶ **PLOT** the data in a meaningful way

# Discovering relationships in data
A note on Machine Learning

## Supervised Learning
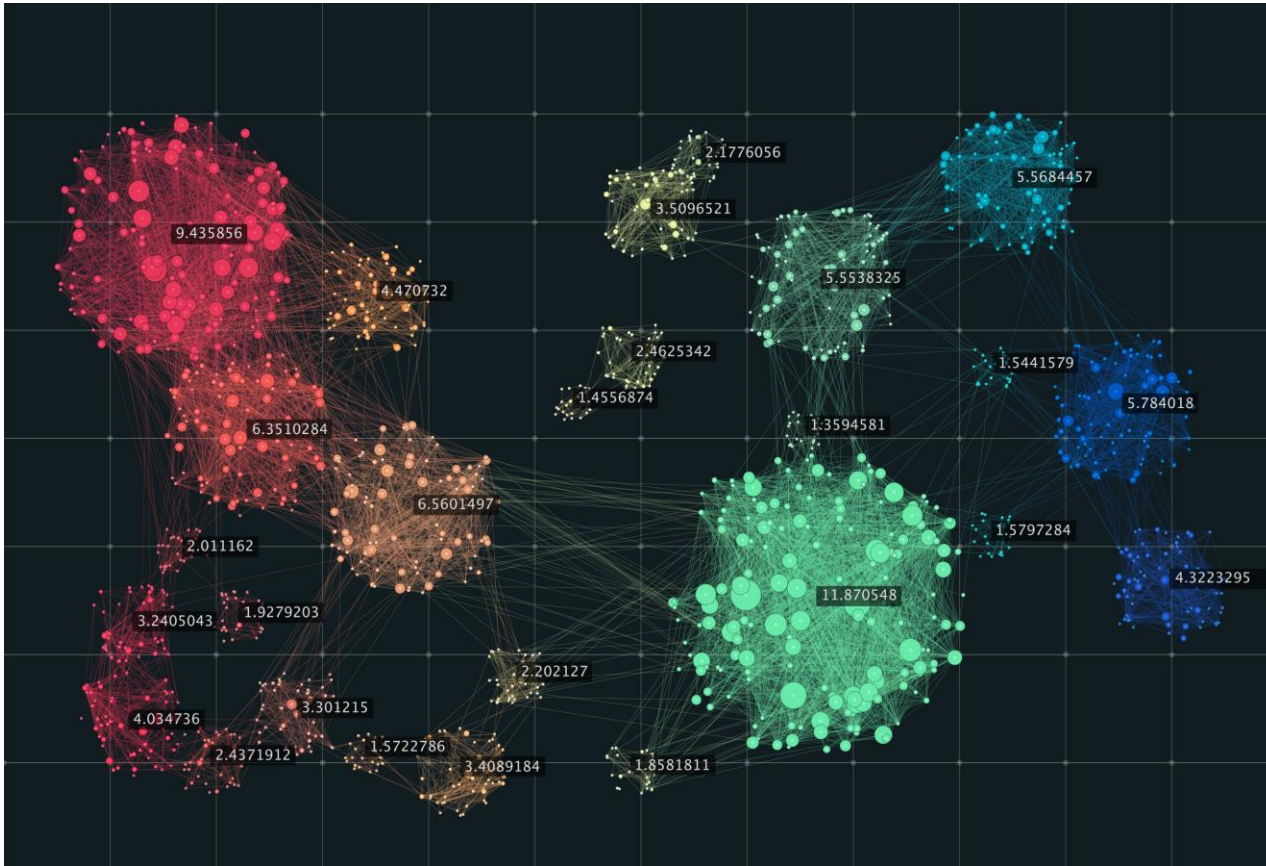
▶ Learn a function that approximates the relationship between input and output observable in the data

▶ **Labeled data**

▶ Examples:

▶ Classification

▶ Regression

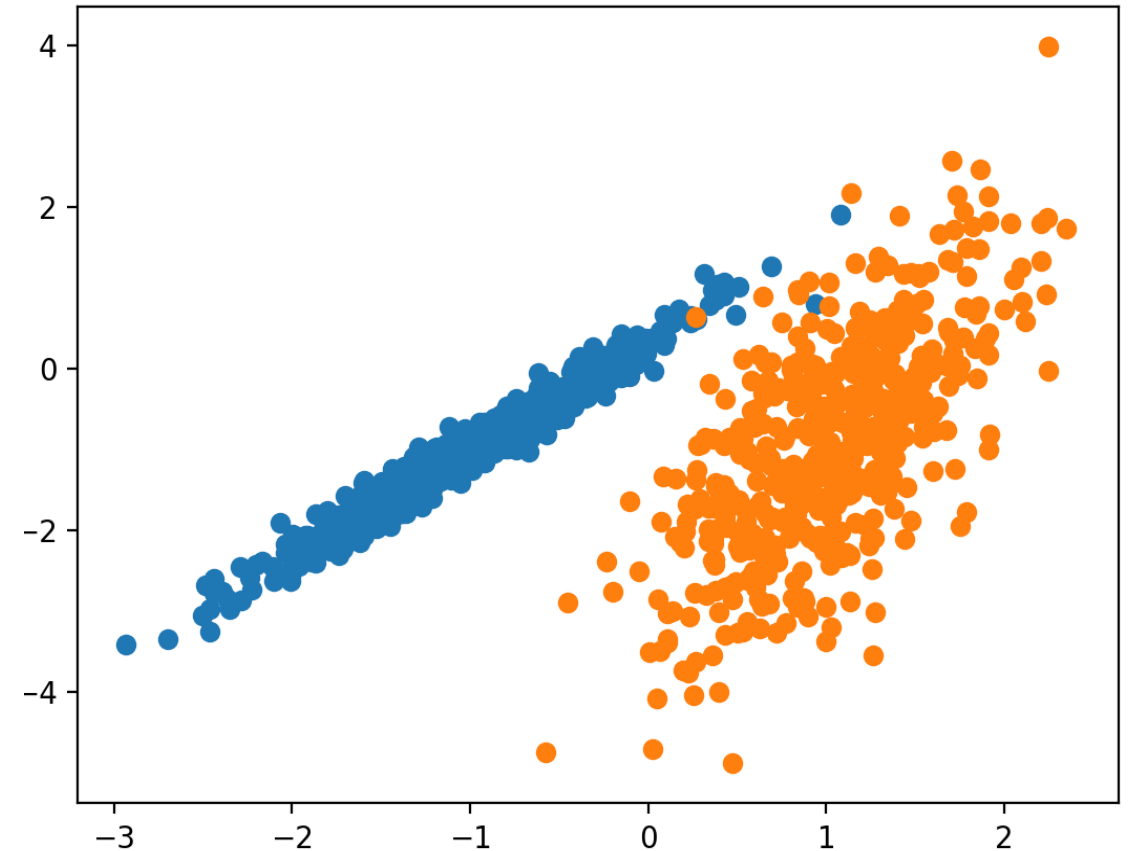## Unsupervised Learning

▶ Uncover the natural structure present within a set of data points

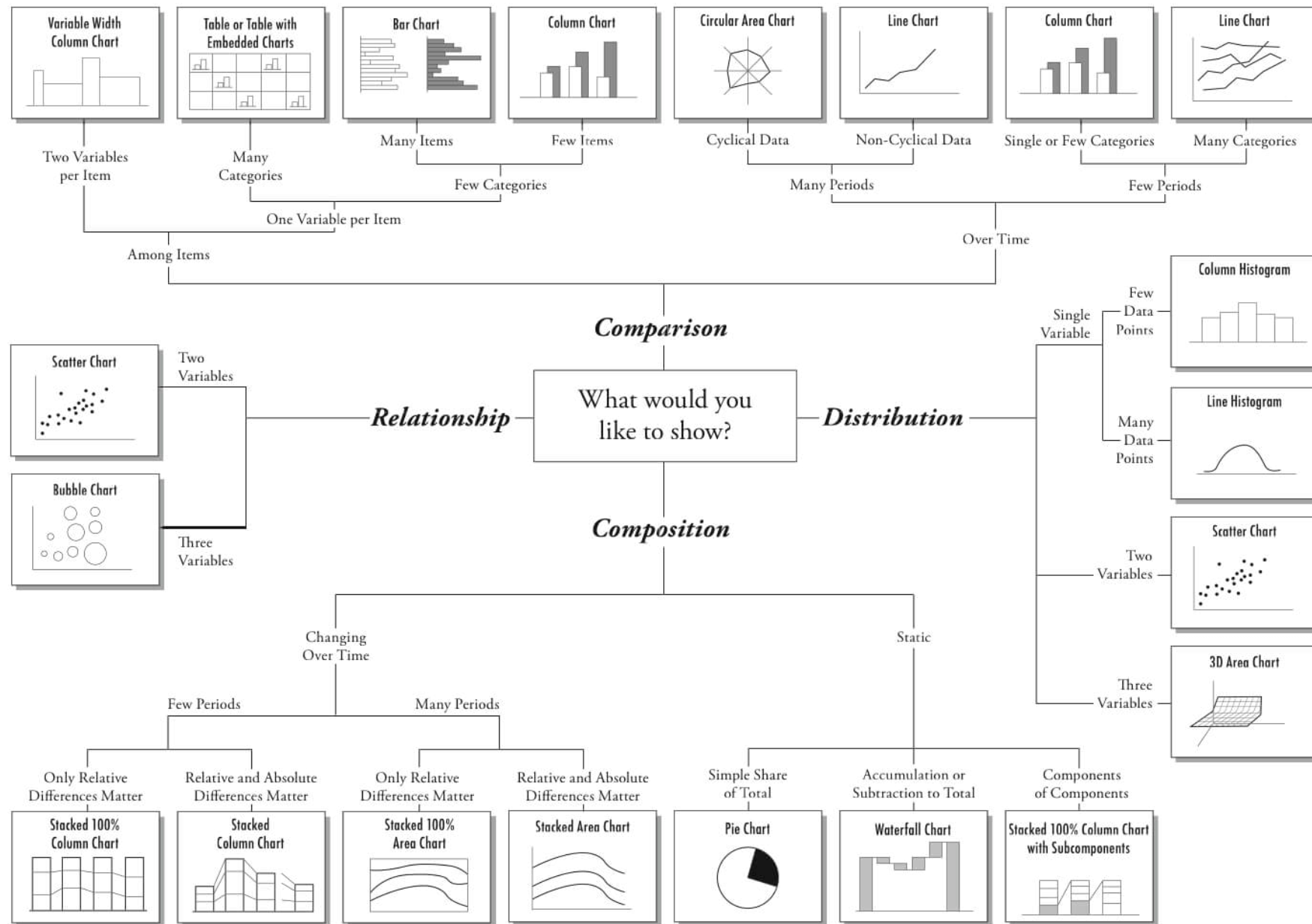▶ **Unlabeled data**

▶ Examples:

▶ Clustering

# Clustering

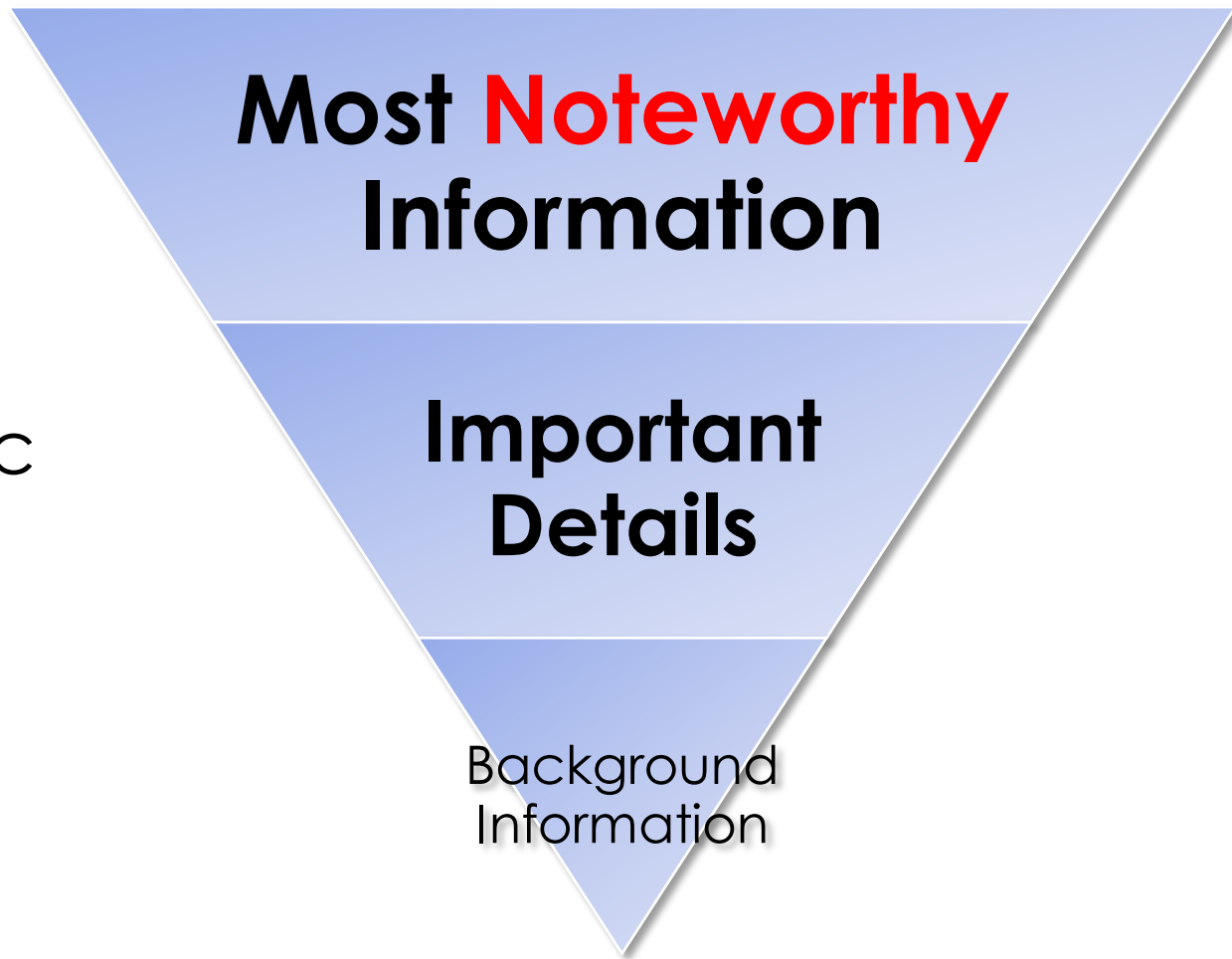# Storytelling with Data

Communicating meaning from complex data

# Chart Suggestions—A Thought-Starter

**Variable Width Column Chart**

**Table or Table with Embedded Charts**

**Bar Chart**

**Column Chart**

**Circular Area Chart**

**Line Chart**

**Column Chart**

**Line Chart**

Two Variables per Item

Many Categories

Many Items

Few Items

Cyclical Data

Non-Cyclical Data

Single or Few Categories

Many Categories

Few Categories

Many Periods

Few Periods

One Variable per Item

Over Time

Among Items

**Comparison**

**Scatter Chart**

Two Variables

**Relationship**

What would you like to show?

**Distribution**

Single Variable

Few Data Points

**Column Histogram**

Many Data Points

**Line Histogram**

**Bubble Chart**

Three Variables

**Composition**

Two Variables

**Scatter Chart**

Changing Over Time

Static

Three Variables

**3D Area Chart**

Few Periods

Many Periods

Only Relative Differences Matter

Relative and Absolute Differences Matter

Only Relative Differences Matter

Relative and Absolute Differences Matter

Simple Share of Total

Accumulation or Subtraction to Total

Components of Components

**Stacked 100% Column Chart**

**Stacked Column Chart**

**Stacked 100% Area Chart**

**Stacked Area Chart**

**Pie Chart**

**Waterfall Chart**

**Stacked 100% Column Chart with Subcomponents**

# Dashboard Design

► Composed of **"elements"** (aka "cards")

► Often laid out on a grid.

► Each element has a specific purpose – tells its own story.

► Choose **size, color, and placement** of elements to draw attention to what is important.

**Most <span style="color:red">Noteworthy</span> Information**

**Important Details**

Background Information

# Types of **Elements**

▶Big Numbers
▶Charts / Graphs
▶Maps
▶Heat map
▶Slicers
▶Tables
▶Sparklines

# Dashboard Best Practices



Consider visual hierarchy

Consider your goal
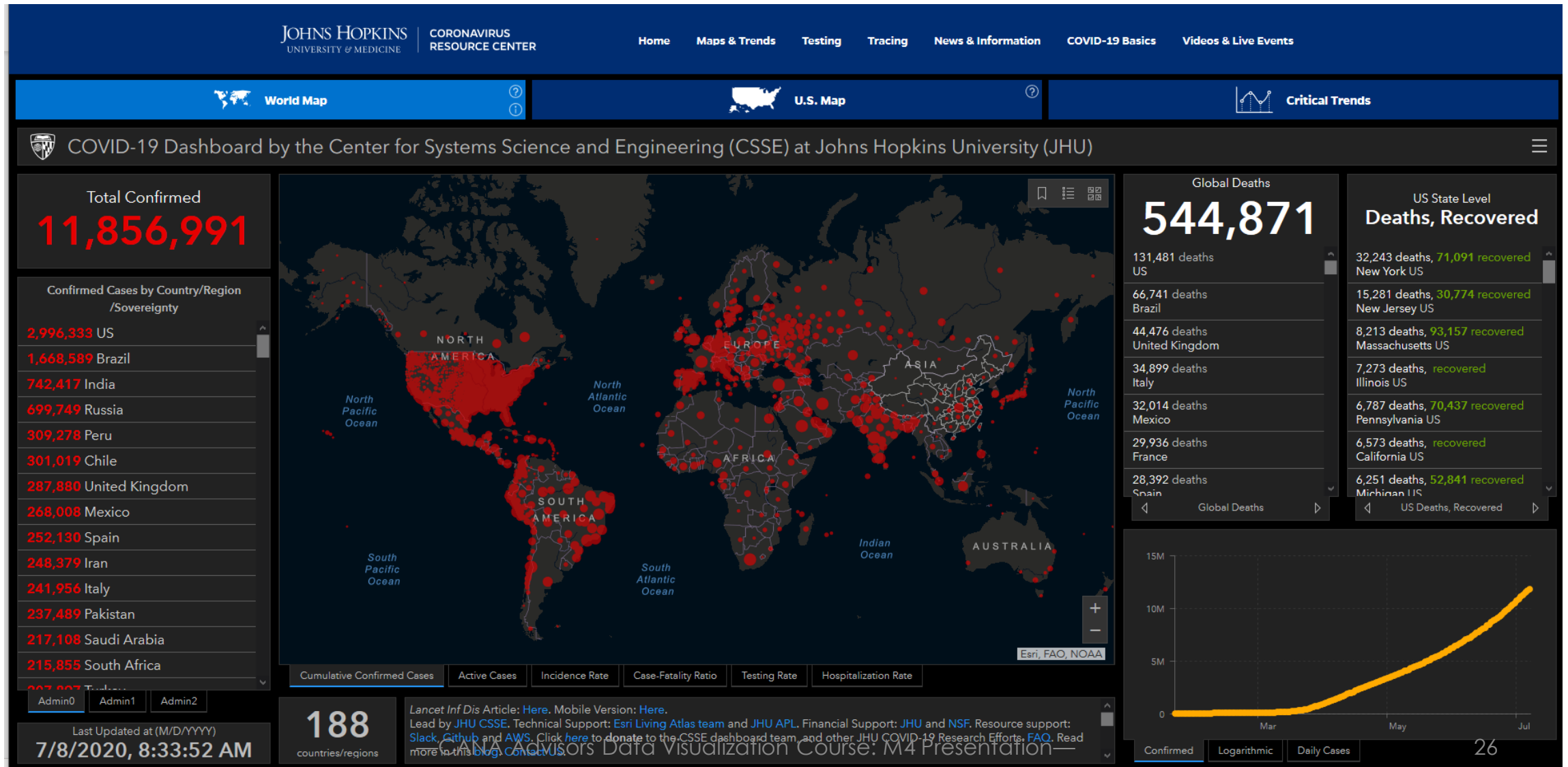
Provide Context

Keep it Simple

Provide user controls

Use the right visuals

Select relevant KPIs

Iterate & Improve

# Dashboard Critique

# **Practical** Examples

Let's see it in practice!