# Course Reminders

- Due Friday (11:59 PM)
  - D4
  - Q4
  - A2
  - Mid-quarter survey (*optional*)
  - Weekly Project Survey (*optional*)

- Scores posted on Canvas: Q3, D3, project review
- Project Proposal grading underway

# Project Reviews

- Really well done overall; thoughtful, clear, and detailed
- Strengths:
    - <u>Premise</u>: interesting questions
    - <u>Writing</u>: background information clear; visualizations, consistency across/throughout; analysis and results explained and interpreted clearly
    - <u>Code</u>: code well commented and explained
- Weaknesses:
    - <u>Flow</u>: code question & hypothesis don't match; typos; poor variable naming; avoid output of large tables/datasets; inconsistencies between explanations and code; mislabeled visualizations
    - <u>Writing</u>: poor flow, wording, clarity; poor communication; confusing as to what was done; too wordy/long-winded; text hard to read/follow; failure to define initialisms/acronyms; used too much jargon; don't overwhelm the reader - subsections; bullet points
    - <u>Data</u>: unclear exactly what was being measured
    - <u>Code</u>: too repetitive, cluttered
    - <u>Analysis</u>: didn't quite have the data they needed (often time-related); explanations lacking for analysis

# Approaches to Analysis
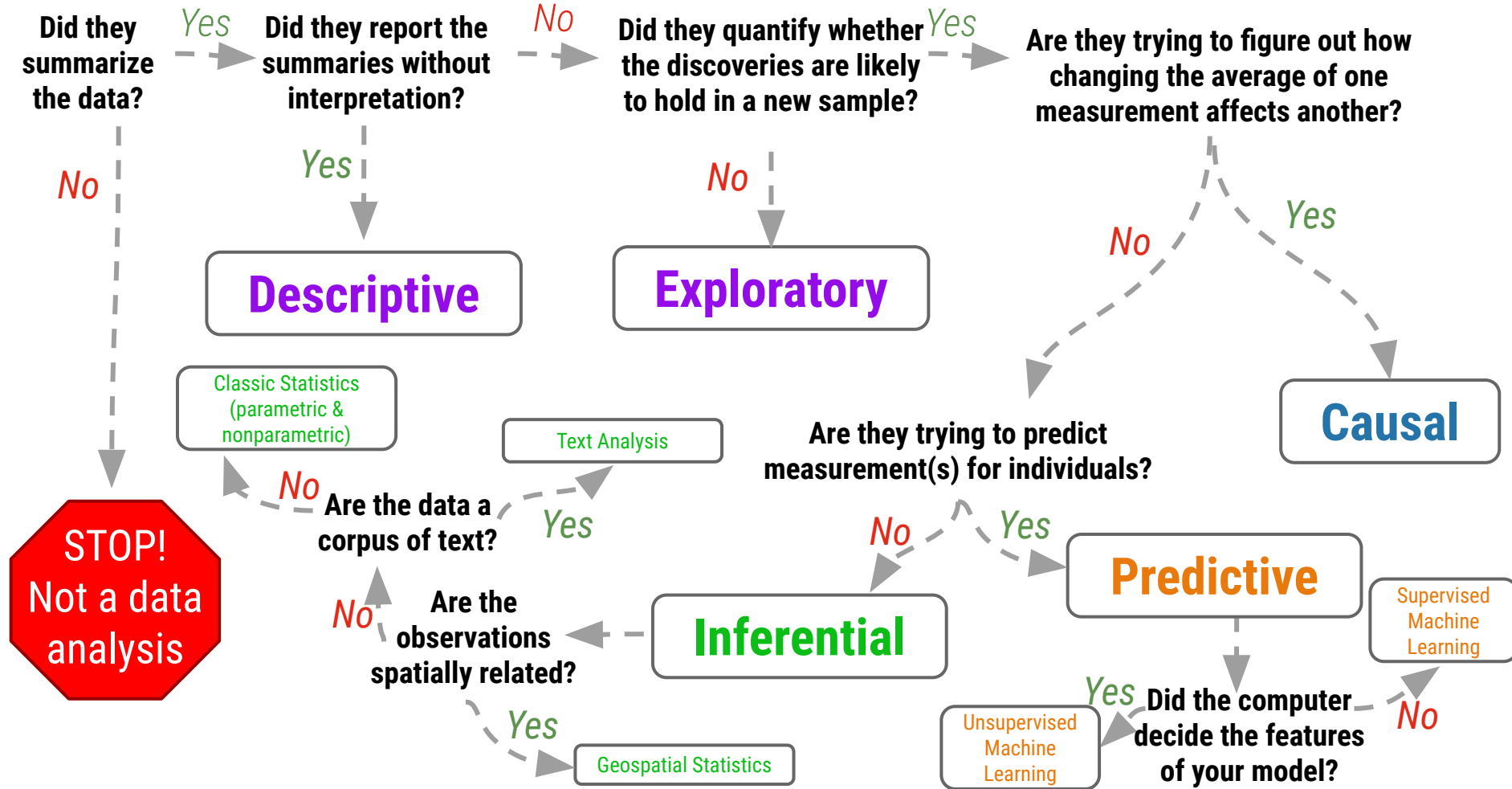
Shannon E. Ellis, Ph.D
UC San Diego

Department of Cognitive Science
sellis@ucsd.edu

*"Data science is the process of formulating a quantitative question that can be answered with data, collecting and cleaning the data, __analyzing the data__, and communicating the answer to the question to a relevant audience."*

To do this, you have to *look at, describe, and explore* the data

# Summary: Analytical Approaches

1. **Descriptive** (and **Exploratory**) Data Analysis are the first step(s)

2. **Inference** establishes relationships
   a. Classic Statistics
   b. Geospatial Analysis
   c. Text Analysis

3. Machine Learning is for **prediction**
   a. Supervised
   b. Unsupervised

4. Experiments best way to establish **causality**

# Exploring Analyses

General question: What impacts politics in America?

Data Science question: Is there a relationship between the sentiment of political words in South Park and America's presidential approval rating?
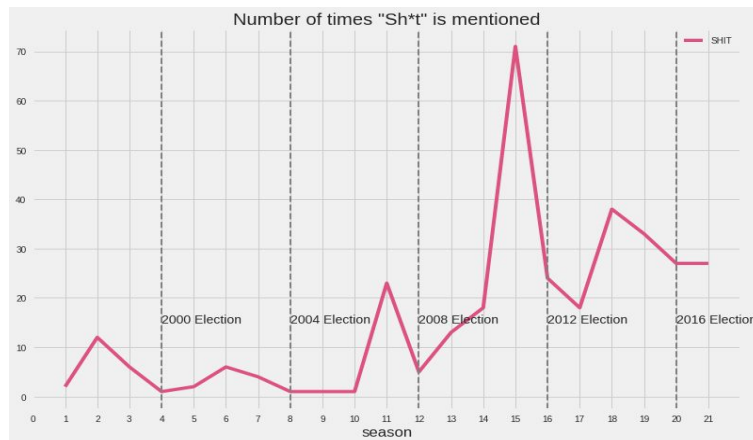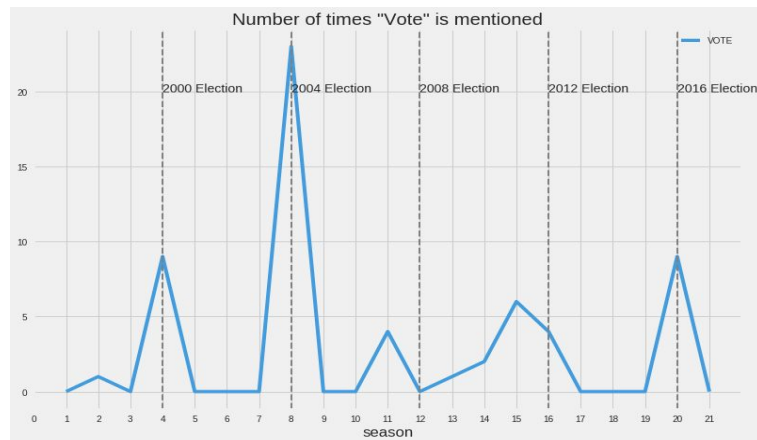
Descriptive

Exploratory

Inferential

Text Analysis

Classic Statistics
(parametric &
nonparametric)



Number of times "Vote" is mentioned

Number of times "Sh*t" is mentioned

General question: What gets too much attention in the news?

Data Science Question: Is there a relationship over time between cause of death terms in the *NYT*, The Guardian, and Google trends data relative to data from the CDC?
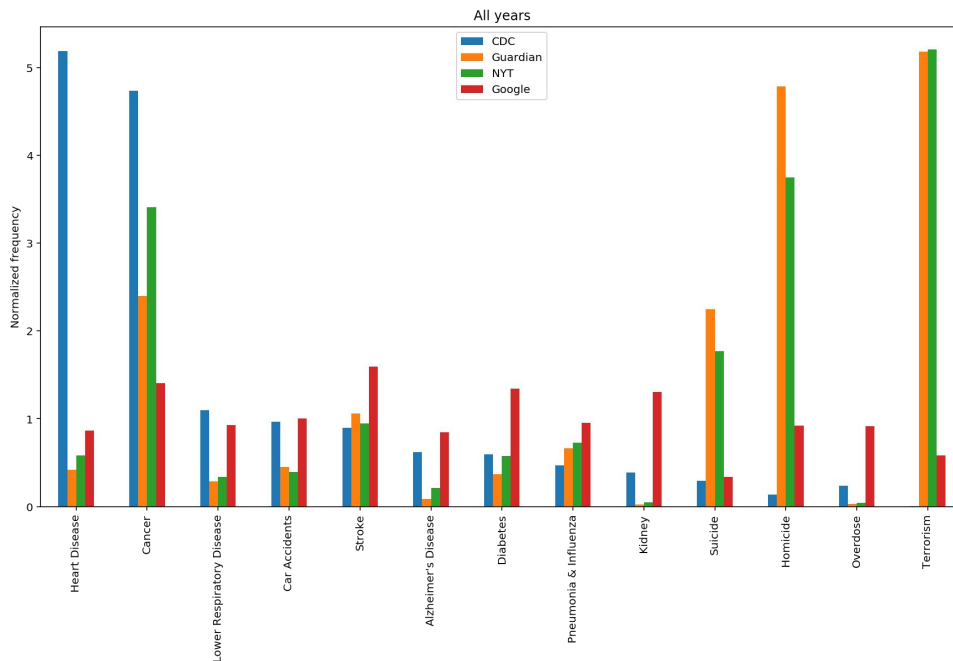


Descriptive

Exploratory

Inferential

Text Analysis

Classic Statistics
(parametric &
nonparametric)

<u>Classification</u>: Often we seek to assign a label to an item from a discrete set of possibilities.

*Can we predict who will win next year's NCAA tournament? The Masters? The Super Bowl? The pennant? A game?*

*Can we predict the genre of a given movie (comedy, drama, or animation?) from just its script?*

**Regression**: A way to forecast a given numerical quantity using other relevant features.

*Can we predict someone's weight given other information?*

*How much snow will the East Coast get this year?*

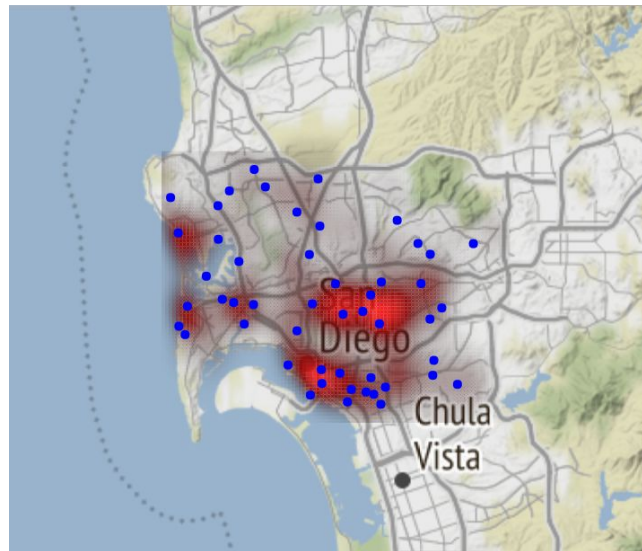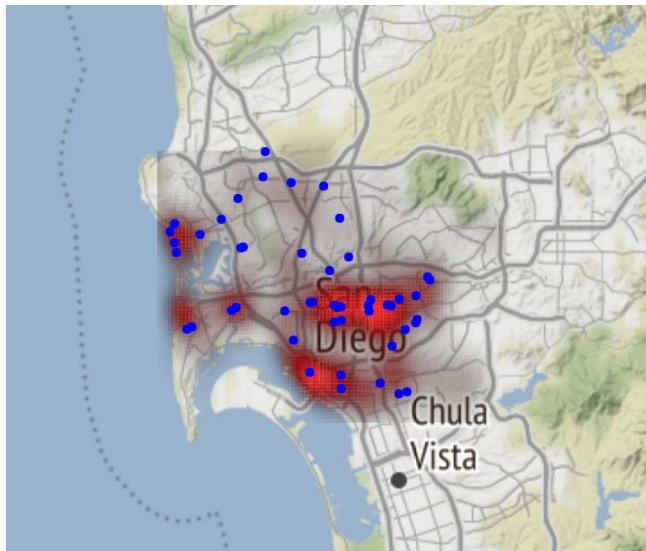General question: Why isn't police response time always the same?

Data Science question: Where should police cars be stationed, accounting for crime levels and time of day, to make police response times equitable throughout San Diego?

Descriptive

Exploratory

Predictive

Inferential

*In case of the total drought in California, how many desalination plant projects we need to supply residential use water for population who live in urban areas in California?*

Average Gallons Per Capita Day for each county in 2016