

Descriptive and Exploratory Analysis

Jason G. Fleischer, Ph.D.

Asst. Teaching Professor

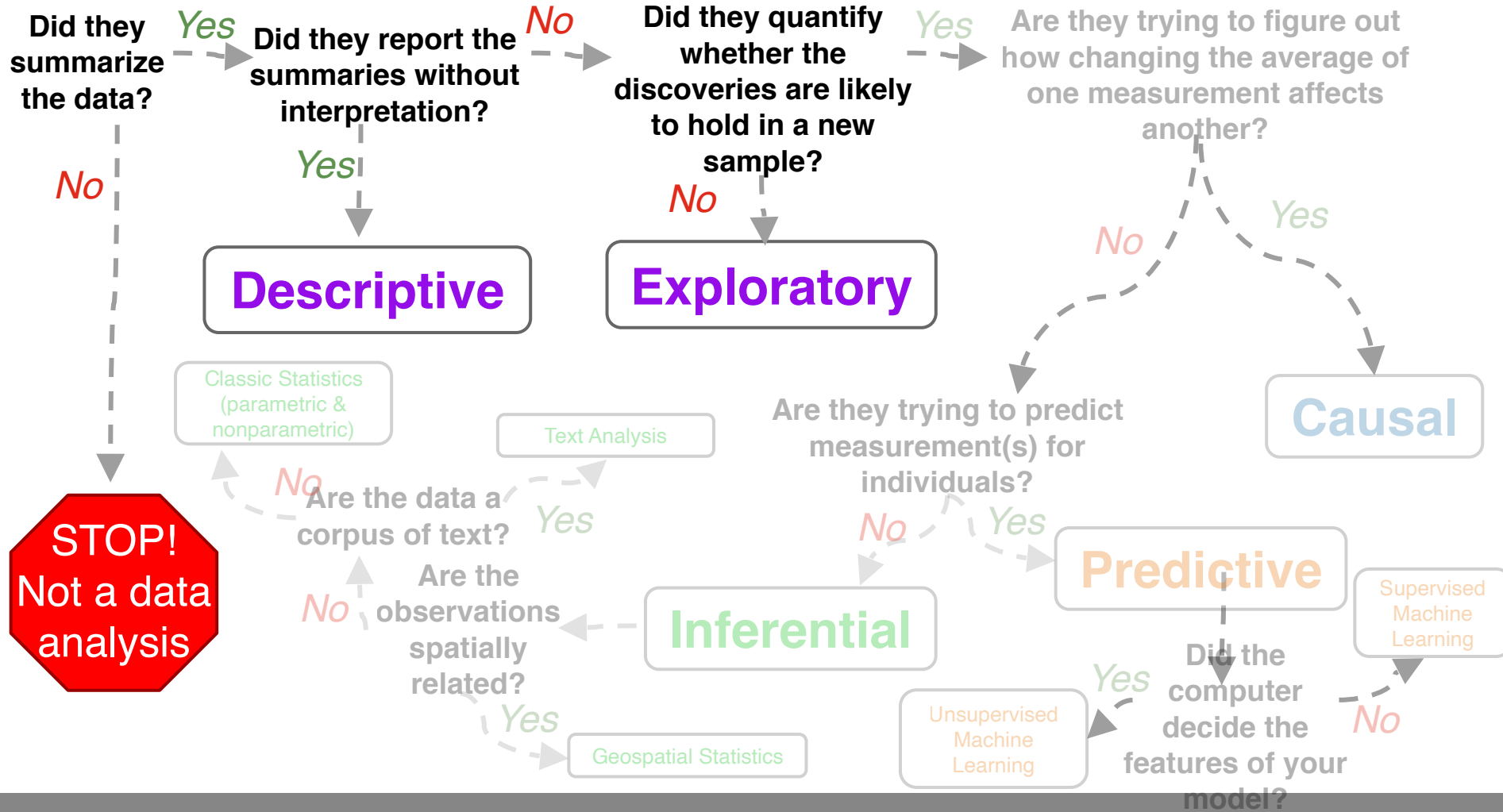
Department of Cognitive Science, UC San Diego

jfleischer@ucsd.edu



@jasongfleischer

<https://jgfleischer.com>



Descriptive: The goal of descriptive analysis is to understand the components of a data set, describe what they are, and explain that description to others who might want to understand the data.

- Problem: Understanding whether users are nice or mean on Youtube
- Data science question: Are the words that people use in their comments more frequently positive words (great, awesome, nice, useful) or negative words (bad, stupid, lame, awful)?
- Type of analysis: Descriptive analysis

To answer this you would calculate statistics about YouTube comments



Statistics

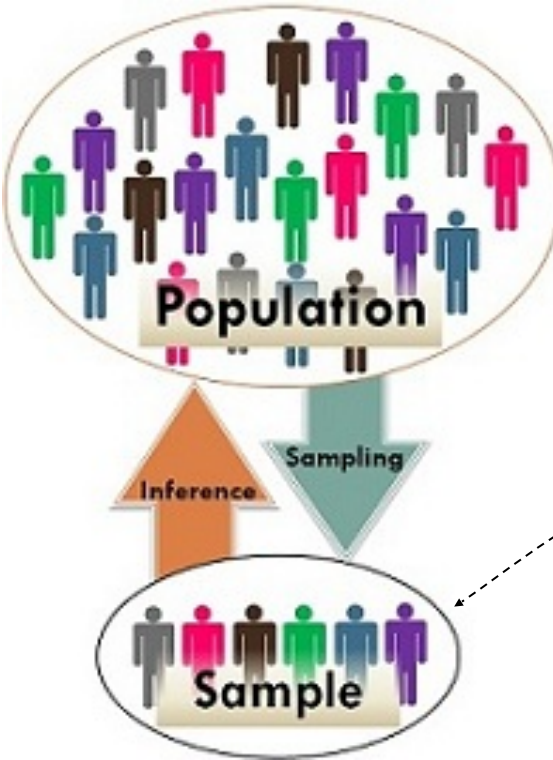
*“the science that deals with the **collection, classification, analysis, and interpretation of numerical facts or data**”*

statistic

“A quantity computed from a sample”

Populations & Samples

We want to learn something about this..



Our population: *all* YouTube comments

Our sample: 100,000 comments

....but we can only *actually* collect data from this

statistic

“A quantity computed from a sample”



For our YouTube analysis, we could take a random sample of comments from YouTube and calculate the following statistic: *the number of positive and the number of negative words in each review.*

Best sampling practices:

- Always think about what your population is
- Collect data from a sample that is representative of your population
- If you have no choice but to work with a dataset that is not collected randomly and is biased, be careful not to generalize your results to the entire population



You'd want to be sure you sample randomly across *all* YouTube comments, making sure not to get more comments from one genre over another, or one location over another, etc.

Examples of bad sampling:

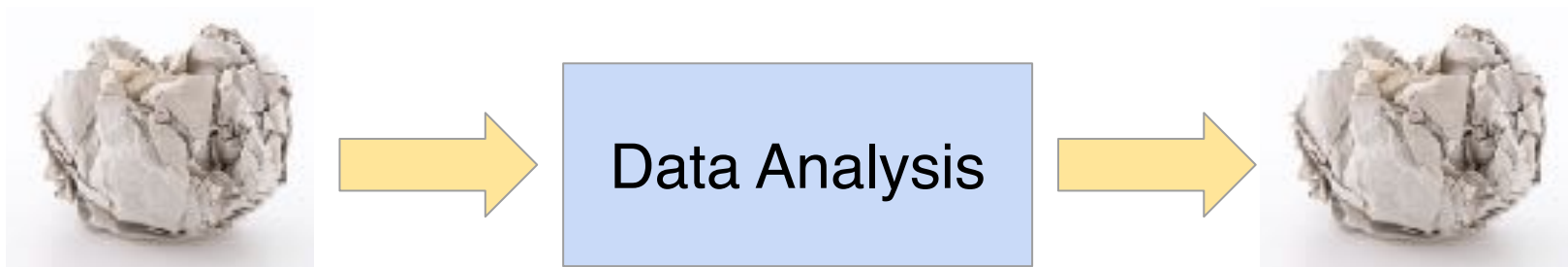
- Surveying subscribers of a gun-related magazine for research on Americans' attitudes toward owning guns
- Randomly sampling Facebook users for what TV shows people like



To understand *all* YouTube comments, you wouldn't just want to sample from one YouTube channel, or videos in a single language.


It's *always* worth spending time at the beginning of a project to determine whether or not the data you have are garbage. Be certain they are actually able to help you answer the question you're interested in.

GIGO : Garbage In. Garbage Out.



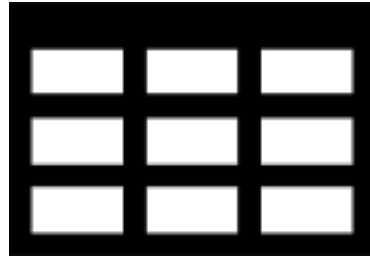


For the survey data I collected from you all, which of the following best describes the population I could generalize findings back to.

- 
- A** Undergraduates
 - B** Undergraduates in the US
 - C** Undergraduates at UCSD
 - D** Students aged 18-25
 - E** UCSD COGS108 students

Descriptive

Descriptive Analysis



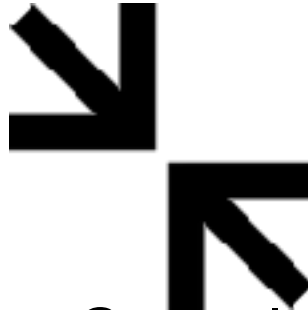
Size



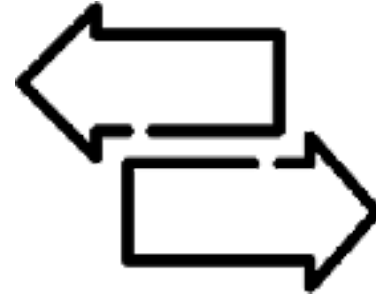
Missingness



Shape



Central
Tendency



Variability



Size

How many observations (rows) and variables (columns) you have is an important first step. You should always be aware of the **size** of your dataset



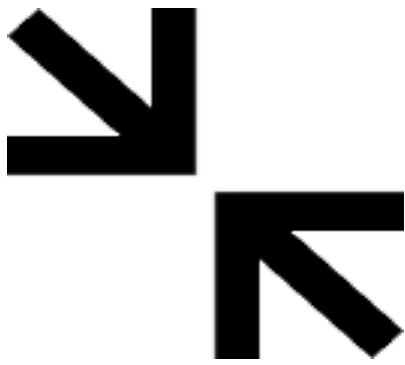
Descriptive

Missingness It's critical to know how many observations have missing data for variables of interest in your data. Knowing *why* their missing is also important.



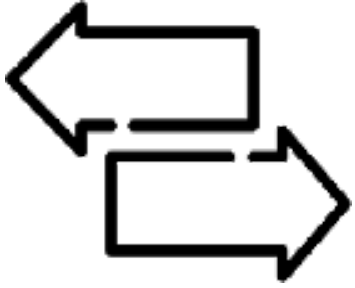
Shape

It's critical to know the distribution of the variables in your dataset. Certain statistical approaches can only be used with certain distributions.



Central Tendency

Knowing the mean, median, and/or mode can help you get an idea of what a typical value is for your variable(s) of interest



Variability

The central tendency tells you part of the story. The **variability in the values** in your observation helps fill in the rest.



Which of the following is NOT something accomplished by a descriptive analysis?

A Describes typical values in your dataset

B Determines the size of your dataset

C Establishes causal relationships between variables

D Identifies missing data

E Determines how variable values in your dataset are

Descriptive Statistics & Summary

“We must suppress some of the truth to communicate the truth... In short, the techniques of descriptive statistics are designed to match the salient features of the data set to human cognitive abilities.”

-I.J. Good (1983)

Descriptive
Analyses are
often included as
“Table 1” in
academic
publications

Table 1. Baseline Characteristics of the Patients.*

Characteristic	Cardiovascular Mortality (N = 305)	Neurovascular Mortality (N = 358)	Ischemic Stroke Mortality (N = 150)	Overall Mortality as Assessed (N = 150)
Age — no. (%)				
50–59 yr	2 (0.7)	1 (0.3)	6 (4.0)	2 (0.7)
60–69 yr	12 (4.0)	28 (8.3)	61 (40.6)	84 (55.3)
70–79 yr	302 (99.3)	329 (93.4)	115 (75.4)	123 (81.3)
80–89 yr	142 (46.2)	118 (33.4)	126 (82.3)	142 (93.7)
≥90 yr	22 (7.1)	23 (6.6)	28 (18.7)	18 (11.9)
Mean ± SD	79.2±7.4	80.1±7.3	78.4±7.5	79.3±7.4
Sex — no. (%)				
Female	163 (53.5)	138 (38.5)	135 (89.3)	184 (121.3)
Male	115 (37.5)	120 (33.7)	113 (73.5)	116 (75.7)
Race — no. (%)†				
White	207 (67.7)	235 (65.9)	128 (85.3)	134 (88.0)
Other	4 (1.3)	5 (1.4)	2 (1.3)	6 (3.9)
History of myocardial infarction — no. (%)	18 (5.9)	12 (3.4)	42 (28.0)	58 (38.0)
History of stroke — no. (%)	16 (5.2)	15 (4.2)	22 (14.7)	18 (11.9)
History of transient ischemic attack — no. (%)	12 (4.0)	15 (4.2)	12 (8.0)	15 (9.9)
Eligible patients — no. (%)				
Spatial	114±18	115±19	126±17	126±17
Duration	75±15	73±12	76±8	75±10
Visual field score and spatial quotient				
60–69 letters, 28/31–40 — no. (%)	111 (36.4)	94 (26.3)	116 (77.3)	133 (87.3)
51–67 letters, 28/39–40 — no. (%)	58 (18.7)	118 (33.2)	108 (71.3)	115 (75.0)
38–50 letters, 28/108–180 — no. (%)	67 (21.6)	13 (3.6)	58 (38.7)	58 (38.0)
21–37 letters, 28/108–180 — no. (%)	25 (8.1)	21 (5.9)	16 (10.7)	26 (17.0)
Mean score	60.1±14.7	58.3±15.1	61.5±11.2	60.4±11.4
Total cholesterol if fasting — mg/dL	418±114	483±136	458±155	461±115
Estimated cholesterol plus calculated total cholesterol at 10 min — mg/dL	211±121	254±132	247±112	252±115
Formal cardiac and vascular — no. (%)				
Coronary atherosclerosis	175 (57.4)	179 (50.3)	176 (117.3)	180 (117.0)
Stroke	85 (27.9)	81 (22.9)	77 (51.3)	72 (47.0)
Hemorrhage	19 (6.2)	24 (6.7)	24 (15.9)	25 (16.3)
Other	15 (4.9)	23 (6.4)	15 (10.0)	18 (11.7)
No coronary atherosclerosis or stroke possible to give	2 (0.7)	8 (2.2)	6 (4.0)	2 (0.7)

* Plus-minus values are means ± SD.
† Race was self-reported.
‡ Total cholesterol at 10 min includes the total, calculated total, cholesterol measured at 10 min, and total protein of apolipoprotein B.

Descriptive

Size

Zooming in on this we see variables stratified by Age, Sex, and Race

Table 1. Baseline Characteristics of the Patients.*

Characteristic	Ranibizumab Monthly (N = 301)	Bevacizumab Monthly (N = 286)	Ranibizumab as Needed (N = 258)	Bevacizumab as Needed (N = 300)
Age — no. (%)				
50–59 yr	2 (0.7)	1 (0.3)	6 (2.0)	2 (0.7)
60–69 yr	33 (11.0)	28 (9.8)	31 (10.4)	34 (11.3)
70–79 yr	102 (33.9)	84 (29.4)	115 (38.6)	103 (34.3)
80–89 yr	142 (47.2)	150 (52.4)	126 (42.3)	142 (47.3)
≥90 yr	22 (7.3)	23 (8.0)	20 (6.7)	19 (6.3)
Mean — yr	79.2 ± 7.4	80.1 ± 7.3	79.4 ± 7.8	79.3 ± 7.6
Sex — no. (%)				
Female	183 (60.8)	180 (62.9)	185 (62.1)	184 (61.3)
Male	118 (39.2)	106 (37.1)	113 (37.9)	116 (38.7)
Race — no. (%)†				
White	297 (98.7)	281 (98.3)	296 (99.3)	294 (98.0)
Other	4 (1.3)	5 (1.7)	2 (0.7)	6 (2.0)

* Plus-minus values are means ± SD.

† Race was self-reported.

‡ Total thickness at the fovea includes the retina, subretinal fluid, choroidal neovascularization, and retinal pigment epithelial elevation.

Shape

Central tendency

variability

Descriptive Statistics & Summary

Calculating descriptive statistics, understanding what they tell you about your data, and reporting them are critical steps in every analysis.

Exploratory: The goal is to find unknown relationships between the variables you have measured in your data set. Exploratory analysis is open ended and designed to verify expected or find unexpected relationships between measurements.

Exploratory



Exploratory Data Analysis (EDA)
detective work answering the question:
“What can the data tell us?”

Why EDA?

- Understand data properties
- Discover Patterns
- Generate & Frame Hypothesis
- Suggest modeling strategies
- Check assumptions (sanity checks)
- Communicate results (present the data)

.....and if you don't, you'll regret it

The
dataset

You

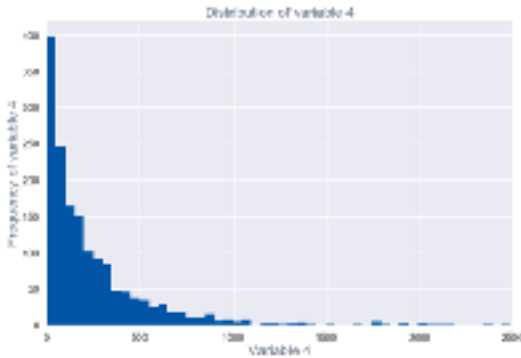


The general principles of exploratory analysis:

- Look for missing values
- Look for outlier values
- Calculate numerical summaries
- Generate plots to explore relationships
- Use tables to explore relationships
- If necessary, transform variables

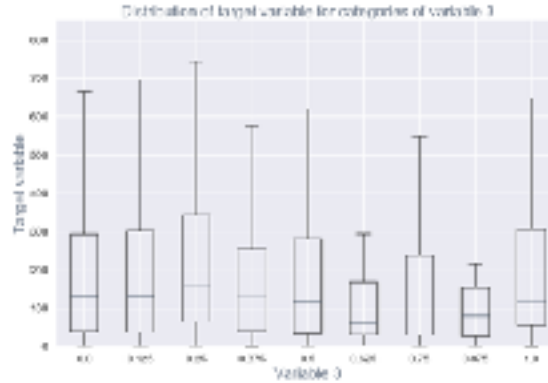
EDA Approaches to “Get a Feel for the Data”

Understanding the relationship between variables in your dataset



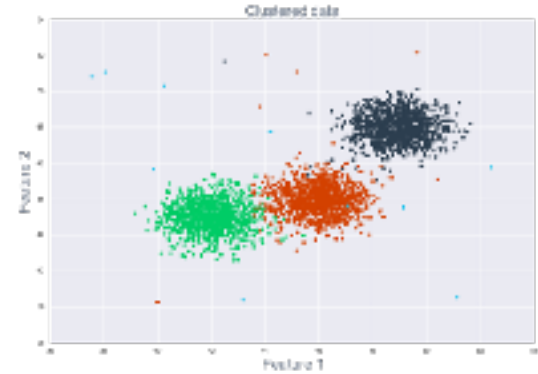
Univariate

understanding a single variable
i.e.: histogram, densityplot, barplot



Bivariate

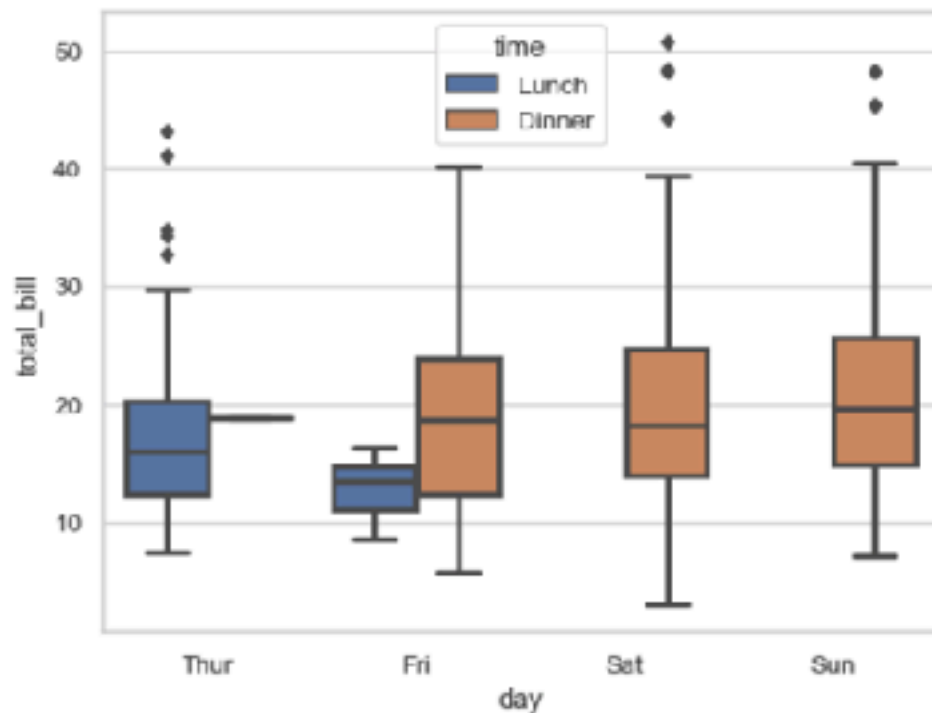
understanding relationship between 2 variables
i.e.: boxplot, scatterplot, grouped barplot, boxplot



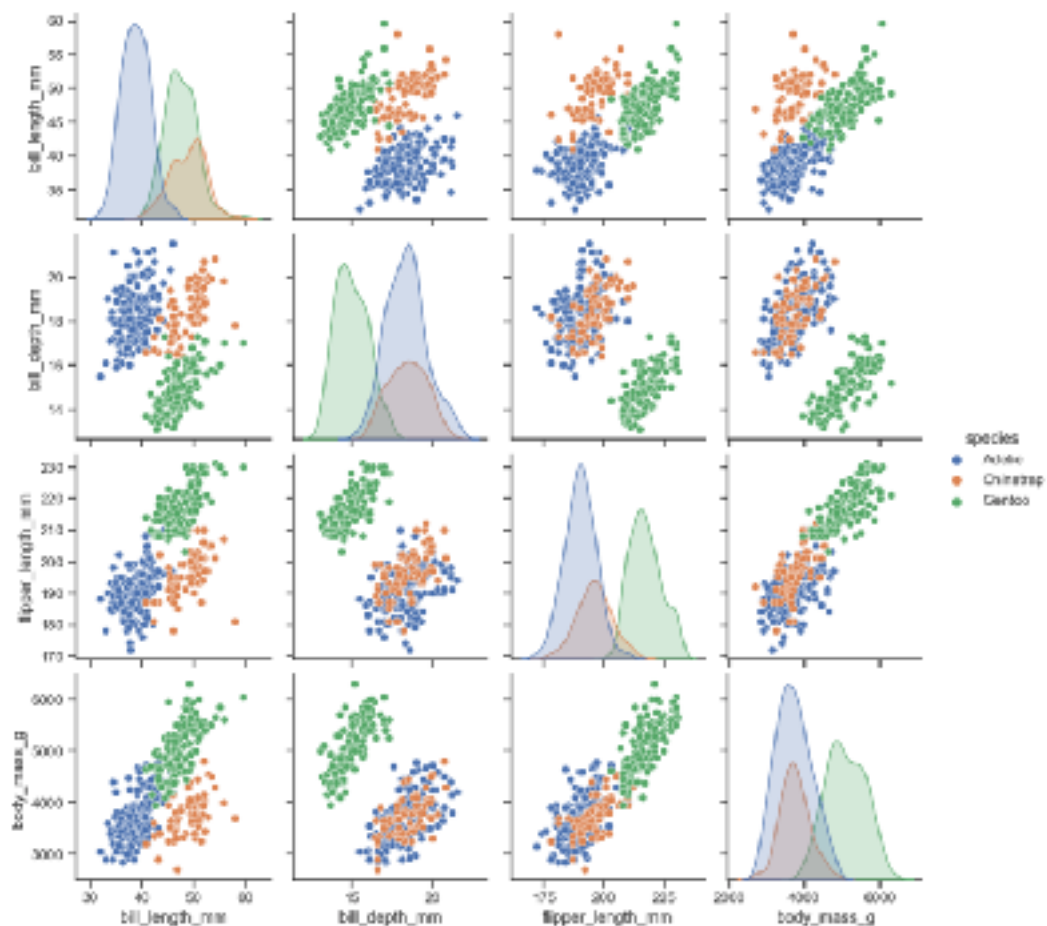
Dimensionality Reduction

projecting high-D data into a lower-D space
i.e.: PCA, ICA, Clustering

```
>>> ax = sns.boxplot(x="day", y="total_bill", hue="time",  
...                  data=tips, linewidth=2.5)
```



```
sns.gpairplot(penguins, hue='species')
```





MNIST data T-SNE projection

