

JSC 270 - LECTURE 2

EXPLORATORY DATA ANALYSIS

<https://jsc270.github.io/>

ANNOUNCEMENT 1

New Office hours for TAs

Matt Edwards: Tues 9a-10a

Lauren Erdman: Thurs 4p-5pm

ANNOUNCEMENT 2

Guest lecturer JSC370: Ben Allison. Tuesday, Jan 19 from 3:10pm to 4:00pm

Ben Allison: Principal Machine Learning Scientist at Amazon. Mr Allison is one of the architects / managers of Amazon's recommender system, which I think is one of the world's largest.

Ben will be talking about the ins and outs of his job, and building and deploying models in real-world settings. Your class is also welcome to stick around for the lecture afterwards.

The zoom link is <https://utoronto.zoom.us/j/89713225921>

Meeting ID: 897 1322 5921

EXPLORATORY DATA ANALYSIS (EDA)

A definition: descriptive and graphical summaries of the data used to understand your data, relationships among variables and problems (missingness, outliers, other data entry errors) that are in the data

Benefit: allow to explore data without making any assumptions

DATA

Subject



Variable/Feature:

Intelligence

Durability

Powers

Strength

Speed

....

Sample



Subject



Sample



DATA

Variable/Feature:

- Cap size
- Cap color
- Habitat
- Poisonous
- Odor

....

Mushroom tabular dataset: <https://archive.ics.uci.edu/ml/datasets/Mushroom>

GETTING DATA

Normally in your DS projects u will be either

Given data

Required to download data

Required to scrape data of the web

Least intensive

Most intensive

TYPES OF DATA

- Categorical
 - Nominal
 - Ordinal
- Numerical
 - Discrete
 - Continuous

CATEGORICAL: NOMINAL

Two or more categories where the ordering doesn't matter

1. province - {Ontario, Quebec, ...}
2. hair color - {black, brunette, blonde,...}

CATEGORICAL: ORDINAL

- Ordinal - order matters but not the difference between values
 1. income {low, medium, high}
 2. education {elementary school, high school, some college, college grad}

BINARY DATA

Data that has two states

Example: bit {0,1}

Surgery: {yes, no}

Is binary data nominal or ordinal?

NUMERICAL DATA: DISCRETE

Numerical data: countable number of values

Can't measure but can count

Example 1:

Number of times apples were bought at a store per day

Example 2:

Number of people in a Data Science program

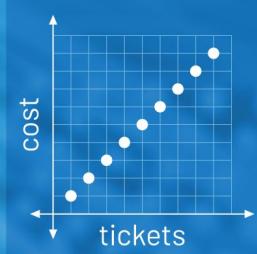
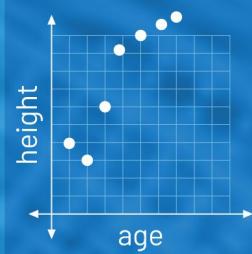
NUMERICAL DATA: CONTINUOUS

Represents measurements – can be measured but not counted

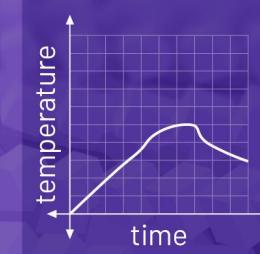
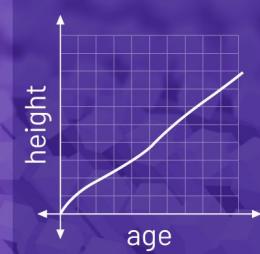
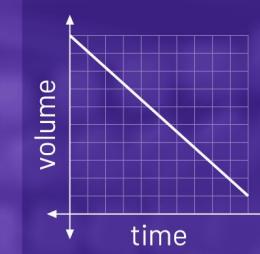
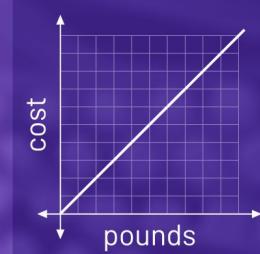
Example: age, weight, height, income

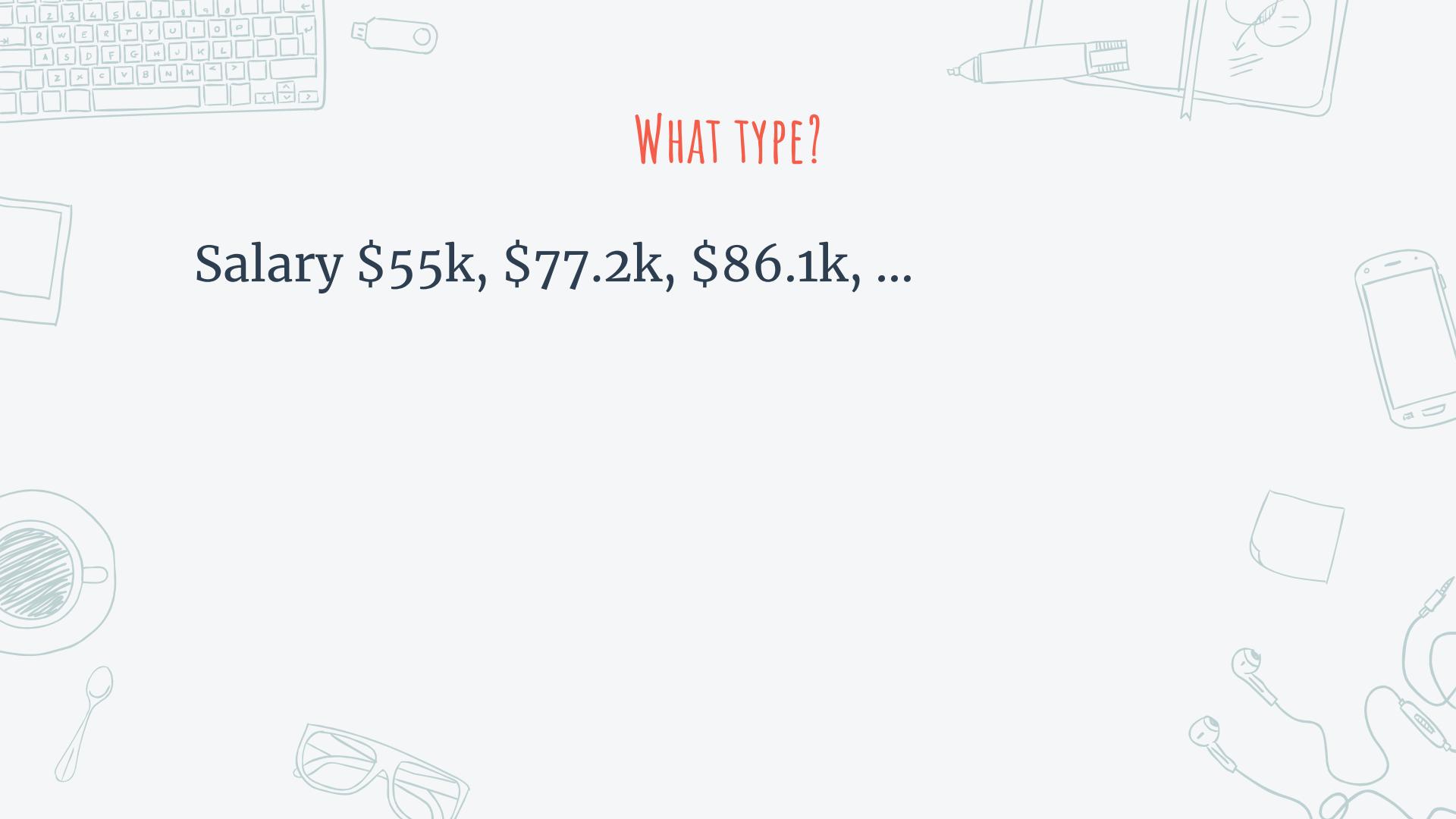
NUMERICAL: DISCRETE VS CONTINUOUS

DISCRETE



CONTINUOUS





WHAT TYPE?

Salary \$55k, \$77.2k, \$86.1k, ...

WHAT TYPE?

Salary \$55k, \$77.2k, \$86.1k, ... - continuous

Salary \$25k, \$26k, \$40k,... (rounded)

WHAT TYPE?

Salary \$55k, \$77.2k, \$86.1k, ... - continuous

Salary \$25k, \$26k, \$40k,... (rounded) - discrete

Salary \$10k-\$19k, \$20k-29k, \$30k-\$39k...

WHAT TYPE?

Salary \$55k, \$77.2k, \$86.1k, ... - continuous

Salary \$25k, \$26k, \$40k,... (rounded) - discrete

Salary \$10k-\$19k, \$20k-29k, \$30k-\$39k... - ordinal

STRUCTURED VS UNSTRUCTURED

Structured Data vs Unstructured Data

Can be displayed in rows, columns and relational databases

XY	1	2
A	A1	A2
B	B1	B2
C	C1	C2
D	D1	D2

Numbers, dates and strings

0, 1, 2,	INV
3, 4, 5,	1234567890
6, 7, 8,	YZ, F-G-H-I
4,2825	

Estimated 20% of enterprise data (Gartner)



Requires less storage



Easier to manage and protect with legacy solutions



vs

Unstructured Data

Cannot be displayed in rows, columns and relational databases



Images, audio, video, word processing files, e-mails



Estimated 80% of enterprise data (Gartner)



Requires more storage



More difficult to manage and protect with legacy solutions

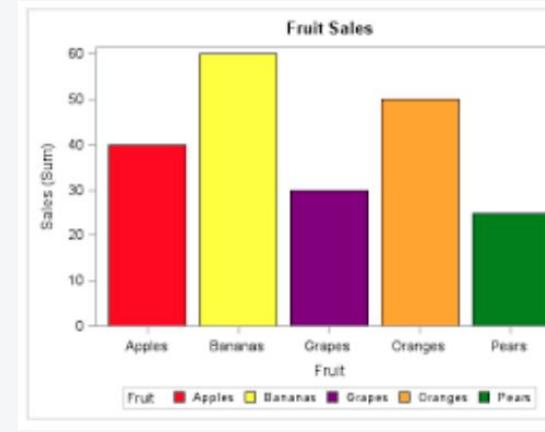
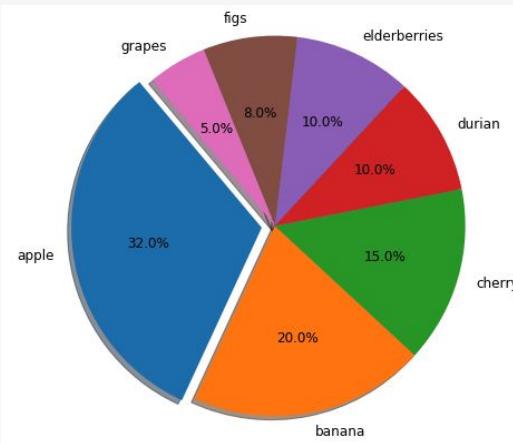


Structured data can be derived from unstructured data and often is

WHY IS IT IMPORTANT?

Types of analysis you could do is different depending on the variable type

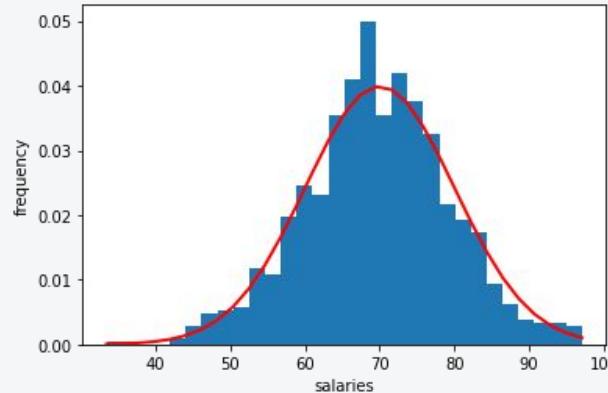
E.g. nominal or ordinal: frequency, proportion or %



WHY IS IT IMPORTANT?

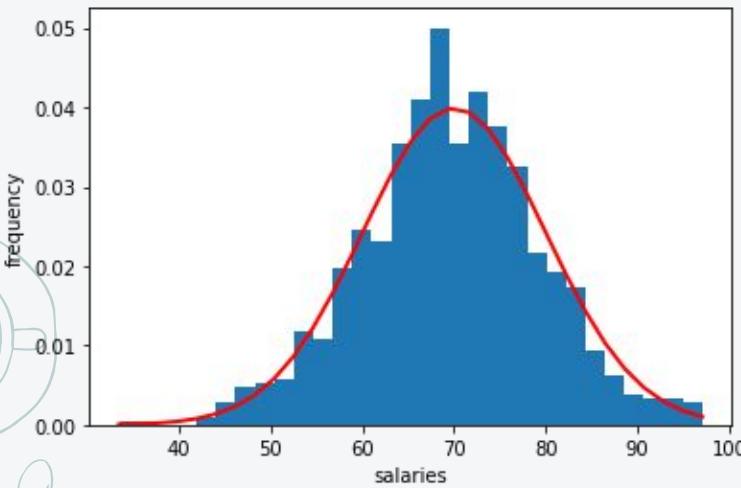
Types of analysis you could do is different depending on the variable type

E.g. distribution for continuous variables (including Gaussian fit)

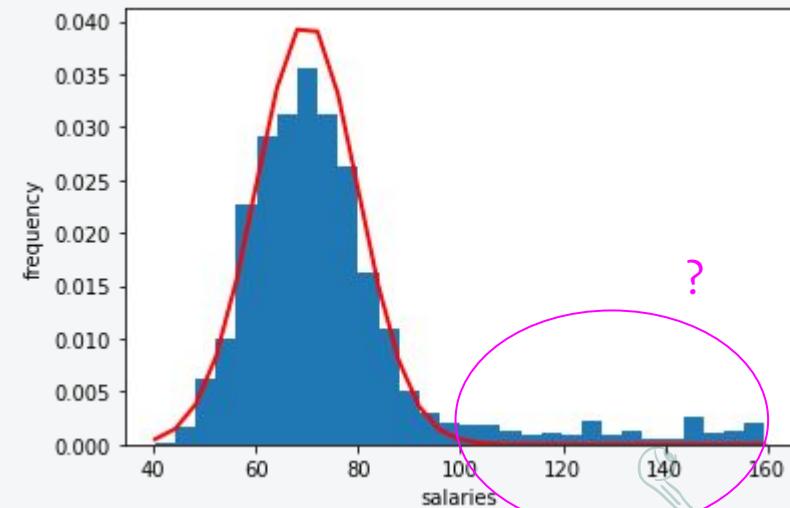


EXAMPLE: SALARIES

Organization A



Organization B



OUTLIER (ANOMALY) DETECTION

Definition ?

OUTLIER (ANOMALY) DETECTION

Definition 1: an observation that deviates markedly from the sample it is found

Definition 2: Observations that are not conforming to the expected behavior

OUTLIER (ANOMALY) DETECTION

Why would we want to detect outliers?

OUTLIER (ANOMALY) DETECTION

Why would we want to detect outliers?
They cause issues with analysis

What do we do?

OUTLIER (ANOMALY) DETECTION

Why would we want to detect outliers?
They cause issues with analysis

What do we do?

1. Detect
2. Remove (?)
3. Impute (?)

AREAS OF APPLICATION

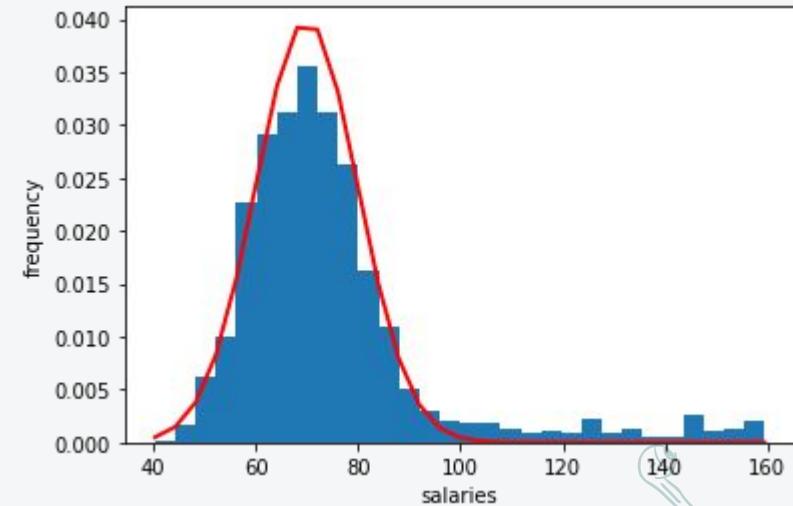


- Fraud Detection
- Bioinformatics/healthcare
- Cyber Security (e.g. networking traffic)
- Trading
- Image Processing
- Many others (can you think of any?)

DETECTION

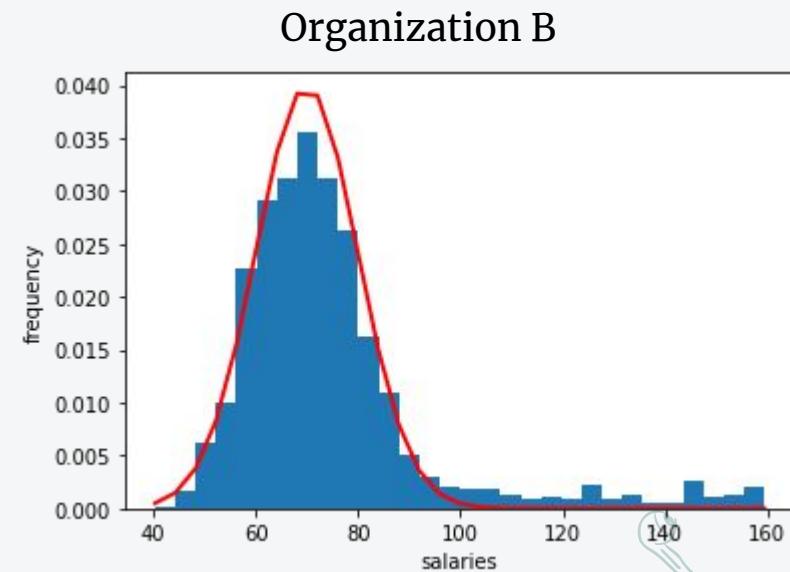
Eyeballing?

Organization B



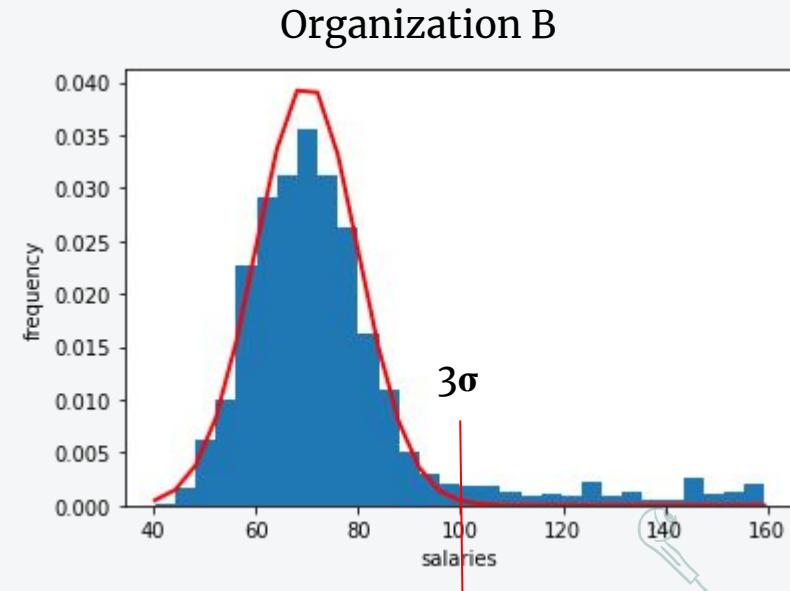
DETECTION

Eyeballing?
Hard to tell...



DETECTION

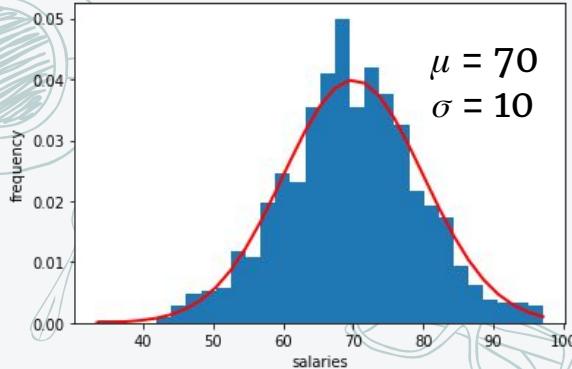
Eyeballing?
Hard to tell...
Statistical Threshold?



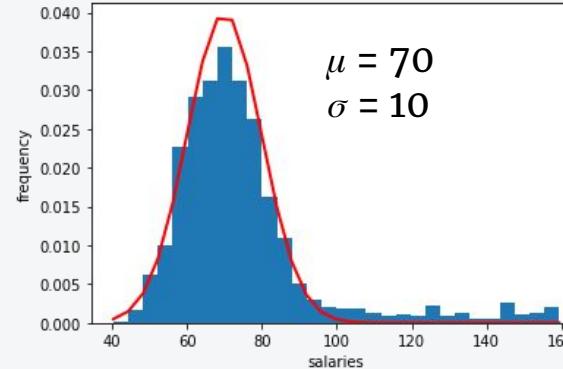
DETECTION

Eyeballing?
Hard to tell...
Statistical Threshold?

Organization A



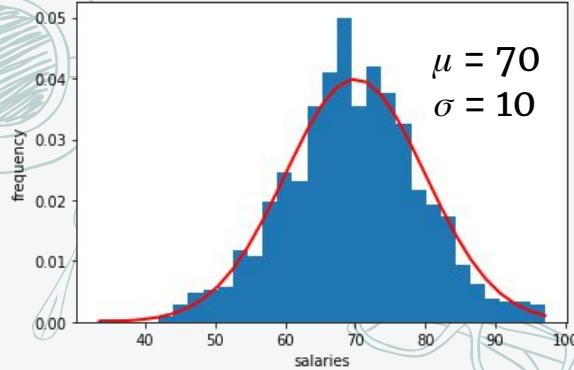
Organization B



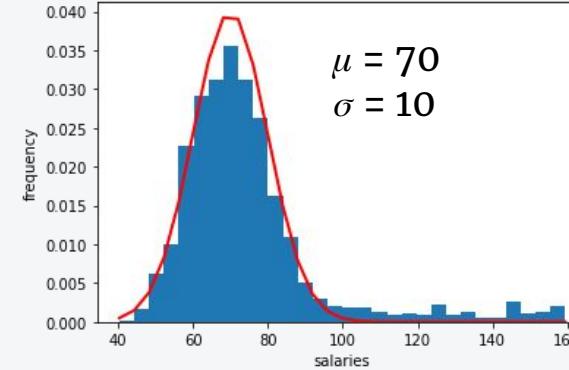
DETECTION

Eyeballing?
Hard to tell...
Statistical Threshold?

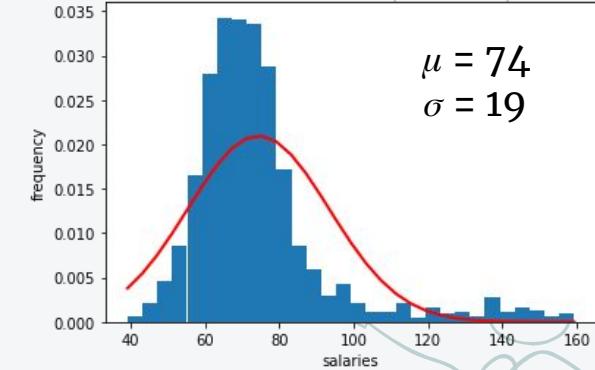
Organization A



Organization B



Organization B

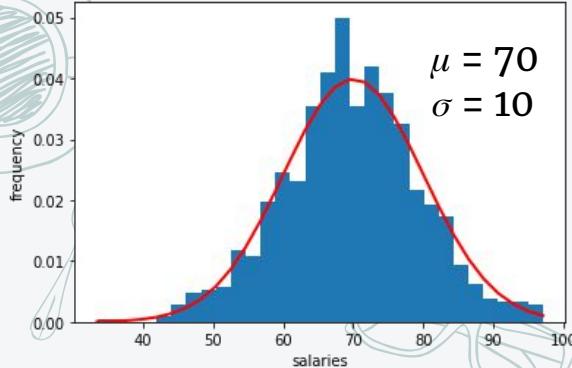


DETECTION

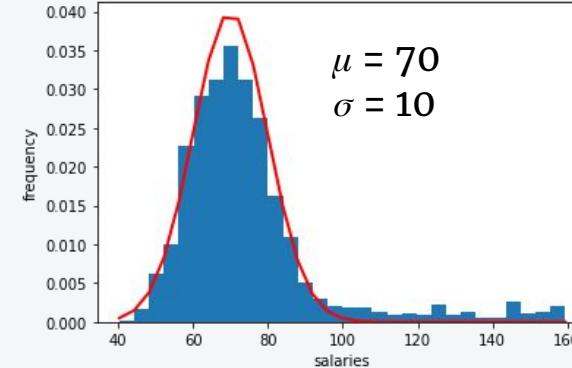
Eyeballing?
Hard to tell...
Statistical Threshold?

Outliers affect the model!!!

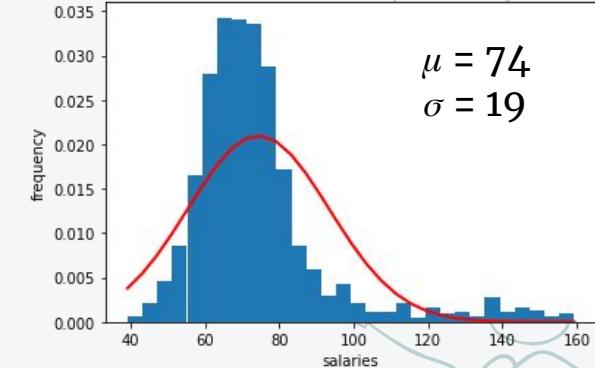
Organization A



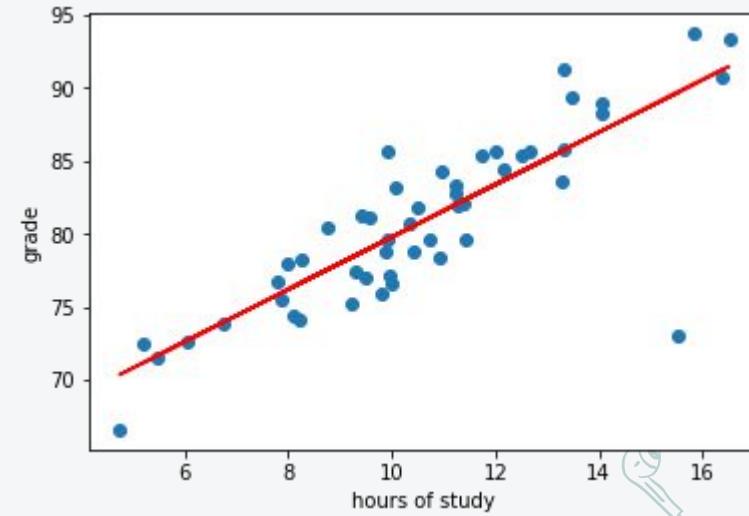
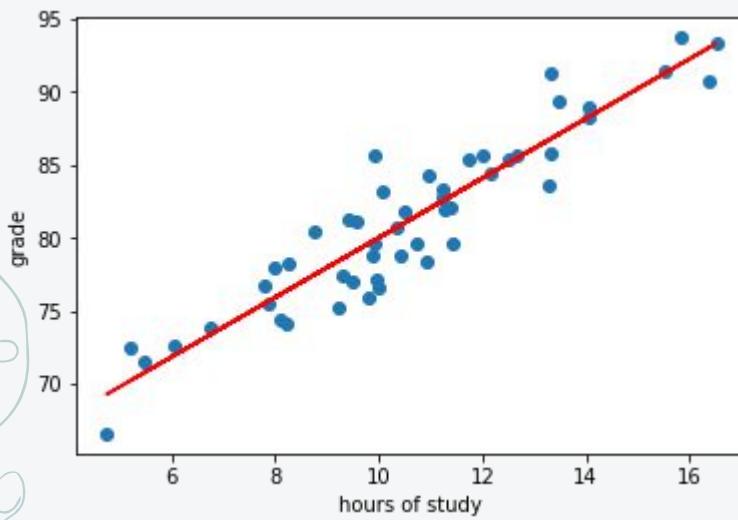
Organization B



Organization B



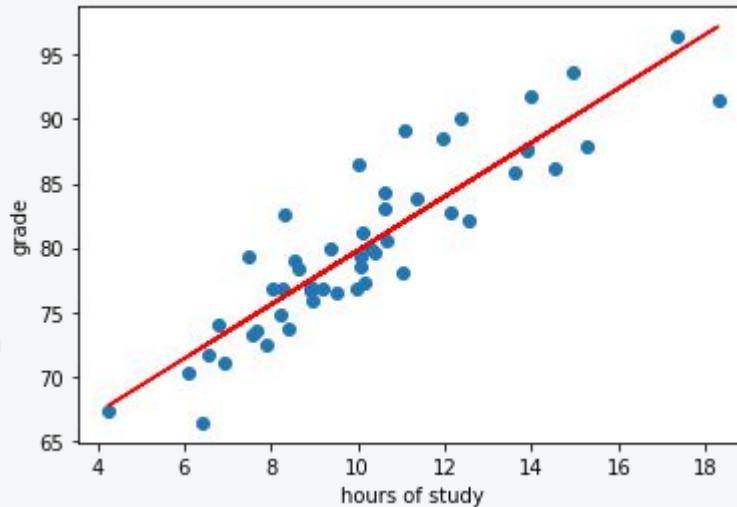
ANOTHER EXAMPLE OF OUTLIERS



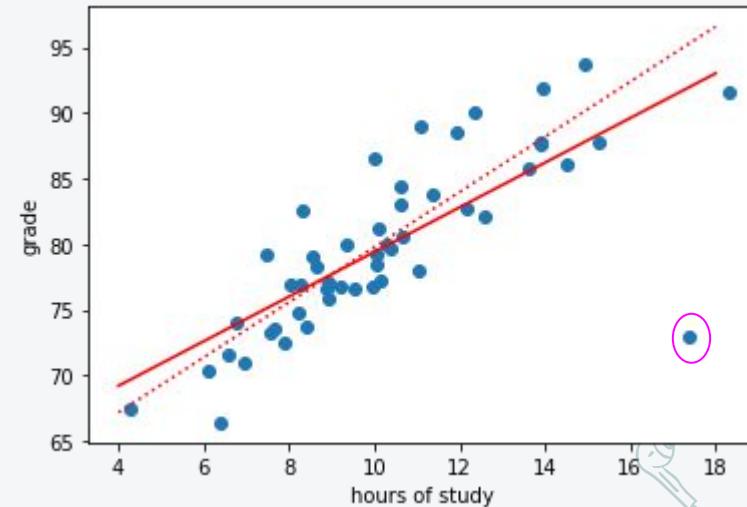
ANOTHER EXAMPLE OF OUTLIERS

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

Pearson correlation $r = 0.9$

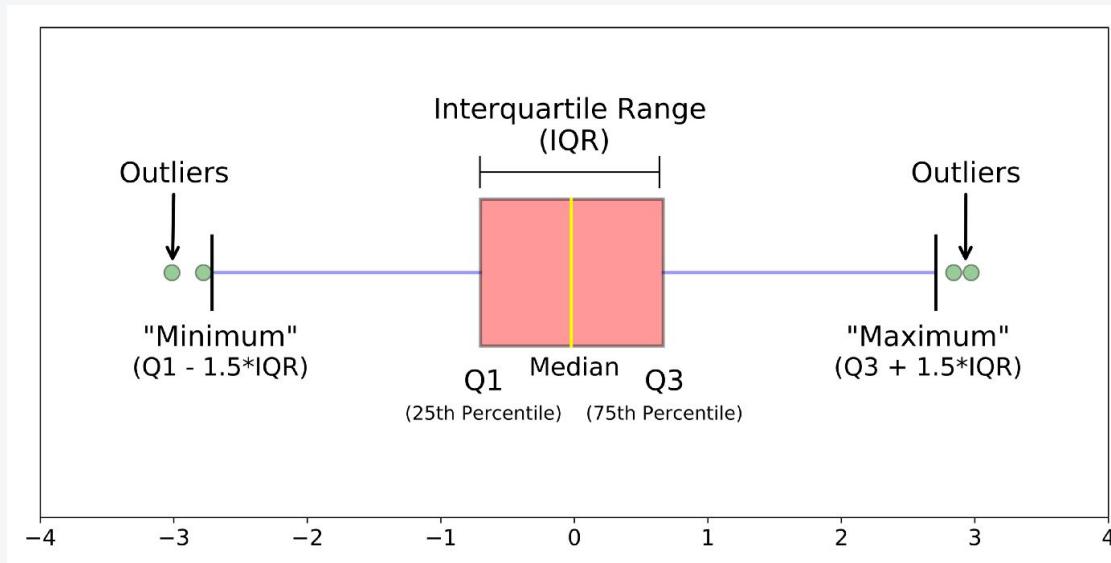


$r = 0.76$



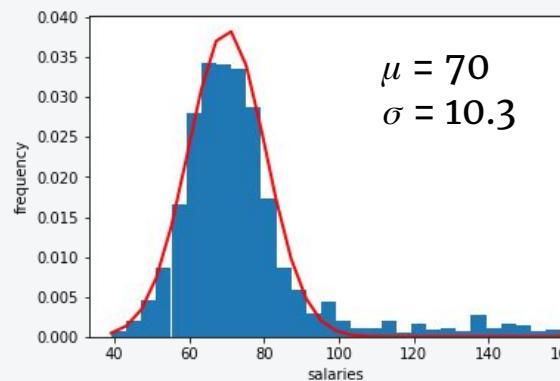
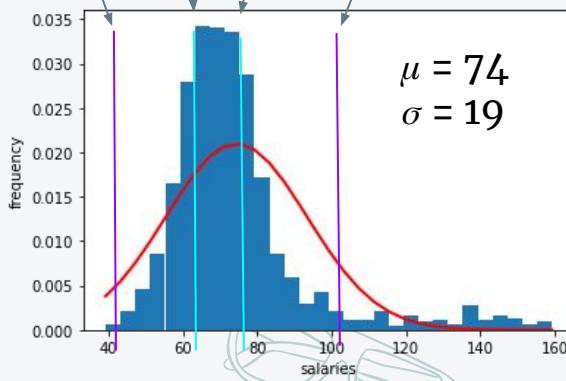
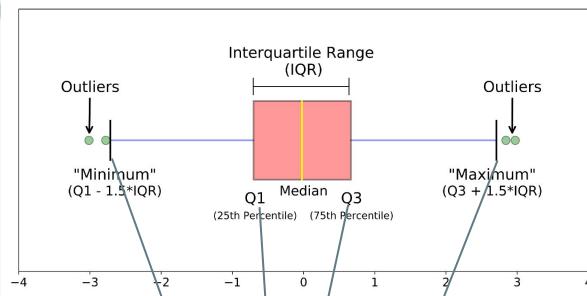
HOW TO MAKE SURE YOUR ANALYSIS IS LESS AFFECTED BY OUTLIERS?

Can estimate parameters using % of data



WHAT CAN YOU DO?

Can estimate parameters with % of data



DATA

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20		NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni		\$40k
Estelle	100	3.2	Festivus Uni	1.7	\$0k
Leo	15	2.4	Festivus Uni	0	
Rachel	50	4.0	Columbia		\$75k

1. OUTLIERS

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20		NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni		\$40k
Estelle	100	3.2	Festivus Uni	1.7	\$0k
Leo	15	2.4	Festivus Uni	0	
Rachel	50	4.0	Columbia		\$75k

1. OUTLIERS

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20		NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni		\$40k
Estelle	100	3.2	Festivus Uni	1.7	\$0k
Leo	15	2.4	Festivus Uni	0	\$35k
Rachel	50	4.0	Columbia		\$75k

1. OUTLIERS

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20		NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni		\$40k
Estelle	100	3.2	Festivus Uni	1.7	\$0k
Leo	15	2.4	Festivus Uni	0	\$35k
Rachel	50	4.0	Columbia		\$75k

remove

1. OUTLIERS

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20		NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni		\$40k
Leo	15	2.4	Festivus Uni	0	\$35k
Rachel	50	4.0	Columbia		\$75k

2. MISSING DATA

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20		NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni		\$40k
Leo	15	2.4	Festivus Uni	0	\$35k
Rachel	50	4.0	Columbia		\$75k

2. MISSING DATA

25% missing data

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20		NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni		\$40k
Leo	15	2.4	Festivus Uni	0	\$35k
Rachel	50	4.0	Columbia		\$75k

2. MISSING DATA

25% missing

Remove?

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20		NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni		\$40k
Leo	15	2.4	Festivus Uni	0	\$35k
Rachel	50	4.0	Columbia		\$75k

2. MISSING DATA

25% missing

R = 65%

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20		NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni		\$40k
Leo	15	2.4	Festivus Uni	0	\$35k
Rachel	50	4.0	Columbia		\$75k

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20		NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni		\$40k
Leo	15	2.4	Festivus Uni	0	\$35k
Rachel	50	4.0	Columbia		\$75k

2. M
R = 84%

25% missing
R = 65%

IMPUTATION: WITH MEAN

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20	3.1	NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni	2.4	\$40k
Leo	15	2.4	Festivus Uni	0	\$35k
Rachel	50	4.0	Columbia	2.4	\$75k

IMPUTATION: WITH MEAN

$$\begin{aligned}\sigma_{\text{original}} &= 1.9 \\ \sigma_{\text{imputed}} &= 1. \\ 6 &\end{aligned}$$

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20	3.1	NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni	2.4	\$40k
Leo	15	2.4	Festivus Uni	0	\$35k
Rachel	50	4.0	Columbia	2.4	\$75k

IMPUTATION: WITH MEAN

$$\begin{aligned}\sigma_{\text{original}} &= 1.9 \\ \sigma_{\text{imputed}} &= 1.6\end{aligned}$$

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20	3.1	NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni	2.4	\$40k
Leo	15	2.4	Festivus Uni	0	\$35k
Rachel	50	4.0	Columbia	2.4	\$75k

$$\begin{aligned}R_{\text{original}} &= 65\% \\ R_{\text{imputed}} &= 56\%\end{aligned}$$

IMPUTATION: FROM SIMILAR INDIVIDUAL

$$\begin{aligned}\sigma_{\text{original}} &= 1.9 \\ \sigma_{\text{imputed}} &= 1.6 \\ \sigma_{\text{similar}} &= 1.9\end{aligned}$$

Student Name	Avg Hours Studied per week	GPA	University	Sense of Humour (0-5)	Salary
George	20	3.1 -> 2.8	NYU	3	\$40k
Jerry	35	3.5	Columbia	5	\$80k
Elaine	55	4.0	Columbia	4.2	\$60k
Cosmo	5	2.0	City College	2	\$25k
Newman	25	2.8	City College	0	\$50k
Frank	35	3	Festivus Uni	2.4 -> 0	\$40k
Leo	15	2.4	Festivus Uni	0	\$35k
Rachel	50	4.0	Columbia	2.4 -> 4.2	\$75k

$$\begin{aligned}R_{\text{original}} &= 65\% \\ R_{\text{imputed}} &= 56\% \\ R_{\text{similar}} &= 71\%\end{aligned}$$

OTHER IMPUTATION METHODS

MICE - draw from distribution of the variable, create multiple rows per individual

Regression

Random Forest

<https://scikit-learn.org/stable/modules/impute.html>

READING MATERIAL

Data Types:

<https://towardsdatascience.com/data-types-in-statistics-347e152e8bee>

Visualization:

<https://jingwen-z.github.io/data-viz-with-matplotlib-series3-pie-chart/>

Reminders:

Data representation, random variables:

https://www.youtube.com/watch?v=vUpuj9K1e4k&list=PLXBDYmaCbeL8efhOZS4g9W6Z3m9_hFSnT&index=5&ab_channel=JeffLeek