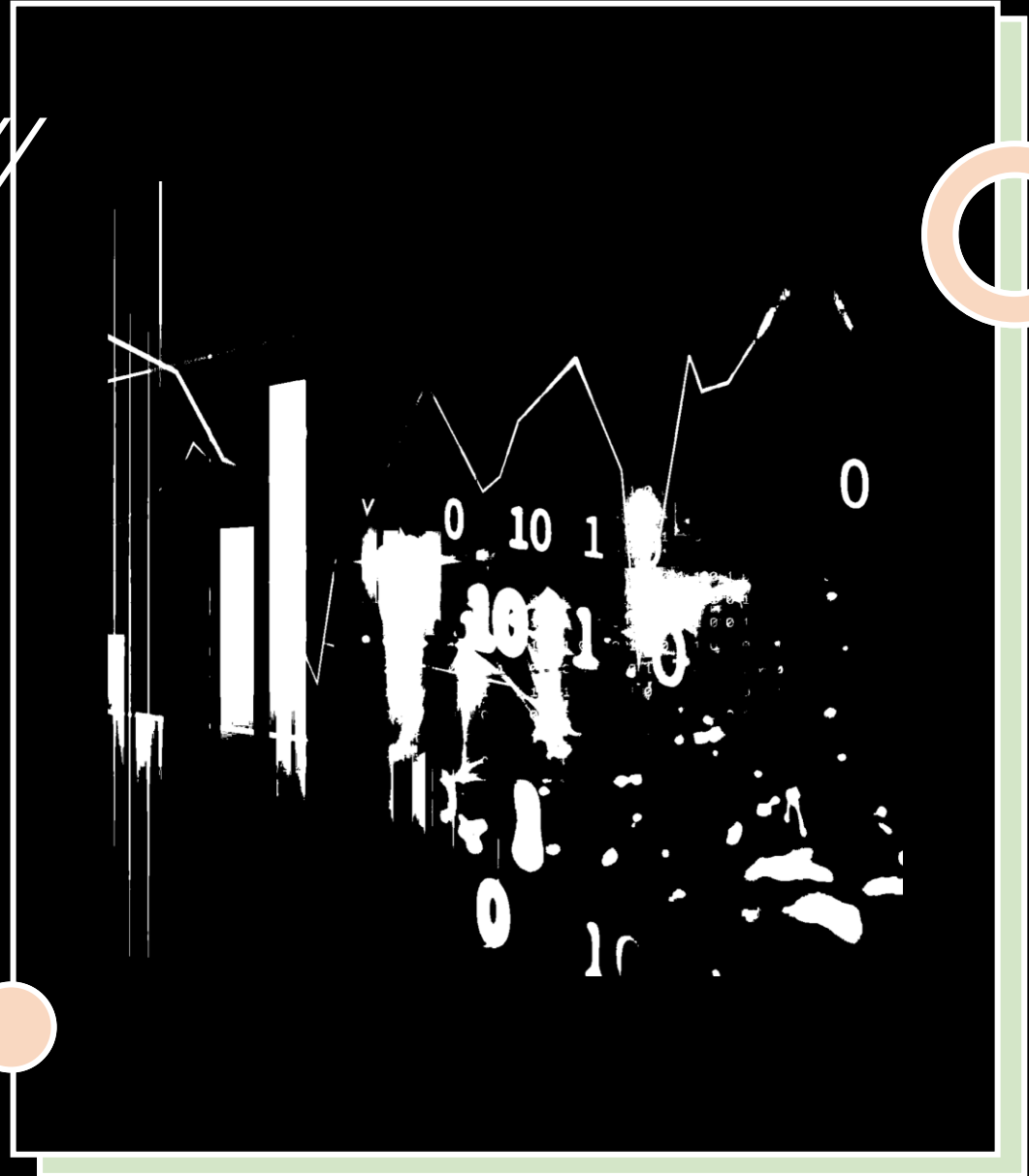


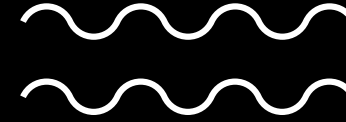
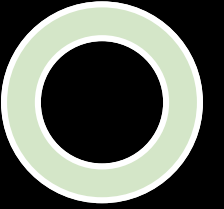
A Brute Force Method of Discovering Information Hidden in Datasets

Kristopher Kurt Honetschlager, Computer Science Major, Winona State University-Rochester '23

kristopher.honetschlager@go.winona.edu

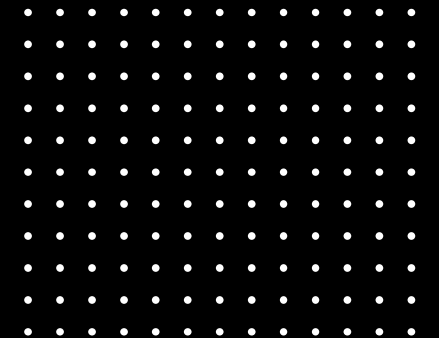
Winona State University at the Rochester Campus





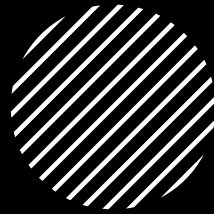
- Computer Science can at least partly be defined as “...the study of computers and computing, including their theoretical and algorithmic foundations, hardware and software, and their uses for processing information.”^{.1}
- Statistics can at least partly be defined as “...the science of collecting, analyzing, presenting, and interpreting data.”^{.2}
- Logic can at least partly be defined as “The science that investigates the principles governing correct or reliable inference.”^{.3}
- Data Science can at least partly be defined as “...[a] field of study that combines domain expertise, programming skills, and knowledge of mathematics and statistics to extract meaningful insights from data.”^{.4}

Computer
Science/Statistics/Logic/
Data Science What are
they?(Background
Information)





Motivation (Mushroom Dataset Theory/Hypothesis)



- Contrarianism: guide states “...no simple rule for determining the edibility of a mushroom...”
 - unlike “...`leaflets three, let it be’ for Poisonous Oak and Ivy.”⁸
- The .names file lays out simple rules for determining whether a given mushroom instance is poisonous or not. Such as: “spore-print-color=green 48 cases missed, 99.41% accuracy”.⁸
- I figured a brute force method of somehow testing every column and pairing of columns against the poisonous column and retrieving a highest score would be a very good way of obtaining these rules.⁸



⁸Schlimmer, Jeff. “Mushroom Data Set.” Edited by Dave Aha, UCI Machine Learning Repository: Mushroom Data Set, University of California, Irvine, 27 Apr. 1987, <https://archive.ics.uci.edu/ml/datasets/Mushroom>.

Example Mushroom Dataset

no_odor	spore-print-color=green	habitat=leaves	cap-color=white	poisonous=yes
0	0	0	0	0
0	0	0	0	0
0	0	0	0	1
1	0	0	1	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	1
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

Number of Columns Selected = 4

no_odor

spore-print-

color=green

habitat=leaves

cap-color-

white

no_odor

spore-print-

color=green

habitat=leaves

cap-color-

white

no_odor

spore-print-

color=green

habitat=leaves

cap-color-

white

no_odor

spore-print-

color=green

habitat=leaves

cap-color-

white

Number of Columns Selected = 1

Number of Columns
Selected = 2

Number of
Columns Selected
= 3

Columns
Selected

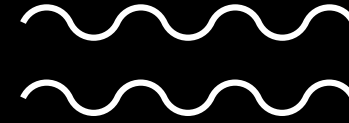
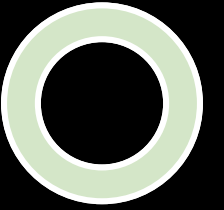
= 4

Number
of

no_odor
spore-print-
color=green
habitat=leaves
cap-color-
white

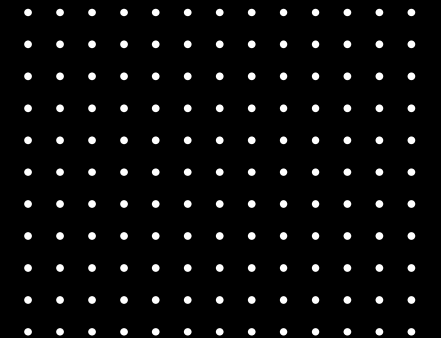
no_odor
spore-print-
color=green
habitat=leaves
cap-color-
white

no_odor
spore-print-
color=green
habitat=leaves
cap-color-
white



- Why not test every possible column and pairing of columns against every possible label to obtain the best rule obtainable for predicting each label in the dataset.
 - Could be a rule that predicts a label that was previously unknown, the hidden aspect, or finds one contrary to current belief, the contrarian aspect.
- Acts like a hypothesis machine
- Selected columns can be referred to as the X or feature vector in data science₇
- Row, also referred to as instance₇
- Column, also referred to as attribute₇
- Poisonous can be referred to as the label in data science₇

Extension of Mushroom Dataset Theory/Hypothesis and Terminology



Example Mushroom Dataset

no_odor	spore-print-color=green	habitat=leaves	cap-color-white	poisonous=yes
0	0	0	0	0
0	0	0	0	0
0	0	0	0	1
1	0	0	1	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	1
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

Number of Columns Selected = 1



Photo by sk8erteck on deviantart.com

Sort of in-Depth Explanation of Method

A
B
C
D
E

A
B
C
D
E

A

1

0

1

0

0

B

0

1

0

1

1

ABC, ABD, and
BCD and BCE
CDE

AB, AC, AD, and AE

BC, BD, and BE

CD and CE

DE

Handwritten notes and diagrams illustrating the method for determining relationships between variables A, B, C, D, and E. The notes include:

- Mooc Dataset (A)**: A table showing the relationship between A and B, with columns for A, B, and C. The table is used to calculate the percentage of instances where A is 1 and B is 1, and where A is 1 and B is 0.
- Value Counts**: A table showing the distribution of values for A, B, and C. The table is used to calculate the percentage of instances where A is 1 and B is 1, and where A is 1 and B is 0.
- Statistics**: A table showing the distribution of values for A, B, and C. The table is used to calculate the percentage of instances where A is 1 and B is 1, and where A is 1 and B is 0.
- Relationships**: A table showing the distribution of values for A, B, and C. The table is used to calculate the percentage of instances where A is 1 and B is 1, and where A is 1 and B is 0.
- Key**: A table showing the distribution of values for A, B, and C. The table is used to calculate the percentage of instances where A is 1 and B is 1, and where A is 1 and B is 0.
- Perfect Rule**: A table showing the distribution of values for A, B, and C. The table is used to calculate the percentage of instances where A is 1 and B is 1, and where A is 1 and B is 0.
- What can we deduce from the following scenarios?**: A table showing the distribution of values for A, B, and C. The table is used to calculate the percentage of instances where A is 1 and B is 1, and where A is 1 and B is 0.

The notes also include a diagram showing the relationships between A, B, C, D, and E, and a table showing the distribution of values for A, B, and C.

Statistics: Dataset Types

- You can use this method on categorical datasets and binary datasets
 - For categorical datasets you must break categorical columns into binary columns

```
3. cap-color:      brown=n, buff=b, cinnamon=c, gray=g, green=r,  
                  pink=p, purple=u, red=e, white=w, yellow=y
```

- My method fundamentally operates on binary datasets



poisonous		poisonous_yes	poisonous_no
No		0	1
No		0	1
Yes		1	0
No		0	1

← CATEGORICAL DATASET		
poisonous	cap-shape	cap-color
No	x	n
No	x	y
Yes	b	w
No	x	w

BINARY DATASET →		
poisonous_yes	cap-color_w	cap-shape_b
0	0	0
0	0	1
1	1	0
0	0	0

Statistics: Binary Dataset Value Counts

Example Binary Dataset

John Doe is Happy	John Doe had ice cream
1	1
1	1
1	1
1	1
1	1
0	0
0	0
0	0
0	0
0	0
0	0

Example Binary Dataset Value Counts

John Doe is Happy	John Doe Had ice cream	Number of Occurences
1	1	5
1	0	0
0	1	0
0	0	5

Value counts represents all possible unique binary combinations given the number of columns in the dataset.

Logic: Conditional Statement Forms

- If p , where p is the hypothesis, then q , where q is the conclusion
 - “If it’s sunny and there’s no snow on the ground, I will go on a walk outside.”
- By traditional logic, “A conditional statement is **not** logically equivalent to its inverse.”⁷. Though perhaps instinctually it may seem that way.
 - Therefore, following from the previous conditional statement it cannot be additionally concluded that “If it’s not sunny and there is snow on the ground, I won’t go on a walk outside.”
- But if we have empirical evidence suggesting the first statement and its inverse occur at a respective frequency of greater than 25% in a given dataset. We may be able to imply that the beginning conditional statement, though not logically equivalent to its inverse, is empirically equivalent to its inverse.
- [TERM] The *denial* of the first point is “If p then not q ”.
 - “If it’s sunny and there’s no snow on the ground, I won’t go on a walk outside.”

Example: Empirical Observation of a Rule

Example Binary Dataset

Example Binary Dataset		
John Doe is Happy	John Doe Had ice cream	
1	1	5
1	0	0
0	1	0
0	0	5
0	0	0
John Doe is happy; did he have ice cream today?		
John Doe is not happy; did he have ice cream today?		

Generation of Disjunctive Rules

- Empirically, rules can be for identifying predictors that a given label's instance is positive.
 - If X is true and Y is true, then Z is true.
- Predictors that the given instance label is negative.
 - If X is false and Y is false, then Z is false.
- Or predictors indicative of the given label's instance being positive or negative based on the inversion of either set of predictors.
 - If X is true and Y is true, then Z is true and if X is false and Y is false, then Z is false.

Statistics: Characteristics of Binary Value Counts

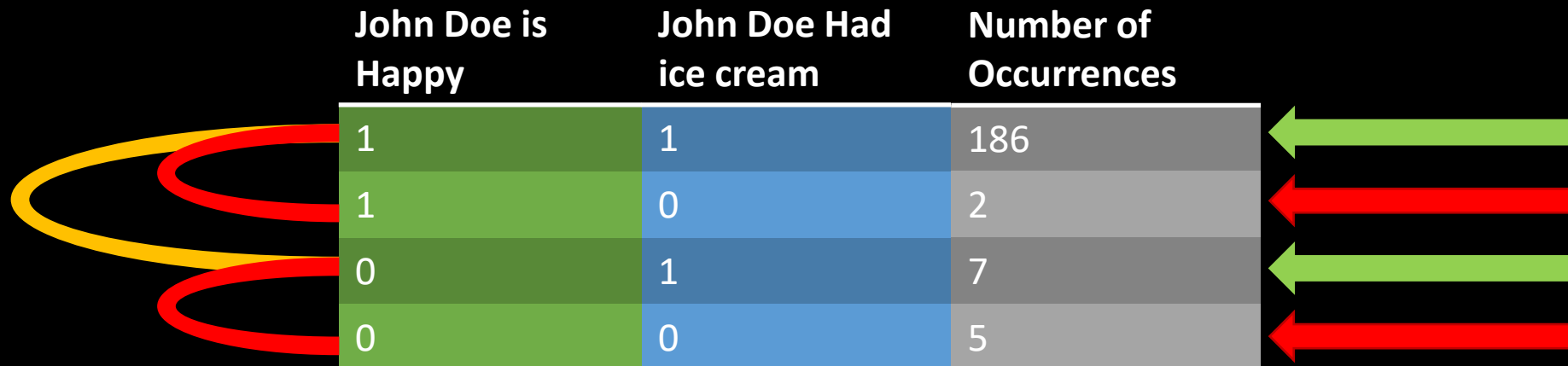
Positive Rule Generation

Occurrence frequency of John Doe is happy, and John Doe had Ice Cream =
 $186/200 = .93 * 100 = 93\%$

Occurrence frequency of John Doe isn't happy, and John Doe had Ice Cream =
 $4/200 = .02 * 100 = 2\%$

Example Binary Dataset Value Counts

John Doe is Happy	John Doe Had ice cream	Number of Occurrences
1	1	186
1	0	2
0	1	7
0	0	5



If one “positive” rule occurs the most frequently in the “positive” rules and in the dataset, greater than 50% respectively; it implies the denials of the rule occur technically less than 50% of the time.

John is happy, did he have ice cream?

Statistics: Characteristics of Binary Value Counts

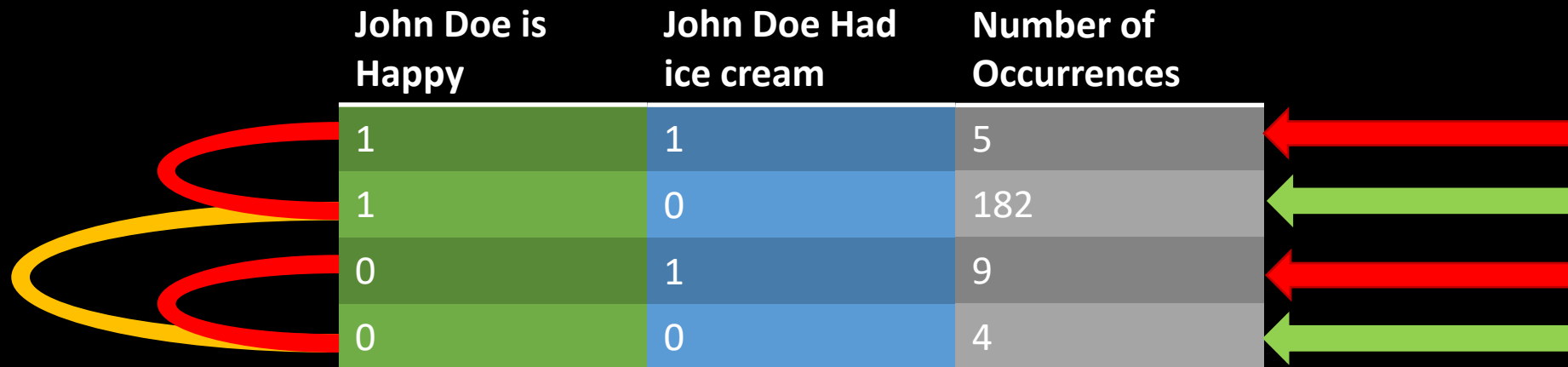
Occurrence frequency of John Doe is happy, and John Doe didn't have Ice Cream =
 $182/200 = .91 * 100 = 91\%$

Negative Rule Generation

Occurrence frequency of John Doe isn't happy, and John Doe didn't have Ice Cream =
 $4/200 = .02 * 100 = 2\%$

Example Binary Dataset Value Counts

John Doe is Happy	John Doe Had ice cream	Number of Occurrences
1	1	5
1	0	182
0	1	9
0	0	4



If one "negative" rule occurs the most frequently in the "negative" rules and in the dataset, greater than 50% respectively; it implies the denials of the rule occur technically less than 50% of the time.

John is happy, did he not have ice cream?

Statistics: Characteristics of Binary Value Counts

Total number of rows in dataset = 200

Combinatorial Rule Generation

Occurrence frequency of John Doe is happy and didn't have ice cream = $95/200 = 0.475 * 100 = 47.5\%$

Occurrence frequency of John Doe is happy and did have ice cream = $95/200 = 0.475 * 100 = 47.5\%$

Example Binary Dataset Value Counts

John Doe is Happy	John Doe Had ice cream	Number of Occurrences
1	1	100
1	0	100
0	1	0
0	0	0

The rule John Doe is happy therefore he had ice cream and John Doe is happy therefore he didn't have ice cream do not deny one another and therefore have no combinatorial effect.

Statistics: Characteristics of Binary Value Counts

Total number of rows in dataset = 200

Combinatorial Rule Generation Cont.

Occurrence frequency of John Doe is happy, he's had ice cream and John Doe isn't happy, he hasn't had ice cream = $100/200 = .5 * 100 = 50\%$

Example Binary Dataset Value Counts

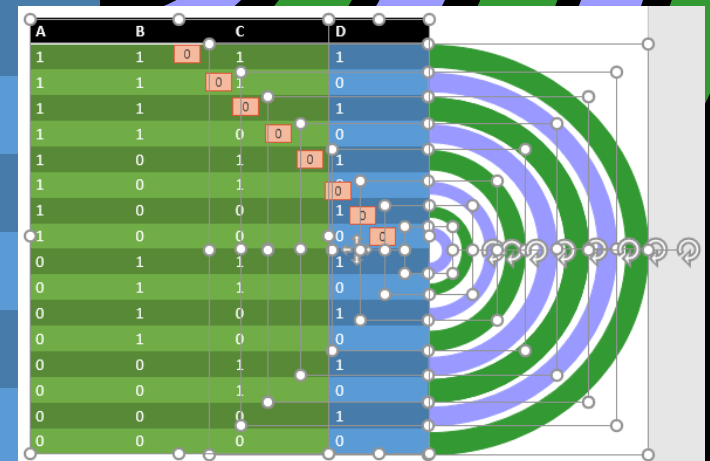
John Doe is Happy	John Doe Had ice cream	Number of Occurrences
1	1	50
1	0	50
0	1	50
0	0	50

Occurrence frequency of John Doe is happy, he hasn't had ice cream and John Doe isn't happy, John Doe has had ice cream = $100/200 = .5 * 100 = 50\%$

“If John Doe is happy, he's had ice cream and if John Doe isn't happy, he hasn't had ice cream.” is the denial of the rule “If John Doe is happy, he hasn't had ice cream and if John Doe isn't happy, John Doe has had ice cream.”.

Value Counts
Dataset from
Hypothetical
Dataset

Value Counts	A	B	C	D
100	1	1	1	1
100	1	1	1	0
100	1	1	0	1
100	1	1	0	0
100	1	0	1	1
100	1	0	1	0
100	1	0	0	1
100	1	0	0	0
100	0	1	1	1
100	0	1	1	0
100	0	1	0	0
100	0	0	1	1
100	0	0	1	0
100	0	0	0	1
100	0	0	1	0
100	0	0	0	1
100	0	0	0	0



Mathematics/Computer Science: Where's the Math and Computer Science?

Value Counts Dataset from Hypothetical Dataset

	Value Counts	A	B	C	D
0	100	1	1	1	1
1	100	1	1	1	0
2	100	1	1	0	1
3	100	1	1	0	0
4	100	1	0	1	1
5	100	1	0	1	0
6	100	1	0	0	1
7	100	1	0	0	0
8	100	0	1	1	1
9	100	0	1	1	0
10	100	0	1	0	1
11	100	0	1	0	0
12	100	0	0	1	1
13	100	0	0	1	0
14	100	0	0	0	1
15	100	0	0	0	0

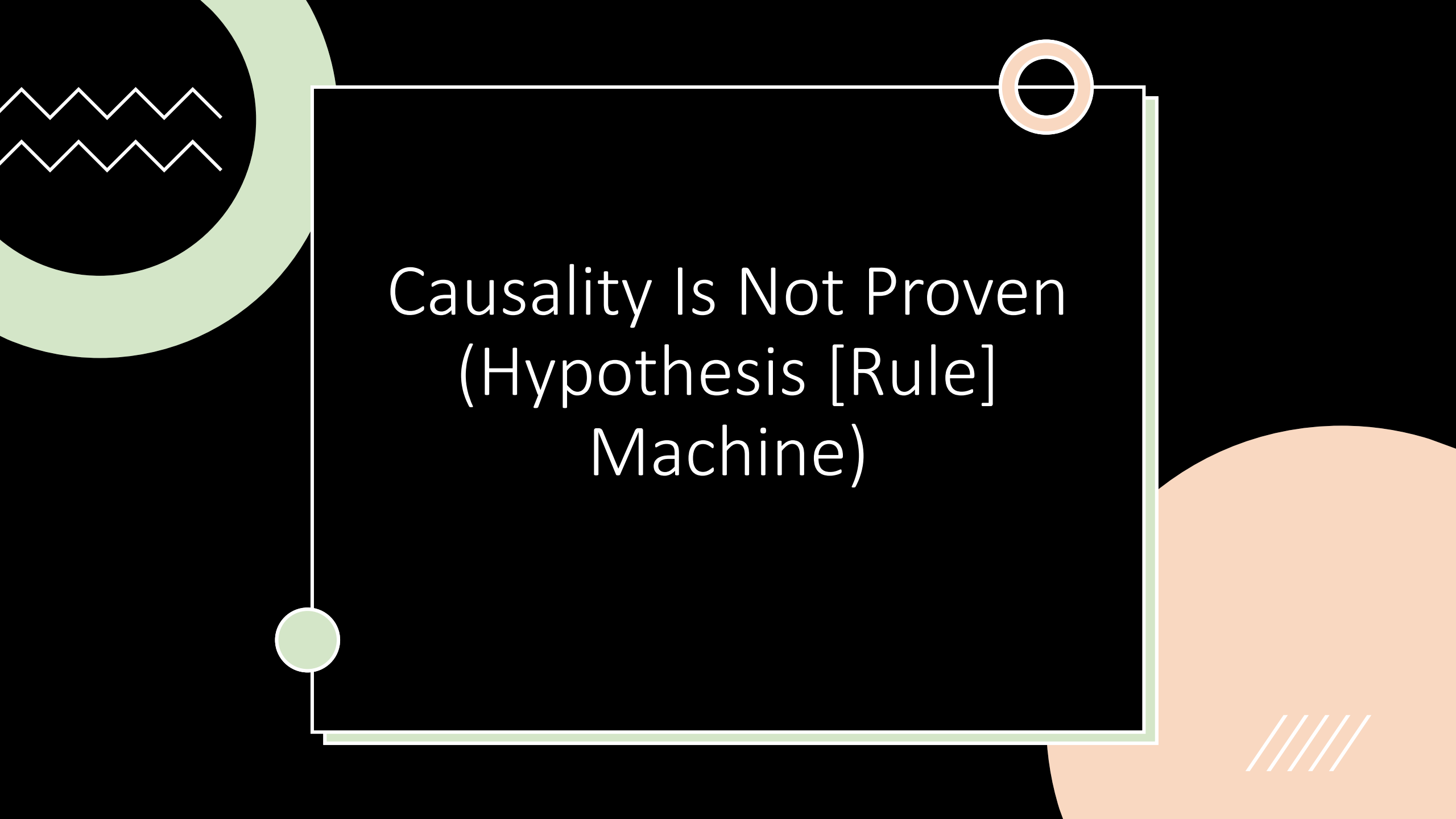
Total number of columns in value counts dataset (not counting value counts column) = 4 = n

Index of inverse of current instance, when the current instance's index is less than $2^n/2 = (2^n - 1) - \text{index}$ current instance

Example using instance at 0: current index $\leq 2^4/2$, therefore inverse of instance index = $(2^4 - 1) - 0 = 15$

Every instance that when it's index modulus 2 is equal to 0 has a positive label

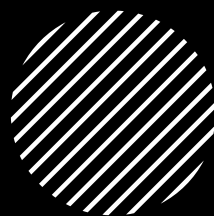
Every instance that when it's index modulus 2 is equal to 1 has a negative label



Causality Is Not Proven
(Hypothesis [Rule]
Machine)



Demonstration



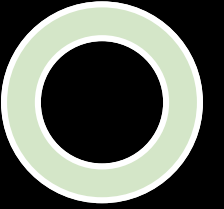
Binary Columns in Dataset (Instances in dataset = 40, Upsampled):

- john-doe_is-happy
- john-doe_had-dinner

If John Doe is happy, John Doe had dinner: {1, 1}

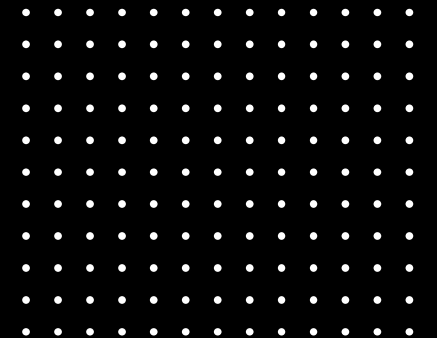
If John Doe is unhappy, John Doe didn't have dinner: {0, 0}

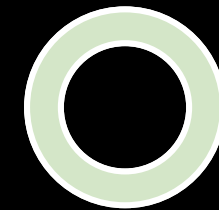
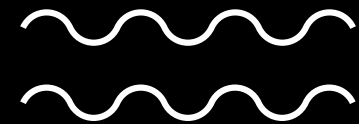
ULTIMATE GOAL is to find a good pairing with machine learning label



- Phenotyping Algorithms “... are special tools that enable researchers to extract phenotypes from complex, and often messy data that get generated during routine interactions within the healthcare system.”₄.
- Induction Logic Programming
- Disjunctive Rule Generation

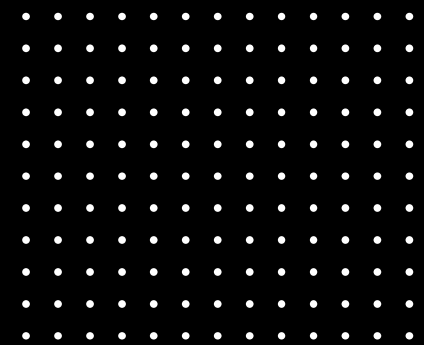
Future Research and Applications





- UC Irvine for hosting mushrooms dataset
- Collin Engstrom my independent study advisor
- Minnstate Undergraduate Scholars Symposium
- Listeners

Thank you





Works Cited

- ¹Belford, Geneva G. and Tucker, Allen. "computer science". Encyclopedia Britannica, Invalid Date, <https://www.britannica.com/science/computer-science>. Accessed 20 March 2022.
- ²Sweeney, Dennis J. , Anderson, David R. and Williams, Thomas A.. "statistics". Encyclopedia Britannica, Invalid Date, <https://www.britannica.com/science/statistics>. Accessed 20 March 2022.
- ³"Logic Definition & Meaning." Dictionary.com, Dictionary.com, <https://www.dictionary.com/browse/logic>.
- ⁴"Data Science." *DataRobot*, 20 Jan. 2022, <https://www.datarobot.com/wiki/data-science/>.
- ⁵"Phenotype Library: HDRUK." Phenotype Library | HDRUK, <https://phenotypes.healthdatagateway.org/>.
- ⁶Sk8erteck. "Light Shining through Windows." *Deviant Art*, 18 Mar. 2004, <https://www.deviantart.com/sk8erteck/art/Light-Shining-through-Windows-5930403>. Accessed 18 Apr. 2022.
- ⁷Kohavi, Ron. "Glossary of Terms." Glossary of Terms Journal of Machine Learning, Kluwer Academic Publishers, 1998, <http://robotics.stanford.edu/~ronnyk/glossary.html>.
- Additional citations on the bottom of some slides