

---

---

---

---

---



THEME O

# ① WHAT IS IT ABOUT

STATISTICS, MACHINE LEARNING,  
PROBABILITY

## • LEARNING FROM DATA

- PREDICTION ←
- INFERENCE
- CAUSALITY
- UNCERTAINTY QUANTIFICATION ←

itself; we discuss how to set these using additional data. Machine learning is essentially a form of applied statistics with increased emphasis on the use of computers to statistically estimate complicated functions and a decreased emphasis on proving confidence intervals around these functions; we therefore present the

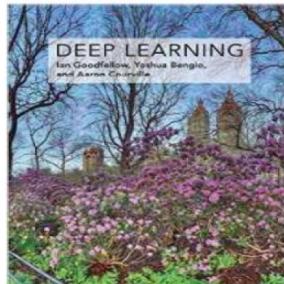
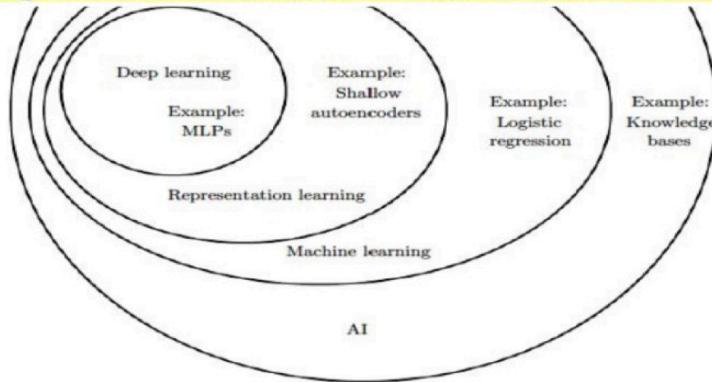


Figure 1.4: A Venn diagram showing how deep learning is a kind of representation learning, which is in turn a kind of machine learning, which is used for many but not all approaches to AI. Each section of the Venn diagram includes an example of an AI technology.

1. DATA VIZ & SUMMARIZATION
2. FROM RANDOMIZATION TO RANDOMNESS
3. WHAT IS PROBABILITY & WHAT IS USEFUL FOR
- \* 4. THE CALCULUS OF PROBABILITY
5. MORE MODELS FOR MORE DATA
- \* 6. JOINT DISTRIBUTIONS, INDEPENDENCE,  
MAXIMUM LIKELIHOOD
- \* 7. EXPECTATION
- \* 8. CONCENTRATION, INEQUALITIES, LIMIT THEOREM
9. STATISTICAL LEARNING

## REFERENCES

- LECTURE NOTES
- ART OF STATISTICS (BOOK)
- VARIOUS ONLINE SOURCES (DATA)
- WIKI

## APPROACH

- START FROM DATA & QUESTIONS
- ONLY TEACH USEFUL THINGS
- COLOUR CODING
  - FUNDAMENTAL RESULT IN STATES
  - FUNDAMENTAL RESULT IN MATH

## EVALUATION

- 1 / 2 / 3 LEVELS
- SEE BLACKBOARD

## THEME 1 : DATA VIZ & SUMMARIZATION

- SHIPMAN'S DEAD PATIENTS
- VOX POPULI
- BAR & BOX PLOT
- MEANS & QUANTILES
- SCATTER PLOT & CORRELATION

## SIMPLEST TYPE OF DATA : BINARY

BINARY

Event	Percentage in 10,267 people allocated placebo	Percentage in 10,269 people allocated statin	% (relative) risk reduction in those allocated statins
Heart attack	11.8	8.7	27%
Stroke	5.7	4.3	25%
Death from any cause	14.7	12.9	13%

→ SUMMARY OF  
BINARY DATA :  
FREQUENCY  
AoS

Table 4.1

The outcomes at five years in the Heart Protection Study, according to treatments allocated to patients. The absolute reduction in the risk of a heart attack was  $11.8 - 8.7 = 3.1\%$ . So out of 1,000 people taking a statin, around 31 heart attacks were prevented – this means that around 30 people had to take a statin for five years to prevent one heart attack.

## SUMMARIZATION (AT THIS POINT) FOR BINARY

- FREQUENCY :

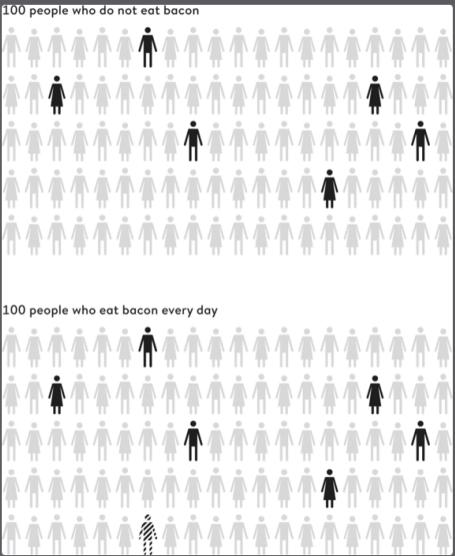
NUMBER OF OBSERVATIONS IN EACH CATEGORY

- PERCENTAGE = FREQUENCY / TOTAL

# VISUALIZATIONS OF SUMMARIES OF BINARY DATA

---

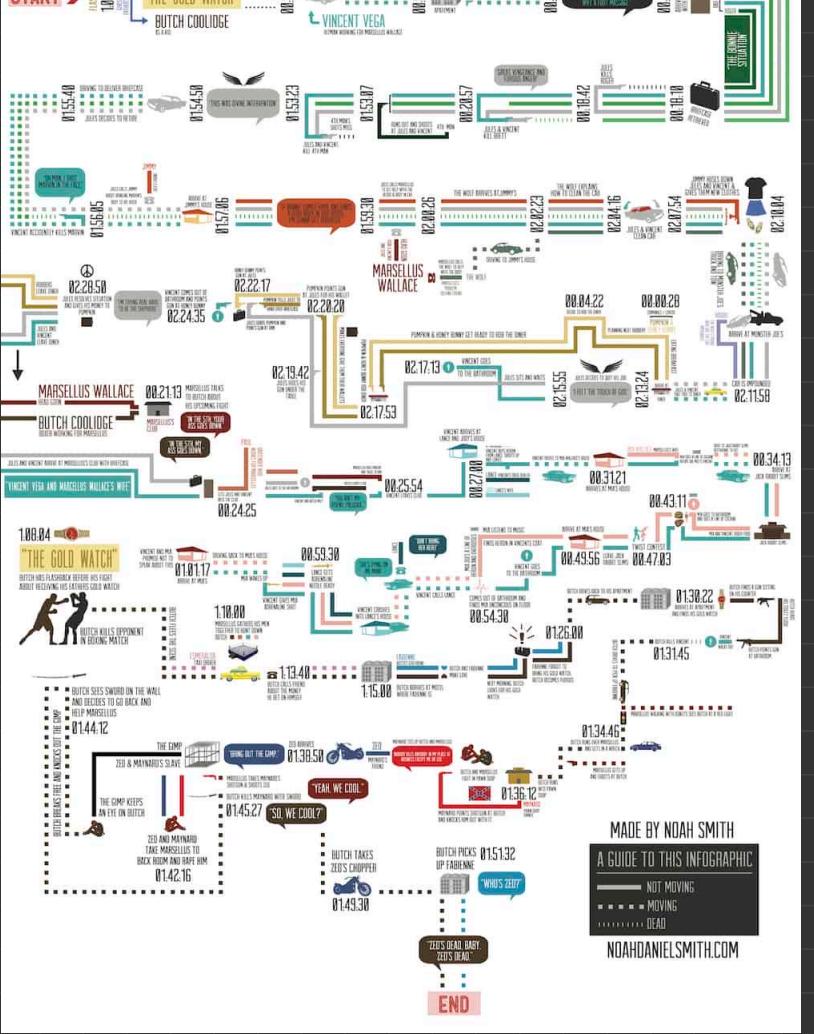
- ALREADY SEEN: TABLE OF FREQ's / PERC's
- INFOGRAMS



**Figure 1.4**  
Bacon sandwich example using a pair of icon arrays, with randomly scattered icons showing the incremental risk of eating bacon every day. Of 100 people who do not eat bacon, 6 (solid icons) develop bowel cancer in the normal run of events. Of 100 people who eat bacon every day of their lives, there is 1 additional (striped) case.<sup>fn5</sup>

Ao S

TWO BINARY JAN'S :



## CATEGORICAL DATA

Event	Percentage in 10,267 people allocated placebo	Percentage in 10,269 people allocated statin	% (relative) risk reduction in those allocated statins
Heart attack	11.8	8.7	27%
Stroke	5.7	4.3	25%
Death from any cause	14.7	12.9	13%

Table 4.1

The outcomes at five years in the Heart Protection Study, according to treatments allocated to patients. The absolute reduction in the risk of a heart attack was  $11.8 - 8.7 = 3.1\%$ . So out of 1,000 people taking a statin, around 31 heart attacks were prevented – this means that around 30 people had to take a statin for five years to prevent one heart attack.

BINARY IS  
CATEGORICAL WITH  
TWO LEVELS  
HERE ALSO SEE  
CATEGORICAL WITH  
4 LEVELS