



INTRO TO DATA SCIENCE





LAW → GOV. & ETHICS
↳ BIAS

STRUCTURE
DATA STORAGE

BIG DATA

PROCESSING DATA

IoT / Events

MACHINE LEARNING

AI

NLP
Comp. Vision

MODELING

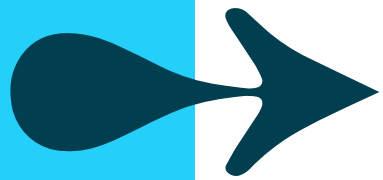
MATHEMATICS
↳ STATISTICS
↳ PROBABILITY

INTERPRETATION
& MEANING

VIZUALING

EXPERIMENTATION
RESEARCH.

DATA SCIENCE



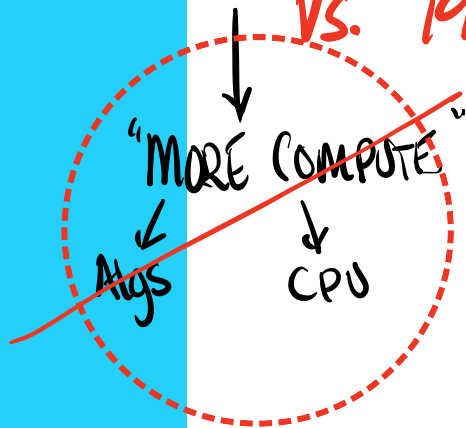


⊛ WHY NOW? → MORE DATA

VS. 1990, 2005... ?

↓
WHERE?

↪ INTERNET
↪ USER BEHAVIOUR (People!)



⊛ WHAT'S DIFFERENT ?

↳ VS. DATA ANALYSIS?

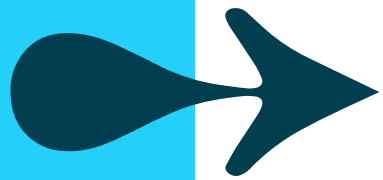
→ SCIENCE (EXPERIMENT, EXPLANATION)

→ DATA ANALYSIS

↳ HISTORICAL ("CERTAIN")

Qx.
ANSW.

↓ E.g. mean(profit₁₉₉₀) = £1m.





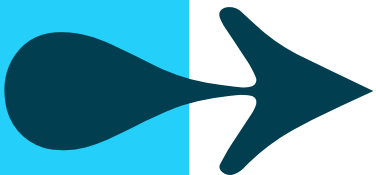
^{All}
DATA SCIENCE

QUESTIONS ARE

PROBABALISTIC

Eg. What profit is most likely
M 2030 - AND -

How confident are we
in that profit?





PROBABILITY

Confidence \rightarrow P (in claim) $\left(\begin{array}{c|c} \text{left} & \text{Assuming (for certain)} \\ \text{C} & \text{A} \end{array} \right)$ Right

$P(\text{"It will rain tomorrow"} | \text{"It has rained today"})$



VARIATIONS IN NOTATION

$$P(Y|X)$$



A Claim about Y



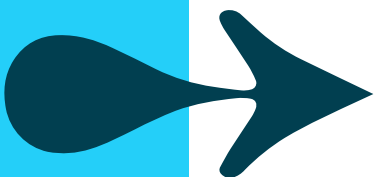
Assumption about X

Y - "amount of rain"
^{Future}

X - "amount of rain"
^{PAST}

$$P(Y) \leftarrow \text{RHS is missing}$$

"ASSUMING IMPLIED
BACKGROUND 'knowledge'"





Eg.

$P(H) \leftarrow$ doesn't make sense...

\downarrow

$$P(H | \{H, T\}) = \frac{1}{2}$$

CLEARER

Eg.

$$P(\text{Fraud} | \text{Age, Amount})$$

$\left| \begin{matrix} 1 \\ 0 \end{matrix} \right.$

 Question

 Assumptions



DATA SCIENCE EXAMPLES

INDUSTRIES & DATASETS

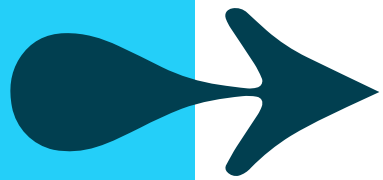
↳ HEALTHCARE - SCREENING

FINANCE - STOCKS

UTILITIES - POWER SUPPLY

INSURANCE - FRAUD

RETAIL - PRODUCTS





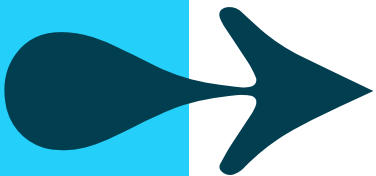
RETAIL - PRODUCTS

2. "HOW ^{← QUANTITATIVE?} WELL WILL A PRODUCT SELL?"
2. y (question)

BASED ON "LOCATION IN STORE" ^X
(assumption)

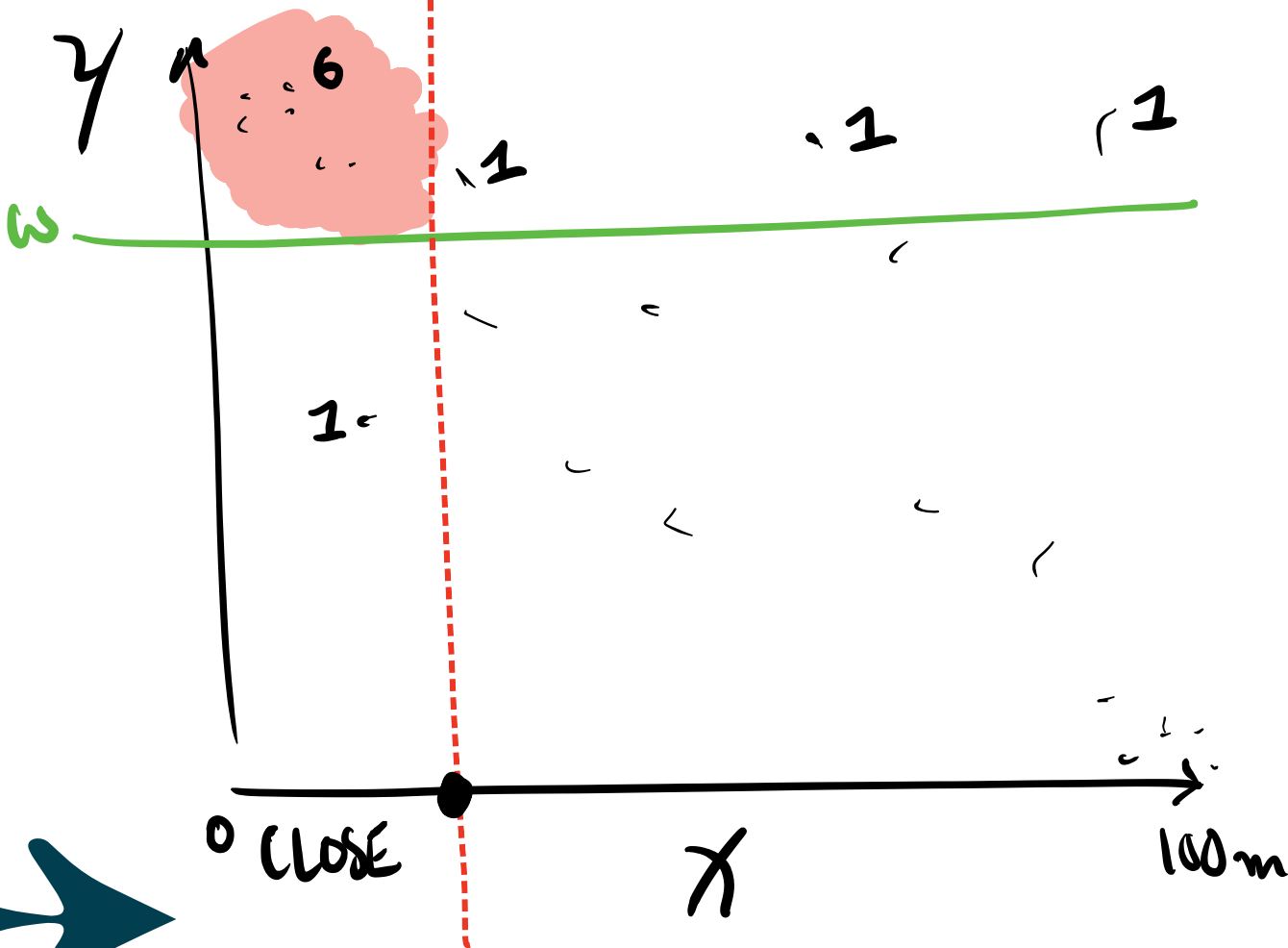
WELL def. $\text{N}^\circ \text{ units / month} \geq w$ ^{← cutoff}

$$P(y \geq w | X)$$

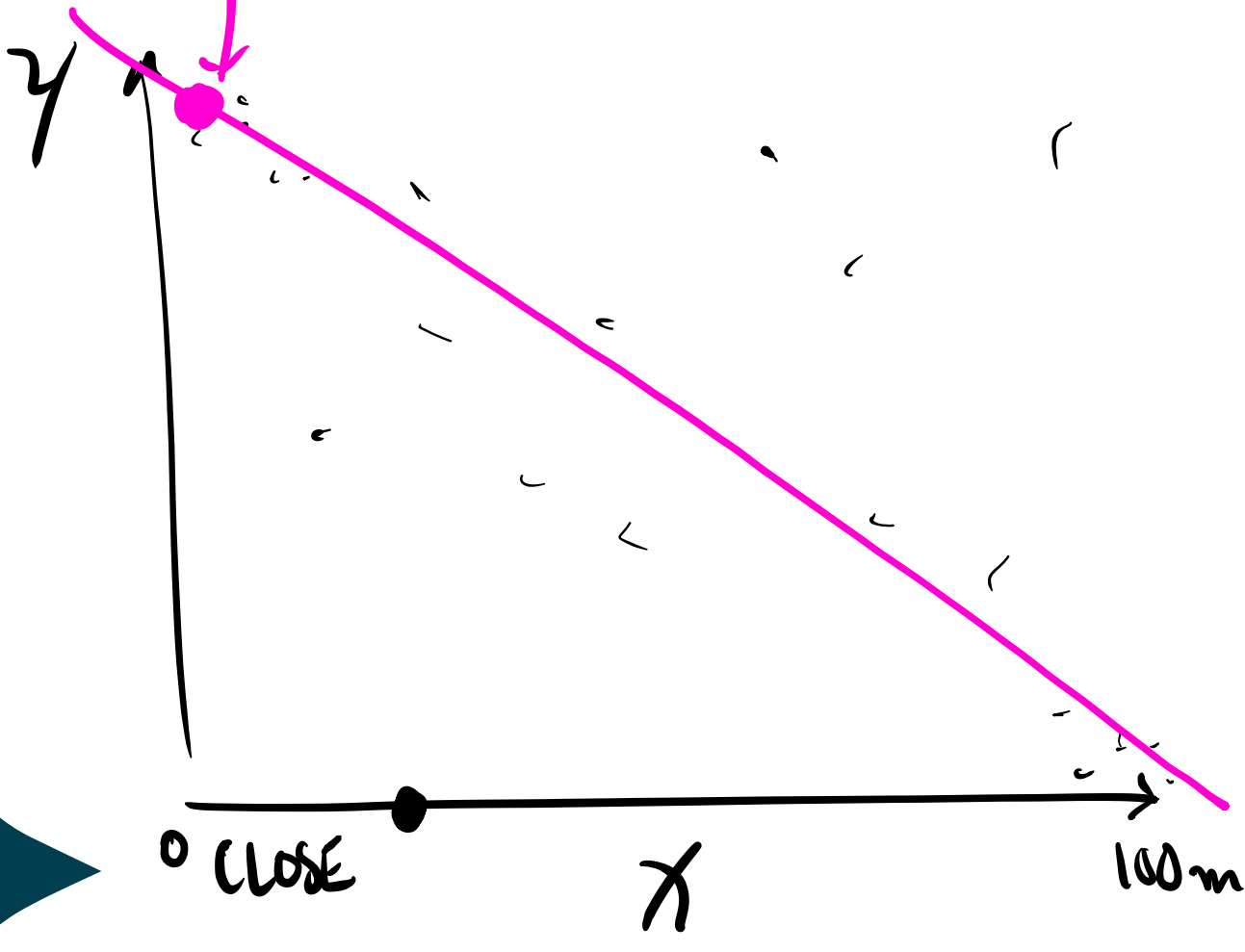




$$P(y \geq w | X)$$



$$\text{MAX } P(Y \geq W | X)$$





ASIDE: EG. Probability Rule

"Emphasizing the significant connections"

$$P(C|E) = \frac{P(E|C) P(C)}{P(E)}$$

↓ "MORE ACCURATE"

$$P(C|E, B) = \frac{P(E|C, B) P(C|B)}{P(E|B)}$$

