

# Data Science by Example

## Problems & Solutions

# Statistics & Machine learning

The Risk of  
Return

# The Problem

---

Argonaut Sell : 100 Toys / yr

Profit / Toy : €10.00 given ~~tx~~ Costs

Sale Price : €100.00

costs [ Delivery : €5.00 ← Paid 2x on Return & Resale  
Admin Costs : €5.00 fixed

Problem: Q8% toys returned ( $r$ )

# Impact

$$\text{Max. Profit: } 100 * \text{£}10 \\ = \text{£}1000$$

$$\text{Actual: } (1 - r) * \text{£}10 = \text{£}800$$

0.2

Goal: Minimize  $r$

# How can we affect $\tau$ ?

- ↳ Modify Prices?
- ↳ Reduce Return Window?

---

↳ Let's Explore...

# Algorithm Proposal

If  $P(\text{Return} | \text{Customer Info}) > \text{Risk}$  :

Return Window := 7 DAYS

else :

Return Window := 30 DAYS

$$P(\text{Return} \mid \underline{x_0 \dots x_n})$$

...?

- \* Age  $\in [0, 100]$
- \* Location  $\in \{\text{UK, FR, ...}\}$
- \* History of Returns  $\xrightarrow{\text{AVG./DAY}}$
- \* Freq. of Visits  $\in [0, 100]$
- \* Gift  $\in \{0, 1\}$   
? Is it a gift?

P( ... )  $\geq$  Risk

At Risk = 0.5,

At Risk = 0.4,

Controls Rates  
Of Mispredictions

X	0.1	X	0.3	✓	0.4	X	0.2	✓	0.7	X	0.6
X	X	X	X	X	X	X	X	✓	X	✓	X
X	X	X	✓	X	X	✓	X	✓	X	✓	✓

Answers:

WAS: False  
-ive

Now: NO ERROR

False  
+ive

How do we find  $P(\text{Return} | X \dots)$ ?

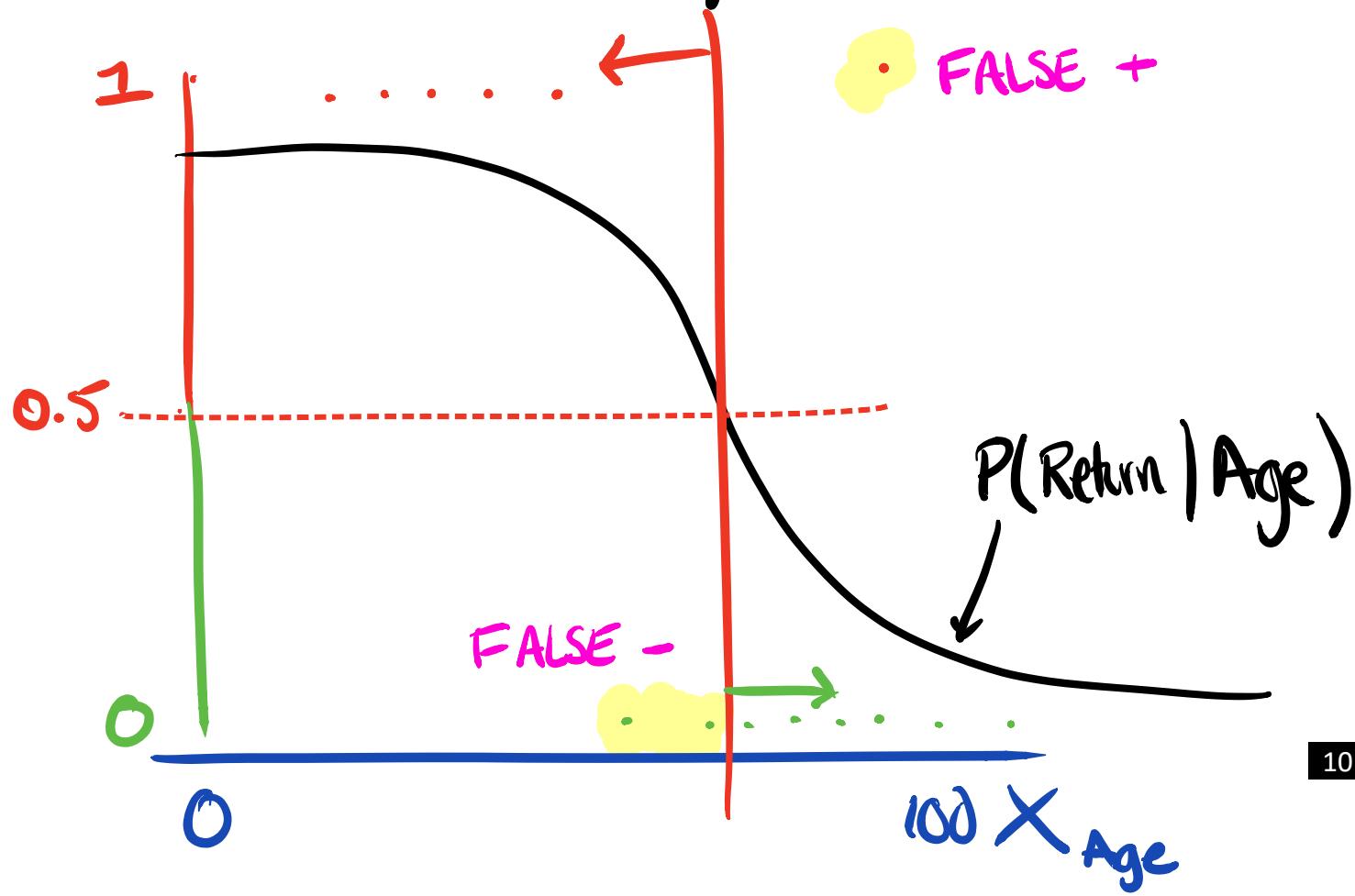
Experiment

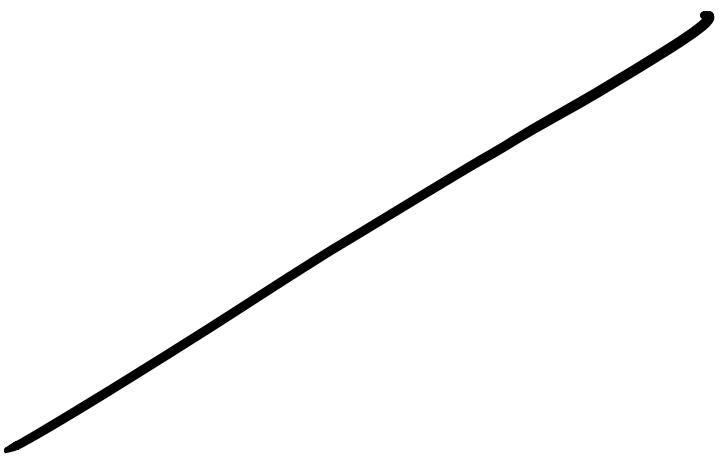
Age	F <sub>q</sub> Return	F <sub>q</sub> Visit	Gift?	...	$\gamma_{\text{Returned}}$
30	0.01	1.2	0	...	0
21	0.2	3.2	0	...	1
...	...	...	...	...	...

**PREDICTIVE FEATURES**

**RESPONSE TARGET**

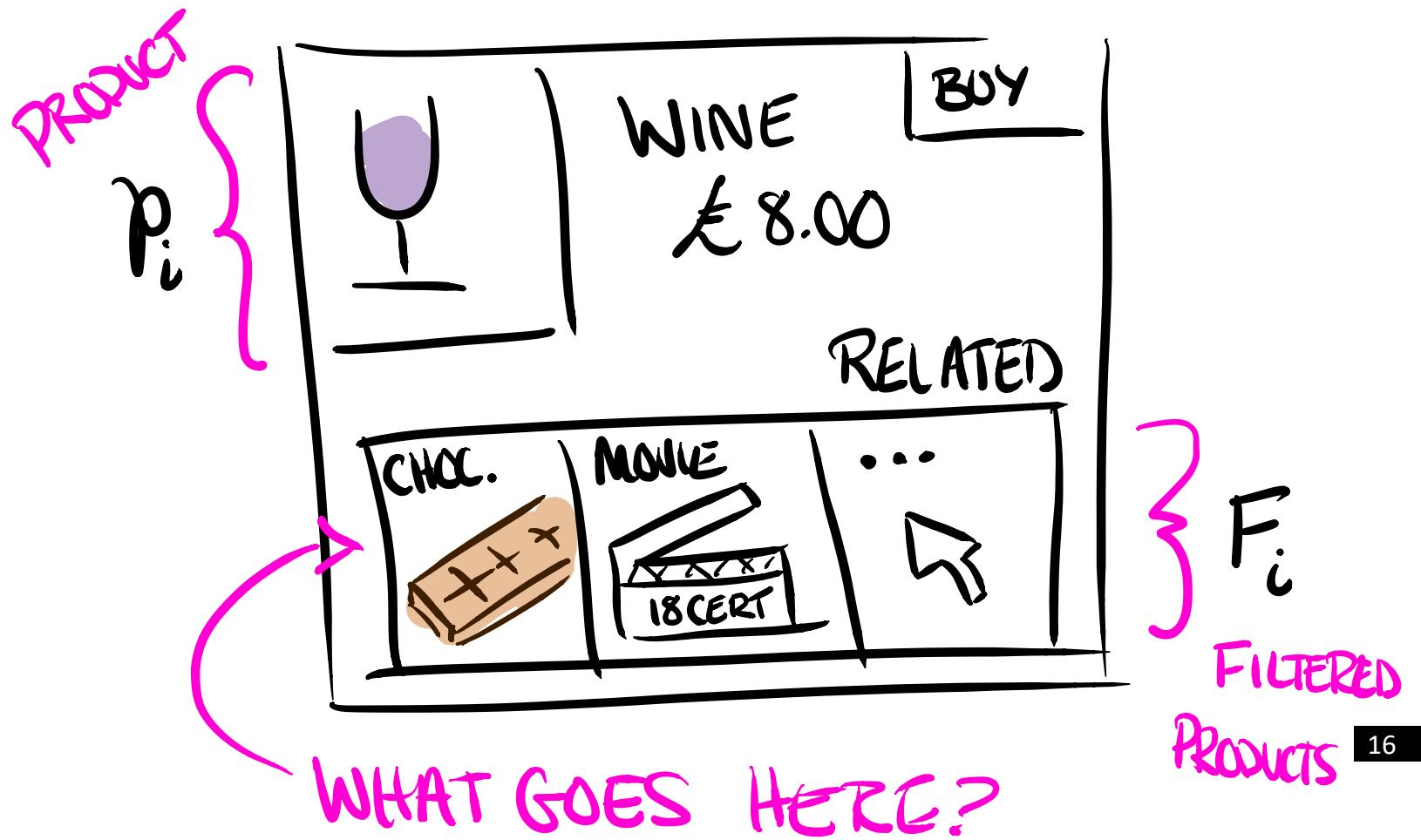
# Find a Relationship





# Algorithms & Graph Theory

Product  
Recommendation



# The Problem

Sales Team sets

filtered set of  
Related products

$$F_i \subset P$$

total set

Products in  
are scored,

$$p_i \underset{\substack{\text{in} \\ \uparrow \text{product}}}{} \in F_i \underset{\uparrow \text{id}}{}$$

Score ( $p_i$ ) = avg. Review

Eg.

$$P = [ \overset{0}{\text{chocolate}}, \overset{1}{\text{Wine}}, \overset{2}{\text{Carrots}} \dots ]$$

$$F_0 = \{ \text{Wine}, \text{Movie}, \dots \}$$

Related Products

$$F_1 = \{ \text{chocolate}, \text{Movie}, \dots \}$$

$$F_2 = \{ \text{Bread}, \text{Parsnips}, \dots \}$$

## Eg. Scoring

$F_1 = \{ \text{chocolate, Movie, ...} \}$

Score (chocolate) = 4.5 #1

Score (movie) = 3.8 #2

∴

When visiting "Wine"

Recommend

# How Can we improve F?

USE PROFILE:

Prob (Purchase | Age, Gender, ...)  
Prob. Customer purchases  
given their age, etc.

USE "SIMILARITY":

SET  $P :=$  PRODUCTS OF SIMILAR  
USERS

LIKED  
How 'alike' is customer to other purchasers?

# Algorithm:

- FOR user  $u$  LOOKING AT  $p_i$  How?
- \* FIND SIMILAR USERS  $S_u$
  - \* Set  $P :=$  PRODUCTS REVIEWED BY  $S_u$
  - \* Set  $F_i :=$  PRODUCTS OF REVIEWERS  
OF  $p_i$  IN  $P$
  - \* SCORE  $F_i$  BY REVIEWS  
& RECOMMEND.

Eg.  $P_i = \text{"WINE"}, u = \text{"Me"}$

PRODUCTS  
RATED  
 $> 4*$

$S_u = \{ \text{Alice}, \text{Bob}, \dots \}$  SIMILAR PEOPLE

$P = \{ \text{Wine}, \text{Cheese}, \text{Bread} \dots \}$

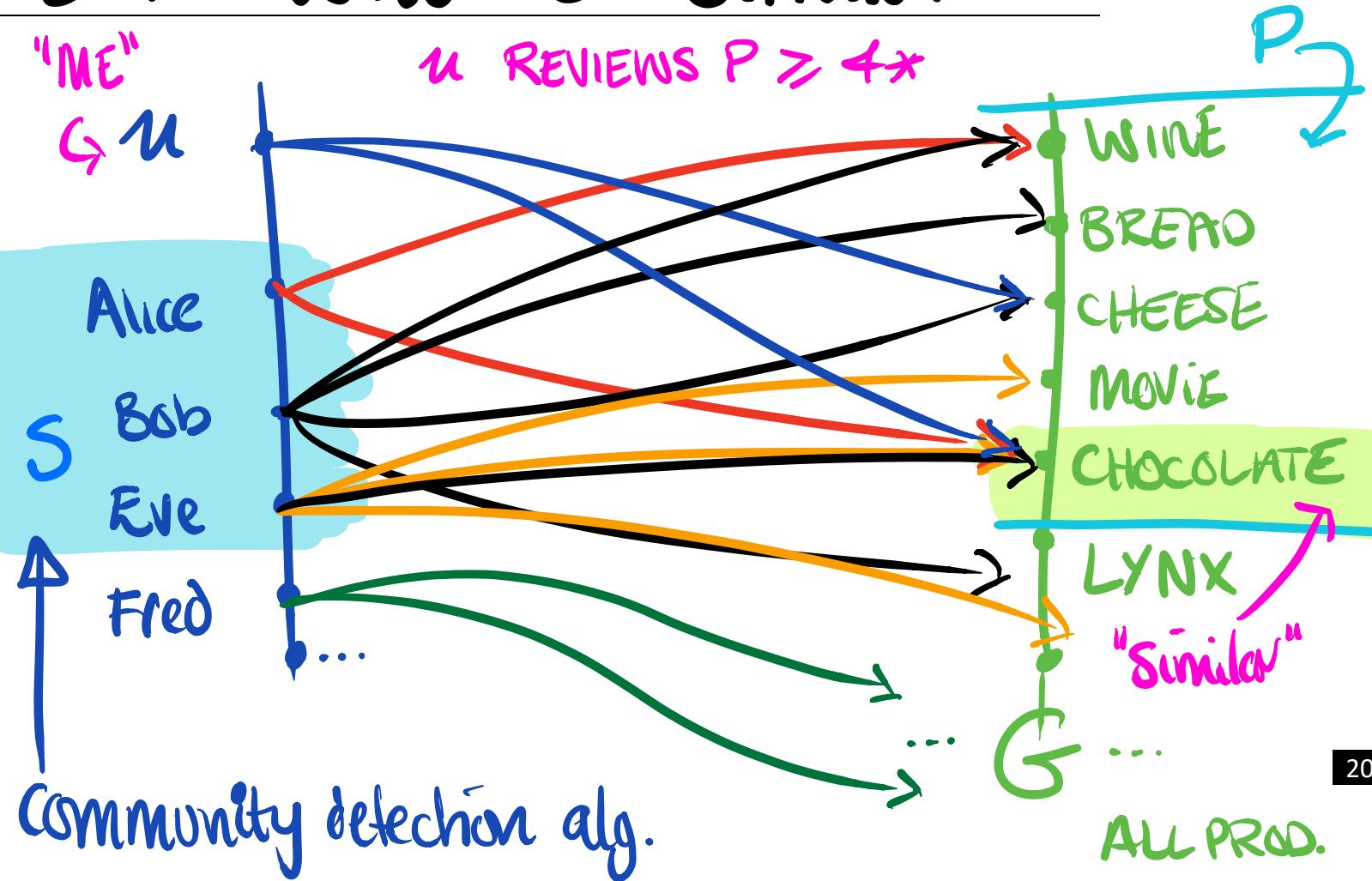
OFTEN REVIEWED } BY SAME PEOPLE

score ↙  
 $F_i = \{ \text{Cheese}, \text{Bread} \}$

[ #2 ↘ ↗ #1  
[ 4.3, 4.7 ... ] ]

RECOMMEND! ↴

# But... Who is 'similar'?



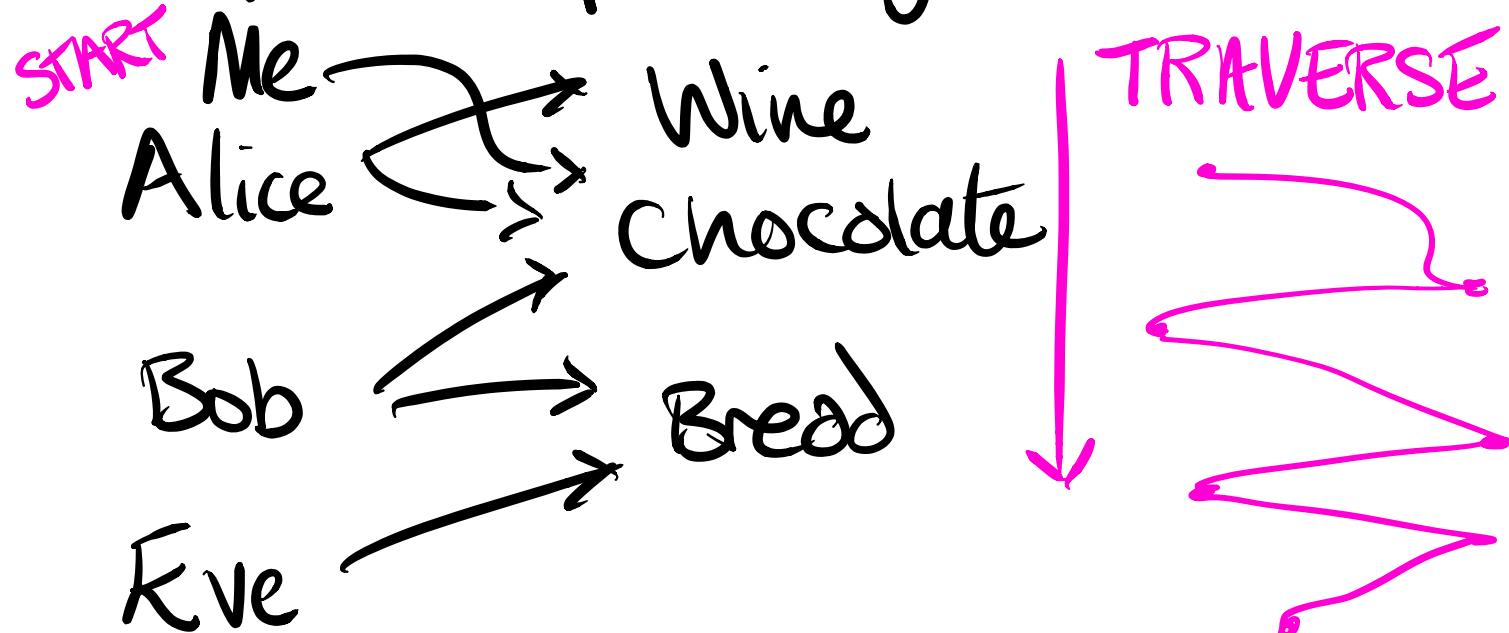
# Community Detection

A 'similar' to B  
by  $J = \frac{\# \text{Shared Reviews}}{\# \text{Reviews}}$

$\{A, B, C, \dots\}$  is a community S

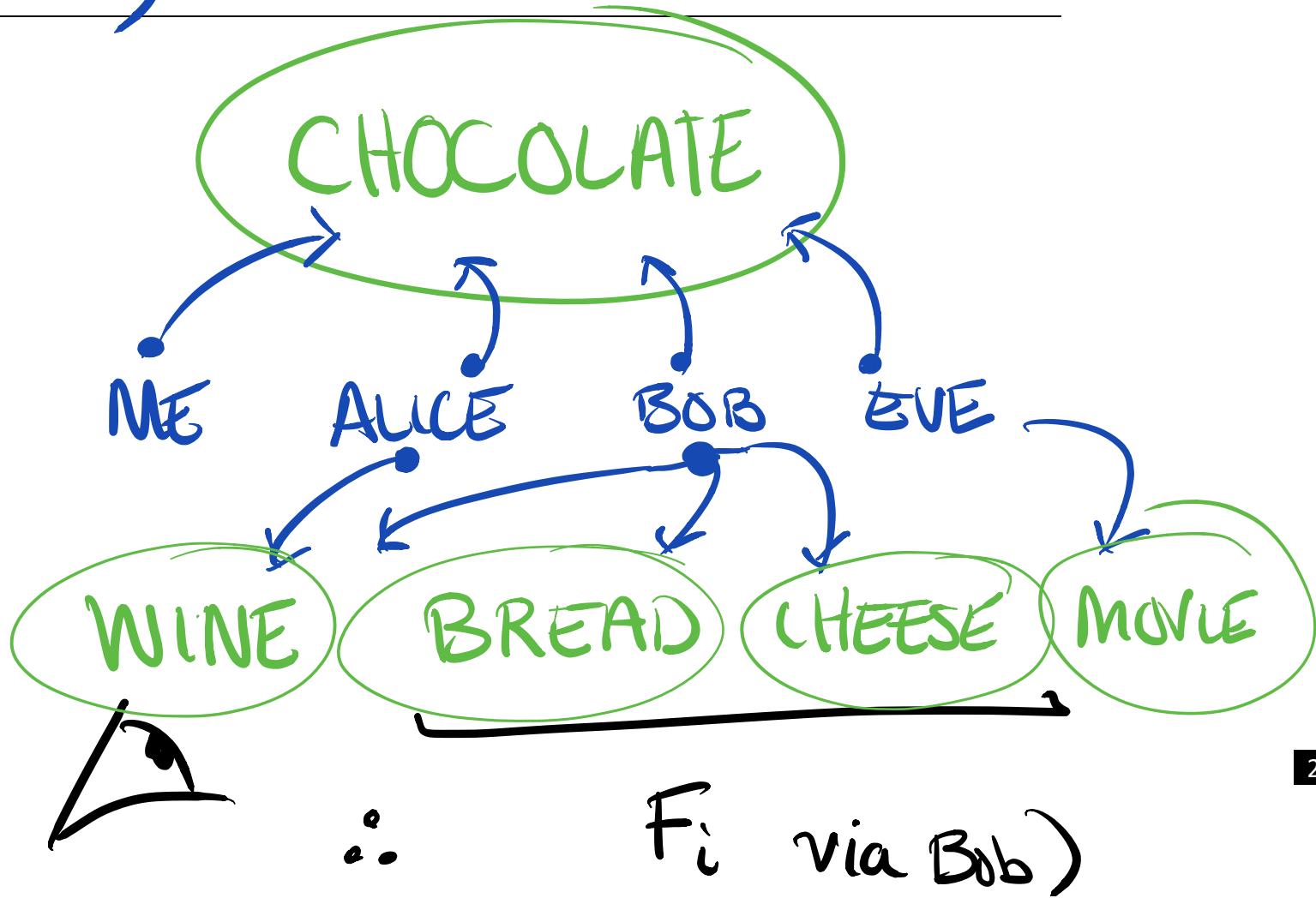
if  $J(S) > \text{minimum}$

# Simple Graph Alg.



Add user  $u$  to  $S$  if  
 $\delta(s) \geqslant \text{minimum}$

So,





# Big Data

---

User Behavior &  
Event Systems

# The Problem

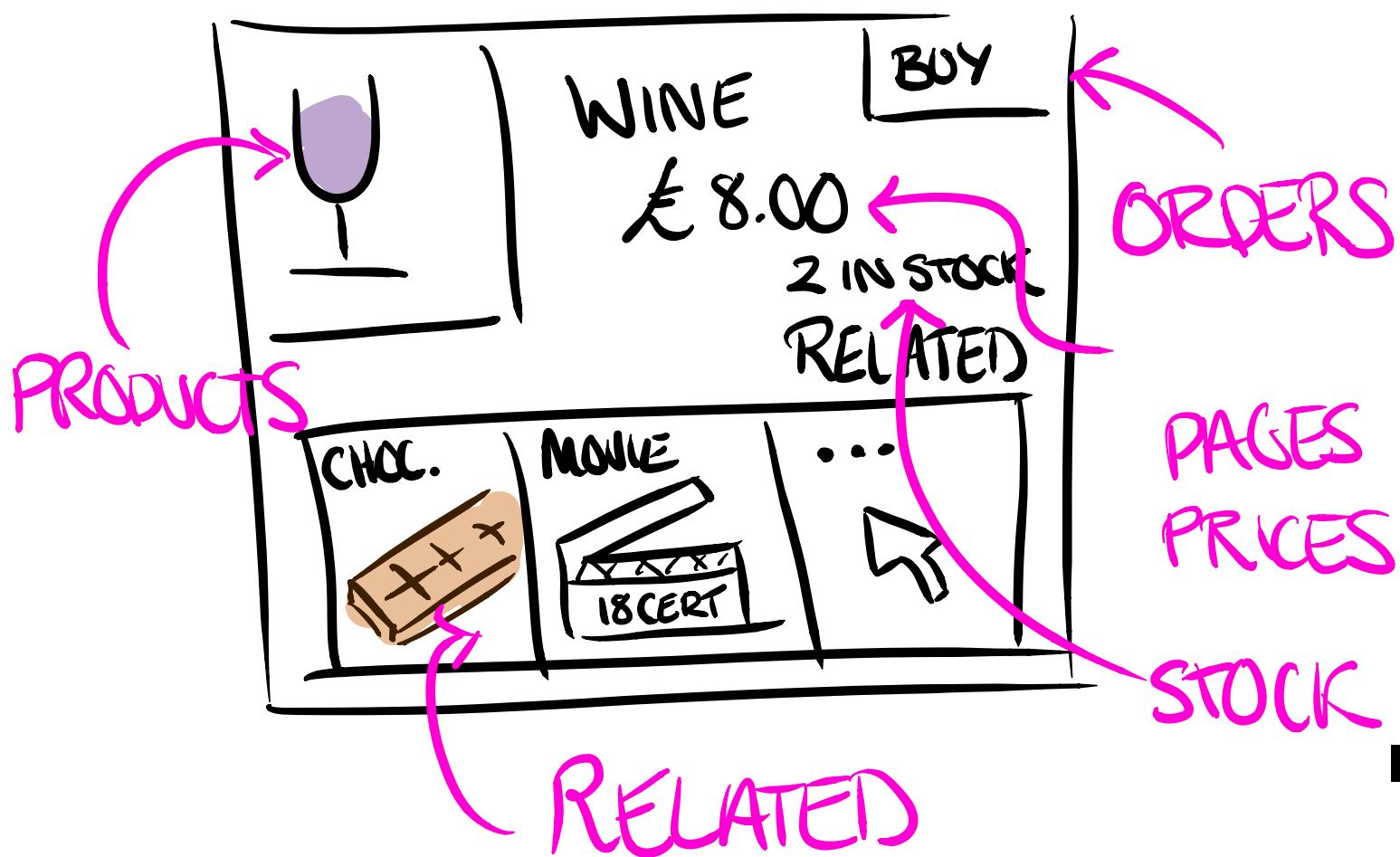
Existing data infrastructure  
does not track

“analytical data”

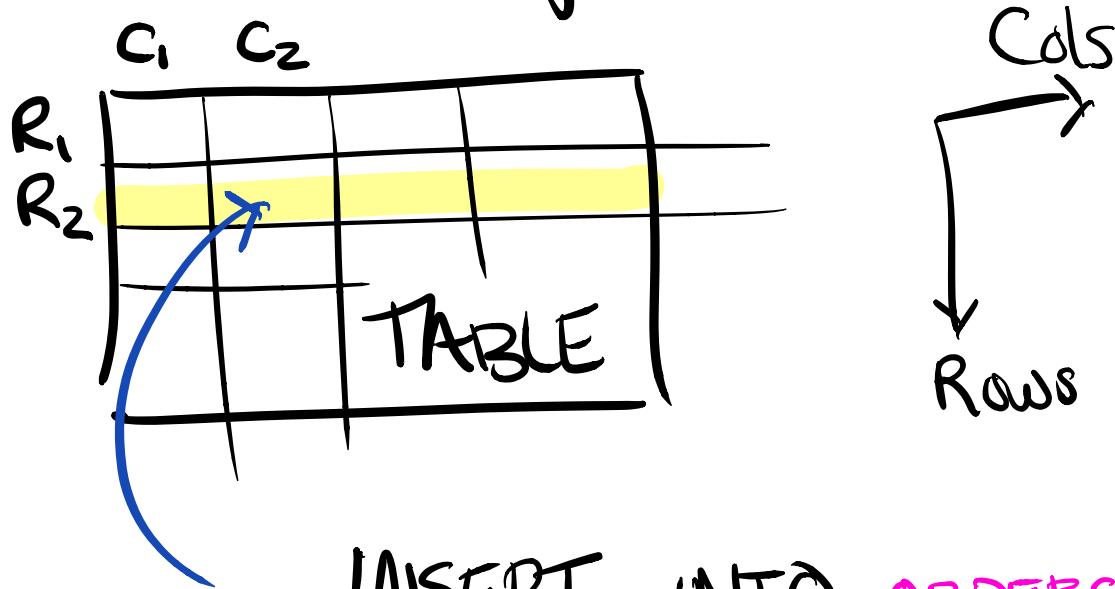
Why?

---

Data infrastructure  
was designed for  
Retail Transactions



# The Existing System



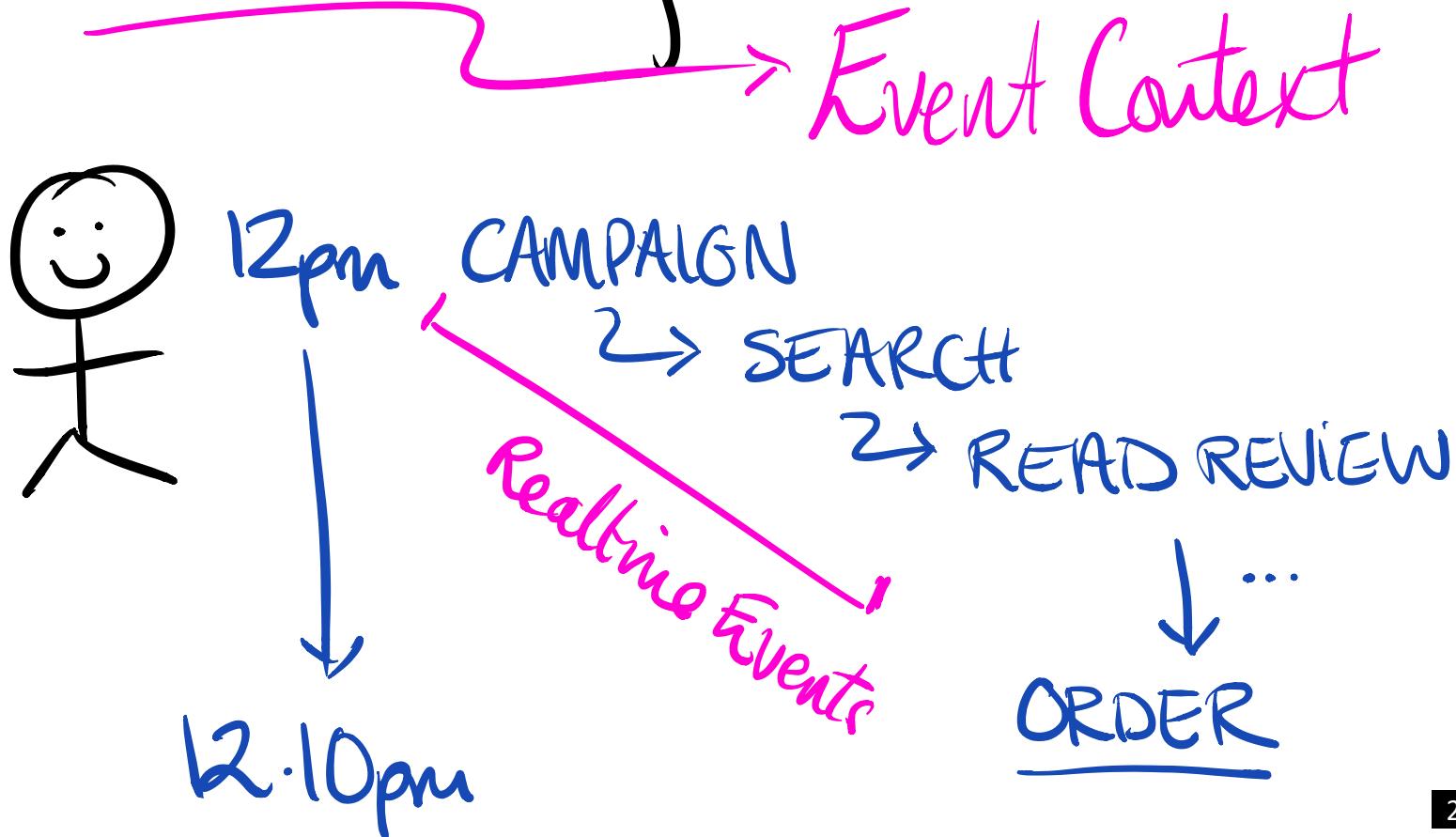
cols  
↓  
rows

INSERT INTO ORDERS VALUES

Row =  
(ID, PRODUCT, QUANTITY, ...)

Essential Retail Information

# What's Missing?



# How do we store 'Events'?

Event =

(Subject, Verb, Object, Context)

→ CUSTOMER CLICKS AD ID: 35  
TIME: 12pm

# Event log

---

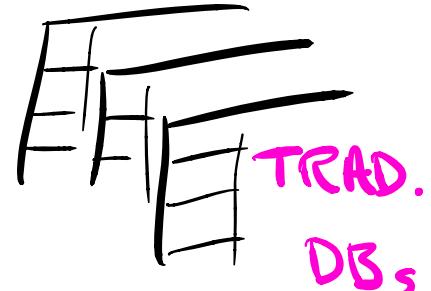
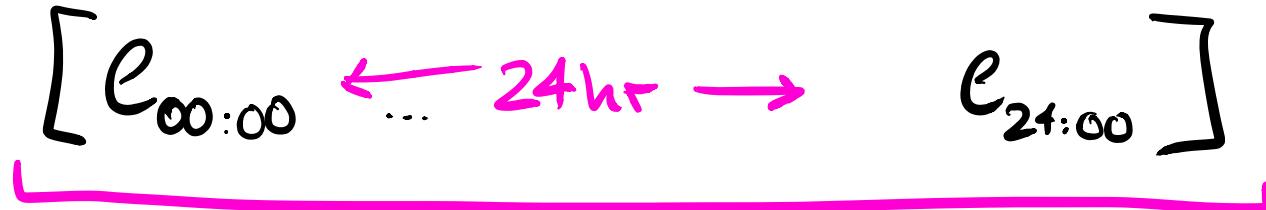
log = [e<sub>1</sub>, e<sub>2</sub>, e<sub>3</sub> ... e<sub>N</sub>]

e <sub>1</sub> ...	CLICK	...	
e <sub>2</sub> ...	READ	...	
e <sub>3</sub> ...	ADD	...	
e <sub>4</sub> ...	ORDER	...	
e <sub>5</sub> ...	COMPARE	...	↓ true

# The New Data System

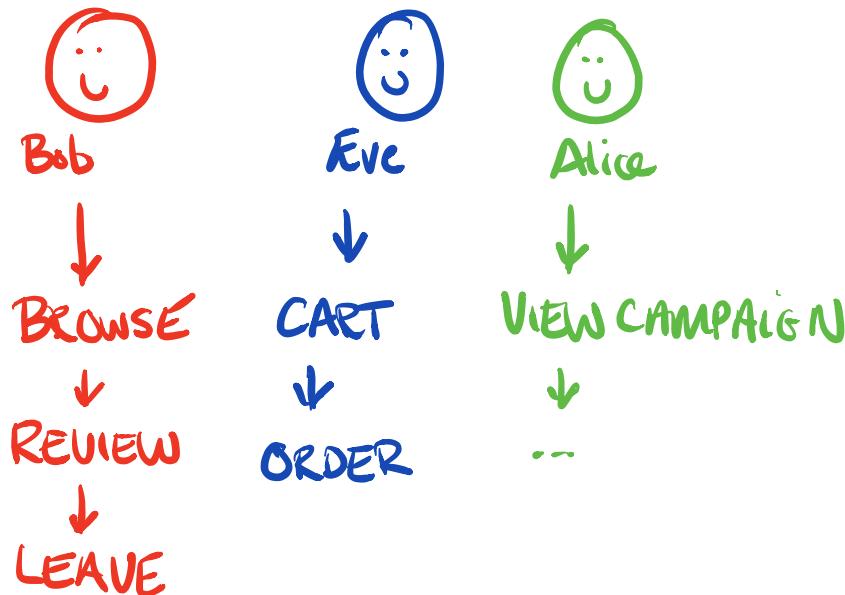
append-only log

CLEARED DAILY



# Fig.

---



$[e_1, e_2, e_3, e_4, e_5, e_6, e_7, \dots]$

Fig.

