

Review 1: Entity Resolution, Blocking, Precision and Recall

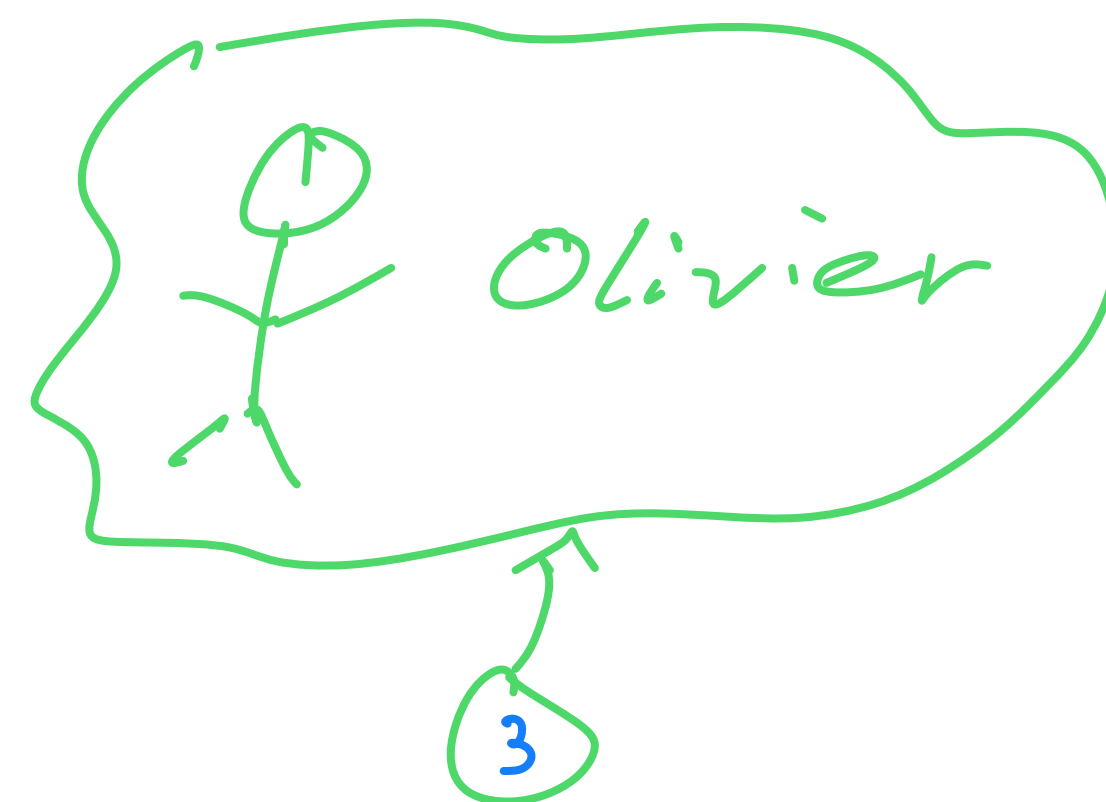
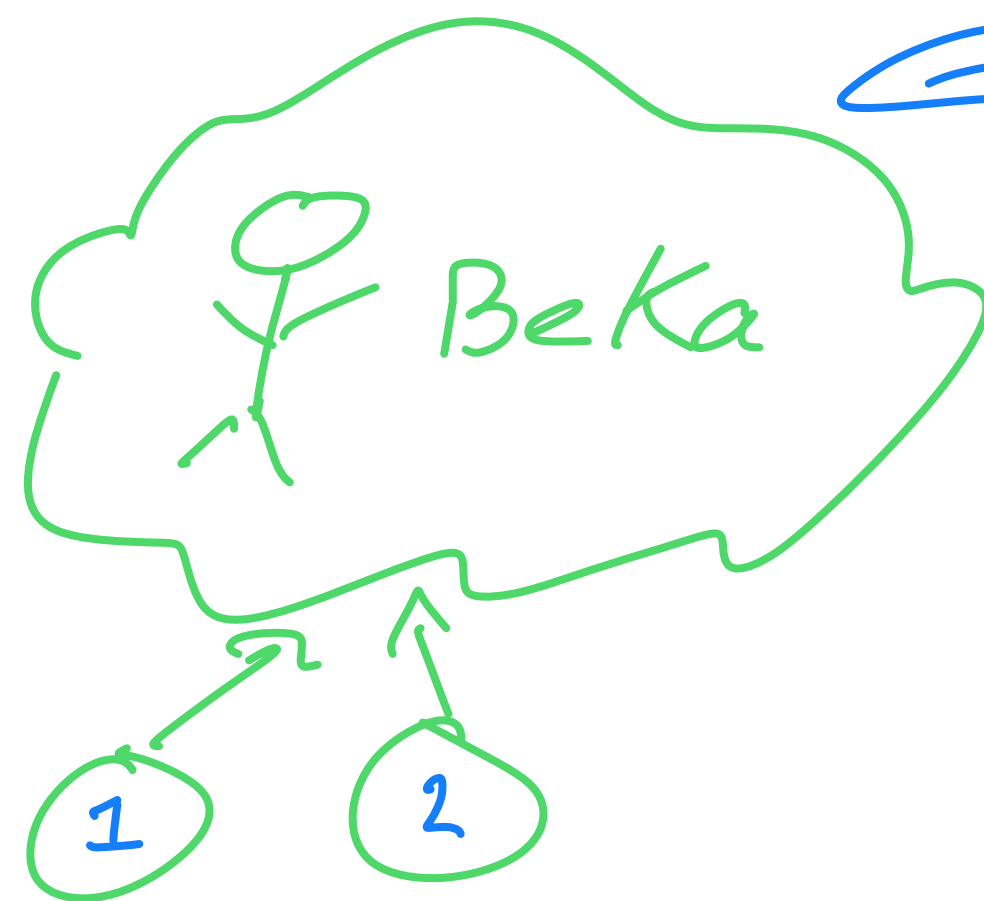
Olivier Binette

Entity Resolution

Databases, Records and Attributes

Database A

	<u>Name</u>	<u>b.y.</u>	<u>occupation</u>	} attributes records
1.	Beka S.	—	Professor	}
2.	Rebecca C.S.	—	Statistician	
3.	Binette, O.	95	PhD Student	



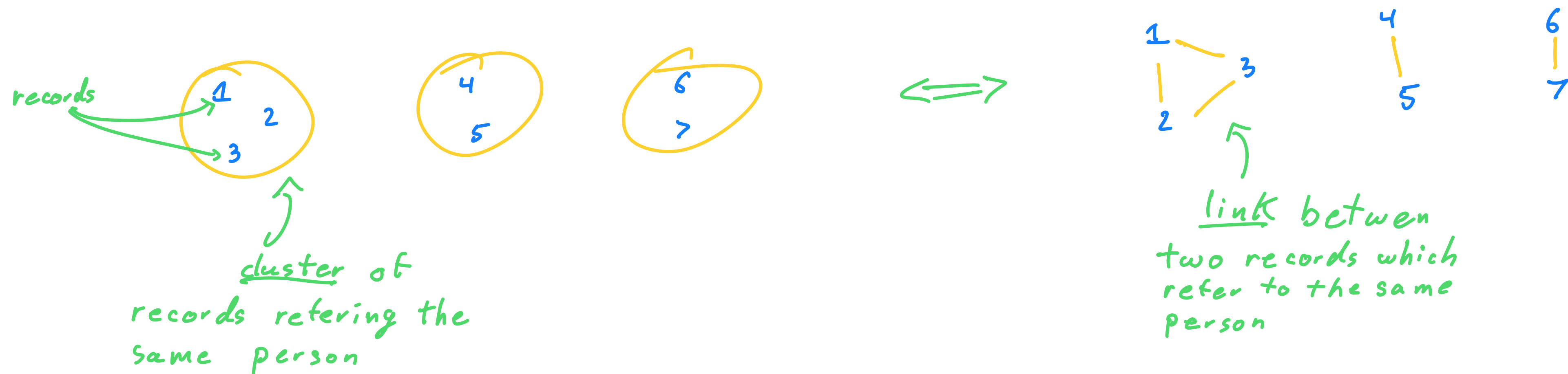
(latent) entities

Entity Resolution

Two records are said to be coreferent or to match if they refer to the same entity.

Goal of Entity Resolution:

Cluster coreferent records \Leftrightarrow identify matching pairs



Entity Resolution

Typical approach to ER:

1. Consider all pairs of records

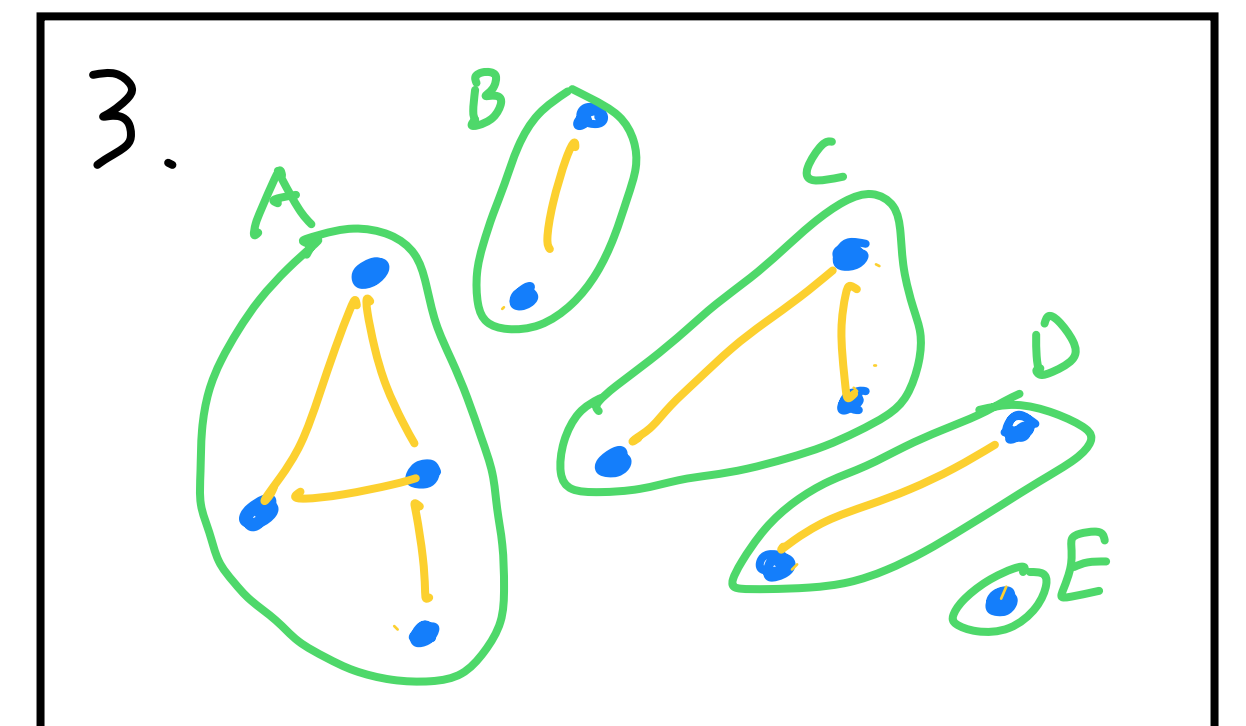
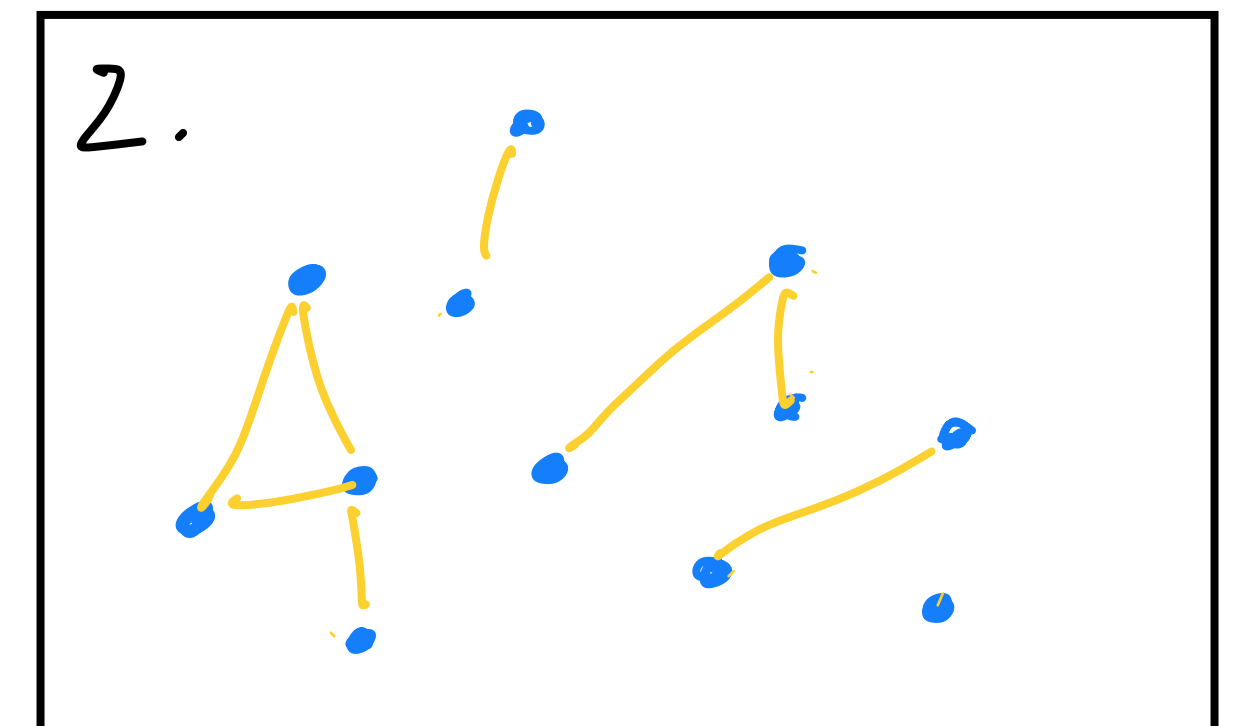
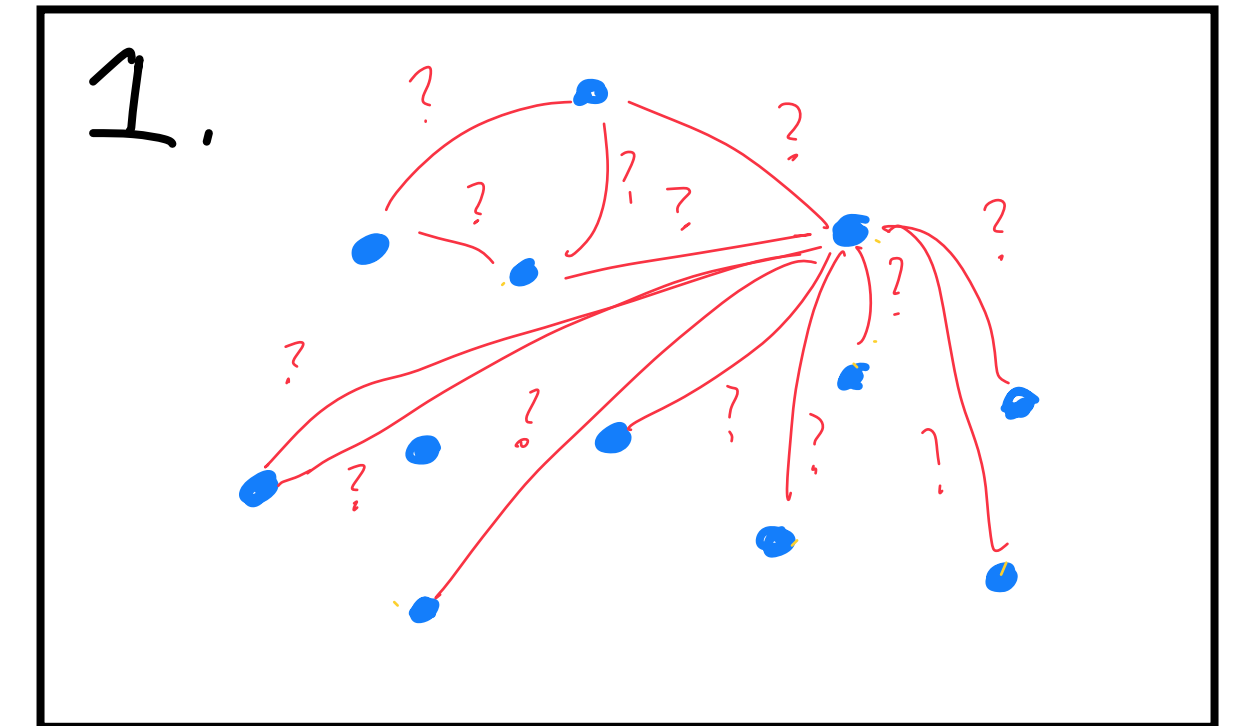
↳ The $\binom{N}{2}$ pairs are the comparison space

2. Classify each pair as being a match (link)
or a non-match (don't link)

↳ This gives a graph

3. Assign a unique entity ID To each
record based on the resulting structure

↳ E.g. use connected components of the graph



Entity Resolution

Problem:

The comparison space is large ($\mathcal{O}(N^2)$)

↳ Many pairs to consider!

Solution:

Blocking...

Blocking

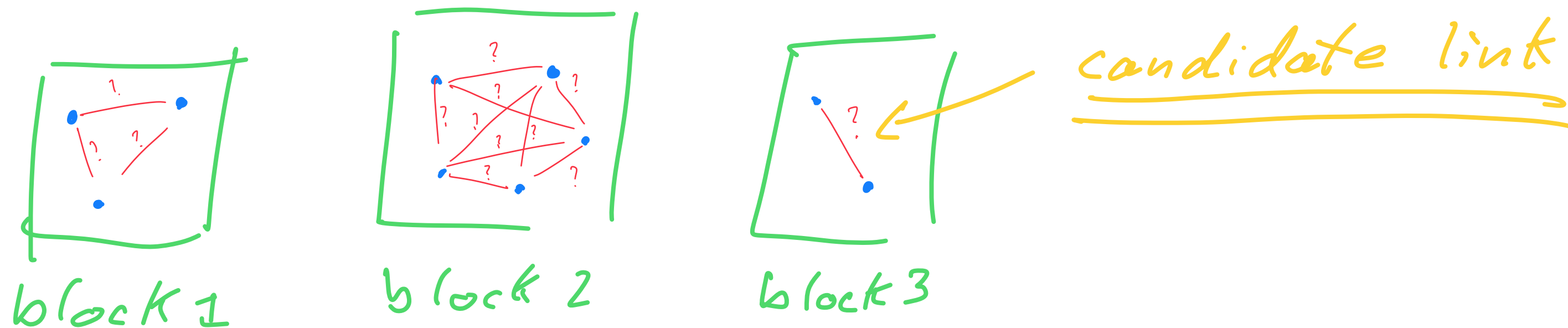
Blocking

Goal of blocking:

- Quickly ($< O(N^2)$) reduce the comparison space

Trick:

- Place each record into a "block"
 - ↳ You want matching records in the same block
- Only make comparisons within blocks
 - ↳ Here we focus on blocks that partition the set of records



Blocking

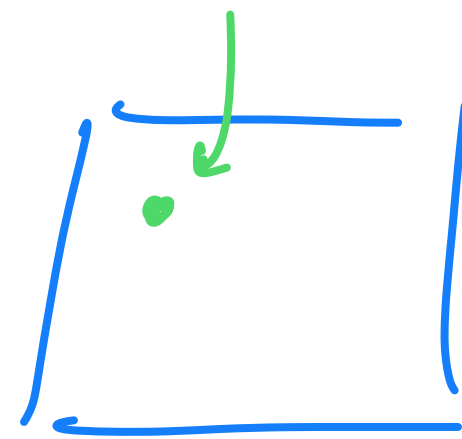
Example:

- Block by last name initial
 - ↳ Each record is placed in a block given the initial of the listed last name.



block A

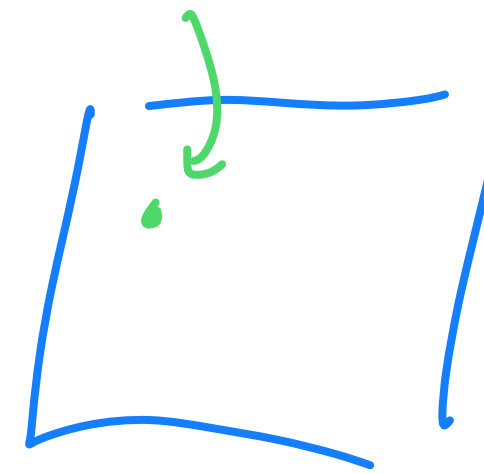
"Olivier Binette"



block B

...

"Olivier Zinette"



block Z

Blocking

The challenge of blocking:

Balance the tradeoff between:

1. Small blocks (fewer comparisons & more efficient)
2. Few mistakes (don't place matching records in different blocks)

How to quantify this tradeoff?

1. Reduction ratio
2. Recall

Reduction ratio

Suppose you have N records.

- Without blocking, you need $\binom{N}{2}$ comparisons
- With k blocks, each of size b_1, b_2, \dots, b_k , you need $\sum_{i=1}^k \binom{b_i}{2}$ comparisons

Reduction ratio:

$$RR = \frac{\binom{N}{2} - \sum_{i=1}^k \binom{b_i}{2}}{\binom{N}{2}} =$$

"% of comparisons that blocking has eliminated"

Recall

Recall = "% of matching pairs that are in the same block"

low recall \Rightarrow many mistakes will necessarily be made

perfect recall (100%) \Rightarrow no accuracy loss in doing blocking.

Precision and Recall

(more generally)

Precision and Recall

Consider a test (e.g. covid test)

↳ Result is positive (P) or negative (N)

→ Patient is sick (C) or not sick (\bar{C})

	C	\bar{C}
P	TP	FP
N	FN	TN

$$\text{precision} = P(C|P) = \frac{P(C \cap P)}{P(P)} = \frac{TP}{TP + FP}$$

" = % of positive patients which are actually sick "

$$\text{recall} = P(P|C) = \frac{P(P \cap C)}{P(C)} = \frac{TP}{TP + FN}$$

" = % of sick patients which have a positive result "

Application to blocking

- We have record pairs
- Each is a candidate pair (P) or not (F)
↳ in the same block
- Each is actually a match (M) or not (\bar{M})

	M	\bar{M}
P	TP	FP
F	FN	TN

$$\text{recall} = \frac{TP}{TP + FN}$$

$$= \frac{\# P \cap M}{\# M}$$

$$= \text{"\% of matches which are candidate pairs"}$$

Computing Recall

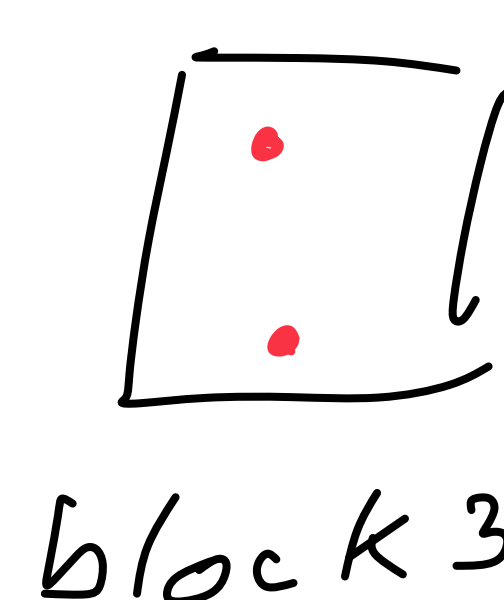
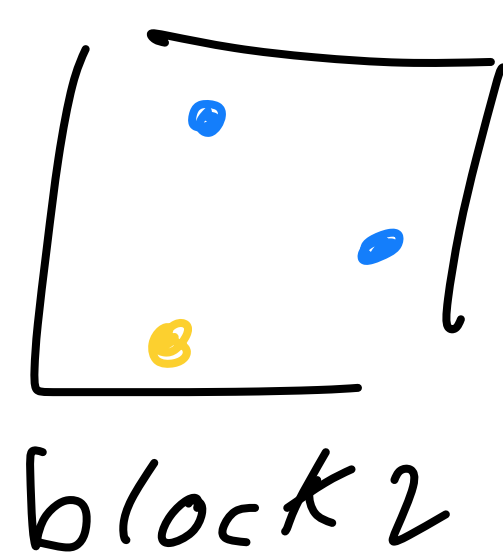
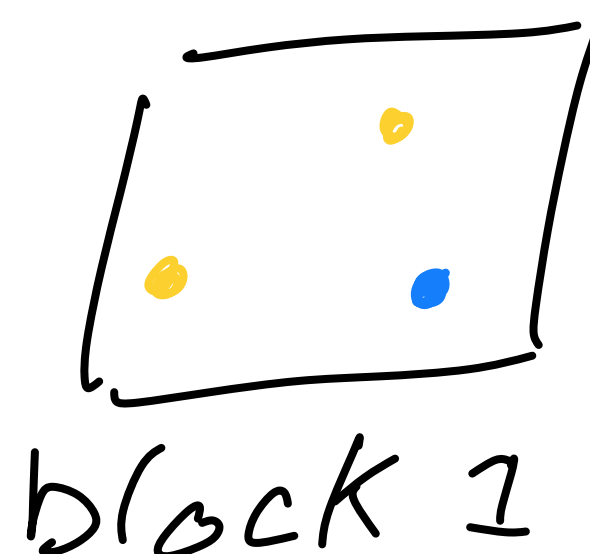
How you compute is tied to your data representation

Here suppose we have:

- A variable "blockID" for each record indicating block membership
- A variable "entityID" representing the true entity ID for the record.

→ These are called membership vectors

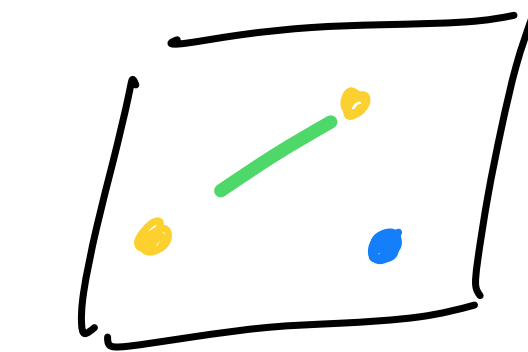
→ Candidate pairs and matching pairs are defined implicitly



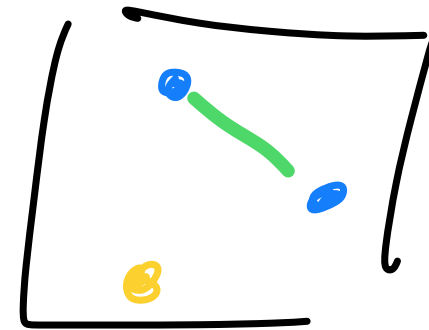
• = entity 1
• = entity 2
• = entity 3

Computing Recall

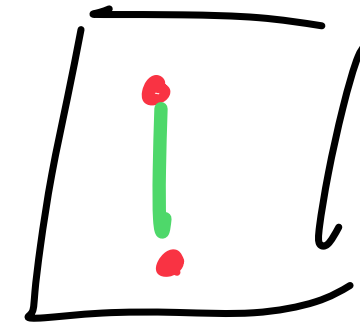
$$TP = 3$$



block 1



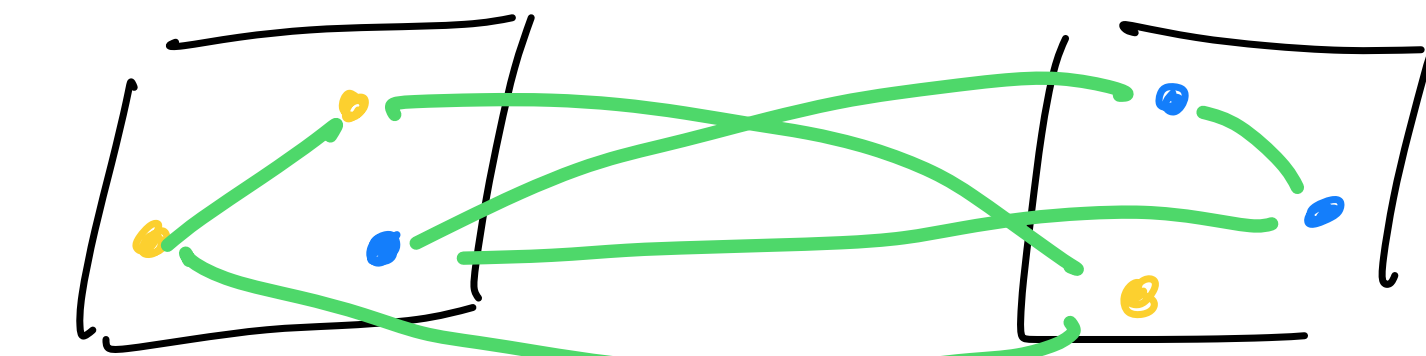
block 2



block 3

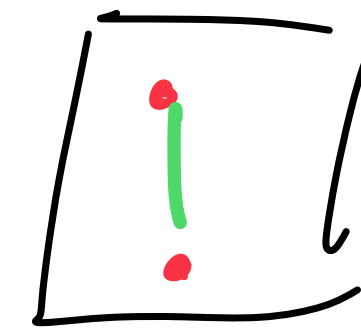
• = entity 1
• = entity 2
• = entity 3

$$\# \text{ matching pairs} = 7$$



block 1

block 2

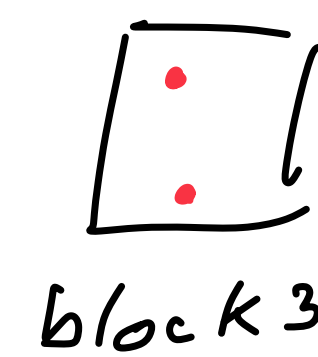
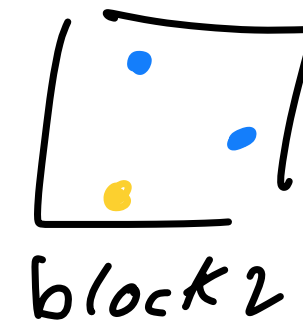
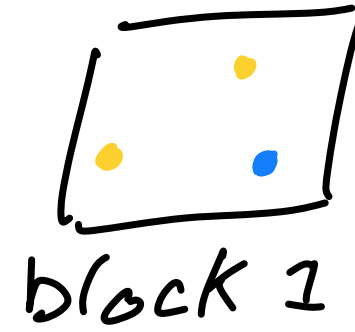


block 3

• = entity 1
• = entity 2
• = entity 3

$$\text{recall} = 3/7$$

Recall in R



• = entity 1
• = entity 2
• = entity 3

$$ct = xtabs(\sim blockID + entityID)$$

$$= \begin{matrix} & \begin{matrix} ent. 1 & ent. 2 & ent. 3 \end{matrix} \\ \begin{matrix} block 1 \\ block 2 \\ block 3 \end{matrix} & \left[\begin{array}{ccc} 2 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 2 \end{array} \right] \end{matrix}$$

$$TP = ?$$

$$\# \text{ matching pairs} = ?$$

$$cs = colSums(ct) = \begin{bmatrix} 3 & 3 & 2 \end{bmatrix}$$

↗
of each
entity

Is this efficient?

Computing precision

$$\text{precision}(\text{block ID}, \text{entity ID}) = \text{recall}(\text{entity ID}, \text{block ID})$$

why?