

Categorical data

- ↳ frequent patterns
- ↳ Association Rule

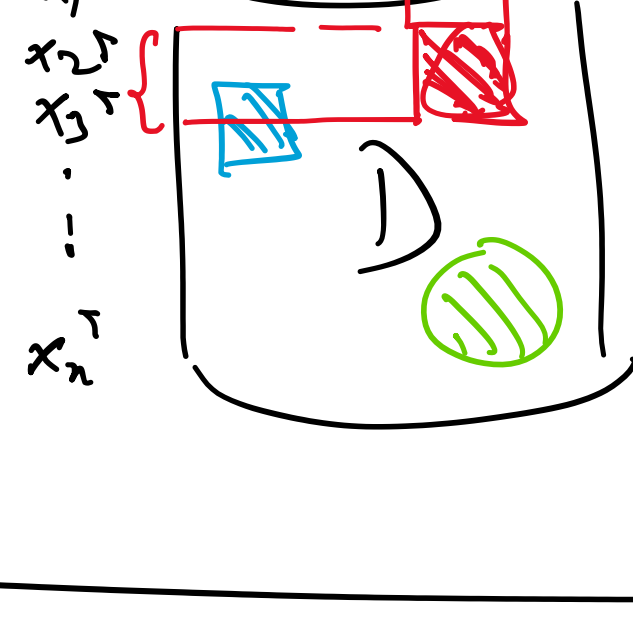
$$X \Rightarrow Y$$

X implies Y , (Confidence, frequency)

Patterns

- ↳ itemset "set" of co occurring events (unordered)
- ↳ sequence \rightarrow temporal (ordered)
- ↳ positional (ordered)

- ↳ graphs (subgraphs that are 'common')



model for the whole data

local model for a subset of the points and a subset of the attributes

Itemset mining

frequently occurring sets of items

$$I = \{ \text{set of possible items} \}$$

$$I = \{ A, B, C, D, E \}$$

$$T = \{ \text{set of instances / transactions} \}$$

$$T = \{ 1, 2, \dots, n \}$$

ids

$$X \subseteq I$$

Itemset is simply some subset of I

$$X = AB$$

$$= \{A, B\}$$

$$X = ACE$$

$$= \{A, C, E\}$$

$Y \subseteq T$ is called a tidset (transaction id set)

binary data matrix

	A	B	C	D	E
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

Sparse matrix

nnz # of non-zero entries

$$nnz \ll |I| \times |T|$$

Task:

Submatrix of all ones

$$X \times Y$$

such that

$$|Y| \geq \text{minsup}$$

then we also call X a frequent itemset

all tids in Y support X (contain all items in X)

BCE is therefore a frequent pattern

minsup \equiv minimum support

e.g.

Combinatorial enumeration task!

find all such frequent patterns X , that occur \geq minsup times

I : $2^{|I|}$ possible subsets / itemsets

$$I = \{ A, B, C, D, E \}$$

2^5 possible subsets \equiv powerset

0 $\rightarrow \emptyset$ empty pattern (always frequent)

1 $\rightarrow A, B, C, D, E$

2 $\rightarrow AB, AC, AD, AE, BC, BD, BE, CD, CE, DE$

3 $\rightarrow ABC, ABD, ABE, \dots$

4 $\rightarrow ABCD, ABCE, \dots$

5 $\rightarrow ABCDE$

then count: find all the tids that support each itemset

Apriori Algorithm!

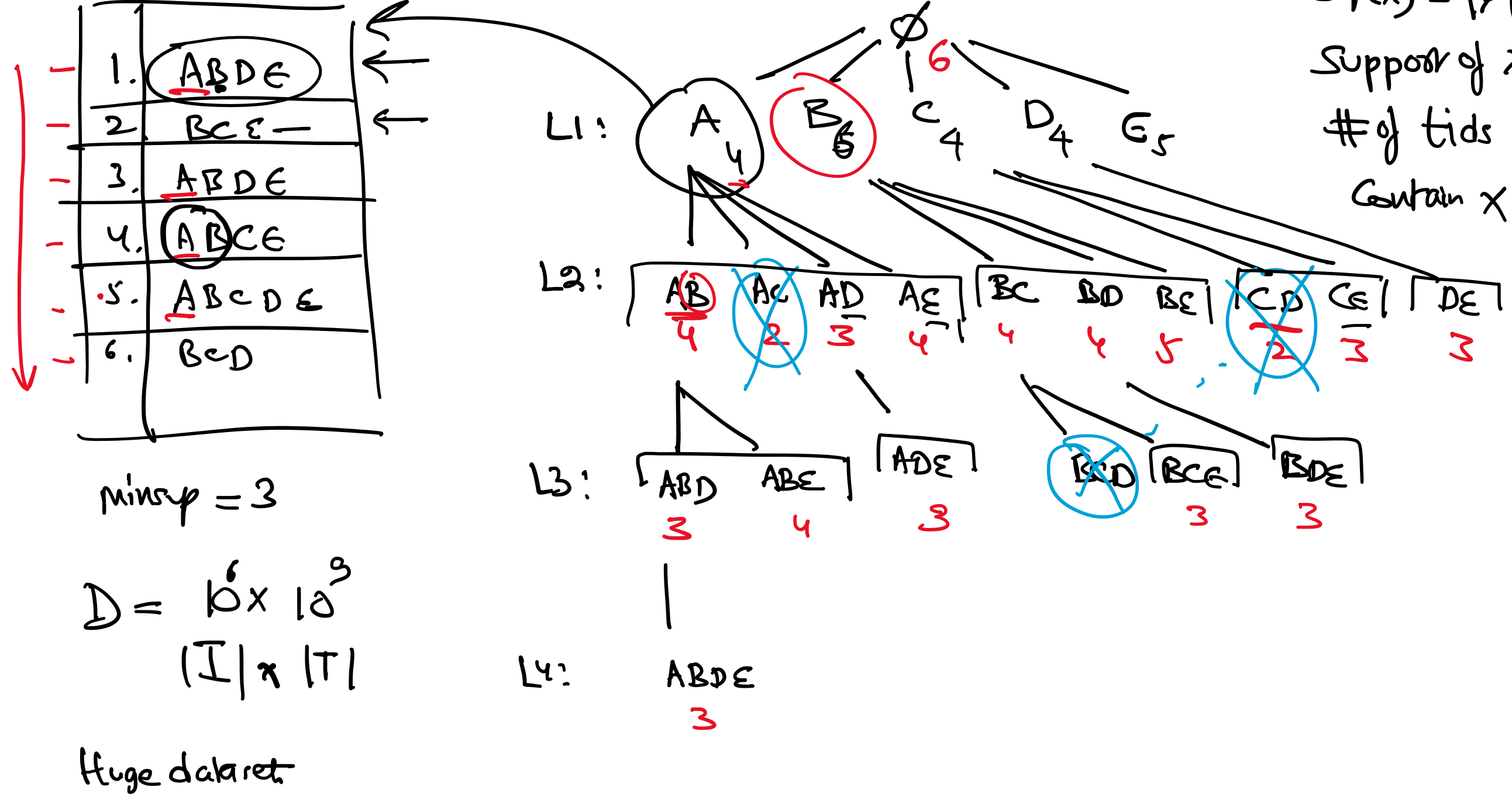
level-wise / BFS

$$(X \times Y)$$

$$\text{sup}(X) = |Y|$$

Support of $X \equiv$

of tids that contain X



$$D = 6 \times 10^3$$

$$|I| \times |T|$$

Huge dataset

$$F = \{ A, B, C, D, E, AB, AD, \dots, ABDE \}$$

Collection of frequent patterns.

1) itemset search space

combinatorial

BFS

DFS

2) counting

data organization

(indexing)

Eclat (Equivalence class transformation)

New db format

A	B	C	D	E
1	2	2	1	1
3	4	5	3	3
5	5	6	6	5
	6			

Inverted index

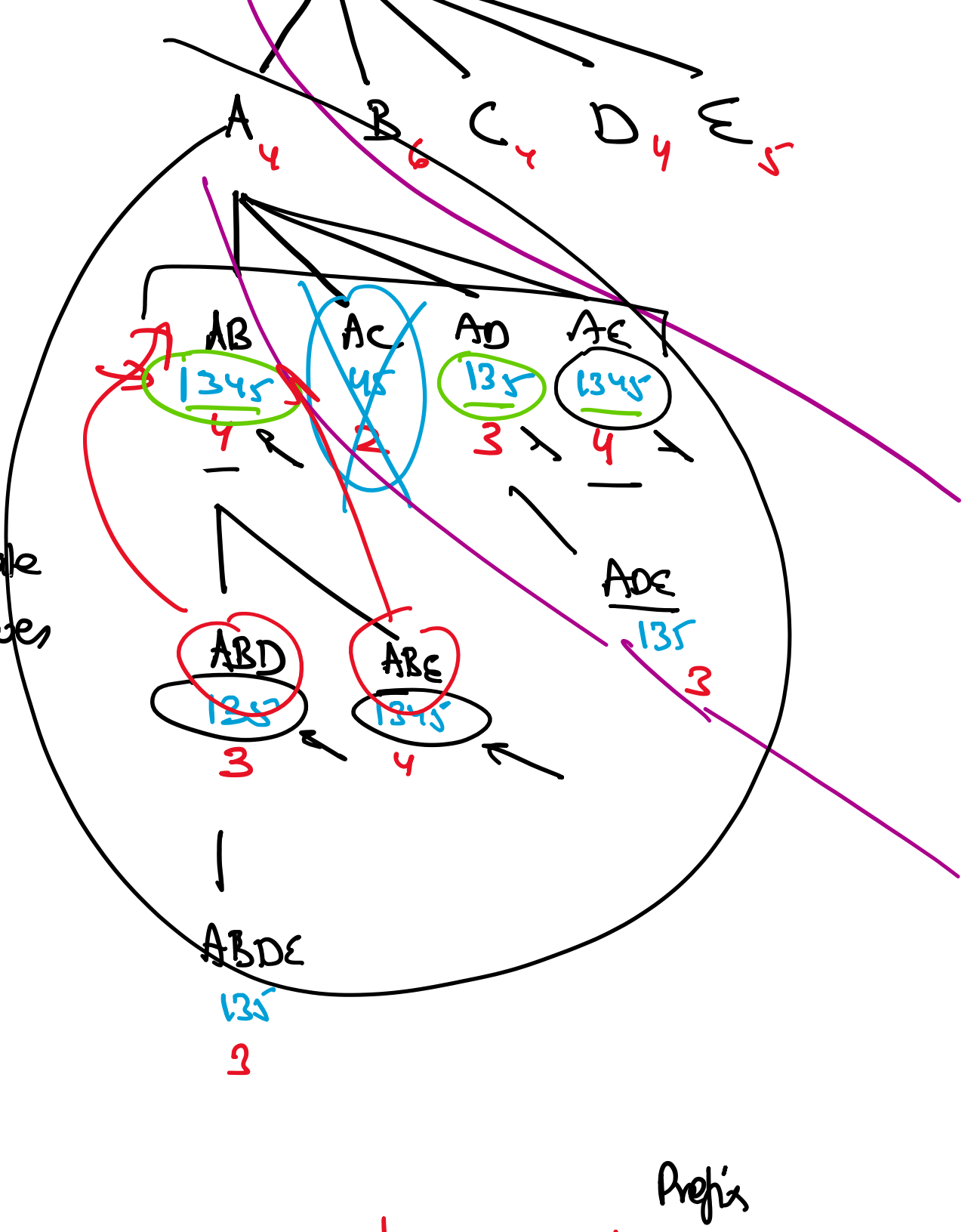
$$t(A) = \text{tidset for } A$$

$$t(X) = \text{tidset for itemset } X$$

$$\text{sup}(X) = |t(X)|$$

$$t(AB) = t(A) \cap t(B)$$

Cur! Intermediate storage goes up



$$X = AB$$

$$t(X) = \text{tidset of } X$$

$$= \{1, 3, 4, 5\} = 1345$$

$$\text{diffset}(X) = t(\text{prefix of } X) - t(X)$$

difference of tidsets

= tids that contain prefix but not the last item in X

$$\text{Prefix } d(ABD) = t(AB) - t(ABD)$$

$$= \{1, 3, 4, 5\} - \{1, 3, 5\}$$

$$= \{4\}$$

dEclat: Eclat on diffsets only

1) Level 1: Use tidsets to create diffsets

2) for other levels use only diffsets!

(propagate only differences)

A	B	C	D	E
1345	123456	2456	1356	12345

dataset

diffset

$$\text{sup}(X) = \text{sup}(\text{prefix}) - \text{len}(\text{diffset})$$

$$d(AB) = d(B) - d(A)$$

$$d(AD) = d(D) - d(A)$$

