

## Numerical tensor $\mathcal{T} \in \mathbb{R}^{\prod_{i=1}^n \mathcal{D}_i}$

Numerical tensors associate  $n$ -tuples with numerical values:

{ (week 1, Botafogo, 3toques)  $\mapsto$  6 retweets  
 (week 1, Botafogo, 4evertton)  $\mapsto$  5 retweets  
 |  
 (week 45, Corinthians, espnagora)  $\mapsto$  185 retweets  
 |  
 (week 49, Goiás, milton neves)  $\mapsto$  100 retweets  
 |  
 (week 49, Vasco, zigugo28)  $\mapsto$  1 retweet }

## Pattern $X \in \prod_{i=1}^n 2^{\mathcal{D}_i}$

A pattern binds subsets of each of the  $n$  dimensions. With an appropriate semantics, it supports the discovery of correlations between elements of all dimensions:

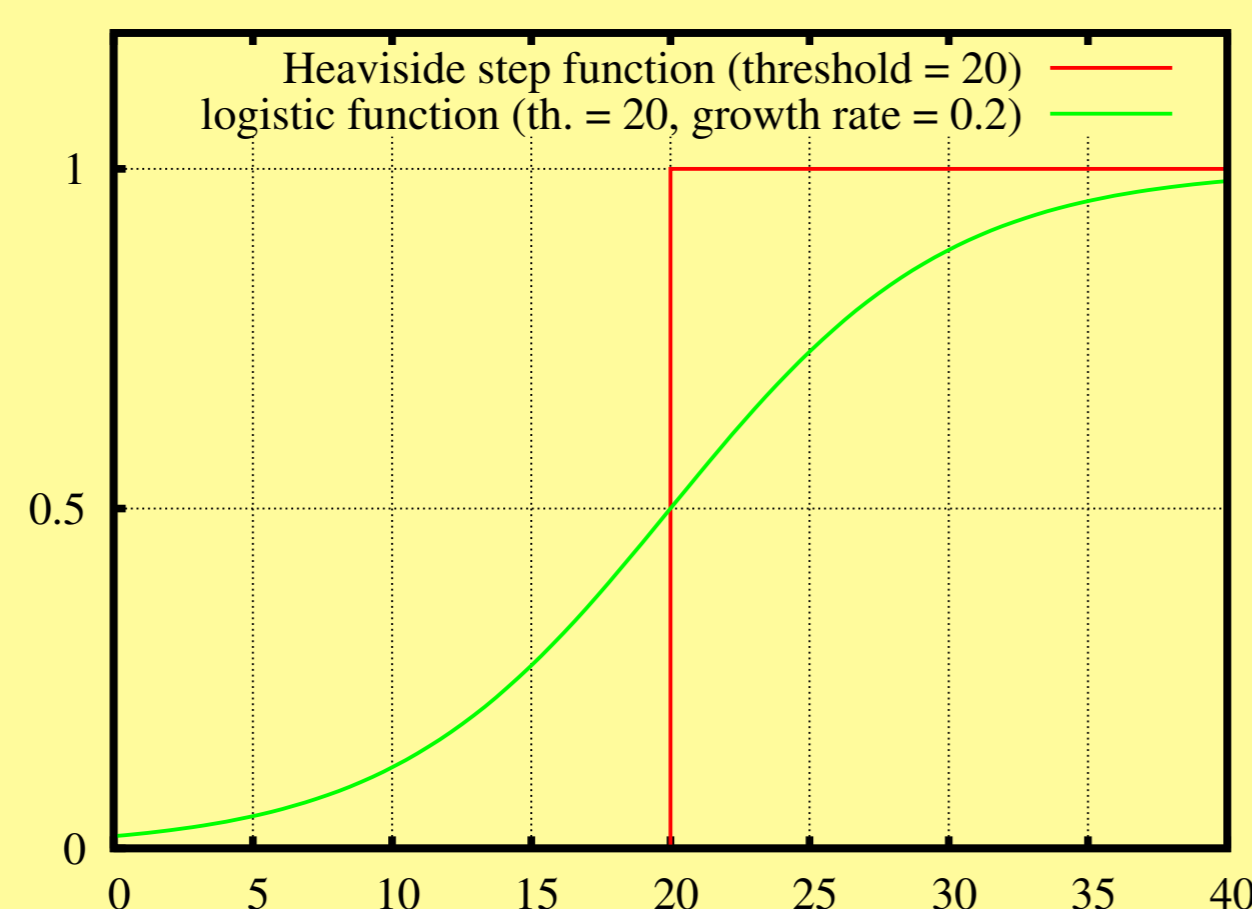
( {week 45, week 47, week 48, week 49}  
 {Corinthians, Fluminense, Goiás}  
 {espnagora, ge.timao, milton neves} )

The pattern is said to *cover* the Cartesian product of its  $n$  subsets. E. g., the pattern above covers this 3-tuple:

(week 45, Corinthians, ge.timao)

## Fuzzy relation $\mathcal{R} \in [0, 1]^{\prod_{i=1}^n \mathcal{D}_i}$

To define patterns covering  $n$ -tuples associated with large numbers (for instance), the numerical tensor is turned into a fuzzy relation. E. g., a logistic function is applied:



The extremal Heaviside step function is commonly used but implies a loss of information w.r.t. the numerical tensor.

## Absolute and per-element tolerance to noise

The *noise* a pattern tolerates is a distance between the membership degrees of its cover in the fuzzy relation and '1' values. In our opinion, noise tolerance must be:

**absolute** so that 1) the sole closed patterns form a lossless representation of all patterns and 2) the tolerance is anti-monotone, hence allows a scalable search of the valid patterns in the pattern space.

**per-element** so that, contrary to a per-pattern tolerance, every element really belongs to the pattern. E. g., in the relation below,  $(\{\alpha, \beta, \gamma\}, \{a, b, c, d\})$  must be invalid:

	a	b	c	d
$\alpha$	0.4	1	1	0
$\beta$	1	1	1	0
$\gamma$	1	1	1	0

## Definition of a valid pattern

Given  $n$  noise tolerance thresholds  $\epsilon = (\epsilon_i)_{i=1..n} \in \mathbb{R}_+^n$ , the pattern  $X = (X_i)_{i=1..n} \in \prod_{i=1}^n 2^{\mathcal{D}_i}$  is valid if and only if:

$$\begin{aligned} \mathcal{C}_{\epsilon\text{-connected}} \quad & \forall i = 1..n, \forall x \in X_i, \sum_{t \in X_1 \times \dots \times \{x\} \times \dots \times X_n} 1 - \mathcal{R}_t \leq \epsilon_i \\ \mathcal{C}_{\epsilon\text{-closed}} \quad & \forall i = 1..n, \forall s \in \mathcal{D}_i \setminus X_i, \\ & \left\{ \begin{array}{l} \sum_{t \in X_1 \times \dots \times \{s\} \times \dots \times X_n} 1 - \mathcal{R}_t > \epsilon_i \\ \text{or} \\ \exists x \in \cup_{j \neq i} X_j \mid \sum_{t \in X_1 \times \dots \times X_i \cup \{s\} \times \dots \times \{x\} \times \dots \times X_n} 1 - \mathcal{R}_t > \epsilon_j \end{array} \right. \end{aligned}$$

Any other relevant property is enforced via an additional user-defined constraint  $\mathcal{C}$ . To have it prune the pattern space, each *occurrence* of the pattern in its expression must behave monotonically or anti-monotonically. E. g., this constraint forces the patterns, with a numerical dimension  $X_i \subset \mathbb{R}$ , to have those elements  $\tau$ -close from each other:

$$\forall t \in [\min(X_i), \max(X_i)], \exists t' \in X_i \text{ s.t. } |t' - t| \leq \tau$$

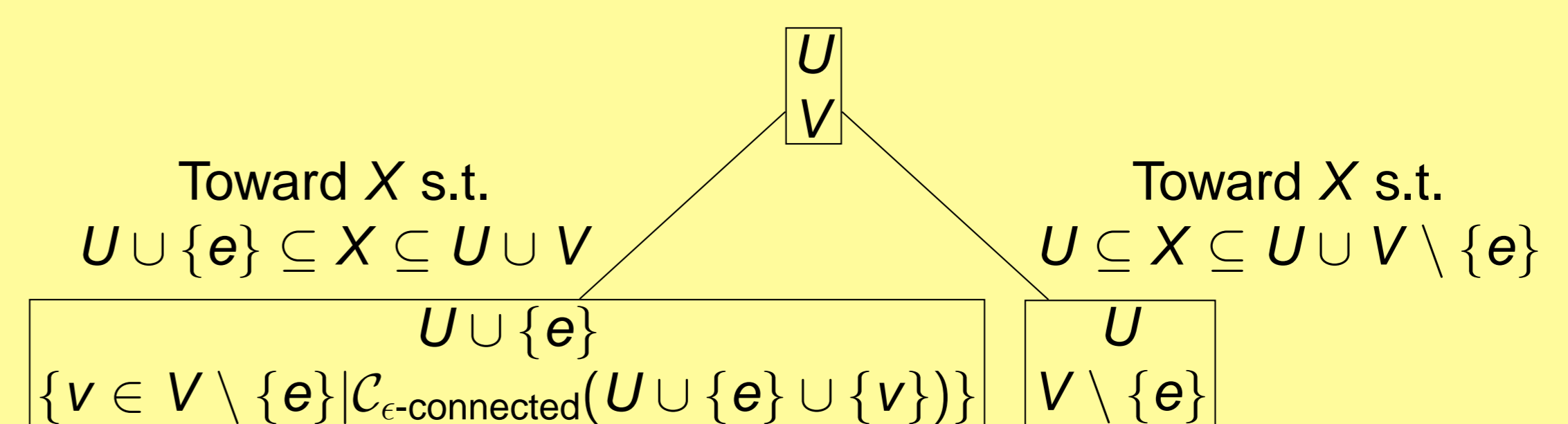
## Traversal of the pattern space

A binary tree structures the pattern space. At every node, root of a sub-tree:

$U \subseteq \cup_{i=1}^n \mathcal{D}_i$  contains the elements that *must* be present in every pattern at the leaves of the sub-tree;

$V \subseteq \cup_{i=1}^n \mathcal{D}_i$  contains the elements that *can* be present in some patterns at the leaves of the sub-tree.

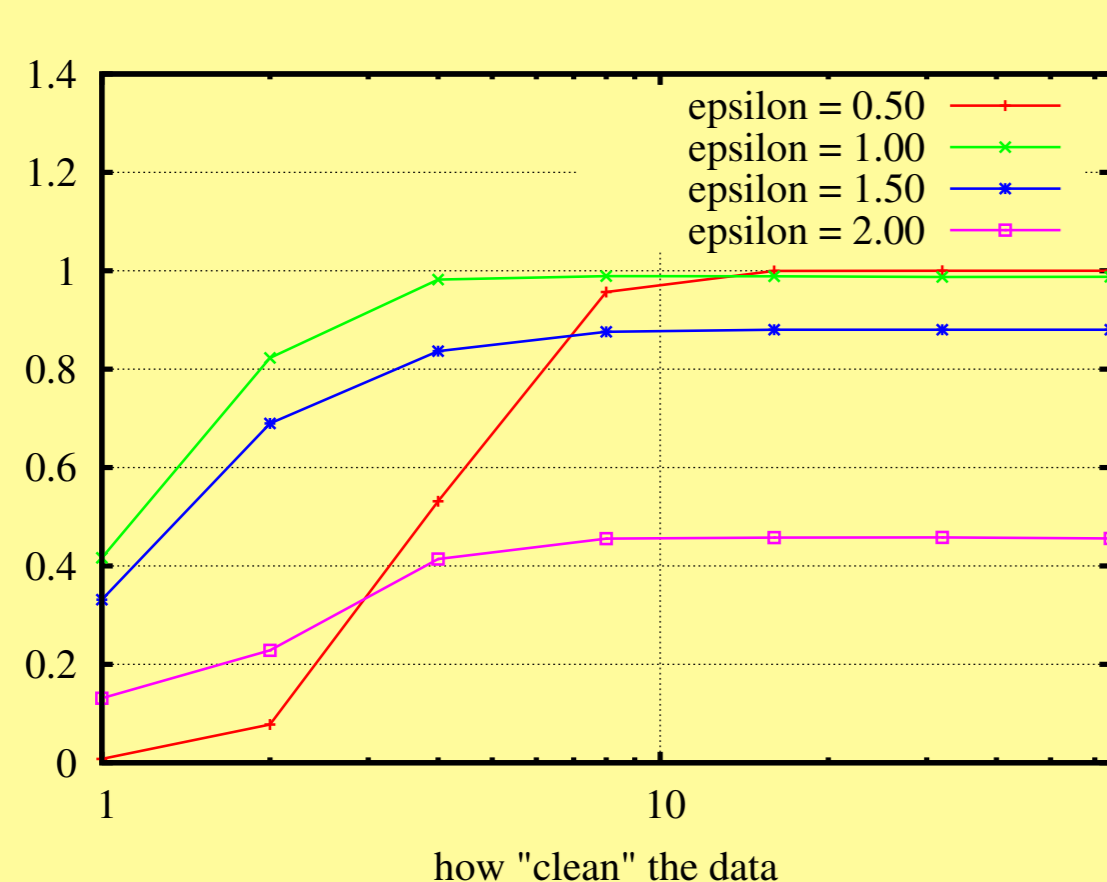
At the root,  $(U, V) = (\emptyset, \cup_{i=1}^n \mathcal{D}_i)$ . A reduced tree, which does not miss any valid pattern, is traversed depth-first. If  $V \neq \emptyset$ ,  $e \in V$  is freely chosen and these rules applied:



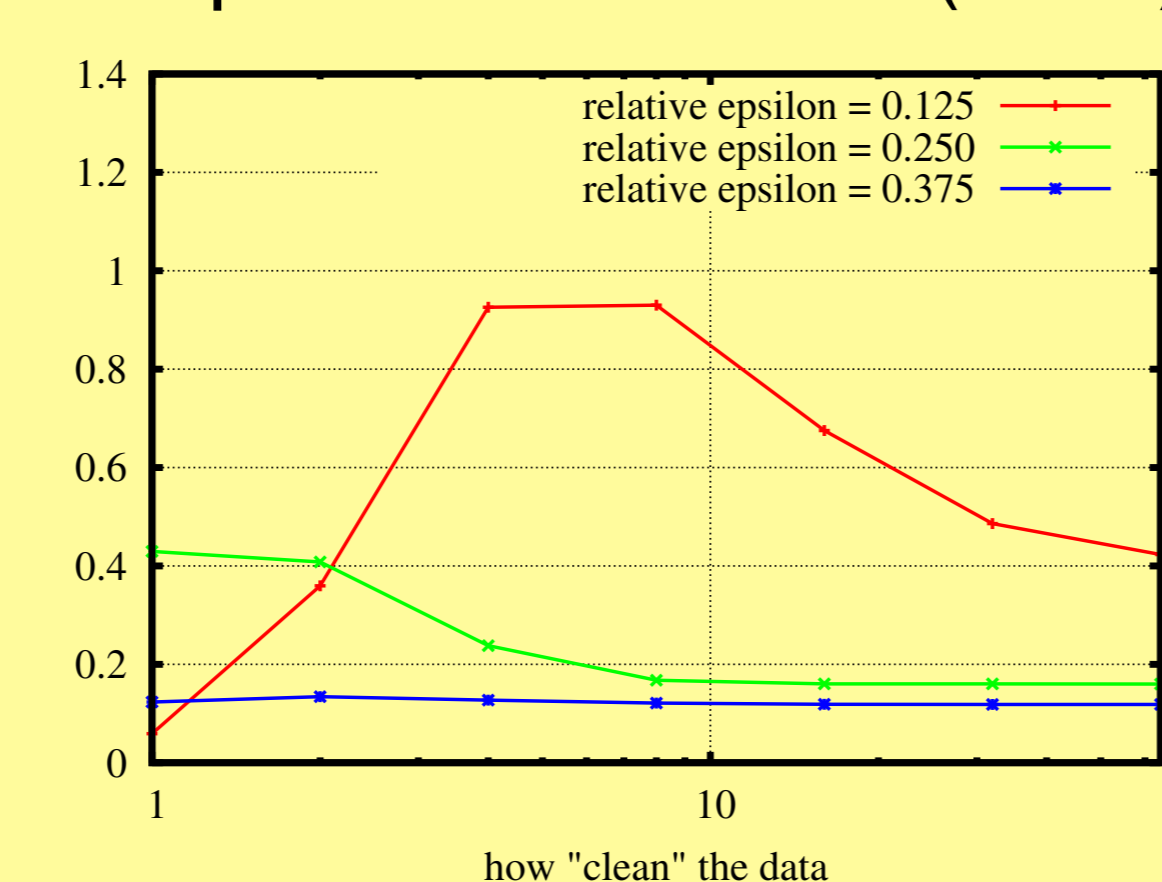
$\mathcal{C}_{\epsilon\text{-closed}}$  and  $\mathcal{C}$  are enforced at every node (pruning).

## Quality of the pattern collection ( $\mathcal{R}$ is $16 \times 16 \times 16$ )

### Per-element tolerance

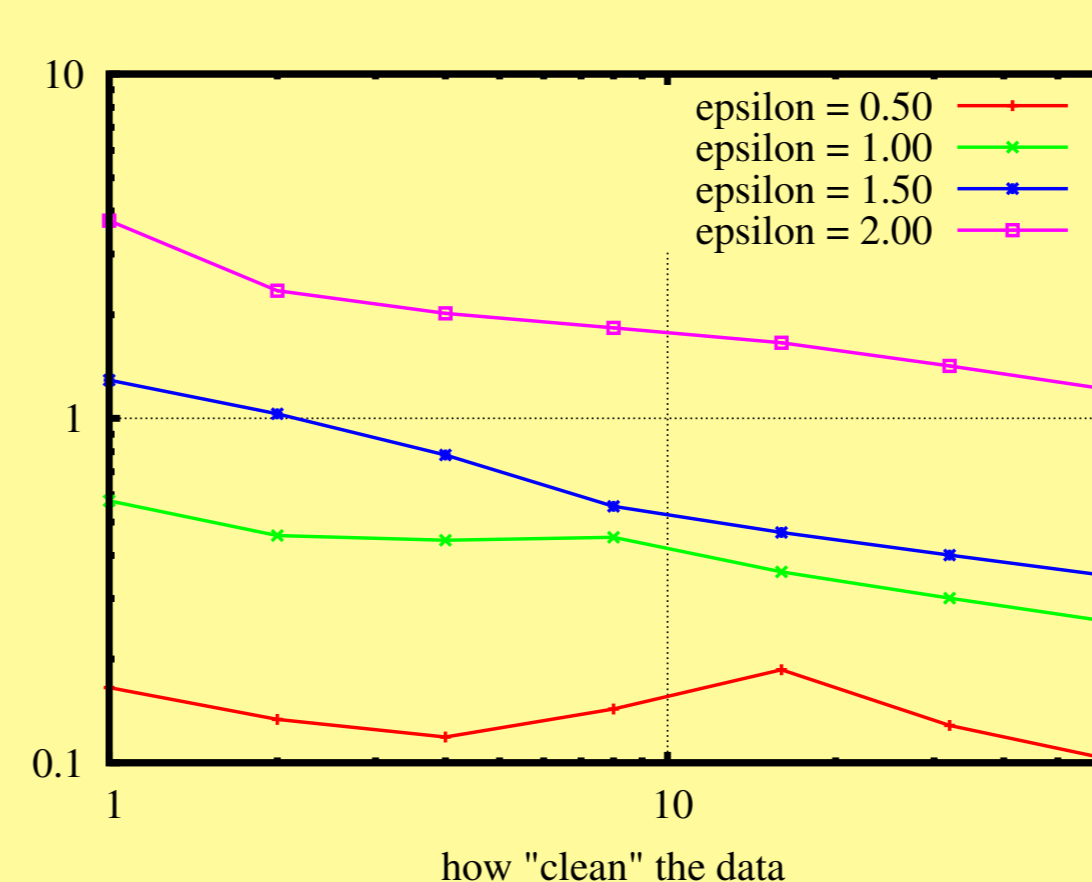


### Per-pattern tolerance (DCE)



## Running time in seconds ( $\mathcal{R}$ is $16 \times 16 \times 16$ )

### Absolute tolerance



### Relative tolerance (DCE)

