

RANKING EPISODES USING A PARTITION MODEL

NIKOLAJ.TATTI@AALTO.FI

Aalto University, Helsinki Institute of Information Technology



Aalto University

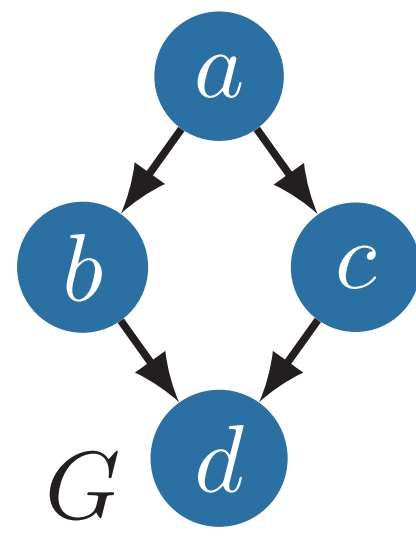
RANKING EPISODES

Episodes are

- (i) patterns occurring in sequences,
- (ii) order of events is described by DAGs,
- (iii) gap events are allowed.

Example:

Episode G occurs in a sequence if and only if (i) a occurs, (ii) then b and c , in any order, (iii) and then d . Gap events are allowed.



For a set of sequences \mathcal{S} , the support is

$$\text{supp}(G) = |\{S \in \mathcal{S} \mid G \text{ occurs in } S\}|$$

Mining using support counter-productive:

- (i) output is humongous,
- (ii) pattern are redundant

Rank patterns based on expectation.

Computing expectation is more intricate than with itemsets:

- (i) models are difficult to compute,
- (ii) depends on the sequence length

Assume we can get

$$p(n) = p(G \text{ occurs in a sequence of length } n)$$

Then the expected support is

$$\mu = \sum_{S \in \mathcal{S}} p(|S|)$$

Compare $\text{supp}(G)$ and μ :

The larger the difference, the more important is the episode.

EXPERIMENTS

Top episodes in *Plant* dataset. The symbols x and y represent noise events.

Independence model			Partition model		
Rank	Episode type	$r_{ind}(G)$	Rank	Episode type	$r_{prt}(G)$
1.	$a \rightarrow b \rightarrow c \rightarrow d$	∞	1.	$a \rightarrow b \rightarrow c \rightarrow d$	10^{308}
2.	$k \begin{matrix} \xrightarrow{n} \\ \xrightarrow{m} \end{matrix} l$	249	2.	$e \rightarrow f$	128
3.-7.	$a \rightarrow b \rightarrow c \rightarrow d \rightarrow x$	184-185	3.	$k \begin{matrix} \xrightarrow{n} \\ \xrightarrow{m} \end{matrix} l$	78
8.	$e \rightarrow f$	128	4.-	$a \rightarrow b \rightarrow c \rightarrow d \rightarrow x$	0-14
9.-	$x \rightarrow y$ or x, y	2-14		or $x \rightarrow y$ or x, y	

Episodes discovered from *JMLR* abstracts:

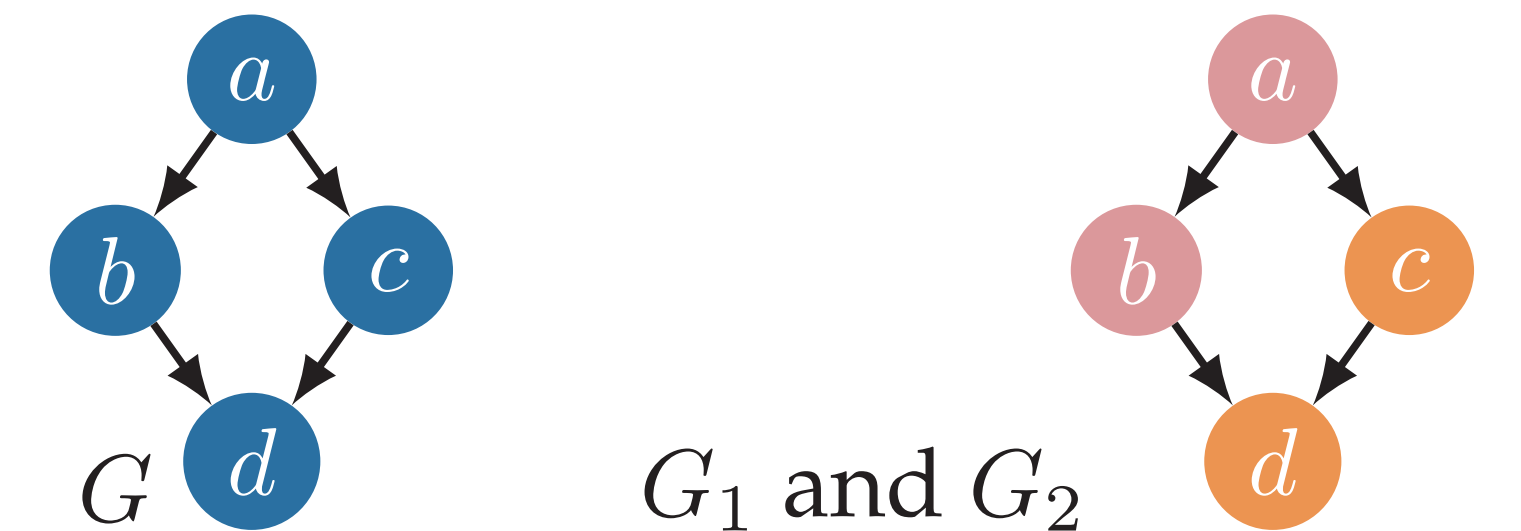
ranked by $r_{ind}(G)$		r_{ind}	r_{prt}	ranked by $r_{prt}(G)$		r_{ind}	r_{prt}
1.	support \rightarrow vector \rightarrow machin	∞	357	support \rightarrow vector	440	440	
2.	support \rightarrow vector	440	440	support \rightarrow vector \rightarrow machin	∞	357	
3.	support \rightarrow vector \rightarrow machin \rightarrow svm	404	90	support \rightarrow machin	324	324	
4.	support \rightarrow vector \rightarrow machin svm	356	10^{-3}	vector \rightarrow machin	306	306	
5.	reproduc \rightarrow kernel \rightarrow hilbert \rightarrow space	341	73	data \rightarrow set	284	284	
6.	support \rightarrow machin	325	325	real \rightarrow world	260	260	
7.	vector \rightarrow machin	306	306	real \rightarrow data	213	213	
8.	data \rightarrow set	284	284	state \rightarrow art	191	191	
9.	real \rightarrow world	260	260	machin \rightarrow learn	190	190	
10.	support \rightarrow vector \rightarrow svm	250	85	bayesian \rightarrow network	166	166	

PARTITION MODEL

General idea:

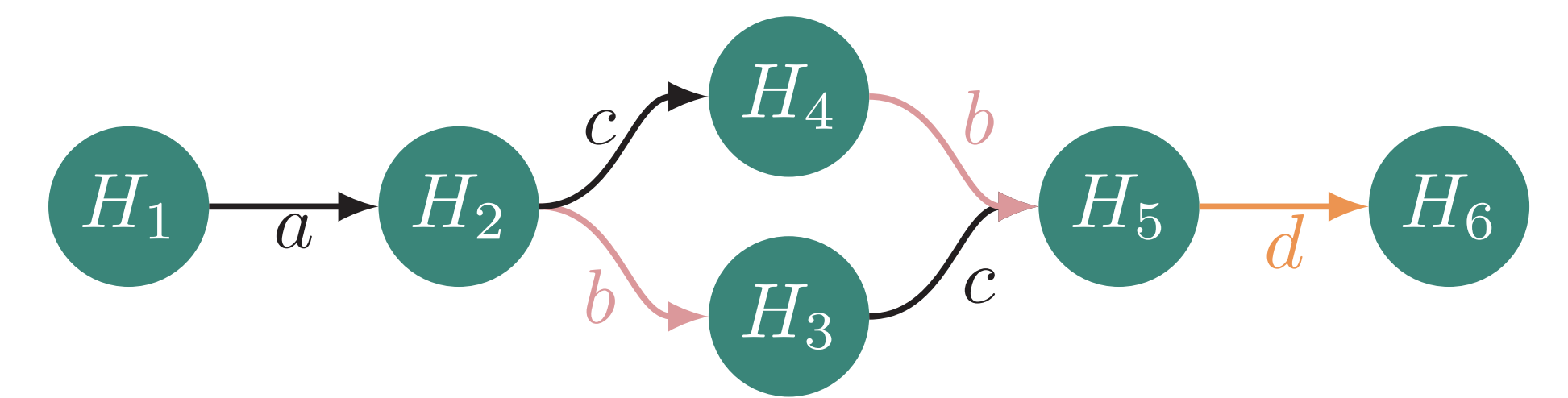
- (i) split the episode into two subepisodes,
- (ii) study how often these episodes occur,
- (iii) incorporate this into a model,
- (iv) try all splits, and use the best.

Example: Split G into G_1 and G_2



Model how likely G_i is discovered once we see at least one event in G_i .

In practice, boost the probabilities $p(e)$ in



Boosted probabilities leads to higher expected support, and lowered rank.

Edge e with label l is modeled as

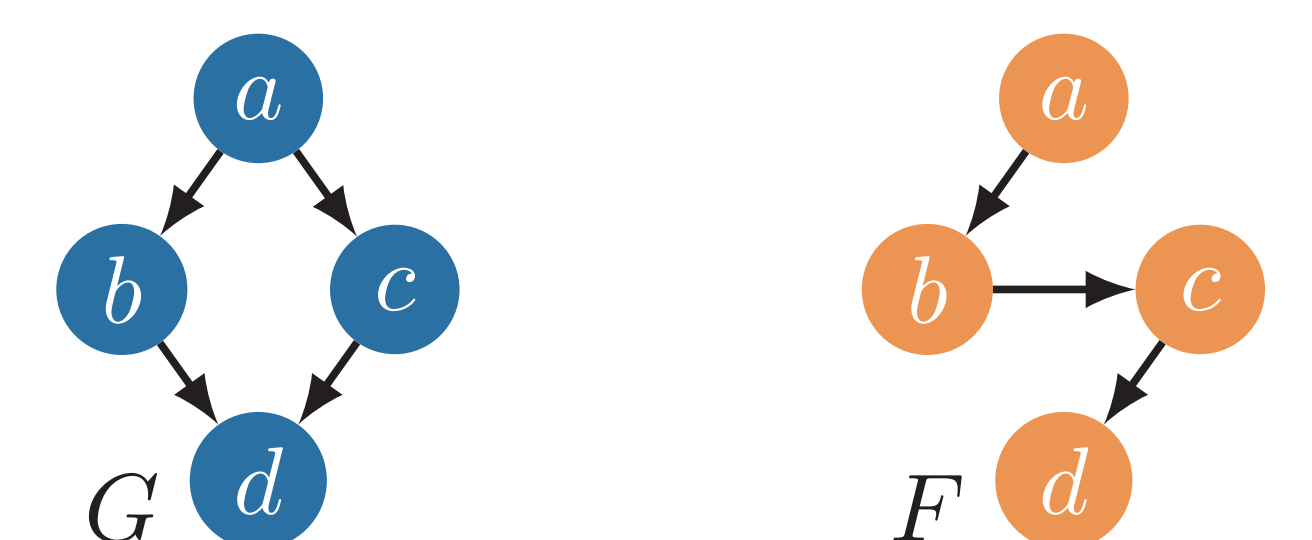
$$p(e) \propto \begin{cases} \exp(u_l + t_1), & \text{if } e \in C_1, \\ \exp(u_l + t_2), & \text{if } e \in C_2, \\ \exp(u_l), & \text{otherwise} \end{cases}$$

Parameters u_l and t_1 can be learned by maximizing likelihood; gradient descent will converge to the global optimum.

SUPEREPISODES

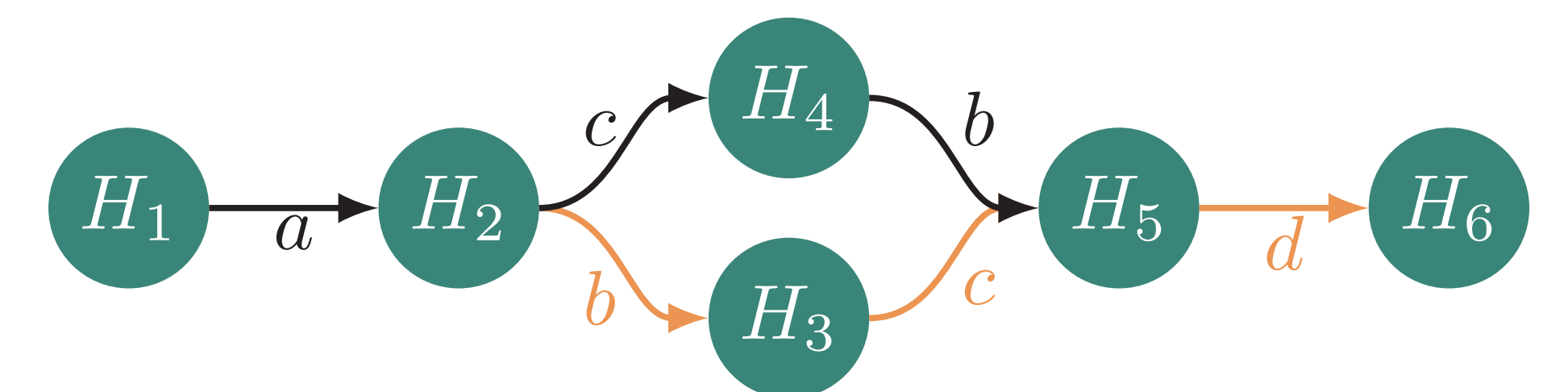
Similar to partition model except now test for superepisodes.

Example: F is a superepisode of G



Model how likely F is discovered once we see at least one event in F .

In practice, boost the probabilities $p(e)$ in



Episodes with high r_{ind} and low r_{prt}

G_1 :	support \rightarrow vector \rightarrow machin	regress	95.1
G_2 :	support \rightarrow vector \rightarrow machin	regress	90.4
G_3 :	support \rightarrow vector \rightarrow machin	number	52.0
G_4 :	support \rightarrow vector \rightarrow machin	regress	86.4
G_5 :	support \rightarrow vector \rightarrow machin	space	51.6