

Oct. 22, 2020
第41回IBISML



Inter-University Research Institute Corporation /
Research Organization of Information and Systems
National Institute of Informatics



統計的に有意な相互作用探索

杉山 磨人 (国立情報学研究所)

Significant Pattern Mining

- Detect **patterns** (feature combinations) that are **statistically significantly** enriched in a class of a dataset

Input:

	x						y
	F1	F2	F3	F4	F5	...	Class
ID1	0	1	1	1	0	...	0
ID2	1	1	0	1	1	...	1
ID3	1	1	0	0	1	...	1
ID4	0	0	1	0	1	...	0
⋮			⋮				⋮

Output:

{F1}, {F3},
{F2, F5},
{F2, F5, F6}, ...



Significant Pattern Mining

- Detect **patterns** (feature combinations) that are **statistically significantly** enriched in a class of a dataset

Input:

	x						y
	F1	F2	F3	F4	F5	...	Class
ID1	0	1	1	1	0	...	0
ID2	1	1	0	1	1	...	1
ID3	1	1	0	0	1	...	1
ID4	0	0	1	0	1	...	0
⋮			⋮				⋮

Output:

{F1}, {F3},
{F2, F5},
{F2, F5, F6}, ...

Significant Pattern Mining

- Detect **patterns** (feature combinations) that are **statistically significantly** enriched in a class of a dataset

Input:

	x						y
	F1	F2	F3	F4	F5	...	Class
ID1	0	1	1	1	0	...	0
ID2	1	1	0	1	1	...	1
ID3	1	1	0	0	1	...	1
ID4	0	0	1	0	1	...	0
⋮			⋮				⋮

Output:

{F1}, {F3},
{F2, F5},
{F2, F5, F6}, ...

Significant Pattern Mining

- Detect **patterns** (feature combinations) that are **statistically significantly** enriched in a class of a dataset

Input:

	x						y
	F1	F2	F3	F4	F5	...	Class
ID1	0	1	1	1	0	...	0
ID2	1	1	0	1	1	...	1
ID3	1	1	0	0	1	...	1
ID4	0	0	1	0	1	...	0
⋮			⋮				⋮

Output:

{F1}, {F3},
{F2, F5},
{F2, F5, F6}, ...

Example: Itemset Mining

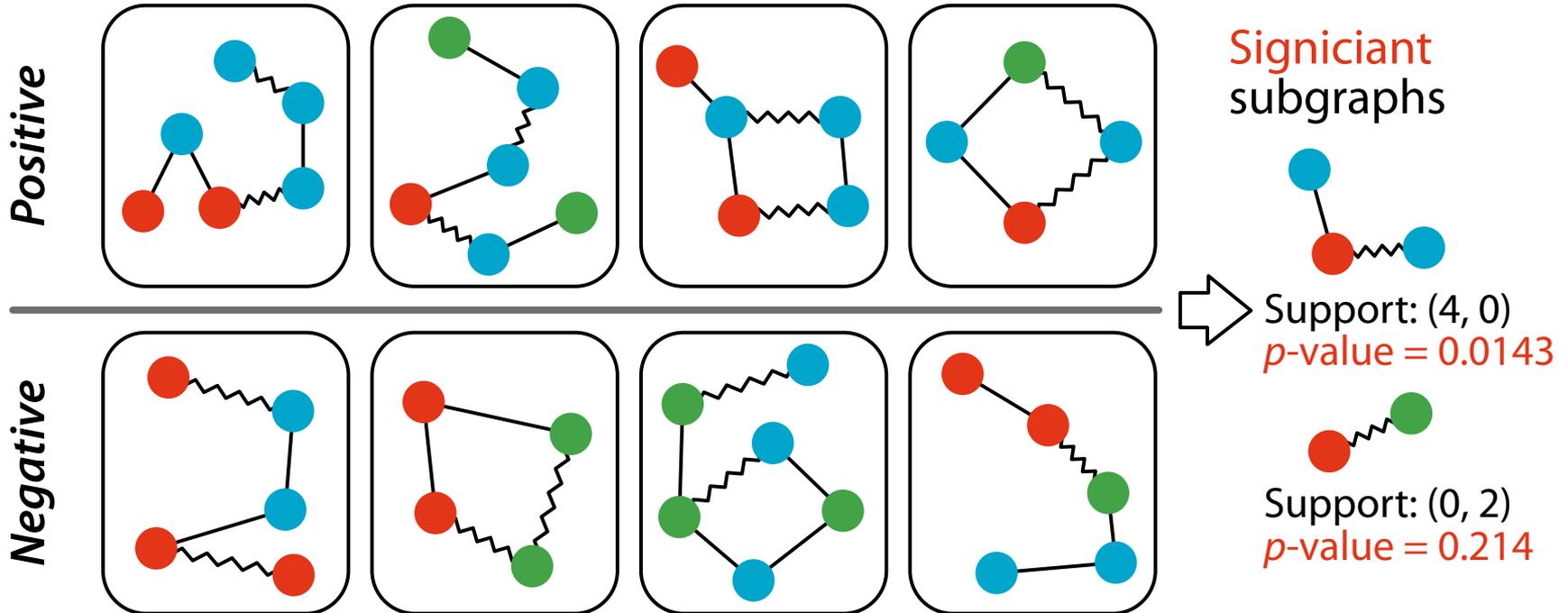
SNPs (items)

 Case	ID 1:	0 0 1 1 0 0 1 1 1 0 0 1 1 1 0 0 1 1 1 0
	ID 2:	1 1 0 0 1 0 1 1 1 1 0 0 0 1 0 1 0 1 0 0
	ID 3:	1 0 1 0 0 0 1 1 1 0 1 0 1 1 0 0 0 0 0 1
	ID 4:	1 1 0 1 1 0 1 1 1 1 1 1 1 1 1 0 1 0 0 1 1
<hr/>		
 Control	ID 5:	0 0 1 1 0 0 0 1 1 0 0 0 1 1 1 1 1 0 0 0
	ID 6:	0 1 0 1 1 0 1 1 0 0 0 0 1 1 0 0 1 0 1 0
	ID 7:	1 0 1 1 0 0 1 0 1 0 0 0 0 0 1 0 1 0 0 0
	ID 8:	1 1 0 0 1 0 0 1 0 1 0 1 0 0 0 1 0 1 0 1

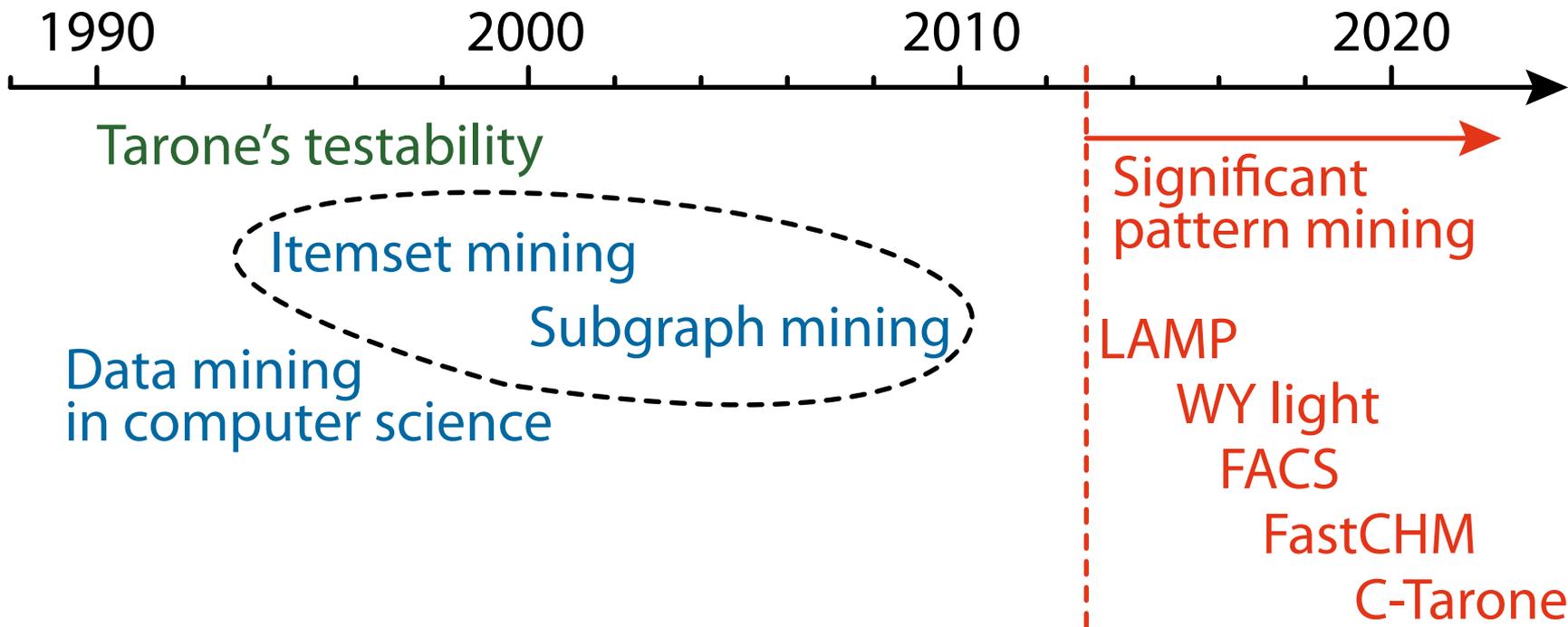
Example: Itemset Mining

		SNPs (items)	Signal
 Case	ID 1:	0 0 1 1 0 0 1 1 1 0 0 1 1 1 0 0 1 1 1 0	1
	ID 2:	1 1 0 0 1 0 1 1 1 1 0 0 0 1 0 1 0 1 0 0	1
	ID 3:	1 0 1 0 0 0 1 1 1 0 1 0 1 1 0 0 0 0 0 1	1
	ID 4:	1 1 0 1 1 0 1 1 1 1 1 1 1 1 0 1 0 0 1 1	1
<hr/>			
 Control	ID 5:	0 0 1 1 0 0 0 1 1 0 0 0 1 1 1 1 1 0 0 0	0
	ID 6:	0 1 0 1 1 0 1 1 0 0 0 0 1 1 0 0 1 0 1 0	0
	ID 7:	1 0 1 1 0 0 1 0 1 0 0 0 0 0 1 0 1 0 0 0	0
	ID 8:	1 1 0 0 1 0 0 1 0 1 0 1 0 0 0 1 0 1 0 1	0

Example: Subgraph Mining



Timeline



Recent Advances

- Webb, G.I., Petitjean, F.: **A Multiple Test Correction for Streams and Cascades of Statistical Hypothesis Tests**, [KDD2016](#)
- Pellegrina, L., Vandin, F.: **Efficient Mining of the Most Significant Patterns with Permutation Testing**, [KDD2018](#)
- Pellegrina, L., Riondato, M., Vandin, F.: **SPuManTE: Significant Pattern Mining with Unconditional Testing**, [KDD2019](#)
- Tran, T.Q., Fukuchi, K., Akimoto, Y., Sakuma, J.: **Statistically Significant Pattern Mining with Ordinal Utility**, [KDD2020](#)

Libraries

- CASMAP
 - Llinares-Lopez, et al.: **CASMAP: Detection of statistically significant combinations of SNPs in association mapping**, [Bioinformatics](#) (2019)
- MP-LAMP (for parallel computation)
 - Yoshizoe, K., Terada, A., Tsuda, K.: **MP-LAMP: parallel detection of statistically significant multi-loci markers on cloud platforms**, [Bioinformatics](#) (2018)

Key Challenges:

1. How to assess the significance for a **multiplicative interaction of variables**?
2. How to perform **multiple testing correction**?
 - How to control the **FWER** (family-wise error rate), the probability to detect one or more false positives?
3. How to manage **combinatorial explosion** (2^d for d variables) of the candidate space?

Problem Formulation

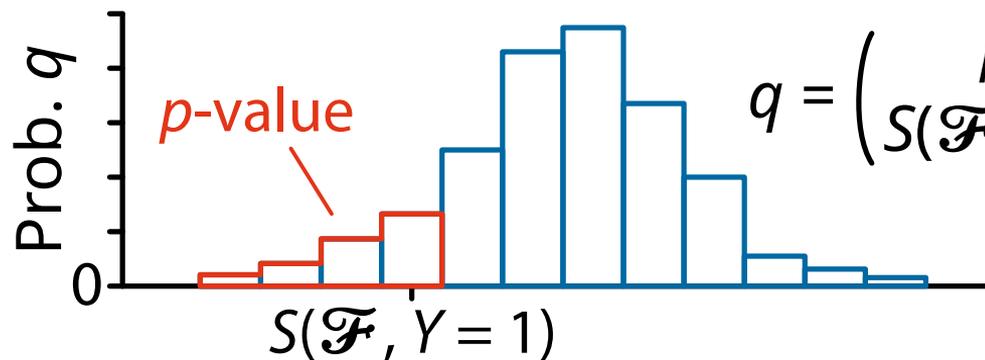
- Define $X_{\mathcal{F}}$ as **binary random variable** of joint occurrence for a feature combination $\mathcal{F} = \{F_i\}_{i \in I}, I \subseteq \{1, \dots, d\}$
 - $X_{\mathcal{F}} = 1$ if \mathcal{F} "occurs", $X_{\mathcal{F}} = 0$ otherwise
- Let Y be an output binary variable

Problem Formulation

- Define $X_{\mathcal{F}}$ as **binary random variable** of joint occurrence for a feature combination $\mathcal{F} = \{F_i\}_{i \in I}, I \subseteq \{1, \dots, d\}$
 - $X_{\mathcal{F}} = 1$ if \mathcal{F} "occurs", $X_{\mathcal{F}} = 0$ otherwise
- Let Y be an output binary variable
- **Task:** Test the null hypothesis $X_{\mathcal{F}} \perp\!\!\!\perp Y$ for *all* $\mathcal{F} \in 2^V$
 - Testing statistical independence between $X_{\mathcal{F}}$ and Y

Fisher's Exact Test

	$X_{\mathcal{F}} = 1$	$X_{\mathcal{F}} = 0$	Total
$Y = 1$	$S(\mathcal{F}, Y = 1)$	$N_1 - S(\mathcal{F}, Y = 1)$	N_1
$Y = 0$	$S(\mathcal{F}, Y = 0)$	$N_0 - S(\mathcal{F}, Y = 0)$	N_0
Total	$S(\mathcal{F})$	$N - S(\mathcal{F})$	N



$$q = \binom{N_1}{S(\mathcal{F}, Y=1)} \binom{N_0}{S(\mathcal{F}, Y=0)} / \binom{N}{S(\mathcal{F})}$$

Multiple Testing Correction

- In each test, [probability of having a false positive] $\leq \alpha$
- If we repeat m tests, αm patterns can be false positives
 - Too many if m is large! For example in itemset mining:
 - For 100000 items, #patterns = 2^{100000}
 - Set significance level $\alpha = 0.01$
 - Number of false positives: $0.01 \cdot 2^{100000} = 10^{30101}$
- The **FWER** should be controlled
 - Probability at least one \mathcal{F} is false positive

Controlling the FWER

- $\text{FWER} = \Pr(\text{FP} > 0)$
 - FP: Number of false positives
- **Objective:** Maximize $\text{FWER}(\delta)$ subject to $\text{FWER}(\delta) \leq \alpha$
 - $\text{FWER}(\delta)$: FWER at corrected significance level δ
 - Cannot be evaluated in closed form (simple but not easy!)
 - Bonferroni correction is popular: $\delta_{\text{Bon}}^* = \alpha/m$

Tarone's Testability Trick

- We use Tarone's testability trick, which requires the minimum achievable p -value $\psi(\mathcal{F})$ for \mathcal{F}

$$\psi(\mathcal{F}) = \binom{N_1}{S(\mathcal{F})} / \binom{N}{S(\mathcal{F})} \quad \text{in Fisher's exact test}$$

Tarone's Testability Trick

$$\mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_3, \dots, \mathcal{F}_{m-1}, \mathcal{F}_m, \mathcal{F}_{m+1}, \dots, \mathcal{F}_{2^d} \quad (\psi(\mathcal{F}_i) \leq \psi(\mathcal{F}_{i+1}))$$

Tarone's Testability Trick

$$m \psi(\mathcal{F}_m) < \alpha \quad \text{and} \quad (m+1)\psi(\mathcal{F}_{m+1}) \geq \alpha$$

$$\mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_3, \dots, \mathcal{F}_{m-1}, \mathcal{F}_m, \mathcal{F}_{m+1}, \dots, \mathcal{F}_{2d} \quad (\psi(\mathcal{F}_i) \leq \psi(\mathcal{F}_{i+1}))$$

Tarone's Testability Trick

$$m \psi(\mathcal{F}_m) < \alpha \quad \text{and} \quad (m+1)\psi(\mathcal{F}_{m+1}) \geq \alpha$$

$\mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_3, \dots, \mathcal{F}_{m-1}, \mathcal{F}_m, \mathcal{F}_{m+1}, \dots, \mathcal{F}_{2^d}$ $(\psi(\mathcal{F}_i) \leq \psi(\mathcal{F}_{i+1}))$

Testable
combinations

Untestable combinations → Prune without testing

↓
 \mathcal{F}_i is significant if: $p\text{-value}(\mathcal{F}_i) < \alpha / \textcircled{m}$ — Correction factor

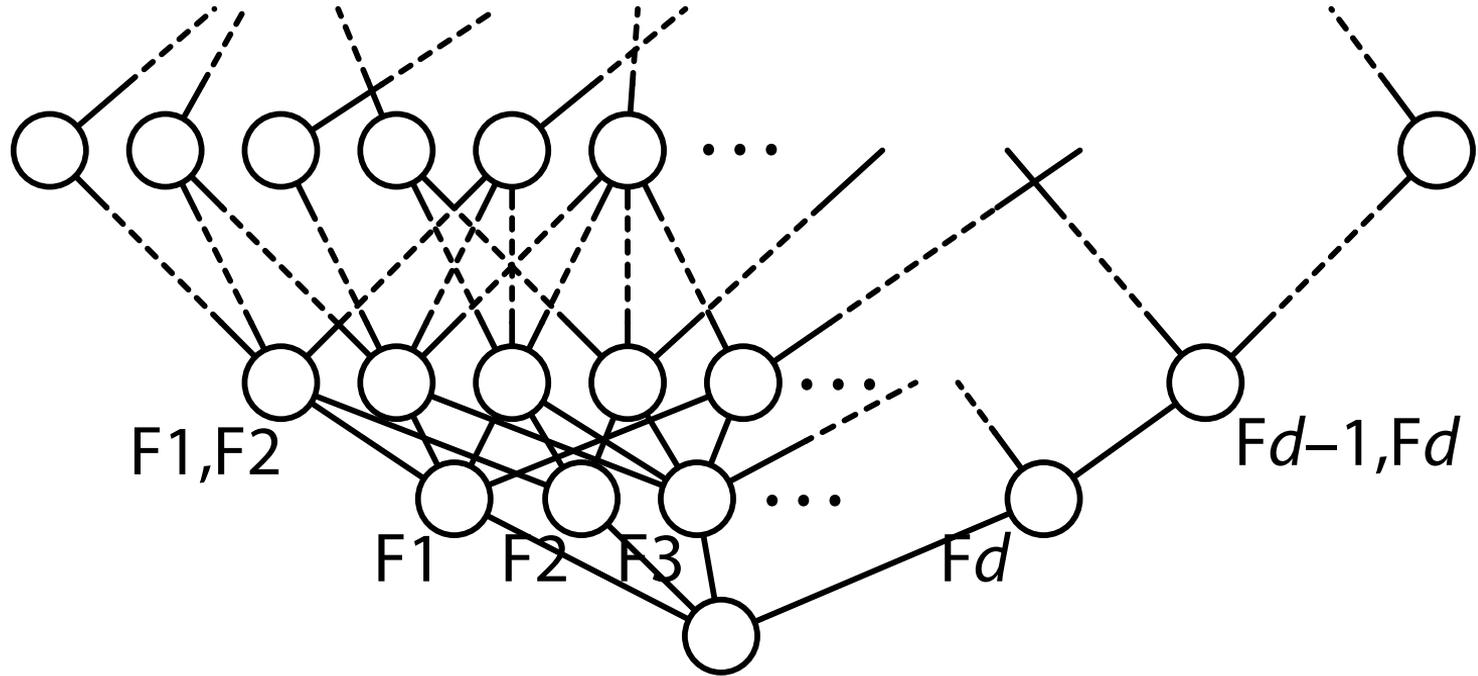
Tarone's Testability Trick with Apriori

- We use Tarone's testability trick, which requires the minimum achievable p -value $\psi(\mathcal{F})$ for \mathcal{F}

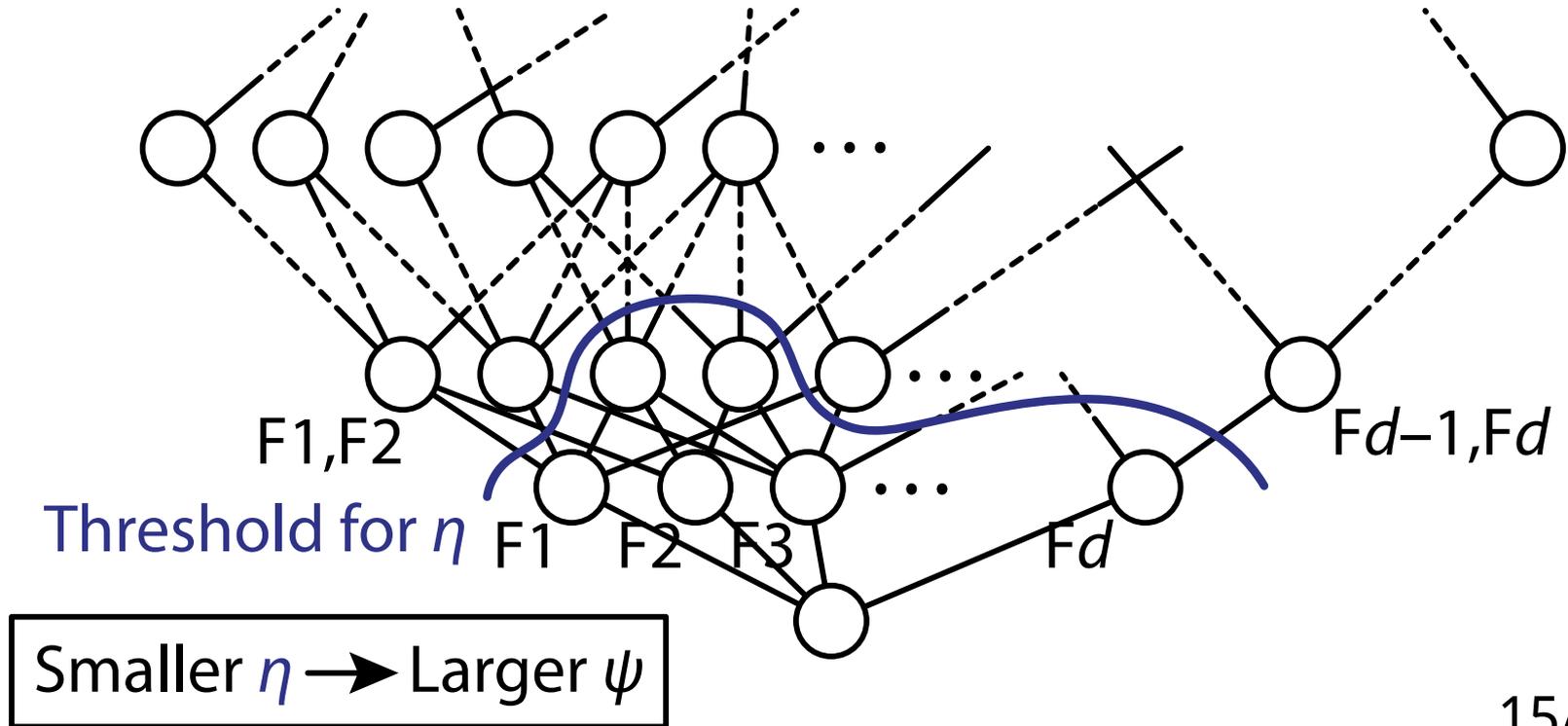
$$\psi(\mathcal{F}) = \binom{N_1}{S(\mathcal{F})} / \binom{N}{S(\mathcal{F})} \quad \text{in Fisher's exact test}$$

- This method is particularly effective if the relationship "Smaller $\eta(\mathcal{F}) \rightarrow$ Larger $\psi(\mathcal{F})$ " holds
 - For each pattern \mathcal{F} ,
 $S(\mathcal{F})$: Support (how many times \mathcal{F} occurs in a dataset)
 $\eta(\mathcal{F}) = S(\mathcal{F})/N$: Frequency

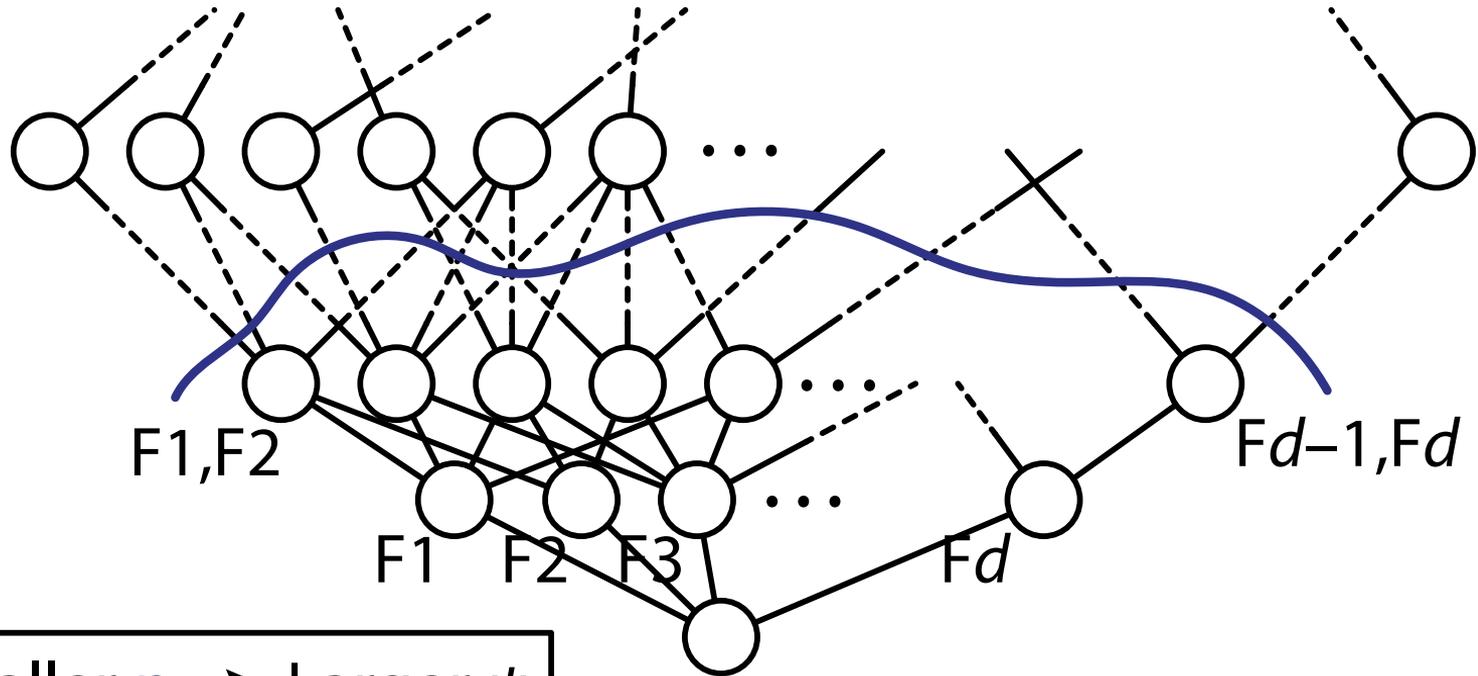
Enumeration Based on Apriori



Enumeration Based on Apriori

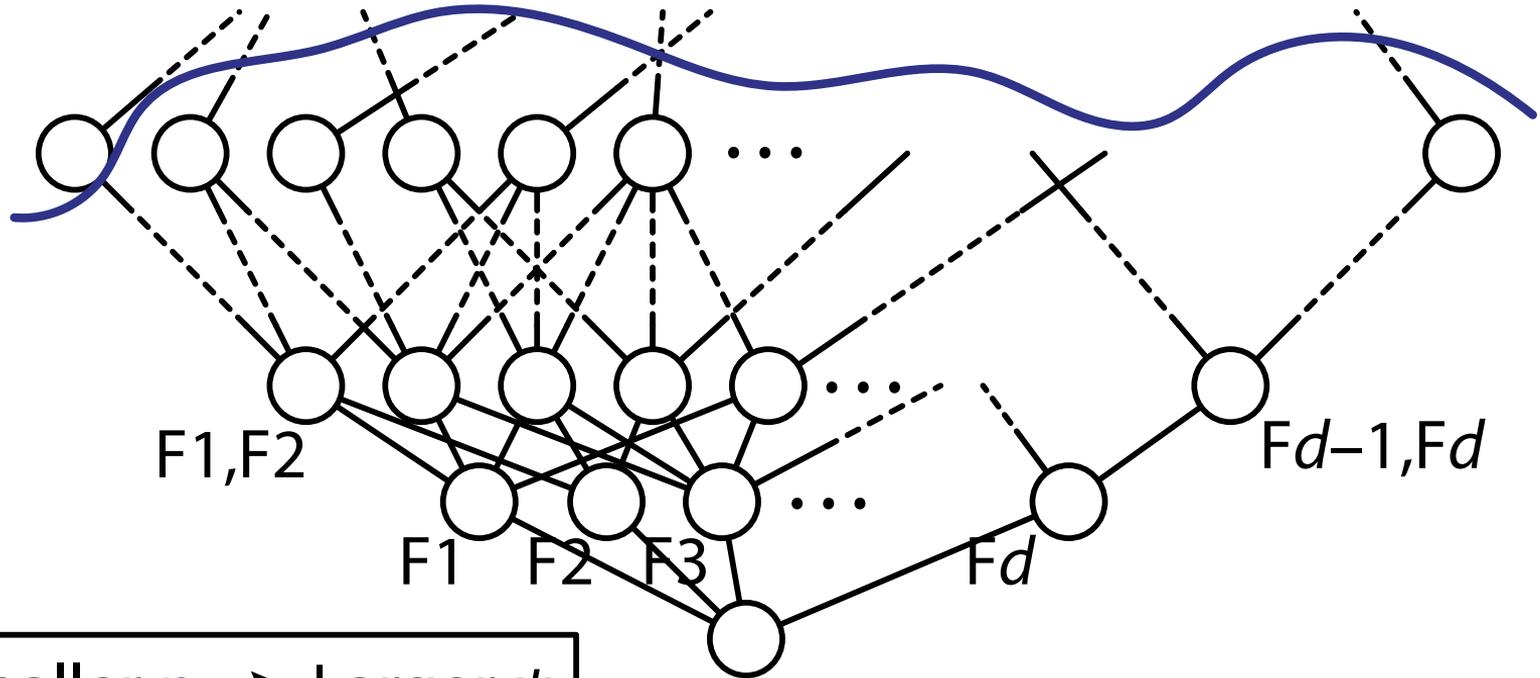


Enumeration Based on Apriori



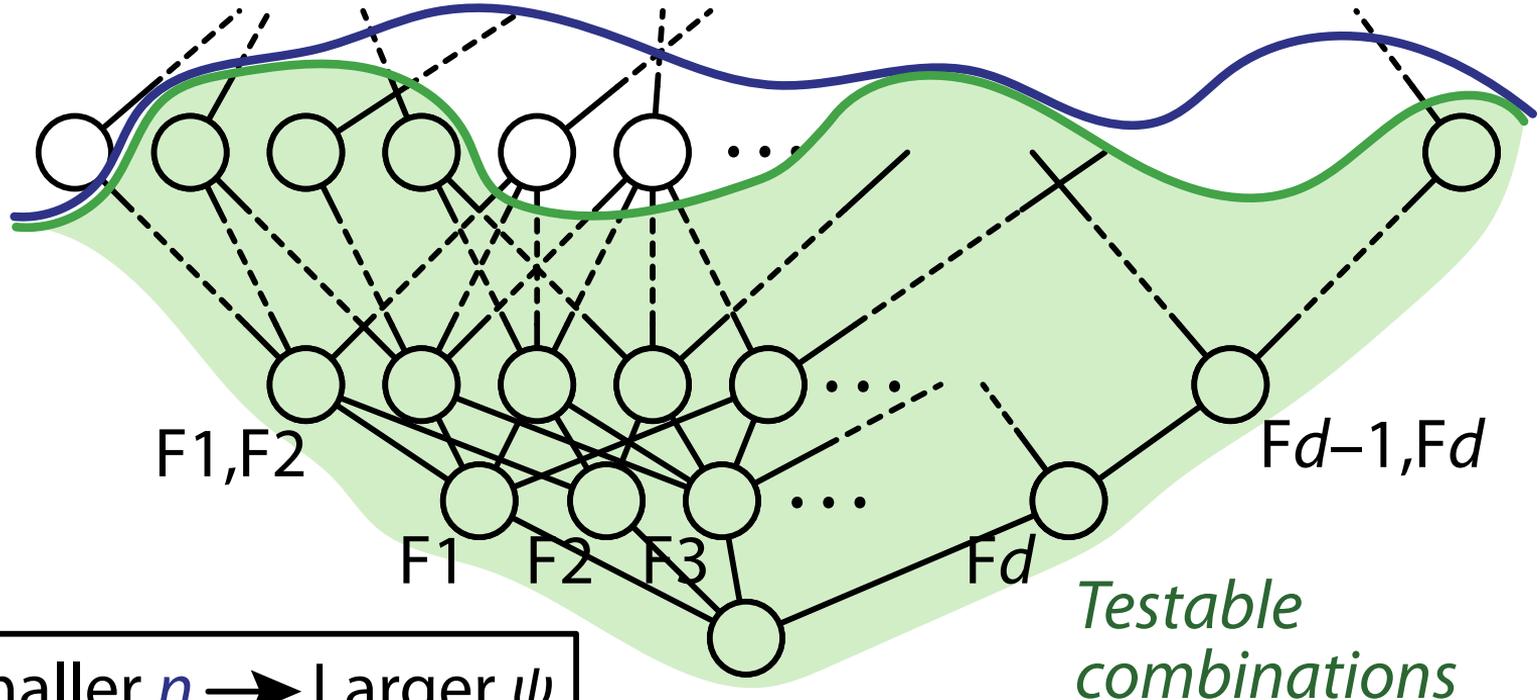
Smaller $\eta \rightarrow$ Larger ψ

Enumeration Based on Apriori



Smaller $\eta \rightarrow$ Larger ψ

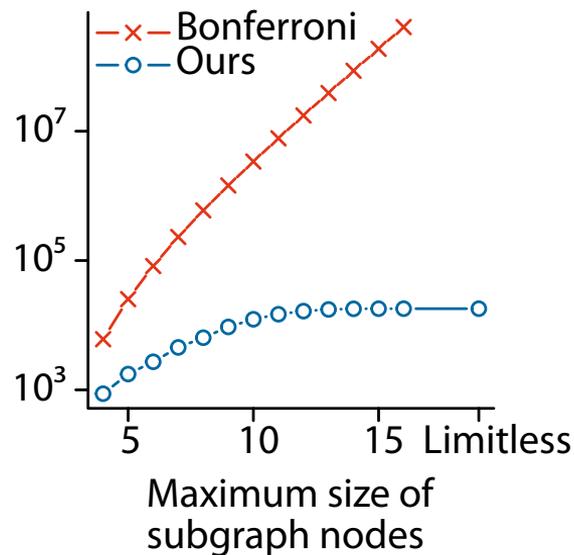
Enumeration Based on Apriori



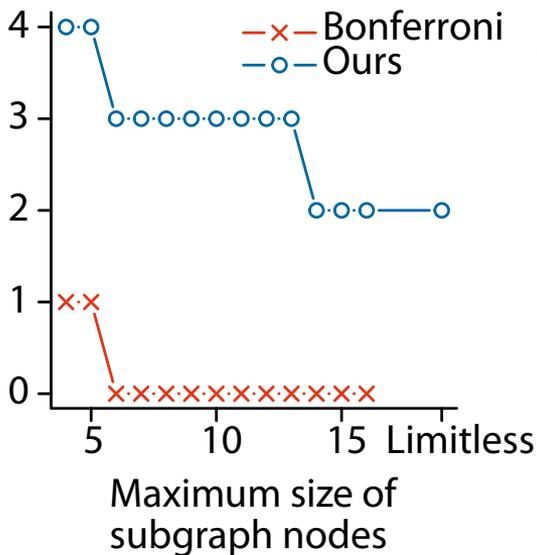
Smaller $\eta \rightarrow$ Larger ψ

Power of Testability

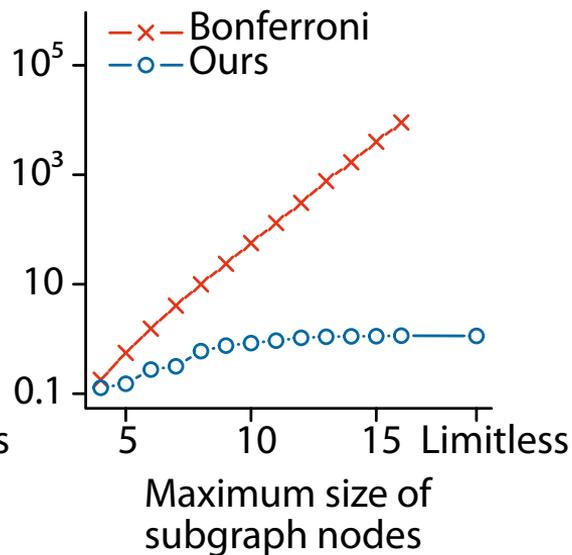
Correction factor



Number of significant subgraphs



Running time (second)



The PTC (Predictive Toxicology Challenge) dataset with 601 chemical compounds

Summary

1. Formulate significance test for each pattern
 - Fisher's exact test is standard, while there are more possibilities
2. Enumerate testable patterns via Tarone's testability + Apriori (DFS)
3. Test each testable pattern

LAMP and WY light

		SNPs (items)	Signal
 Case	ID 1:	0 0 1 1 0 0 1 1 1 0 0 1 1 1 0 0 1 1 1 0	1
	ID 2:	1 1 0 0 1 0 1 1 1 1 0 0 0 1 0 1 0 1 0 0	1
	ID 3:	1 0 1 0 0 0 1 1 1 0 1 0 1 1 0 0 0 0 0 1	1
	ID 4:	1 1 0 1 1 0 1 1 1 1 1 1 1 1 0 1 0 0 1 1	1
<hr/>			
 Control	ID 5:	0 0 1 1 0 0 0 1 1 0 0 0 1 1 1 1 1 0 0 0	0
	ID 6:	0 1 0 1 1 0 1 1 0 0 0 0 1 1 0 0 1 0 1 0	0
	ID 7:	1 0 1 1 0 0 1 0 1 0 0 0 0 0 1 0 1 0 0 0	0
	ID 8:	1 1 0 0 1 0 0 1 0 1 0 1 0 0 0 1 0 1 0 1	0

FACS [Papaxanthos et al. 2016] for Covariates

		SNPs (items)	Cov.	Signal
 Case	ID 1:	0 0 1 1 0 0 1 1 1 0 0 1 1 1 0 0 1 1 1 0	Europe	1
	ID 2:	1 1 0 0 1 0 1 1 1 1 0 0 0 1 0 1 0 1 0 0	Europe	1
	ID 3:	1 0 1 0 0 0 1 1 1 0 1 0 1 1 0 0 0 0 0 1	Asia	1
	ID 4:	1 1 0 1 1 0 1 1 1 1 1 1 1 1 0 1 0 0 1 1	Asia	1
<hr/>				
 Control	ID 5:	0 0 1 1 0 0 0 1 1 0 0 0 0 1 1 1 1 1 0 0 0	Europe	0
	ID 6:	0 1 0 1 1 0 1 1 0 0 0 0 0 1 1 0 0 1 0 1 0	Europe	0
	ID 7:	1 0 1 1 0 0 1 0 1 0 0 0 0 0 0 1 0 1 0 0 0	Asia	0
	ID 8:	1 1 0 0 1 0 0 1 0 1 0 1 0 0 0 0 1 0 1 0 1	Asia	0

FAIS [Llinares-López et al. 2015] for Intervals

		SNPs (items)	Signal
 Case	ID 1:	0 0 1 1 0 0 1 1 1 0 0 1 1 1 0	1
	ID 2:	1 1 0 0 1 0 1 1 1 1 0 0 0 1 0 1 0 1 0 0	1
	ID 3:	1 0 1 0 0 0 1 1 1 0 1 0 1 1 0 0 0 0 0 1	1
	ID 4:	1 1 0 1 1 0 1 1 1 1 1 1 1 1 0 1 0 0 1 1	1
<hr/>			
 Control	ID 5:	0 0 1 1 0 0 0 1 1 0 0 0 1 1 1 1 1 0 0 0	0
	ID 6:	0 1 0 1 1 0 1 1 0 0 0 0 1 1 0 0 1 0 1 0	0
	ID 7:	1 0 1 1 0 0 1 0 1 0 0 0 0 0 1 0 1 0 0 0	0
	ID 8:	1 1 0 0 1 0 0 1 0 1 0 1 0 0 0 1 0 1 0 1	1

FastCHM for Intervals + Cov.

		SNPs (items)	Cov.	Signal
 Case	ID 1:	0 0 1 1 0 0 1 1 1 0 0 1 1 1 0	Europe	1
	ID 2:	1 1 0 0 1 0 1 1 1 1 0 0 0 1 0 1 0 1 0 0	Europe	1
	ID 3:	1 0 1 0 0 0 1 1 1 0 1 0 1 1 0 0 0 0 0 1	Asia	1
	ID 4:	1 1 0 1 1 0 1 1 1 1 1 1 1 1 0 1 0 0 1 1	Asia	1
<hr/>				
 Control	ID 5:	0 0 1 1 0 0 0 1 1 0 0 0 1 1 1 1 1 0 0 0	Europe	0
	ID 6:	0 1 0 1 1 0 1 1 0 0 0 0 1 1 0 0 1 0 1 0	Europe	0
	ID 7:	1 0 1 1 0 0 1 0 1 0 0 0 0 0 1 0 1 0 0 0	Asia	0
	ID 8:	1 1 0 0 1 0 0 1 0 1 0 1 0 0 0 1 0 1 0 1	Asia	1

C-Tarone [Sugiyama & Borgwardt, 2019]

- Find all feature interactions from **continuous** data

Input:

	x						y
	F1	F2	F3	F4	F5	...	Class
ID1	-0.96	-3.03	3.38	2.57	-6.06	...	0
ID2	-1.80	4.45	-4.35	0.82	8.90	...	1
ID3	-3.29	1.39	-4.44	-0.77	2.78	...	1
ID4	-0.53	-1.96	-3.43	-4.42	-3.92	...	0
⋮			⋮				⋮

C-Tarone [Sugiyama & Borgwardt, 2019]

- Find all feature interactions from **continuous** data

Input:

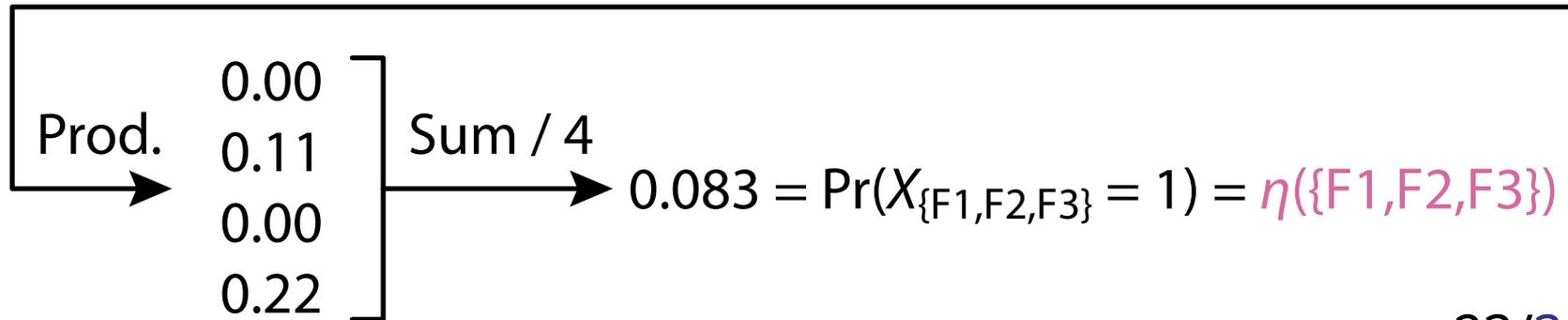
	x						y
	F1	F2	F3	F4	F5	...	Class
ID1	-0.96	-3.03	3.38	2.57	-6.06	...	0
ID2	-1.80	4.45	-4.35	0.82	8.90	...	1
ID3	-3.29	1.39	-4.44	-0.77	2.78	...	1
ID4	-0.53	-1.96	-3.43	-4.42	-3.92	...	0
⋮			⋮				⋮

Output:

{F1}, {F3},
{F2, F5},
{F2, F5, F6}, ...

Use Copula Support [Tatti, 2013]

	F1	F2	F3		R(F1)	R(F2)	R(F3)		$\pi(F1)$	$\pi(F2)$	$\pi(F3)$
x_1	-0.96	-3.03	3.38		2	0	3		0.67	0.00	1.00
x_2	-1.80	4.45	-4.35	Rank →	1	3	1	Norm. →	0.34	1.00	0.34
x_3	-3.29	1.39	-4.44		0	2	0		0.00	0.67	0.00
x_4	-0.53	-1.96	-3.43		3	1	2		1.00	0.34	0.67



Contingency Tables

- For **each** pattern (variable combination), we construct two types of **contingency tables**
 - One is from the **expected** situation under null
 - The other is from the **observed** situation from data
- Significance is assessed by comparison of the two tables
 - Each table is represented as a four-dimensional vector

Expected for p_E	$X_{\mathcal{F}} = 1$	$X_{\mathcal{F}} = 0$	Total
$Y = 1$	$\eta(\mathcal{F}) r_1$	$r_1 - \eta(\mathcal{F}) r_1$	r_1
$Y = 0$	$\eta(\mathcal{F}) r_0$	$r_0 - \eta(\mathcal{F}) r_0$	r_0
Total	$\eta(\mathcal{F})$	$1 - \eta(\mathcal{F})$	1

Observed for p_O	$X_{\mathcal{F}} = 1$	$X_{\mathcal{F}} = 0$	Total
$Y = 1$	$\eta(\mathcal{F}, Y = 1)$	$r_1 - \eta(\mathcal{F}, Y = 1)$	r_1
$Y = 0$	$\eta(\mathcal{F}, Y = 0)$	$r_0 - \eta(\mathcal{F}, Y = 0)$	r_0
Total	$\eta(\mathcal{F})$	$1 - \eta(\mathcal{F})$	1

Significance Test

- The independence $X_{\mathcal{F}} \perp\!\!\!\perp Y$ is translated into:

$$H_0 : D_{\text{KL}}(\mathbf{p}_O, \mathbf{p}_E) = 0, \quad H_1 : D_{\text{KL}}(\mathbf{p}_O, \mathbf{p}_E) \neq 0$$

- \mathbf{p}_E and \mathbf{p}_O are vectorized contingency tables:

$$\mathbf{p}_E = (\eta(\mathcal{F})r_1, \eta(\mathcal{F})r_0, r_1 - \eta(\mathcal{F})r_1, r_0 - \eta(\mathcal{F})r_0)$$

$$\mathbf{p}_O = (\eta(\mathcal{F}, Y=1), \eta(\mathcal{F}, Y=0), r_1 - \eta(\mathcal{F}, Y=1), r_0 - \eta(\mathcal{F}, Y=0))$$

- We apply **G-test**: the statistic $\lambda = 2ND_{\text{KL}}(\mathbf{p}_O, \mathbf{p}_E)$ follows the χ^2 -distribution with the d.f. 1

KL Divergence Bound

- **Theorem** (tight upper bound of KL divergence):

$$D_{\text{KL}}(\mathbf{p}, \mathbf{p}_{\text{E}})$$

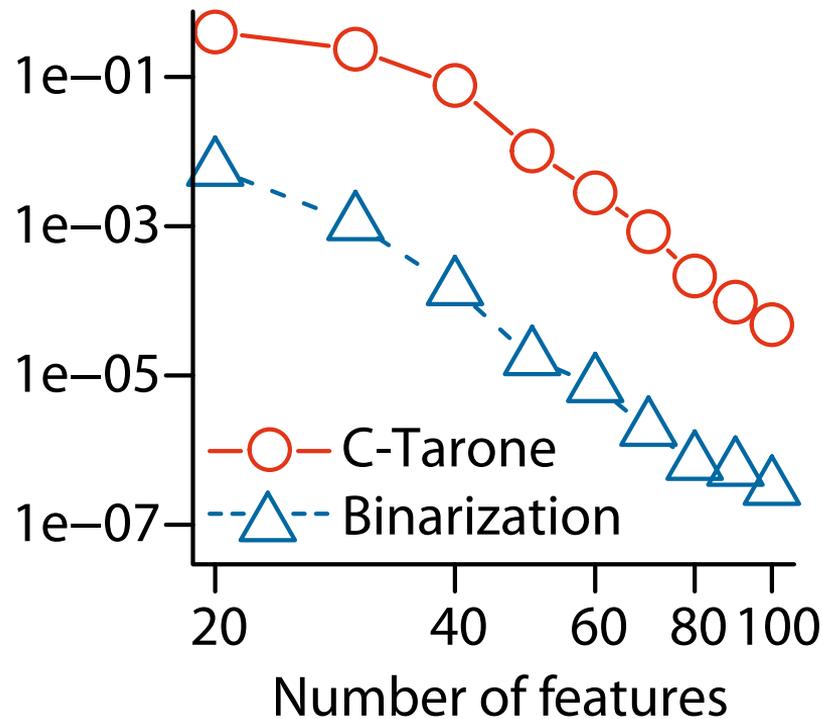
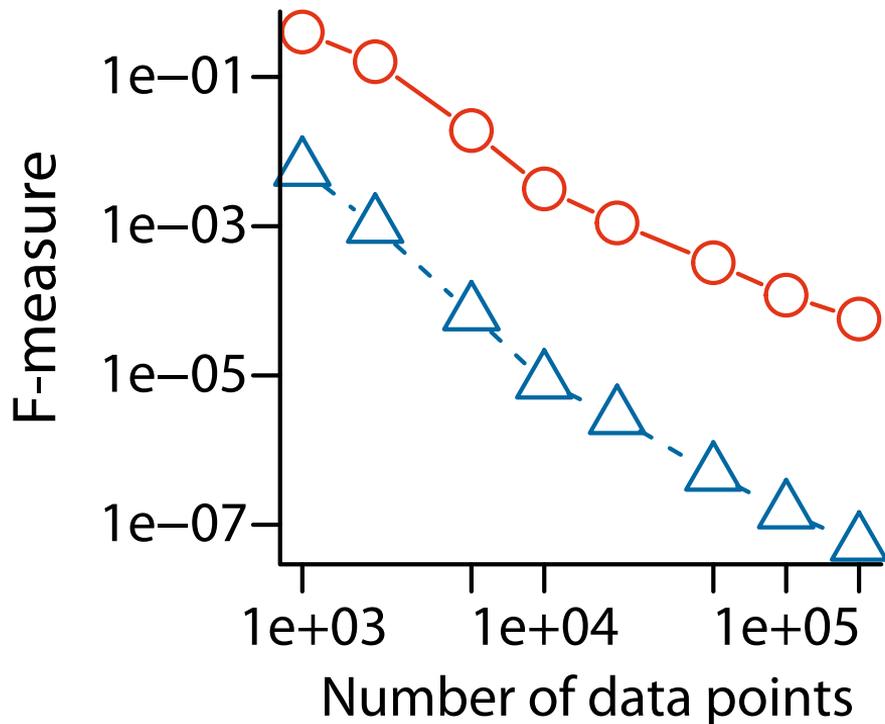
$$< a \log \frac{1}{b} + (b - a) \log \frac{b - a}{(1 - a)b} + (1 - b) \log \frac{1}{(1 - a)}$$

$$- \mathbf{p}_{\text{E}} = (ab, a(1 - b), (1 - a)b, (1 - a)(1 - b)),$$

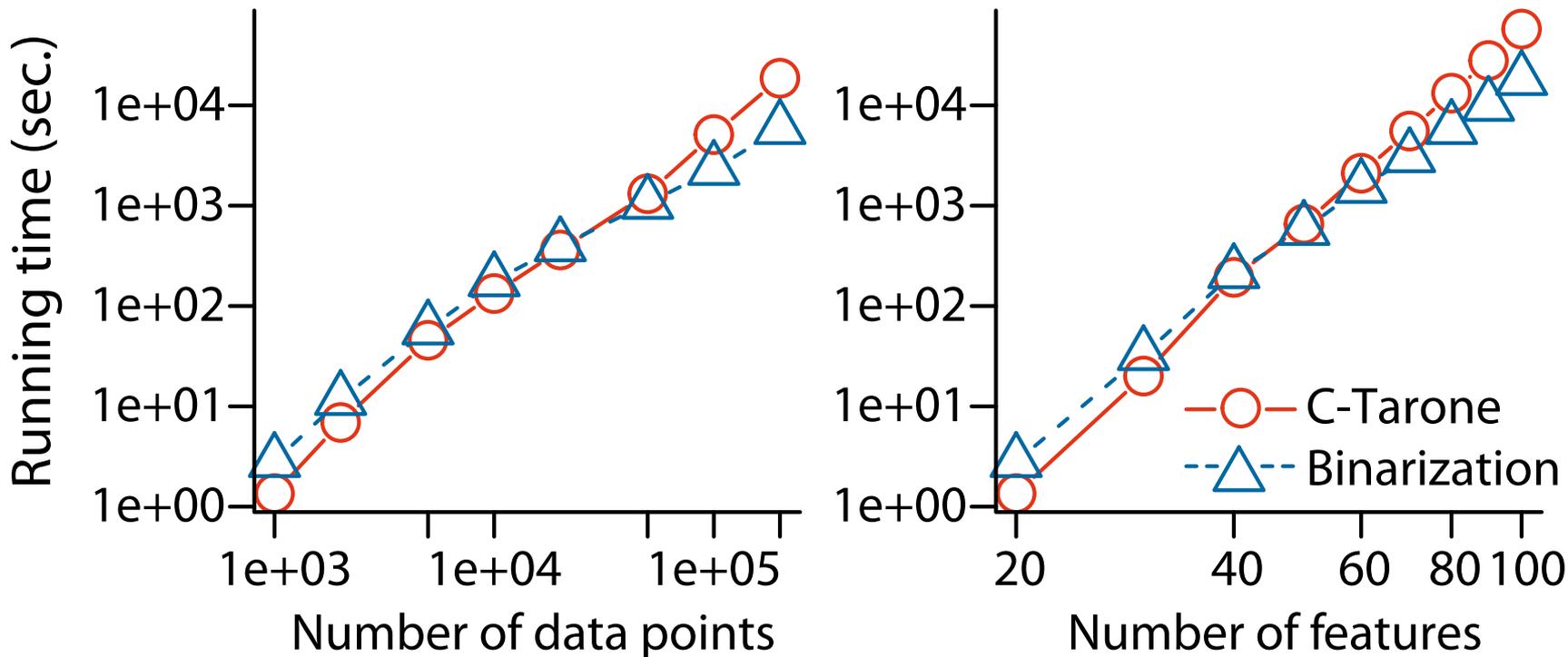
$$\mathbf{p} \in \{ \mathbf{p} \in \mathcal{P} \mid p_1 + p_2 = a, p_1 + p_3 = b \}$$

- The p -value for this upper bound is the minimum achievable p -value

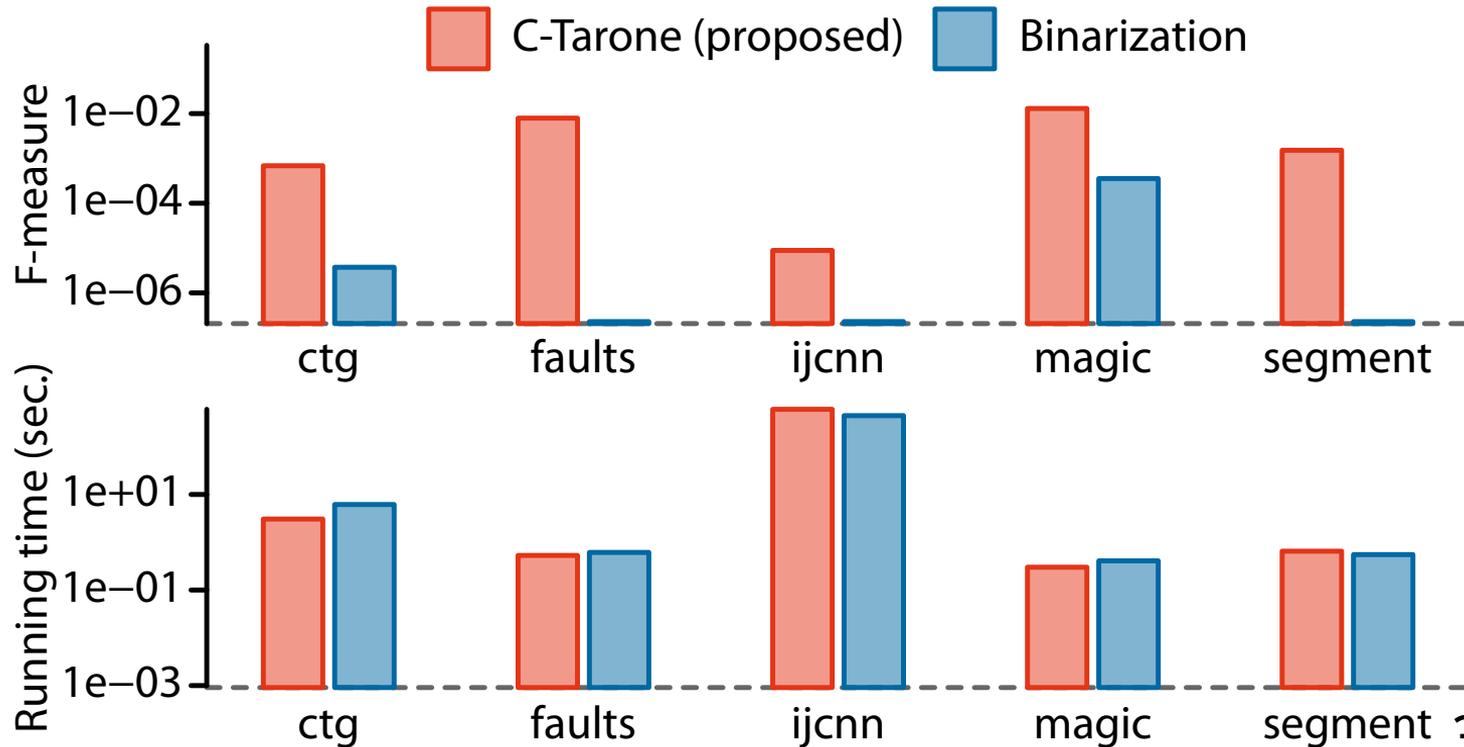
Exp. Results on Synthetic Data



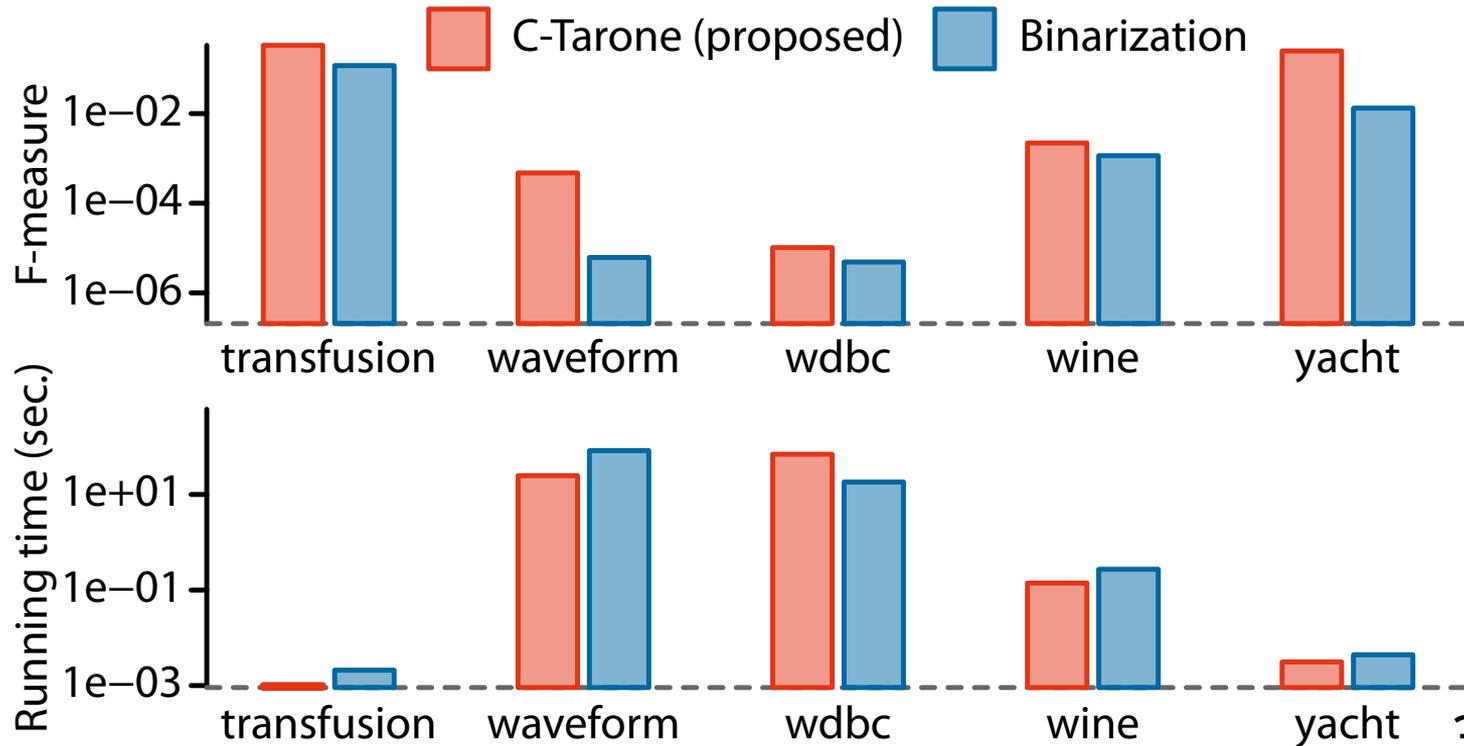
Exp. Results on Synthetic Data



Experimental Results on Real Data



Experimental Results on Real Data



Conclusion

- Significant pattern mining is introduced
 - Find significant interactions while controlling the FWER
 - pattern mining (data mining) + multiple testing correction (statistics)
- Key to solve the problem is Tarone's testability trick
 - This method can be used if the minimum achievable p -value ψ exists
 - If we have the relationship "Smaller $\eta \rightarrow$ Larger ψ ", Apriori can be used to efficiently enumerate testable patterns