

August 15, 2019

IJCAI 2019



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

National Institute of Informatics



# Finding Statistically Significant Interactions between Continuous Features

---

Mahito Sugiyama (National Institute of Informatics)

Karsten Borgwardt (ETH Zürich)

# Our Proposal: *C-Tarone*

---

- Find all feature interactions that are significantly associated with class labels from multivariate data with controlling the FWER

**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
⋮			⋮				⋮

# Our Proposal: *C-Tarone*

- Find all feature interactions that are significantly associated with class labels from multivariate data with controlling the FWER

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}, ...

# Existing Method: Significant Pattern Mining

---

- So far only binary (or discrete) data can be used  
→ Results obtained by SPM via binarization can be uninformative!

**Input:**

	<b>x</b>						<b>y</b>
	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}, ...

# We solve:

---

1. How to assess the significance for a **multiplicative interaction of continuous features**?
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** of the candidate space?
  - The number of possible interactions is  $2^d$  for  $d$  features

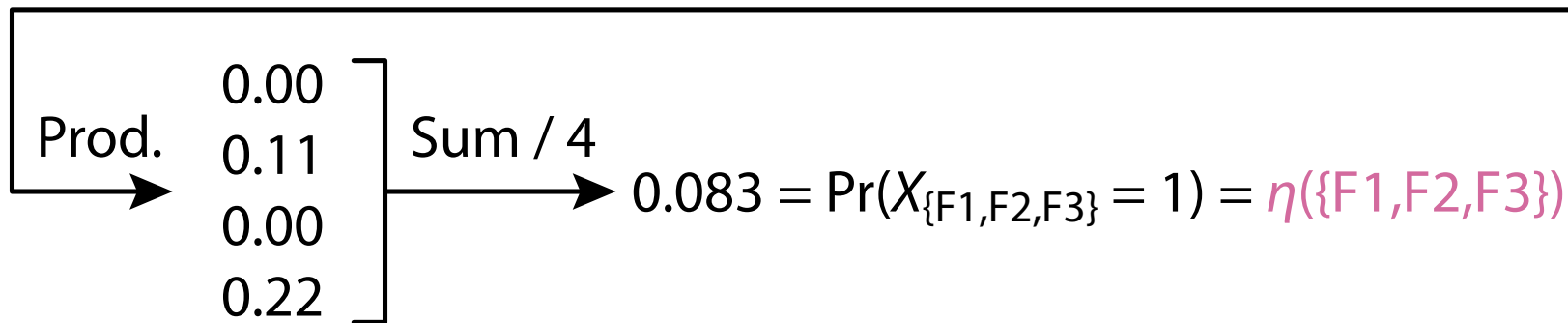
# Problem Formulation

---

- Define  $X_{\mathcal{F}}$  as the **binary random variable** of joint occurrence for a feature combination  $\mathcal{F} = \{F_i, F_{i+1}, \dots, F_{i+k}\}$ 
  - $X_{\mathcal{F}} = 1$  if  $\mathcal{F}$  “occurs”,  $X_{\mathcal{F}} = 0$  otherwise
- Let  $Y$  be an output binary variable
- **Our 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  $Y$
- We need to estimate the probability  $\Pr(X_{\mathcal{F}})$  from data

# Copula Support [Tatti, 2013] for $\Pr(X_{\mathcal{F}} = 1)$

	F1	F2	F3		$R(F1)$	$R(F2)$	$R(F3)$		$\pi(F1)$	$\pi(F2)$	$\pi(F3)$
$\mathbf{x}_1$	-0.96	-3.03	3.38		2	0	3		0.67	0.00	1.00
$\mathbf{x}_2$	-1.80	4.45	-4.35	Rank	1	3	1	Norm.	0.34	1.00	0.34
$\mathbf{x}_3$	-3.29	1.39	-4.44	→	0	2	0	→	0.00	0.67	0.00
$\mathbf{x}_4$	-0.53	-1.96	-3.43		3	1	2		1.00	0.34	0.67



# Contingency Tables

---

Expected (under null) for $\mathbf{p}_E$	$X_{\mathcal{F}} = 1$	$X_{\mathcal{F}} = 0$	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 $\mathbf{p}_O$	$X_{\mathcal{F}} = 1$	$X_{\mathcal{F}} = 0$	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 the condition:

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

- $\mathbf{p}_E$  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  $\chi^2$ -distribution with the d.f. 1

# Multiple Testing Correction

---

- The **FWER** should be controlled
  - Probability that at least one feature combination is a false positive
  - If we naively test all combinations,  $a2^d$  **false positives** could occur!!
- We use **Tarone's testability trick**, which requires the minimum achievable  $p$ -value  $\psi(\mathcal{F})$  for  $\mathcal{F}$
- **Theorem** (tight upper bound of KL divergence):

$$D_{\text{KL}}(\mathbf{p}, \mathbf{p}_E) < a \log \frac{1}{b} + (b - a) \log \frac{b - a}{(1 - a)b} + (1 - b) \log \frac{1}{(1 - a)}$$

- $\mathbf{p}_E = (ab, a(1 - b), (1 - a)b, (1 - a)(1 - b))$ ,  
 $\mathbf{p} \in \{ \mathbf{p} \in \mathcal{P} \mid p_1 + p_2 = a, p_1 + p_3 = b \}$

# Tarone's Testability Trick

---

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

# Tarone's Testability Trick

---

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

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

# Tarone's Testability Trick

---

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

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

Testable  
combinations

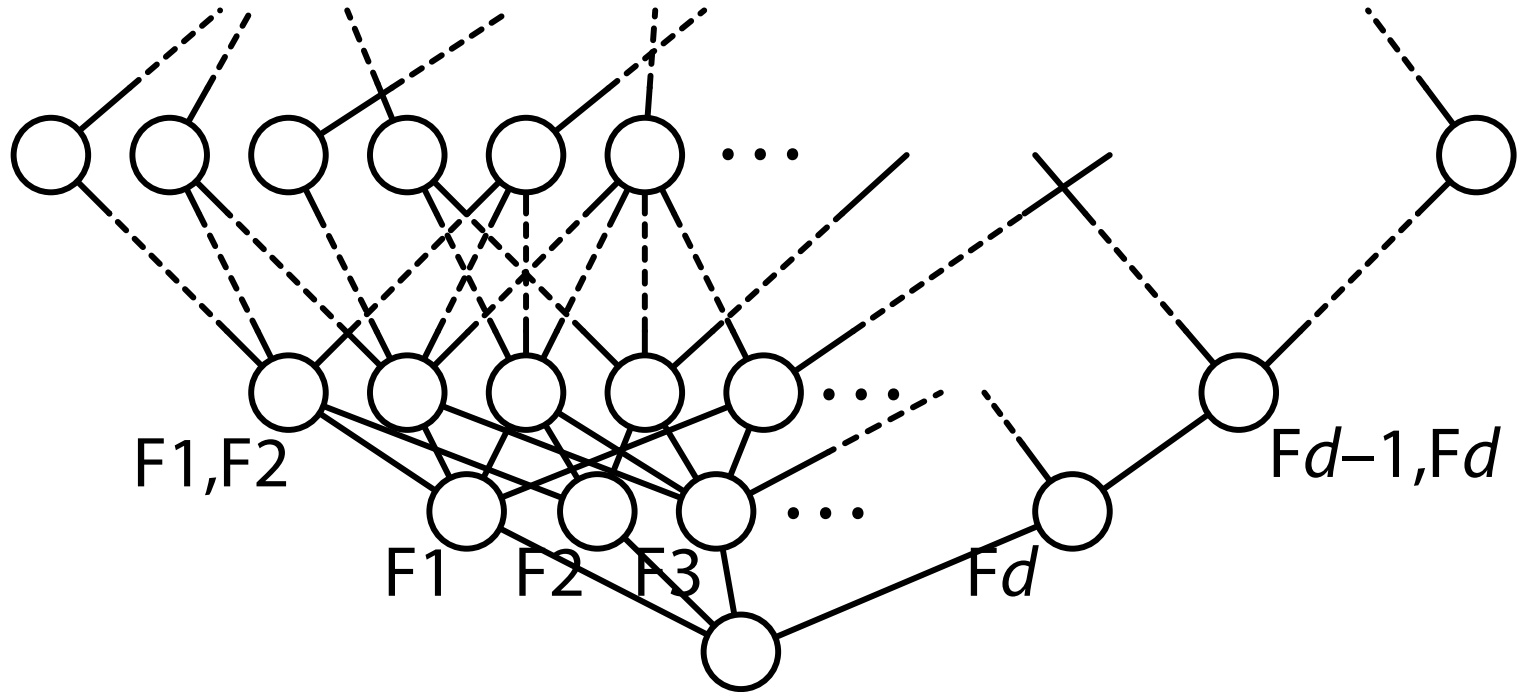
Untestable  
combinations

→ Prune without testing

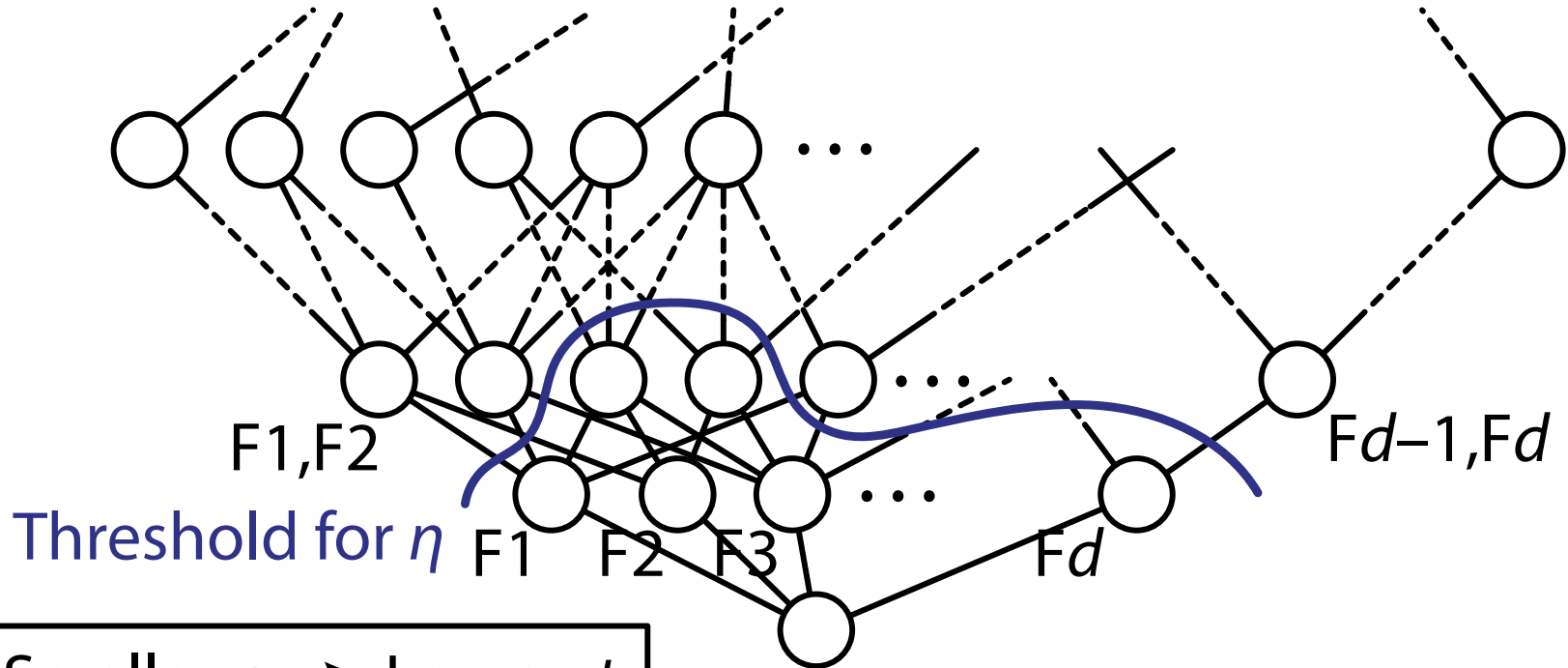
$\mathcal{F}_i$  is significant if:  $p\text{-value}(\mathcal{F}_i) < \alpha / m$  — Correction factor

# Enumeration Algorithm Based on Apriori

---

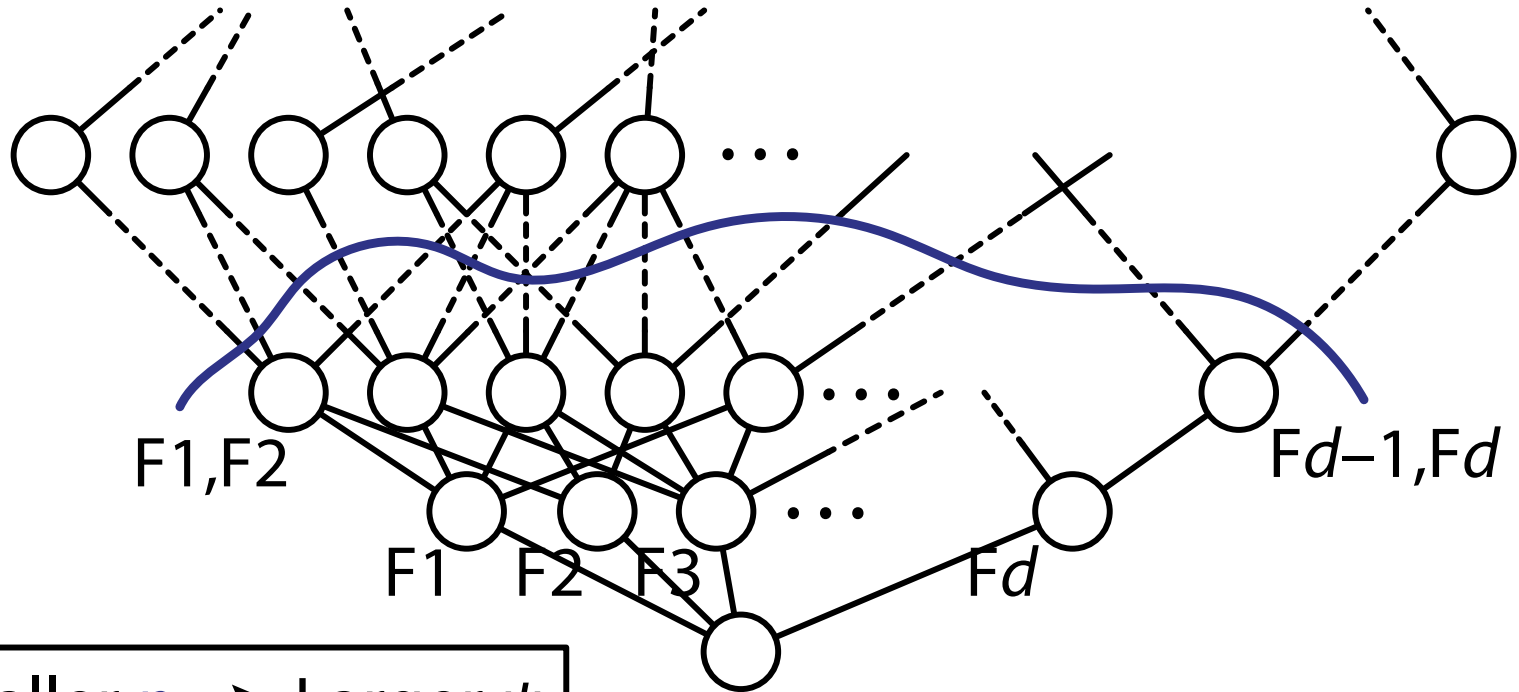


# Enumeration Algorithm Based on Apriori



# Enumeration Algorithm Based on Apriori

---

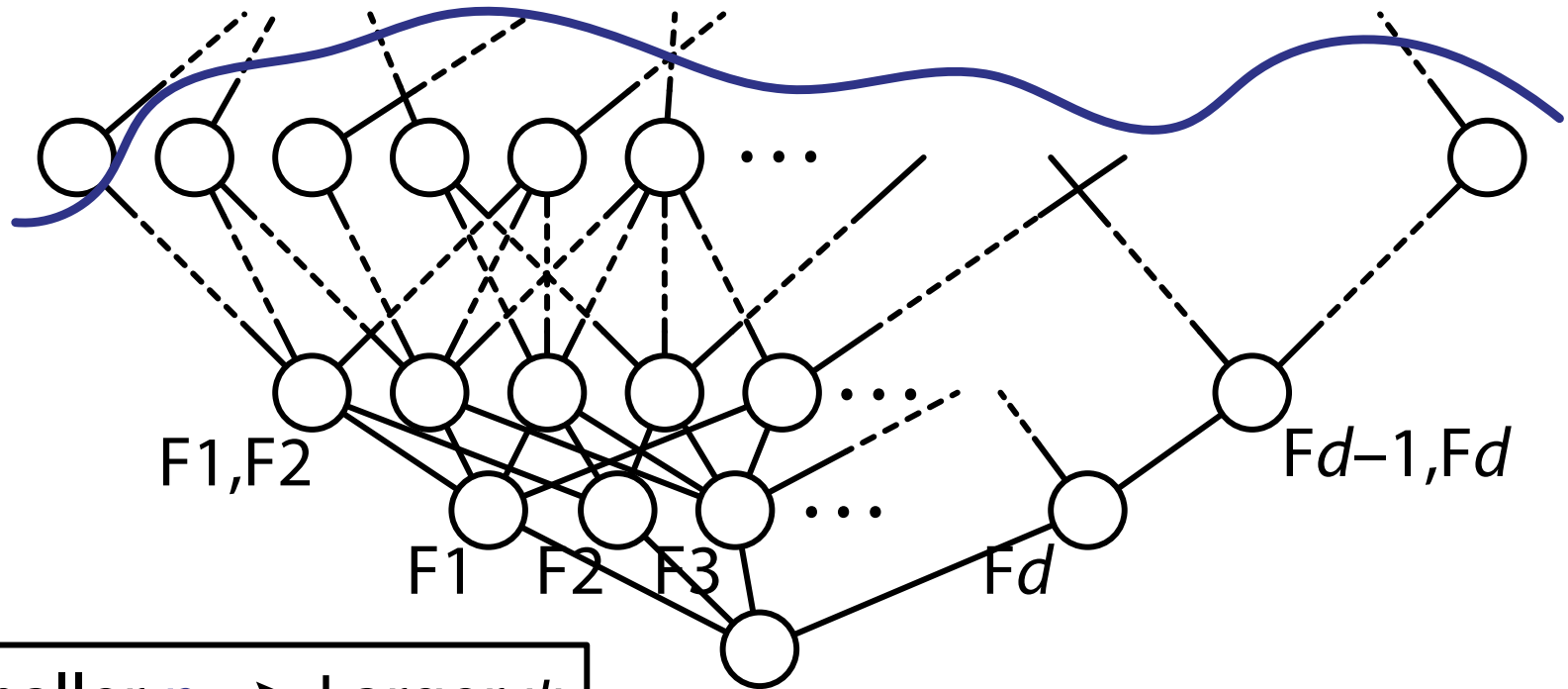


Smaller  $\eta \rightarrow$  Larger  $\psi$



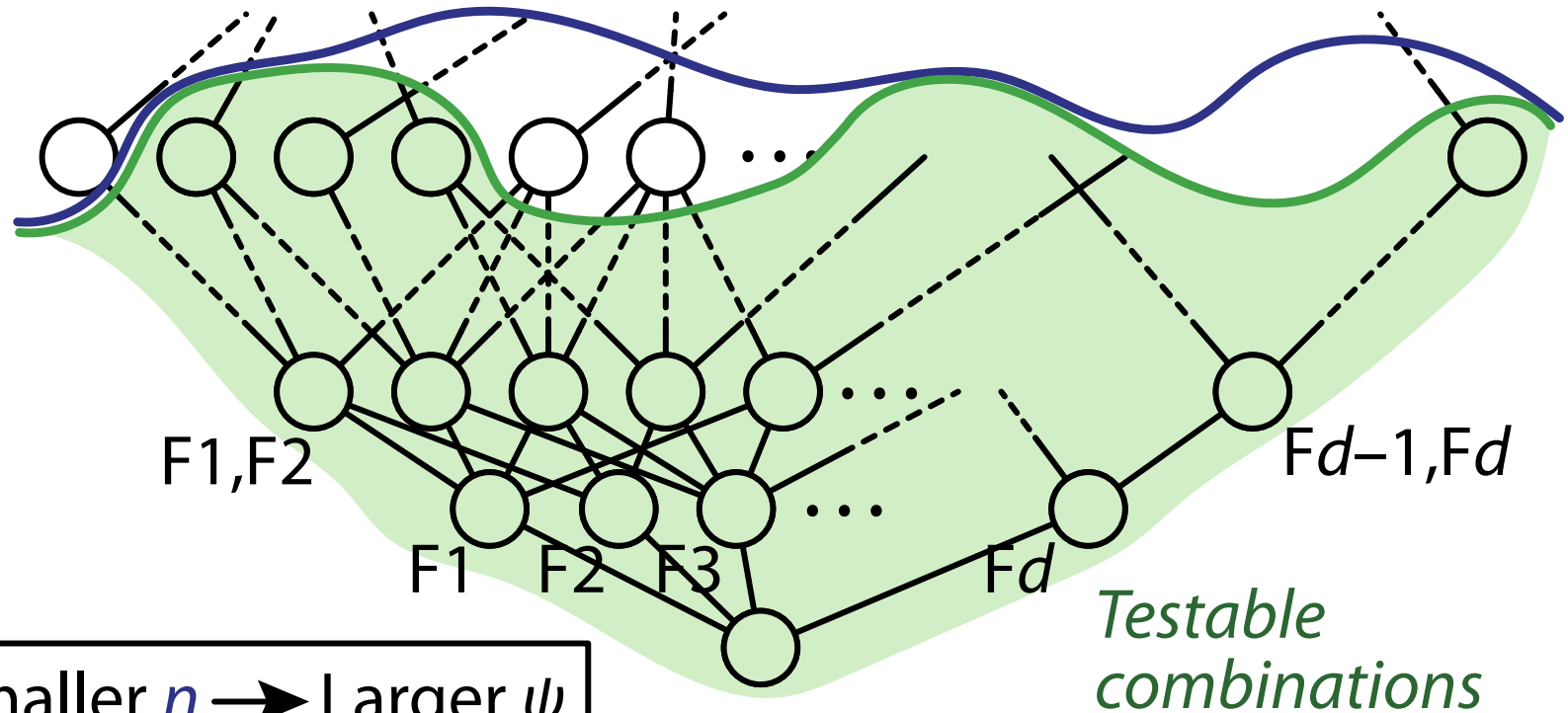
# Enumeration Algorithm Based on Apriori

---

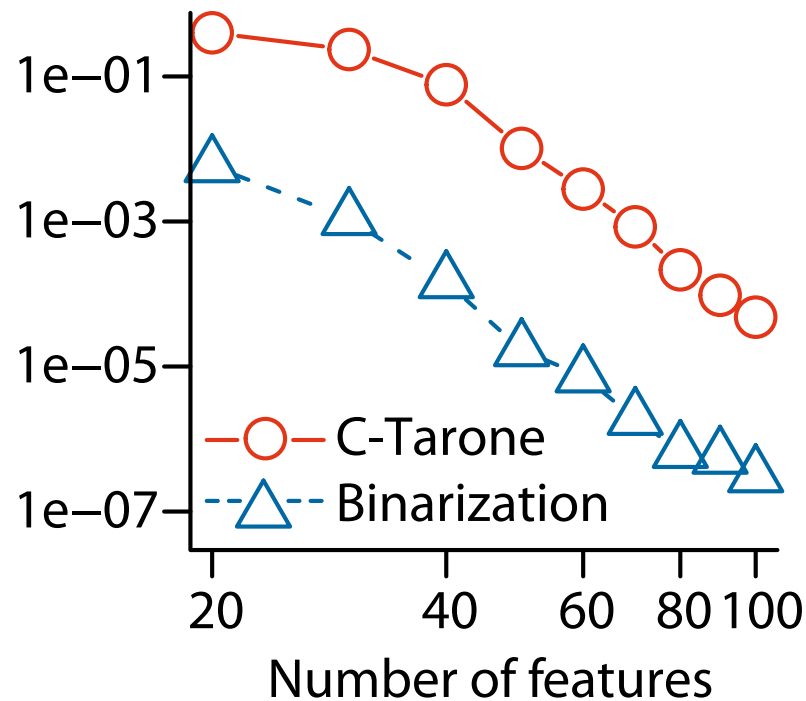
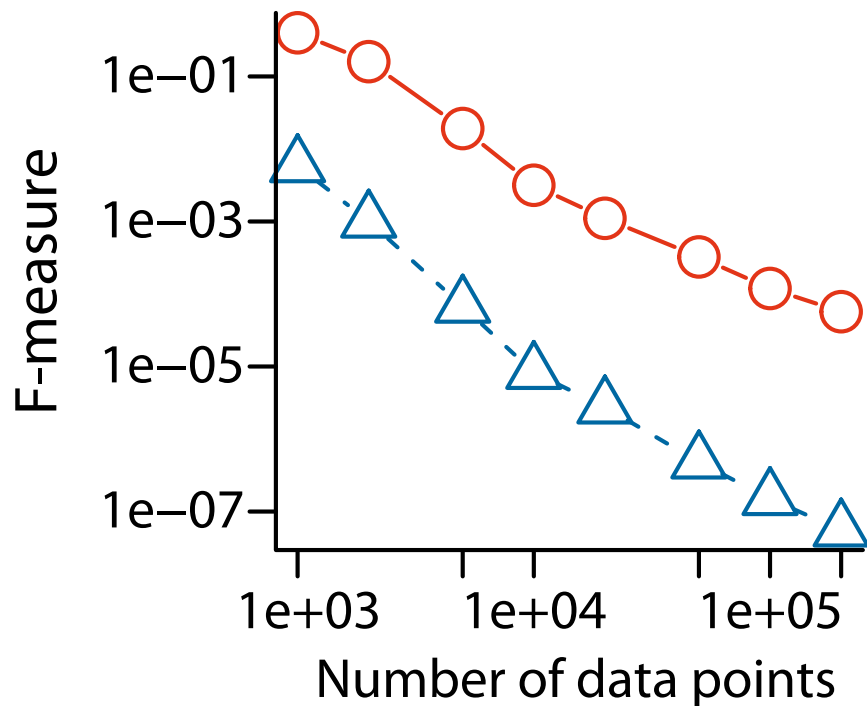


Smaller  $\eta \rightarrow$  Larger  $\psi$

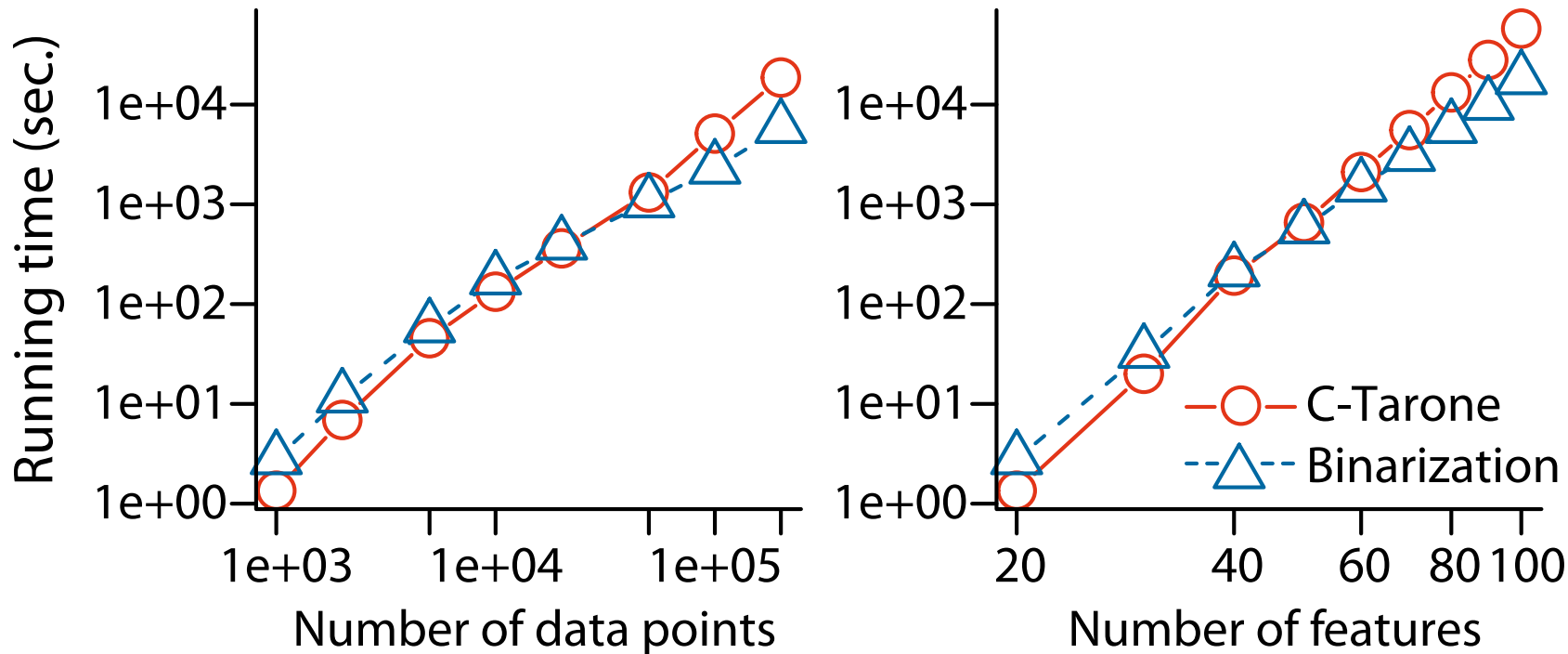
# Enumeration Algorithm Based on Apriori



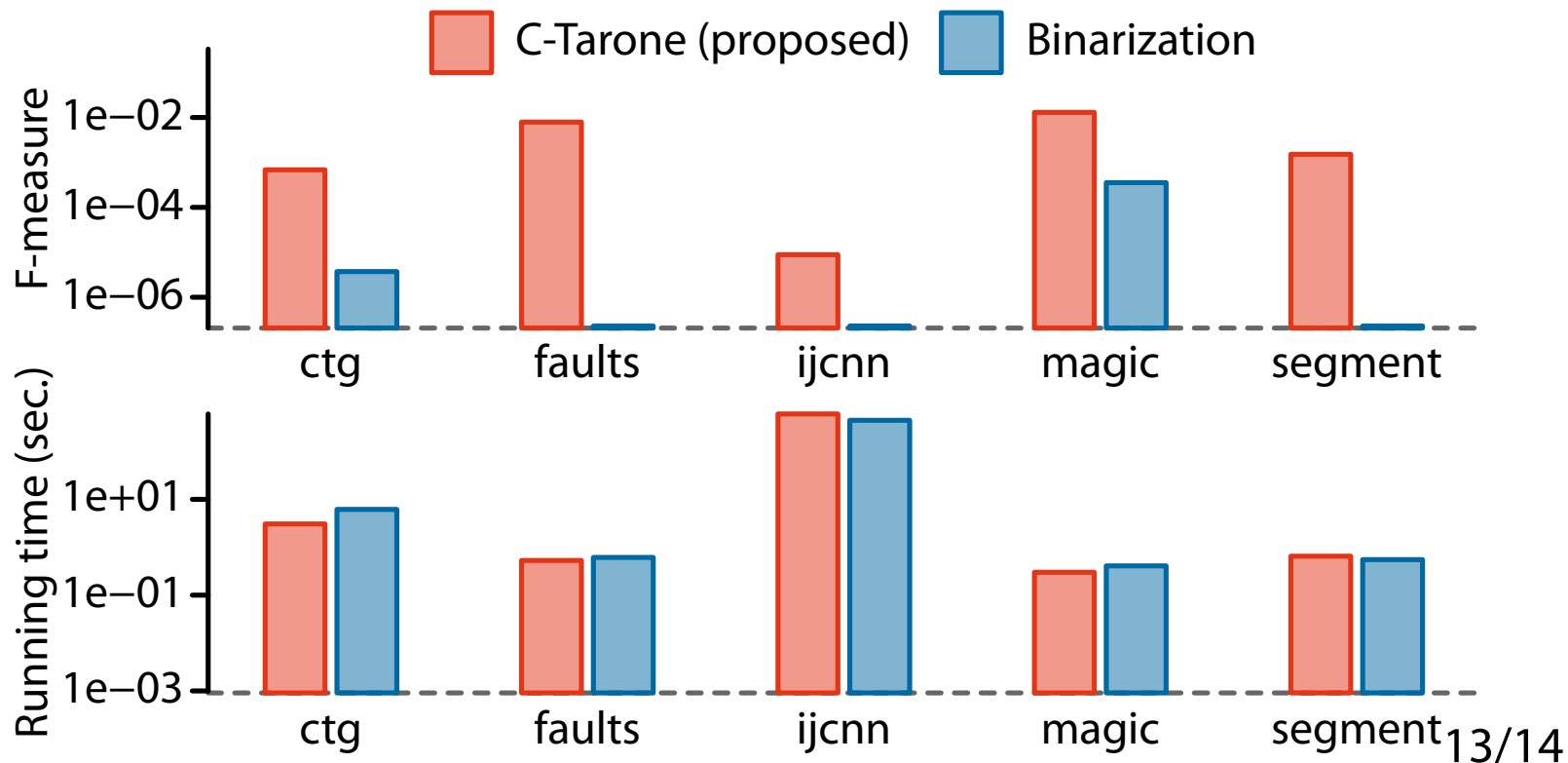
# Experimental Results on Synthetic Data



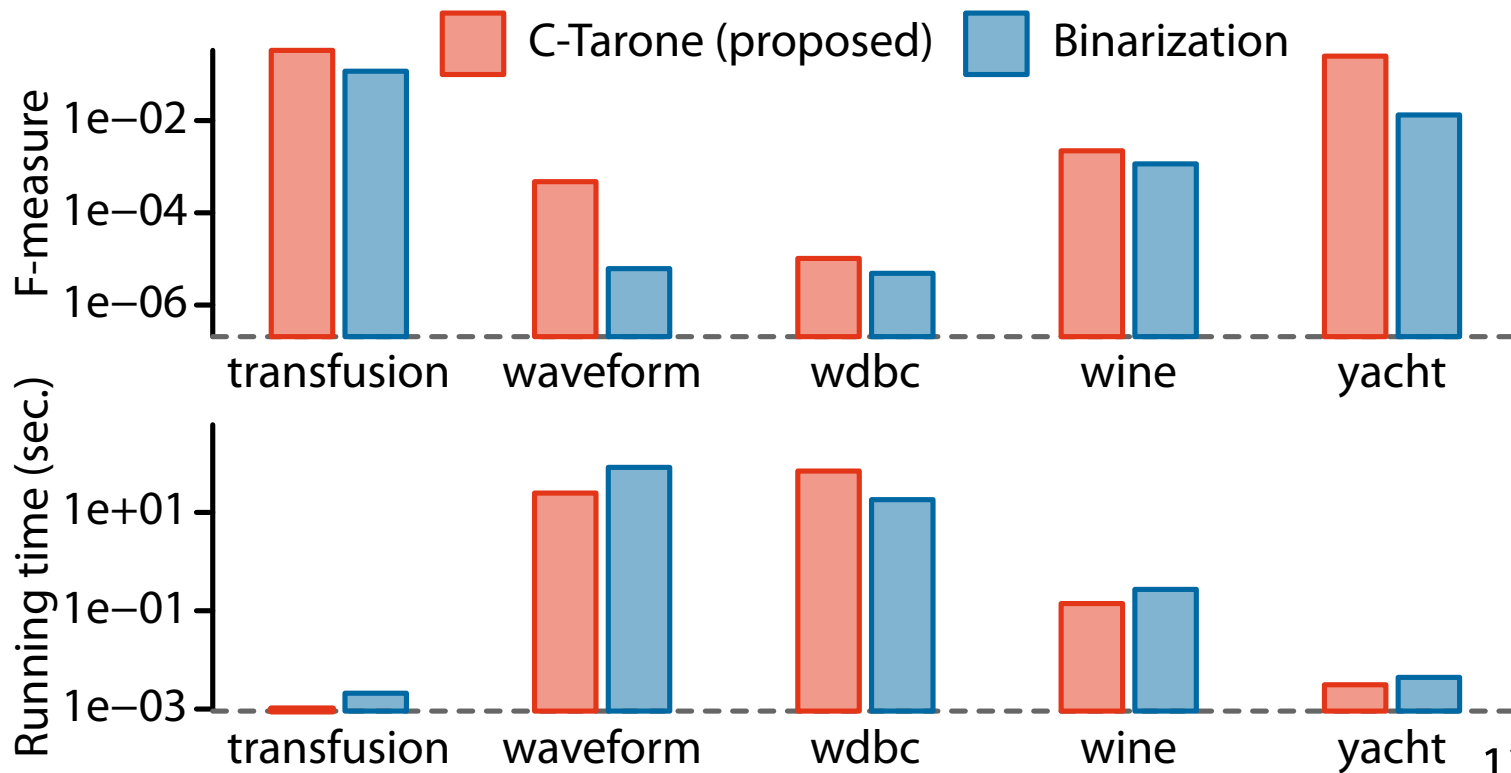
# Experimental Results on Synthetic Data



# Experimental Results on Real Data



# Experimental Results on Real Data



# Conclusion

---

- We have proposed **C-Tarone**, a solution to the open problem of finding *all* multiplicative interactions between **continuous** features significantly associated with an output variable
  - Significance is rigorously controlled for multiple testing
- Our work opens the door to many applications of searching significant feature combinations, in which the data is not adequately described by binary features