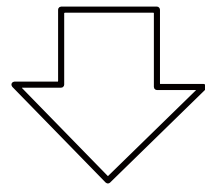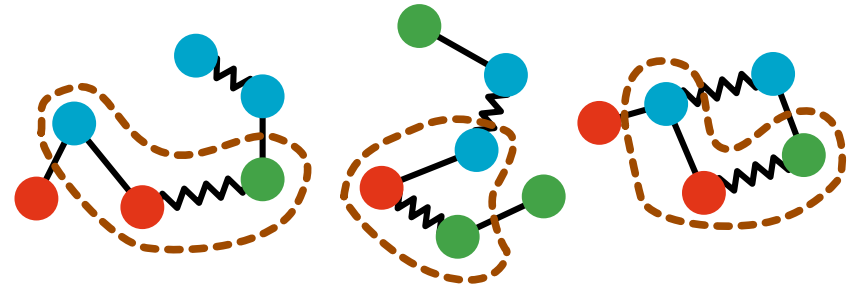# Multiple Testing Correction in Graph Mining

Mahito Sugiyama (Osaka University, JST PRESTO)

Joint work with Felipe Llinares López[1], Niklas Kasenburg[2], Karsten Borgwardt[1] ([1]ETH Zürich, [2]Univ. Copenhagen)
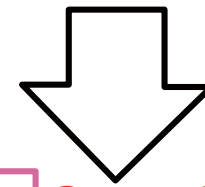
# Binary data

ID a b c d e f g h i j
1  0 0 1 1 0 0 1 1 1 0
2  1 1 0 1 1 0 1 1 1 0
3  1 0 1 1 0 0 1 1 1 0
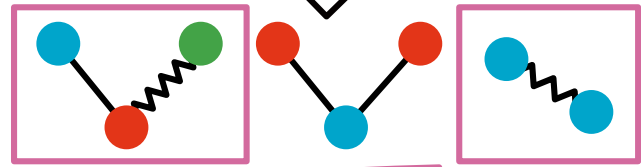4  1 1 0 0 1 0 0 1 0 1
5  1 1 0 1 1 0 1 1 1 0

# Graphs



Pattern mining

{a, b, e}  {d, g, h, i}



Support:       3          4                           3        2        3

(Statistically) Significant patterns
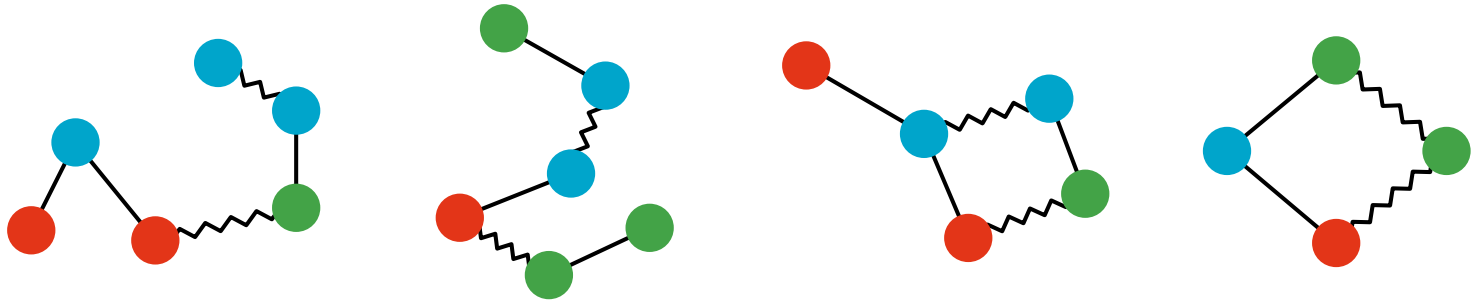($P$ value $< 0.05$)

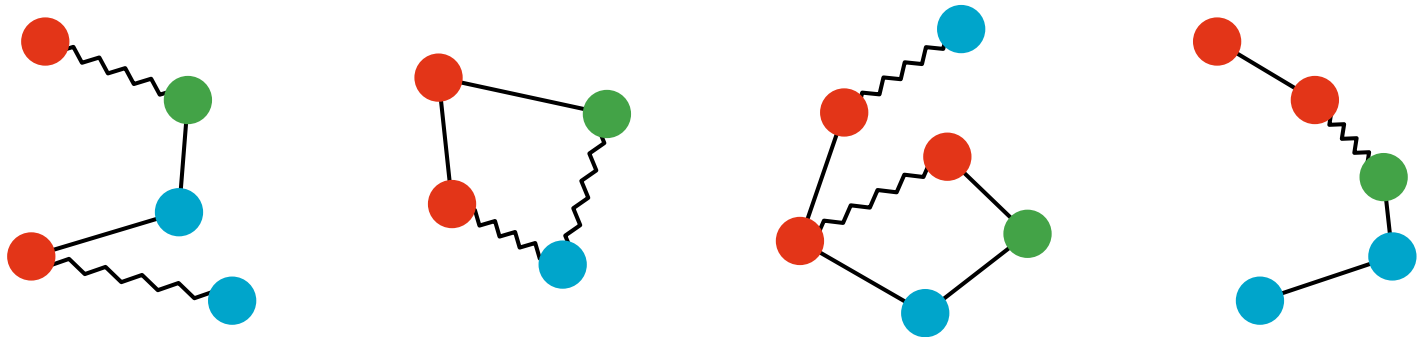$P$ value:     0.06       0.01                      0.02     0.07     0.02

The $P$ value is crucial for scientific discovery!
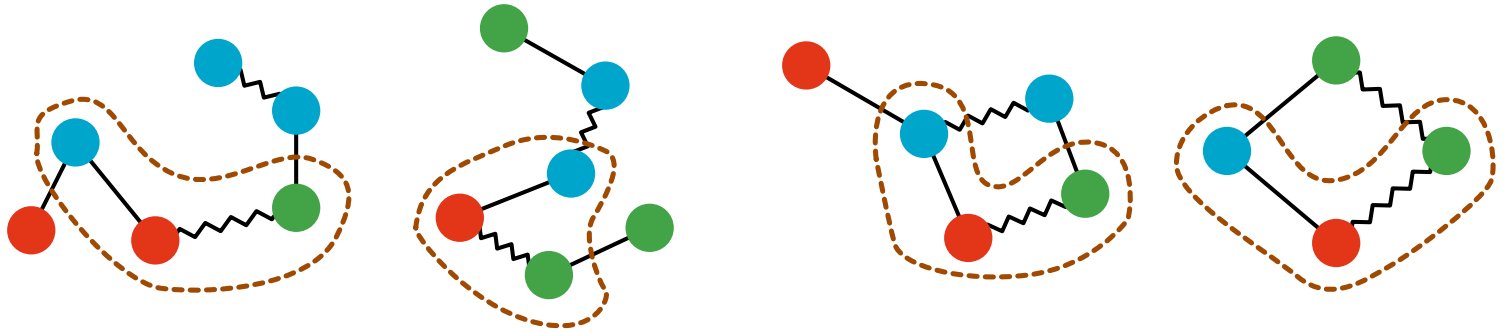
# Find Subgraphs



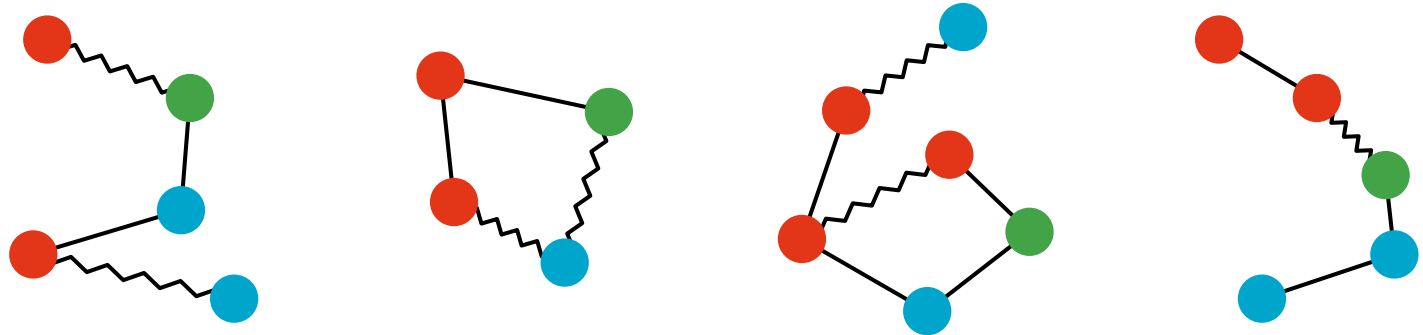Active

Inactive

# Find Subgraphs

Active

Inactive

# Find Subgraphs



Active

Inactive

# Find Subgraphs



Active

Inactive

# Hypothesis Test for Each Subgraph

|  | Alternative hypothesis is true | Null hypothesis is true |
|---|---|---|
| Declared significant | True Positive | False Positive (Type I Error) |
| Declared non-significant | False Negative (Type II Error) | True Negative |

**Null hypothesis**: The occurence of the subgraph is independent from the activity

**Alternative hypothesis**: The occurence of the subgraph is associated with the activity

# Testing the Independence of Subgraph

- Given two sets of graphs $\mathcal{G}$ and $\mathcal{G}'$
  - $|\mathcal{G}| = n$, $|\mathcal{G}'| = n'$ ($n \leq n'$)

- The *P value* of each subgraph $H \sqsubseteq G$ with $G \in \mathcal{G} \cup \mathcal{G}'$ is determined by the Fisher's exact test

|  | Occ. | Non-occ. | Total |
|---|---|---|---|
| $\mathcal{G}$ | $x$ | $n - x$ | $n$ |
| $\mathcal{G}'$ | $x'$ | $n' - x'$ | $n'$ |
| Total | $x + x'$ | $(n - x) + (n' - x')$ | $n + n'$ |

# Fisher's Exact Test

- The probability $q(x)$ of obtaining $x$ and $x'$ is given by the hypergeometric distribution:

$$q(x) = \binom{n}{x}\binom{n'}{x'} \bigg/ \binom{n+n'}{x+x'}$$

|  | Occ. | Non-occ. | Total |
|---|---|---|---|
| $\mathcal{G}$ | $x$ | $n - x$ | $n$ |
| $\mathcal{G}'$ | $x'$ | $n' - x'$ | $n'$ |
| Total | $x + x'$ | $(n - x) + (n' - x')$ | $n + n'$ |

$P$ value
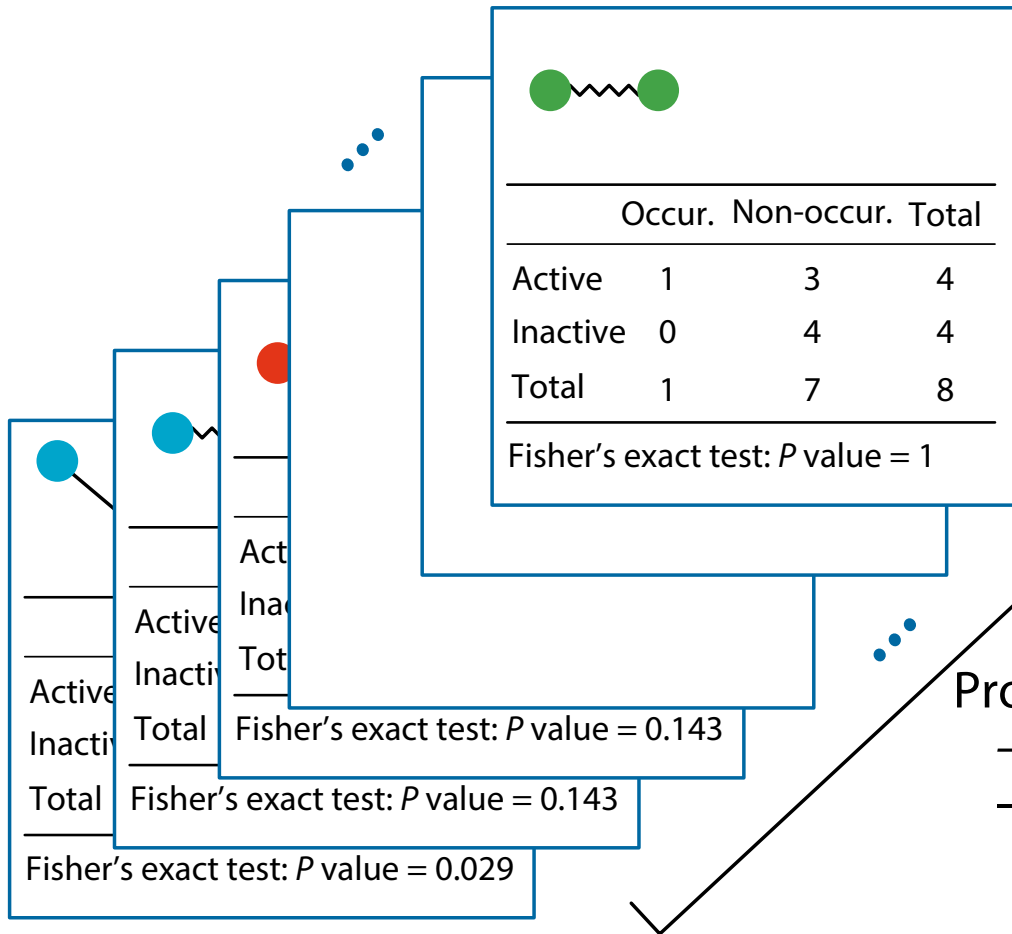
$x_{min} = \max\{0, x + x' - n'\}$    $x_{max} = \min\{x + x', n\}$

# Multiple Testing



Task: Detect all significant subgraphs

We need multiple testing correction! Otherwise, too many false positives:

$$\text{FWER} = 1 - (1 - \alpha)^m$$

$m$ subgraphs

Problems:
  – $m$ is massive
  – The significance level $\alpha / m$ in Bonferroni correction becomes too conservative

|  | Occur. | Non-occur. | Total |
|---|---|---|---|
| Active | 1 | 3 | 4 |
| Inactive | 0 | 4 | 4 |
| Total | 1 | 7 | 8 |

Fisher's exact test: $P$ value = 1

Fisher's exact test: $P$ value = 0.143

Fisher's exact test: $P$ value = 0.143

Fisher's exact test: $P$ value = 0.029

Active

Inactive

# Counting the Frequency of Subgraphs

# Counting the Frequency of Subgraphs

Frequency

$f($  $) = 7$

# Counting the Frequency of Subgraphs

Frequency

$f(\ \ ) = 6$

# The Minimum *P* Value

- The minimum achievable *P* value for the frequency $f(H)$ of a subgraph $H$ is

$$P_{\min} = \binom{n}{f(H)} \bigg/ \binom{n+n'}{f(H)}$$

|          | Occ.   | Non-occ.          | Total   |
|----------|--------|-------------------|---------|
| Active   | $f(H)$ | $n - f(H)$        | $n$     |
| Inactive | $0$    | $n'$              | $n'$    |
| Total    | $f(H)$ | $(n - f(H)) + n'$ | $n + n'$ |

Most biased case $(f(H) < n)$



Minimum *P* value

$x_{\min} = \max\{0, f(H) - n'\}$

$x_{\max} = \min\{f(H), n\}$

# Testability

- The minimum achievable *P* value
  for the frequency $f(H)$ of a subgraph $H$ is

$$P_{\min} = \binom{n}{f(H)} \Big/ \binom{n + n'}{f(H)}$$

- Tarone (1990) pointed out (and Terada et al. (2013) revisited):
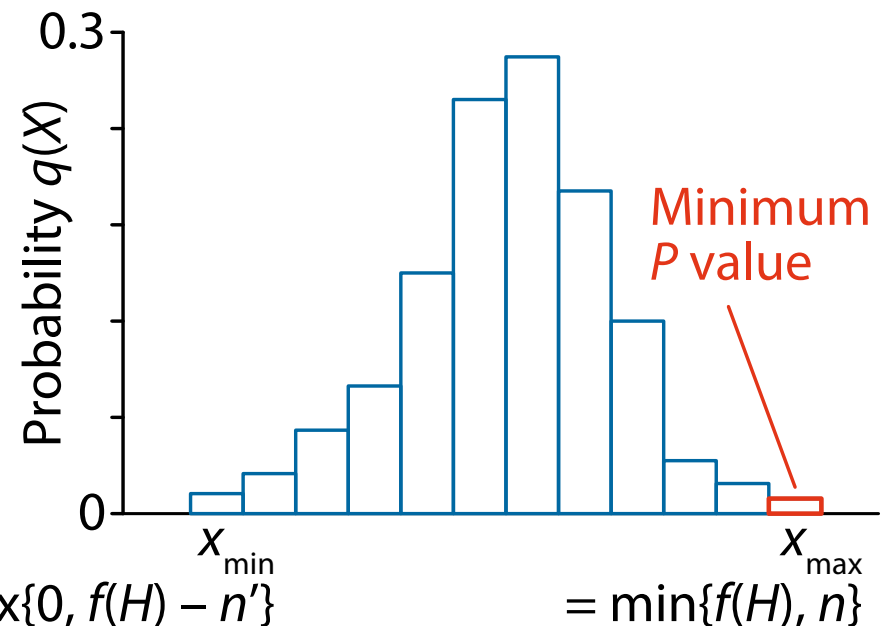  *For a hypothesis H, if its minimum P value is smaller than the significance threshold, this is untestable and we can ignore it*
  - Untestable hypotheses (subgraphs) do not increase the FWER
  - The Bonferroni factor reduces to the number of testable hypotheses

# Finding the Optimal Correction Factor

- $m(k)$: # of subgraphs whose minimum $P$ values $< \alpha/k$
  - $k$: the correction factor, $\alpha/k$: the corrected significance level

- For each $k$, FWER is controlled as (Tarone 1990):

$$\text{FWER} \leq m(k)\frac{\alpha}{k} = \frac{m(k)}{k}\alpha$$

- Our task:
  - Find the smallest $k$ while controlling FWER $\leq \alpha$
    - Coincides with the "root" $k_{rt}$ of the function $m(k) - k$
    - $m(k) \leq k$ for all $k \geq k_{rt}$ and $m(k) > k$ for all $k < k_{rt}$
  - Enumerate testable subgraphs whose min. $P$ values $< \alpha/k_{rt}$

# Testable Subgraphs

# Testable Subgraphs

# Testable Subgraphs



$k = 9$,     $m(9) = 4$

Frequency is large

Minimum $P$ value

$\alpha / 9$

Significance level

Untestable subgraphs

Testable subgraphs

# Testable Subgraphs



$k = 8, \quad m(8) = 6$

Frequency is large

Minimum $P$ value

Significance level

$\alpha / 8$

Untestable subgraphs

Testable subgraphs

# Testable Subgraphs



$k = 7, \quad m(7) = 8$

Frequency is large

Minimum $P$ value

$\alpha / 7$

Significance level

Untestable subgraphs

Testable subgraphs

# Testable Subgraphs

# Subgraphs Are Testable Iff Frequent

- Our task:

Find $k$ such that

(# of subgraphs whose minimum $P$ values $< \alpha/k$) $= k$

$$\Downarrow$$

Find $\sigma$ such that

(# of subgraphs whose frequency $\geq \sigma$) $= \alpha/\psi(\sigma)$

Testable subgraphs = Frequent subgraphs

# Use Frequent Subgraph Mining

- Testable subgraphs can be enumerated by frequent subgraph mining algorithms

- **Proposition:**
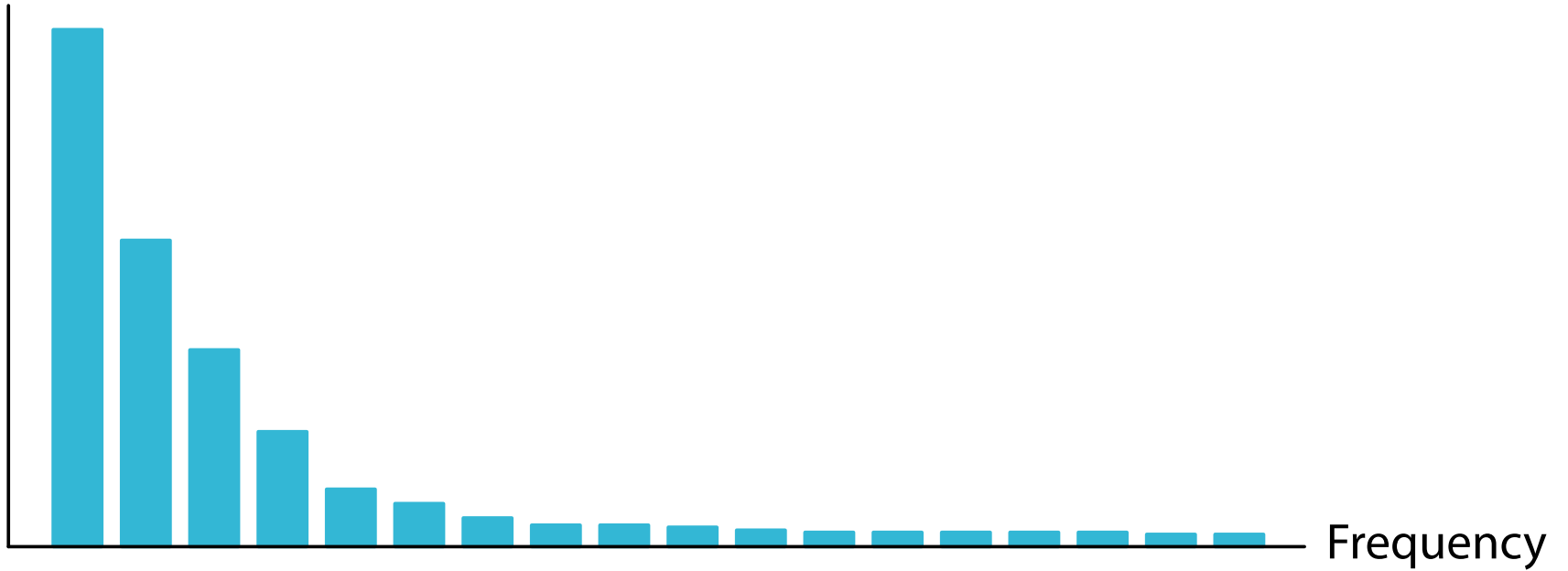  The set of testable subgraphs $\tau(\mathcal{H})$ coincides with the set of frequent subgraphs with the threshold $\sigma_{rt}$ s.t.

    # of subgraphs with minfreq $\sigma_{rt} - 1 > \alpha / \psi(\sigma_{rt} - 1)$,

    # of subgraphs with minfreq $\sigma_{rt} \leq \alpha / \psi(\sigma_{rt})$,

  – $\alpha / \psi(\sigma)$ shows the admissible number of subgraphs at $\sigma$

    ○ $\psi(\sigma) = \binom{n}{\sigma} / \binom{n+n'}{\sigma}$ (Minimum $P$ value at $\sigma$)
    ○ For $k_{rt} = \alpha / \psi(\sigma_{rt})$, if $\psi$ is monotonically decreasing, $m(k_{rt}) = |\{ H \in \mathcal{H} \mid \psi(f(H)) \leq \psi(\sigma_{rt}) \}| = |\{ H \in \mathcal{H} \mid f(H) \geq \sigma_{rt} \}|$

# How to Use Subgraph Mining

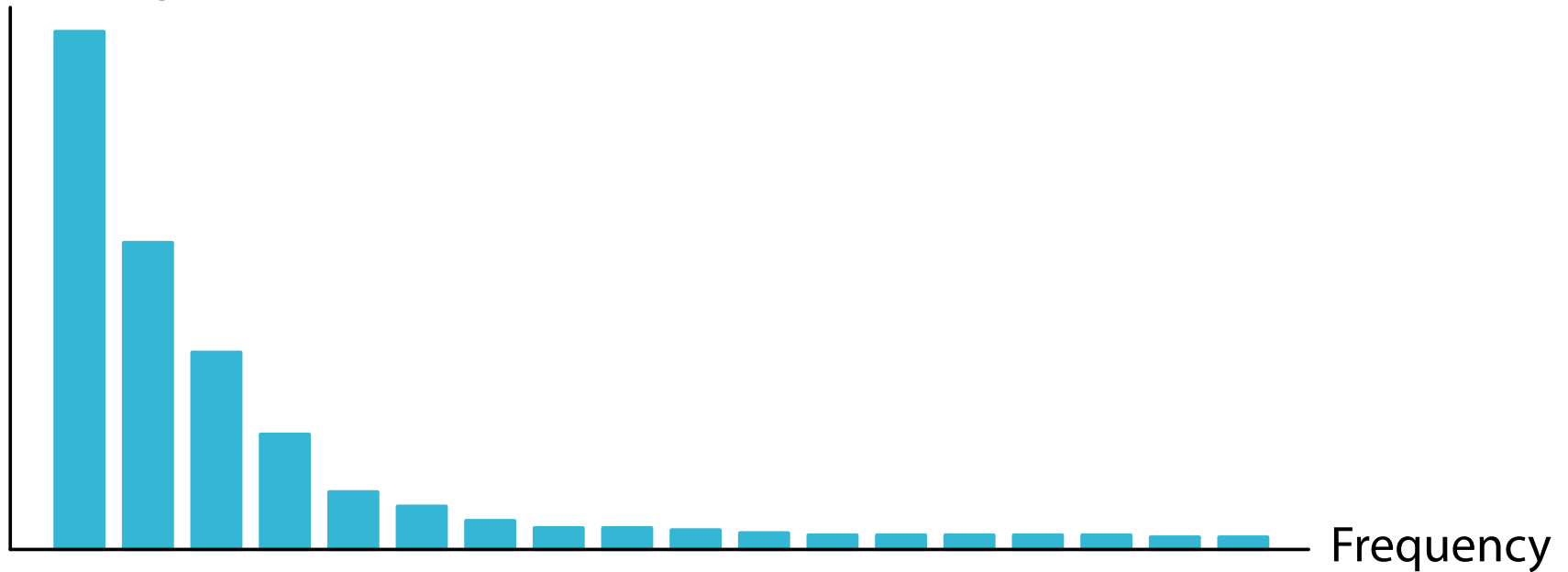# Brute-Force Search (Bonferroni)



# of subgraphs

Frequency

Freq. threshold is 1

Brute-force
(Bonferroni method)

# Decremental Search (LAMP)



# of subgraphs

Frequency

Decremental search

Terminate if # of subgraphs is larger than $\alpha / \psi(\sigma)$

# Incremental Search

# of subgraphs

Frequency

Terminate if # of subgraphs detected so far exceeds $\alpha / \psi(\sigma)$

Incremental search

# Datasets

| Dataset | Size | #positive | avg.$|V|$ | avg.$|E|$ | max$|V|$ | max$|E|$ |
|---|---|---|---|---|---|---|
| PTC (MR) | 584 | 181 | 31.96 | 32.71 | 181 | 181 |
| MUTAG | 188 | 125 | 17.93 | 39.59 | 28 | 66 |
| D&D | 1178 | 691 | 284.32 | 715.66 | 5748 | 14267 |
| NCI1 | 4208 | 2104 | 60.12 | 62.72 | 462 | 468 |
| NCI167 | 80581 | 9615 | 39.70 | 41.05 | 482 | 478 |
| NCI220 | 900 | 290 | 46.87 | 48.52 | 239 | 255 |

# Correction Factor

# Number of Significant Subgraphs

# Running Time (second)

# Running Time Summary

- RMSD (root mean square deviation) of running time (seconds) to the best (fastest) running time on all datasets

| Brute-force | Decremental (LAMP) | Incremental |
|:---:|:---:|:---:|
| $6.994 \times 10^4$ | $2.410 \times 10^4$ | $1.230 \times 10^2$ |

- Incremental search is the fastest
  - More than two orders of magnitude faster than brute-force
  - Much faster than decremental (LAMP) as the final minimum frequency is usually small (~20)

# Final Minimum Frequency

| Dataset | Maximum size of subgraph nodes | | | | | | | $n$ |
|---|---|---|---|---|---|---|---|---|
| | 5 | 7 | 9 | 11 | 13 | 15 | Limitless | |
| PTC(MR) | 9 | 10 | 11 | 11 | 11 | 11 | 11 | 181 |
| MUTAG | 8 | 10 | 11 | 12 | 14 | — | — | 125 |
| D&D | 20 | 22 | 22 | 22 | 22 | 22 | 22 | 691 |
| NCI1 | 17 | 20 | 22 | 25 | 27 | 29 | — | 2104 |
| NCI167 | 7 | 8 | 9 | 10 | 11 | — | — | 9615 |
| NCI220 | 10 | 11 | 13 | 14 | 15 | 16 | 18 | 290 |

# Detected Significant Subgraphs

PTC (MR)
(carcinogenicity)

NCI 220
(anti-cancer activity)

# FWER Is still Too Low!

# Related work: LAMP version 2

- Minato et al. proposed a faster version of LAMP in itemset mining
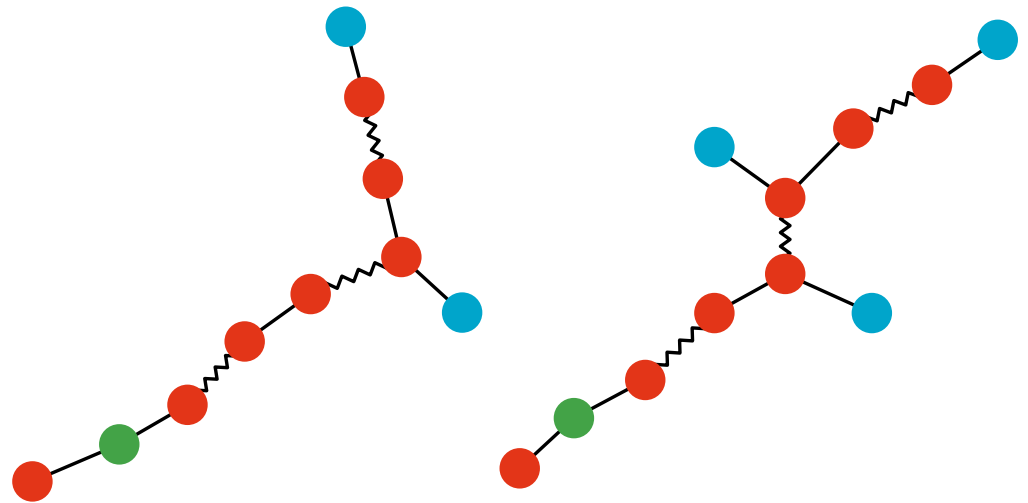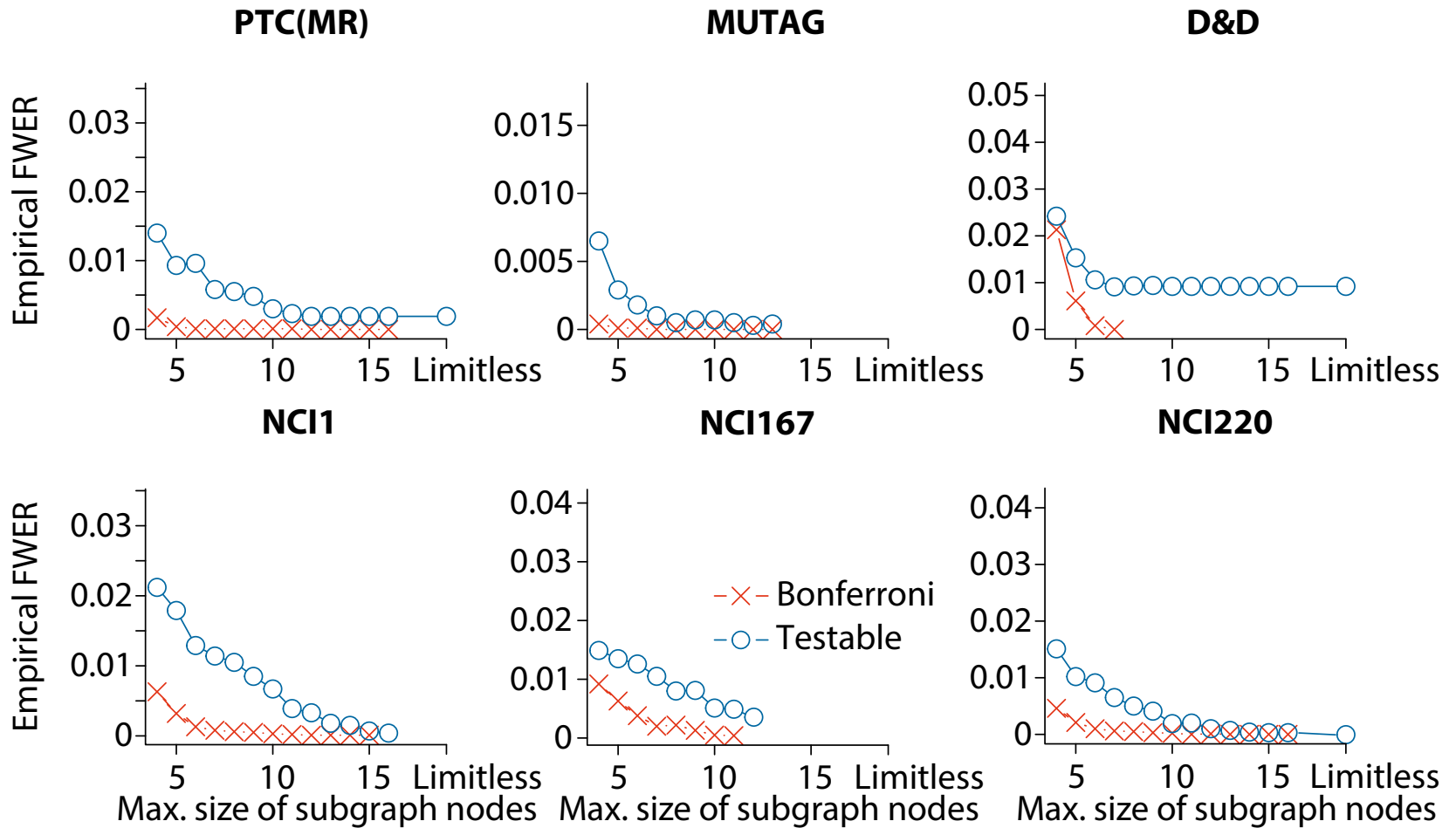  - Minato, S., Uno, T., Tsuda, K., Terada, A. and Sese, J.:
    **Fast Statistical Assessment for Combinatorial Hypotheses Based on Frequent Itemset Mining**
    ECML PKDD 2014

- The idea is almost the same with our incremental search
  - Start from $\sigma = 1$, every time an item is added, the condition $|\mathcal{I}(\sigma)| \leq \alpha/\psi(\sigma)$ is checked
    - $\mathcal{I}(\sigma)$: the set of itemsets found so far with the frequency $\geq \sigma$
  - As soon as $|\mathcal{I}(\sigma)| > \alpha/\psi(\sigma)$, the current $\sigma$ is too large and we decrement it

# Conclusion

- Significant subgraphs mining with multiple testing correction is achieved
    - The first work that considers multiple testing correction in graph mining

- Efficient and effective (less false negatives) using testability

- Future work
    - Increase the FWER with keeping $\leq \alpha$
        - Currently we ignore correlations between subgraphs

# Papers about Testability

- Tarone, R.E.:
  **A modified Bonferroni method for discrete data**
  Biometrics (1990)

- Terada, A., Okada-Hatakeyama, M., Tsuda, K., Sese, J.:
  **Statistical significance of combinatorial regulations**,
  *Proc. Natl. Acad. Sci. USA* (2013).

- Minato, S., Uno, T., Tsuda, K., Terada, A., Sese, J.:
  **Fast Statistical Assessment for Combinatorial Hypotheses
  Based on Frequent Itemset Mining**
  ECML PKDD 2014

- Sugiyama, M., Llinares López, F., Kasenburg, N., Borgwardt, K.M.:
  **Significant Subgraph Mining with Multiple Testing Correction**,
  SIAM SDM 2015 (`http://arxiv.org/abs/1407.0316`)
  - Code: `http://git.io/N126`