

# Exploiting discrete test statistics for significant pattern mining

## Theory and applications

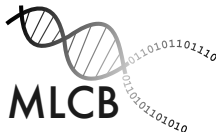
Felipe Llinares López

D-BSSE, ETH Zürich

February 16th, 2015



**D-BSSE**  
Department of Biosystems  
Science and Engineering



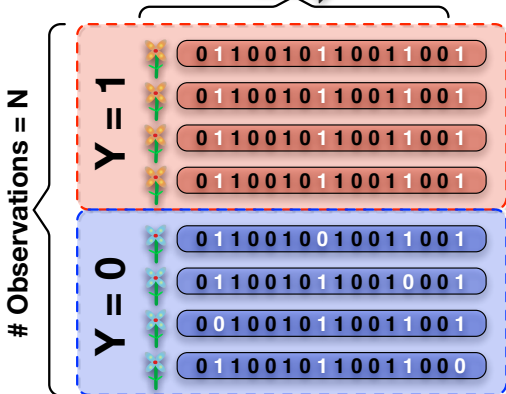
# Outline

- 1 Introduction
- 2 Statistical association testing in pattern mining
  - Background
  - The minimum attainable  $p$ -value
  - Tarone's method and LAMP
- 3 Westfall-Young Light
  - Preliminaries
  - State-of-the-art
  - Our proposal
  - Results
- 4 Conclusions

# Application: significant itemset mining

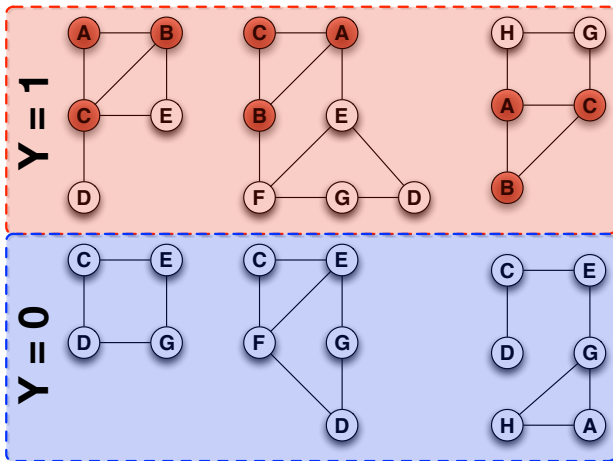
- *Example:* Find high-order (multiplicative) combinations of binary predictors associated with class membership

# Features = P  $\longrightarrow$  # Patterns =  $2^P$



# Application: significant subgraph mining

- *Example:* Find molecular motifs statistically associated with drug activity



# Problem statement

## OUR GOAL

Find **all** patterns whose occurrence within an object is **statistically associated** with class membership, **after correction for multiple testing**

## Statistical association testing

### Goal

Given  $\{(x_i, y_i)\}_{i=1}^N$  sampled iid from  $p_{\mathbf{X}\mathbf{Y}}(x, y)$ , determine if  $\mathbf{X} \not\perp \mathbf{Y}$ ; i.e. are dependent RVs.

- 1 Assume  $\mathbf{X} \perp \mathbf{Y}$  unless *proven* otherwise
- 2 Choose a *test statistic*  $\mathbf{T}$  to measure the strength of the association between  $\mathbf{X}$  and  $\mathbf{Y}$  exhibited by the sample  $\{(x_i, y_i)\}_{i=1}^N$
- 3 Derive the distribution of  $\mathbf{T}$  under the assumption  $\mathbf{X} \perp \mathbf{Y}$
- 4 Given  $t = T(\{(x_i, y_i)\}_{i=1}^N)$ , compute the *p-value* as  $p = \Pr(\mathbf{T} \geq t | \mathbf{X} \perp \mathbf{Y})$
- 5 *Reject* the assumption  $\mathbf{X} \perp \mathbf{Y}$  if  $p \leq \alpha$

## Statistical association testing in pattern mining

- In pattern mining, every  $p$ -value can be obtained from a  $2 \times 2$  contingency table:

Variables	$\mathbf{X} = 1$	$\mathbf{X} = 0$	Row totals
$\mathbf{Y} = 1$	$a$	$b$	$n$
$\mathbf{Y} = 0$	$c$	$d$	$N - n$
Col totals	$x$	$N - x$	$N$

- Common test statistics  $\mathbf{T}$  for this case are:
  - Fisher's exact test [Fisher, 1922]
  - $\chi^2$  test [Pearson, 1900]

## The multiple hypothesis testing problem

- For a dataset with  $D$  patterns  $\Rightarrow$  Test  $\mathbf{X}_i \perp \mathbf{Y} \forall i = 1, \dots, D$ 
  - We have  $D$  different contingency tables
  - Margins  $n$  and  $N$  the same  $\forall i = 1, \dots, D$
  - Margin  $x$  depends on  $i \Rightarrow x_i$
- When testing the hypothesis  $\mathbf{X} \perp \mathbf{Y}$  via  $p \leq \alpha$ , the probability of a false discovery occurring is  $\alpha$
- What if we test  $D$  hypotheses  $\mathbf{X}_i \perp \mathbf{Y} \ i = 1, \dots, D$ ?
  - *Remark:*  $\mathbb{E}[\text{FP}] = \alpha D$ 
    - $(\alpha = 0.05, D = 2 \cdot 10^5) \Rightarrow 10^4$  false positives on average
    - $(\alpha = 0.05, D = 2 \cdot 10^{10}) \Rightarrow 10^8$  false positives on average
    - $(\alpha = 0.05, D = 2^{2 \cdot 10^5}) \Rightarrow 10^{60205}$  false positives on average



## The multiple hypothesis testing problem

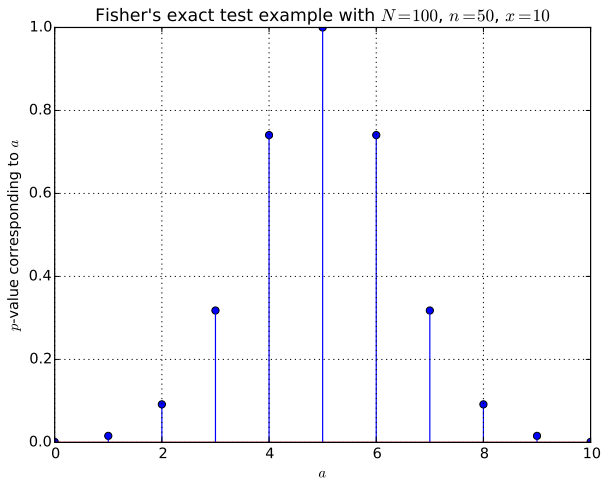
### Family-Wise Error Rate (FWER)

The FWER is defined as the probability of producing one or more false discoveries. If one can *guarantee* that  $\text{FWER} \leq \alpha$ , then the multiple hypothesis testing procedure is said to *control the FWER at level  $\alpha$*

- *Solution*: Reject each hypothesis  $\mathbf{X}_i \perp \mathbf{Y}$   $i = 1, \dots, D$  iff  $p_i \leq \delta$ , where  $\delta$  is chosen to ensure  $\text{FWER} \leq \alpha$
- *Remark*:  $\text{FWER} \ll \alpha$  is **not** beneficial
- *Bonferroni correction [Bonferroni, 1936]*: Let  $\delta = \frac{\alpha}{D}$
- **What if, as in pattern mining,  $D$  is a gigantic number?**

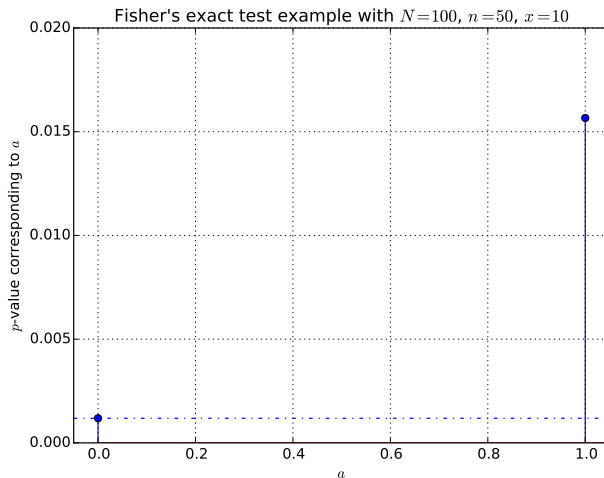
## The minimum attainable $p$ -value

In pattern mining, test statistics and attainable  $p$ -values are discrete...



## The minimum attainable $p$ -value

Thus, a minimum attainable  $p$ -value exists...



## The concept of testability

- Given  $x$ ,  $n$  and  $N$ , let  $\Psi(x) = \min_a p(a, x, n, N)$  be the minimum  $p$ -value attainable by the discrete test
  - **Remark:** Well-defined since  $(x, n, N)$  are assumed fixed!
- For each pattern  $i = 1, \dots, D$  the minimum attainable  $p$ -value  $\Psi(x_i)$  is a function of the pattern support  $x_i$
- If  $\Psi(x_i) > \delta$ , the  $i$ -th pattern can never be significant at corrected level  $\delta$
- Define  $\mathcal{I}_T(\delta) = \{ i \in \{ 1, \dots, D \} \mid \Psi(x_i) \leq \delta \}$ , the set of *testable hypotheses at level  $\delta$*
- **What are the implications?**

## An improved Bonferroni Correction for discrete data

- This phenomenon was first exploited in [Tarone, 1990]
- Tarone showed that  $\text{FWER} \leq \delta |\mathcal{I}_{\mathcal{T}}(\delta)|$
- To ensure  $\text{FWER} \leq \alpha$  choose  $\delta_{\text{tar}}^* = \max \{ \delta \mid \delta |\mathcal{I}_{\mathcal{T}}(\delta)| \leq \alpha \}$
- Usually,  $\delta_{\text{tar}}^* \gg \frac{\alpha}{D}$ , leading to greatly increased statistical power
- Computing  $\delta_{\text{tar}}^*$  as proposed by Tarone is unfeasible computationally
- In [Terada et al., 2013a], Terada et. al. link Tarone's method to frequent itemset mining, proposing the Limitless-Arity Multiple Testing Procedure (LAMP)

# Westfall-Young Light

Felipe Llinares, Mahito Sugiyama, Laetitia  
Papaxanthos, Karsten Borgwardt

## Addressing the dependence between test statistics

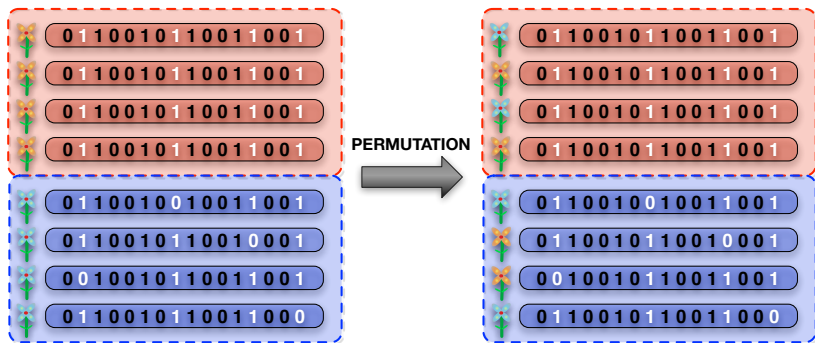
- Tarone's method ignores the dependence structure between patterns:  $\text{FWER}(\delta_{\text{tar}}^*) \ll \alpha$  frequently
- The optimal FWER-controlling method would use  $\delta^*$  such that:

$$\delta^* = \underset{\delta}{\operatorname{argmax}} \delta \text{ s.t. } \text{FWER}(\delta) \leq \alpha$$

- Evaluating  $\text{FWER}(\delta)$  in closed-form is not possible
- *Solution:* Use resampling methods, like Westfall-Young (WY) permutation testing [Westfall and Young, 1993]

## The Westfall-Young permutation testing procedure

- **Step 1:** Randomly permute class labels  $\{y_i\}_i^N$  to obtain  $\{\tilde{y}_i\}_i^N$



- By construction,  $\mathbf{X}_i \perp \tilde{\mathbf{Y}} \forall i = 1, \dots, D$  (i.e. **no** pattern is associated with the class labels)



## The Westfall-Young permutation testing procedure

- **Step 2:** Compute the  $p$ -values  $\tilde{p}_i$  for each pattern  $i = 1, \dots, D$  using the permuted labels  $\tilde{\mathbf{Y}}$ 
  - Since  $\mathbf{X}_i \perp \tilde{\mathbf{Y}} \forall i = 1, \dots, D$ , any pattern for which  $\tilde{p}_i \leq \delta$  would be a false positive at level  $\delta$
- **Step 3:** Compute  $p_{\min} = \min_{i=1, \dots, D} \tilde{p}_i$ 
  - $\text{FP} > 0 \Leftrightarrow p_{\min} \leq \delta$
- **Step 4:** Repeat steps 1-3  $J$  times, obtaining  $\left\{ p_{\min}^{(j)} \right\}_{j=1}^J$ 
  - $\text{FWER} = \Pr(\text{FP} > 0) \approx \frac{1}{J} \sum_{j=1}^J \mathbb{1} \left[ p_{\min}^{(j)} \leq \delta \right]$
- **Step 5:**  $\delta^*$  can be found as the  $\alpha$ -quantile of  $\left\{ p_{\min}^{(j)} \right\}_{j=1}^J$

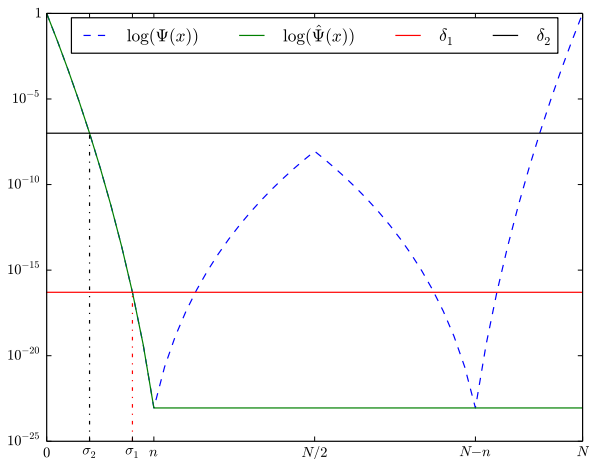
## FastWY

- Computing  $p_{\min}$  naively requires enumerating and computing  $J$   $p$ -values for all  $D$  patterns
- Terada et. al. propose in [Terada et al., 2013b] the FastWY algorithm as an extension of LAMP to WY permutation testing
- FastWY provides a way to speedup the computation of  $p_{\min}$  over the naive approach

### KEY CONCEPT

The target is computing  $p_{\min} = \min_{i=1, \dots, D} \tilde{p}_i$ . If  $p'_{\min} = \min_{i \in \mathcal{I}(\delta)} \tilde{p}_i$  satisfies  $p'_{\min} \leq \delta$ , then  $p'_{\min} = p_{\min}$  and the search can be stopped early.

# Linking testability and frequent pattern mining



# FastWY algorithm

---

**Algorithm 1** FASTWY as proposed in [Terada et al., 2013b]

---

```

1: function FASTWY
2:   for  $j = 1, \dots, J$  do
3:      $\mathbf{y}^{(j)} \leftarrow \text{permute}(\mathbf{y})$ 
4:      $\sigma \leftarrow n + 1$ 
5:     repeat
6:        $\sigma \leftarrow \sigma - 1, \delta_\sigma \leftarrow \Psi(\sigma)$ 
7:        $\hat{\mathcal{I}}_T(\sigma) \leftarrow \text{FPM}(\{1, \dots, P\}, \sigma)$ 
8:       Compute  $p_i \forall i \in \hat{\mathcal{I}}_T(\sigma)$ 
9:        $p_{\min}^{(j)} \leftarrow \min_{i \in \hat{\mathcal{I}}_T(\sigma)} p_i$ 
10:    until  $p_{\min}^{(j)} \leq \delta_\sigma$ 
11:   end for
12:    $\delta^* \leftarrow \alpha\text{-quantile of } \left\{ p_{\min}^{(j)} \right\}_{j=1}^J$ 
13: end function

```

---

## Limitations of FastWY

- 1 Like LAMP, relies on using a monotonically decreasing surrogate  $\hat{\Psi}(x) \leq \Psi(x)$
- 2 Uses a decremental (in support threshold) search strategy
- 3 Needs to either repeat pattern mining  $J \approx 10^4$  times or store the occurrence list of every frequent pattern
- 4 Requires computing the whole set  $\left\{ p_{\min}^{(j)} \right\}_{j=1}^J$  exactly
- 5 As a consequence of (3), with overwhelming probability, some  $p_{\min}^{(j)}$  will require mining patterns with very low supports

Westfall-Young light removes all these limitations!



## Removing limitations (1)-(4): the WY-light algorithm

---

### Algorithm 2 WY-light core

---

```

1: function PROCESSPATTERN
2:   if  $x_i \in \Sigma_k$  then
3:     Compute  $\tilde{p}_i^{(j)} \forall j = 1, \dots, J$ 
4:      $p_{\min}^{(j)} \leftarrow \min(p_{\min}^{(j)}, \tilde{p}_i^{(j)}) \forall j = 1, \dots, J$ 
5:      $\text{FWER} \leftarrow \frac{1}{J} \sum_{j=1}^J \mathbb{1} \left[ p_{\min}^{(j)} \leq \delta \right]$ 
6:     while  $\text{FWER} > \alpha$  do
7:        $k \leftarrow k + 1$ 
8:       Update  $\delta_k, \Sigma_k$  and  $\sigma_k = \min \{x | x \in \Sigma_k\}$ 
9:     end while
10:  end if
11:  Enumerate all patterns  $j \in \text{Children}(i) | x_j \geq \sigma_k$ 
12: end function
    
```

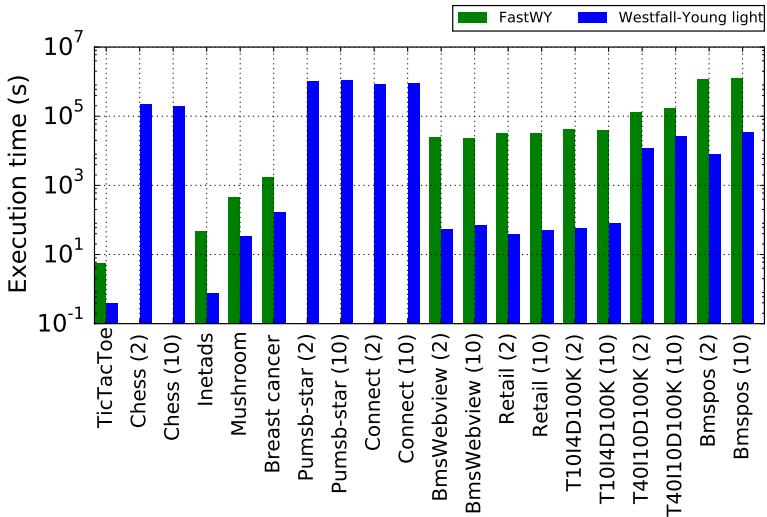
---

## Removing limitations (1)-(4): the WY-light algorithm

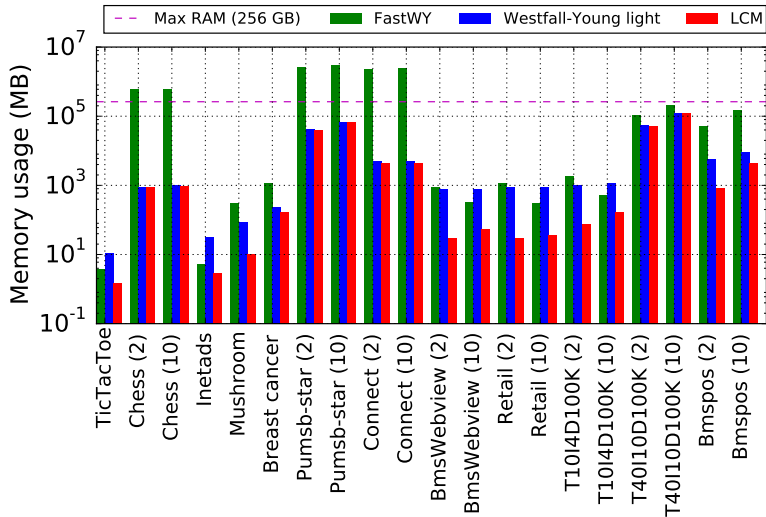
- Property 1** Whenever a new pattern  $i$  is processed, the updated empirical FWER estimate can never decrease
- Property 2**  $\text{FWER}(\delta)$  for all  $\delta \in [0, \delta_k]$  can be evaluated exactly using only the  $p$ -values of patterns in  $\mathcal{I}_T(\Sigma_k)$ .
- Property 3** For fixed  $x_i$ ,  $n$ , and  $N$ , the computational complexity of evaluating Fisher's exact test  $p$ -value  $p_i(\gamma)$  for a single value of  $\gamma$  or for all possible values of  $\gamma$  in  $[a_{i,\min}, a_{i,\max}]$  is the same and equal to  $O(\min\{x_i, n\})$



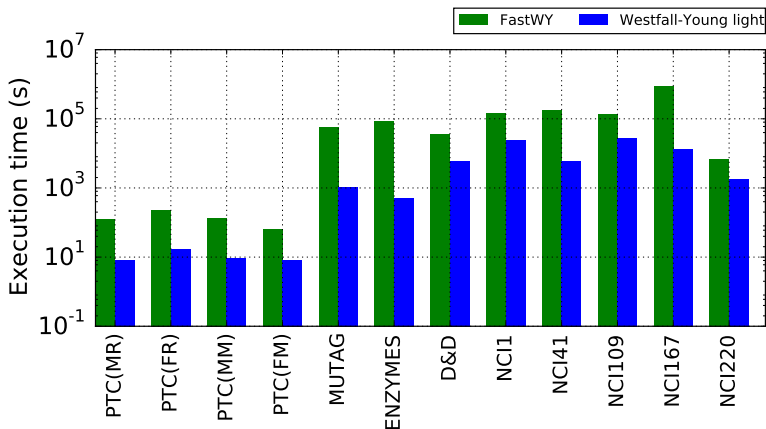
# Runtime in itemset mining



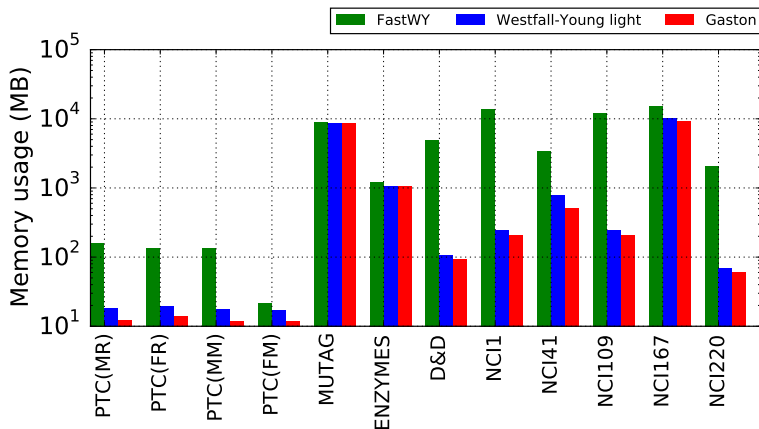
# Memory usage in itemset mining



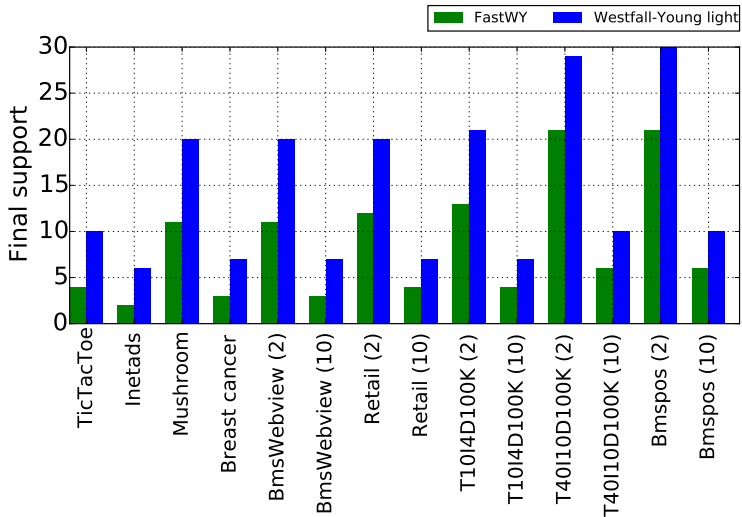
## Runtime in subgraph mining



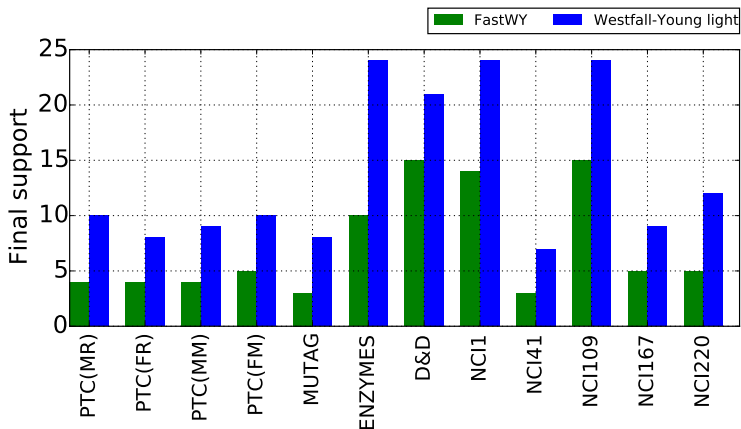
# Memory usage in subgraph mining



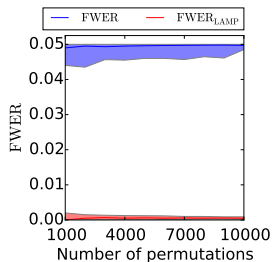
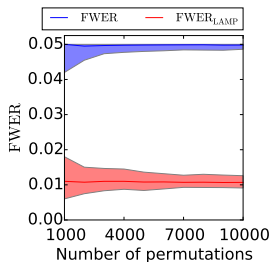
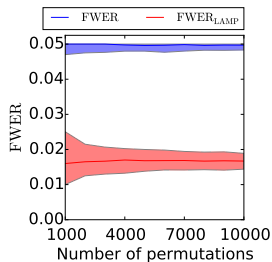
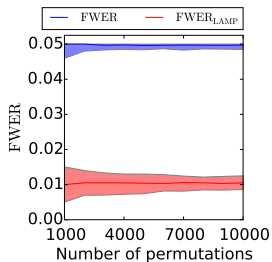
# Final support in itemset mining



## Final support in subgraph mining



# Power comparison: LAMP vs WY



## Conclusions

- Discovering patterns significantly associated with class membership is a fundamental problem in data mining
- Rigorous correction for multiple hypothesis testing is mandatory if statistically reliable results are needed
- The discrete nature of test statistics in pattern mining can be exploited to get great gains in statistical power
- Westfall-Young light allows applying the Westfall-Young permutation testing procedure to large-scale datasets
  - Scalable pattern mining under optimal FWER-control!



Thank you!

## References I



Bonferroni, C. E. (1936).

Teoria statistica delle classi e calcolo delle probabilità.

*Pubblicazioni del R Istituto Superiore di Scienze Economiche e Commerciali di Firenze*, 8:3–62.



Fisher, R. A. (1922).

On the Interpretation of  $\chi^2$  from Contingency Tables, and the Calculation of P.

*Journal of the Royal Statistical Society*, 85(1):87–94.

## References II



Pearson, K. (1900).

X. On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling.

*Philosophical Magazine Series 6*, 50:157–175.






Tarone, R. E. (1990).

A modified bonferroni method for discrete data.

*Biometrics*, 46(2):515–522.

## References III

-  Terada, A., Okada-Hatakeyama, M., Tsuda, K., and Sese, J. (2013a).  
Statistical significance of combinatorial regulations.  
*Proceedings of the National Academy of Sciences*,  
110(32):12996–13001.
-  Terada, A., Tsuda, K., and Sese, J. (2013b).  
Fast westfall-young permutation procedure for combinatorial  
regulation discovery.  
In *IEEE International Conference on Bioinformatics and  
Biomedicine*, pages 153–158.
-  Westfall, P. H. and Young, S. S. (1993).  
Resampling-based multiple testing.  
*Statistics in Medicine*, 13(10):1084–1086.

## References IV