



# Clipper: p-value-free FDR control on high-throughput data from two conditions

---

Jingyi Jessica Li

Joint work with Xinzhou Ge and Yiling Elaine Chen (Ph.D. students)

Department of Statistics  
University of California, Los Angeles

<http://jsb.ucla.edu>

# Introduction

---

# Background

- ▶ High-throughput biological data
  - Small sample size (number of replicates, often  $\leq 3$ )
  - Huge number of features (often  $\sim 10^4$ )
- ▶ Two conditions
  - Experimental
  - Background / negative control
- ▶ Identification of “interesting” features

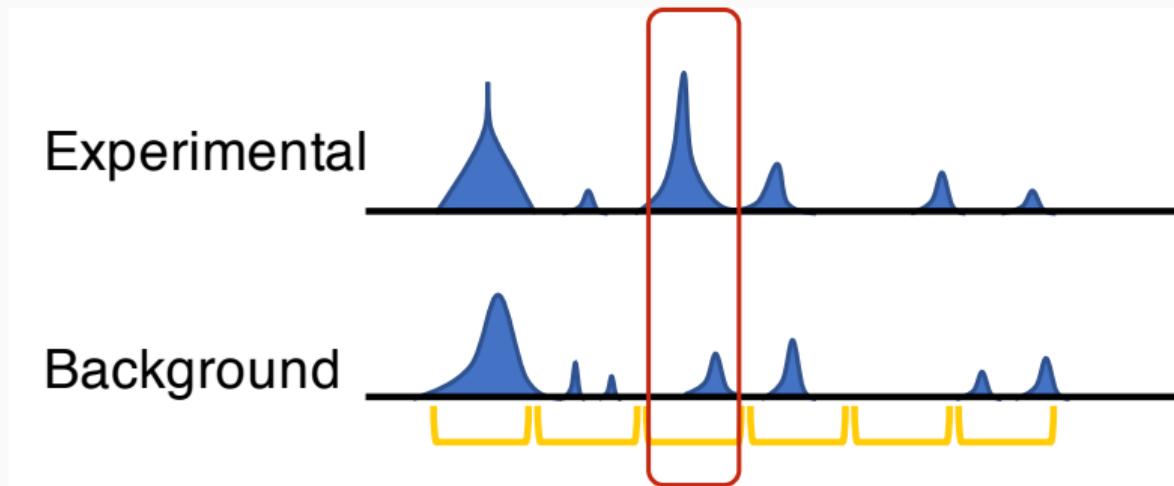


## Four example high-throughput biological applications

- ▶ Peak calling from ChIP-seq data
  - Protein-DNA binding sites
- ▶ Peptide identification from mass spectrometry (MS) data
  - Peptide-spectrum matches
- ▶ Differential analysis of RNA-seq data
  - Differentially expressed genes
- ▶ Differential analysis of Hi-C data
  - Differentially interacting chromatin regions

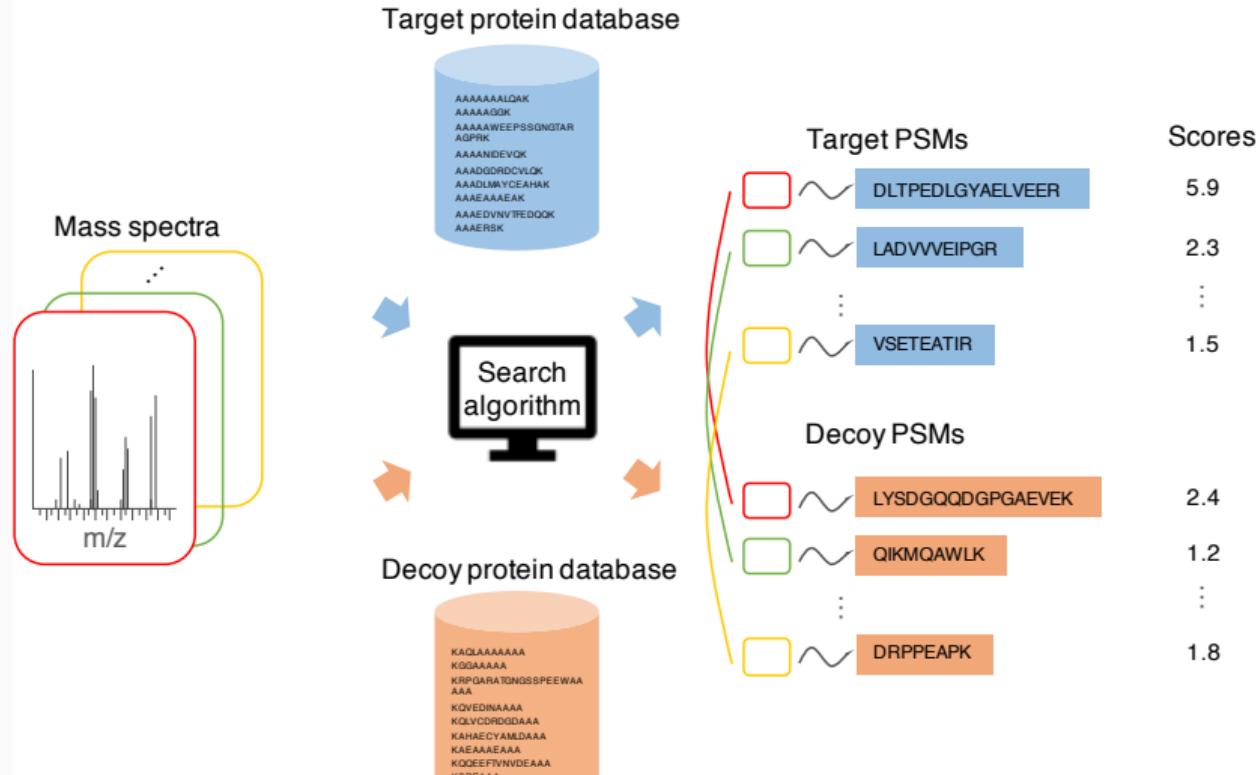


## Peak calling from ChIP-seq data



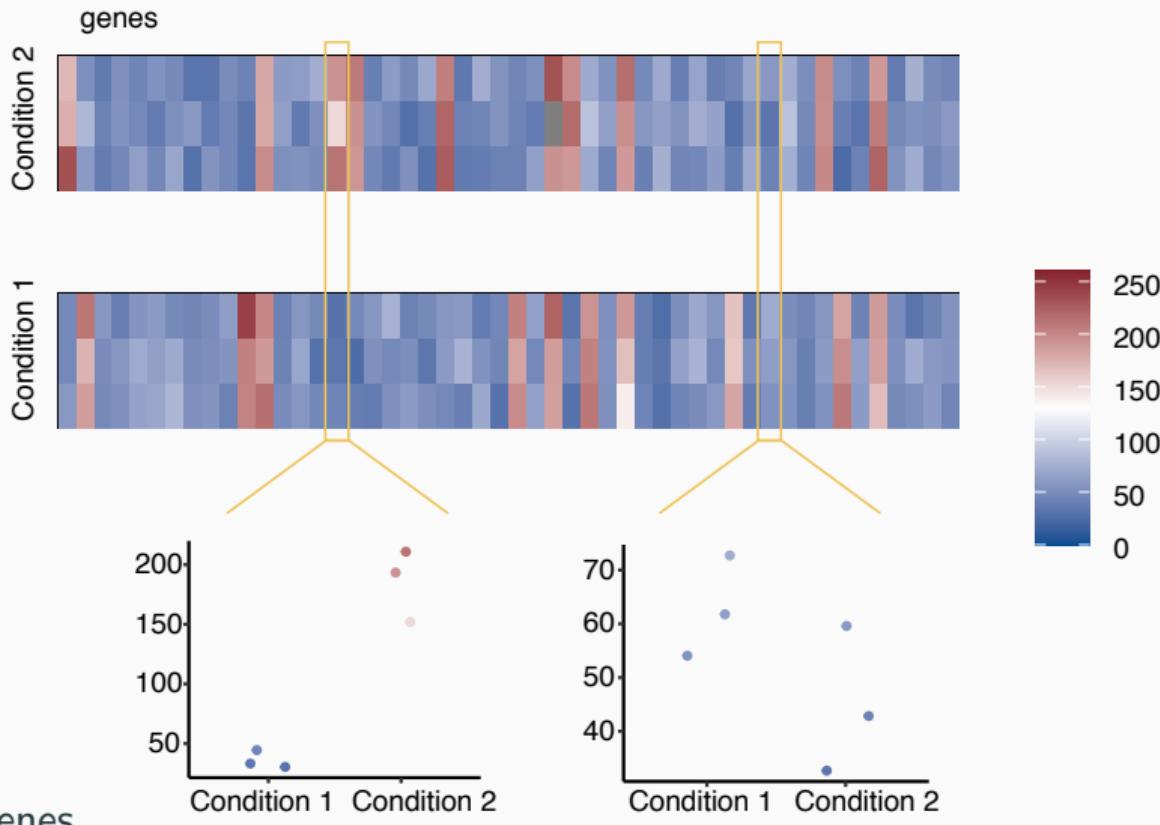
Features: genomic regions

# Peptide identification from MS data

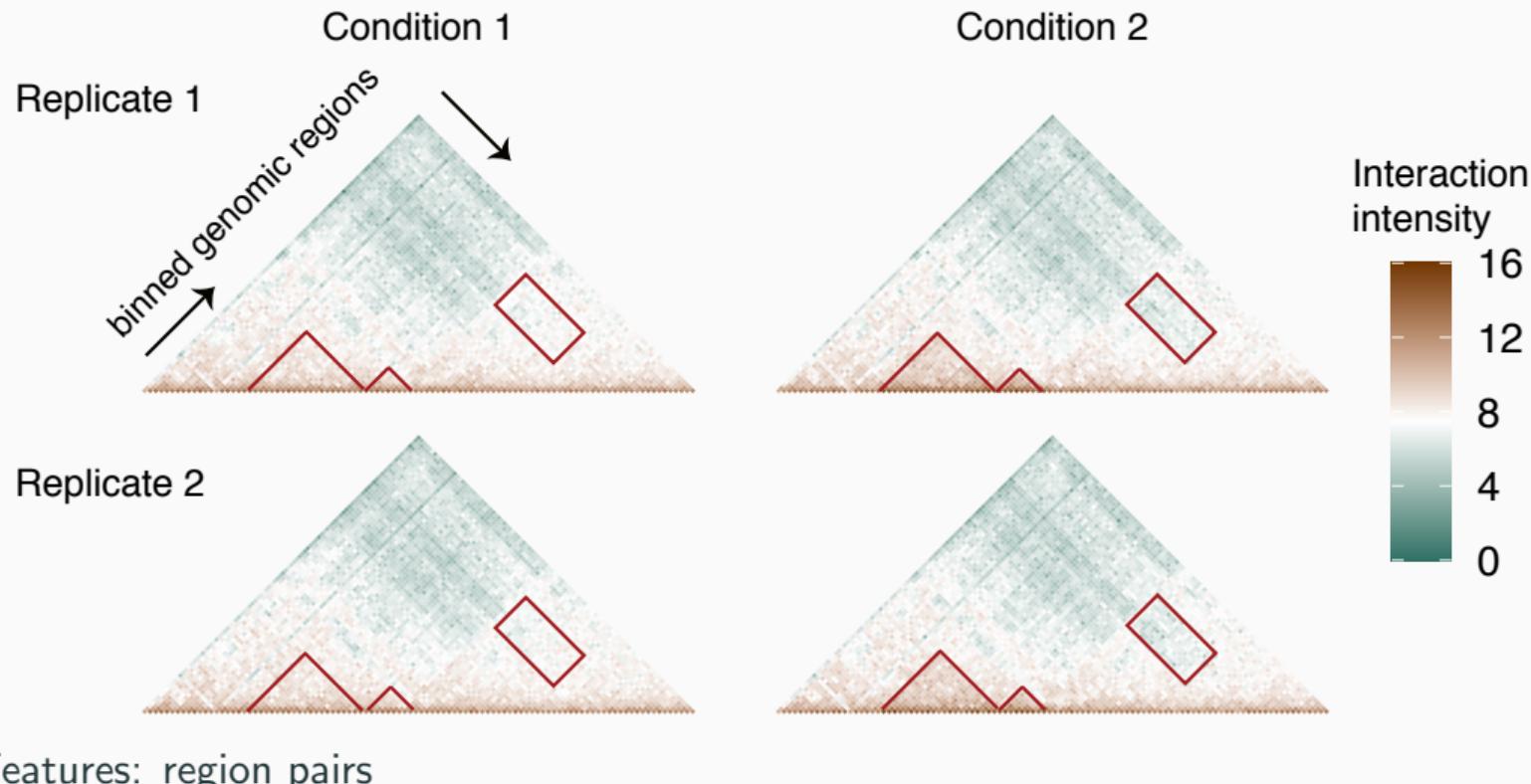


Features: Peptide-spectrum matches (PSMs)

# Identification of differentially expressed genes from RNA-seq data



# Identification of differentially interacting regions from Hi-C data



## Enrichment and differential analyses

- ▶ Interesting means “enriched” or “differential”

- ▶ Enriched features:

$$\mathbb{E}[\text{experimental}] > \mathbb{E}[\text{background}]$$

- Protein binding regions (peaks)
- Peptide-spectrum matches

- ▶ Differential features:

$$\mathbb{E}[\text{experimental}] \neq \mathbb{E}[\text{background}]$$

- Differentially expressed genes
- Differentially interacting regions

## False discovery rate (FDR)

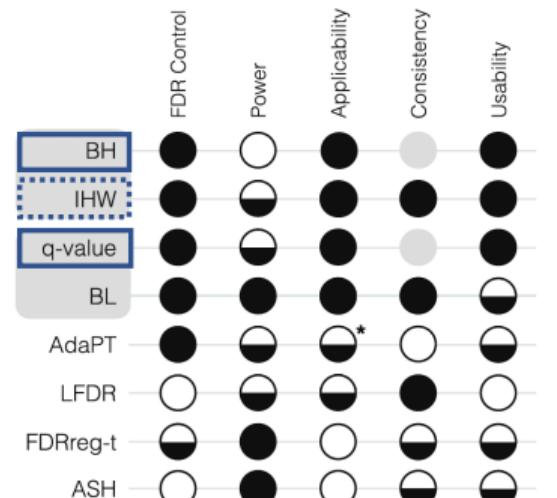
- ▶ Criterion for controlling false discoveries
- ▶ **Frequentist FDR** [Benjamini and Hochberg, *JRSSB*, 1995]

$$\text{FDR} := \mathbb{E} \left[ \frac{\# \text{ false discoveries}}{\# \text{ discoveries} \vee 1} \right]$$

- ▶ Bayesian paradigm:
  - Bayesian false discovery rate [Efron and Tibshirani, *Genet Epidemiol*, 2002]
  - Local false discovery rate (**fdr**) [Efron et al., *JASA*, 2001]
  - Local false sign rate [Stephens, *Biostatistics*, 2017]

# Existing FDR control methods

	Input	Assumptions	Output	R package
BH	p-values	exchangeability	adjusted p-values	<code>stats</code>
IHW	(1) p-values (2) independent & informative covariate	exchangeability within covariate groups		<code>ihw</code>
q-value	p-values	exchangeability	q-values	<code>qvalue</code>
BL			adjusted p-values	<code>swfdr</code>
AdaPT	(1) p-values (2) independent & informative covariate	exchangeability conditional on covariate(s)	q-values	<code>adaptMT</code>
LFDR		exchangeability within covariate groups	adjusted p-values	none
FDRreg	(1) z-scores (2) independent & informative covariate	exchangeability conditional on covariate(s); normal test statistics	Bayesian FDRs	<code>FDRreg</code>
ASH	(1) effect sizes (2) standard errors of (1)	effects are unimodal; test statistics have normal or t mixture components	q-values	<code>ash</code>



“A practical guide to methods controlling false discoveries in computational biology”  
 [Korthauer et al., *Genome Biol.*, 2019]

## Generic FDR control methods

---

- ▶ **P-value based methods** (exact)

- Benjamini-Hochberg (**BH**) procedure [Benjamini and Hochberg, *JRSSB*, 1995]
- Storey's **q-value** procedure [Storey, *JRSSB*, 2002]

- ▶ **Local fdr based method** (approximate)

- Thresholding local fdr to  $q$  (e.g., 5%) approximately controls the FDR

# Calculation of p-values

---

## ► Requirements

- Distributional assumptions (parametric)
- Large number of replicates (nonparametric)

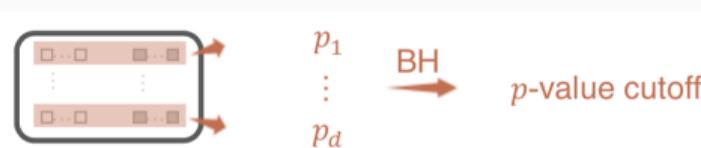
## ► Approaches

- “**Paired**” approach
- “**Pooled**” approach



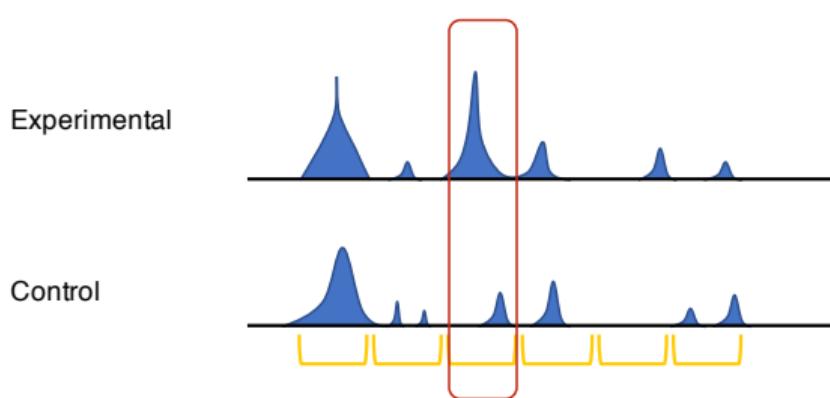
## Paired approach

- ▶ Used in
  - Peak calling from ChIP-seq data
  - Identification of differentially expressed genes from RNA-seq data
  - Identification of differentially interacting regions from Hi-C data
- ▶ One feature at a time, two-sample test



- ▶ Issues
  - Mis-formulation (e.g., two-sample test as one-sample test)
  - Mis-specification (e.g., negative binomial as Poisson)

## Example of paired mis-formulation: MACS [Zhang et al., Genome Biol, 2008]



In the control samples, we often observe tag distributions with local fluctuations and biases. For example, at the FoxA1 candidate peak locations, tag counts are well correlated between ChIP and control samples (Figure 1c,d). Many possible sources for these biases include local chromatin structure, DNA amplification and sequencing bias, and genome copy number variation. Therefore, instead of using a uniform  $\lambda_{BG}$  estimated from the whole genome, MACS uses a dynamic parameter,  $\lambda_{local}$ , defined for each candidate peak as:

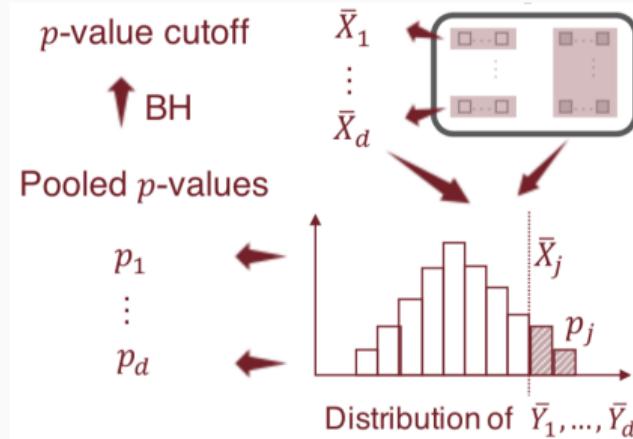
$$\lambda_{local} = \max(\lambda_{BG}, [\lambda_{1k}, \lambda_{5k}, \lambda_{10k}])$$

where  $\lambda_{1k}$ ,  $\lambda_{5k}$  and  $\lambda_{10k}$  are  $\lambda$  estimated from the 1 kb, 5 kb or 10 kb window centered at the peak location in the control sample, or the ChIP-Seq sample when a control sample is not available (in which case  $\lambda_{1k}$  is not used).  $\lambda_{local}$  captures the influence of local biases, and is robust against occasional low tag counts at small local regions. MACS uses  $\lambda_{local}$  to calculate the *p*-value of each candidate peak and removes potential false positives due to local biases (that is, peaks significantly under  $\lambda_{BG}$ , but not under  $\lambda_{local}$ ). Candidate peaks with *p*-values below a user-defined threshold *p*-value (default  $10^{-5}$ ) are called, and the ratio between the ChIP-Seq tag count and  $\lambda_{local}$  is reported as the *fold\_enrichment*.

Cited for more than 8,000 times

## Pooled approach

- ▶ Used in
  - Peptide identification from MS data
- ▶ Pools all features' background measurements to form a null distribution
  - Assumes a **homogeneous** background: features are i.i.d. under background



## Generic FDR control methods

---

- ▶ **P-value based methods** (exact)
  - Benjamini-Hochberg (**BH**) procedure [Benjamini and Hochberg, *JRSSB*, 1995]
  - Storey's **q-value** procedure [Storey, *JRSSB*, 2002]
- ▶ **Local fdr based method** (approximate)
  - Thresholding local fdr to  $q$  (e.g., 5%) approximately controls the FDR

## Local fdr [Efron et al., JASA, 2001]

- ▶ Local fdr of feature  $j$

$$\text{local fdr}_j := \mathbb{P}(\text{feature is uninteresting} \mid Z = z_j)$$

vs. Bayesian false discovery rate

$$\text{Fdr}(z) := \mathbb{P}(\text{feature is uninteresting} \mid Z \geq z)$$

- ▶ It can be shown that (assuming local  $\text{fdr}_j$  is monotone decreasing in  $z_j$ )

$$\text{Fdr}(z^*) \leq q \text{ if } z^* := \min\{z_j : \text{local fdr}_j \leq q\}$$

- ▶ With **discoveries** :=  $\{\text{feature } j : \text{local fdr}_j \leq q\}$

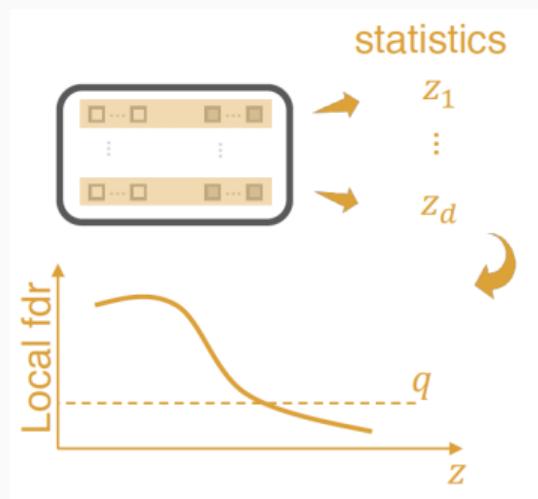
$$\text{FDR}(\text{discoveries}) \approx \text{Fdr}(z^*)$$

- ▶ So FDR is approximately controlled



## Local fdr [Efron et al., JASA, 2001]

- ▶ Requires estimating the null distribution of test statistic by
  - Normal distributional assumption
  - or
  - Swapping replicates between conditions



# Our proposal: Clipper

- ▶ Does not
  - use p-values
  - assume parametric distributions
  - require many replicates
- ▶ Two components
  - Contrast scores
  - Cutoff
    - ▶ Enrichment analysis with equal numbers of replicates: **BC procedure** [Barber and Candès, *Ann Stat*, 2015]
    - ▶ Differential analysis and other enrichment analysis: **GZ procedure** [Gimenez and Zou, *PMLR*, 2019]
- ▶ Robust to
  - distributions
  - numbers of replicates
  - outliers



# Illustration of methods

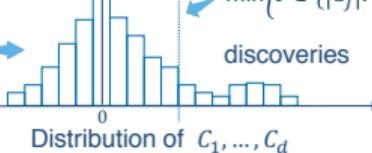
	Experiment	Background
Feature	1	
1	$X_{11}$ ... $X_{1m}$	$Y_{11}$ ... $Y_{1n}$
:	⋮	⋮
$d$	$X_{d1}$ ... $X_{dm}$	$Y_{d1}$ ... $Y_{dn}$

**a Clipper** Contrast scores

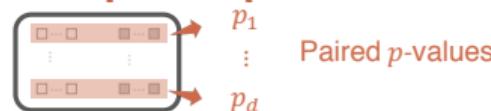


Contrast score cutoff (BC procedure):  

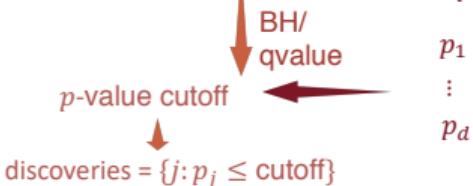
$$\min\{t \in \{|C_j| : C_j \neq 0\} : \frac{1+\#\{j : C_j \leq -t\}}{\#\{j : C_j \geq t\} v_1} \leq q\}$$



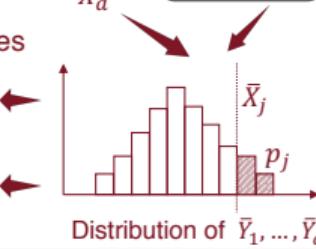
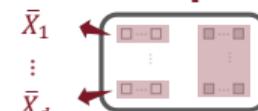
**b BH/qvalue-pair**



Pooled p-values



**c BH/qvalue-pool**



**d locfdr**



## Clipper Method

---

## Notations

- $d$ : number of features
- $\mathbf{X}_j = (X_{j1}, \dots, X_{jm})^\top$ ,  $\mathbf{Y}_j = (Y_{j1}, \dots, Y_{jn})^\top$ : measurements under two conditions
  - $X_{j1}, \dots, X_{jm} \geq 0$  are i.i.d.;  $\mu_{Xj} = \mathbb{E}[X_{j1}]$
  - $Y_{j1}, \dots, Y_{jn} \geq 0$  are i.i.d.;  $\mu_{Yj} = \mathbb{E}[Y_{j1}]$
  - $j = 1, \dots, d$

	Experiment			Background			
Feature	1	$X_{11}$	$\dots$	$X_{1m}$	$Y_{11}$	$\dots$	$Y_{1n}$
	$\vdots$	$\vdots$		$\vdots$	$\vdots$		$\vdots$
$d$		$X_{d1}$	$\dots$	$X_{dm}$	$Y_{d1}$	$\dots$	$Y_{dn}$

Feature  $j$  is **interesting**

- Enrichment analysis:  $\mu_{Xj} > \mu_{Yj}$
- Differential analysis:  $\mu_{Xj} \neq \mu_{Yj}$



# Assumptions

Conditioning on  $\{\mu_{Xj}\}_{j=1}^d$  and  $\{\mu_{Yj}\}_{j=1}^d$ ,

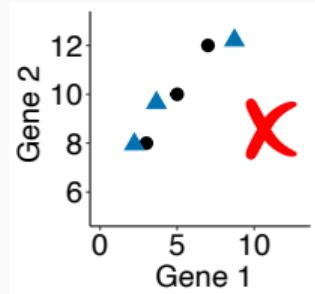
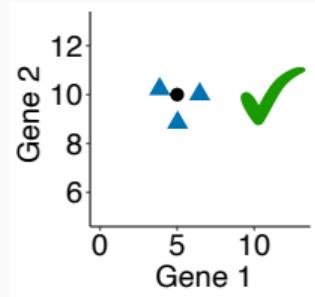
- ▶ **Independence:**

$X_1, \dots, X_m, Y_1, \dots, Y_n$  are mutually independent (1)

$\mathbf{X}_1, \dots, \mathbf{X}_d, \mathbf{Y}_1, \dots, \mathbf{Y}_d$  are mutually independent

- ▶ For **uninteresting** feature  $j$ ,

$X_1, \dots, X_m, Y_1, \dots, Y_n$  are identically distributed (2)



- Expectation
- ▲ Measurement

## Clipper for three analysis tasks

---

Calculation of contrast scores depends on analysis tasks:

- ▶ Enrichment analysis with  $m = n$
- ▶ Enrichment analysis with  $m \neq n$
- ▶ Differential analysis

## Enrichment analysis with $m = n$ : contrast scores

Two summary statistics:

$$t^{\text{diff}}(\mathbf{x}, \mathbf{y}) := \bar{x} - \bar{y} \quad (3)$$

$$t^{\max}(\mathbf{x}, \mathbf{y}) := \max(\bar{x}, \bar{y}) \cdot \text{sign}(\bar{x} - \bar{y}) \quad (4)$$

In enrichment analysis with  $m = n$ , **contrast score** of feature  $j$ :

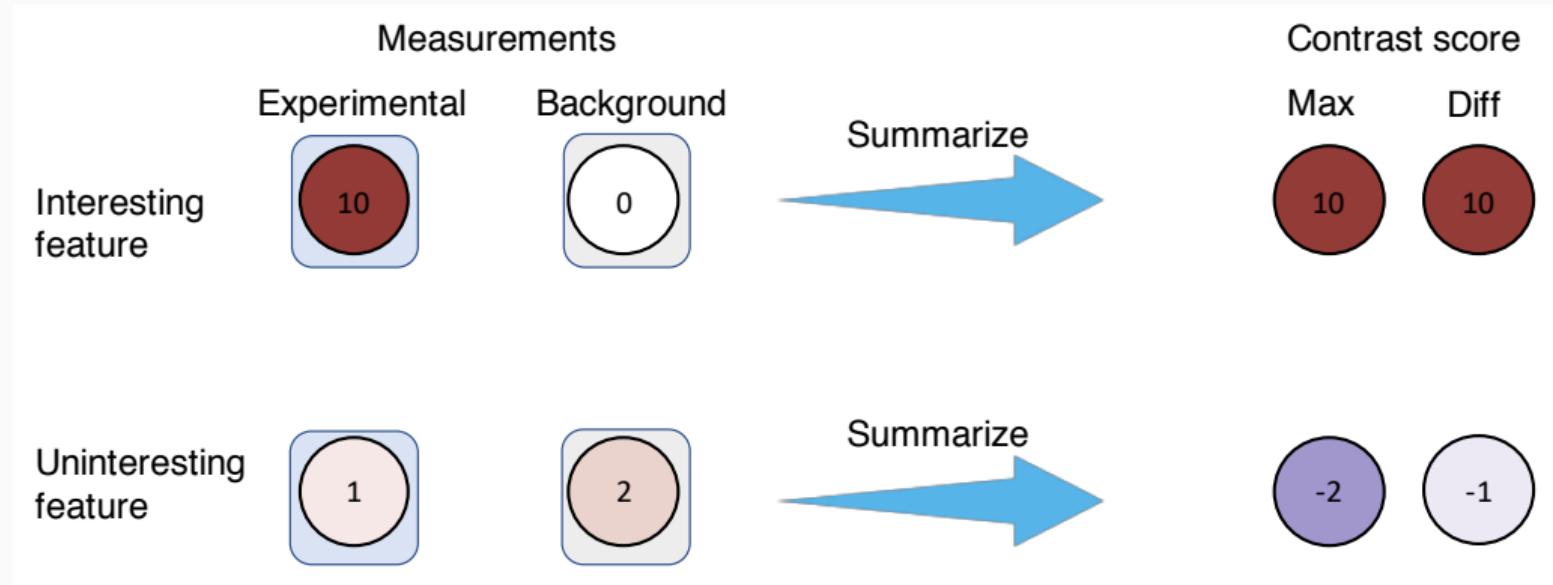
$$C_j := t^{\text{diff}}(\mathbf{X}_j, \mathbf{Y}_j) \quad \text{difference contrast score} \quad (5)$$

or

$$C_j := t^{\max}(\mathbf{X}_j, \mathbf{Y}_j) \quad \text{maximum contrast score} \quad (6)$$



## Enrichment analysis with $m = n = 1$ : contrast scores

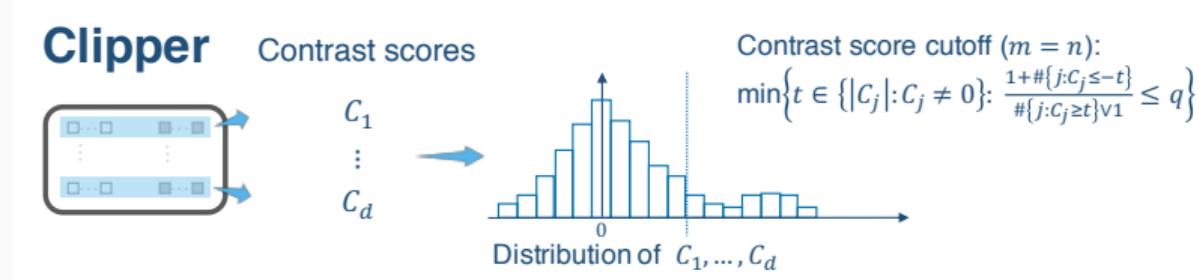


## Enrichment analysis with $m = n$ : cutoff

### Definition 1 BC procedure [Barber and Candès, Ann Stat, 2015]

- Given contrast scores  $\{C_j\}_{j=1}^d$ , define  $\mathcal{C} = \{|C_j| : C_j \neq 0 ; j = 1, \dots, d\}$
- Based on the target FDR threshold  $q \in (0, 1)$ , contrast-score cutoff  $T^{\text{BC}}$ :

$$T^{\text{BC}} := \min \left\{ t \in \mathcal{C} : \frac{\text{card}(\{j : C_j \leq -t\}) + 1}{\text{card}(\{j : C_j \geq t\}) \vee 1} \leq q \right\} \quad (7)$$



- Discoveries:**  $\{j : C_j \geq T^{\text{BC}}\}$
- BC vs. BH [Arias-Castro and Chen, *Electronic J Stat*, 2016]

## Key lemma to guarantee the theoretical FDR control by the BC procedure

Define  $S_j = \text{sign}(C_j) \in \{-1, 0, 1\}$

$\mathcal{N}$ : the set of uninteresting features

Then

1.  $S_1, \dots, S_d$  are mutually independent
2.  $\mathbb{P}(S_j = 1) = \mathbb{P}(S_j = -1)$  for all  $j \in \mathcal{N}$
3.  $\{S_j\}_{j \in \mathcal{N}} \perp \mathcal{C}$

When  $m \neq n$ , 2 and 3 are not guaranteed to hold



## Enrichment analysis with $m \neq n$ : permutation

Data:

$$\mathbf{W} := \begin{bmatrix} X_{11} & \cdots & X_{1m} & Y_{11} & \cdots & Y_{1n} \\ \vdots & & & & & \vdots \\ X_{d1} & \cdots & X_{dm} & Y_{d1} & \cdots & Y_{dn} \end{bmatrix}$$

- $\sigma$ : a permutation function to permute the columns of  $\mathbf{W}$  and output  $\mathbf{W}^\sigma$

$$\mathbf{W}^\sigma := \begin{bmatrix} X_{11}^\sigma & \cdots & X_{1m}^\sigma & Y_{11}^\sigma & \cdots & Y_{1n}^\sigma \\ \vdots & & & & & \vdots \\ X_{d1}^\sigma & \cdots & X_{dm}^\sigma & Y_{d1}^\sigma & \cdots & Y_{dn}^\sigma \end{bmatrix}$$

- The permuted measurements  $\left\{(\mathbf{X}_j^\sigma, \mathbf{Y}_j^\sigma)\right\}_{j=1}^d$

## Enrichment analysis with $m \neq n$ : permutation

- ▶ The identity permutation  $\sigma_0$
- ▶ Sample  $h$  non-identity permutations  $\sigma_1, \dots, \sigma_h$
- ▶  $\left\{(\mathbf{X}_j^{\sigma_0}, \mathbf{Y}_j^{\sigma_0}), (\mathbf{X}_j^{\sigma_1}, \mathbf{Y}_j^{\sigma_1}), \dots, (\mathbf{X}_j^{\sigma_h}, \mathbf{Y}_j^{\sigma_h})\right\}_{j=1}^d$
- ▶ Degree of “**interestingness**” of feature  $j$  given  $\sigma_\ell$ :  $T_j^{\sigma_\ell} := t^{\text{diff}}(\mathbf{X}_j^{\sigma_\ell}, \mathbf{Y}_j^{\sigma_\ell})$
- ▶ Sort  $\{T_j^{\sigma_\ell}\}_{\ell=0}^h$  so that

$$T_j^{(0)} \geq T_j^{(1)} \geq \dots \geq T_j^{(h)}$$

## Enrichment analysis with $m \neq n$ : contrast scores

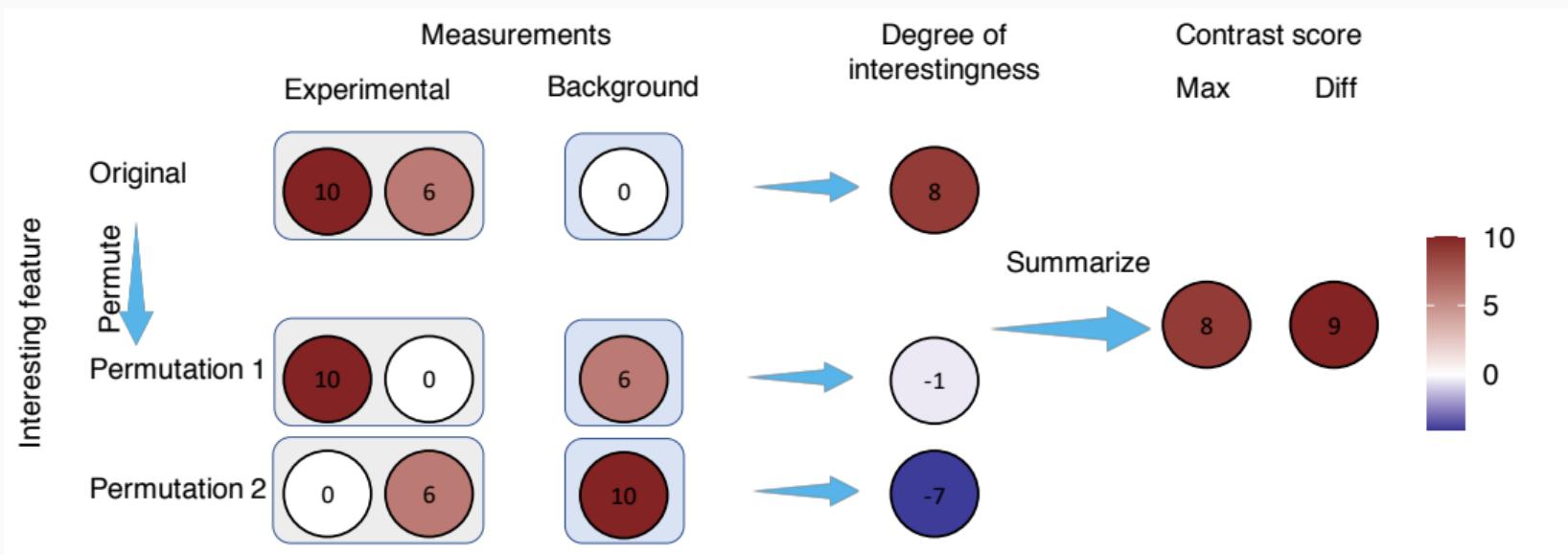
Contrast score of feature  $j$ :

$$C_j := \begin{cases} T_j^{(0)} - T_j^{(1)} & \text{if } T_j^{(0)} = T_j^{\sigma_0} \\ T_j^{(1)} - T_j^{(0)} & \text{otherwise} \end{cases} \quad \text{difference contrast score} \quad (8)$$

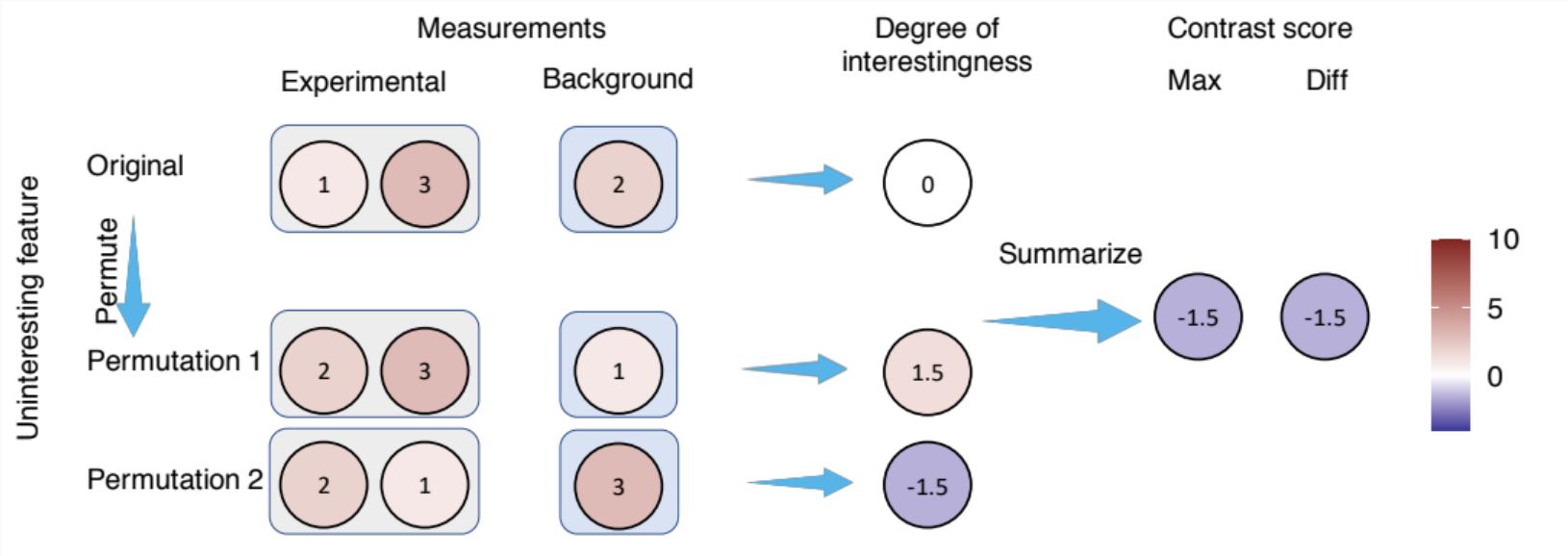
or

$$C_j := \begin{cases} |T_j^{(0)}| & \text{if } T_j^{(0)} = T_j^{\sigma_0} > T_j^{(1)} \\ 0 & \text{if } T_j^{(0)} = T_j^{(1)} \\ -|T_j^{(0)}| & \text{otherwise} \end{cases} \quad \text{maximum contrast score} \quad (9)$$

## Enrichment analysis with $m = 2$ , $n = 1$ : contrast scores



## Enrichment analysis with $m = 2$ , $n = 1$ : contrast scores



## Enrichment analysis with $m \neq n$ : cutoff

### Definition 2 GZ procedure [Gimenez and Zou, PMLR, 2019]

- ▶ Given contrast scores  $\{C_j\}_{j=1}^d$ , define  $\mathcal{C} = \{|C_j| : C_j \neq 0 ; j = 1, \dots, d\}$
- ▶ Based on the target FDR threshold  $q \in (0, 1)$ , contrast-score cutoff  $T^{GZ}$ :

$$T^{GZ} := \min \left\{ t \in \mathcal{C} : \frac{\frac{1}{h} + \frac{1}{h} \text{card}(\{j : C_j \leq -t\})}{\text{card}(\{j : C_j \geq t\}) \vee 1} \leq q \right\} \quad (10)$$

- ▶ **Discoveries:**  $\{j : C_j \geq T^{GZ}\}$

## Key lemma to guarantee the theoretical FDR control by the GZ procedure

Define  $S_j = \text{sign}(C_j) \in \{-1, 0, 1\}$

$\mathcal{N}$ : the set of uninteresting features

Then

1.  $S_1, \dots, S_d$  are mutually independent ;
2.  $\mathbb{P}(S_j = 1) \leq \frac{1}{h+1}$  for all  $j \in \mathcal{N}$ ;
3.  $\{S_j\}_{j \in \mathcal{N}} \perp \mathcal{C}$ .

## Differential analysis

---

- ▶ Almost the same as enrichment analysis with  $m \neq n$
- ▶ Only difference: the degree of “interestingness” of feature  $j$ :

- Differential:

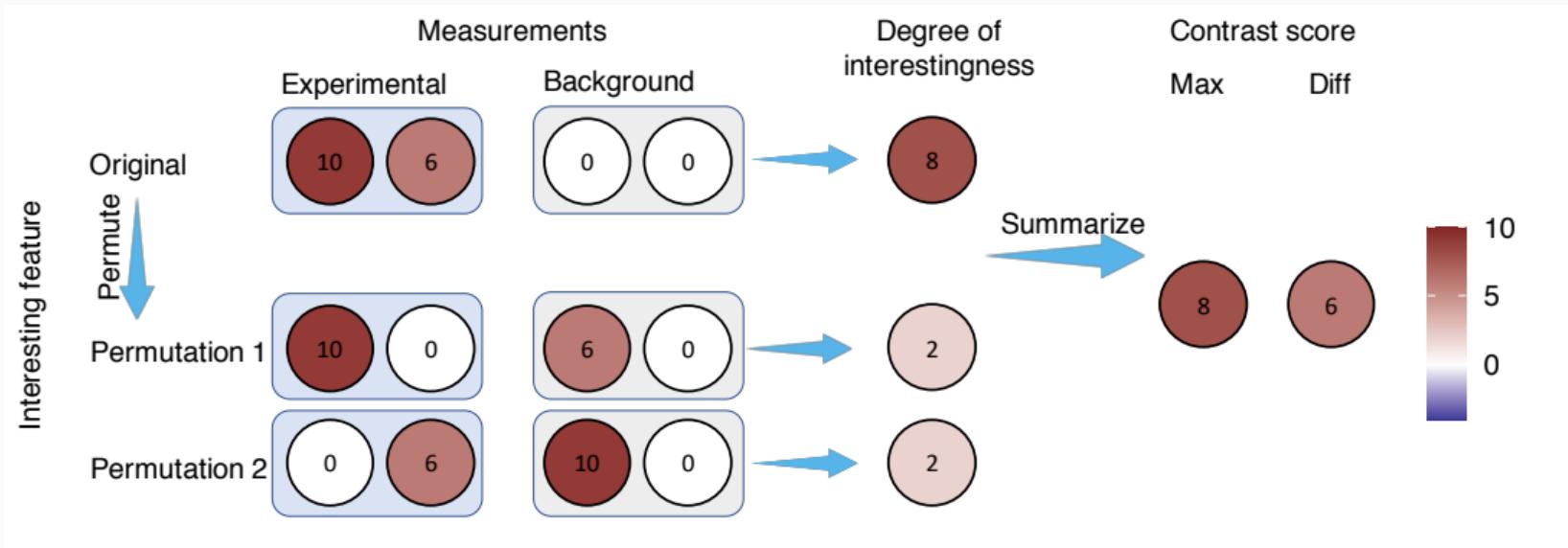
$$T_j^{\sigma_\ell} := \left| t^{\text{diff}}(\mathbf{X}_j^{\sigma_\ell}, \mathbf{Y}_j^{\sigma_\ell}) \right|$$

- Enrichment:

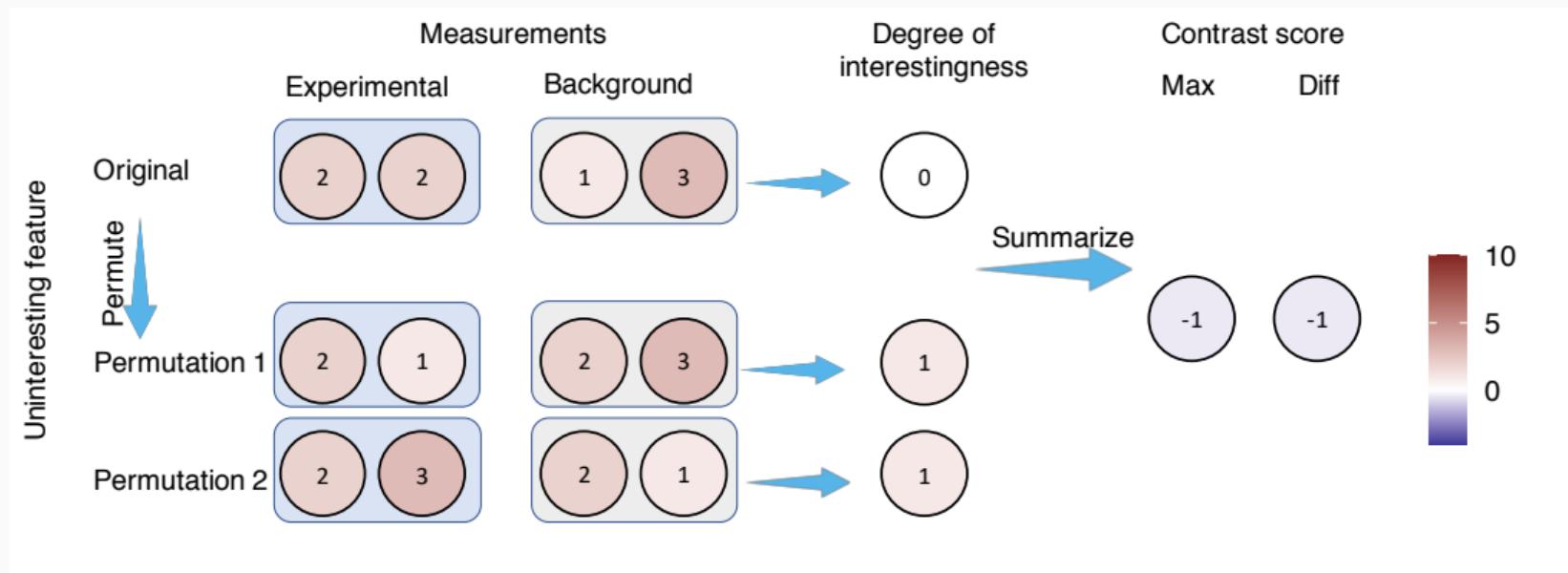
$$T_j^{\sigma_\ell} := t^{\text{diff}}(\mathbf{X}_j^{\sigma_\ell}, \mathbf{Y}_j^{\sigma_\ell})$$

- ▶ Contrast-score cutoff also by the GZ procedure

## Differential analysis with $m = n = 2$ : contrast scores



## Differential analysis with $m = n = 2$ : contrast scores



## Simulation analysis

---

## Simulation design

---

- ▶ **Analysis task:** enrichment or differential
- ▶ **Numbers of replicates**  $mvsn$ : 1vs1, 2vs1, 3vs3, or 10vs10
- ▶ **Distribution:** Gaussian, Poisson, or negative binomial
- ▶ **Background:** homogeneous or heterogeneous
- ▶ **Number of features:**  $d = 1,000$  or  $10,000$

## ► P-value based methods

- p-value calculation approach (paired or pooled)
- p-value thresholding procedure (BH or Storey's qvalue)

⇒ BH-pair  
BH-pool  
qvalue-pair  
qvalue-pool

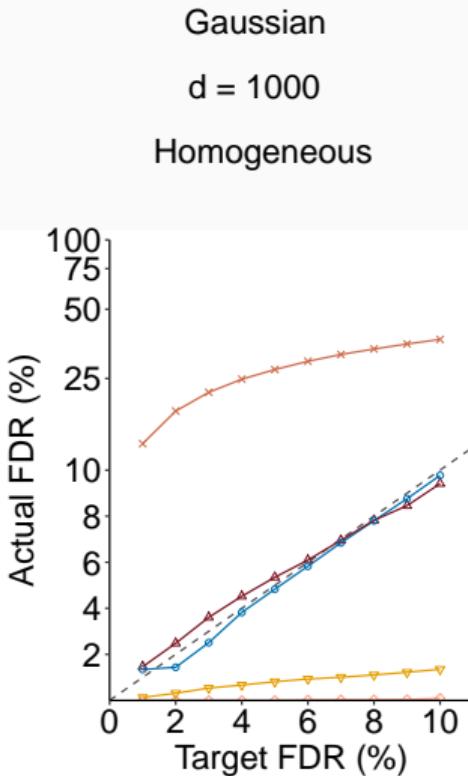
## ► Local fdr based methods

- empirical null (Gaussian)
- swapping

⇒ locfdr-emp  
locfdr-swap

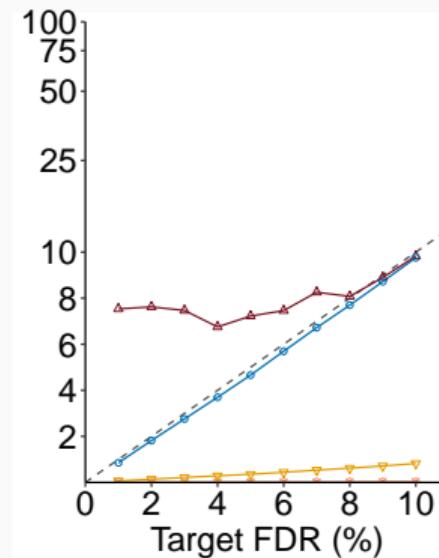
# Simulation results

## 1vs1 Enrichment



## 1vs1 Enrichment

Gaussian  
 $d = 10000$   
Heterogeneous



- Clipper
- BH-pool
- BH-pair-2as1
- BH-pair-mis
- locfdr-emp

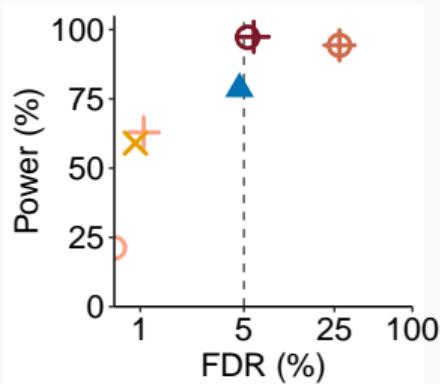
# Simulation results: target FDR threshold $q = 5\%$

## 1vs1 Enrichment

Gaussian

$d = 1000$

Homogeneous

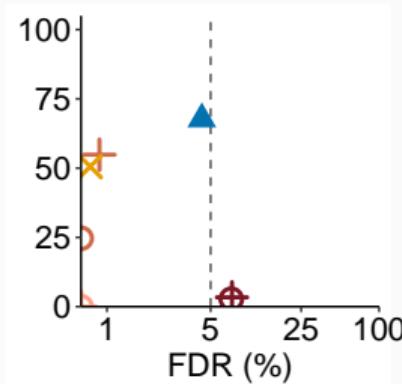


## 1vs1 Enrichment

Gaussian

$d = 10000$

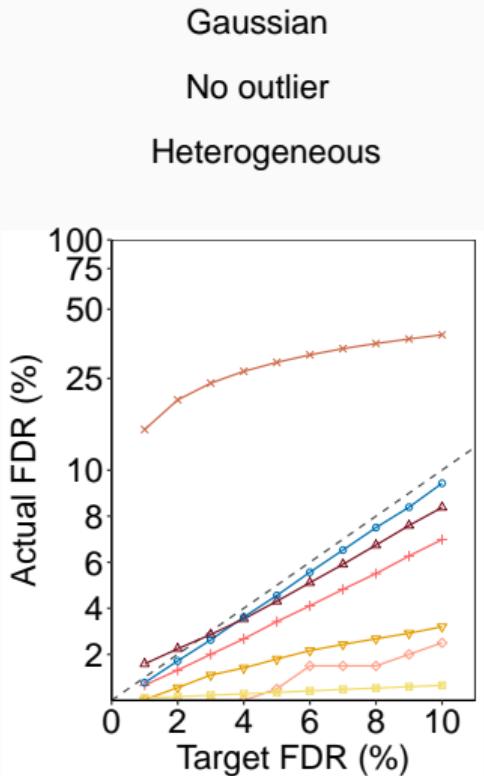
Heterogeneous



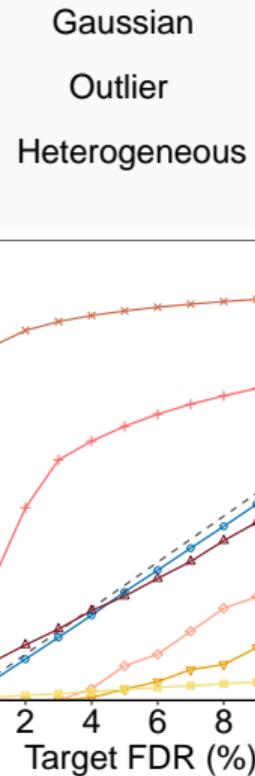
- ▲ Clipper
  - BH-pool
  - BH-pair-2as1
  - BH-pair-mis
  - +
  - +
  - +
  - +
- qvalue-pool  
qvalue-pair-2as1  
qvalue-pair-mis  
locfdr-emp

# Simulation results

## 3vs3 Enrichment



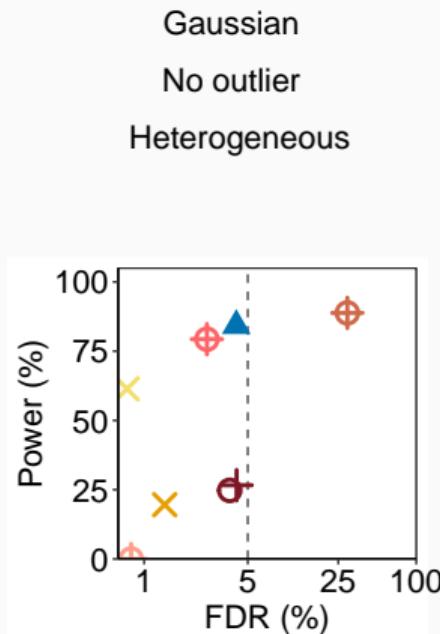
## 3vs3 Enrichment



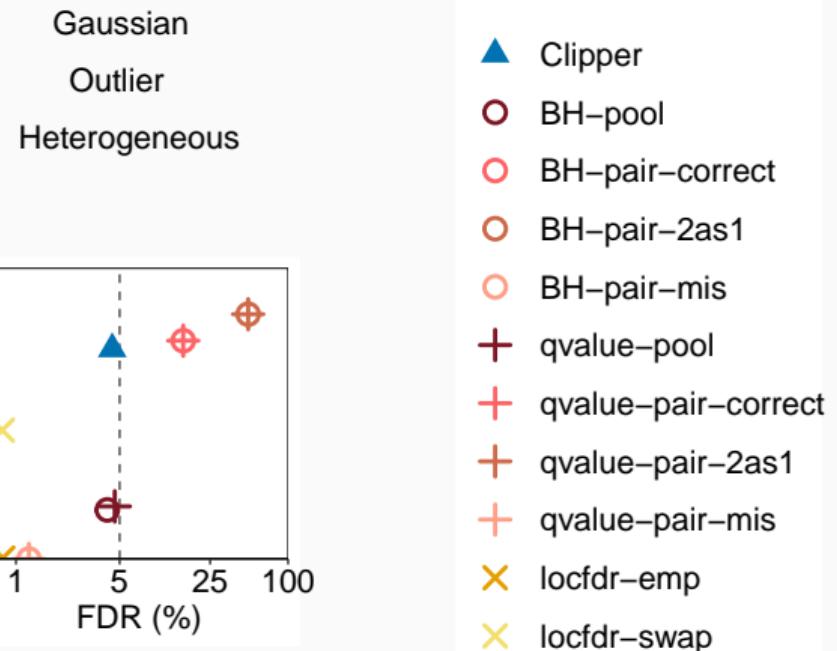
- Clipper
- BH-pool
- BH-pair-correct
- BH-pair-2as1
- BH-pair-mis
- locfdr-emp
- locfdr-swap

# Simulation results: target FDR threshold $q = 5\%$

## 3vs3 Enrichment

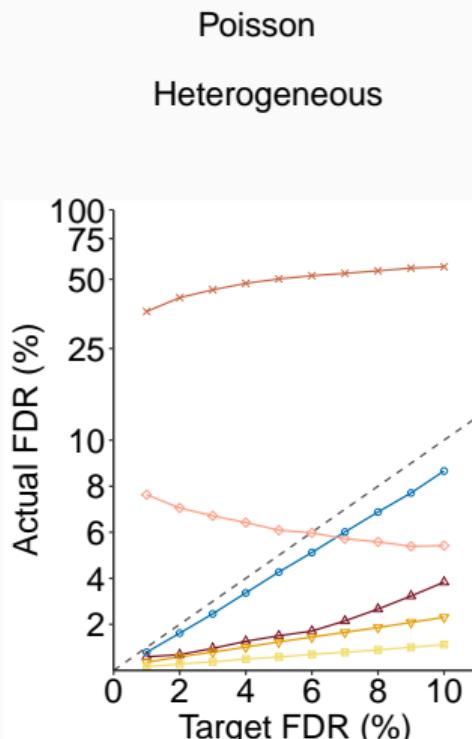


## 3vs3 Enrichment



# Simulation results

## 2vs1 Enrichment

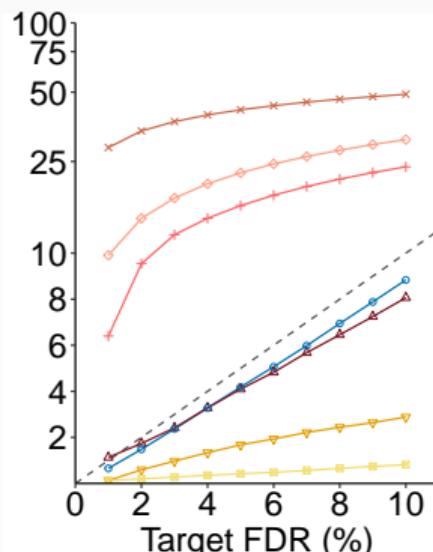


## 3vs3 Differential

Negative binomial

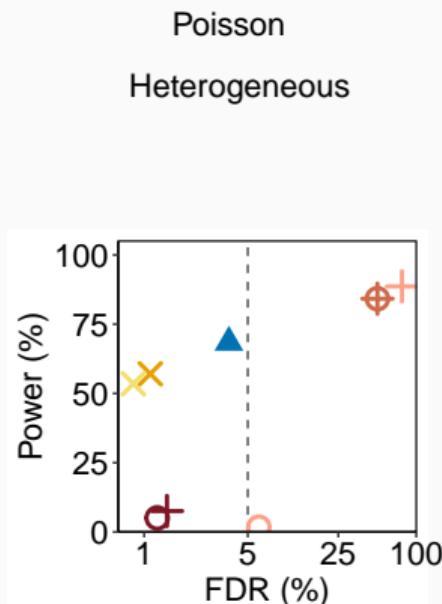
Heterogeneous

- 
- A legend listing eight methods with their corresponding line colors and markers:
- Clipper (blue circle)
  - BH-pool (red triangle)
  - BH-pair-correct (orange diamond)
  - BH-pair-2as1 (brown cross)
  - BH-pair-mis (pink circle)
  - locfdr-emp (yellow inverted triangle)
  - locfdr-swap (yellow square)

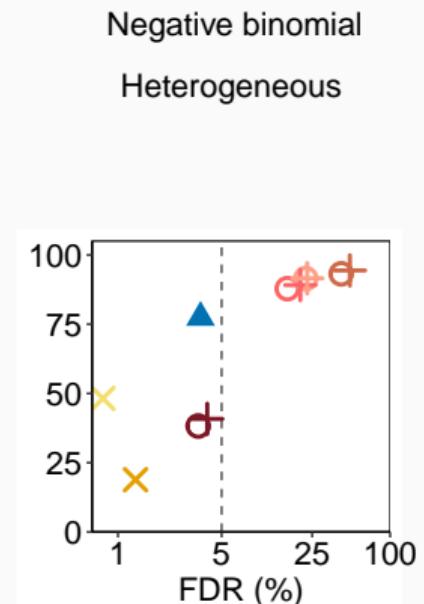


# Simulation results: target FDR threshold $q = 5\%$

## 2vs1 Enrichment



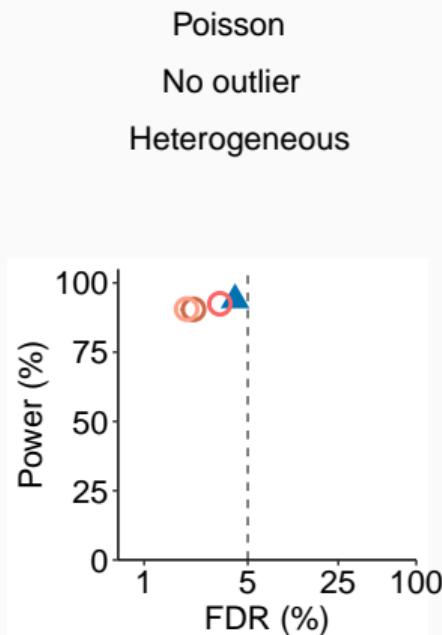
## 3vs3 Differential



- ▲ Clipper
- BH-pool
- BH-pair-correct
- BH-pair-2as1
- BH-pair-mis
- + qvalue-pool
- + qvalue-pair-correct
- + qvalue-pair-2as1
- + qvalue-pair-mis
- X locfdr-emp
- X locfdr-swap

# Simulation results: target FDR threshold $q = 5\%$

## 10vs10 Enrichment

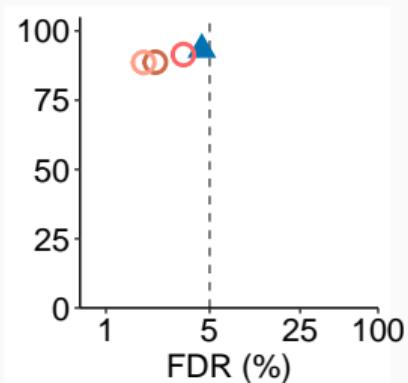


## 10vs10 Enrichment

Negative Binomial

No Outlier

Heterogeneous



- ▲ Clipper
- BH-pair-Wilcoxon
- BH-pair-parametric
- BH-pair-permutation

## **High-throughput biological data applications**

---

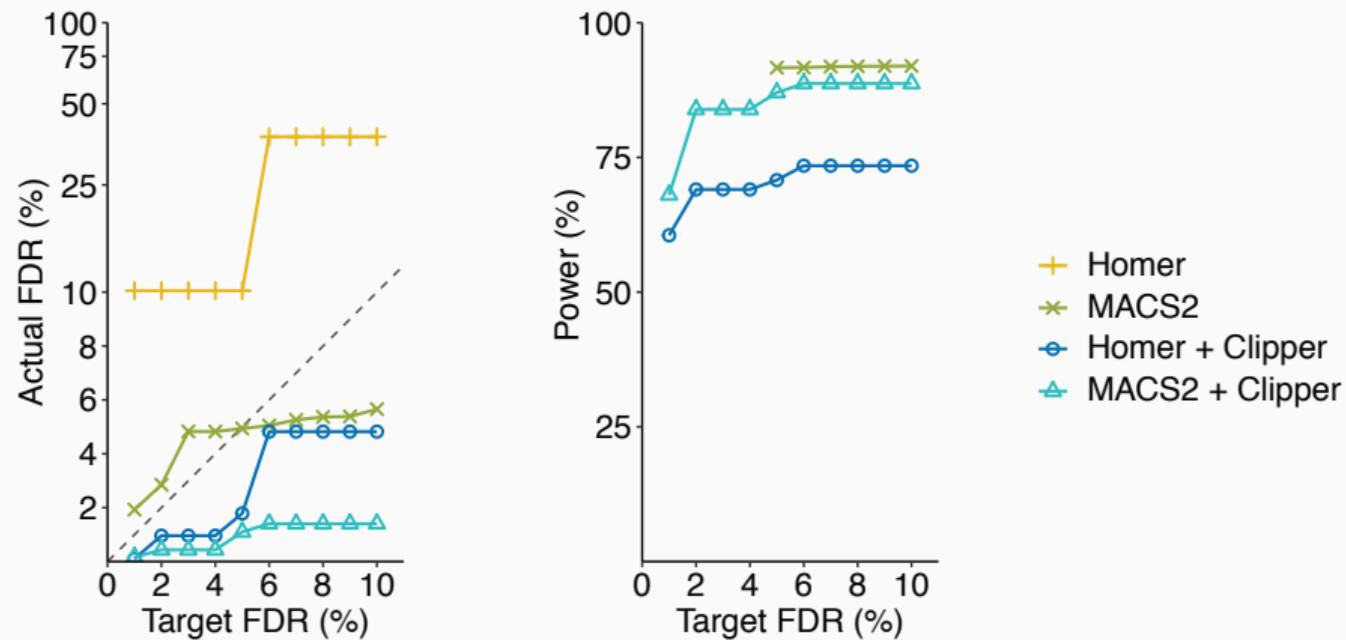
# High-throughput biological data applications

We compare Clipper with mainstream bioinformatics methods

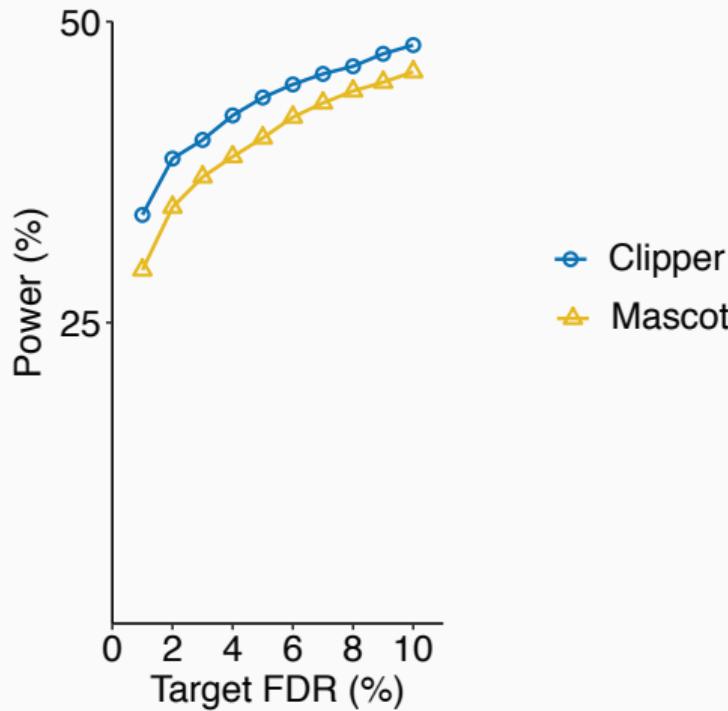
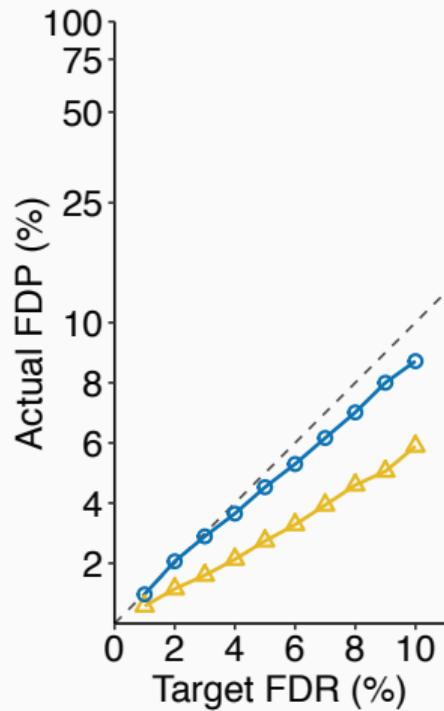
- ▶ Peaking calling from ChIP-seq data:
  - **MACS2** [Zhang et al., *Genome Biol*, 2008]
  - **Homer** [Heinz et al., *Mol Cell*, 2010]
- ▶ Peptide identification from MS data
  - **Mascot** [Spivak et al., *J Proteome Res*, 2009]
- ▶ Identification of differentially expressed genes from RNA-seq data
  - **DESeq2** [Love et al., *Genome Biol*, 2014]
  - **edgeR** [Robinson et al., *Bioinformatics*, 2010]
  - Covariate-based p-value weighting: **IHW** [Ignatiadis et al., *Nat Methods*, 2016]
- ▶ Identification of differentially interacting regions from Hi-C data
  - **diffHic** [Lun et al., *BMC Bioinformatics*, 2015]
  - **FIND** [Djekidel et al., *Genome Res*, 2018]
  - **multiHiCcompare** [Stansfield et al., *Bioinformatics*, 2019]



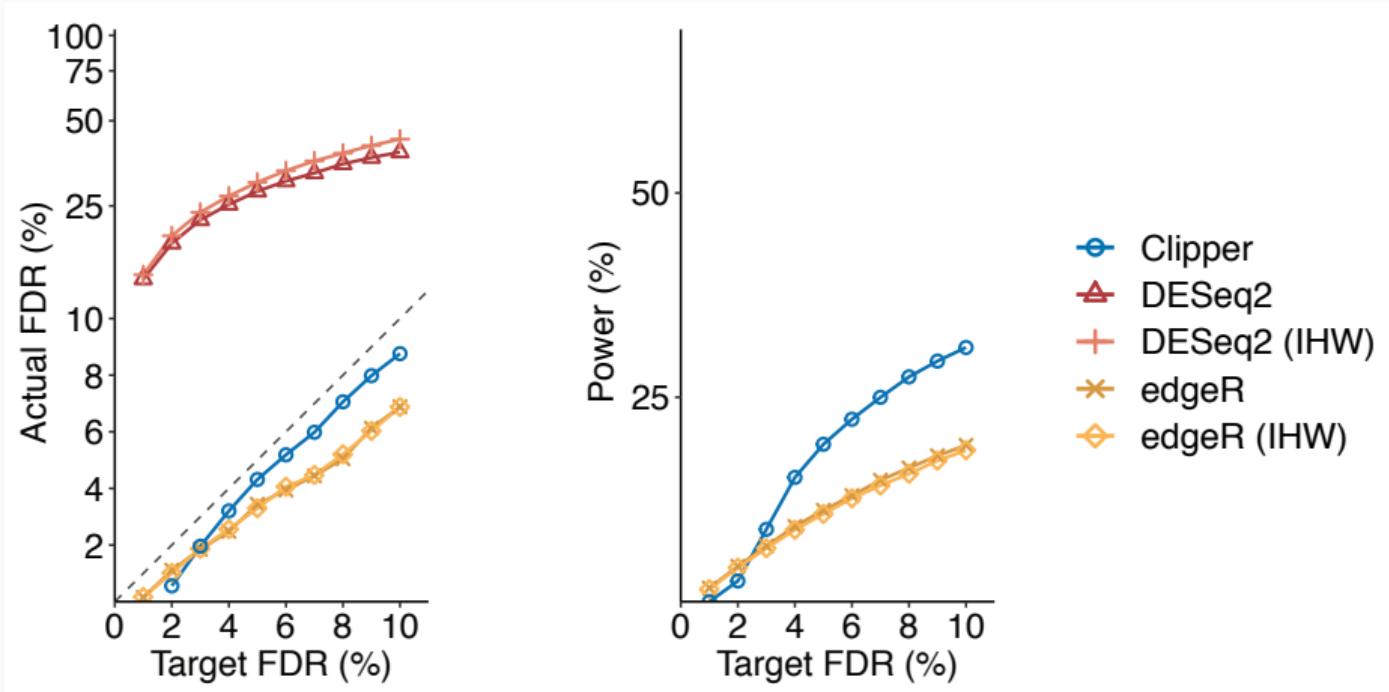
## Peaking calling from ChIP-seq data



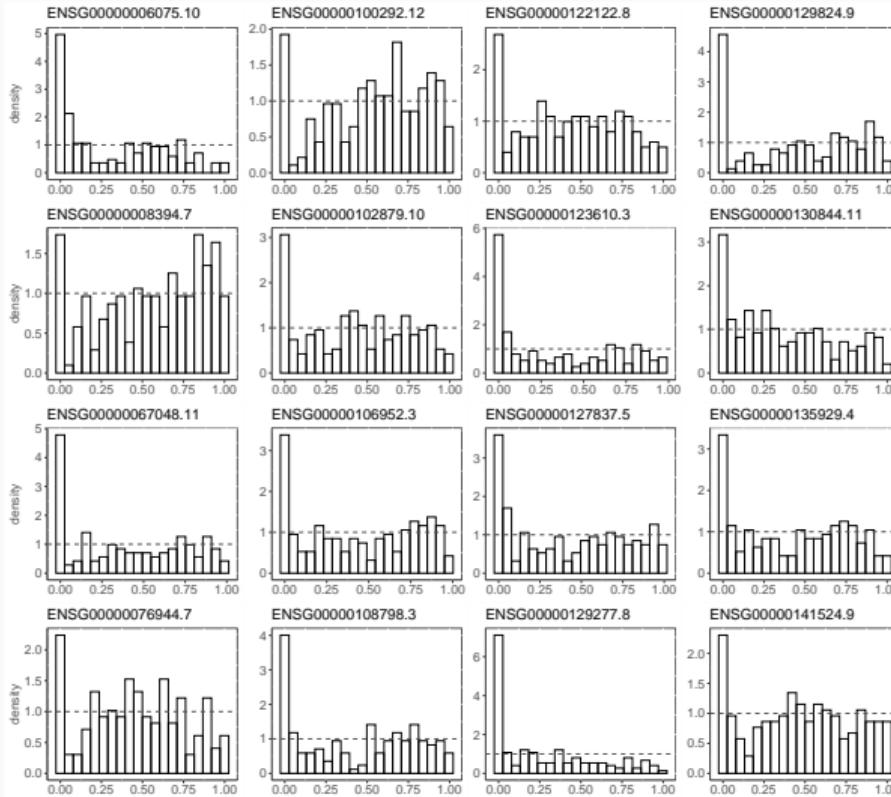
# Peptide identification from MS data



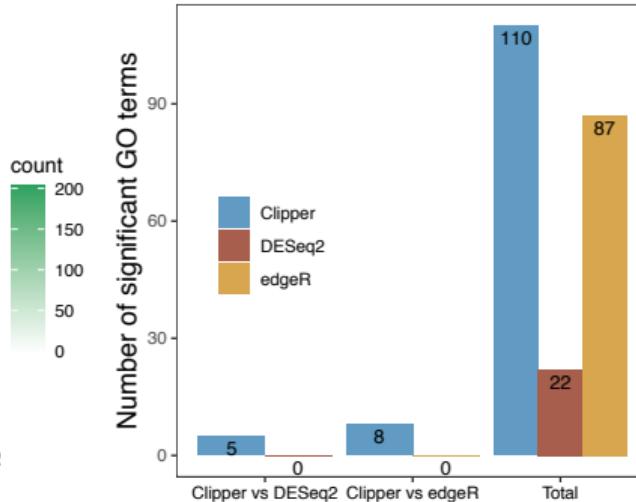
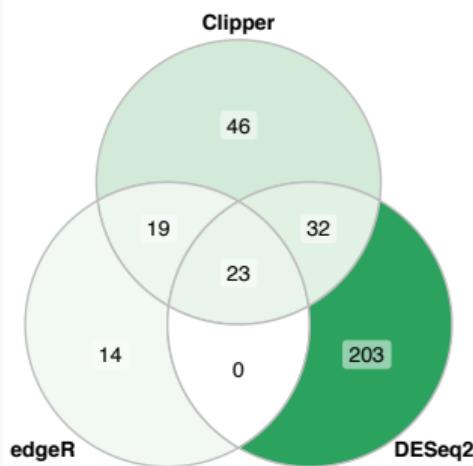
# Identification of differentially expressed genes from RNA-seq data



# The p-value distributions of 16 non-DEGs most frequently identified by DESeq2



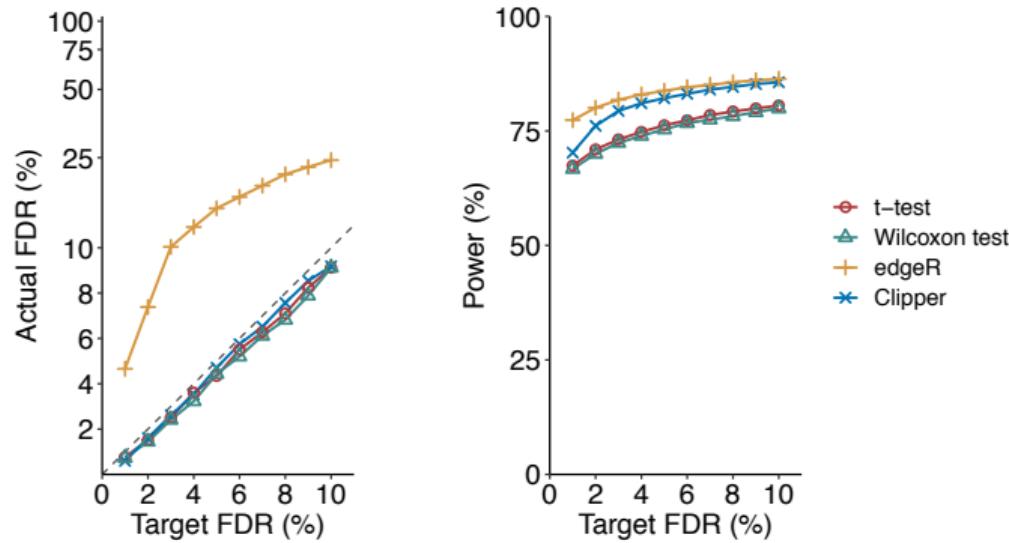
## Classical (inflammatory) vs. non-classical (chronical) monocytes



### Most significant GO term from DESeq2

GO term (ID)	qvalue (DESeq2)	qvalue (edgeR)	qvalue (Clipper)
leukocyte chemotaxis (GO:0030595)	9.930044e-06	9.594885e-09	3.104557e-10
myeloid leukocyte migration (GO:0097529)	1.107612e-05	2.921486e-08	5.740217e-10
granulocyte chemotaxis (GO:0071621)	2.698853e-05	1.008808e-08	1.167108e-09
neutrophil chemotaxis (GO:0030593)	2.698853e-05	2.921486e-08	2.691033e-09

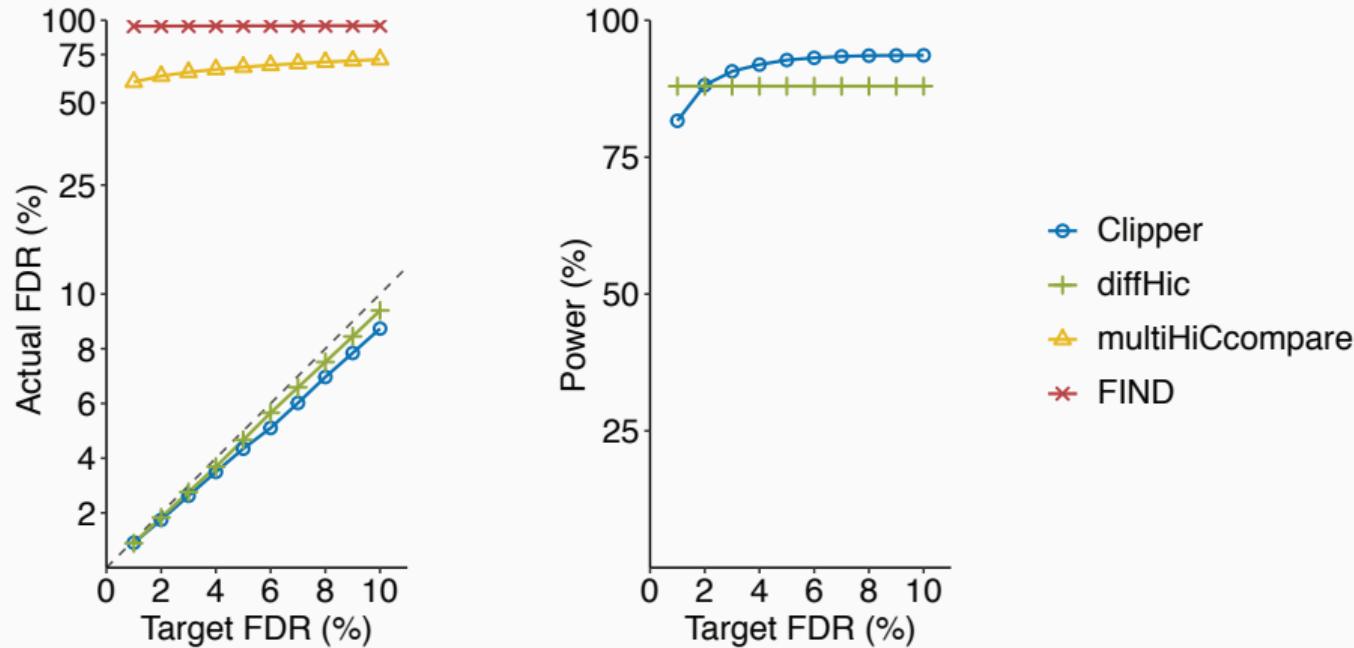
# Identification of differentially expressed genes from single-cell RNA-seq data



10x Genomics data generated by **scDesign2** <https://github.com/JSB-UCLA/scDesign2>

 “Bias, robustness and scalability in single-cell differential expression analysis” [Soneson and Robinson *Nat Methods*, 2018]

# Identification of differentially interacting regions from Hi-C data



## Discussion

---

## Discussion

---

- ▶ Clipper **avoids** the need for
  - valid (finite-sample or asymptotic) null distribution
  - high-resolution p-values
- ▶ Broad applications in high-throughput data analysis
  - key: **contrast scores**
- ▶ Importance of validating the FDR control
  - FDR control is bluntly assumed but **rarely validated** in most bioinformatics methods

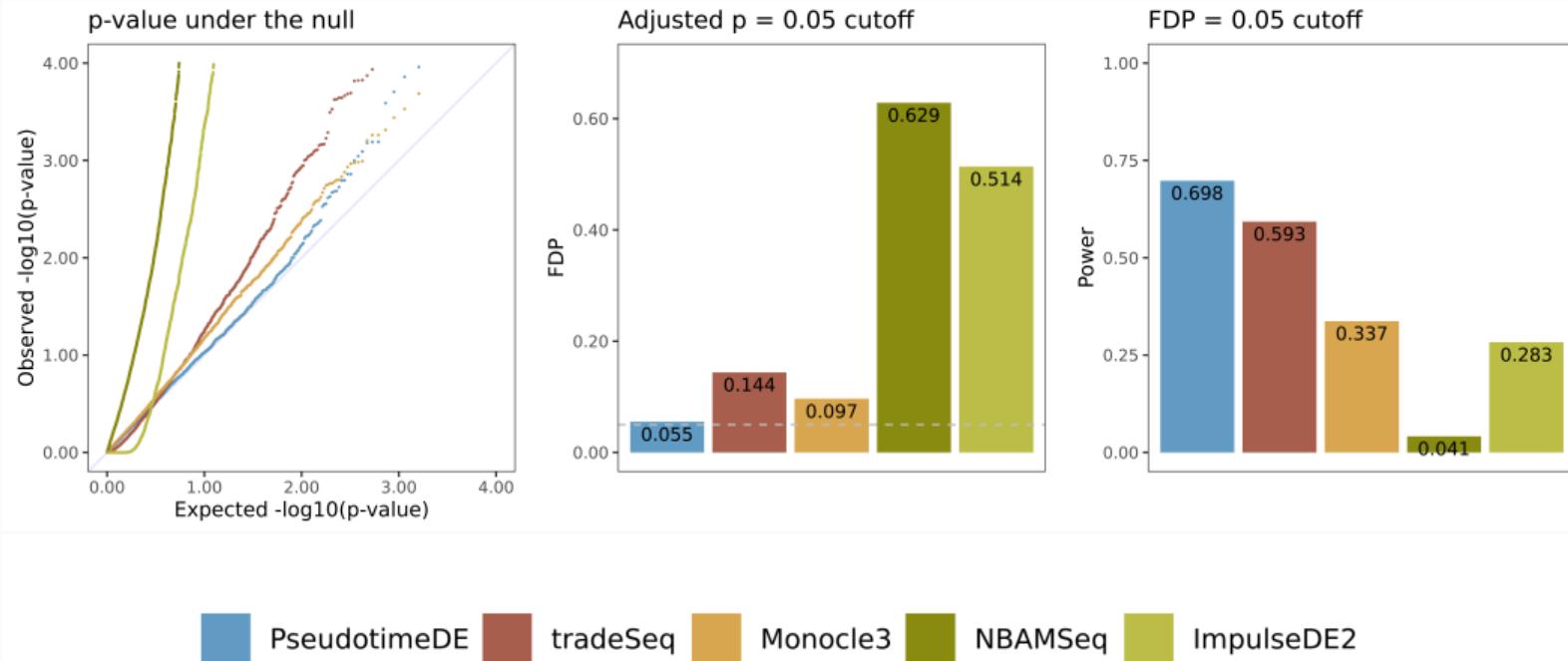
## Discussion

---

- ▶ Clipper avoids the need for
  - valid (finite-sample or asymptotic) null distribution
  - high-resolution p-values
- ▶ Broad applications in high-throughput data analysis
  - key: **contrast scores**
- ▶ Importance of validating the FDR control
  - FDR control is bluntly assumed but **rarely validated** in most bioinformatics methods
- ▶ Future work
  - ▶ incorporate **feature covariates** (e.g., gene variance)
  - ▶ choose contrast score & # of permutations: **power**
  - ▶ implement Clipper in bioinformatics tools



# PseudotimeDE: Identification of DEGs along Pseudotime with Valid p-values



Dongyuan Song R package:

<https://github.com/SONGDONGYUAN1994/PseudotimeDE>

# Acknowledgements

## UCLA JSB:

- ▶ Dongyuan Song
- ▶ Elaine Huang
- ▶ Tianyi Sun
- ▶ Kexin Li
- ▶ Other lab members

## City of Hope:

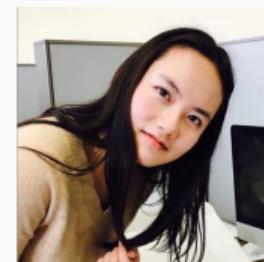
- ▶ Dr. Leo Wang
- ▶ MeiLu McDermott
- ▶ Kyla Woyshner
- ▶ Antigoni Manousopoulou

## UCI

- ▶ Dr. Wei Li



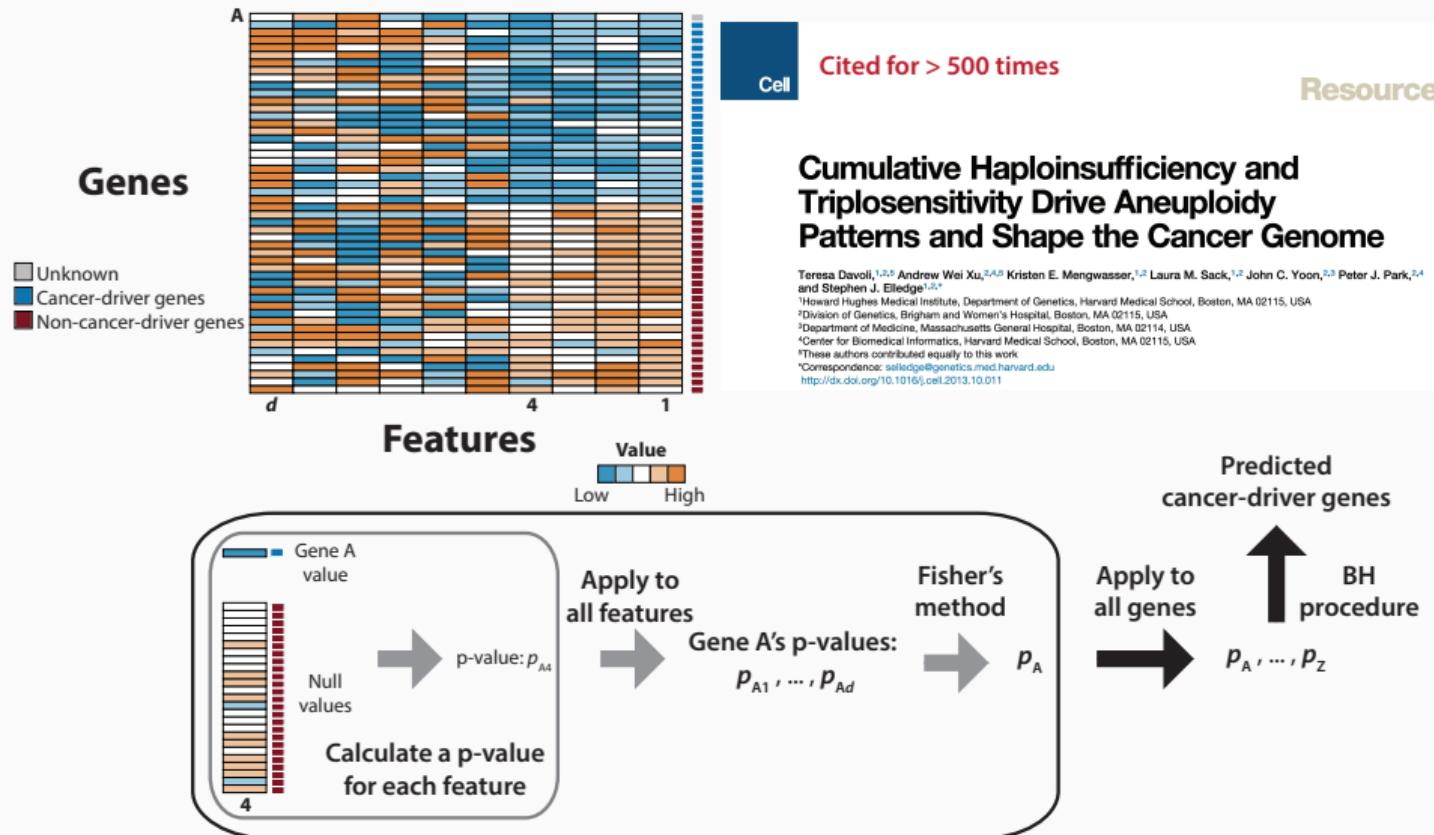
Xinzhou Ge



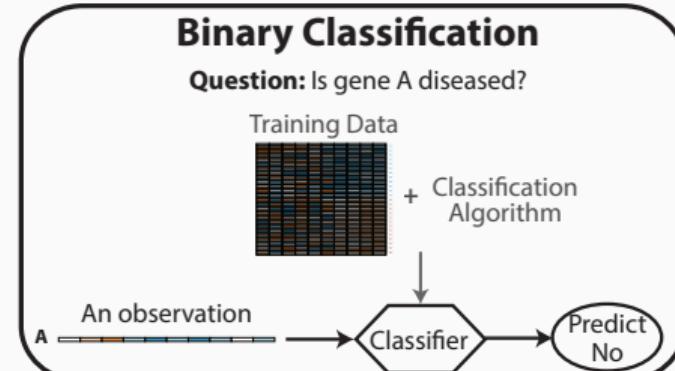
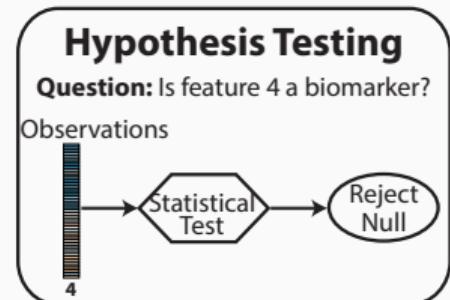
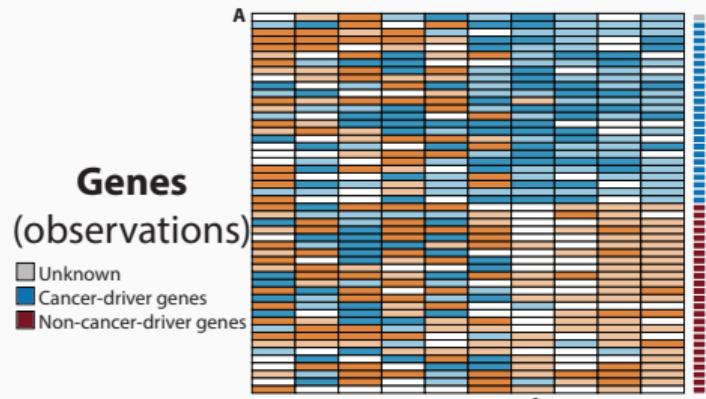
Yiling Elaine Chen



# Cancer-driver gene prediction: multiple testing?



# Multiple testing vs. binary classification



## CANCER

### DORGE: Discovery of Oncogenes and tumor suppressor genes using Genetic and Epigenetic features

Jie Lyu<sup>1\*</sup>, Jingyi Jessica Li<sup>2\*†</sup>, Jianzhong Su<sup>3</sup>, Fanglue Peng<sup>3</sup>, Yiling Elaine Chen<sup>2</sup>, Xinzhou Ge<sup>2</sup>, Wei Li<sup>1†</sup>

Patterns

CellPress  
OPEN ACCESS

## Perspective

### Statistical Hypothesis Testing versus Machine Learning Binary Classification: Distinctions and Guidelines

Jingyi Jessica Li<sup>1,\*</sup> and Xin Tong<sup>2</sup>

<sup>1</sup>Department of Statistics, University of California, Los Angeles, CA 90095-1554, USA

<sup>2</sup>Department of Data Sciences and Operations, Marshall School of Business, University of Southern California, Los Angeles, CA 90089, USA

\*Correspondence: jli@stat.ucla.edu

<https://doi.org/10.1101/patter.2020.100115>