
False Discovery Rate : enjeux et nouveaux défis

E. Roquain¹
Joint work with S. Delattre²

¹Laboratory LPMA, Université Pierre et Marie Curie (Paris 6), France

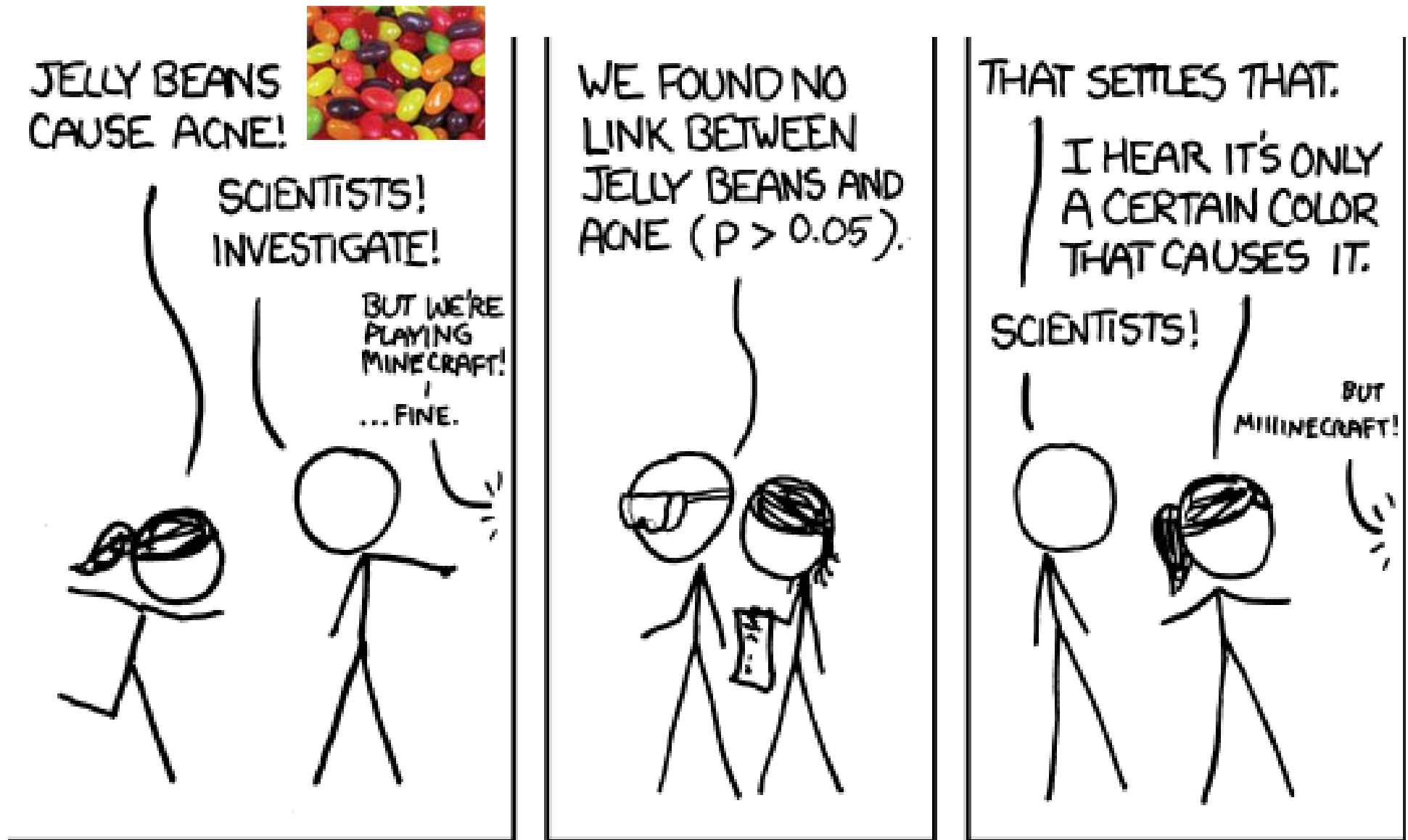
²Laboratory LPMA, Université Paris Diderot (Paris 7), France

Journées MAS - SMAI, 27 Août 2014

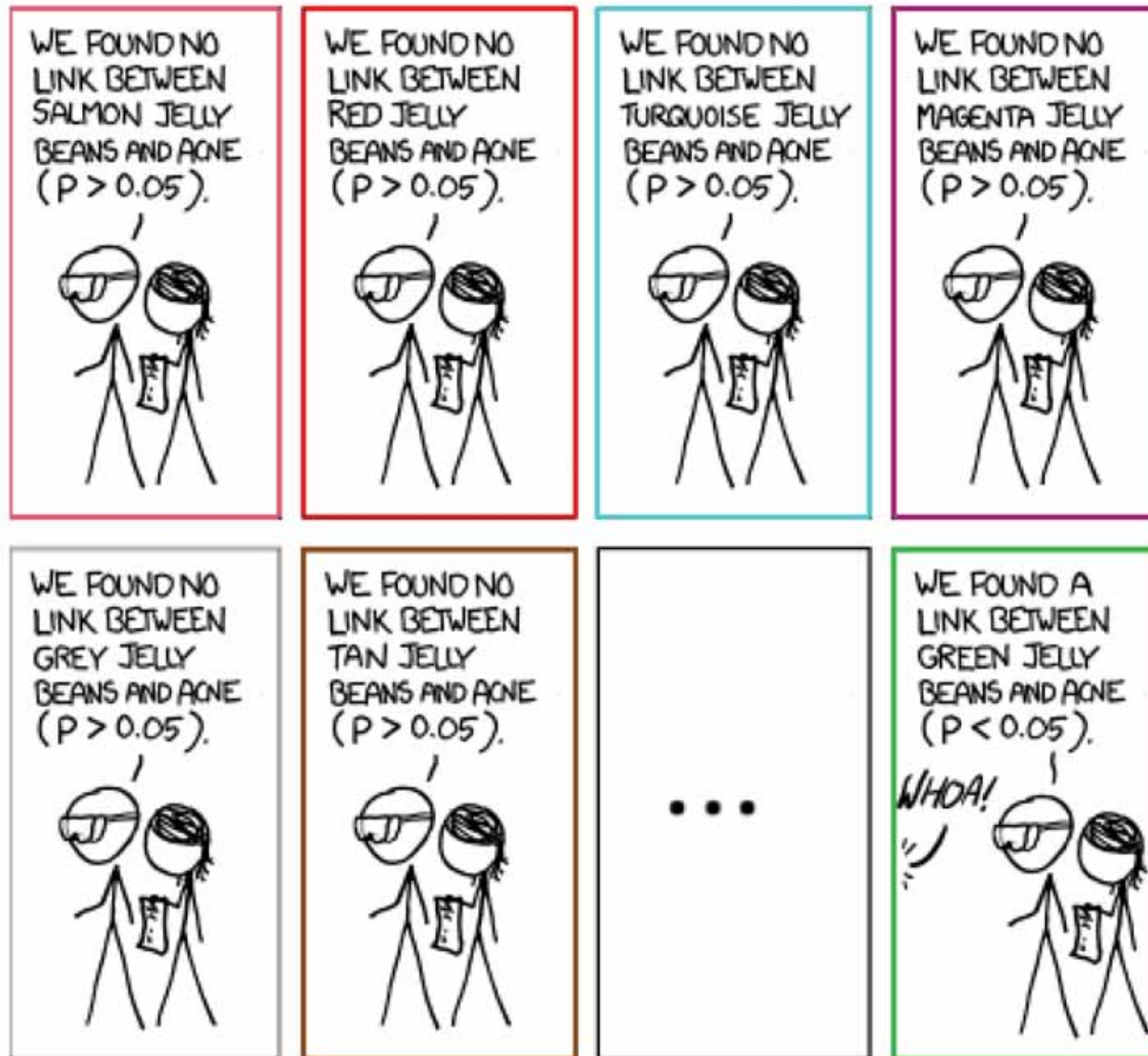
- 1 Introduction
- 2 False discovery rate control
- 3 Dependence and limitations
- 4 New challenges and results

- 1 Introduction
- 2 False discovery rate control
- 3 Dependence and limitations
- 4 New challenges and results

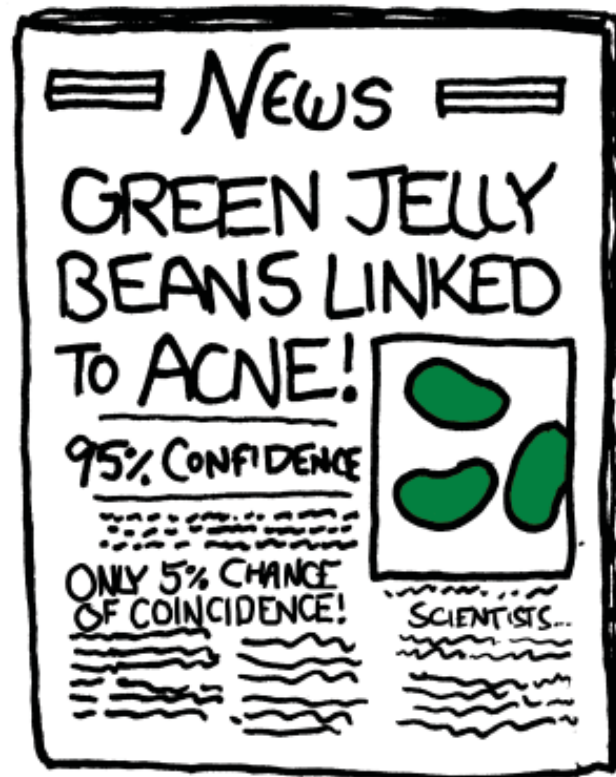
A “multiple testing joke” (<http://xkcd.com>)



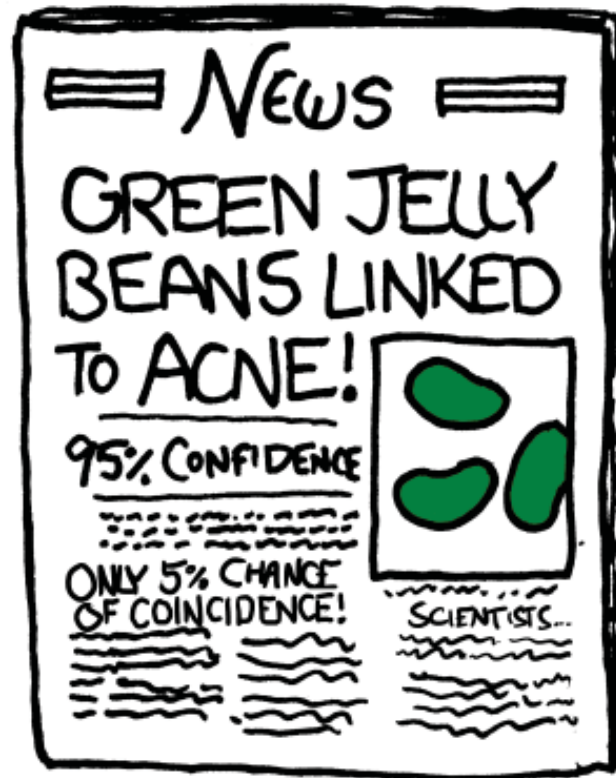
A “multiple testing joke” (<http://xkcd.com>)



A “multiple testing joke” (<http://xkcd.com>)



A “multiple testing joke” (<http://xkcd.com>)



Multiplicity problem

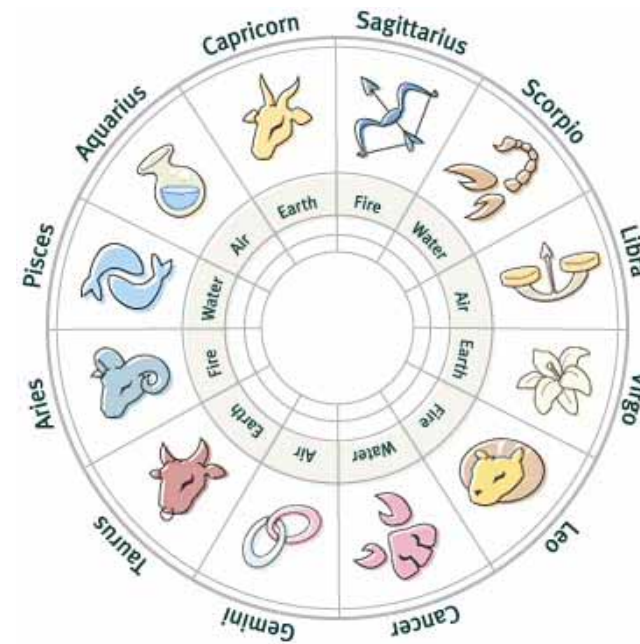
$\mathbf{P}(\text{make at least one false discovery}) \gg \mathbf{P}(\text{the } i\text{-th is a false discovery})$

A correction is needed to assess significance!

Some other examples

Paradoxes due to large scale experiments

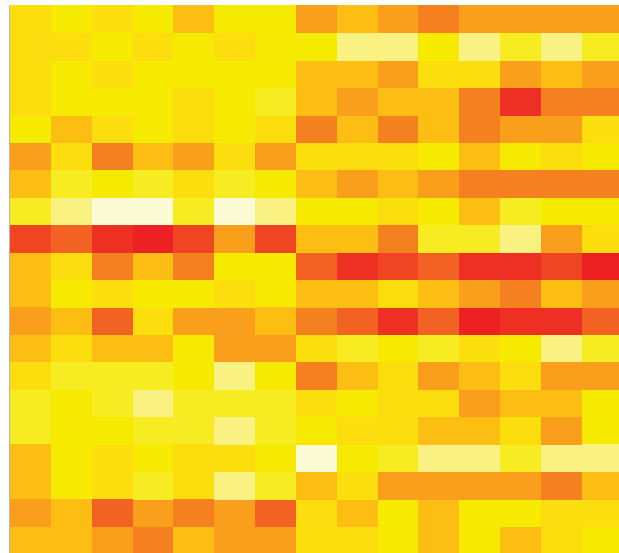
Probable facts appear significant



Multiplicity in microarray [Hedenfalk et al. (2001)]

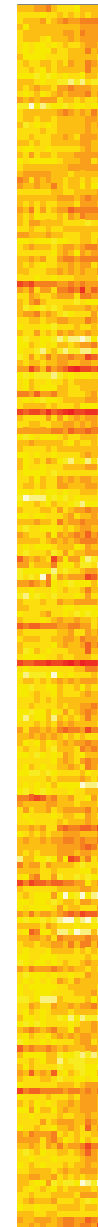
BRCA1 vs BRCA2

genes

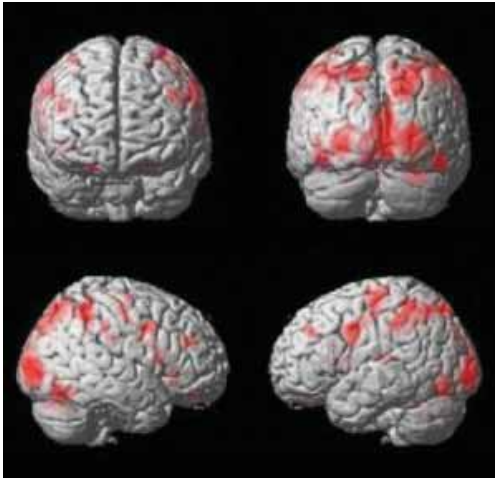


- ▶ expression level (activity)
- ▶ genes differentially activated?
- ▶ 1 test for each gene
- ▶ thousands of genes

- ▶ nb replications \ll dimension
- ▶ correlations



Other applications



- ▶ Neuroimaging (fMRI)
activated regions?
- ▶ Econometrics
winning strategies?
- ▶ Astronomy
directions with stars?



Canonical setting

- ▶ $X_i = \text{avg group 2} - \text{avg group 1}$ (rescaled) for genes i
- ▶ Gaussian model :

$$\begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_m \end{pmatrix} = \mu \begin{pmatrix} H_1 \\ H_2 \\ \vdots \\ H_m \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_m \end{pmatrix},$$

with $\mu > 0$, $H \in \{0, 1\}^m$ (fixed) and $\varepsilon \sim \mathcal{N}(0, \Gamma)$ ($\Gamma_{i,i} = 1$).

- ▶ $\Gamma = \text{dependence structure} = I_m$ for now

Question: for each i , $H_i = 0$ or $H_i = 1$?

Multiple testing : favors the "0" decision

Individual decision and errors

- ▶ Test statistic: X_i
- ▶ p -value: $p_i = \bar{\Phi}(X_i)$, with $\bar{\Phi}(z) = \mathbf{P}(Z \geq z)$, $Z \sim \mathcal{N}(0, 1)$

p_i such that

if $H_i = 0$, $p_i \sim U(0, 1)$

if $H_i = 1$, $p_i \sim \bar{\Phi}(\bar{\Phi}^{-1}(\cdot) - \mu)$

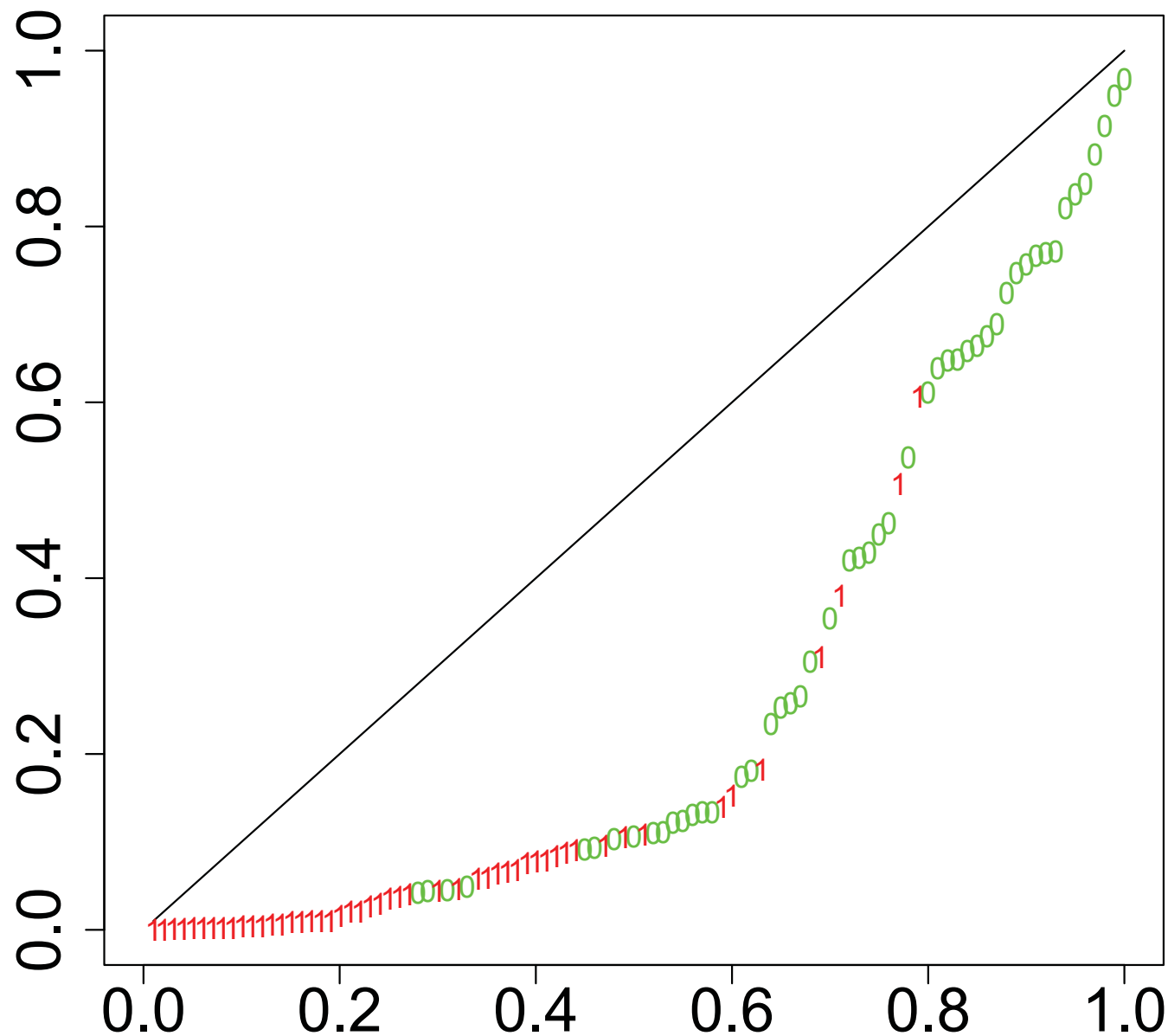
- ▶ Choose $\hat{H}_i = \mathbf{1}\{p_i \leq t\}$ for some threshold t

- ▶ Two errors:

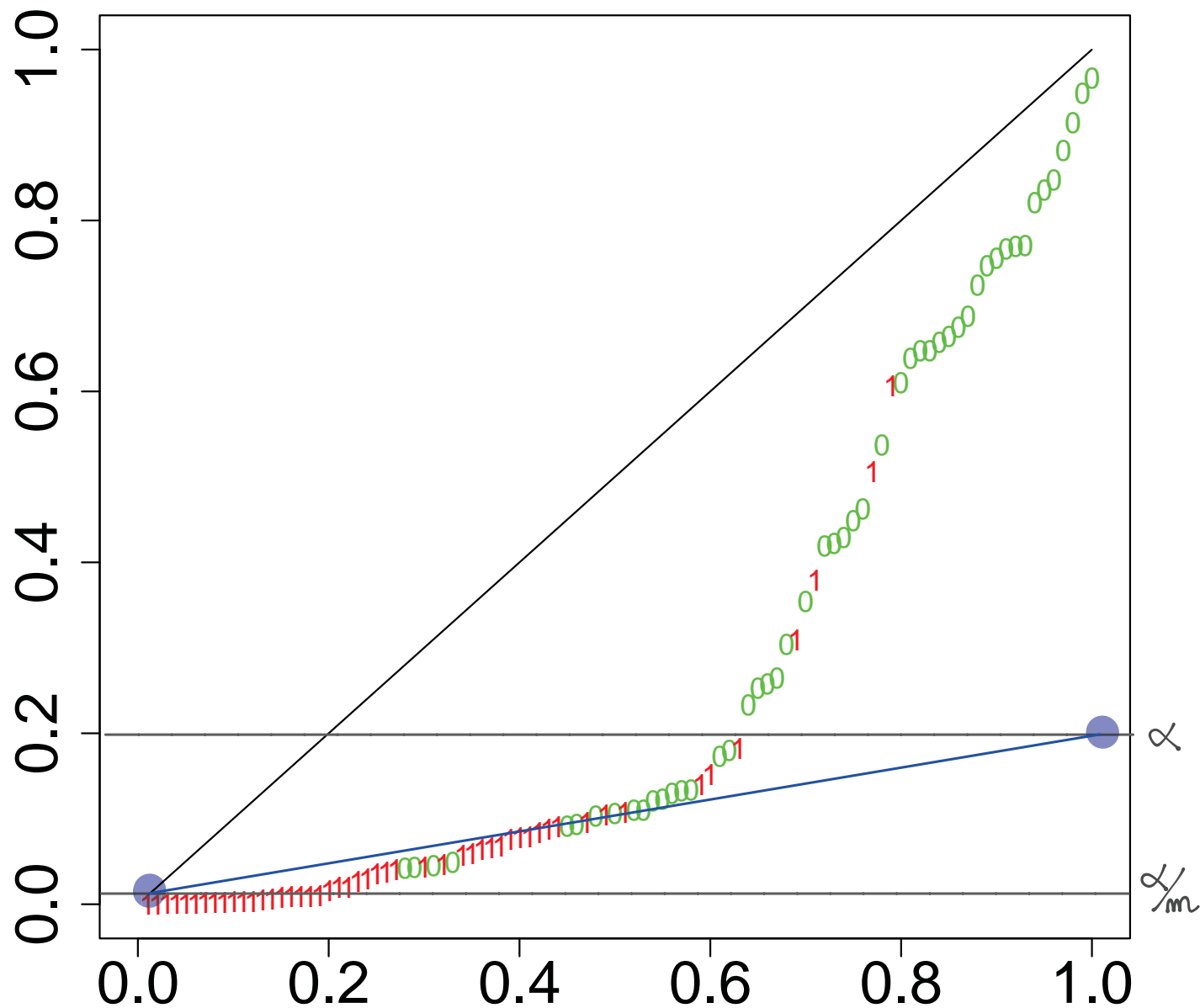
	$\hat{H}_i = 0$	$\hat{H}_i = 1$
$H_i = 0$	true negative	false positive
$H_i = 1$	false negative	true positive

- ▶ False positive more annoying

Data: $m = 100$; $m_0 = 50$; $\mu = 2$; $\Gamma = I_m$



BH thresholding



- 1 Introduction
- 2 False discovery rate control**
- 3 Dependence and limitations
- 4 New challenges and results

False discovery rate control

For a decision $\hat{H}_i = \mathbf{1}\{p_i \leq \hat{t}\}$ ($\forall i$),

$$\text{FDP}(\hat{t}) = \frac{\#\{i : H_i = 0, \hat{H}_i = 1\}}{\#\{i : \hat{H}_i = 1\}} \quad \left(\frac{0}{0} = 0 \right)$$

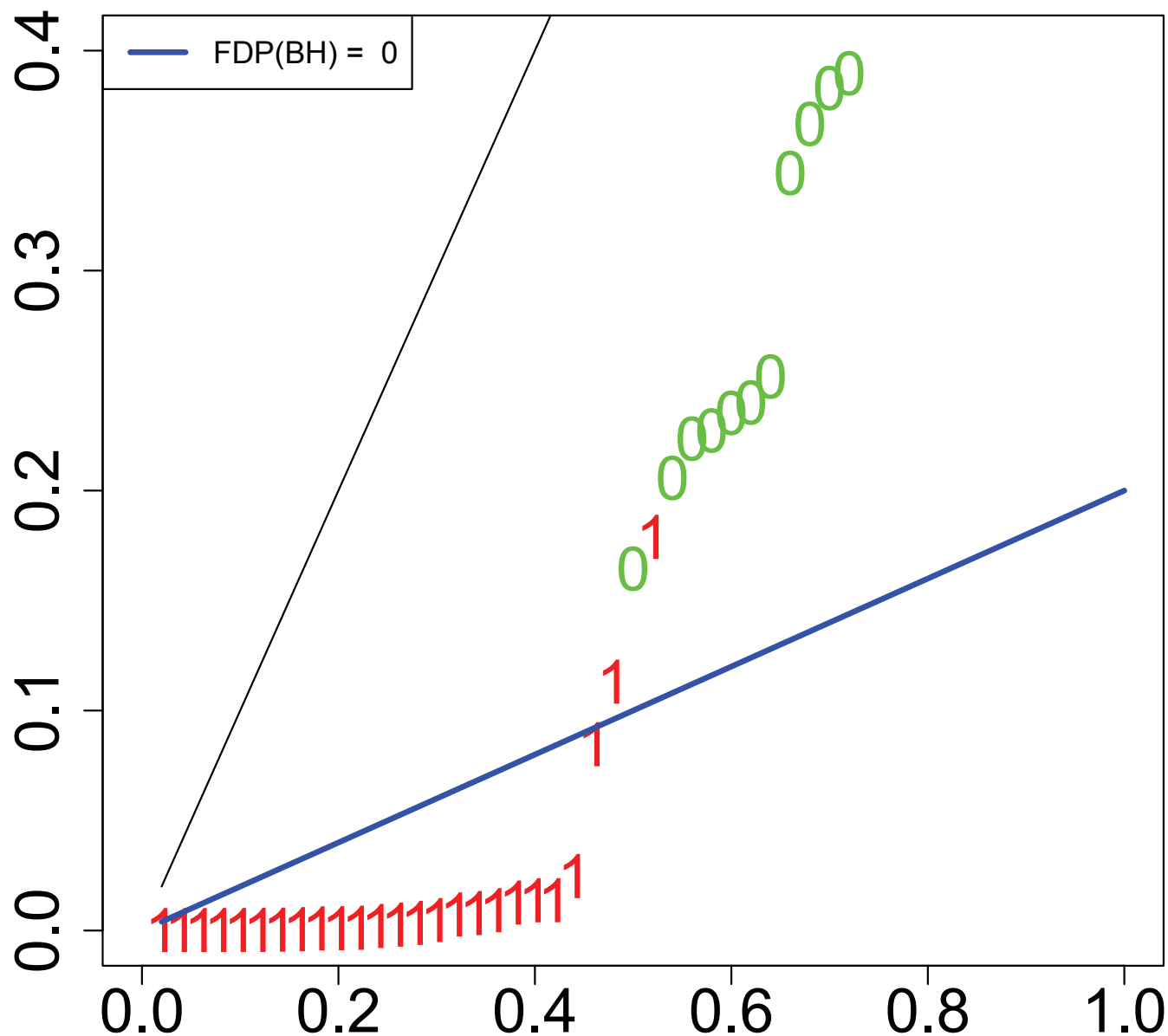
$$\text{FDR}(\hat{t}) = \mathbf{E}[\text{FDP}(\hat{t})]$$

Theorem [Benjamini and Hochberg (1995)] [Benjamini and Yekutieli (2001)]

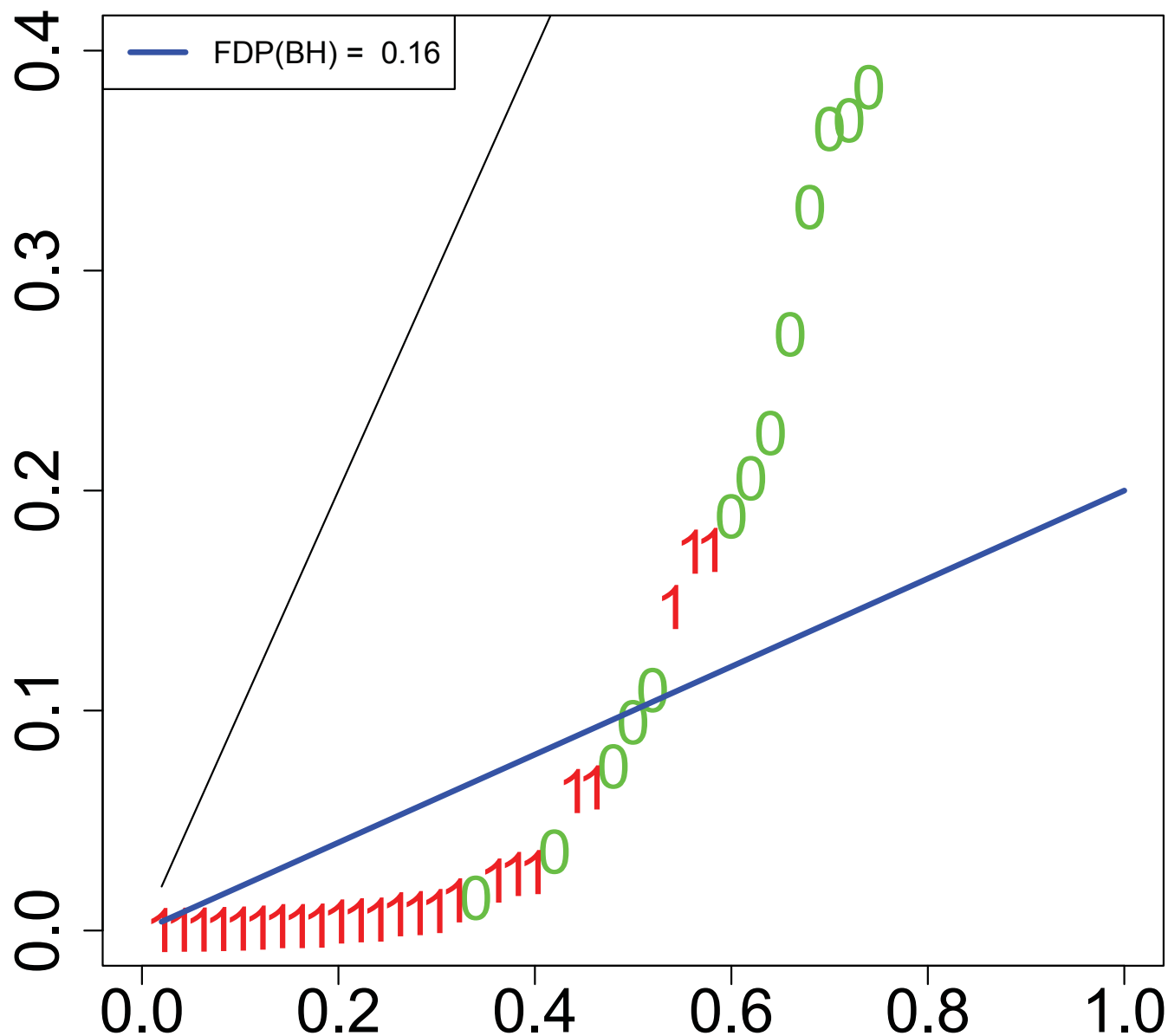
If $\Gamma = I_m$ and \hat{t} threshold of BH procedure, $\forall \mu, H$,

$$\text{FDR}(\hat{t}) = (m_0/m)\alpha \leq \alpha$$

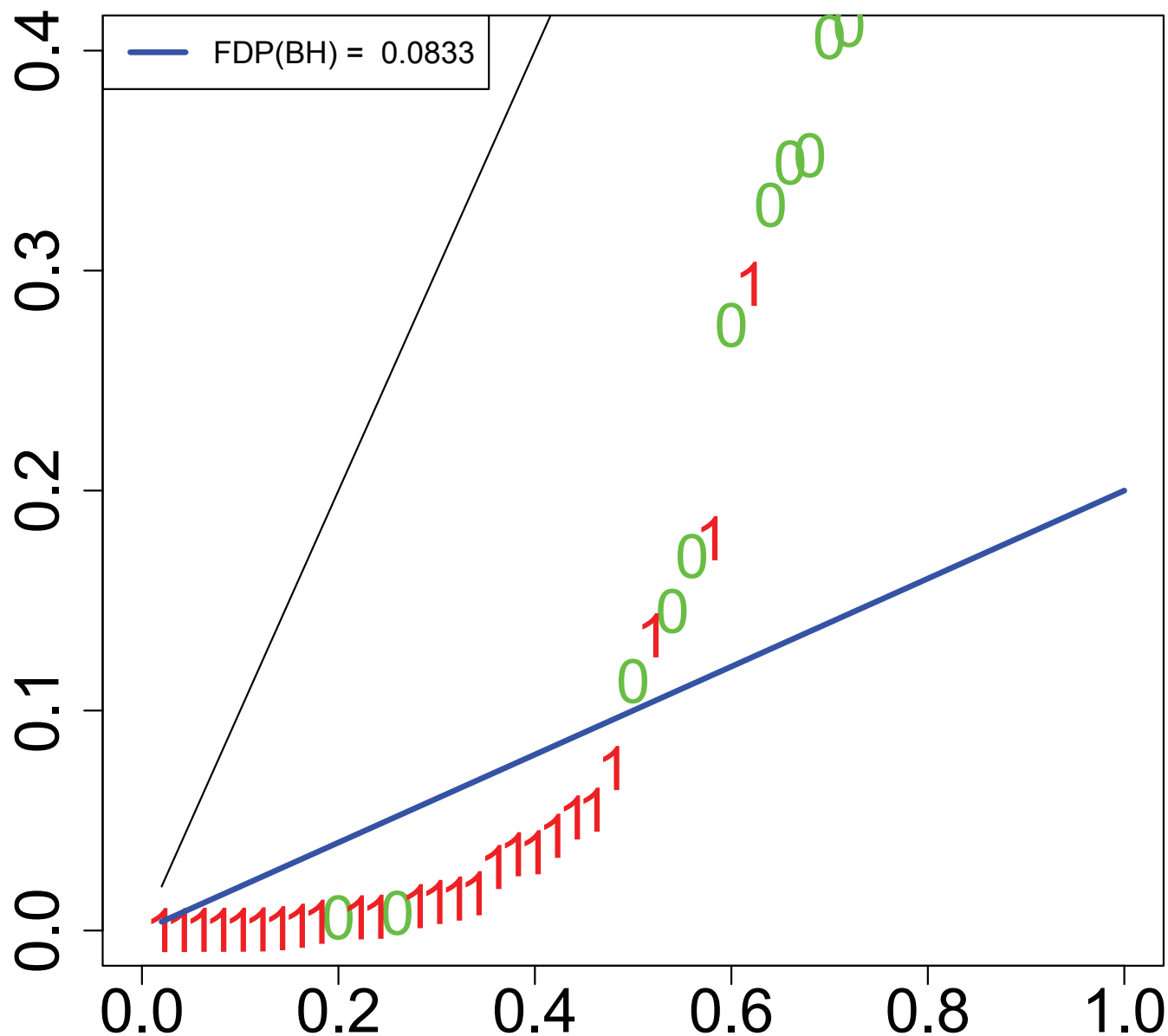
Simulations. $m = 50$; $m_0 = 25$; $\mu = 3$; $\Gamma = I_m$



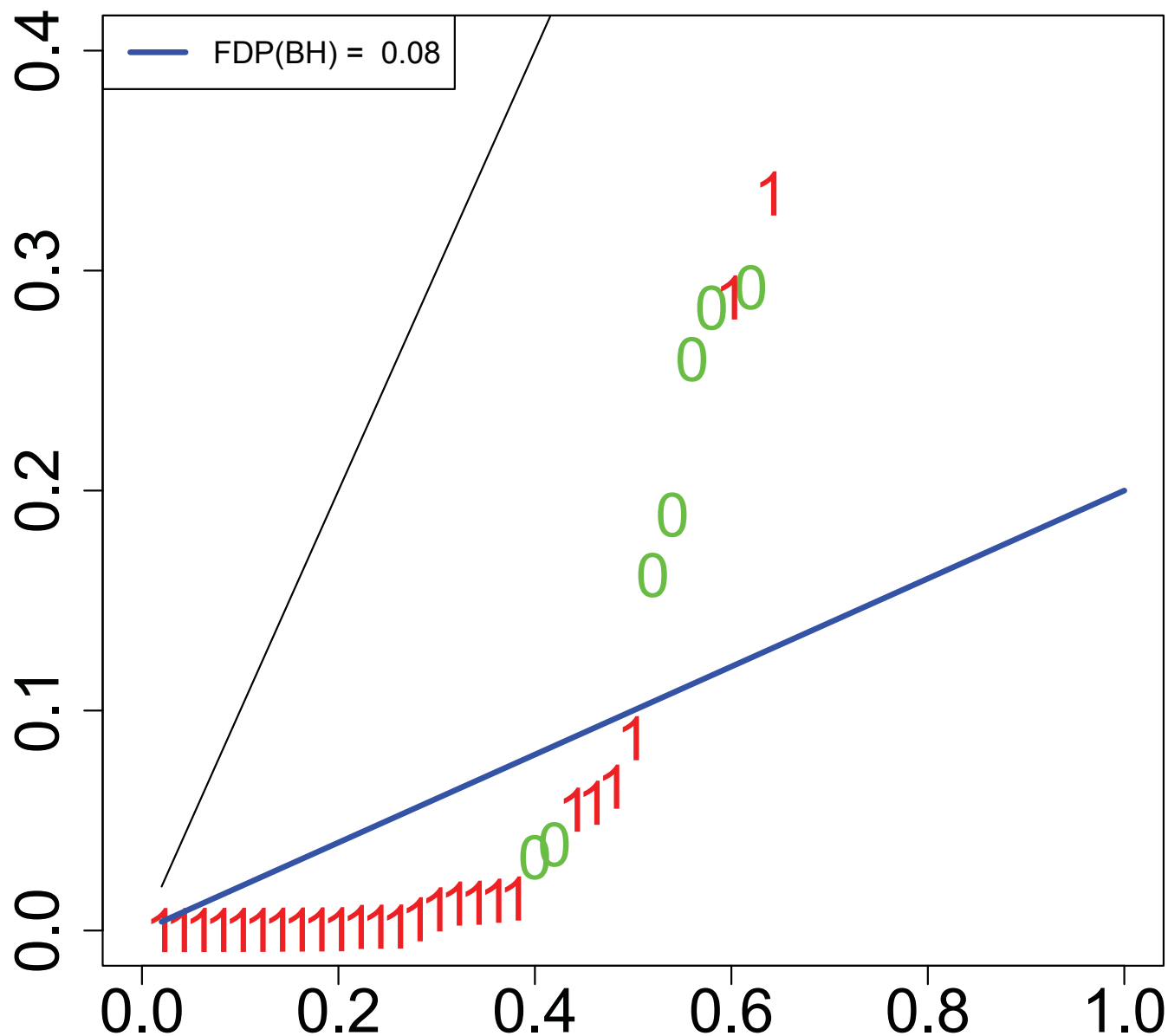
Simulations. $m = 50$; $m_0 = 25$; $\mu = 3$; $\Gamma = I_m$



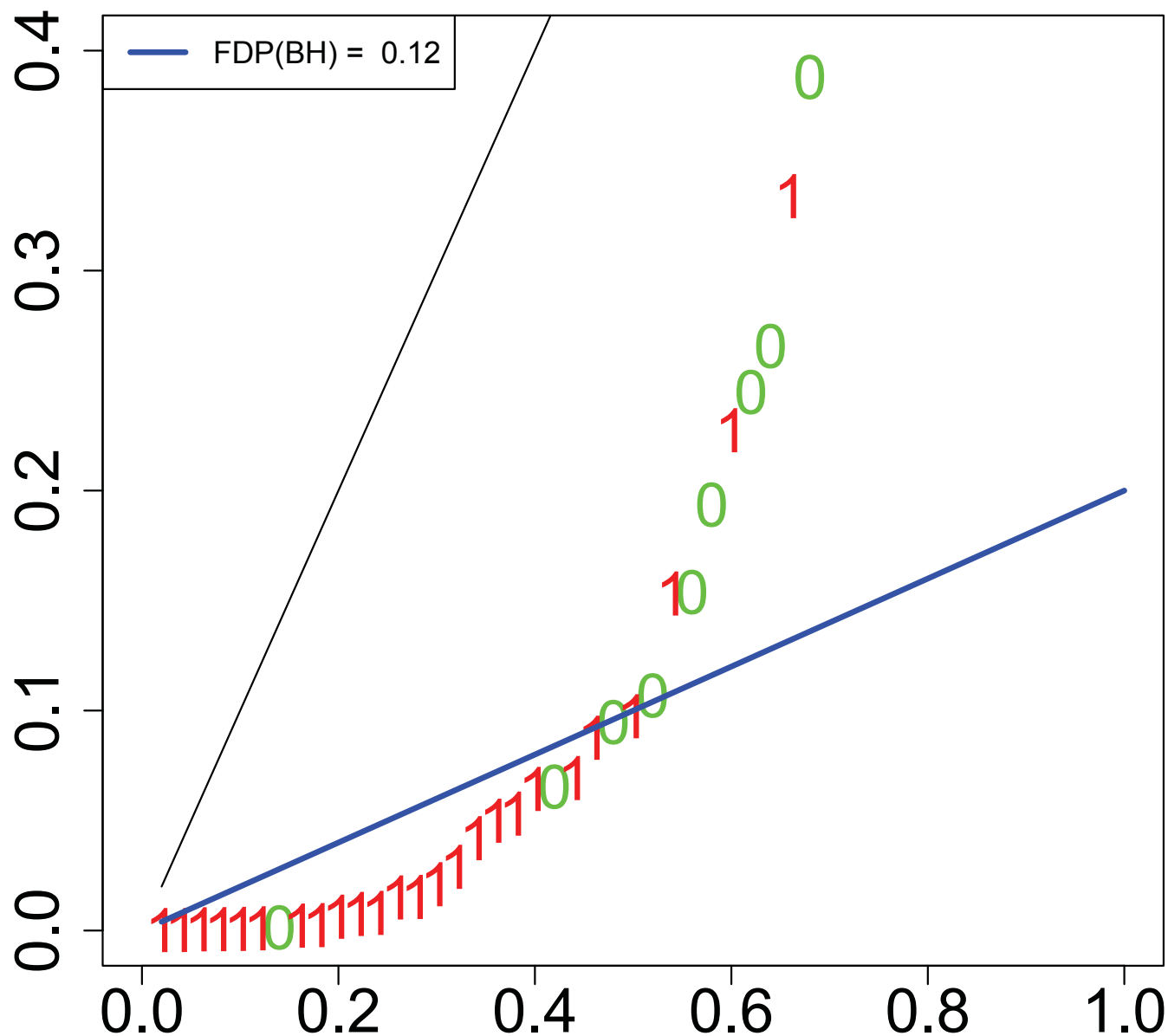
Simulations. $m = 50$; $m_0 = 25$; $\mu = 3$; $\Gamma = I_m$



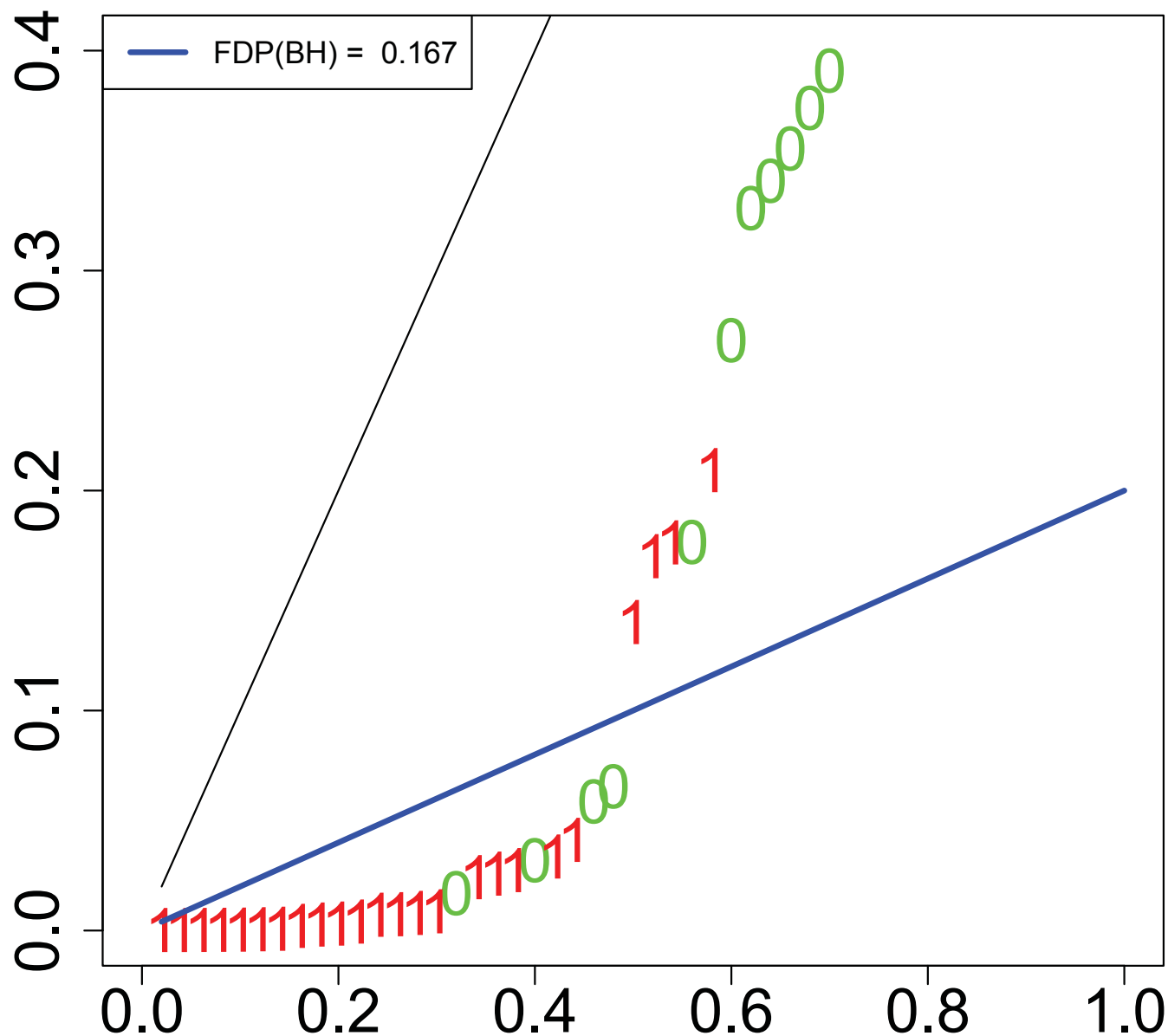
Simulations. $m = 50$; $m_0 = 25$; $\mu = 3$; $\Gamma = I_m$



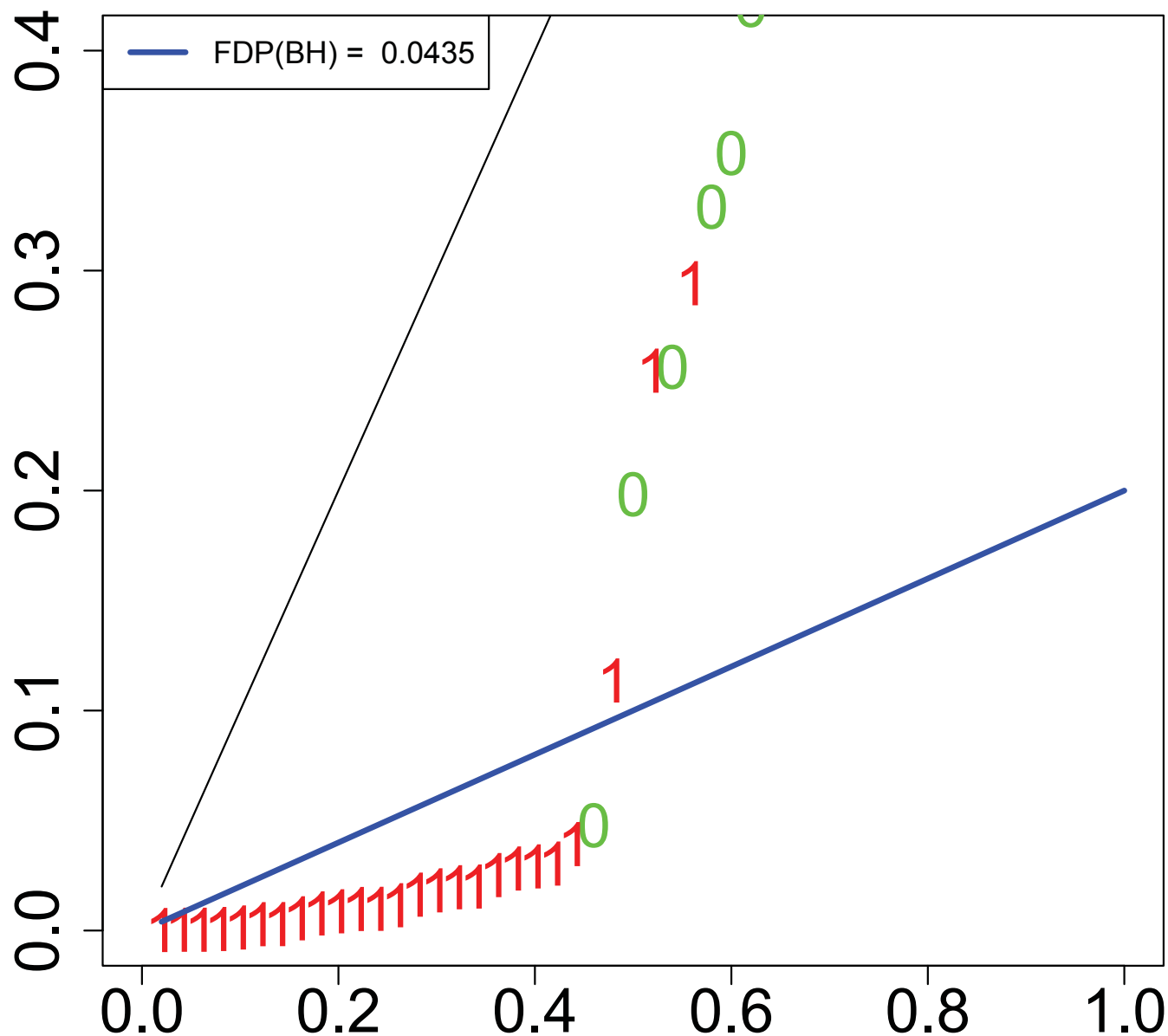
Simulations. $m = 50$; $m_0 = 25$; $\mu = 3$; $\Gamma = I_m$



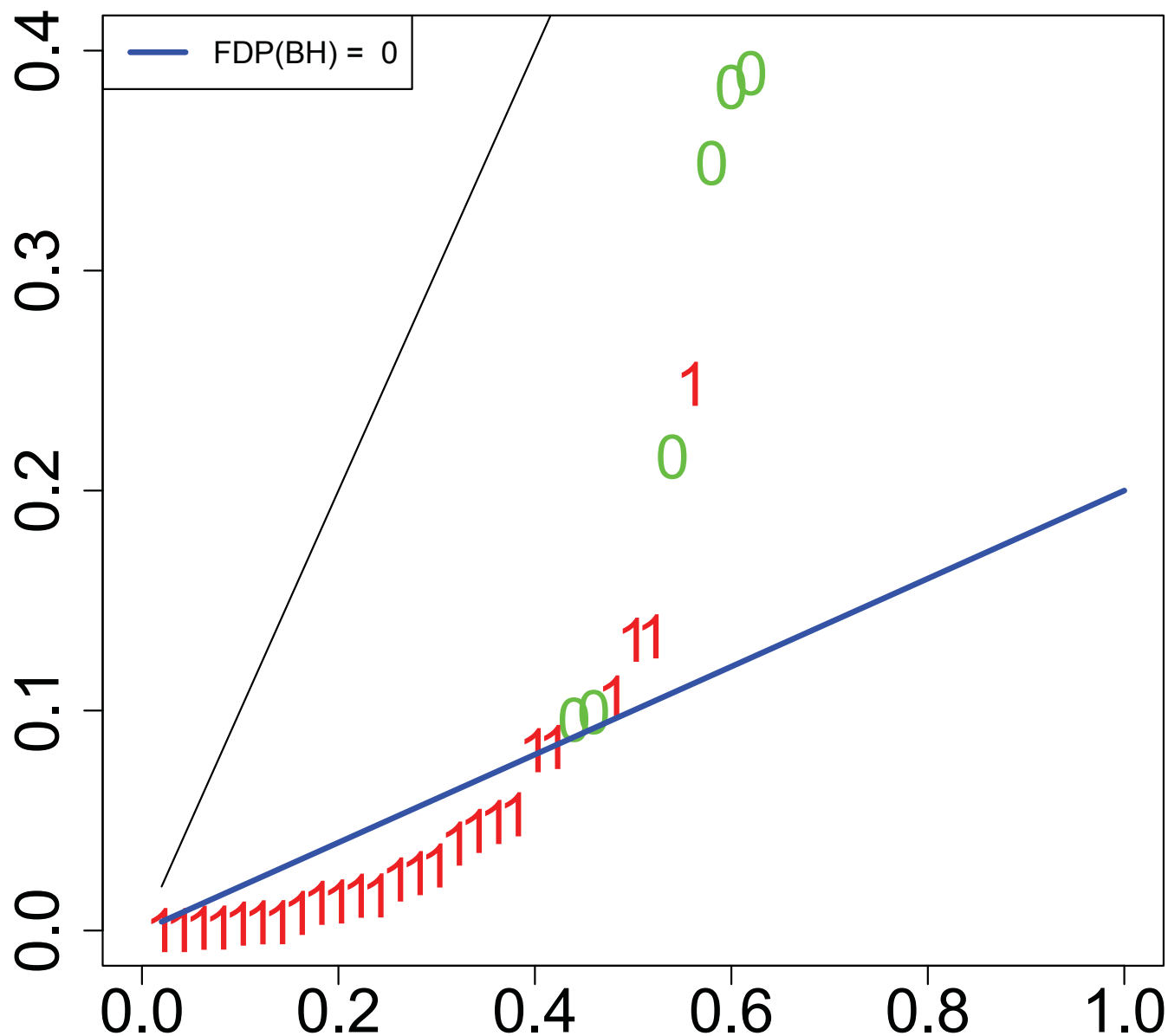
Simulations. $m = 50$; $m_0 = 25$; $\mu = 3$; $\Gamma = I_m$



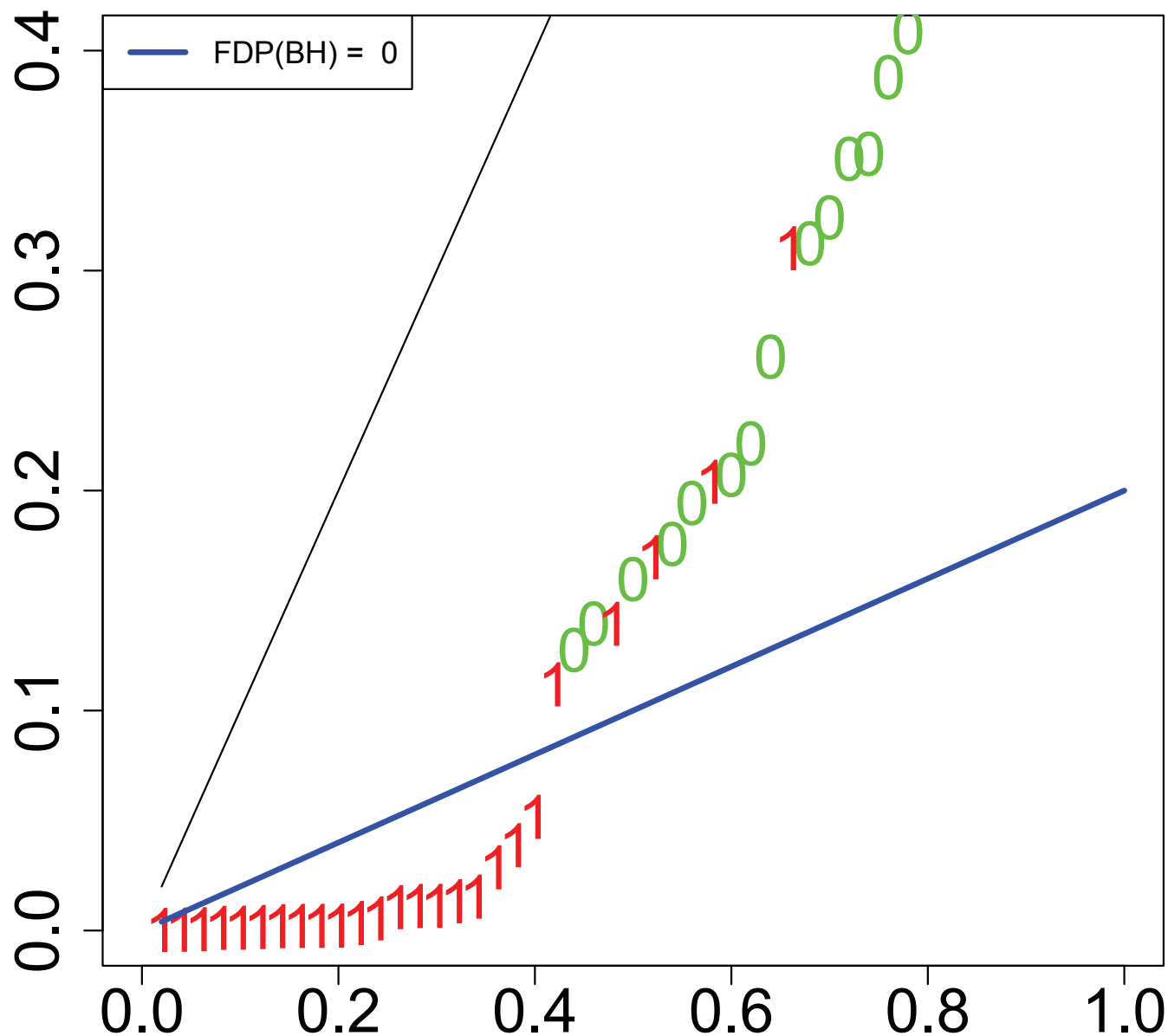
Simulations. $m = 50$; $m_0 = 25$; $\mu = 3$; $\Gamma = I_m$



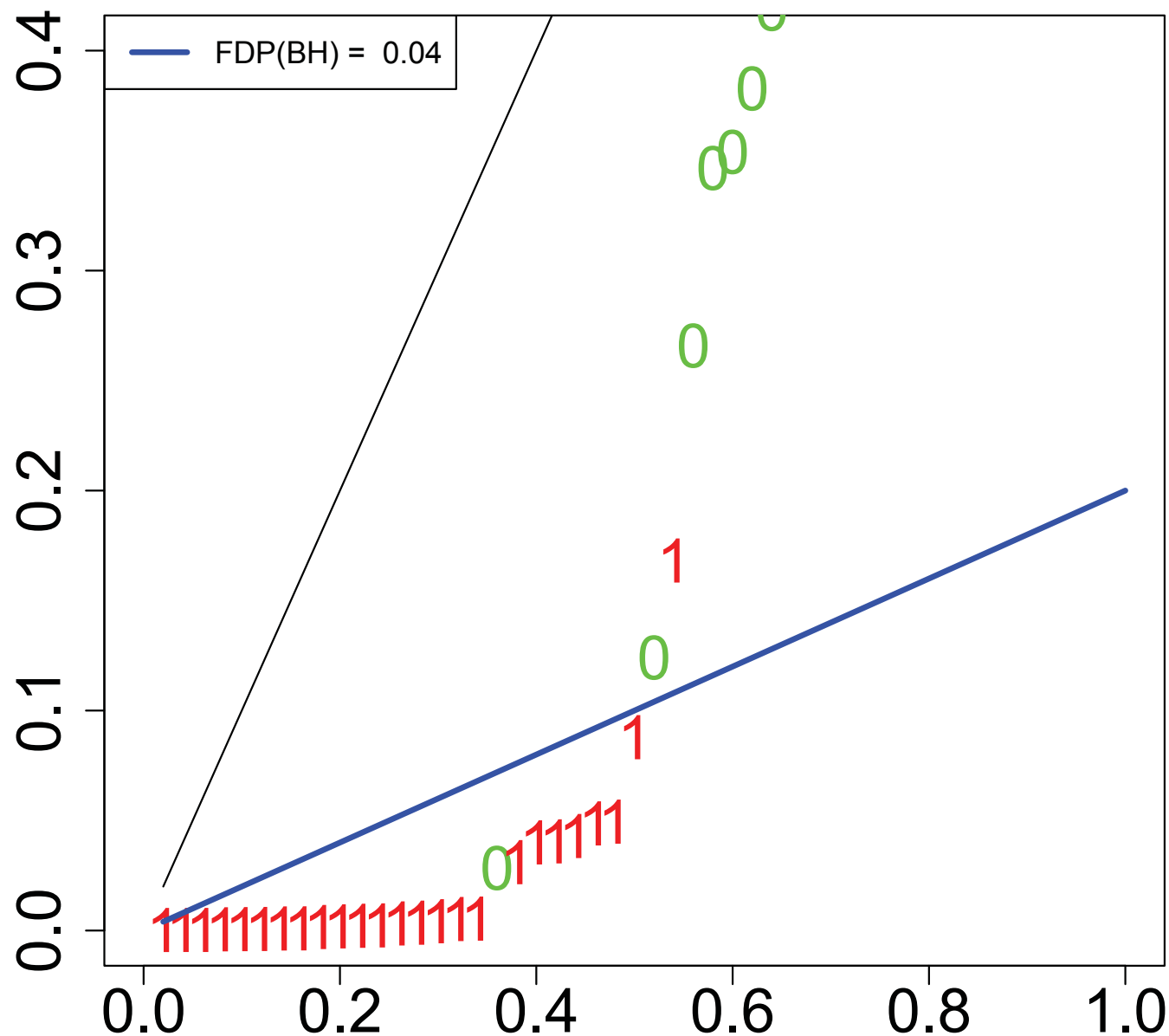
Simulations. $m = 50$; $m_0 = 25$; $\mu = 3$; $\Gamma = I_m$



Simulations. $m = 50$; $m_0 = 25$; $\mu = 3$; $\Gamma = I_m$

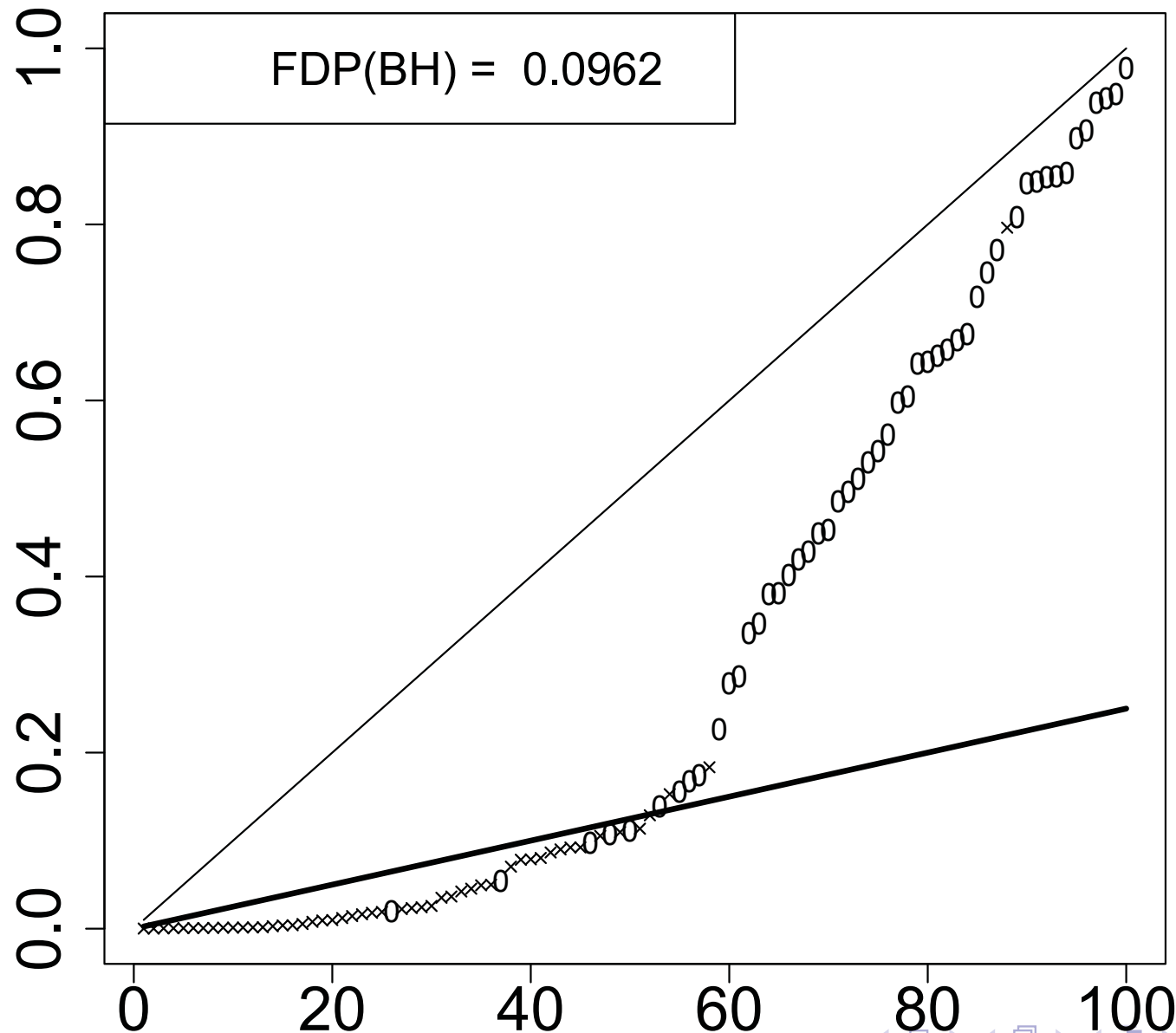


Simulations. $m = 50$; $m_0 = 25$; $\mu = 3$; $\Gamma = I_m$



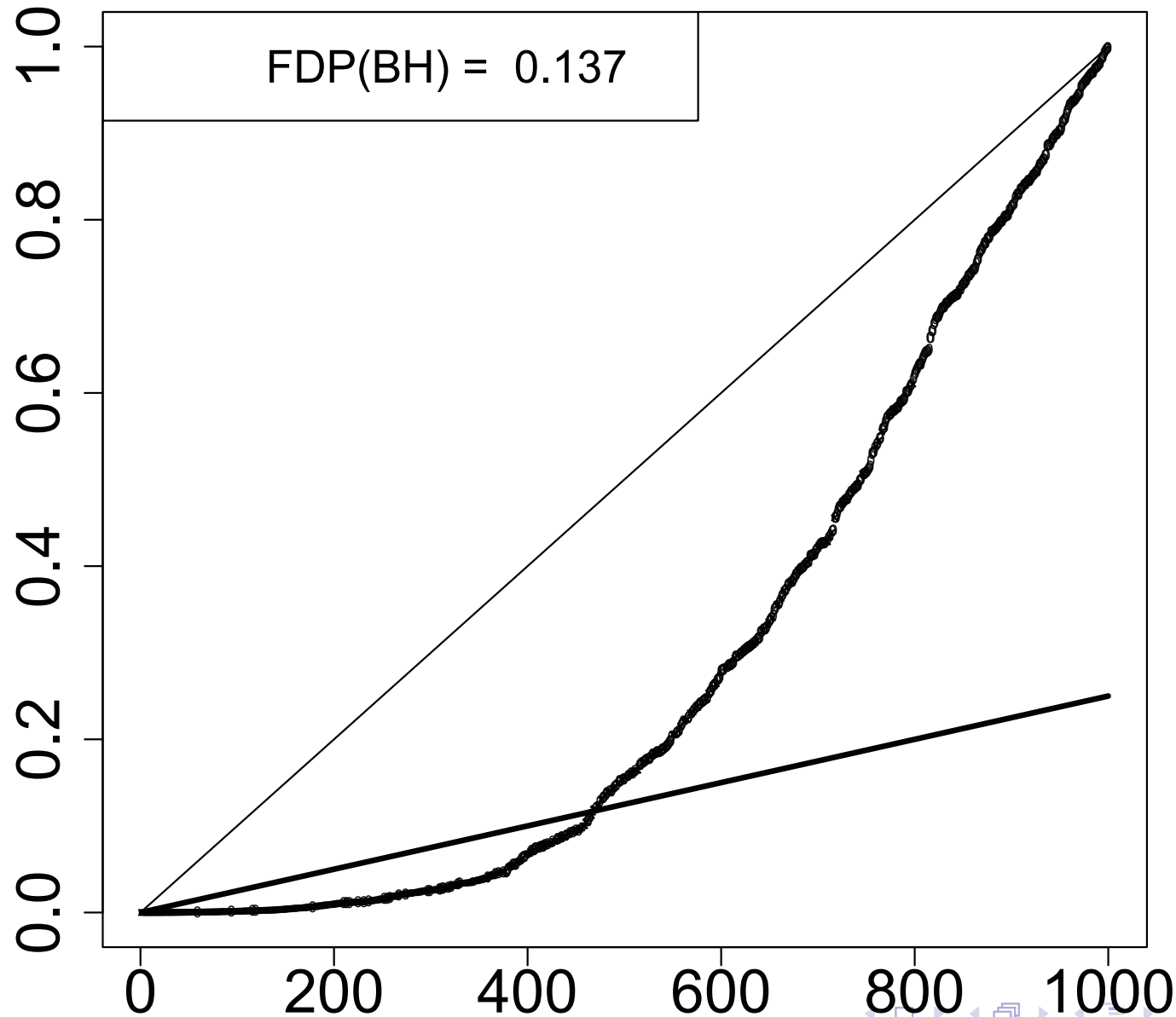
Scale invariance property $\pi_0 = 0.5; \mu = 2$

$m = 100$



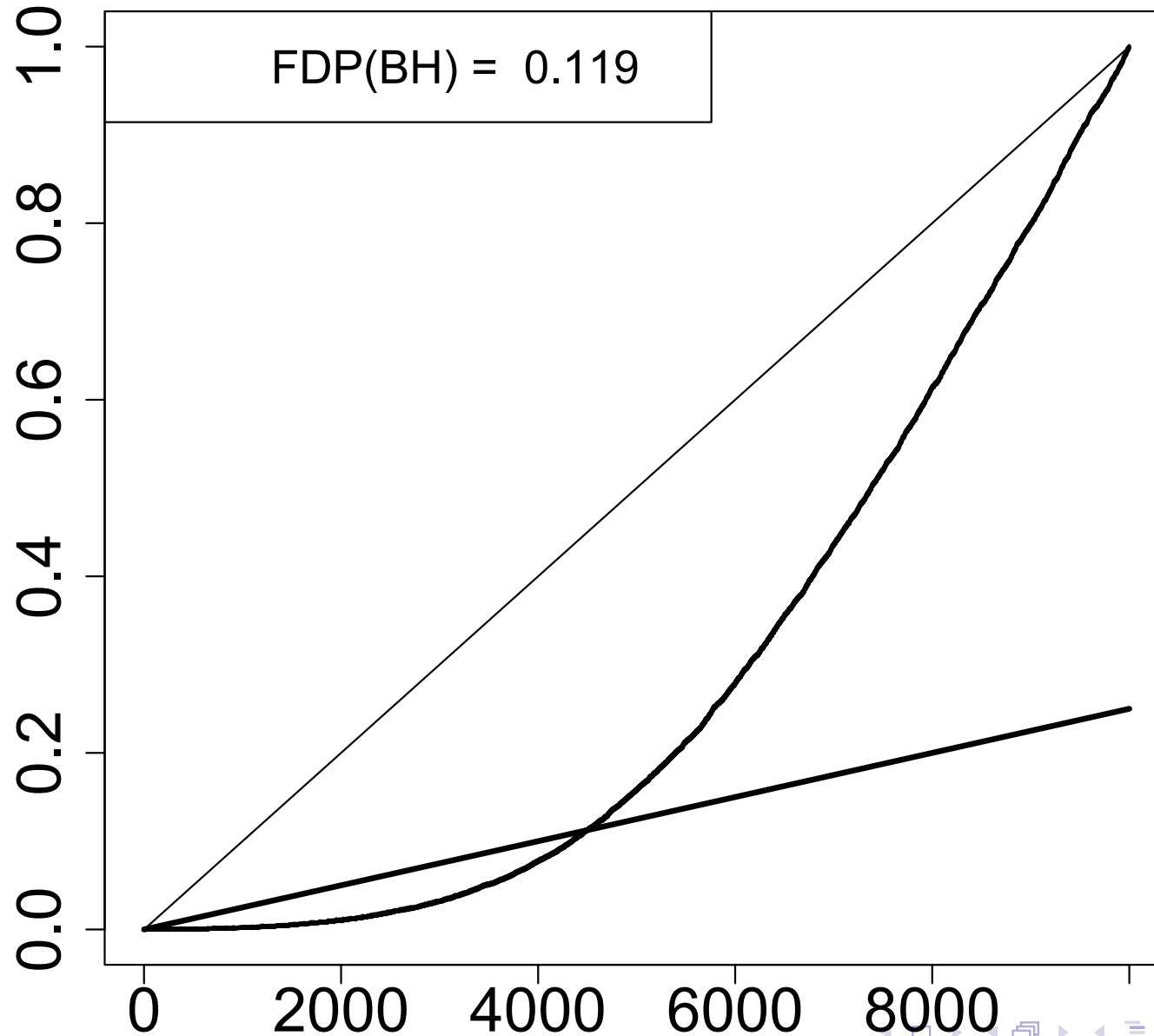
Scale invariance property $\pi_0 = 0.5; \mu = 2$

$m = 1000$



Scale invariance property $\pi_0 = 0.5; \mu = 2$

$m = 10000$





Benjamini and Hochberg (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing

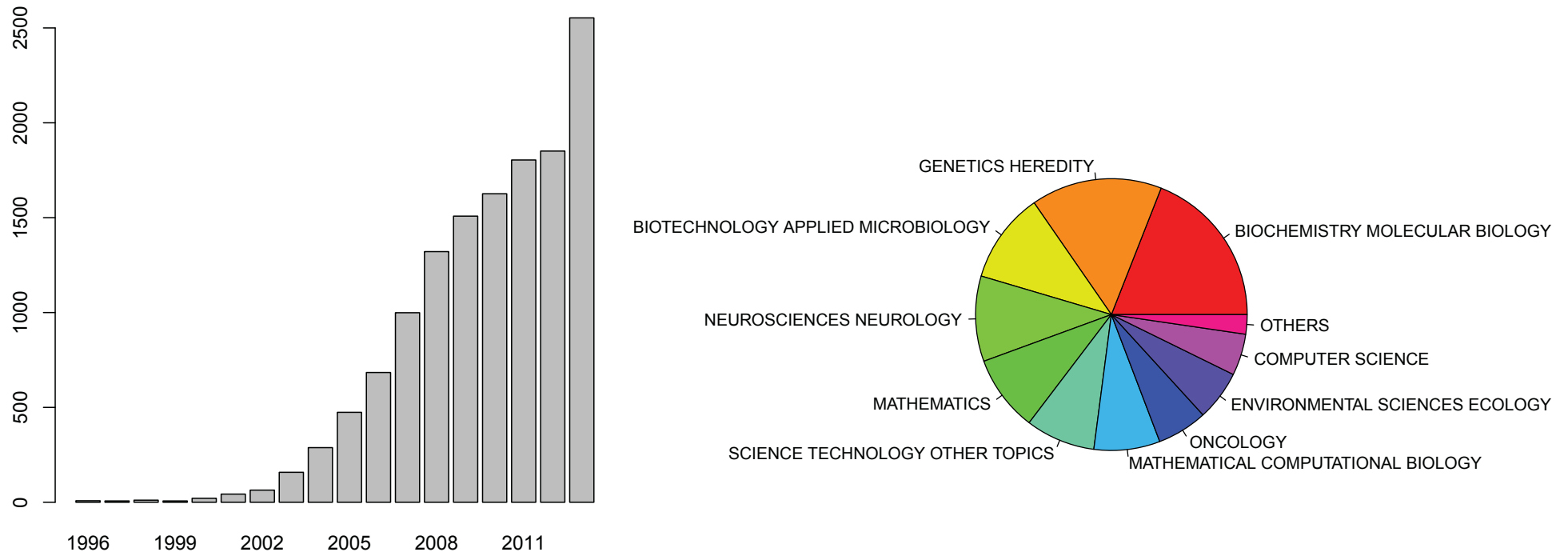


Figure: Statistics for the 13,427 papers citing this work from 1996 to 2013 and according to "the web of science". Left: per year. Right: per research fields.

- 1 Introduction
- 2 False discovery rate control
- 3 Dependence and limitations**
- 4 New challenges and results

FDR control under dependence

Theorem [Benjamini and Yekutieli (2001)]

If $\Gamma_{i,j} \geq 0$ for $i \neq j$ (PRDS) and \hat{t} threshold of BH procedure, $\forall \mu, H$,

$$\text{FDR}(\hat{t}) \leq (m_0/m)\alpha \leq \alpha$$

Well-accepted property [Farcomeni (2006)], [Kim and van de Wiel (2008)],...

$\text{FDR}(BH) \leq$ or $\simeq \alpha$ for any “realistic” dependencies

End of the story?

Consider ρ -equicorrelation: $\Gamma = \begin{pmatrix} 1 & \rho & \dots & \rho \\ \rho & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \rho \\ \rho & \dots & \rho & 1 \end{pmatrix}$ with $\rho \geq 0$

FDR control under dependence

Theorem [Benjamini and Yekutieli (2001)]

If $\Gamma_{i,j} \geq 0$ for $i \neq j$ (PRDS) and \hat{t} threshold of BH procedure, $\forall \mu, H$,

$$\text{FDR}(\hat{t}) \leq (m_0/m)\alpha \leq \alpha$$

Well-accepted property [Farcomeni (2006)], [Kim and van de Wiel (2008)],...

$\text{FDR}(BH) \leq$ or $\simeq \alpha$ for any “realistic” dependencies

End of the story?

Consider ρ -equicorrelation: $\Gamma = \begin{pmatrix} 1 & \rho & \dots & \rho \\ \rho & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \rho \\ \rho & \dots & \rho & 1 \end{pmatrix}$ with $\rho \geq 0$

FDR control under dependence

Theorem [Benjamini and Yekutieli (2001)]

If $\Gamma_{i,j} \geq 0$ for $i \neq j$ (PRDS) and \hat{t} threshold of BH procedure, $\forall \mu, H$,

$$\text{FDR}(\hat{t}) \leq (m_0/m)\alpha \leq \alpha$$

Well-accepted property [Farcomeni (2006)], [Kim and van de Wiel (2008)],...

$\text{FDR}(BH) \leq$ or $\simeq \alpha$ for any “realistic” dependencies

End of the story?

Consider ρ -equicorrelation: $\Gamma = \begin{pmatrix} 1 & \rho & \dots & \rho \\ \rho & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \rho \\ \rho & \dots & \rho & 1 \end{pmatrix}$ with $\rho \geq 0$

FDR control under dependence

Theorem [Benjamini and Yekutieli (2001)]

If $\Gamma_{i,j} \geq 0$ for $i \neq j$ (PRDS) and \hat{t} threshold of BH procedure, $\forall \mu, H$,

$$\text{FDR}(\hat{t}) \leq (m_0/m)\alpha \leq \alpha$$

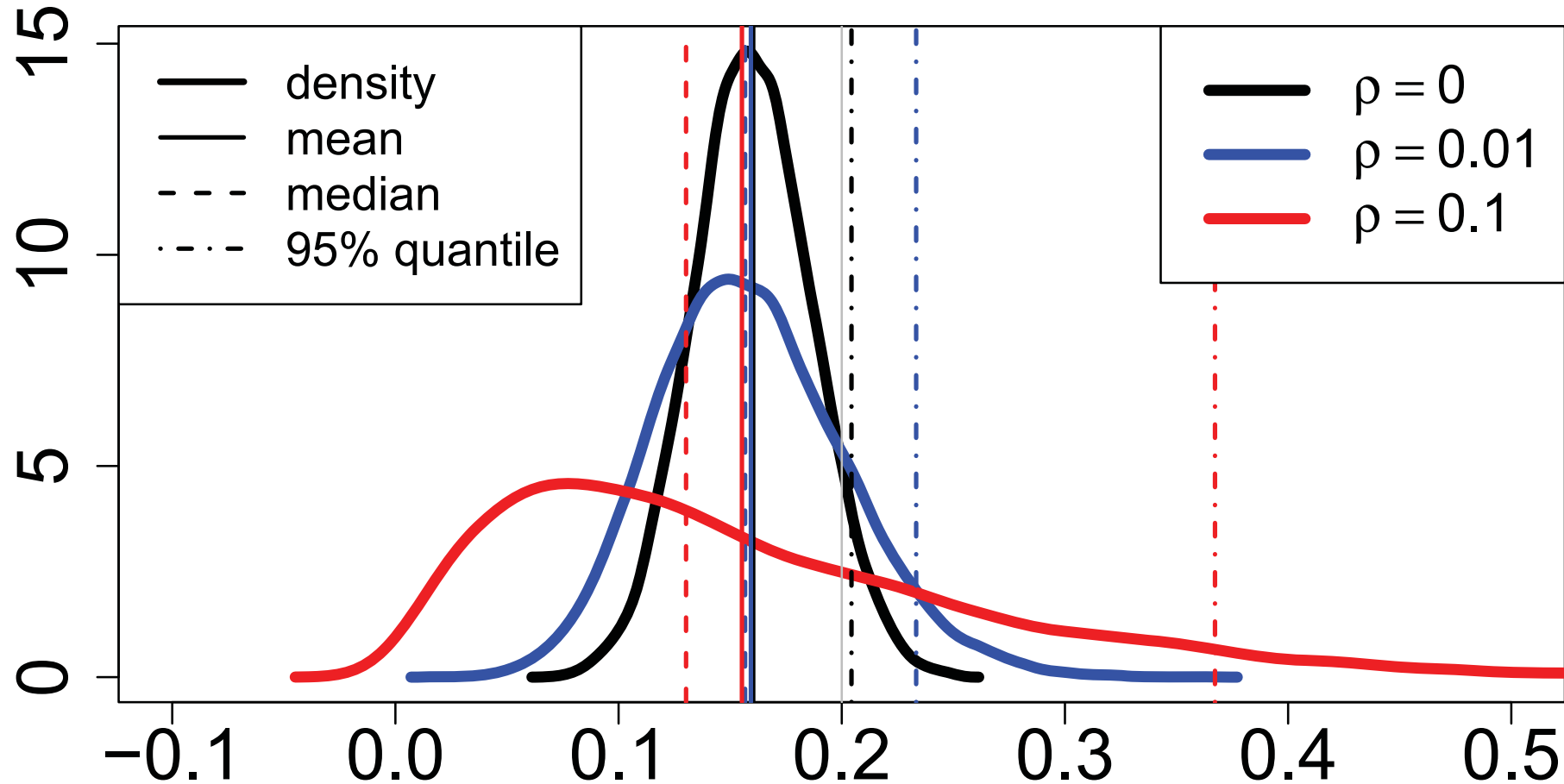
Well-accepted property [Farcomeni (2006)], [Kim and van de Wiel (2008)],...

$\text{FDR}(BH) \leq$ or $\simeq \alpha$ for any “realistic” dependencies

End of the story?

Consider ρ -equicorrelation: $\Gamma = \begin{pmatrix} 1 & \rho & \dots & \rho \\ \rho & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \rho \\ \rho & \dots & \rho & 1 \end{pmatrix}$ with $\rho \geq 0$

True underlying FDP(BH) $m = 1000, \pi_0 = 0.8, \alpha = 0.2$



True underlying FDP(BH) $m = 1000, \pi_0 = 0.8, \alpha = 0.2$



$FDR = \mathbf{E}[FDP] \leq \alpha$
does not provide
 $FDP \leq \alpha$ or even $FDP \simeq \alpha$!

- 1 Introduction
- 2 False discovery rate control
- 3 Dependence and limitations
- 4 New challenges and results**

New challenges under dependence

- **Task 1:** behavior of FDP(BH) under “weak” dependence

[Genovese Wasserman (2004)], [Farcomeni (2007,2008)], [Neuviel (2008)]

$$\text{FDP(BH)} = \Psi \left(\widehat{F}_0, \widehat{F}_1 \right)$$

- **Task 2:** control the FDP: for $\alpha, \zeta \in (0, 1)$, get $\widehat{t} = \widehat{t}(\Gamma)$ s.t.

$$\mathbf{P}(\text{FDP}(\widehat{t}) \leq \alpha) \geq 1 - \zeta.$$

[Lehmann and Romano (2005)], [Romano and Shaikh (2006a, 2006b)], [Romano and Wolf (2007)], [Guo and Romano (2007)], [Guo, He and Sarkar (2013)]

- **Task 3:** Build new test statistics $X_i^* \neq X_i$ that try to remove Γ .

[Friguet et al. (2009), Fan et al. (2012)]

Our starting point: ρ -equicorrelation

New challenges under dependence

- **Task 1:** behavior of FDP(BH) under “weak” dependence

[Genovese Wasserman (2004)], [Farcomeni (2007,2008)], [Neuviel (2008)]

$$\text{FDP(BH)} = \Psi \left(\widehat{F}_0, \widehat{F}_1 \right)$$

- **Task 2:** control the FDP: for $\alpha, \zeta \in (0, 1)$, get $\widehat{t} = \widehat{t}(\Gamma)$ s.t.

$$\mathbf{P}(\text{FDP}(\widehat{t}) \leq \alpha) \geq 1 - \zeta.$$

[Lehmann and Romano (2005)], [Romano and Shaikh (2006a, 2006b)], [Romano and Wolf (2007)], [Guo and Romano (2007)], [Guo, He and Sarkar (2013)]

- **Task 3:** Build new test statistics $X_i^* \neq X_i$ that try to remove Γ .

[Friguet et al. (2009), Fan et al. (2012)]

Our starting point: ρ -equicorrelation

New results under **weak** dependence

Theorem (weak equicorrelation)

$\rho = \rho_m \rightarrow 0$, $m_0/m \rightarrow \pi_0 \in (0, 1)$, $r_m = (m^{-1} + |\rho_m|)^{-1/2}$ then

$$r_m(\text{FDP}(BH) - \pi_0\alpha) \rightsquigarrow \mathcal{N}(0, \nu) ; \nu \in (0, \infty)$$

In particular,

$$\mathbf{P}(\text{FDP}(BH) \leq \alpha) \rightarrow 1 \geq 1 - \zeta$$

- ▶ BH works well under weak dep (asymp)
- ▶ convergence rate can be slow
- ▶ Relaxation to general weak dependence

$$\frac{1}{m(m-1)} \sum_{i \neq j} \Gamma_{i,j}^2 \rightarrow 0$$

New results under **weak** dependence

Theorem (weak equicorrelation)

$\rho = \rho_m \rightarrow 0$, $m_0/m \rightarrow \pi_0 \in (0, 1)$, $r_m = (m^{-1} + |\rho_m|)^{-1/2}$ then

$$r_m(\text{FDP}(BH) - \pi_0\alpha) \rightsquigarrow \mathcal{N}(0, \nu) ; \nu \in (0, \infty)$$

In particular,

$$\mathbf{P}(\text{FDP}(BH) \leq \alpha) \rightarrow 1 \geq 1 - \zeta$$

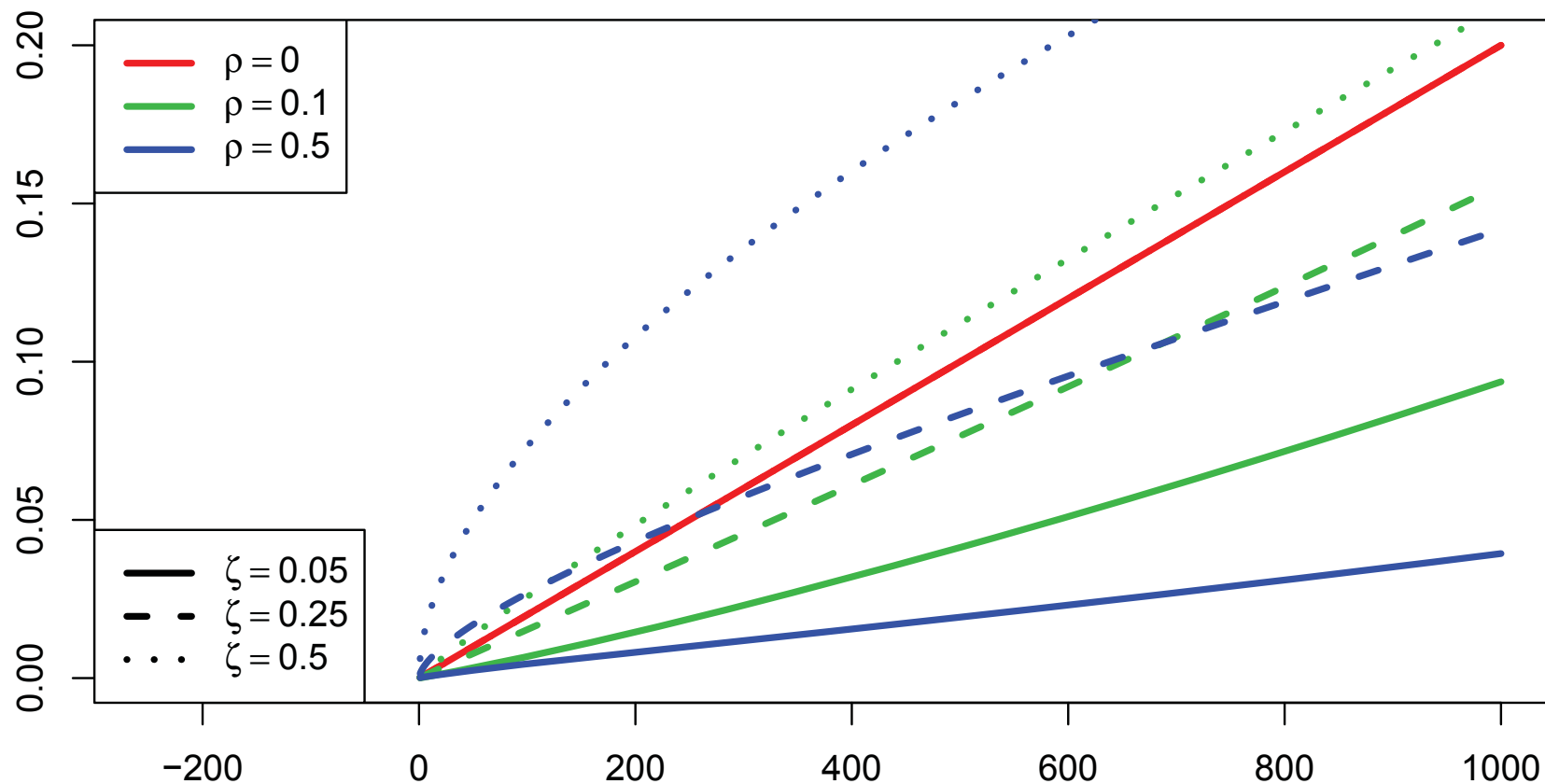
- ▶ BH works well under weak dep (asymp)
- ▶ convergence rate can be slow
- ▶ Relaxation to general weak dependence

$$\frac{1}{m(m-1)} \sum_{i \neq j} \Gamma_{i,j}^2 \rightarrow 0$$

New results under **strong** dependence

New BH type critical values under ρ -equicorrelation:

$$\tau_\ell = \Phi \left(\rho^{1/2} \Phi^{-1}(\zeta) + (1 - \rho)^{1/2} \Phi^{-1}(\alpha \ell / m) \right), \quad 1 \leq \ell \leq m.$$



New results under **strong** dependence

Theorem (strong equicorrelation)

$\rho \in (0, 1)$ (**fixed!**), $m_0/m \rightarrow \pi_0 \in (0, 1)$, new critical values

$$\tau_\ell = \Phi \left(\rho^{1/2} \Phi^{-1}(\zeta) + (1 - \rho)^{1/2} \Phi^{-1}(\alpha \ell / m) \right), \quad 1 \leq \ell \leq m.$$

Then

$$\liminf_m \mathbf{P}(\text{FDP}(\text{New}) \leq \alpha) \geq 1 - \zeta.$$

- ▶ Plug-in possible: $(\log m)(\widehat{\rho}_m - \rho_m)^2 = o_P(1)$
- ▶ Relaxation to general dependence (Romano-Wolf's heuristic)

New results under **strong** dependence

Theorem (strong equicorrelation)

$\rho \in (0, 1)$ (**fixed!**), $m_0/m \rightarrow \pi_0 \in (0, 1)$, new critical values

$$\tau_\ell = \Phi \left(\rho^{1/2} \Phi^{-1}(\zeta) + (1 - \rho)^{1/2} \Phi^{-1}(\alpha \ell / m) \right), \quad 1 \leq \ell \leq m.$$

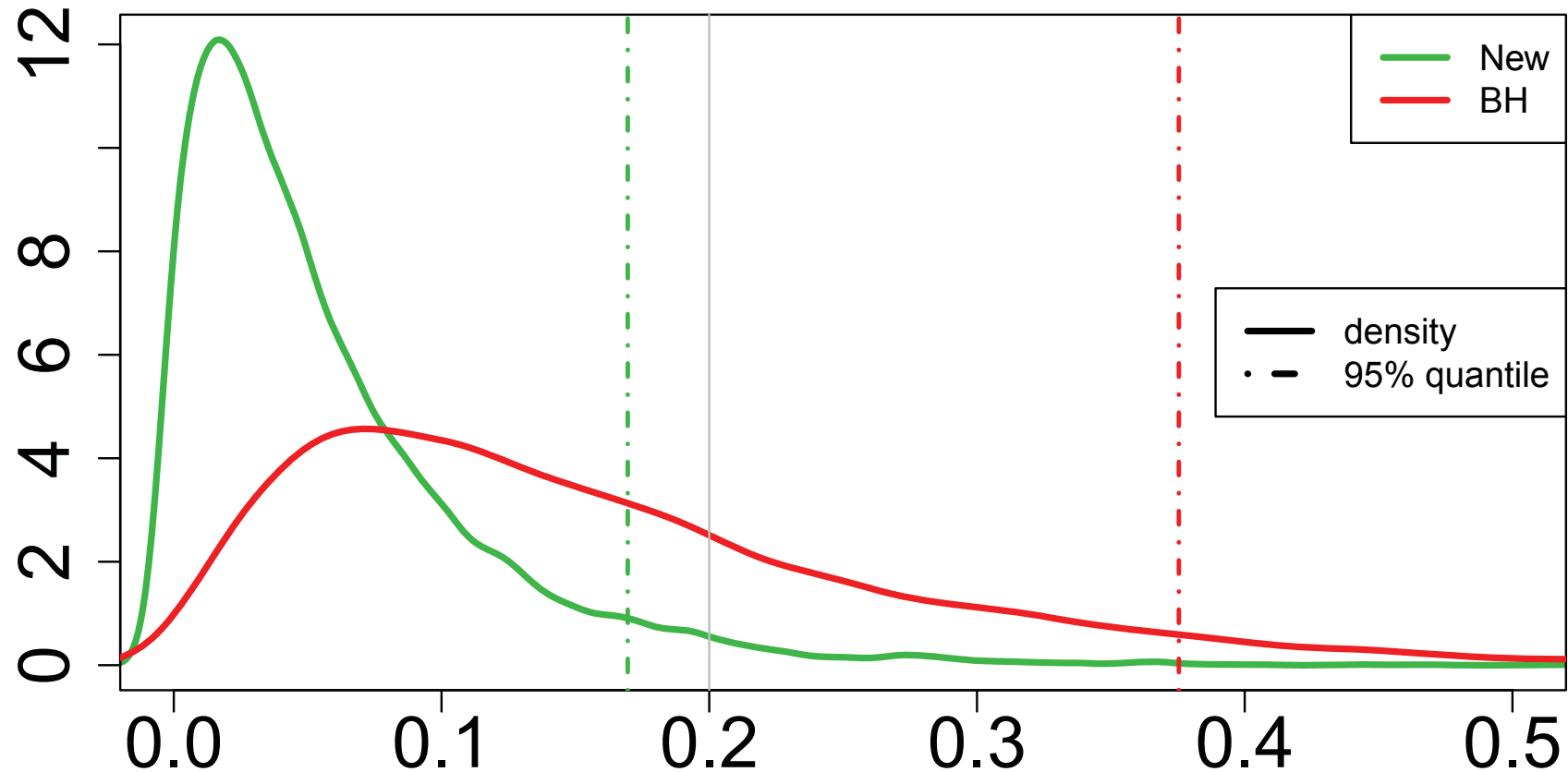
Then

$$\liminf_m \mathbf{P}(\text{FDP}(\text{New}) \leq \alpha) \geq 1 - \zeta.$$

- ▶ Plug-in possible: $(\log m)(\widehat{\rho}_m - \rho_m)^2 = o_P(1)$
- ▶ Relaxation to general dependence (Romano-Wolf's heuristic)

Distribution of FDP(New)

$m = 1000, \alpha = 0.2, \zeta = 0.05, \rho = 0.1$



Some take home messages

- ▶ BH procedure does not suffer from the curse of dimensionality
- ▶ “works” under weak dependence
- ▶ when equicorrelated strong: new critical values using ρ

Open problems

- ▶ FDP control in general, when Γ unknown (resampling)?
- ▶ when dependence strong and “well structured”: change test statistics?
- ▶ optimality issues

Outlook

Some take home messages

- ▶ BH procedure does not suffer from the curse of dimensionality
- ▶ “works” under weak dependence
- ▶ when equicorrelated strong: new critical values using ρ

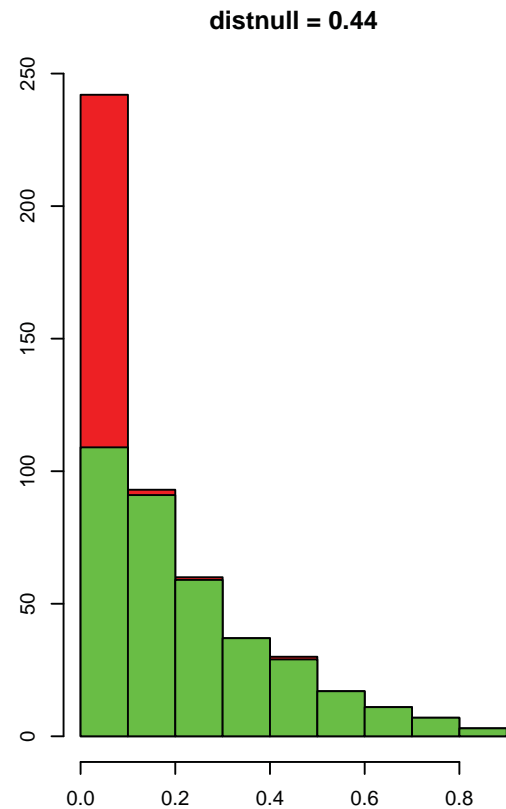
Open problems

- ▶ FDP control in general, when Γ unknown (resampling)?
- ▶ when dependence strong and “well structured”: change test statistics?
- ▶ optimality issues

Removing strong dependence

Equicorrelation

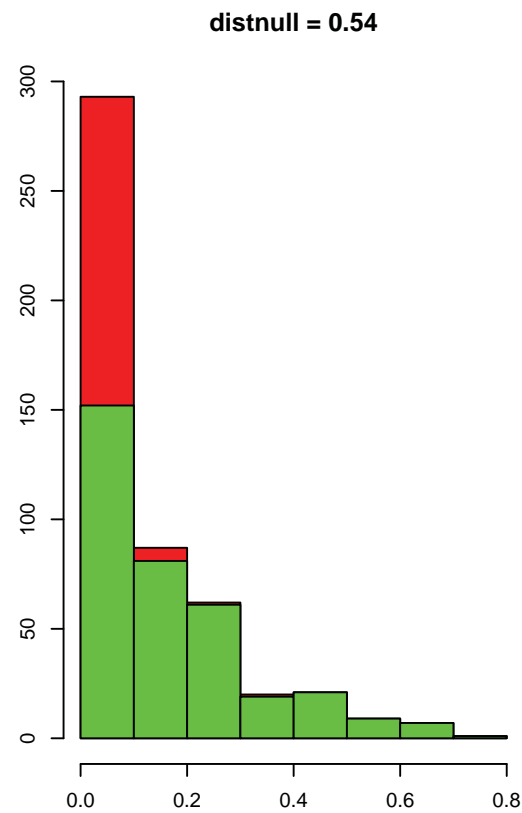
$$X_i = \mu H_i + U + \sqrt{1 - \rho} \xi_i,$$



Removing strong dependence

Equicorrelation

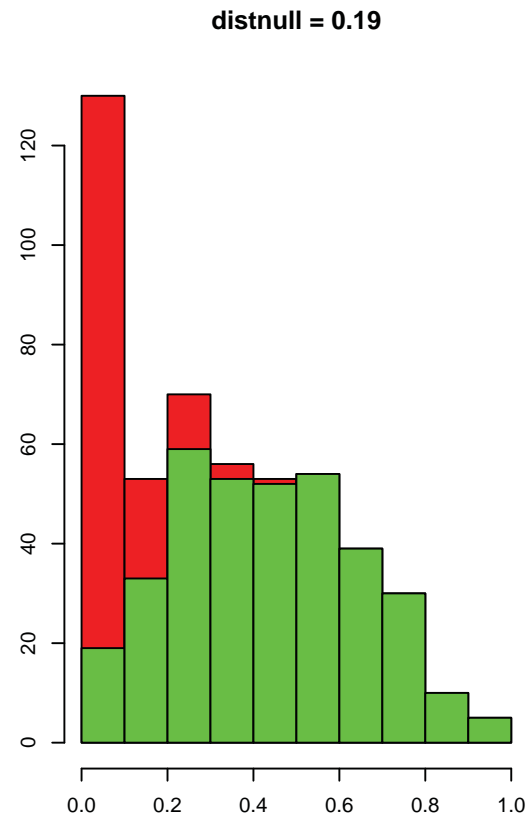
$$X_i = \mu H_i + U + \sqrt{1 - \rho} \xi_i,$$



Removing strong dependence

Equicorrelation

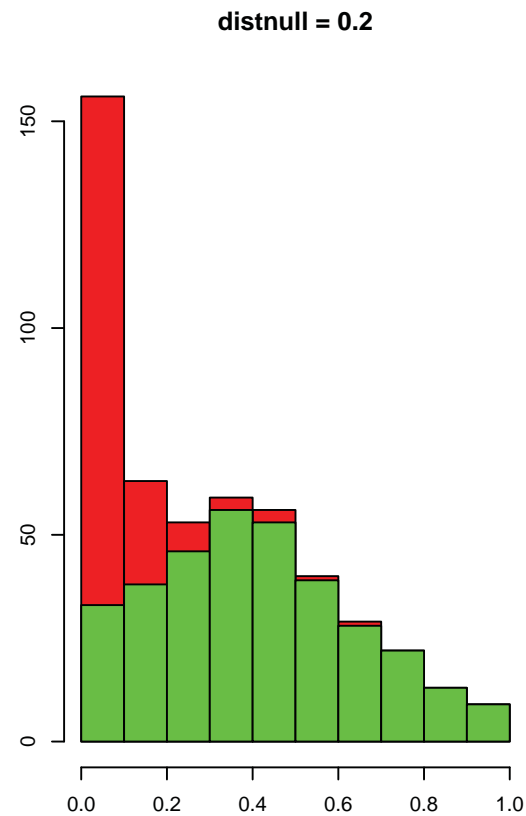
$$X_i = \mu H_i + U + \sqrt{1 - \rho} \xi_i,$$



Removing strong dependence

Equicorrelation

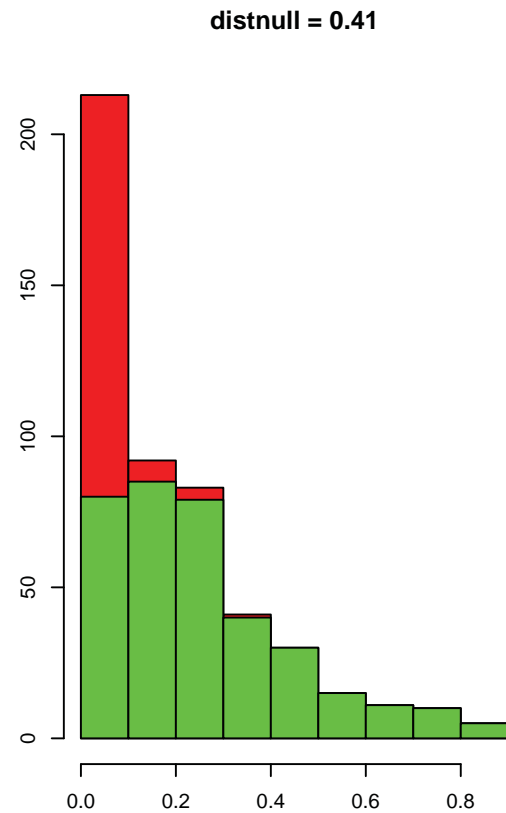
$$X_i = \mu H_i + U + \sqrt{1 - \rho} \xi_i,$$



Removing strong dependence

Equicorrelation

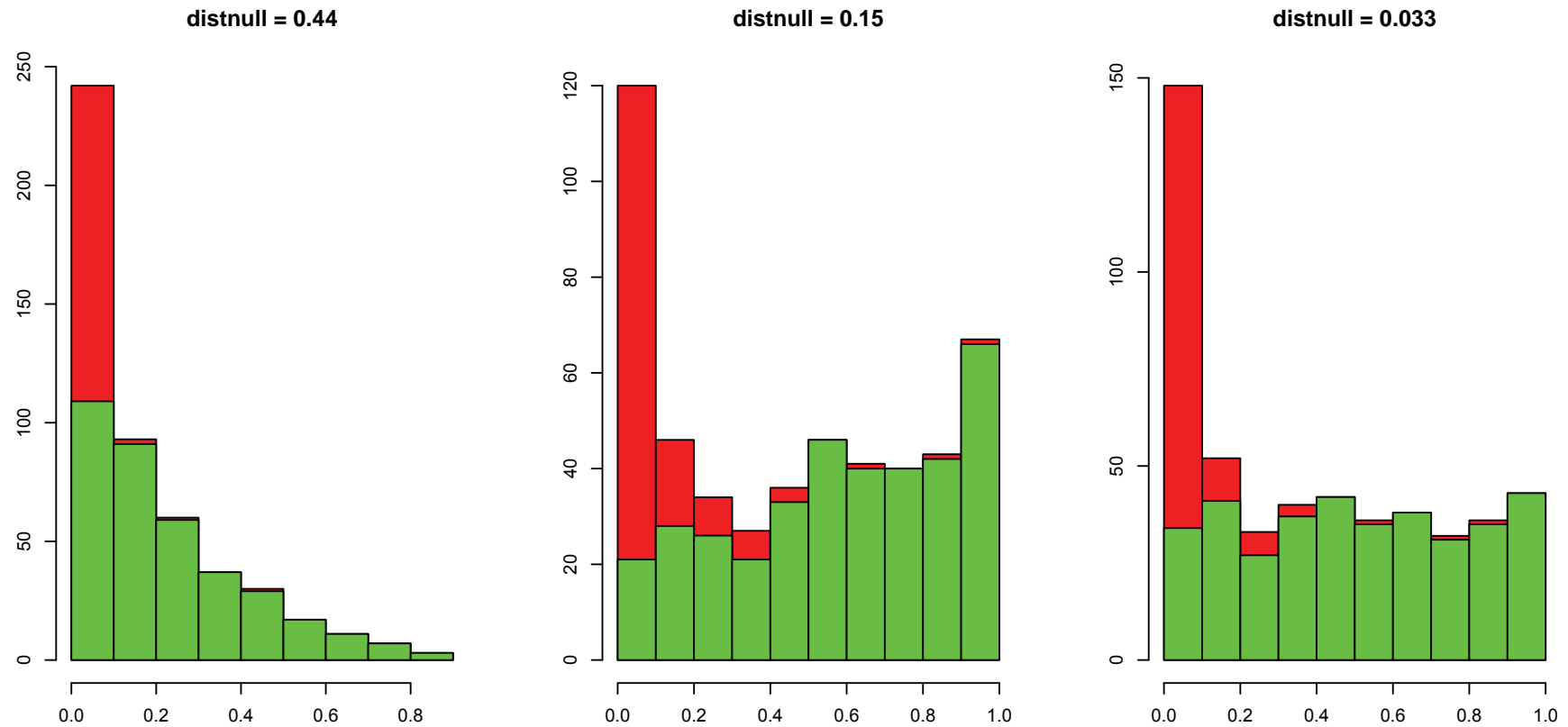
$$X_i = \mu H_i + U + \sqrt{1 - \rho} \xi_i,$$



Removing strong dependence

Equicorrelation

$$X_i = \mu H_i + U + \sqrt{1 - \rho} \xi_i, \quad X_i^* = X_i - \hat{U}$$



Original

$$\hat{U}_1 = X_{(0.475m)}$$

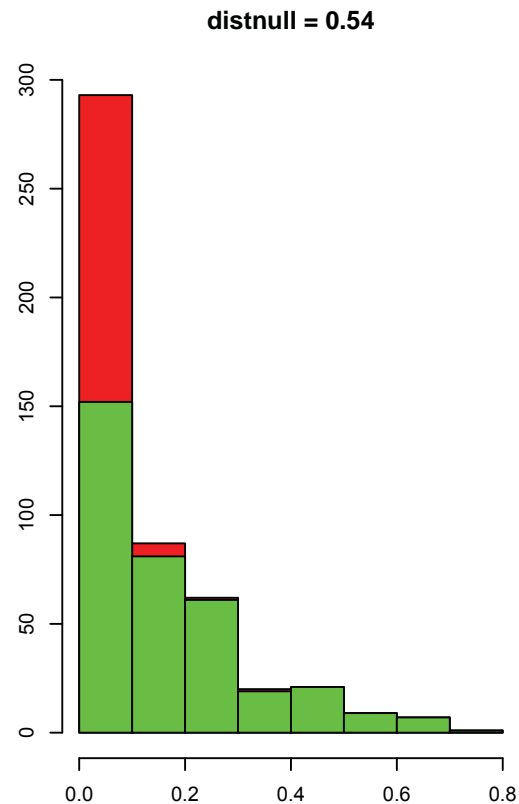
[Fan et al. (2012)]

$$\hat{U}_2 = \text{Mean}_{k \leq m^{1/2}}(X_{(k)} - e_k)$$

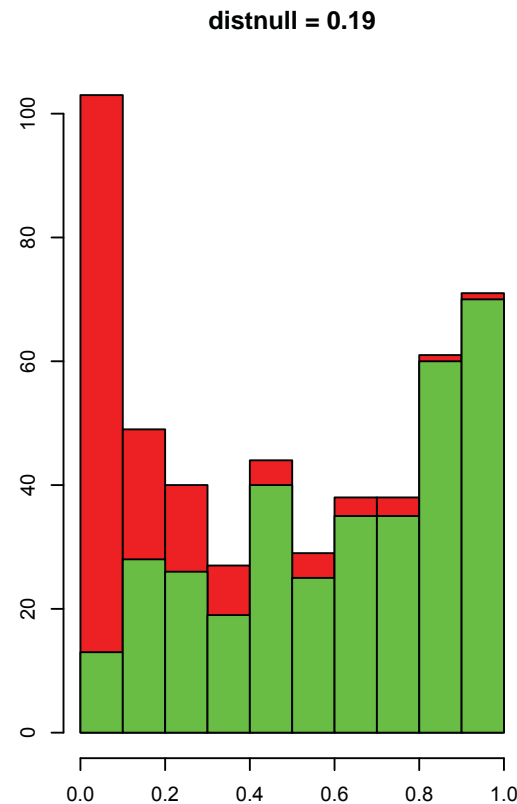
Removing strong dependence

Equicorrelation

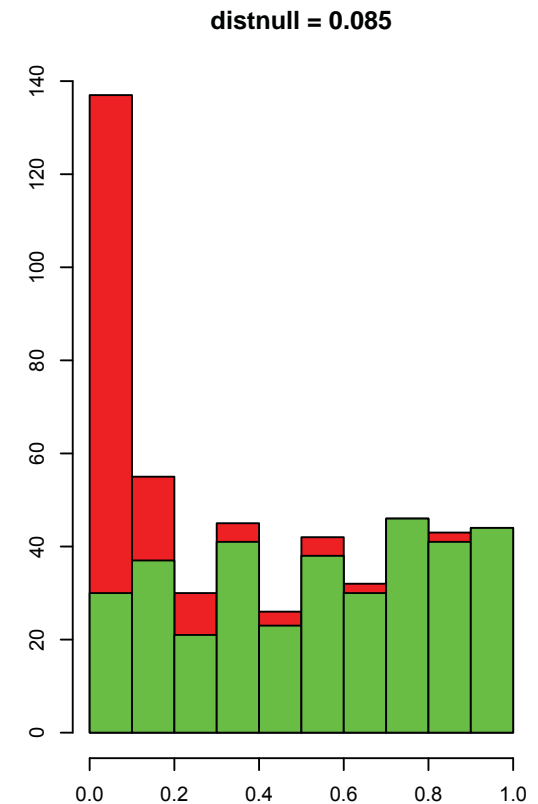
$$X_i = \mu H_i + U + \sqrt{1 - \rho} \xi_i, \quad X_i^* = X_i - \hat{U}$$



Original



$\hat{U}_1 = X_{(0.475m)}$
[Fan et al. (2012)]

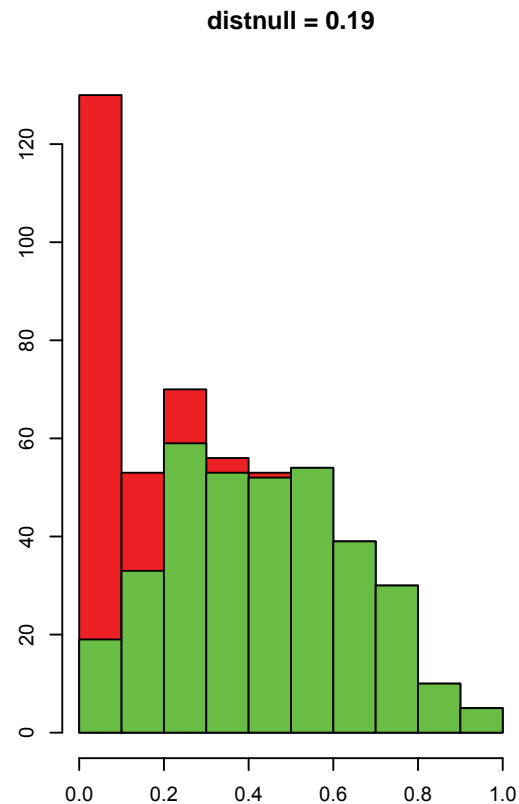


$\hat{U}_2 = \text{Mean}_{k \leq m^{1/2}}(X_{(k)} - e_k)$

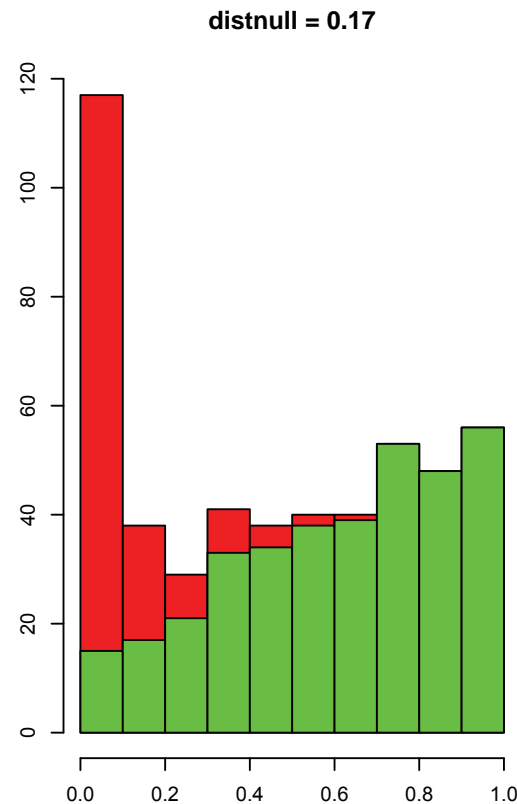
Removing strong dependence

Equicorrelation

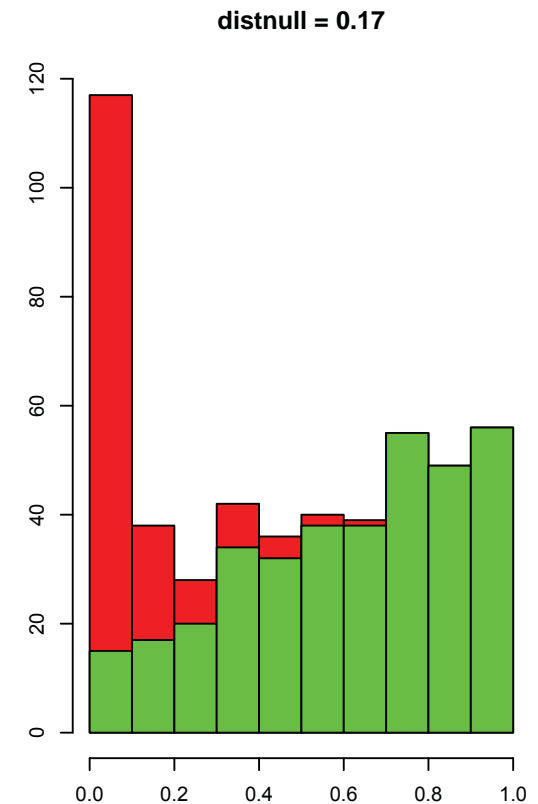
$$X_i = \mu H_i + U + \sqrt{1 - \rho} \xi_i, \quad X_i^* = X_i - \hat{U}$$



Original



$\hat{U}_1 = X_{(0.475m)}$
[Fan et al. (2012)]



$\hat{U}_2 = \text{Mean}_{k \leq m^{1/2}}(X_{(k)} - e_k)$

Removing strong dependence

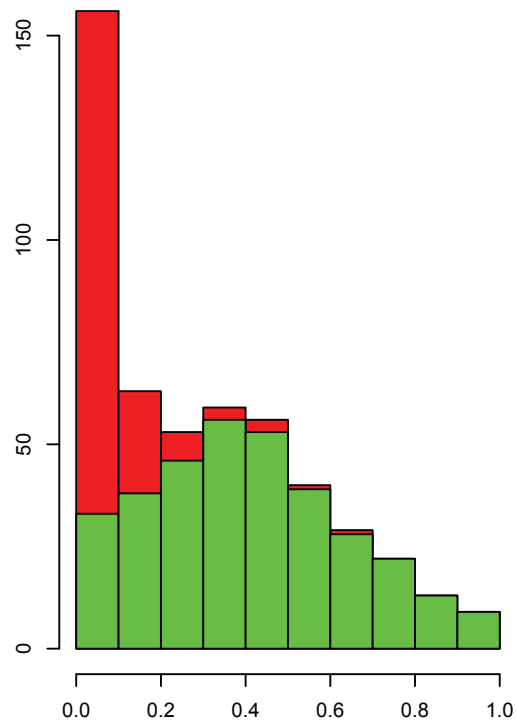
Equicorrelation

$$X_i = \mu H_i + U + \sqrt{1 - \rho} \xi_i, \quad X_i^* = X_i - \hat{U}$$

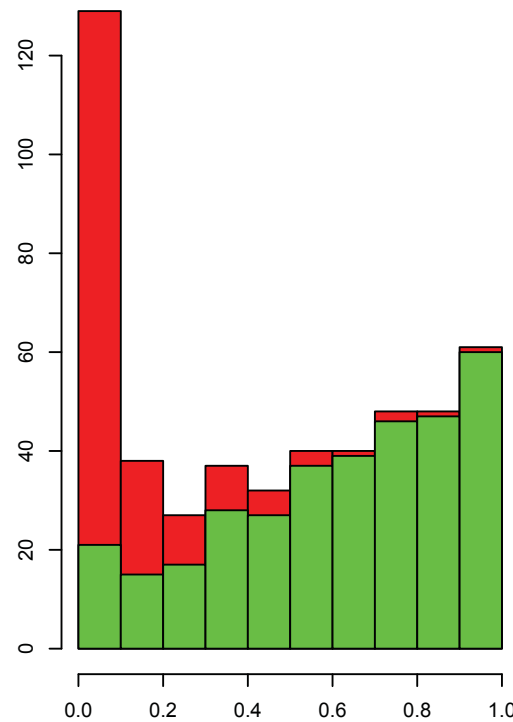
distnull = 0.2

distnull = 0.19

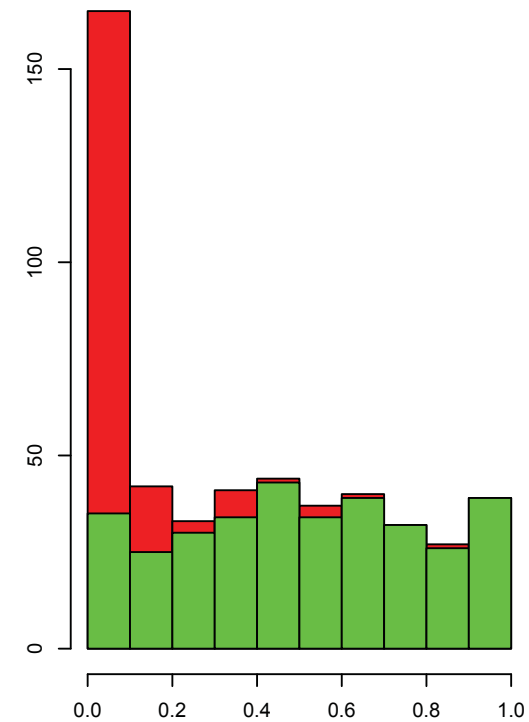
distnull = 0.036



Original



$\hat{U}_1 = X_{(0.475m)}$
[Fan et al. (2012)]



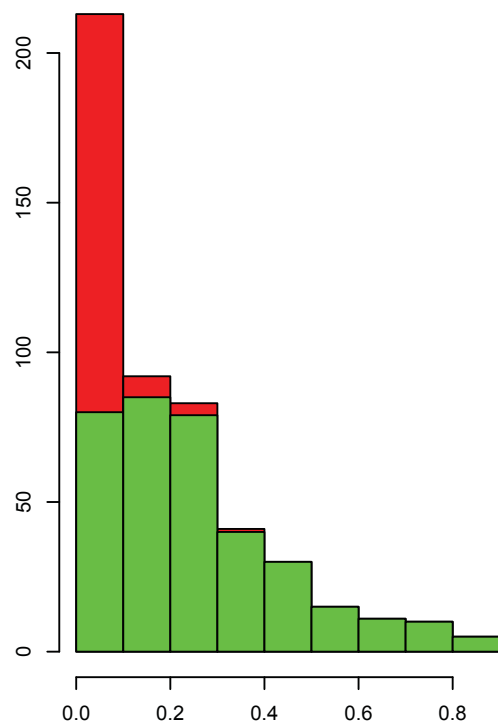
$\hat{U}_2 = \text{Mean}_{k \leq m^{1/2}}(X_{(k)} - e_k)$

Removing strong dependence

Equicorrelation

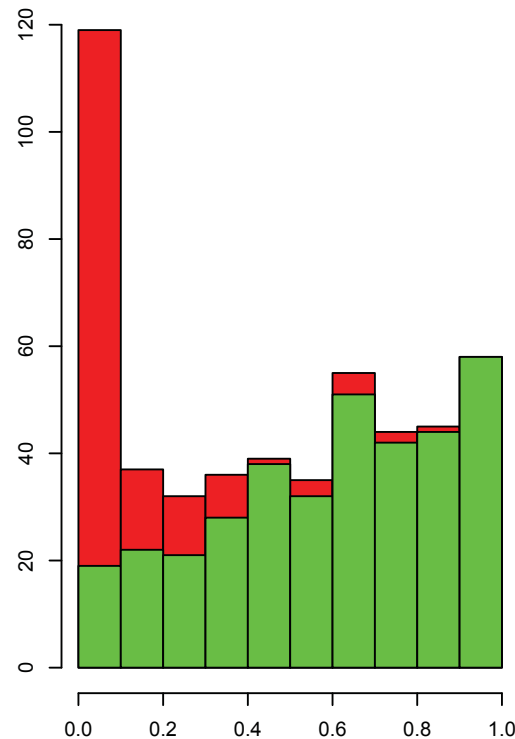
$$X_i = \mu H_i + U + \sqrt{1 - \rho} \xi_i, \quad X_i^* = X_i - \hat{U}$$

distnull = 0.41



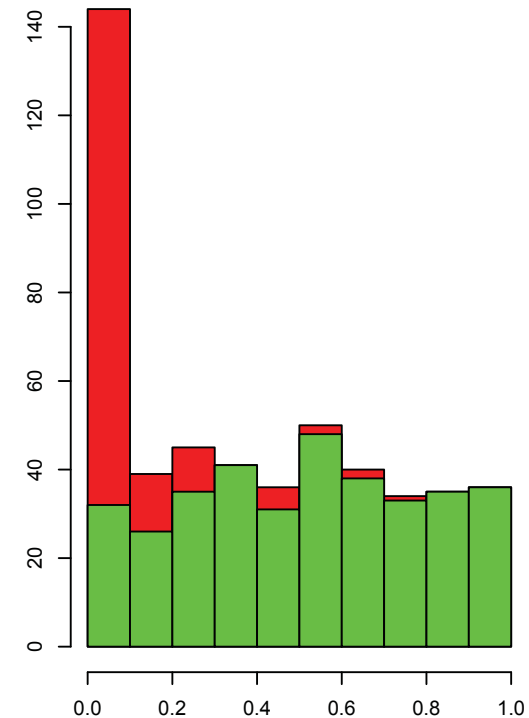
Original

distnull = 0.16



$\hat{U}_1 = X_{(0.475m)}$
[Fan et al. (2012)]

distnull = 0.046



$\hat{U}_2 = \text{Mean}_{k \leq m^{1/2}}(X_{(k)} - e_k)$