# FDR and Online FDR

## Adel Javanmard and Andrea Montanari

USC and Stanford

December 11, 2015

# Outline

# Large-scale Hypothesis Testing

# Assume

▶ I am the CTO of a big web company

▶ ≈ 1000 data scientists

▶ ≈ 1000 *'brilliant ideas'* per day
  ▶ Users are more likely to click on the first search result
  ▶ Users are more likely to on top right ads
  ▶ Users are more engaged with page layout A

▶ How to avoid wasting company resources?

Compute 'significance level' from data!

## Assume

▶ I am the CTO of a big web company

▶ $\approx$ 1000 data scientists

▶ $\approx$ 1000 *'brilliant ideas'* per day
  ▶ Users are more likely to click on the first search result
  ▶ Users are more likely to on top right ads
  ▶ Users are more engaged with page layout A

▶ How to avoid wasting company resources?

Compute 'significance level' from data!

# Assume

- I am the CTO of a big web company

- $\approx 1000$ data scientists

- $\approx 1000$ *'brilliant ideas'* per day
  - Users are more likely to click on the first search result
  - Users are more likely to on top right ads
  - Users are more engaged with page layout A

- How to avoid wasting company resources?

Compute 'significance level' from data!

# Assume

- I am the CTO of a big web company

- $\approx 1000$ data scientists

- $\approx 1000$ *'brilliant ideas'* per day
  - Users are more likely to click on the first search result
  - Users are more likely to on top right ads
  - Users are more engaged with page layout A

- How to avoid wasting company resources?

Compute 'significance level' from data!

# Assume

- I am the CTO of a big web company

- $\approx 1000$ data scientists

- $\approx 1000$ *'brilliant ideas'* per day
  - Users are more likely to click on the first search result
  - Users are more likely to on top right ads
  - Users are more engaged with page layout A

- How to avoid wasting company resources?

Compute 'significance level' from data!

# Assume

- I am the CTO of a big web company

- $\approx 1000$ data scientists

- $\approx 1000$ *'brilliant ideas'* per day
  - Users are more likely to click on the first search result
  - Users are more likely to on top right ads
  - Users are more engaged with page layout A

- How to avoid wasting company resources?

Compute 'significance level' from data!

# Example

**Idea:** *Users click more on the first search result than on the second*

Null $H_0$: Users are equaly likely to click on first and second

Data:

- $n$ events
- $n_1$ clicks on the *first* result
- $n_2 = n - n_1$ clicks on the *second* result

Idea

$$H_0 \quad \Rightarrow \quad z \equiv \frac{n_1 - n_2}{\sqrt{n}} \approx N(0, 1)$$

- If $z \gg 1$, then declare it significant

# Example

**Idea:** *Users click more on the first search result than on the second*

**Null $H_0$:** Users are equaly likely to click on first and second

Data:

- $n$ events
- $n_1$ clicks on the *first* result
- $n_2 = n - n_1$ clicks on the *second* result

Idea

$$H_0 \quad \Rightarrow \quad z \equiv \frac{n_1 - n_2}{\sqrt{n}} \approx \mathrm{N}(0, 1)$$

- If $z \gg 1$, then declare it significant

# Example

**Idea:** *Users click more on the first search result than on the second*

**Null $H_0$:** Users are equaly likely to click on first and second

**Data:**

- $n$ events
- $n_1$ clicks on the *first* result
- $n_2 = n - n_1$ clicks on the *second* result

Idea

$$H_0 \quad \Rightarrow \quad z \equiv \frac{n_1 - n_2}{\sqrt{n}} \approx \mathrm{N}(0, 1)$$

- If $z \gg 1$, then declare it significant

# Example

**Idea:** *Users click more on the first search result than on the second*

**Null $H_0$:** Users are equaly likely to click on first and second

**Data:**

- $n$ events
- $n_1$ clicks on the *first* result
- $n_2 = n - n_1$ clicks on the *second* result

**Idea**

$$H_0 \quad \Rightarrow \quad z \equiv \frac{n_1 - n_2}{\sqrt{n}} \approx \mathrm{N}(0, 1)$$

- If $z \gg 1$, then declare it significant

# Formally

$$z \equiv \frac{n_1 - n_2}{\sqrt{n}} \approx \mathrm{N}(0, 1)$$

**p-value** $(G \sim \mathrm{N}(0, 1))$

$$p \equiv \mathbb{P}(G \geq z) = \int_z^\infty \frac{e^{-x^2/2}}{\sqrt{2\pi}} \, \mathrm{d}x$$

- Null:      $p \sim \mathrm{Uniform}([0, 1])$ (Definition)
- Small $p$:   significant

# Formally

$$z \equiv \frac{n_1 - n_2}{\sqrt{n}} \approx \mathrm{N}(0, 1)$$

**p-value** $(G \sim \mathrm{N}(0, 1))$

$$p \equiv \mathbb{P}(G \geq z) = \int_z^\infty \frac{e^{-x^2/2}}{\sqrt{2\pi}} \, \mathrm{d}x$$

- Null:  $p \sim \mathrm{Uniform}([0, 1])$  (Definition)
- Small $p$:  significant

# Company policy

Bring your idea up only if $p \leq \alpha$

[$\alpha = 0.05$, Fisher's rule of thumb]

# Company policy

Bring your idea up only if $p \leq \alpha$

[$\alpha = 0.05$, Fisher's rule of thumb]

# Problem

- $M \approx 1000$ hypotheses per day
- $M\alpha \approx 1000 \cdot 0.05 = 50$ pass the test
- Still too much waste

New company policy (Bonferroni):

Bring up your idea only if $p \le \alpha_M = \alpha/M$

# Problem

- $M \approx 1000$ hypotheses per day
- $M\alpha \approx 1000 \cdot 0.05 = 50$ pass the test
- Still too much waste

New company policy (Bonferroni):

Bring up your idea only if $p \leq \alpha_M = \alpha/M$

# Problem

- $M \approx 1000$ hypotheses per day
- $M\alpha \approx 1000 \cdot 0.05 = 50$ pass the test
- Still too much waste

**New company policy (Bonferroni):**

Bring up your idea only if $p \leq \alpha_M = \alpha/M$

# Problem with Bonferroni

Bring up your idea only if $p \leq \alpha_M = \alpha/M$

- More data scientists $\Rightarrow$ Less sensitive
- $\alpha$ false positives per day $\Rightarrow$ Does not scale with $M$

# Problem with Bonferroni

Bring up your idea only if $p \leq \alpha_M = \alpha/M$

- More data scientists $\Rightarrow$ Less sensitive
- $\alpha$ false positives per day $\Rightarrow$ Does not scale with $M$

What do we want to achieve?

# FDR (Benjamini, Hochberg, 1995)

- $M$ hypotheses
- $D \equiv$ Total number of discoveries (positives)
- $FD \equiv$ Number of false discoveries

$$FDR = \mathbb{E}\Big\{ \frac{FD}{\max(D, 1)} \Big\}$$

**Interpretation:** $FDR \leq 0.1 \Rightarrow$ At most 10% of the discoveries is false.

# FDR (Benjamini, Hochberg, 1995)

- $M$ hypotheses
- D ≡ Total number of discoveries (positives)
- FD ≡ Number of false discoveries

$$\text{FDR} = \mathbb{E}\left\{ \frac{\text{FD}}{\max(\text{D}, 1)} \right\}$$

**Interpretation:** FDR $\leq 0.1 \Rightarrow$ At most 10% of the discoveries is false.

# Controlling FDR

# Setting

**Null hypotheses:**

$$H_{0,1}, H_{0,2}, \ldots, H_{0,M}$$

**p-values:**

$$p_1, p_2, \ldots, p_M$$

**Ground truth:**

$$\theta_1, \theta_2, \ldots, \theta_M \, [H_{0,i} : \theta_i = 0]$$

**Test ouput $(\boldsymbol{p} = (p_i)_{1 \leq i \leq M}$:**

$$T_1(\boldsymbol{p}), T_2(\boldsymbol{p}), \ldots, T_M(\boldsymbol{p}) \in \{0, 1\}$$

$$\theta_i = 0 \Rightarrow p_i \sim \mathrm{Uniform}([0, 1])$$

# Setting

**Null hypotheses:**

$$H_{0,1}, H_{0,2}, \ldots, H_{0,M}$$

**p-values:**

$$p_1, p_2, \ldots, p_M$$

**Ground truth:**

$$\theta_1, \theta_2, \ldots, \theta_M \, [H_{0,i} : \theta_i = 0]$$

**Test ouput $(\boldsymbol{p} = (p_i)_{1 \le i \le M}$:**

$$T_1(\boldsymbol{p}), T_2(\boldsymbol{p}), \ldots, T_M(\boldsymbol{p}) \in \{0, 1\}$$

$$\theta_i = 0 \Rightarrow p_i \sim \mathrm{Uniform}([0, 1])$$

# Benjamini-Hochberg procedure

▶ Order the p-values

$$p_{(1)} \leq p_{(2)} \leq \cdots \leq p_{(M)}$$

▶ Set threshold

$$I = \max \left\{ i \in [M] : \ p_{(i)} \leq \frac{i\alpha}{M} \right\}$$

▶ Reject at level $p_{(I)}$:

$$T_\ell(p) = \begin{cases} 1 & \text{if } p_\ell \leq p_{(I)}, \\ 0 & \text{otherwise.} \end{cases}$$

Theorem (Benjamini, Hochberg, 1995)

*If the p-values are independent, and BH is used, then*

$$\mathrm{FDR} \leq \alpha$$

# Benjamini-Hochberg procedure

- Order the p-values

$$p_{(1)} \leq p_{(2)} \leq \cdots \leq p_{(M)}$$

- Set threshold

$$I = \max \left\{ i \in [M] : \ p_{(i)} \leq \frac{i\alpha}{M} \right\}$$

- Reject at level $p_{(I)}$:

$$T_\ell(p) = \begin{cases} 1 & \text{if } p_\ell \leq p_{(I)}, \\ 0 & \text{otherwise.} \end{cases}$$

---

Theorem (Benjamini, Hochberg, 1995)

*If the p-values are independent, and BH is used, then*

$$\text{FDR} \leq \alpha$$

# Interpretation

- $M_0$ true nulls, $M_1 = M - M_0$ true non-null
- Reject $H_{0,i}$ if $p_i \leq q$

$$\text{FD} \approx M_0 q$$
$$\text{D} = J(q) \equiv \max\{i : p_{(i)} < q\}$$

$$\text{FDR} \approx \widehat{\text{FDR}}(q) \equiv \frac{M_0 q}{J(q)} \leq \frac{Mq}{J(q)}$$
$$\widehat{\text{FDR}}(p_{(I)}) \leq \frac{Mp_{(I)}}{I} \leq \alpha$$

# Interpretation

- $M_0$ true nulls, $M_1 = M - M_0$ true non-null
- Reject $H_{0,i}$ if $p_i \leq q$

$$\mathrm{FD} \approx M_0 q$$
$$\mathrm{D} = J(q) \equiv \max\{i : p_{(i)} < q\}$$

$$\mathrm{FDR} \approx \widehat{\mathrm{FDR}}(q) \equiv \frac{M_0 q}{J(q)} \leq \frac{Mq}{J(q)}$$
$$\widehat{\mathrm{FDR}}(p_{(I)}) \leq \frac{Mp_{(I)}}{I} \leq \alpha$$

# Interpretation

- $M_0$ true nulls, $M_1 = M - M_0$ true non-null
- Reject $H_{0,i}$ if $p_i \leq q$

$$\mathrm{FD} \approx M_0 q$$
$$\mathrm{D} = J(q) \equiv \max\{i : \ p_{(i)} < q\}$$

$$\mathrm{FDR} \approx \widehat{\mathrm{FDR}}(q) \equiv \frac{M_0 q}{J(q)} \leq \frac{Mq}{J(q)}$$
$$\widehat{\mathrm{FDR}}(p_{(I)}) \leq \frac{M p_{(I)}}{I} \leq \alpha$$

# Controlling Online FDR

# Back to our company

**BH policy:** Collect $M$ p-values every day, and run BH

**Problems:**
- Centralized
- Controls end-of-day FDR
  Not end-of-year FDR

$\rightarrow$ Online FDR control

# Back to our company

**BH policy:** Collect $M$ p-values every day, and run BH

**Problems:**

▶ Centralized

▶ Controls end-of-day FDR
Not end-of-year FDR

$\rightarrow$ Online FDR control

# Back to our company

**BH policy:** Collect $M$ p-values every day, and run BH

**Problems:**

- Centralized

- Controls end-of-day FDR
  Not end-of-year FDR

$\rightarrow$ Online FDR control

## Setting

**Null hypotheses:**

$$H_{0,1}, H_{0,2}, \ldots, H_{0,M}$$

**Sequence of p-values: one at each time**

$$p_1, p_2, p_3, \ldots$$

**Ground truth:**

$$\theta_1, \theta_2, \theta_3, \ldots [H_{0,i} : \theta_i = 0]$$

**Test ouput $(p_1^t = (p_1, \ldots, p_t)$:**

$$T_1(p_1^1), T_2(p_1^2), T_3(p_1^3), \cdots \in \{0, 1\}$$

[Foster, Stine, 2007]

# Setting

**Null hypotheses:**

$$H_{0,1}, H_{0,2}, \ldots, H_{0,M}$$

**Sequence of p-values: one at each time**

$$p_1, p_2, p_3, \ldots$$

**Ground truth:**

$$\theta_1, \theta_2, \theta_3, \ldots [H_{0,i} : \theta_i = 0]$$

**Test ouput $(p_1^t = (p_1, \ldots, p_t)$:**

$$T_1(p_1), T_2(p_2; T_1), T_3(p_3; T_1, T_2), \cdots \in \{0, 1\}$$

[Foster, Stine, 2007]

# What do we want to control?

- FD($n$) $\equiv$ False discoveries up to time $n$
- D($n$) $\equiv$ Total number of discoveries up to time $n$

$$\text{FDR}(n) \equiv \mathbb{E}\left\{ \frac{\text{FD}(n)}{\max(\text{D}(n), 1)} \right\}$$

Want FDR($n$) $\leq \alpha$ for all $n$, $\theta$

# What do we want to control?

- $\text{FD}(n) \equiv$ False discoveries up to time $n$
- $\text{D}(n) \equiv$ Total number of discoveries up to time $n$

$$\text{FDR}(n) \equiv \mathbb{E}\left\{\frac{\text{FD}(n)}{\max(\text{D}(n), 1)}\right\}$$

Want $\text{FDR}(n) \leq \alpha$ for all $n$, $\theta$

# Trivial approach (Bonferroni)

- Choose $\beta_i \in [0, 1]$, $\sum_{i=1}^{\infty} \beta_i \leq \alpha$
- Set

$$T_i = \begin{cases} 1 & \text{if } p_i \leq \beta_i, \\ 0 & \text{otherwise.} \end{cases}$$

Indeed

$$\text{FDR}(n) \leq \mathbb{E}\{\text{FD}(n)\} \leq \sum_{i:\,\theta_i=0} \mathbb{P}(p_i \leq \beta_i) = \sum_{i:\,\theta_i=0} \beta_i \leq \alpha$$

Very conservative!

# Trivial approach (Bonferroni)

- Choose $\beta_i \in [0, 1]$, $\sum_{i=1}^{\infty} \beta_i \leq \alpha$
- Set

$$T_i = \begin{cases} 1 & \text{if } p_i \leq \beta_i, \\ 0 & \text{otherwise.} \end{cases}$$

Indeed

$$\text{FDR}(n) \leq \mathbb{E}\{\text{FD}(n)\} \leq \sum_{i:\,\theta_i=0} \mathbb{P}(p_i \leq \beta_i) = \sum_{i:\,\theta_i=0} \beta_i \leq \alpha$$

Very conservative!

# Trivial approach (Bonferroni)

- Choose $\beta_i \in [0, 1]$, $\sum_{i=1}^{\infty} \beta_i \leq \alpha$
- Set

$$T_i = \begin{cases} 1 & \text{if } p_i \leq \beta_i, \\ 0 & \text{otherwise.} \end{cases}$$

Indeed

$$\text{FDR}(n) \leq \mathbb{E}\{\text{FD}(n)\} \leq \sum_{i:\,\theta_i=0} \mathbb{P}(p_i \leq \beta_i) = \sum_{i:\,\theta_i=0} \beta_i \leq \alpha$$

Very conservative!

# A simple rule

**LORD**     (Levels based On Recent Discovery)

- Choose $\beta_i \in [0, 1]$, $\sum_{i=1}^{\infty} \beta_i \leq \alpha$
- $\tau_i \equiv$ Time of the last discovery before $i$
- Set

$$T_i = \begin{cases} 1 & \text{if } p_i \leq \beta_{i-\tau_i}, \\ 0 & \text{otherwise.} \end{cases}$$

Each discovery resets everything.

# A simple rule

**LORD**    (Levels based On Recent Discovery)

- Choose $\beta_i \in [0, 1]$, $\sum_{i=1}^{\infty} \beta_i \leq \alpha$
- $\tau_i \equiv$ Time of the last discovery before $i$
- Set

$$T_i = \begin{cases} 1 & \text{if } p_i \leq \beta_{i - \tau_i}, \\ 0 & \text{otherwise.} \end{cases}$$

Each discovery resets everything.

# A theorem

**Theorem** (Javanmard, Montanari, 2015)

*If the null p-values are indepenent, then* LORD *satifies*

$$\sup_{\theta} \sup_{n} \mathrm{FDR}(n) \leq \alpha \,.$$

# Remarks

- Foster, Stine 2007:
  - Introduced model
  - Introduced *alpha investing rules*
  - Proved they control mFDR *(see next)*

- Last theorem applies to *generalized alpha investing*

- LORD uses very little information on the past!

# Remarks

- Foster, Stine 2007:
    - Introduced model
    - Introduced *alpha investing rules*
    - Proved they control mFDR *(see next)*
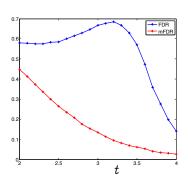
- Last theorem applies to *generalized alpha investing*

- LORD uses very little information on the past!

# Remarks

- Foster, Stine 2007:
    - Introduced model
    - Introduced *alpha investing rules*
    - Proved they control mFDR *(see next)*


- Last theorem applies to *generalized alpha investing*


- LORD uses very little information on the past!

# FDRvs mFDR

$$\mathrm{mFDR}_\eta(n) = \frac{\mathbb{E}_\theta\{\mathrm{FD}(n)\}}{\mathbb{E}_\theta\{\mathrm{D}(n)\} + \eta}$$

mFDR control $\not\Rightarrow$ FDR control

# Example



## Data

- $Z_1, \ldots, Z_{n_0} \sim_{iid} \mathrm{N}(0, 1)$, $(Z_{n_0+1}, \ldots, Z_n) \sim \mathrm{N}(\theta_* \mathbf{1}, \rho \mathbf{1}\mathbf{1}^\mathsf{T} + \overline{\rho}\mathbf{I})$
- $n = 3000$, $n_0 = 2700$, $\theta_* = 2$, $\rho = 0.9$

## Rule

$$T_i = \begin{cases} 1 & \text{if } |Z_i| \geq t, \\ 0 & \text{otherwise.} \end{cases}$$
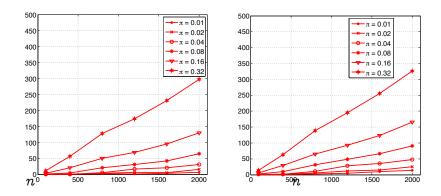
# Statistical power?

**Two-groups model**

$$\theta_i \sim_{iid} \text{Bernoulli}(\pi),$$

$$\mathbb{P}_{\theta_i}(p_i \leq x) = \begin{cases} F(x) = x & \text{if } \theta_i = 0, \\ G(x) & \text{otherwise.} \end{cases}$$

**'Discoveries should keep coming'**

- A good rule should have $D(n) = \Theta(n)$.

# Two experiments



- **Left:** $\theta_i \sim_{iid} (1 - \pi)\delta_0 + \pi N(0, \sigma^2)$, $Z_i \sim N(0, \theta_i)$
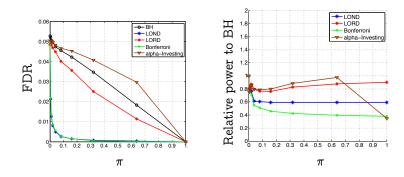- **Right:** $\theta_i$ re-ordered, decreasing $|\theta_i|$

# A theorem

**Theorem** (Javanmard, Montanari, 2015)

*Assume the two-groups model, and use of* LORD. *Then, almost surely*

$$\lim_{n \to \infty} \frac{1}{n} D(n) \geq \mathcal{A}(G, \beta),$$

$$\mathcal{A}(G, \beta) \equiv \left( \sum_{k=1}^{\infty} e^{-\sum_{\ell=1}^{k} G(\beta_\ell)} \right)^{-1}.$$

- $\mathcal{A}(G, \beta) > 0$ strictly if $G(\beta_\ell) > (1 + \varepsilon)/\ell$ for all $\ell$ large enough.
- Sufficient $G(x) \approx G_0 x^{1+\delta}$ as $x \to 0$.
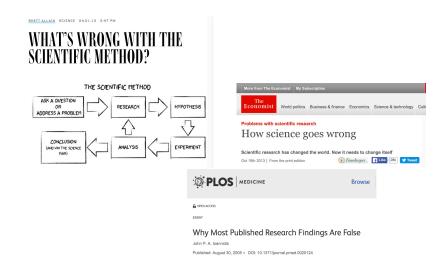
# Comparison under the Gaussian two-groups model



$$\mathrm{TD}(n) = \text{ True discoveries}$$

$$\mathrm{RelativePower}(n) \equiv \mathbb{E}\Big\{ \frac{\mathrm{TD}(n)}{\max(\mathrm{TD_{BH}}(n), 1)} \Big\}.$$

▶ $n = 1000$, $\sigma^2 = 2 \log n$

Conclusion

What if I am not CTO of a big-data company?

# Take the "company" as a metaphor for science

# Conclusion

▶ FDR control is fundamental for reasoning about data

▶ Online FDR is likely more realistic

Thanks!

# Conclusion

▶ FDR control is fundamental for reasoning about data

▶ Online FDR is likely more realistic

Thanks!

# Conclusion

- FDR control is fundamental for reasoning about data

- Online FDR is likely more realistic

Thanks!