

Controlling the False Discovery Rate: Understanding and Extending the Benjamini-Hochberg Method

Christopher R. Genovese

Department of Statistics

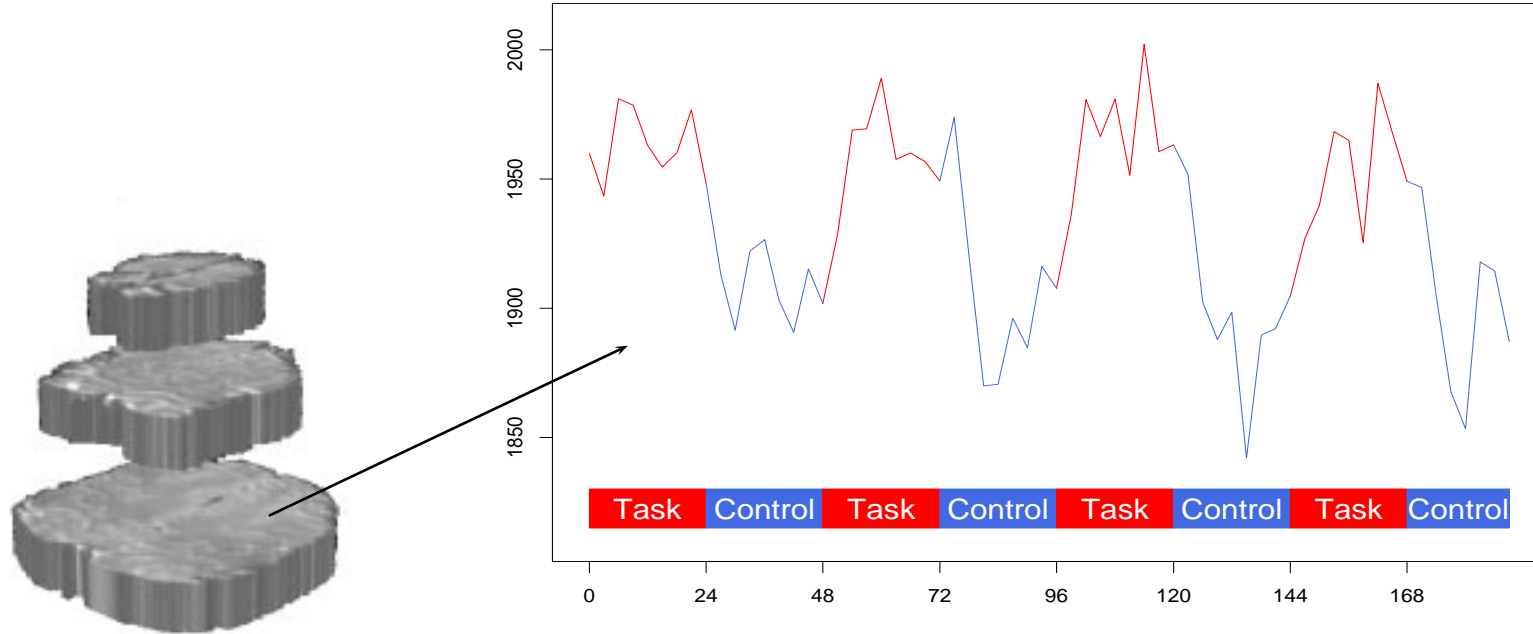
Carnegie Mellon University

joint work with Larry Wasserman

This work partially supported by NSF Grant SES 9866147.

Motivating Example #1: fMRI

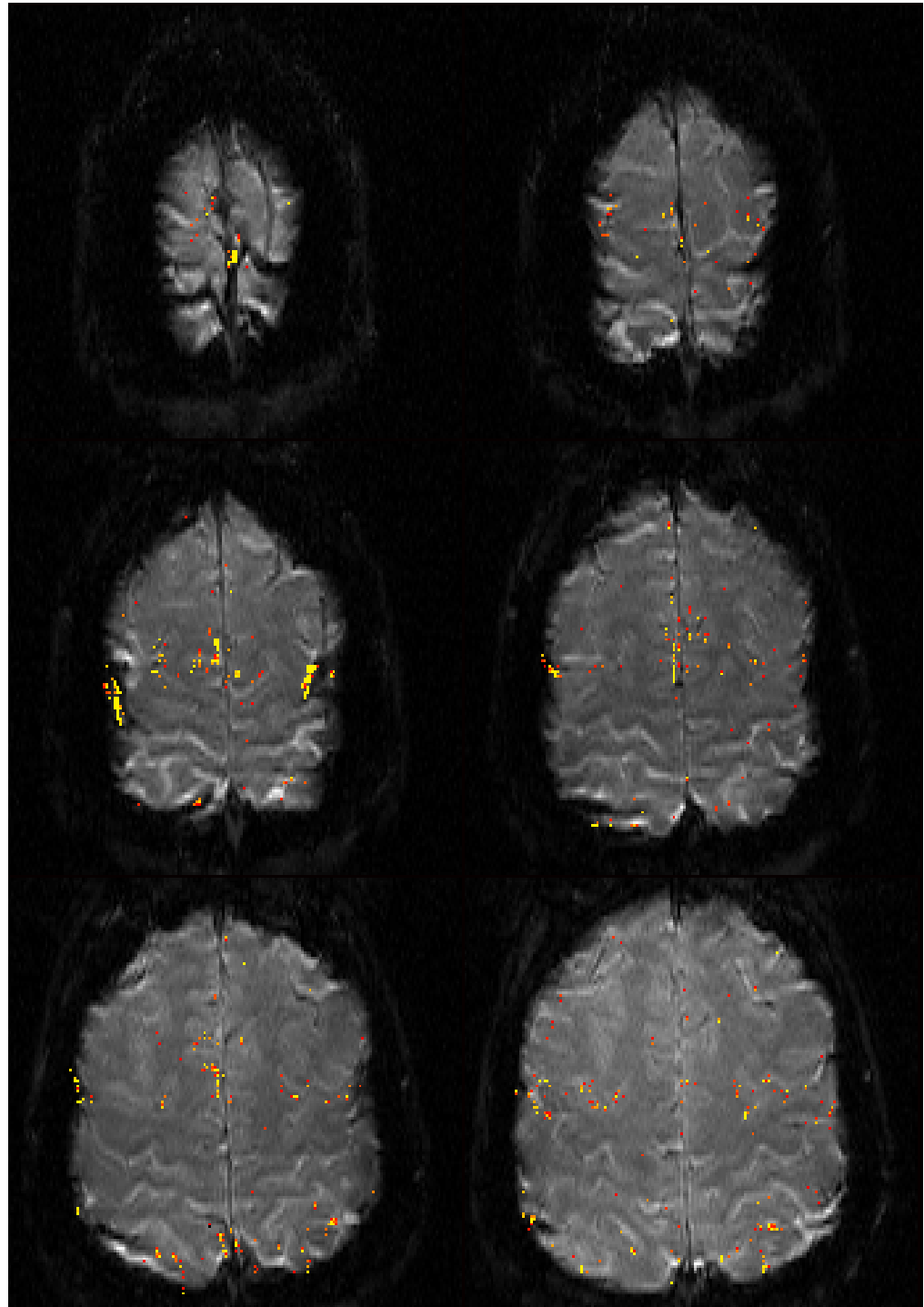
- fMRI Data: Time series of 3-d images acquired while subject performs specified tasks.



- Goal: Characterize task-related signal changes caused (indirectly) by neural activity. [See, for example, Genovese (2000), *JASA* 95, 691.]

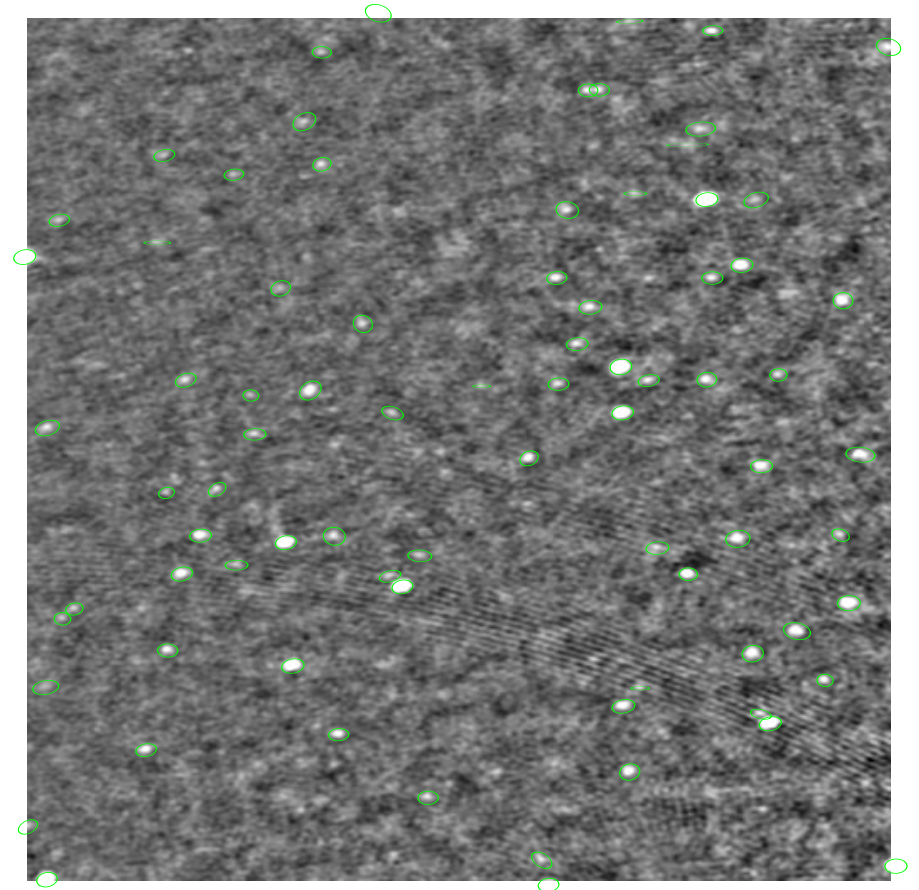
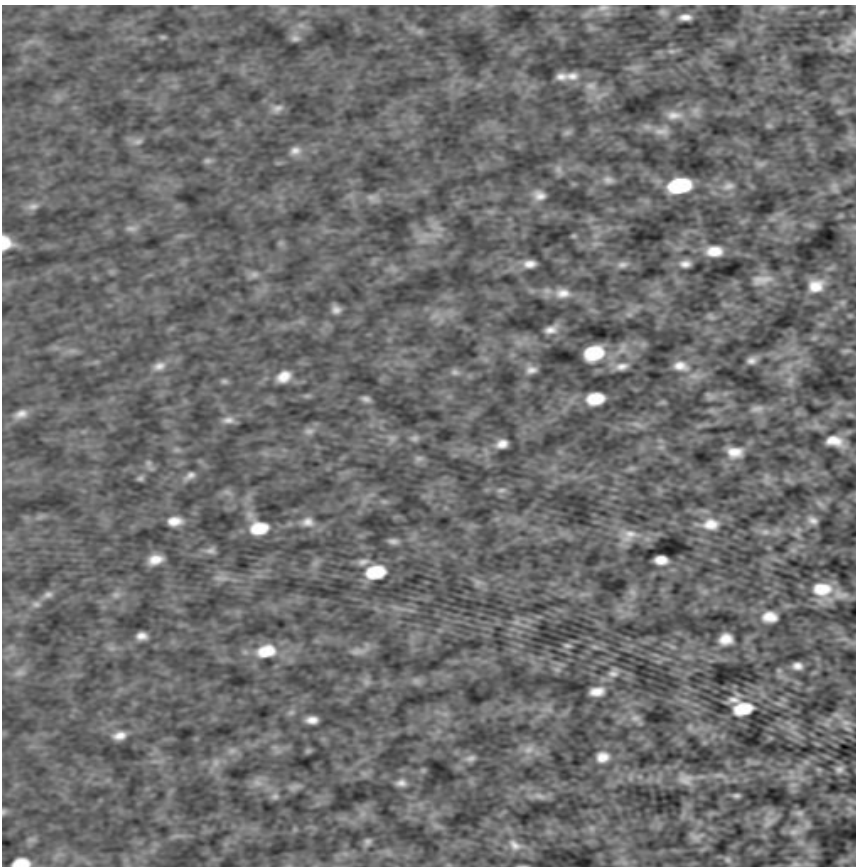
fMRI (cont'd)

Perform hypothesis tests at many thousands of volume elements to identify loci of activation.



Motivating Example #2: Source Detection

- Interferometric radio telescope observations processed into digital image of the sky in radio frequencies.
- Signal at each pixel is a mixture of source and background signals.



Motivating Example #3: DNA Microarrays

- New technologies allow measurement of gene expression for thousands of genes simultaneously.

		Subject				Subject			
		1	2	3	...	1	2	3	...
Gene	1	X_{111}	X_{121}	X_{131}	...	X_{112}	X_{122}	X_{132}	...
	2	X_{211}	X_{221}	X_{231}	...	X_{212}	X_{222}	X_{232}	...
	3	⋮	⋮	⋮	...	⋮	⋮	⋮	...
	4								
	5								
	6								
	⋮								
		<u>Condition 1</u>				<u>Condition 2</u>			

- Goal: Identify genes associated with differences among conditions.
- Typical analysis: hypothesis test at each gene.

The Multiple Testing Problem

- Perform m simultaneous hypothesis tests.

Classify results as follows:

	H_0 Retained	H_0 Rejected	Total
H_0 True	$N_{0 0}$	$N_{1 0}$	M_0
H_0 False	$N_{0 1}$	$N_{1 1}$	M_1
Total	$m - R$	R	m

Only R is observed here.

- Assess outcome through combined error measure.
This binds the separate decision rules together.

Multiple Testing (cont'd)

- Traditional methods seek strong control of familywise Type I error (FWER).
 - Weak Control: If all nulls true, $P\{N_{1|0} > 0\} \leq \alpha$.
 - Strong Control: Corresponding statement holds for any subset of tests for which all nulls are true.

For example, Bonferroni correction provides strong control but is quite conservative.

- Can power be improved while maintaining control over a meaningful measure of error?

Enter Benjamini & Hochberg ...

FDR and the BH Procedure

- Define the *realized* False Discovery Rate (FDR) by

$$\text{FDR} = \begin{cases} \frac{N_{1|0}}{R} & \text{if } R > 0, \\ 0, & \text{if } R = 0. \end{cases}$$

- Benjamini & Hochberg (1995) define a sequential p-value procedure that controls *expected* FDR.

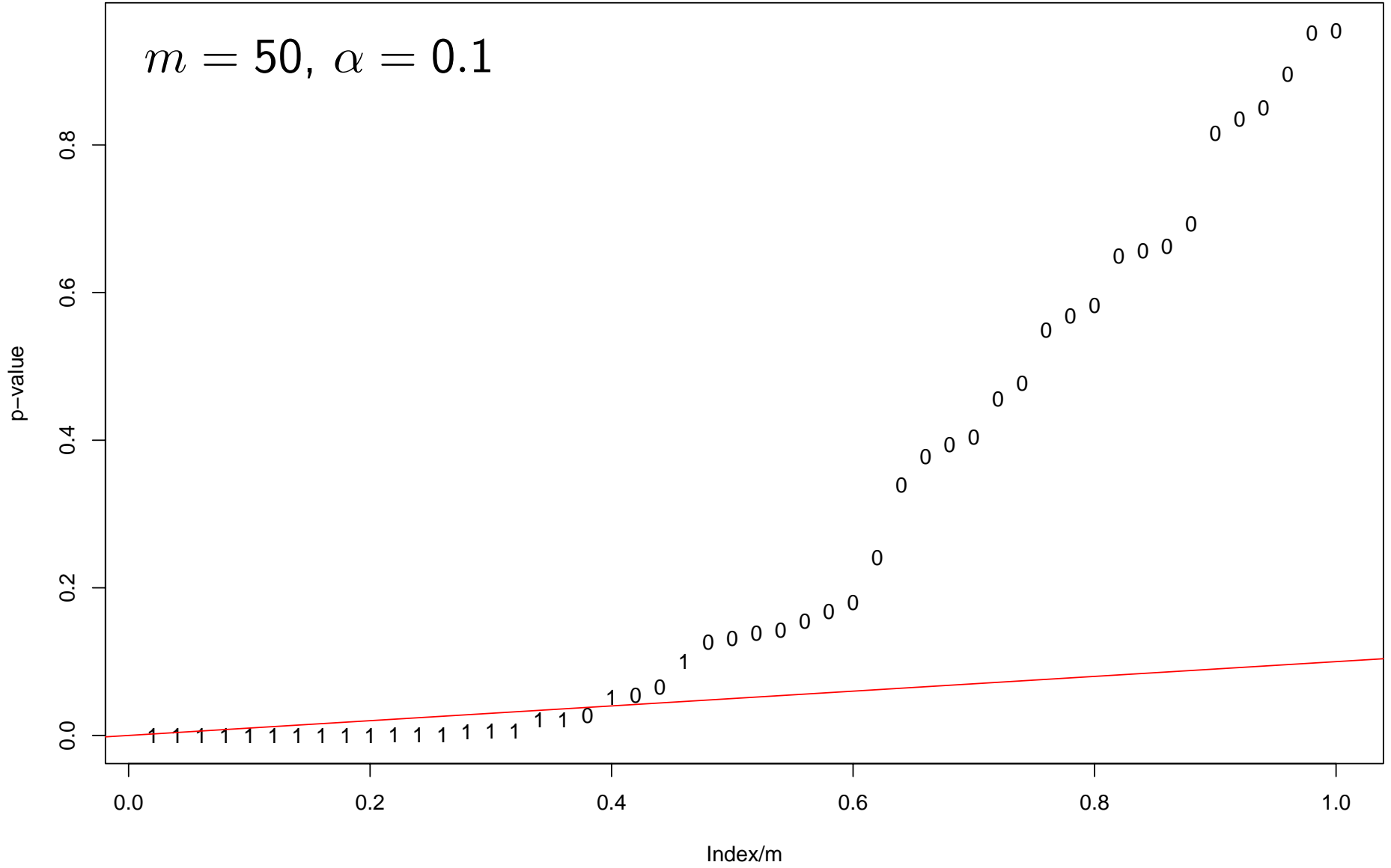
Specifically, the BH procedure guarantees

$$E(\text{FDR}) \leq \frac{M_0}{m} \alpha \leq \alpha$$

for a pre-specified $0 < \alpha < 1$.

(The first inequality is an equality in the continuous case.)

$m = 50, \alpha = 0.1$



- The BH procedure for p-values P_1, \dots, P_m :

0. Select $0 < \alpha < 1$.

1. Define $P_{(0)} \equiv 0$ and

$$R_{\text{BH}} = \max \left\{ 0 \leq i \leq m: P_{(i)} \leq \alpha \frac{i}{m} \right\}.$$

2. Reject H_0 for every test where $P_j \leq P_{(R_{\text{BH}})}$.

- Several variant procedures also control E(FDR).
- Bound on E(FDR) holds if p-values are independent or positively dependent (Benjamini & Yekutieli, 2001). Storey (2001) shows it holds under a possibly weaker condition.
- By replacing α with $\alpha / \sum_{i=1}^m 1/i$, control E(FDR) at level α for any joint distribution on the p-values. (Very conservative!)

Road Map

1. Preliminaries

- Considering both types of errors: The False Nondiscovery Rate (FNR)
- Models for realized FDR and FNR
- FDR and FNR as stochastic processes

2. Understanding BH

- Re-express BH procedure as plug-in estimator
- Asymptotic behavior of BH
- Improving the power – more general plug-ins
- Asymptotic risk comparisons

3. Extensions to BH

- Conditional risk
- FDR control as an estimation problem
- Confidence intervals for realized FDR
- Confidence thresholds

Recent Work on FDR

Benjamini & Hochberg (1995)

Benjamini & Liu (1999)

Benjamini & Hochberg (2000)

Benjamini & Yekutieli (2001)

Storey (2001a,b)

Efron, et al. (2001)

Storey & Tibshirani (2001)

Tusher, Tibshirani, Chu (2001)

Abromovich, et al. (2000)

Genovese & Wasserman (2001a,b)

See also technical reports 735, 737, 747, 752, 754
at <http://lib.stat.cmu.edu/www/cmu-stats/tr/>.

The False Nondiscovery Rate

- Controlling FDR alone only deals with Type I errors.
- Define the *realized* False Nondiscovery Rate as follows:

$$\text{FNR} = \begin{cases} \frac{N_{0|1}}{m - R} & \text{if } R < m, \\ 0 & \text{if } R = m. \end{cases}$$

This is the proportion of false non-rejections among those tests whose null hypothesis is not rejected.

- Idea: Combine FDR and FNR in assessment of procedures.

Basic Models

- Let $H_i = 0$ (or 1) if the i^{th} null hypothesis is true (or false).
These are unobserved.
- Let P_i be the i^{th} p-value.
- We assume that $(P_1, H_1), \dots, (P_m, H_m)$ are independent with $P_i | \{H_i = 0\} \sim \text{Uniform}\langle 0, 1 \rangle$, and $P_i | \{H_i = 1\} \sim F \in \mathcal{F}$,
a class of alternative p-value distributions.
 - Under the *conditional model*, H_1, \dots, H_m are fixed, unknown.
 - Under the *mixture model*, we assume each $H_i \sim \text{Bernoulli}\langle a \rangle$.
- Define $M_0 = \sum_i (1 - H_i)$ and $M_1 = \sum_i H_i = m - M_0$.
Under the *mixture model*, $M_1 \sim \text{Binomial}\langle m, a \rangle$.
Under the *conditional model*, these are fixed.

Basic Models (cont'd)

- Typical examples:

- Parametric family: $\mathcal{F}_\Theta = \{F_\theta: \theta \in \Theta\}$

- Concave, continuous distributions

$$\mathcal{F}_C = \{F: F \text{ concave, continuous cdf with } F \prec U\}.$$

- Remark: The assumption of the mixture model does not require the same alternative for each test. For example, suppose that

$$\begin{aligned} P_i | \Psi_i = \psi &\sim F_\psi \\ \Psi_i &\sim H \end{aligned}$$

Then, $F = \int F_\psi dH(\psi)$.

Recurring Notation

$m, M_0, N_{1 0}$	# of tests, true nulls, false discoveries
a	Mixture weight on a lternative
$H^m = (H_1, \dots, H_m)$	Unobserved true classifications
$P^m = (P_1, \dots, P_m)$	Observed p-values
$P_{()}^m = (P_{(1)}, \dots, P_{(m)})$	Sorted p-values (define $P_{(0)} \equiv 0$)
U	CDF of Uniform $\langle 0, 1 \rangle$
F, f	Alternative CDF and density
$G = (1 - a)U + aF$	Marginal CDF of P_i (mixture model)
\hat{G}	Empirical CDF of P^m
$\epsilon_m = \sqrt{\frac{1}{2m} \log \left(\frac{2}{\beta} \right)}$	DKW bound $1 - \beta$ quantile of $\ \hat{G} - G\ _\infty$

Multiple Testing Procedures

- A multiple testing procedure T is a map $[0, 1]^m \rightarrow [0, 1]$, where the null hypotheses are rejected in all those tests for which $P_i \leq T(P^m)$.
- Examples:

Uncorrected testing	$T_U(P^m) = \alpha$
Bonferroni	$T_B(P^m) = \alpha/m$
Benjamini-Hochberg	$T_{\text{BH}}(P^m) = P_{(R_{\text{BH}})}$
Fixed Threshold	$T_t(P^m) = t$
First- r	$T_{(r)}(P^m) = P_{(r)}$

FDR and FNR as Stochastic Processes

- Define the realized FDR and FNR processes, respectively, by

$$\text{FDR}(t) \equiv \text{FDR}(t; P^m, H^m) = \frac{\sum_i \mathbf{1}\{P_i \leq t\} (1 - H_i)}{\sum_i \mathbf{1}\{P_i \leq t\} + \prod_i \mathbf{1}\{P_i > t\}}$$
$$\text{FNR}(t) \equiv \text{FNR}(t; P^m, H^m) = \frac{\sum_i \mathbf{1}\{P_i > t\} H_i}{\sum_i \mathbf{1}\{P_i > t\} + \prod_i \mathbf{1}\{P_i \leq t\}}.$$

- For procedure T , the realized FDR and FNR are obtained by evaluating these processes at $T(P^m)$.
- Both these processes converge to Gaussian processes outside a neighborhood of 0 and 1 respectively.

FDR and FNR as Stochastic Processes (cont'd)

- For example, define

$$Z_m(t) = \sqrt{m} (\text{FDR}(t) - Q(t)), \quad \delta \leq t \leq 1,$$

where $0 < \delta < 1$ and $Q(t) = (1 - a)U/G$.

- Let Z be a mean 0 Gaussian process on $[\delta, 1]$ with covariance kernel

$$K(s, t) = \frac{a(1 - a) [(1 - a)stF(s \wedge t) + aF(s)F(t)(s \wedge t)]}{G^2(s)G^2(t)}.$$

- Then, $Z_m \rightsquigarrow Z$.

BH as a Plug-in Procedure

- Let \hat{G} be the empirical cdf of P^m under the mixture model. Ignoring ties, $\hat{G}(P_{(i)}) = i/m$, so BH equivalent to

$$T_{\text{BH}}(P^m) = \arg \max \left\{ t: \hat{G}(t) = \frac{t}{\alpha} \right\}.$$

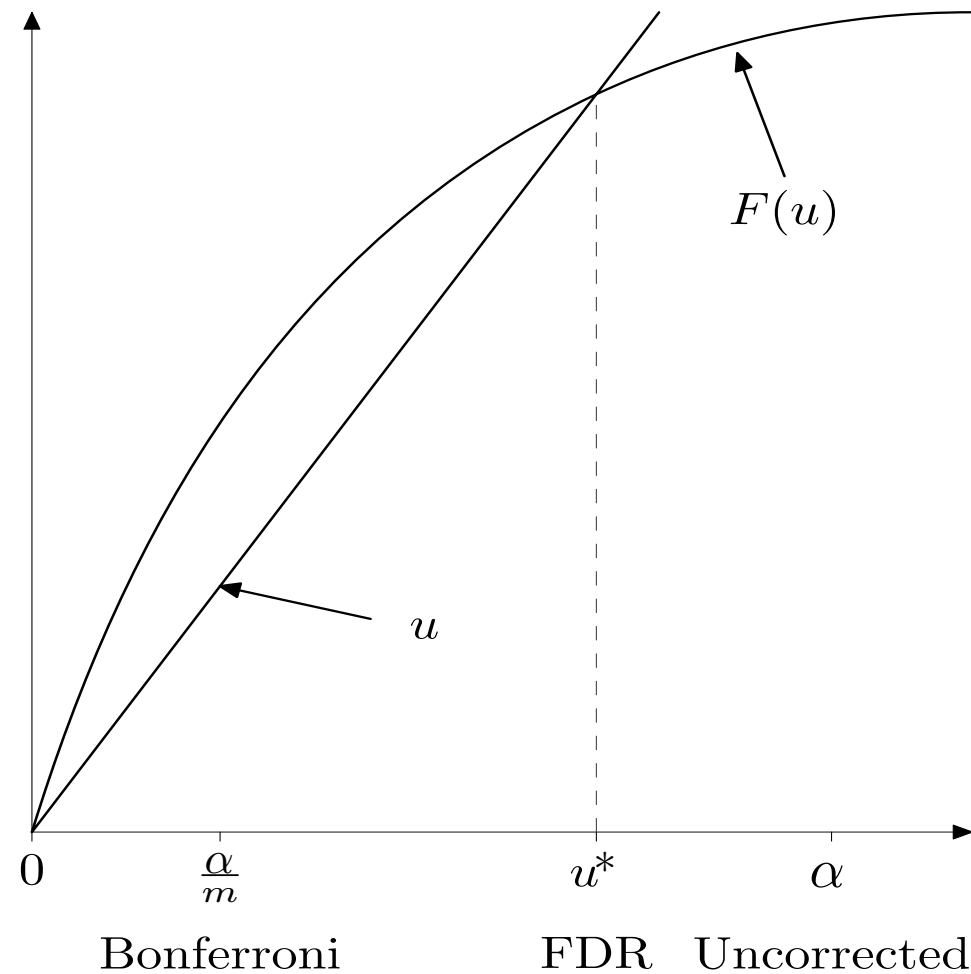
- We can think of this as a plug-in procedure for estimating

$$\begin{aligned} u^*(a, F) &= \arg \max \left\{ t: G(t) = \frac{t}{\alpha} \right\} \\ &= \arg \max \{ t: F(t) = \beta t \}, \end{aligned}$$

where $\beta = (1 - \alpha + \alpha a)/\alpha a$.

Asymptotic Behavior of BH Procedure

This yields the following picture:



Optimal Thresholds

- Under the mixture model and in the continuous case,

$$E(\text{FDR}(T_{\text{BH}}(P^m))) = (1 - a)\alpha.$$

- The BH procedure overcontrols $E(\text{FDR})$ and thus will not in general minimize $E(\text{FNR})$.
- This suggests finding a plug-in estimator for

$$\begin{aligned} t^*(a, F) &= \arg \max \left\{ t: G(t) = \frac{(1 - a)t}{\alpha} \right\} \\ &= \arg \max \{ t: F(t) = (\beta - 1/\alpha)t \}, \end{aligned}$$

where $\beta - 1/\alpha = (1 - a)(1 - \alpha)/a\alpha$.

- Note that $t^* \geq u^*$.

Optimal Thresholds (cont'd)

- For each $0 \leq t \leq 1$,

$$E(\text{FDR}(t)) = \frac{(1-a)t}{G(t)} + O\left(\frac{1}{\sqrt{m}}\right)$$

$$E(\text{FNR}(t)) = a \frac{1-F(t)}{1-G(t)} + O\left(\frac{1}{\sqrt{m}}\right).$$

- Ignoring $O(m^{-1/2})$ terms and choosing t to minimize $E(\text{FNR}(t))$ subject to $E(\text{FDR}(t)) \leq \alpha$, yields $t^*(a, F)$ as the optimal threshold.
- Can the potential improvement in power be achieved when estimating t^* ?

Yes, if F sufficiently far from U .

Operating Characteristics of the BH Method

- Define the misclassification risk of a procedure T by

$$R_M(T) = \frac{1}{m} \sum_{i=1}^m \mathbf{E} \left| \mathbf{1} \{P_i \leq T(P^m)\} - H_i \right|.$$

This is the average fraction of errors of both types.

- Then $R_M(T_{\text{BH}}) \sim R(a, F)$ as $m \rightarrow \infty$, where

$$R(a, F) = (1 - a)u^* + a(1 - F(u^*)) = (1 - a)u^* + a(1 - \beta u^*).$$

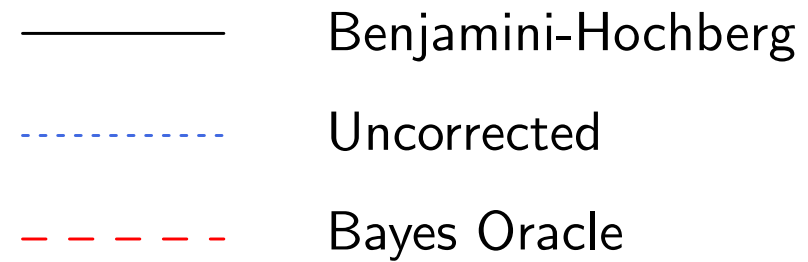
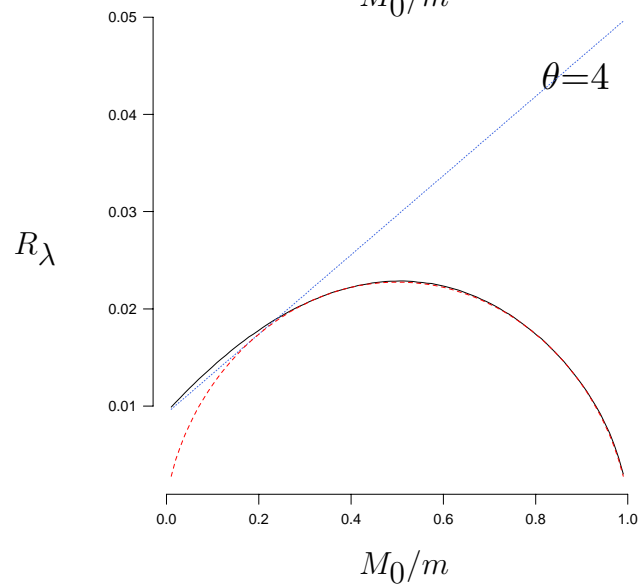
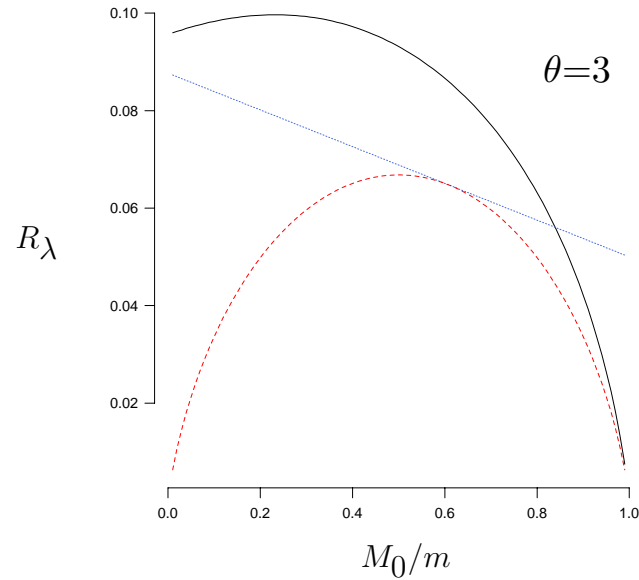
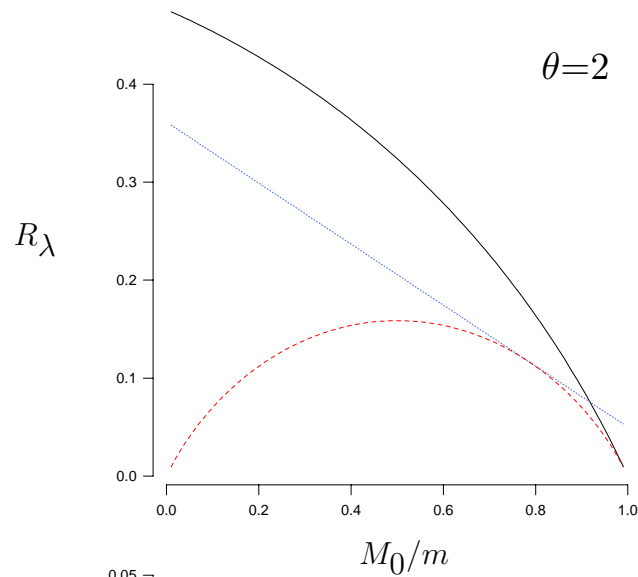
- Compare this to Uncorrected and Bonferroni and the oracle rule $T_O(P^m) = b$ where b solves $f(b) = (1 - a)/a$.

$$R_M(T_U) = (1 - a)\alpha + a(1 - F(\alpha))$$

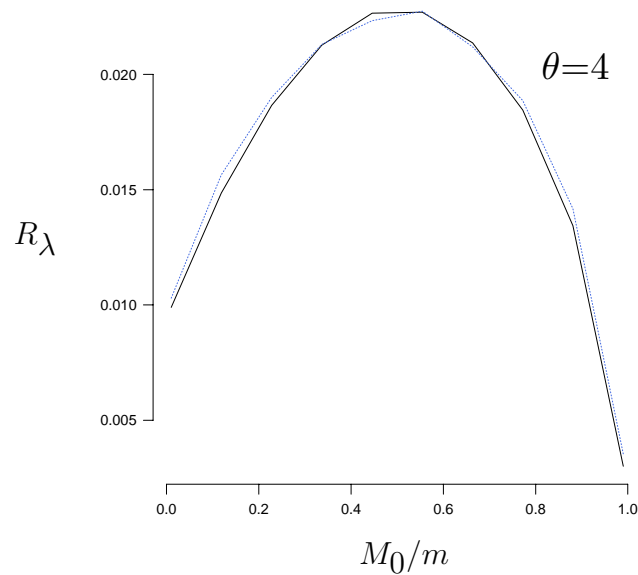
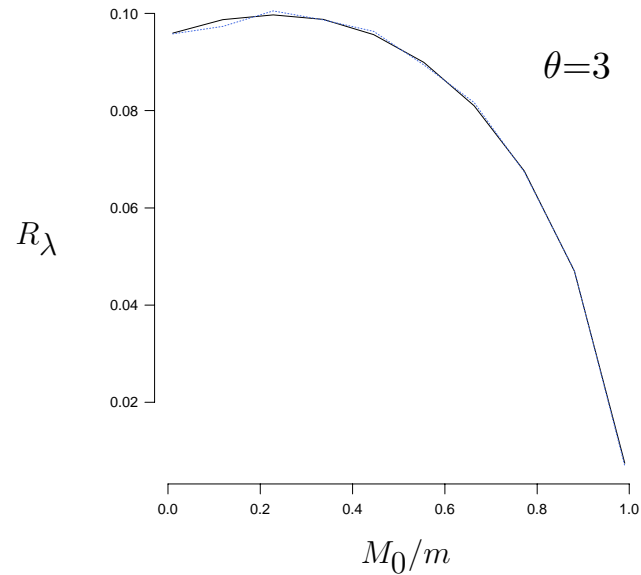
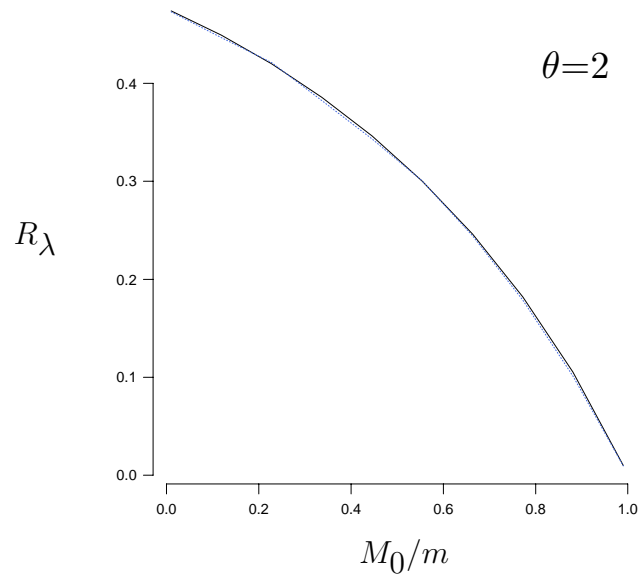
$$R_M(T_B) = (1 - a)\frac{\alpha}{m} + a\left(1 - F\left(\frac{\alpha}{m}\right)\right)$$

$$R_M(T_O) = (1 - a)b + a(1 - F(b)).$$

Normal $\langle\theta, 1\rangle$ Model, $\alpha = 0.05$



Check Approximation Accuracy $m = 100$



— Asymptotic
- - - Exact

Extension 1: Conditional Risk

- It is intuitively appealing (cf. Kiefer, 1977) to assess the performance of a procedure conditionally given the ordered p-values.
- When conditioning, we need only consider the $m + 1$ procedures $T_{(r)}(P^m) = P_{(r)}$ for $r = 0, \dots, m$.
- Under the conditional model, once $P_{()}^m$ is observed, only the randomness in the labelling of the true classifications remains.
- Consider a parametric family $\mathcal{F} = \{F_\theta: \theta \in \Theta\}$ of alternative p-value distributions.

Then, (M_0, θ) becomes the unknown parameter.

Begin by treating this as known.

Conditional Risk (cont'd)

- Define a conditional risk for $\lambda \geq 0$ by

$$R_\lambda(r; M_0, \theta \mid P_{()}^m) = \mathbf{E}_{M_0, \theta} \left[\text{FNR}(P_{(r)}) + \lambda \text{FDR}(P_{(r)}) \mid P_{()}^m \right],$$

where M_0 and r are in $\{0, \dots, m\}$ and $\theta \in \Theta$.

- Here λ determines the balance between the two error types.

It also serves as a Lagrange multiplier for the optimization problem:

$$r_* = \arg \min_{0 \leq r \leq m} \mathbf{E}_{M_0, \theta}(\text{FNR}(P_{(r)}) \mid P_{()}^m)$$

subject to

$$\mathbf{E}_{M_0, \theta}(\text{FDR}(P_{(r)}) \mid P_{()}^m) \leq \alpha.$$

Conditional Risk (cont'd)

- This problem can be solved exactly:
 - Closed form for conditional distribution of FDR and FNR based on expressions for

$$P_{M_0, \theta} \{ N_{1|0} = k \mid P_{()}^m \} \quad \text{and} \quad E_{M_0, \theta} (N_{1|0} \mid P_{()}^m)$$

derived via generating function methods.

- Find R_λ -minimizer explicitly.
 - Select λ to satisfy the constraint.
- Remark: The R_λ minimizing conditional procedure also minimizes the unconditional R_λ risk, but the constrained optimization problem is harder to solve unconditionally.

Conditional Risk (cont'd)

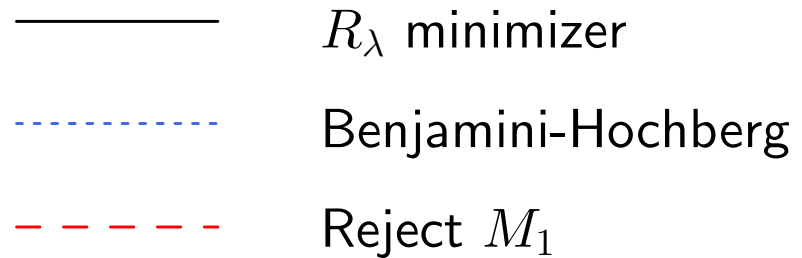
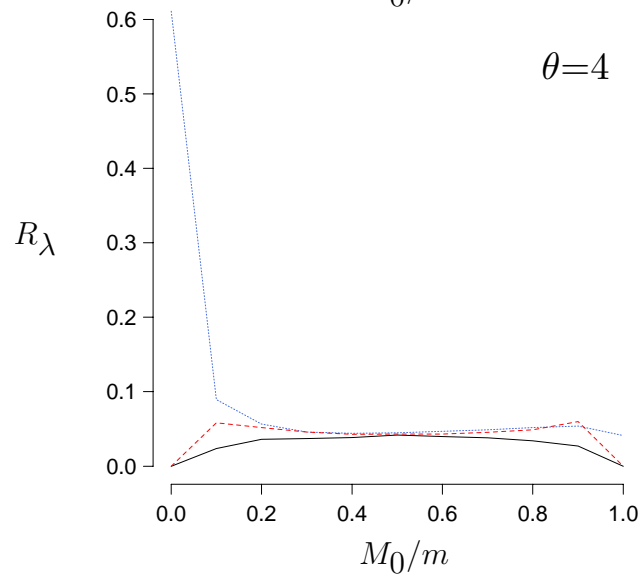
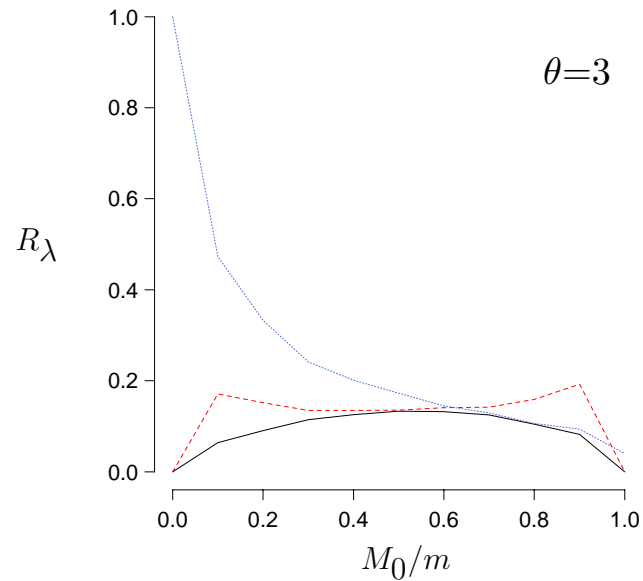
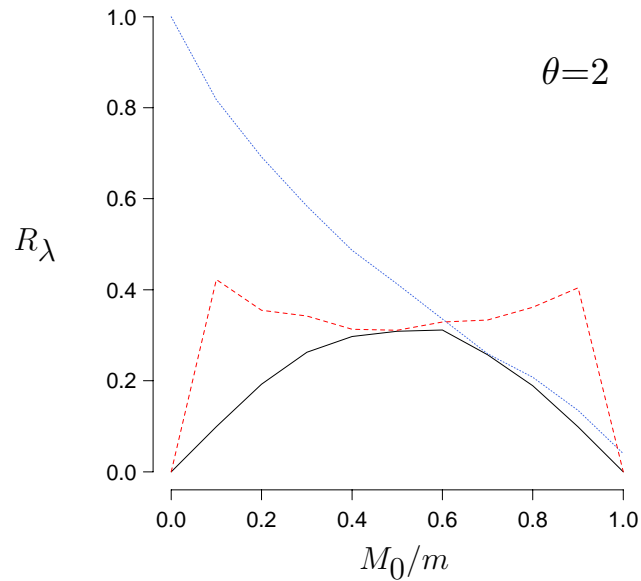
- For M_0 unknown case, R_λ dominated by extremes, $M_0 = 0$ or $M_0 = m$.
- One approach: minimize conditional Bayes risk based on R_λ

$$R_\lambda(r; \theta | P_{()}^m) = \sum_{m_0=0}^m R_\lambda(r; m_0, \theta | P_{()}^m) p_\theta(m_0 | P_{()}^m),$$

where $p_\theta(m_0 | P_{()}^m)$ is derived from a specified (e.g., Uniform) prior on $\{0, \dots, m\}$.

- This minimizes the unconditional Bayes risk.

Normal $\langle\theta, 1\rangle$ Model, $m = 100$, $\alpha = 0.05$



Bayesian FDR

- These conditional results yield the posterior distribution of FDR and FNR (and related quantities).

No simulation necessary: can compute full posterior directly.

- Suggests the procedure $T_{\text{Bayes}}(P^m) = P_{(r_*)}$, where

$$r_* = \arg \min_{0 \leq r \leq m} \text{E}(\text{FNR}(P_{(r)}) \mid P_{()}^m)$$

subject to

$$\text{E}(\text{FDR}(P_{(r)}) \mid P_{()}^m) \leq \alpha.$$

- This procedure has good asymptotic frequentist performance.

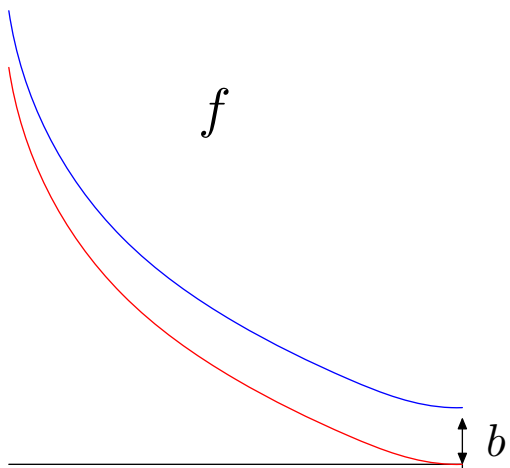
Extension 2: Estimating a and F

- To compute plug-in estimates that approximate the optimal threshold, we need a good estimate of a .

For instance,

$$\hat{t}^* = \arg \max \left\{ t: \hat{G}(t) = \frac{(1 - \hat{a})t}{\alpha} \right\}.$$

- For confidence thresholds, need estimate of a and F .
- Identifiability



If $\min f = b > 0$, can write $F = (1-b)U + bF_0$, so many (a, F) pairs yield the same G .

If $f = F'$ is decreasing with $f(1) = 0$, then (a, F) is identifiable.

Estimating a and F (cont'd)

- Even when non-identifiable, a can be bounded from below by \underline{a} .
 $a - \underline{a}$ is typically small. For example, $a - \underline{a} = ae^{-n\theta^2/2}$ in the two-sided test of $\theta = 0$ versus $\theta \neq 0$ in the Normal $\langle\theta, 1\rangle$ model.
- Parametric Case: (a, θ) typically identifiable; use MLE.
- Non-parametric case:
 - Derived a $1 - \beta$ confidence interval for \underline{a} and thus a .
 - When F concave, get $\hat{a}_{\text{LCM}} = \underline{a} + O_P(m^{-1/3})$.
Can do better with further smoothness assumptions.
 - In general, requires density estimate of g .
 - Can estimate F by: $\hat{F}_m = \operatorname{argmin}_H \|\hat{G} - (1 - \hat{a})U - \hat{a}H\|_\infty$.
Consistent for reduced F if \hat{a} consistent for \underline{a} .
- Note: Assumption of concavity has a big effect.

Extension 3: Confidence Intervals

- Beyond controlling FDR and FNR on average, we would like to be able to make inferences about the realized quantities.
- Want to find $c(P^m, T)$, for any procedure T , such that

$$P_{a,F} \left\{ \text{FDR}(T(P^m)) \leq c(P^m, T) \right\} \geq 1 - \alpha,$$

at least asymptotically.

- Let $r(P^m, T) = \sum_i 1 \{P_i \leq T(P^m)\}$ be the number of rejections.
- Template: $c(P^m, T)$ is a $1 - \beta$ quantile of the sum of $r(P^m, T)$ independent Bernoulli $\langle q_i \rangle$ variables.

Here, the q_i bound $q(P_{(i)})$ with high probability, where $q(t) = (1 - a)/g(t)$ gives the conditional distribution of H_1 given P_1 .

The q_i depend on the assumed class \mathcal{F} of alternative p-value distributions.

Confidence Intervals (cont'd)

- Case 1: $\mathcal{F} = \{F_\theta: \theta \in \Theta\}$

– Asymptotic: $\beta = \alpha$ and $q_i = \frac{1 - \hat{a}}{1 - \hat{a} + \hat{a}f_{\hat{\theta}}(P_{(i)})}$.

– Exact: Let $\beta = 1 - \sqrt{1 - \alpha}$ and let Ψ_m be a $1 - \beta$ confidence set for (a, θ) .

$$q_i = \sup_{\Psi_m} \frac{1 - a}{1 - a + af_\theta(P_{(i)})}$$

Example: Invert DKW Envelope

$$\Psi_m = \{(a, \theta): \|G_{a,\theta} - \hat{G}\|_\infty \leq \epsilon_m\}.$$

Confidence Intervals (cont'd)

- Case 2: $\mathcal{F} = \{F: F \text{ concave, continuous cdf and } F \prec U\}$.

Can find a minimal concave cdf \underline{G} in DKW envelope. Define

$$q_i = \frac{1 - \hat{a}}{\underline{g}(P_i)},$$

and use $\beta = 1 - (1 - \alpha)^{1/3}$.

- May be possible to obtain nonparametric results in non-concave case, but the intervals appear to be hopelessly wide in practice.
- Bayesian posterior intervals also have asymptotically valid frequentist coverage.
- All these results extend to give joint confidence intervals for FDR and FNR.

Extension 4: Confidence Thresholds

- In practice, it would be useful to have a procedure T_C that guarantees

$$P_G\{\text{FDR}(T_C) > c\} \leq \alpha$$

for some specified c and α .

We call this a $(1 - \alpha, c)$ *confidence threshold procedure*.

- Two approaches: an asymptotic threshold using the Bootstrap, and an exact (small-sample) threshold requiring numerical search.
- Here, I'll discuss the case where a is known.

In general, can use an estimate of a , but this introduces additional complexity.

Bootstrap Confidence Thresholds

- First guess: Choose T such that

$$P_{\hat{G}}\{ \text{FDR}^*(T) \leq c \} \geq 1 - \alpha.$$

Unfortunately, this fails.

- The problem is an additional bias term:

$$\begin{aligned} 1 - \alpha &= P_{\hat{G}}\{ \text{FDR}^*(T) \leq c \} \\ &\approx P_G\{ \text{FDR}(T) \leq c + (Q(T) - \hat{Q}(T)) \} \\ &\neq P_G\{ \text{FDR}(T) \leq c \}, \end{aligned}$$

where $Q = (1 - \alpha)U/G$ and $\hat{Q} = (1 - \alpha)U/\hat{G}$.

Bootstrap Confidence Thresholds (cont'd)

- Let $\beta = \alpha/2$ and $\epsilon_m = \sqrt{\frac{1}{2m} \log \left(\frac{2}{\beta} \right)}$.
- Procedure
 1. Draw $H_1^* \dots, H_m^*$ iid Bernoulli $\langle a \rangle$
 2. Draw $P_i^* | H_i^*$ from $(1 - H_i^*)U + H_i^* \hat{F}$.
 3. Define $\Omega_c^*(t) = \sum_i I\{P_i^* \leq t\}(1 - H_i^* - c)$.
 4. Use threshold defined by

$$T_C = \max \left\{ t: P_{\hat{G}} \left\{ \Omega_c^*(t) \leq -c \epsilon_m \right\} \geq 1 - \beta \right\}.$$

- Then,

$$P_G \left\{ \text{FDR}(T_C) \leq c \right\} \geq 1 - \alpha + O \left(\frac{1}{\sqrt{m}} \right).$$

Exact Confidence Thresholds

- Let \mathcal{M}_β be a $1 - \beta$ confidence set for M_0 , derived from the Binomial $\langle m, 1 - \alpha \rangle$.

- Define

$$S(t; h^m, p^m) = \frac{\sum_i 1 \{p_i \leq t\} (1 - h_i)}{\sum_i (1 - h_i)},$$

$$\mathcal{U} = \left\{ (h^m, p^m) : \sum_i (1 - h_i) \in \mathcal{M}_\alpha \text{ and } \|S(\cdot; h^m, p^m) - U\|_\infty \leq \epsilon_{m_0} \right\},$$

where $\epsilon_m = \sqrt{\log(2/\beta)/2m}$ as above.

- Then, if $\beta = 1 - \sqrt{1 - \alpha}$, $P_G\{(H^m, P^m) \in \mathcal{U}\} \geq 1 - \alpha$ and

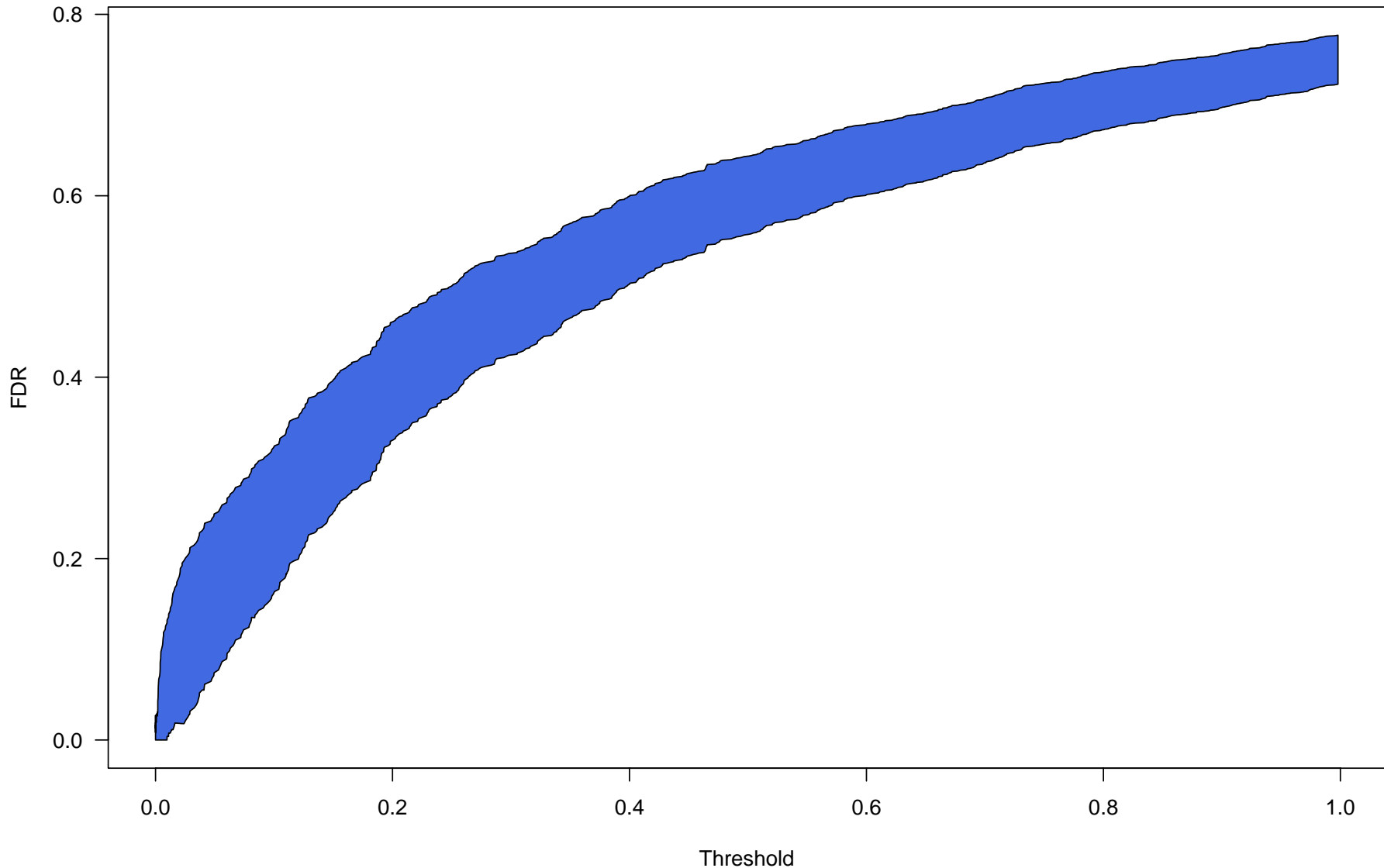
$$T_C = \sup \{t : \text{FDR}(t; h^m, P^m) \leq c \text{ and } h^m : (h^m, P^m) \in \mathcal{U}\}$$

is a $(1 - \alpha, c)$ confidence threshold procedure.

That is, $P_G\{\text{FDR}(T_C) \leq c\} \geq 1 - \alpha$.

Exact Confidence Thresholds (cont'd)

\mathcal{U} yields a confidence envelope for $FDR(t)$ sample paths.



Take-Home Points

- Realized versus Expected FDR
- Considering both FDR and FNR yields greater power
- Multiple testing problem is transformed to an estimation problem.
- Must control FDR and FNR as stochastic processes.

In general, the threshold and the FDR are coupled, and these correlations can have a large effect.