

A Large-Sample Approach to Controlling the False Discovery Rate

Christopher R. Genovese

Department of Statistics
Carnegie Mellon University

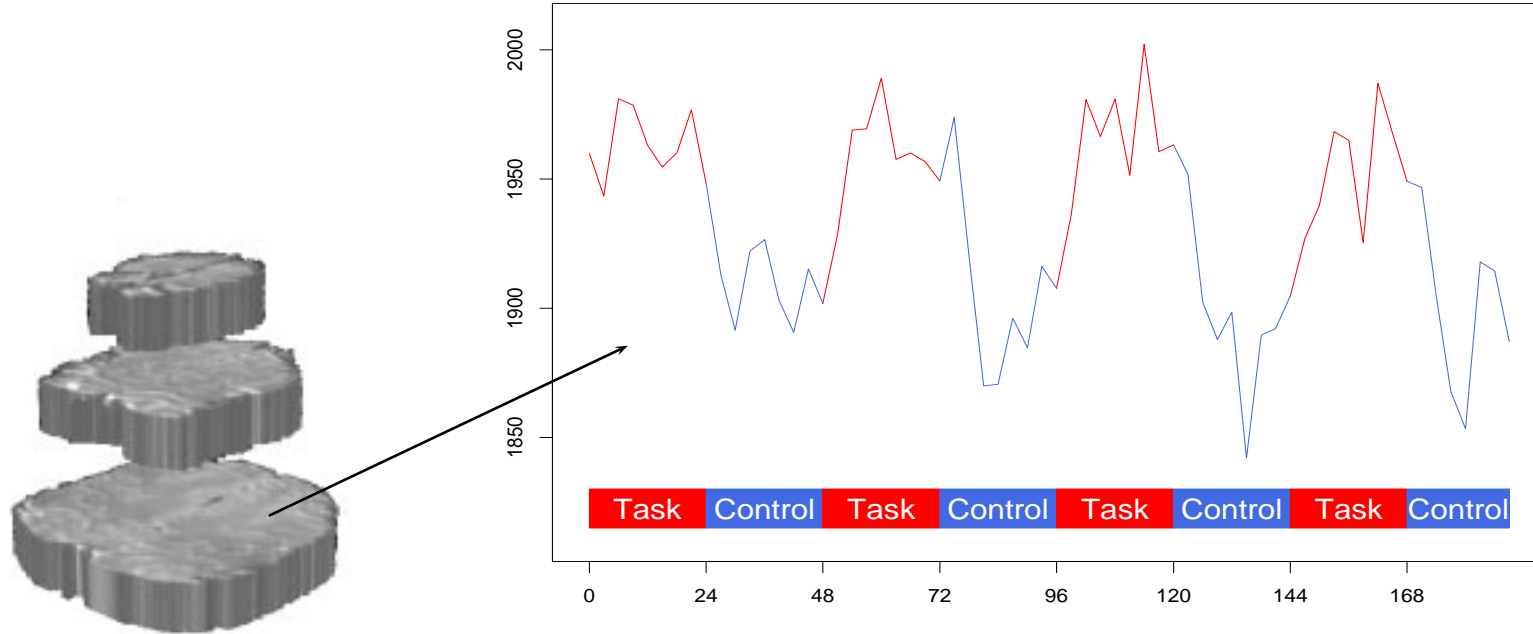
Larry Wasserman

Department of Statistics
Carnegie Mellon University

This work partially supported by NSF Grant SES 9866147.

Motivating Example #1: fMRI

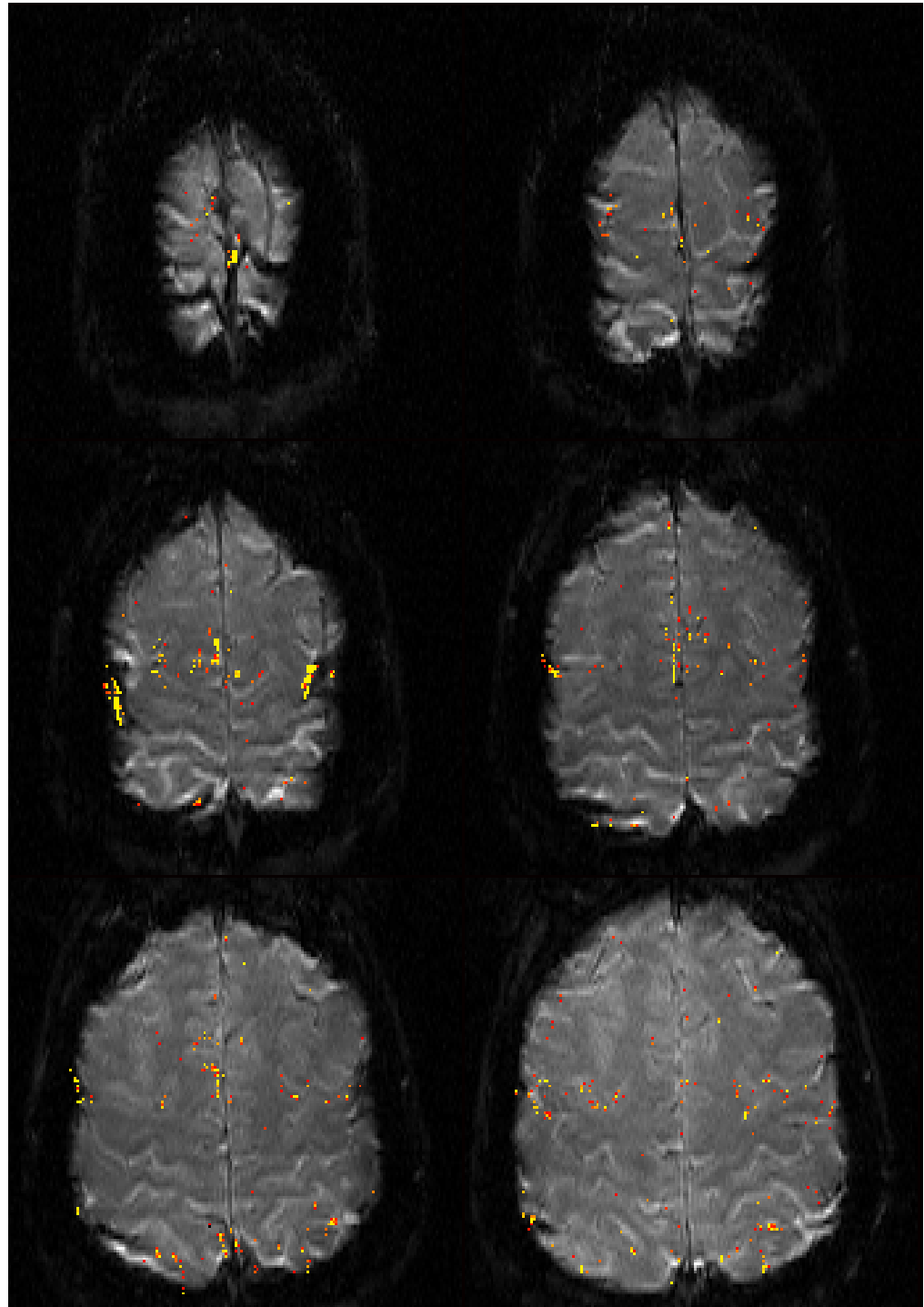
- fMRI Data: Time series of 3-d images acquired while subject performs specified tasks.



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

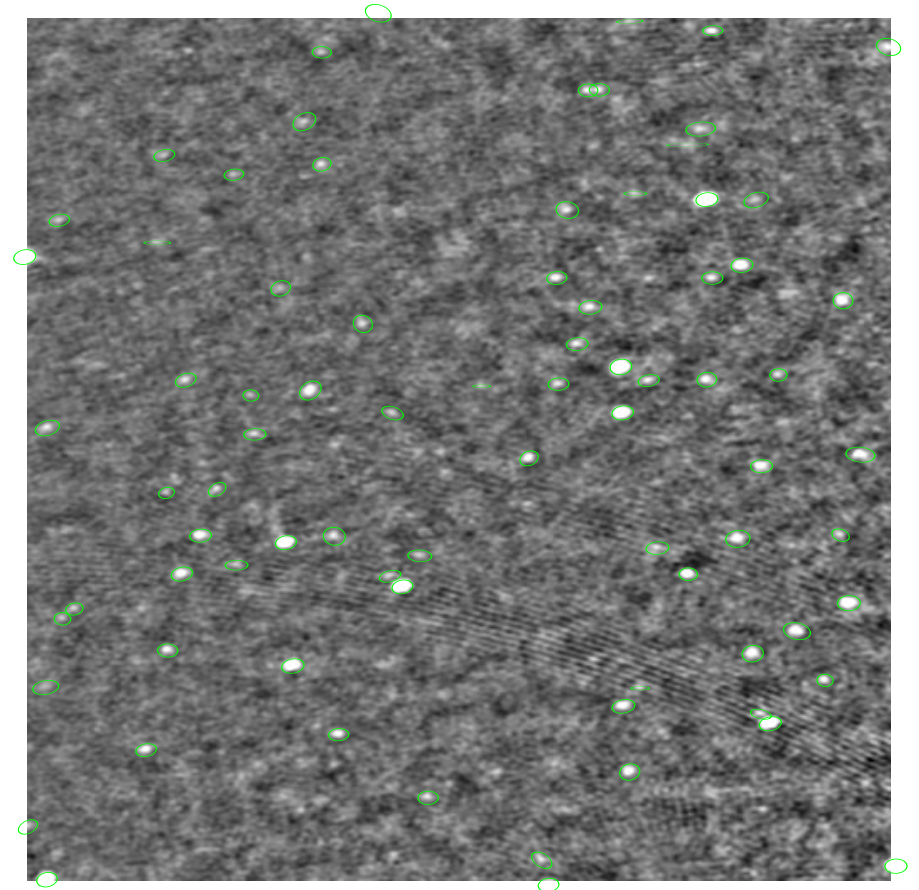
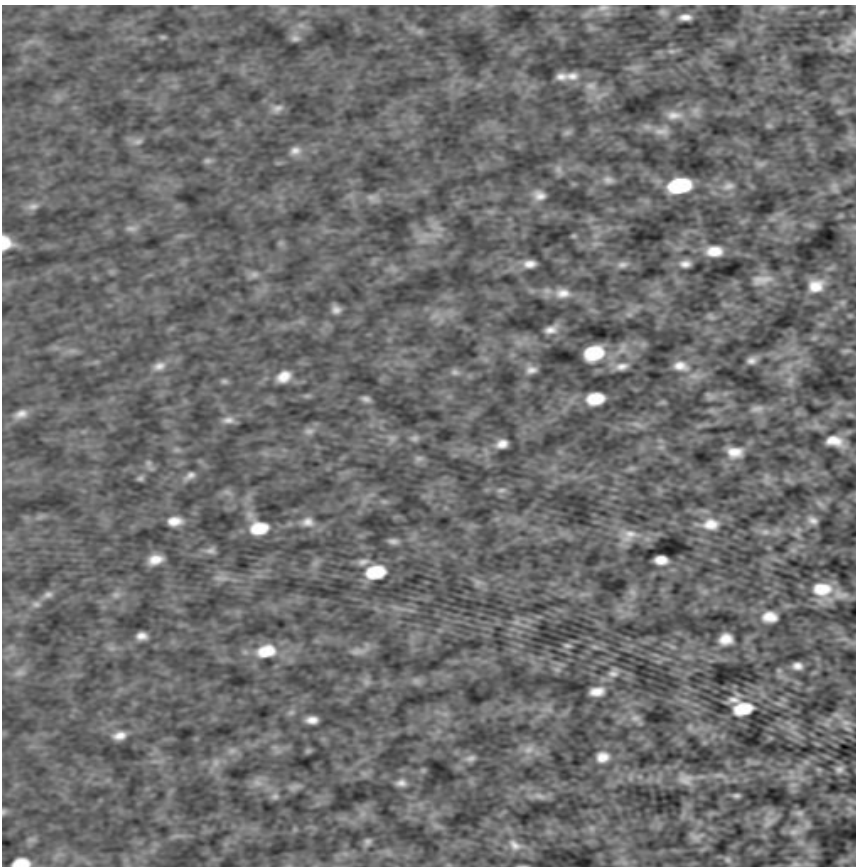
fMRI (cont'd)

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.

Recent Work on FDR

Abromovich, et al. (2000)

Benjamini & Hochberg (1995)

Benjamini & Liu (1999)

Benjamini & Hochberg (2000)

Benjamini & Yekutieli (2001)

Efron, et al. (2001)

Finner and Roters (2001, 2002)

Genovese & Wasserman (2001,2002)

Sarkar (2002)

Storey (2001,2002)

Storey & Tibshirani (2001)

Tusher, Tibshirani, Chu (2001)

The Multiple Testing Problem

- Perform m simultaneous hypothesis tests.

Classify results as follows:

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

Here, $M_{i|j}$ is the number of H_i chosen when H_j true.

Only R and m are observed.

False Discovery and Nondiscovery Proportions

- Define the False Discovery Proportion (FDP) and the False Nondiscovery Proportion (FNP) as follows:

$$\text{FDP} = \begin{cases} \frac{M_{1|0}}{R} & \text{if } R > 0, \\ 0, & \text{if } R = 0. \end{cases} \quad \text{FNP} = \begin{cases} \frac{M_{0|1}}{m - R} & \text{if } R < m, \\ 0, & \text{if } R = m. \end{cases}$$

- Then, the False Discovery Rate (FDR) and the False Nondiscovery Rate (FNR) are given by

$$\text{FDR} = E(\text{FDP}) \quad \text{FNR} = E(\text{FNP}).$$

Road Map

1. Preliminaries

- Models for FDP and FNP
- FDP and FNP as stochastic processes

2. Plug-in Procedures

- Asymptotic behavior of BH procedure
- Optimal Thresholds

3. Confidence Thresholds

- Controlling probability of exceeding specified proportion of false discoveries

4. Estimating the p -value distribution

Basic Models

- Let $P^m = (P_1, \dots, P_m)$ be the p-values for the m tests.
- Let $H^m = (H_1, \dots, H_m)$ where $H_i = 0$ (or 1) if the i^{th} null hypothesis is true (or false).
- We assume the following model:

$$H_1, \dots, H_m \text{ iid Bernoulli}\langle a \rangle$$

$$\Xi_1, \dots, \Xi_m \text{ iid } \mathcal{L}_{\mathcal{F}}$$

$$P_i \mid H_i = 0, \Xi_i = \xi_i \sim \text{Uniform}\langle 0, 1 \rangle$$

$$P_i \mid H_i = 1, \Xi_i = \xi_i \sim \xi_i.$$

where $\mathcal{L}_{\mathcal{F}}$ denotes a probability distribution on a class \mathcal{F} of distributions on $[0, 1]$.

Basic Models (cont'd)

- Marginally, P_1, \dots, P_m are drawn iid from

$$G = (1 - a)U + aF,$$

where U is the Uniform $\langle 0, 1 \rangle$ cdf and

$$F = \int \xi d\mathcal{L}_{\mathcal{F}}(\xi).$$

- Typical examples:
 - Parametric family: $\mathcal{F}_{\Theta} = \{F_{\theta} : \theta \in \Theta\}$
 - Concave, continuous distributions

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

- Can also work under what we call the *conditional model* where H_1, \dots, H_m are fixed, unknown.

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 $P_i \leq T(P^m)$. Often call T a *threshold*.
- Examples:
 - Uncorrected testing $T_U(P^m) = \alpha$
 - Bonferroni $T_B(P^m) = \alpha/m$
 - Fixed threshold at t $T_t(P^m) = t$
 - First r $T_{(r)}(P^m) = P_{(r)}$
 - Benjamini-Hochberg $T_{BH}(P^m) = \sup\{t: \hat{G}(t) = t/\alpha\}$
 - Oracle $T_O(P^m) = \sup\{t: G(t) = (1 - a)t/\alpha\}$
 - Plug In $T_{PI}(P^m) = \sup\{t: \hat{G}(t) = (1 - \hat{a})t/\alpha\}$
 - Regression Classifier $T_{Reg}(P^m) = \sup\{t: \hat{P}\{H_1=1|P_1=t\} > 1/2\}$

FDP and FNP as Stochastic Processes

- Inherent difficulty: FDP, FNP, and a general threshold all depend on the same data.
- Define the FDP and FNP processes, respectively, by

$$\text{FDP}(t) \equiv \text{FDP}(t; P^m, H^m) = \frac{\sum_i 1\{P_i \leq t\} (1 - H_i)}{\sum_i 1\{P_i \leq t\} + 1\{\text{all } P_i > t\}}$$

$$\text{FNP}(t) \equiv \text{FNP}(t; P^m, H^m) = \frac{\sum_i 1\{P_i > t\} H_i}{\sum_i 1\{P_i > t\} + 1\{\text{all } P_i \leq t\}}.$$

- For procedure T , the FDP and FNP are obtained by evaluating these processes at $T(P^m)$.

FDP and FNP as Stochastic Processes (cont'd)

- Both these processes converge to Gaussian processes outside a neighborhood of 0 and 1 respectively.
- For example, define

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

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

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

- Then, $Z_m \rightsquigarrow Z$.

Plug-in Procedures

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

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

as Storey (2002) first noted.

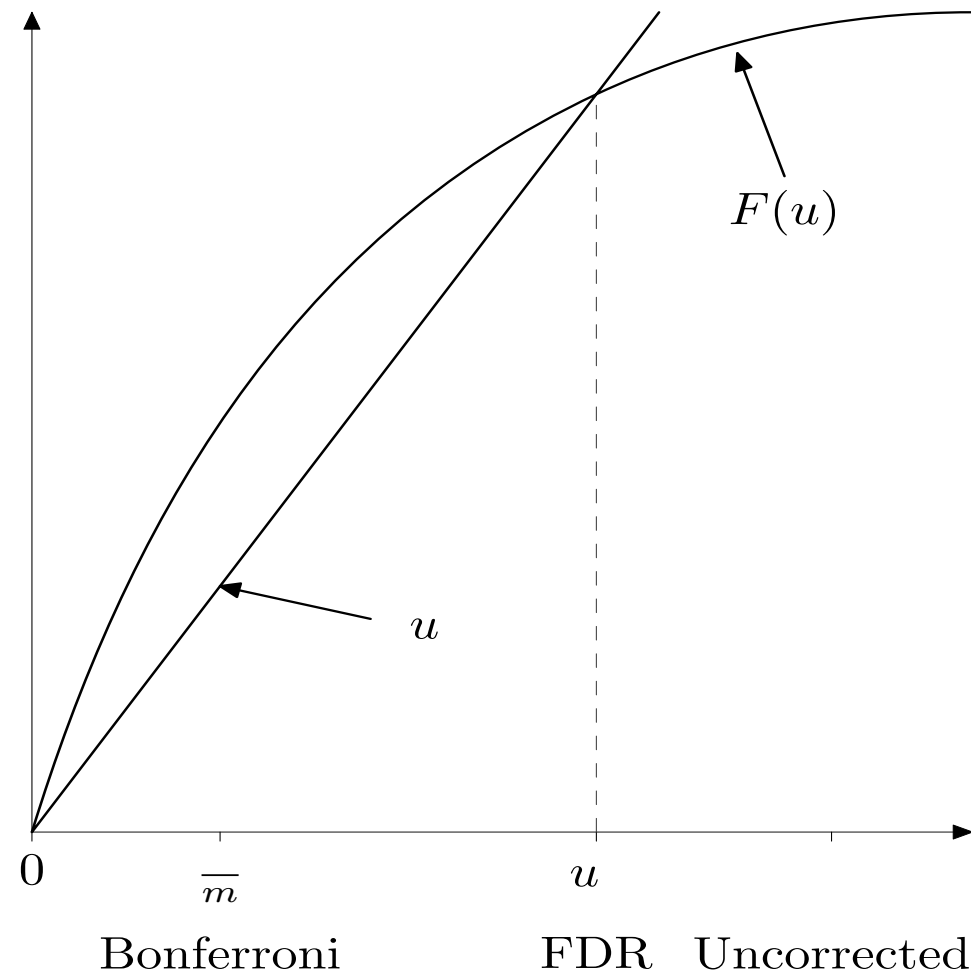
- One can think of this as a plug-in procedure for estimating

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

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,

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

- The BH procedure overcontrols FDR and thus will not in general minimize FNR.
- This suggests using T_{PI} , the plug-in estimator for

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

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

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

Optimal Thresholds (cont'd)

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

$$E(\text{FDP}(t)) = \frac{(1-a)t}{G(t)} + O((1-t)^m)$$

$$E(\text{FNP}(t)) = a \frac{1-F(t)}{1-G(t)} + O((a+(1-a)t)^m).$$

- Ignoring $O()$ terms and choosing t to minimize $E(\text{FNP}(t))$ subject to $E(\text{FDP}(t)) \leq \alpha$, yields $t^*(a, G)$ as the optimal threshold.
- GW (2002) show that

$$E(\text{FDP}(t^*(\hat{a}, \hat{G}))) \leq \alpha + O(m^{-1/2}).$$

Confidence Thresholds

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

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

for some specified c and α .

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

- Four approaches: (i) an asymptotic Bootstrap threshold, (ii) an asymptotic closed-form threshold, (iii) an exact (small-sample) threshold requiring numerical search, and (iv) a Bayesian threshold.
- Here, I'll discuss the case where α is known.

In general, all of this works using an estimator, but this introduces additional complexity.

Bootstrap Confidence Thresholds

- First guess: Choose T such that

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

- This fails. The problem is an additional bias term:

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

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

- Can fix this with double bootstrap (harder) or DKW correction (easier).

Bootstrap Confidence Thresholds (cont'd)

- Let $\beta = \alpha/2$ and $\epsilon_m \equiv \epsilon_m(\beta) = \sqrt{\frac{1}{2m} \log \left(\frac{2}{\beta} \right)}$.

- 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 $\Omega_c^*(t) = \sum_i 1\{P_i^* \leq t\} (1 - H_i^* - c)$.

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,

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

Closed-Form Asymptotic Confidence Thresholds

- Let

$$t_0 = Q^{-1}(c) \quad \hat{t}_0 = \hat{Q}^{-1}(c).$$

- Then define

$$T_C = \hat{t}_0 + \frac{\hat{\Delta}_{m,\alpha}}{\sqrt{m}},$$

where $\hat{\Delta}_{m,\alpha}$ is depends on a density estimate of $g = G'$.

- Then,

$$P_G\{ \text{FDP}(T_C) \leq c \} \geq 1 - \alpha + o(1).$$

Closed-Form Asymptotic Confidence Thresholds

- Details:

$$\hat{\Delta}_{m,\alpha} = \frac{z_{\alpha/2} \left(\sqrt{\hat{K}_{Q^{-1}}(\hat{t}_0, \hat{t}_0)} + \hat{g}(\hat{t}_0) \right) + 2\sqrt{\log m}}{1 - \hat{a} - c\hat{g}(\hat{t}_0)}$$

$$\hat{K}_{Q^{-1}}(s, t) = \frac{\hat{K}_Q(\hat{Q}^{-1}(s), \hat{Q}^{-1}(t))}{\hat{Q}'(\hat{Q}^{-1}(s))\hat{Q}'(\hat{Q}^{-1}(t))}$$

$$\hat{K}_Q(s, t) = \frac{(1 - \hat{a})^2 st}{\hat{G}^2(s)\hat{G}^2(t)} \left[\hat{G}(s \wedge t) - \hat{G}(s)\hat{G}(t) \right].$$

- This requires no bootstrapping but does require density estimation. This is analogous to the situation faced when estimating the standard error of a median.

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

$$S(t; h^m, p^m) = \frac{\sum_i \mathbf{1}\{p_i \leq t\} (1 - h_i)}{\sum_i (1 - h_i)} \quad [\text{EDF of null p-values}]$$

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

where $m_0 = \sum_i (1 - h_i)$ and $\epsilon_{m_0}(\beta) = \sqrt{\log(2/\beta)/2m_0}$.

- Take $\beta = 1 - \sqrt{1 - \alpha}$.

Exact Confidence Thresholds (cont'd)

- Let

$$T_C = \sup \left\{ t : \text{FDP}(t; h^m, P^m) \leq c \text{ and } h^m \in \mathcal{U}_\beta(P^m) \right\}$$
$$\mathcal{G} = \left\{ \text{FDP}(\cdot; h^m, P^m) : h^m \in \mathcal{U}_\beta(P^m) \right\}.$$

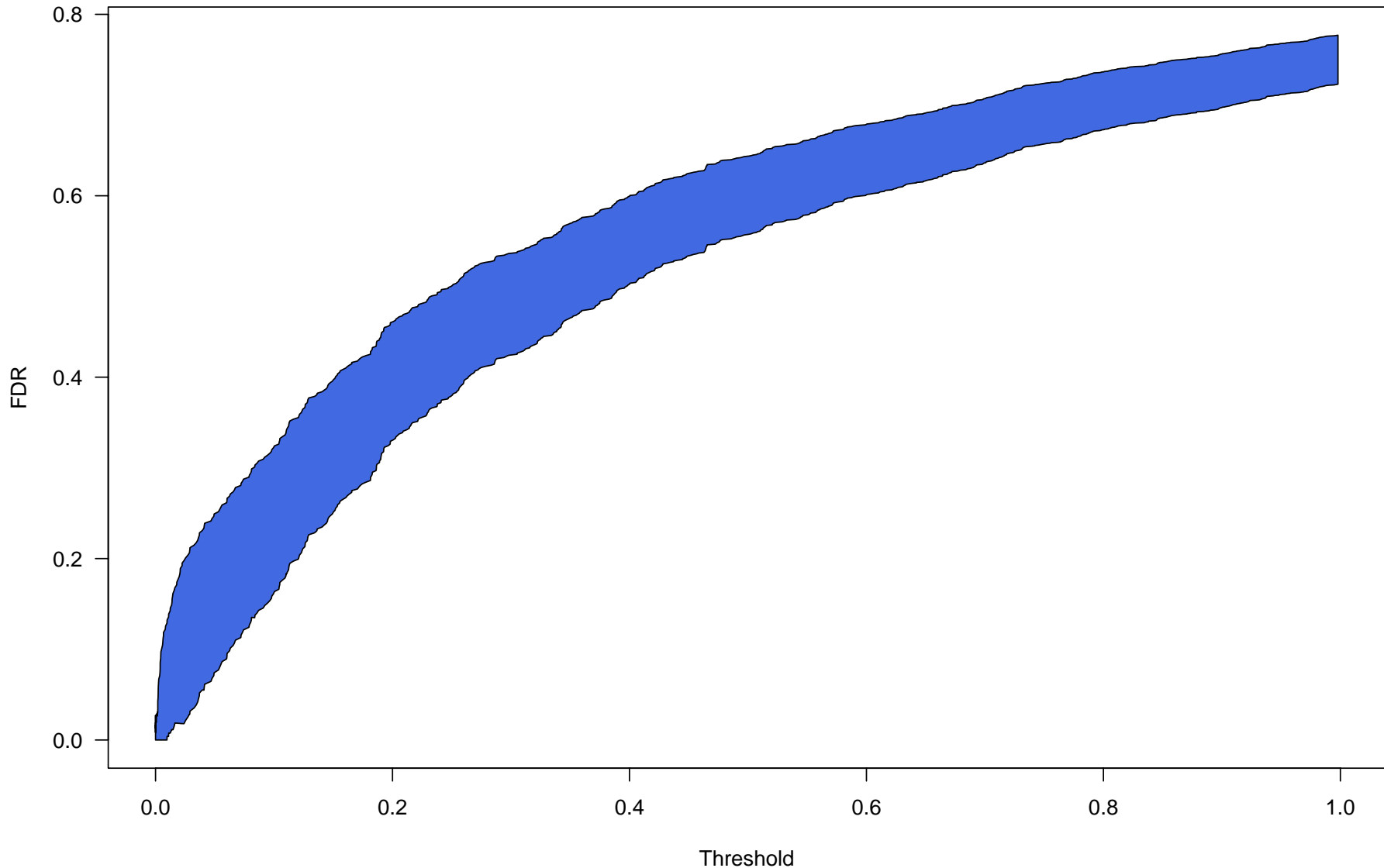
- Then,

$$\begin{aligned} P_G \left\{ H^m \in \mathcal{U}_\beta(P^m) \right\} &\geq 1 - \alpha, \\ P_G \left\{ \text{FDP}(\cdot; H^m, P^m) \in \mathcal{G} \right\} &\geq 1 - \alpha, \\ P_G \left\{ \text{FDP}(T_C) \leq c \right\} &\geq 1 - \alpha. \end{aligned}$$

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

Exact Confidence Thresholds (cont'd)

\mathcal{G} gives a confidence envelope for $FDP(t)$ sample paths.



Estimating a and F

- Recall that the p-value distribution $G = (1 - a)U + aF$ where a and F are unknown.
- We need a good estimate of a for plug-in estimates,

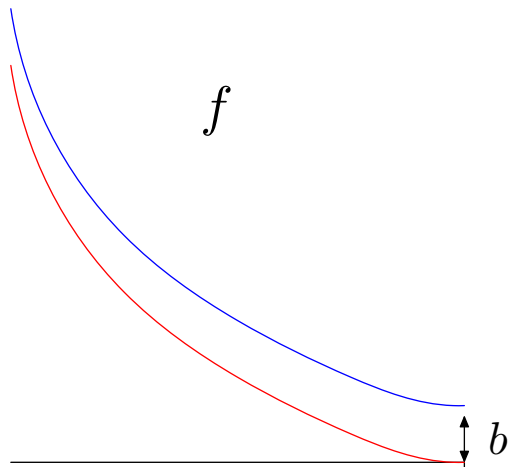
$$T_{\text{PI}}(P^m) = \max \left\{ t: \hat{G}(t) = \frac{(1 - \hat{a})t}{\alpha} \right\},$$

that approximate the optimal threshold.

- We need good estimates of a and F for confidence thresholds.

Estimating a and F (cont'd)

- Identifiability and Purity



If $\min f = b > 0$, can write $F = (1-b)U + bF_0$,
 $\mathcal{O}_G = \{(\tilde{a}, \tilde{F}) : \tilde{F} \in \mathcal{F}, G = (1 - \tilde{a})U + \tilde{a}\tilde{F}\}$
 may contain more than one element.

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

- In general, let $\underline{a} \leq a$ be the smallest mixing weight in the orbit:
 $\underline{a} = 1 - \min_t g(t)$. This is identifiable.

Storey (2002) notes that $0 \leq \sup_{0 < t < 1} \frac{G(t) - t}{1 - t} \leq \underline{a} \leq a \leq 1$.

- $a - \underline{a}$ is typically small: $a - \underline{a} = ae^{-n\theta^2/2}$ in the two-sided test of $\theta = 0$ versus $\theta \neq 0$ in the Normal $\langle \theta, 1 \rangle$ model.

Estimating a and F (cont'd)

- Parametric Case

- Derived a $1 - \beta$ one-sided conf. int. for \underline{a} and thus a .
 (a, θ) typically identifiable even if $a > \underline{a}$; use MLE.

- Non-parametric case:

- Derived a $1 - \beta$ one-sided conf. int. for \underline{a} and thus a .
- When F concave, get $\hat{a}_{\text{LCM}} = \underline{a} + O_P(m^{-1/3})$.
- When F smooth enough, get $\hat{a}_{\text{S}} = \underline{a} + O_P(m^{-2/5})$.
- Consistent estimate for F_0 if \hat{a} consistent for \underline{a} :

$$\hat{F}_m = \operatorname{argmin}_{H \in \mathcal{F}} \|\hat{G} - (1 - \hat{a})U - \hat{a}H\|_{\infty}.$$

Estimating a and F (cont'd)

- \hat{a}_S uses “spacings” estimator (Swanepoel, 1999) to estimate $\min g(t)$. This yields

$$\frac{m^{2/5}}{(\log m)^\delta} (\hat{a} - \underline{a}) \rightsquigarrow \text{Normal}\langle 0, (1 - \underline{a})^2 \rangle$$

- In the concave case, take $\hat{g} = G'_{LCM}$ and $\hat{a}_{LCM} = 1 - \hat{g}(1)$.
A $1 - \alpha$ confidence interval for a is

$$\hat{a}_{LCM} \pm 4q_\alpha |\hat{g}(1)|^{1/3} n^{-1/3}$$

where $P\{\operatorname{argmax}_h (W(h) - h^2) \geq q_\alpha\} = \alpha$ and W_h is a 2-sided Brownian motion tied down at 0.

Estimating a and F (cont'd)

- Confidence interval for a given by

$$\mathcal{A}_m = \left[\max_t \frac{\widehat{G}_m(t) - t - \epsilon_m(\alpha)}{1 - t}, 1 \right],$$

where \widehat{G}_m is EDF and $\epsilon_m(\alpha) = \sqrt{\log(2/\alpha)/2m}$.

Then,

$$1 - \alpha \leq \inf_{a, F} \mathbb{P}\{a \in \mathcal{A}_m\} \leq 1 - \alpha + R_m$$

where

$$R_m = \sum_j (-1)^j \frac{\alpha^{j^2}}{2^{j^2-1}} + O\left(\frac{(\log m)^2}{\sqrt{m}}\right)$$

Take-Home Points

- Asymptotic view motivated by particular applications, but asymptotics appear to kick in rather quickly.
- Confidence thresholds address a question that collaborating scientists frequently raise.
- Helpful to think of FDP (FDR) and FNP (FNR) as stochastic processes.

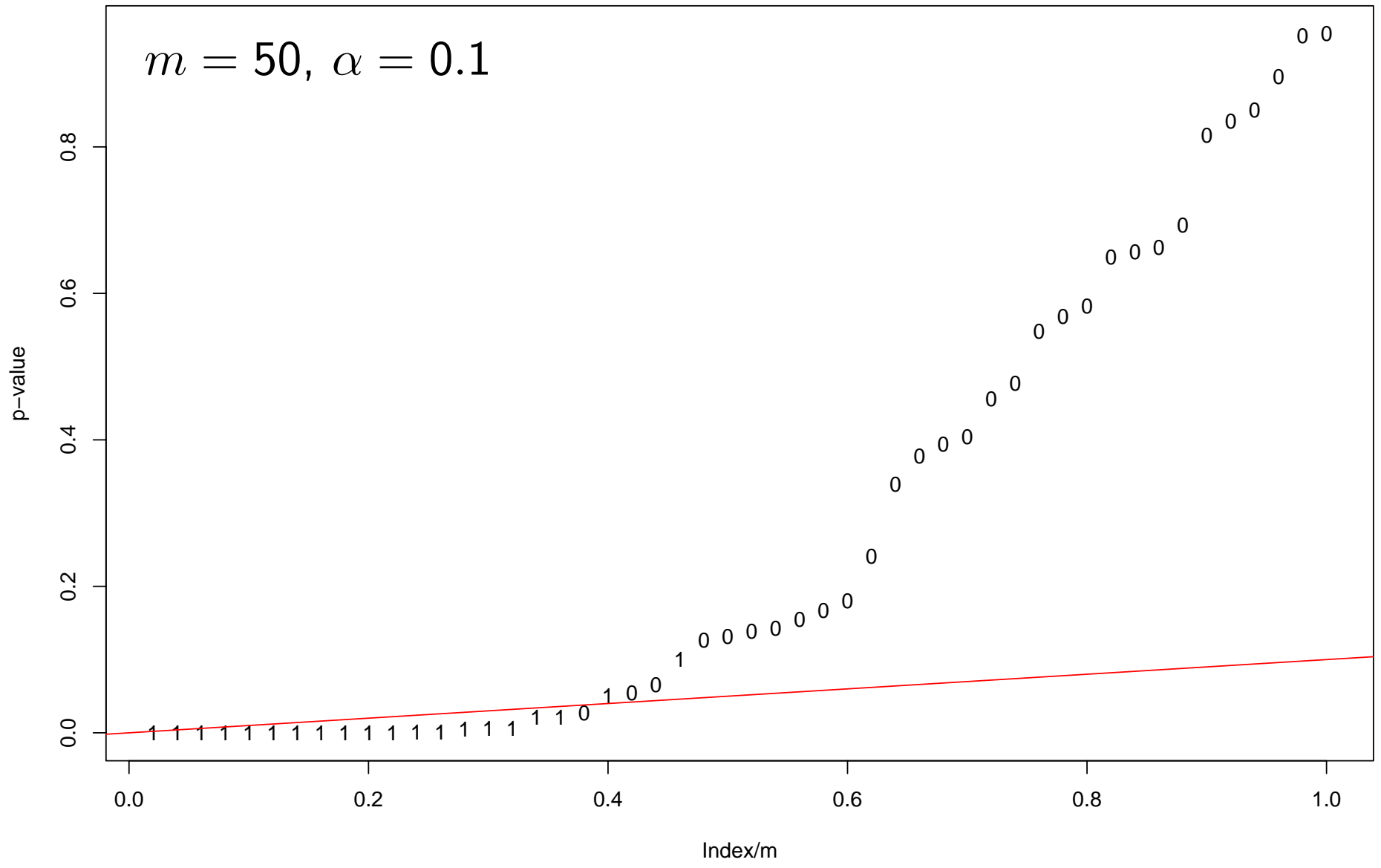
In general, the threshold and the FDP are coupled, and these correlations can have a large effect.

- Dependence

Recurring Notation

$m, M_0, M_{1 0}$	# of tests, true nulls, false discoveries
a	Mixture weight on <i>a</i> lternative
$H^m = (H_1, \dots, H_m)$	Unobserved true classifications
$P^m = (P_1, \dots, P_m)$	Observed p-values
U	CDF of Uniform $\langle 0, 1 \rangle$
F, f	Alternative CDF and density
$G = (1 - a)U + aF$	Marginal CDF of P_i
$g = G'$	Marginal density of P_i
\hat{G}_m	Estimate of G (e.g., empirical CDF of P^m)
$\epsilon_k(\beta) = \sqrt{\frac{1}{2k} \log \left(\frac{2}{\beta} \right)}$	DKW bound $1 - \beta$ quantile of $\ \hat{G}_k - G\ _\infty$

$m = 50, \alpha = 0.1$



Bayesian Thresholds

- Bayesian Threshold bounds posterior FDR:

$$T_{\text{Bayes}} = \sup\{t : E(\text{FDP}(t) \mid P^m) \leq \alpha\}$$

- Similarly, can construct a posterior (c, α) confidence threshold $T_{\text{Bayes},c}$ by

$$T_{\text{Bayes},c} = \sup\{t : P\{\text{FDP}(t) \leq c \mid P^m\} \leq \alpha\}$$

EBT (Empirical Bayes Testing)

- Efron et al (2001) note that

$$P\{H_i = 0 \mid P^m\} = \frac{(1 - a)}{g(P_i)} \equiv q(P_i)$$

- Reject whenever $q(p) \leq \alpha$?
- For a, f unknown, $f \geq 0$ implies that

$$a \geq 1 - \min_p g(p) \implies \hat{a} = 1 - \min_p \hat{g}(p).$$

- Then,
- $$\hat{q}(p) = \frac{1 - \hat{a}}{\hat{g}(p)} = \frac{\min_s \hat{g}(s)}{\hat{g}(p)}$$

EBT versus FDR

- If we reject when $P\{H_i = 0 \mid P^m\} \leq \alpha$,
how many errors are we making?
- Under weak conditions, can show that

$$q(t) \leq \alpha \text{ implies } Q(t) < \alpha$$

So EBT is conservative.

Behavior of \hat{q}

- THEOREM. Let $\hat{q}(t) = \frac{(1-a)}{\hat{g}(t)}$. Suppose that

$$m^\alpha(\hat{g}(t) - g(t)) \rightsquigarrow W$$

for some $\alpha > 0$, where W is a mean 0 Gaussian process with covariance kernel $\tau(v, w)$. Then

$$m^\alpha(\hat{q}(t) - q(t)) \rightsquigarrow Z$$

where Z is a Gaussian process with mean 0 and covariance kernel

$$K_q(v, w) = \frac{(1-a)^2 \tau(v, w)}{g(v)^4 g(w)^4}.$$

Behavior of \hat{q} (cont'd)

- Parametric Case: $g \equiv g_\theta = (1 - a) + af_\theta(v)$ Then,

$$\text{rel}(v) = \frac{\widehat{\text{se}}(\hat{q}(v))}{q(v)} \approx O\left(\frac{1}{\sqrt{m}}\right) \left| \frac{\partial \log g_\theta}{\partial d\theta} \right| = O\left(\frac{1}{\sqrt{m}}\right) |v - \theta| \quad \text{Normal case}$$

- Nonparametric Case

$$\hat{g}(t) = \frac{1}{m} \sum_{i=1}^m \frac{1}{h_m} K\left(\frac{t - P_i}{h_m}\right)$$

$h_m = cm^{-\beta}$ where $\beta > 1/5$ (undersmooth). Then

$$\text{rel}_v = \frac{c}{m^{(1-\beta)/2} \sqrt{g(v)}}.$$