

Bayes and Empirical Bayes Information (Learning from the experience of others)

Bradley Efron
Stanford University

A Jar of Biased Coins

- **Coin A** 45 flips, 18 heads: $\hat{p}_A = 18/45 (= .400)$
- **Other coins** $\hat{p}_B, \hat{p}_C, \hat{p}_D, \dots$
- Coins are different and all flips independent
- **Question** Can $\hat{p}_B, \hat{p}_C, \hat{p}_D, \dots$ help you estimate p_A ?
- **Answer** Yes!
(*Empirical Bayes* circa 1950: Robbins, Stein, Good, Turing)

Bayes Theorem

- *Prior density* $g(\mu)$ for unknown parameter μ
- *Sampling density* $f_\mu(x)$ for observed data x

Bayes Rule $g(\mu|x) = g(\mu)f_\mu(x)/f(x)$

where $f(x) = \int g(\mu)f_\mu(x) d\mu$, the *marginal density* of x

- **Empirical Bayes** Try to learn $g(\mu)$ from “other” data

The 18 Baseball Players

	Player	hits/AB	\hat{p} (“MLE”)	\tilde{p} (“James–Stein”)	p (“Truth”)
1.	Clemente	18/45	.400	.290	.346
2.	Robinson	17/45	.378	.286	.298
3.	Howard	16/45	.356	.281	.276
⋮	⋮	⋮	⋮	⋮	⋮
17.	Munson	8/45	.178	.244	.316
18.	Alvis	7/45	.156	.239	.200
	<i>Grand Average</i>		.265	.265	.265

Total Prediction Error

- *MLE* $\sum_1^{18} (p_i - \hat{p}_i)^2 = 0.075$
- *James–Stein* $\sum_1^{18} (p_i - \tilde{p}_i)^2 = 0.021$
- *Ratio* 3.5
- No fluke!

The James–Stein Rule (1961)

- *Observe* $x_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mu_i, \sigma_0^2)$ for $i = 1, 2, \dots, N \geq 4$,
i.e., independent normals, with $f_{\mu_i}(x_i) = \left(\frac{1}{\sqrt{2\pi\sigma_0^2}} \right) \exp \left\{ -\frac{1}{2} \left(\frac{x_i - \mu_i}{\sigma_0} \right)^2 \right\}$
- **James–Stein** $\bar{x} = \frac{\sum x_i}{N}$ $\hat{V} = \frac{\sum (x_i - \bar{x})^2}{N - 3}$
- *MLE* $\hat{\mu}_i = x_i$ • *JS* $\tilde{\mu}_i = \bar{x} + \left[1 - \sigma_0^2 / \hat{V} \right] (x_i - \bar{x})$

Theorem *JS always beats MLE in expected squared error:*

$$E \left\{ \sum (\mu_i - \tilde{\mu}_i)^2 \right\} < E \left\{ \sum (\mu_i - \hat{\mu}_i)^2 \right\}$$

Bayesian Justification

- If *a priori* $\mu_i \stackrel{\text{ind}}{\sim} \mathcal{N}(M, A)$ and $x_i | \mu_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mu_i, \sigma_0^2)$:
then $x_i \sim \mathcal{N}(M, V)$ with $V = A + \sigma_0^2$ and

$$E\{\mu_i | x_i\} = M + \left[1 - \sigma_0^2 / V\right] (x_i - M)$$

- **James–Stein**: \bar{x} unbiased for M , and $1/\hat{V}$ unbiased for $1/V$
- “Empirical Bayes”

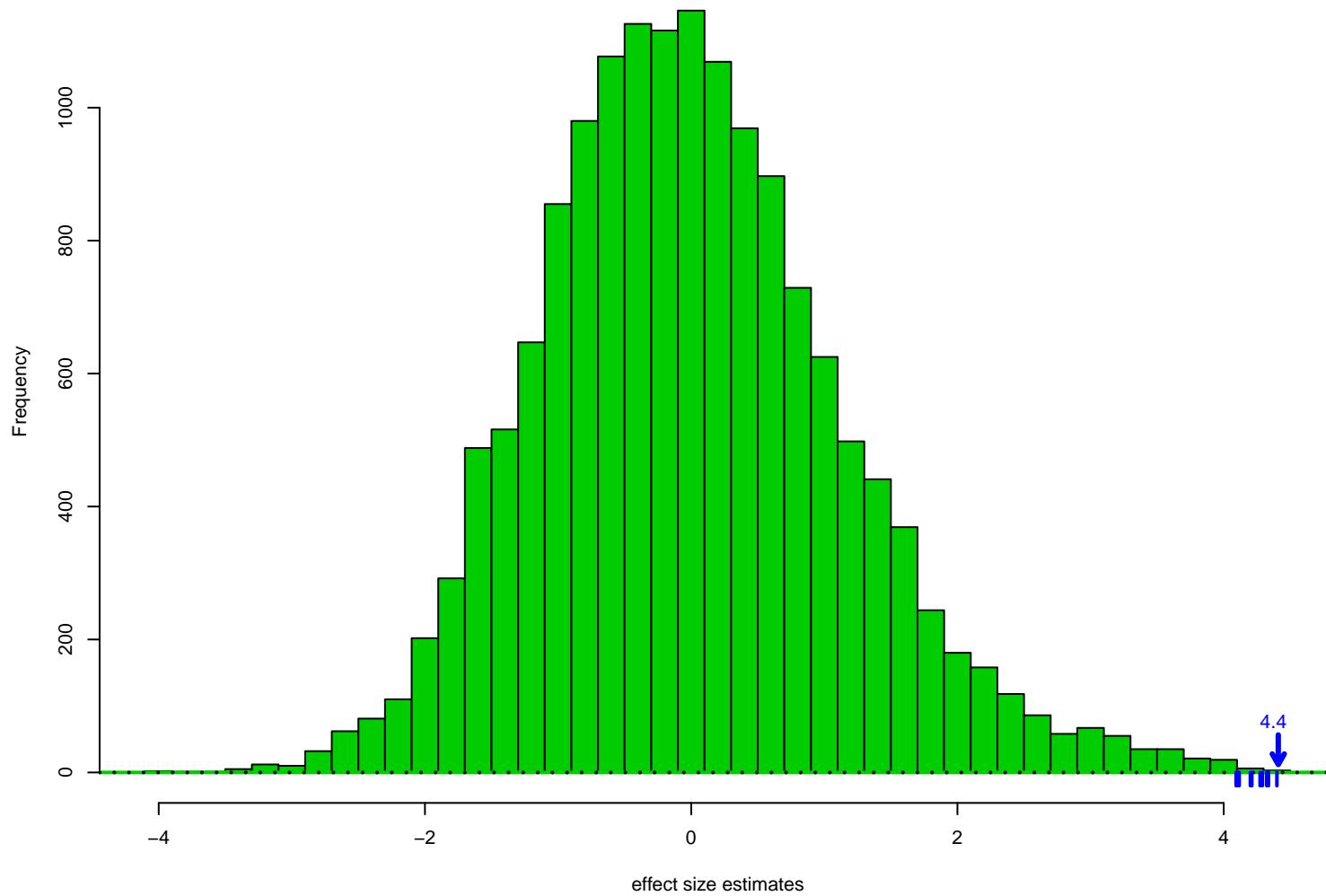
A Diffusion Tensor Imaging Study

- 6 dyslexic children and 6 normal controls
- Each measured at $N = 15443$ brain locations (*voxels*)
- x_i compares dyslexics with controls at voxel i :

$$x_i \sim \mathcal{N}(\mu_i, \sigma_0^2)$$

- μ_i the “effect size” • $\sigma_0 = 1.06$
- x_i 's range from -4.0 to 4.4

DTI Study comparing 6 dyslexic with 6 control children;
effect size estimates for 15443 brain voxels.
Dashes are 10 largest x's



The Winner's Curse

(“Selection bias,” “Regression to the mean”)

- *10 largest x_i 's* 4.40, 4.32, 4.30, ..., 4.09
- **Question** What are the 10 corresponding effect sizes μ_i ?
- **Curse** μ_i 's usually (much) smaller than x_i 's
- Bayes estimates are immune to Winner's Curse
- Empirical Bayes?

Tweedie's Formula (1956)

- *Bayes model* $\mu \sim g(\cdot)$ *a priori* and $x|\mu \sim \mathcal{N}(\mu, \sigma_0^2)$

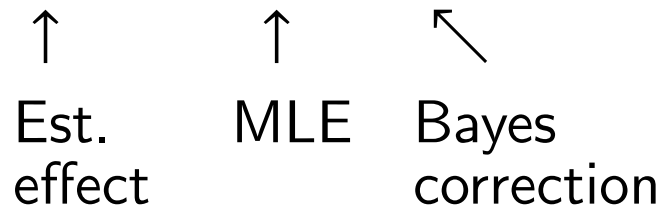
- *Marginal density* $f(x) = \int_{-\infty}^{\infty} f_0(x - \mu)g(\mu) d\mu$

with $f_0(x)$ the $\mathcal{N}(0, \sigma_0^2)$ density

- Tweedie's Formula

$$E\{\mu|x\} = x + \sigma_0^2 l'(x)$$

- $l'(x) = \frac{d}{dx} \log f(x)$



- *Advantage* works directly with f , not g

Empirical Bayes Implementation

- We don't know prior $g(\mu)$ or mixture $f(x)$ but can estimate $\hat{f}(x)$ and $\hat{l} = \log \hat{f}(x)$ from histogram of all N x 's:

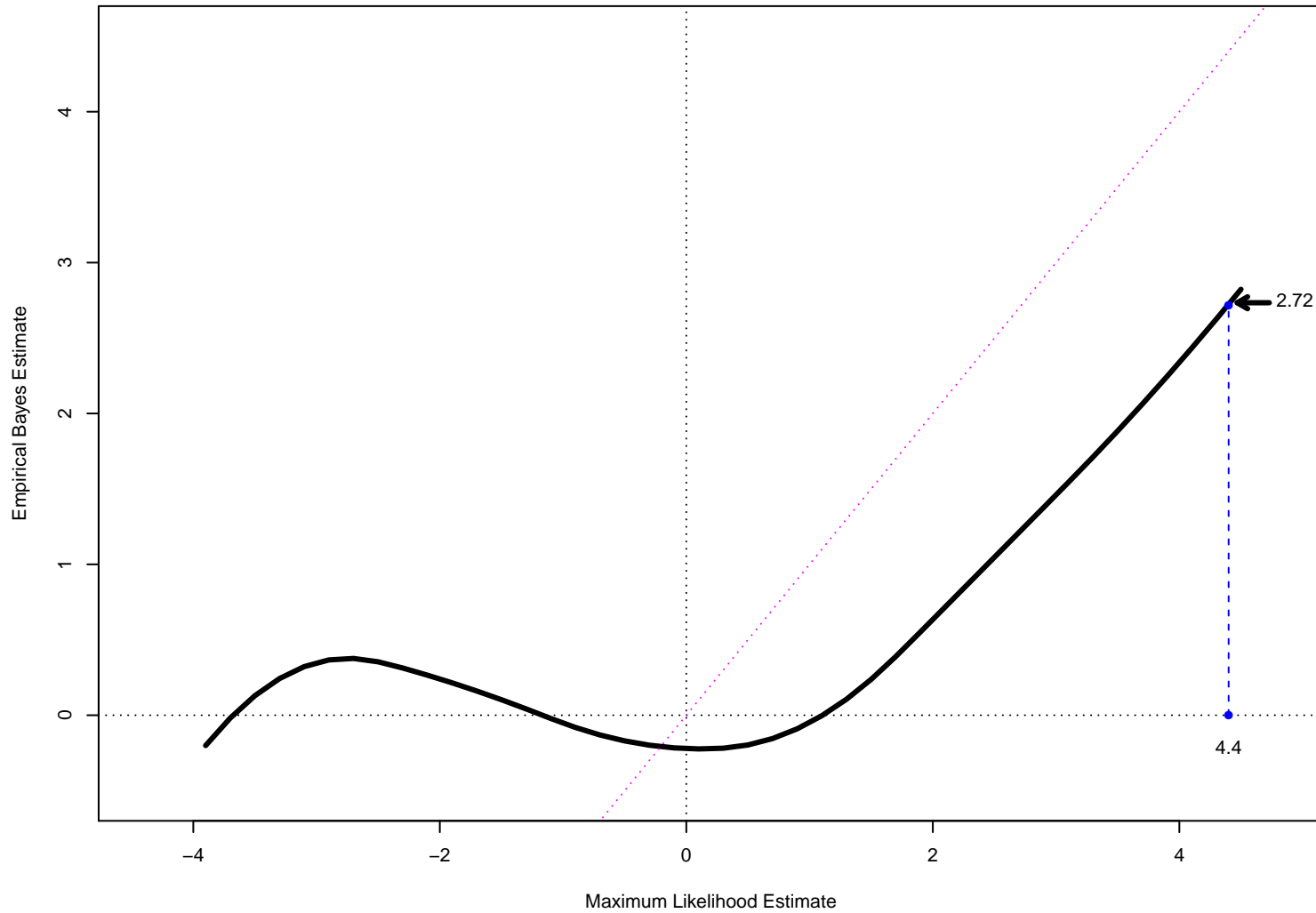
$$\hat{E}\{\mu_i|x_i\} = x_i + \sigma_0^2 \hat{l}'(x_i)$$

- *Parametric model* $\log f(x) = \sum_{j=0}^J \beta_j x^j$ [$J = 2$: normal]

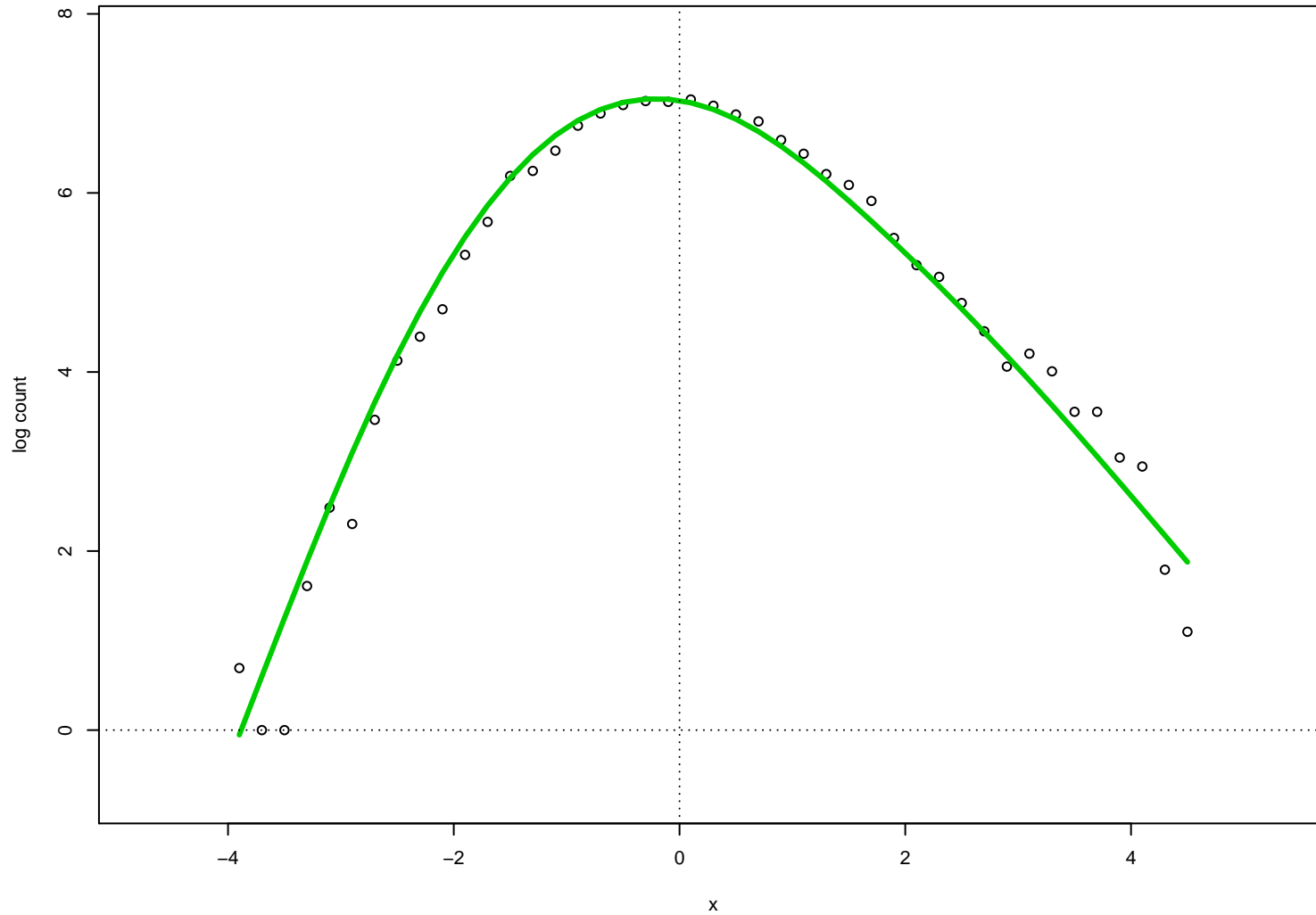
- *MLE* $\hat{l}(x) = \sum_{j=0}^J \hat{\beta}_j x^j$ so $\hat{l}'(x) = \sum_{j=1}^J j \hat{\beta}_j x^{j-1}$

[$J = 2$: James–Stein]

Empirical Bayes estimation curve for the DTI data
(Using Tweedie's formula, $J=6$)



Estimating $\log\{f(x)\}$ from the DTI data; points are log bin counts from histogram; curve is J=6 degree polynomial



Does “Tweedie” Lift the Winner’s Curse?

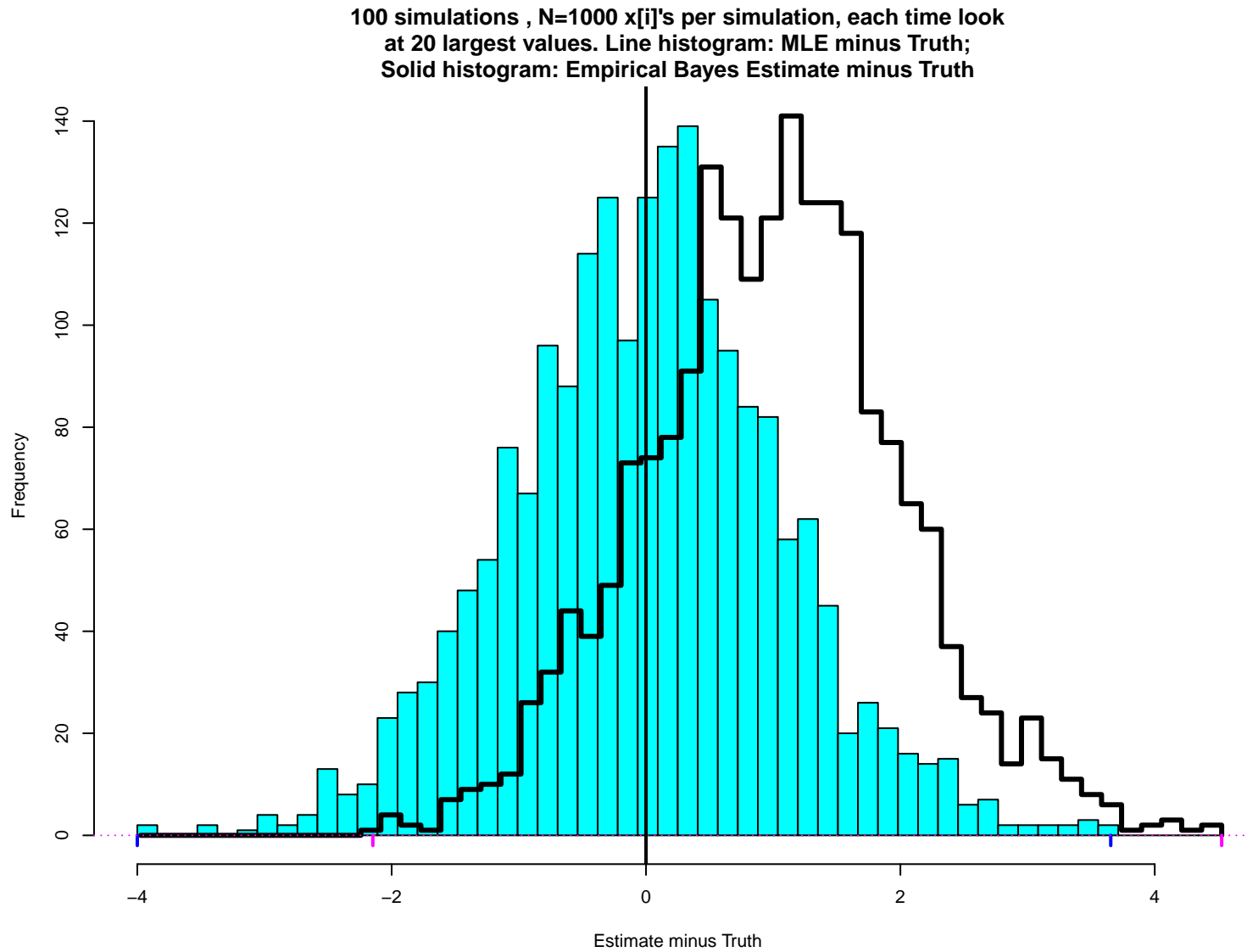
- *Simulations* prior $g(\mu) = e^{-\mu}$ [$\mu > 0$] and $x|\mu \sim \mathcal{N}(\mu, 1)$:

$$\{(\mu_i, x_i), i = 1, 2, \dots, N = 1000\}$$

- Compute Tweedie estimates $\{\tilde{\mu}_i\}$
- Repeat 100 times, each time comparing

$$x_i - \mu_i \quad \text{with} \quad \tilde{\mu}_i - \mu_i$$

for 20 largest x_i values



Empirical Bayes Regret

- For a given observed value x_0 :

true Bayes estimate $\mu^\dagger = x_0 + l'(x_0)$

must be more accurate than EB estimate $\tilde{\mu} = x_0 + \hat{l}'_x(x_0)$

- $\text{Regret}(x_0) = E \left\{ (\mu - \tilde{\mu})^2 - (\mu - \mu^\dagger)^2 \mid x_0 \right\}$

$$= E \left[\hat{l}'_x(x_0) - l'(x_0) \right]^2$$

- *Asymptotically* $\text{Reg}(x_0) \doteq \frac{c(x_0)}{N}$ $\left[\hat{l}'_x \text{ method determines } c(x_0) \right]$

Empirical Bayes Information

- $N = 1$: $\hat{\mu}_i = x_i$ (MLE)
- $N = \infty$: $\mu_i^\dagger = x_i - l'(x_i)$ (Bayes)
- N : $\tilde{\mu}_i = x_i - \hat{l}'_{\mathbf{x}}(x_i)$ (EB)

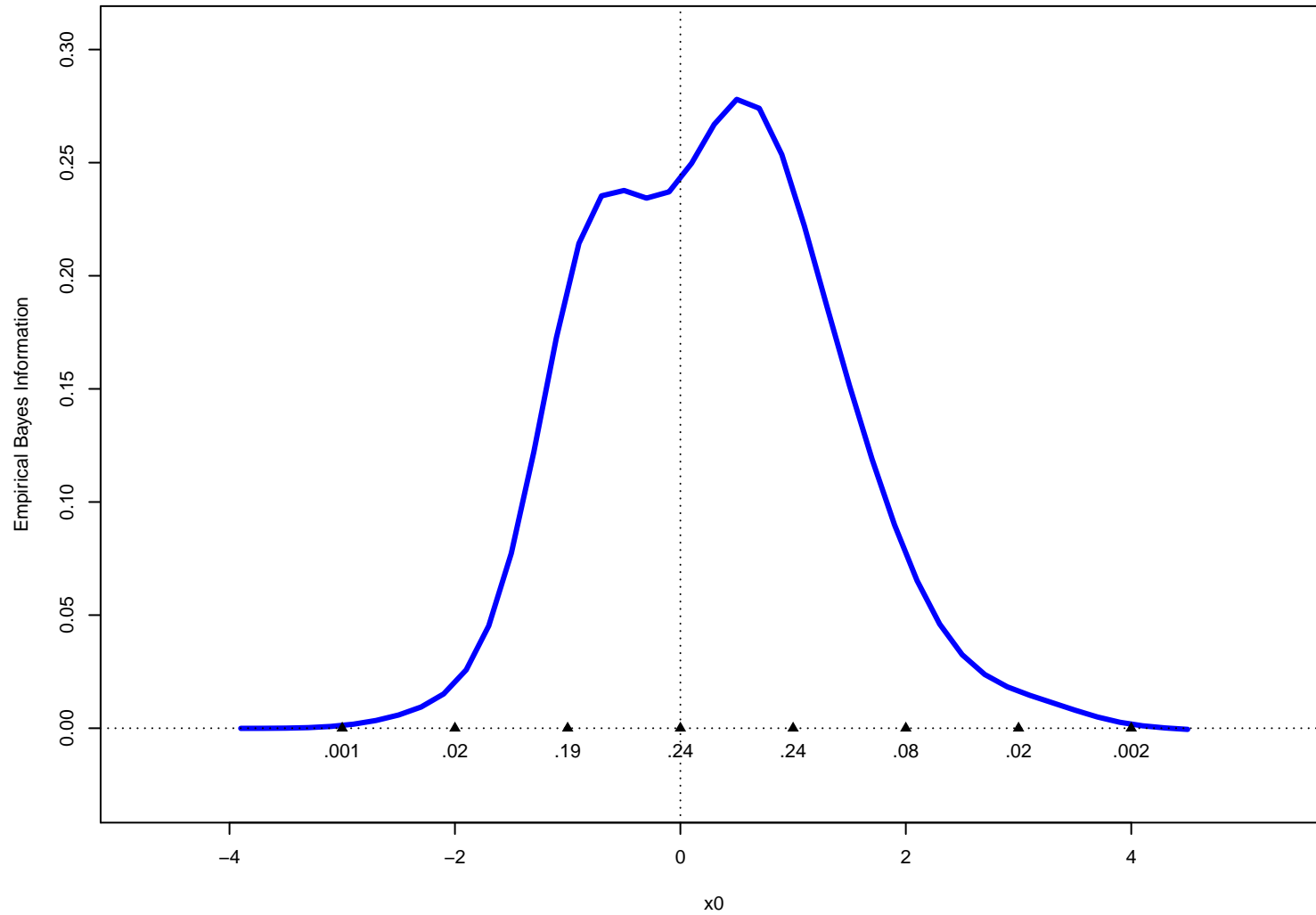
- Empirical Bayes information

$$\mathcal{I}(x_0) = \frac{1}{c(x_0)}$$

(amount of info for estimating $E\{\mu|x_0\}$ per other observation)

- *Bayes regret* $\text{Reg}(x_0) \doteq \frac{1}{N\mathcal{I}(x_0)}$

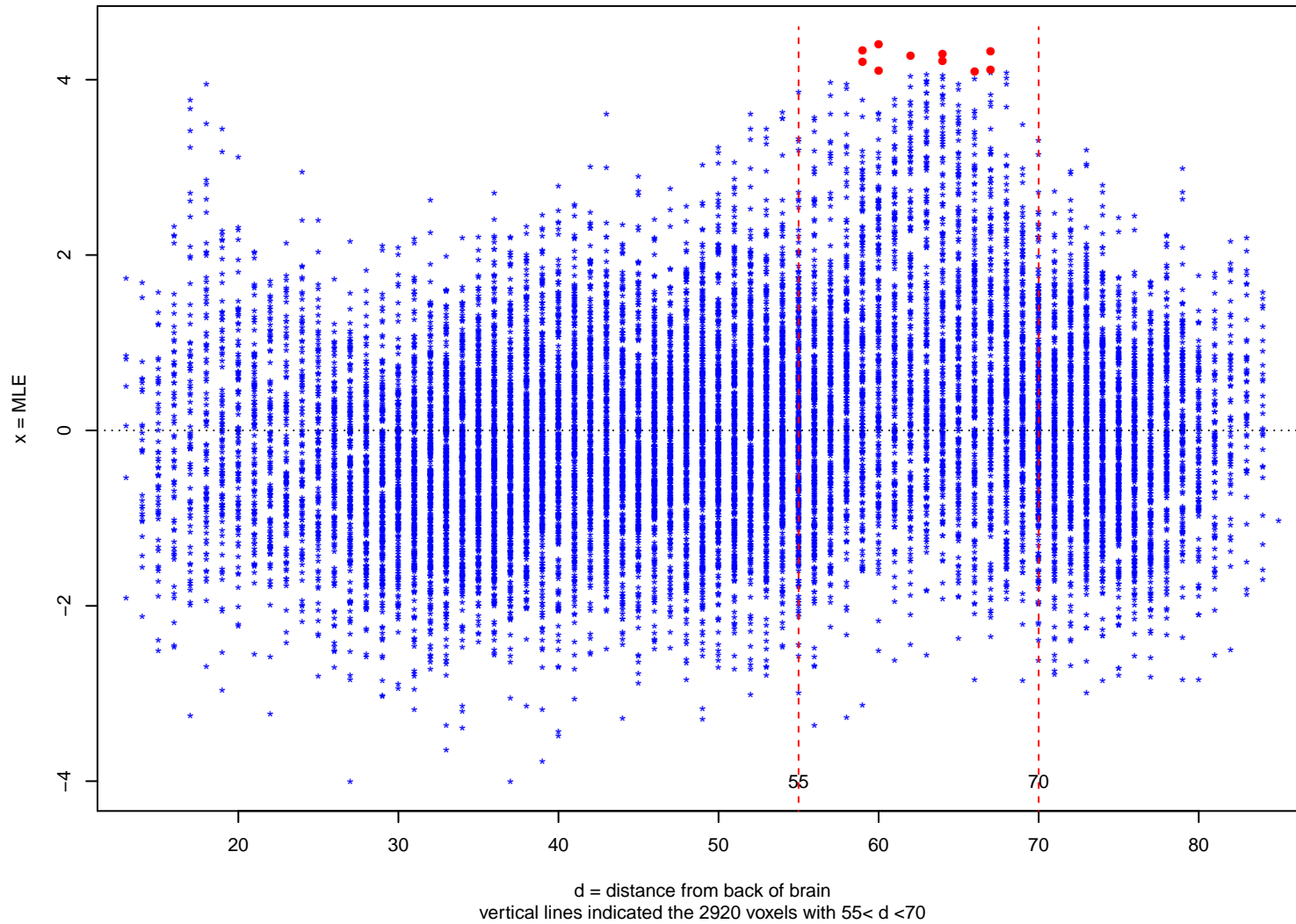
Empirical Bayes Information per other observation
for the DTI data



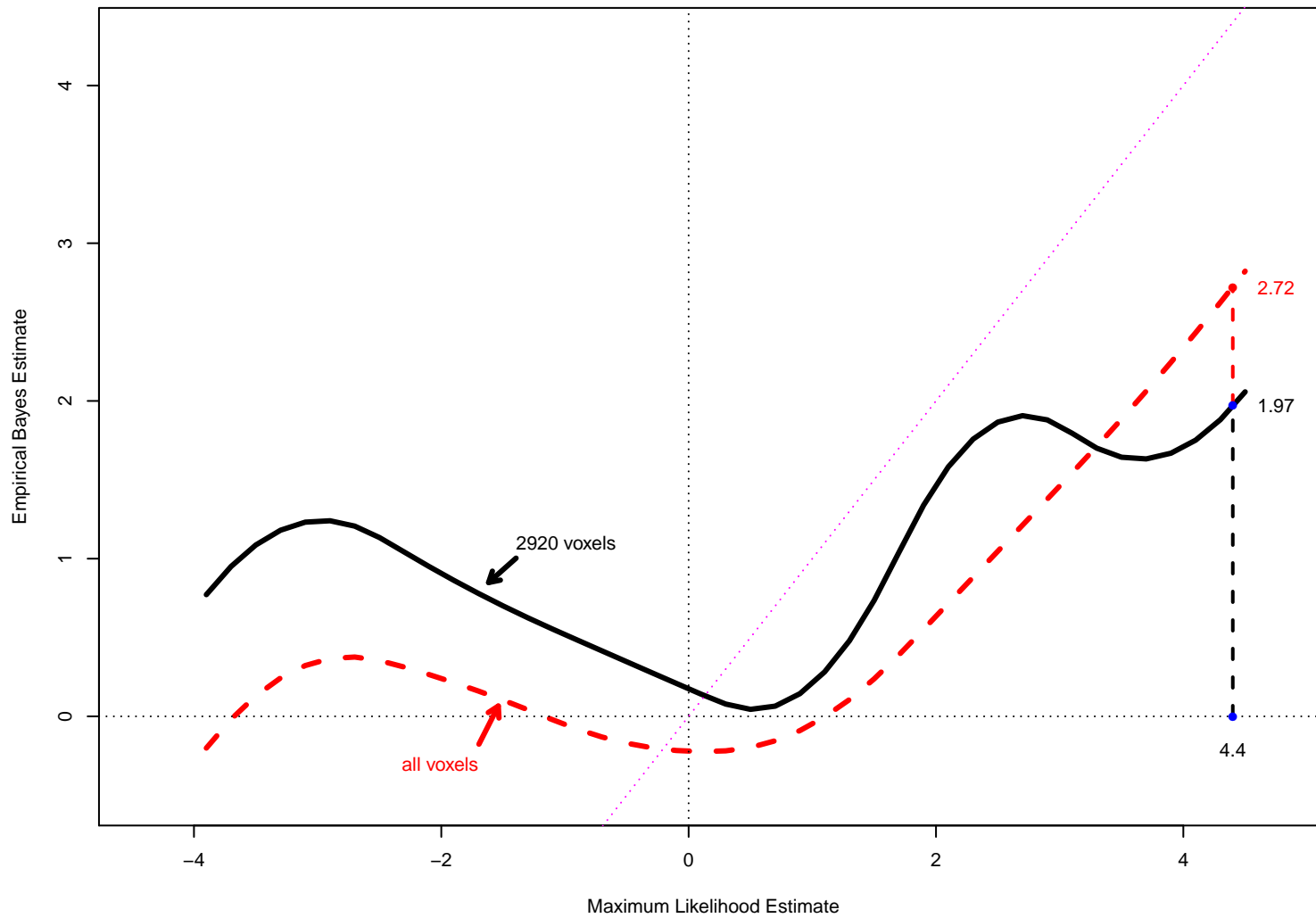
Learning from the Experience of Others

- *Which “others” ?*
 - All the other baseball players?
 - All the other voxels?
- *Next* DTI effect size MLEs x_i plotted versus d_i , distance from back of brain

MLE estimates x vs distance d from back of brain;
Red dots show the 10 largest x values



Empirical Bayes estimates based on only the 2920 voxels having $55 < d < 70$ (Solid); Red curve shows original estimates.



Relevance Functions

- Let $\rho_0(i)$ be *relevance* of “other” case i to target case i_0 [e.g., $\rho_0(i) = e^{-|d_i-65|/10}$ for target case having $d = 65$]
- Extended Tweedie Formula

$$\tilde{\mu}_{i_0} = x_{i_0} + \hat{l}'(x_{i_0}) + \frac{d}{dx} \log \hat{R}_0(x_{i_0})$$

where $\hat{R}_0(x) = \hat{E}\{\rho_0(i)|x\}$.

False Discovery Rates (1995)

- *Tweedie* $\mu \sim g(\cdot)$ observe $x|\mu \sim \mathcal{N}(\mu, \sigma_0^2)$
- **Testing** $\pi_0 = \Pr\{\mu = 0\}$ [“null”; often π_0 near 1]
- Local false discovery rate = $\Pr\{\text{null} | x\}$

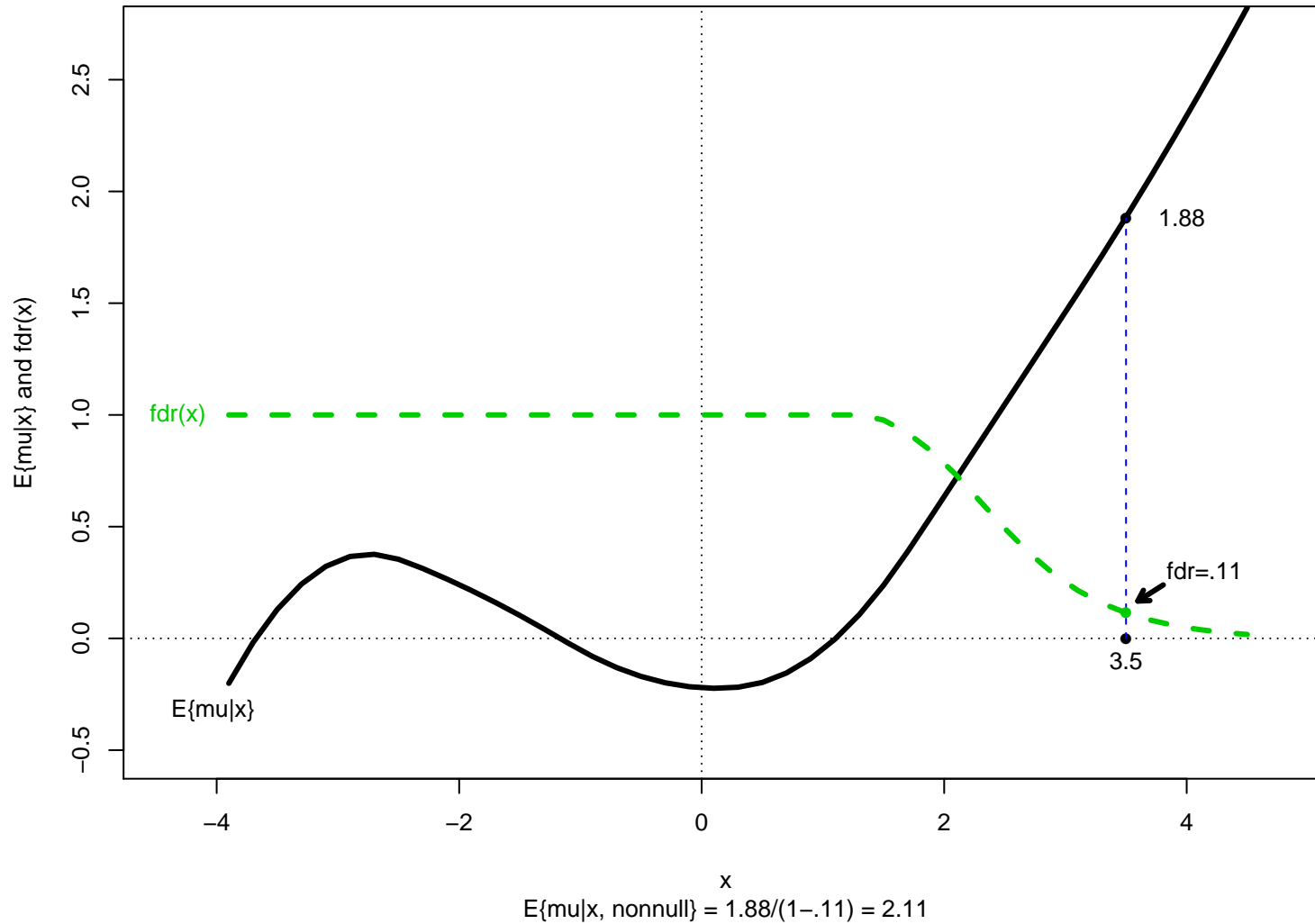
$$\text{fdr}(x) = \pi_0 f_0(x) / f(x)$$

- $f_0(x) = \frac{e^{-\frac{1}{2}\left(\frac{x}{\sigma_0}\right)^2}}{\sqrt{2\pi\sigma_0^2}}$ and $f(x) =$ mixture density

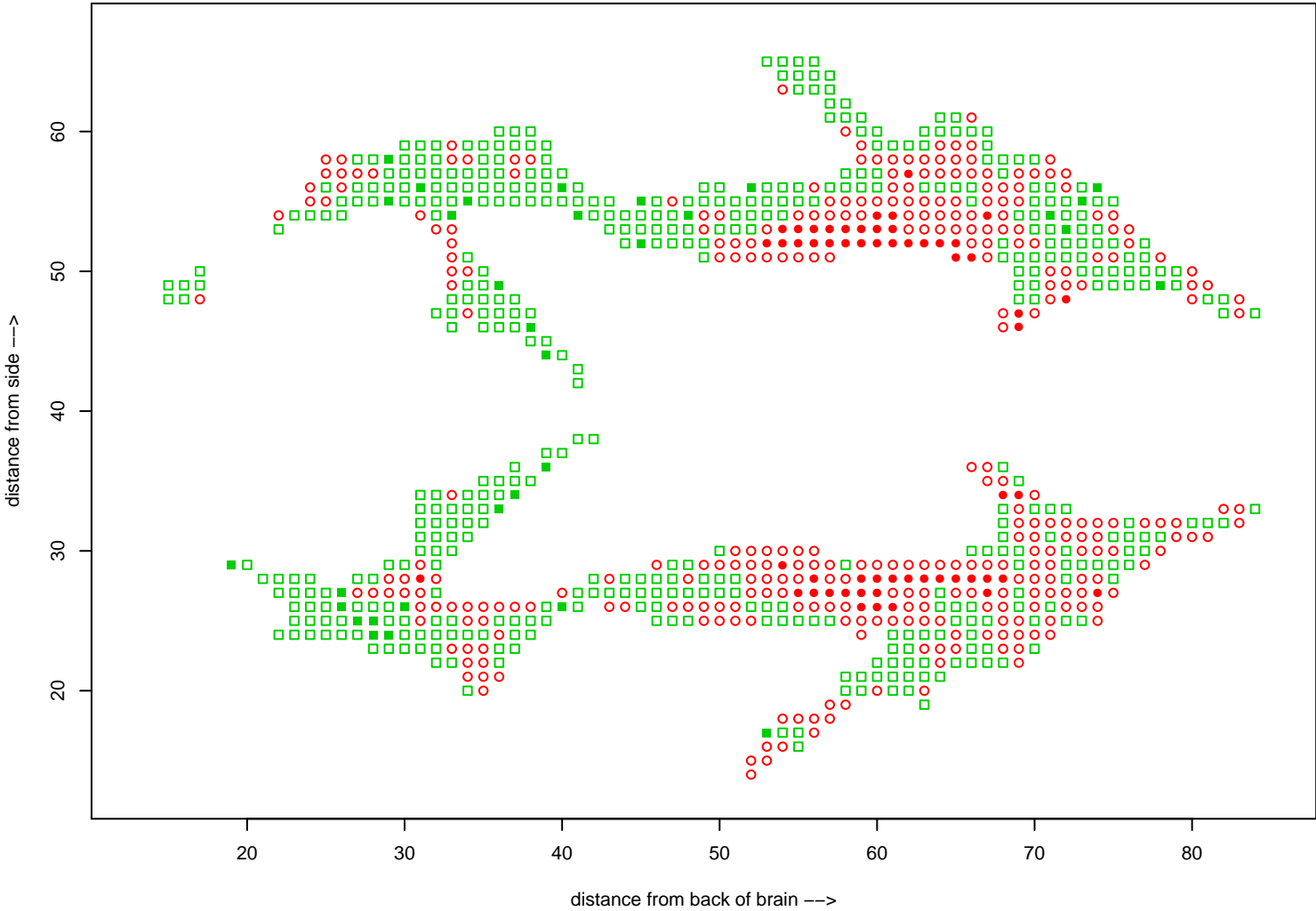
Tweedie and the fdr

- $-\log \text{fdr}(x) = \frac{x^2}{2\sigma_0^2} + l(x) + \text{constant}$ [$l(x) = \log f(x)$]
- $\frac{d}{dx} -\log \text{fdr}(x) = \frac{1}{\sigma_0^2} [x + \sigma_0^2 l'(x)] = \frac{E\{\mu|x\}}{\sigma_0^2}$
- `locfdr` (language R) estimates $\hat{l}(x)$ as before $\Rightarrow \widehat{\text{fdr}}(x)$
- fdr requires \hat{l} ; Tweedie requires \hat{l}'

Empirical Bayes estimation curve for the DTI data $E\{\mu|x\}$ (solid) and also false discovery rate $fdr(x)$ (dashed); at $x=3.5$, $fdr(x)=.11$ and $E\{\mu|x\}=1.88$



Horizontal slice of DTI brain data (848 voxels);
Red for $x > 0$, Green for $x < 0$; solid for $\text{abs}(x) > 2$



References

<http://stat.stanford.edu/~brad/papers/>

- “Tweedie’s formula and selection bias”
- “Microarrays, empirical Bayes, and the two-groups model”
- “The future of indirect evidence”
- “Empirical Bayes estimates for large-scale prediction problems”