

The Bayes Deconvolution Problem

Bradley Efron

Stanford University

Bayes Deconvolution Problem

- Unknown prior density $g(\theta)$ gives *unobserved* realizations

$$\Theta_1, \Theta_2, \dots, \Theta_N \stackrel{\text{iid}}{\sim} g(\theta)$$

- Each Θ_k gives *observed* $X_k \sim p_{\Theta_k}(x)$ [$p_{\Theta}(x)$ known]
- Marginal density

$$f(x) = \int p_{\theta}(x)g(\theta) d\theta$$

- Wish to estimate $g(\theta)$ from X_1, X_2, \dots, X_N

Two Familiar Cases

- **Poisson** $X_k \sim \text{Poi}(\Theta_k)$, $p_\theta(x) = e^{-\theta}\theta^x/x!$

$$\left(\text{Robbins: } \hat{E}\{\Theta|X = x\} = (x + 1) \frac{\hat{f}(x + 1)}{\hat{f}(x)}; \quad \hat{f}(x) = \frac{\#\{X_k = x\}}{N} \right)$$

- **Normal** $X_k \sim \mathcal{N}(\Theta_k, 1)$, $p_\theta(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\theta)^2}$

$$\left(\text{Tweedie, Stein, Brown: } \hat{E}\{\Theta|X = x\} = x + \frac{d}{dx} \log \hat{f}(x) \right)$$

- Neither requires \hat{g}

More Ambitious Goal

- Estimate entire prior density $g(\theta)$
- Why?
- *Ensemble properties* of Θ 's

$$\Pr\{\Theta > 2\}, \text{ etc.}$$

- *Empirical Bayes*: estimate full posterior distribution

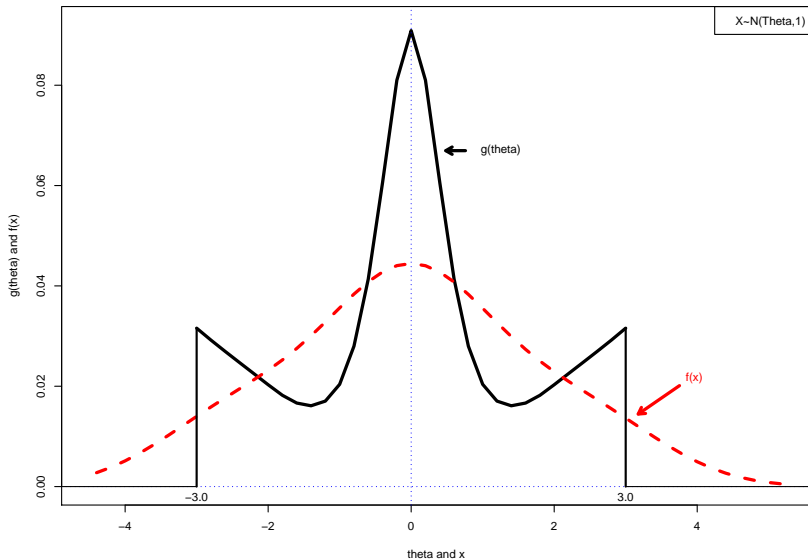
$$\Pr\{\Theta > 2|X = x\}, \text{ etc.}$$

- Asymptotics discouraging (Carroll and Hall, 1988)

A Test Case

- **Normal model** $X_k \sim \mathcal{N}(\Theta_k, 1)$ with $\Theta \in [-3, 3]$
- $g(\theta)$ an equal mixture of $\mathcal{N}(0, 0.5^2)$ and a symmetric density proportional to $|\theta|$
- Gives triangular-shaped marginal $f(x)$
- *Goal* Sample from $f(\cdot)$, estimate g

Test Case: prior density $g(\theta)$ (black) gives marginal density $f(x)$ (red); Goal: sample from f , estimate g



Fourier Method (Stefanski and Carroll, 1990)

- $X \sim \mathcal{N}(\Theta, 1)$: $\mathcal{F}(f) = \mathcal{F}(g)e^{-t^2/2}$ (\mathcal{F} = Fourier transform)

- Smoothing

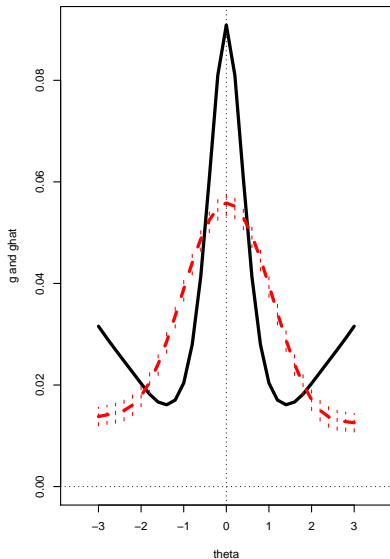
$$\hat{f}(x) = \frac{1}{N} \sum_{k=1}^N \sin\left(\frac{X_k - x}{\lambda}\right) / (X_k - x)$$

- Stef-Carroll: $\hat{g}(\theta) = \mathcal{F}^{-1}\{\mathcal{F}(\hat{f})e^{t^2/2}\}$

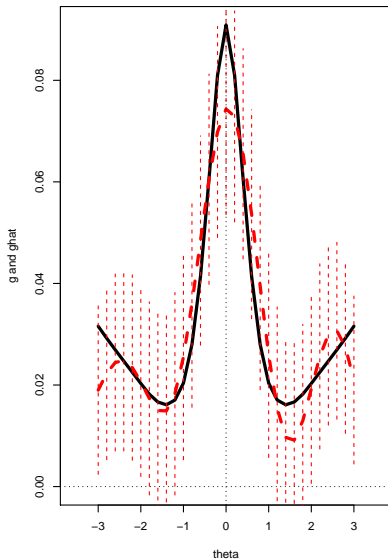
- Kernel form $\hat{g}(\theta) = \frac{1}{N} \sum_{k=1}^N k_\lambda(X_k - \theta)$ where

$$k_\lambda(x) = \frac{1}{\pi} \int_0^{1/\lambda} e^{t^2/2} \cos(tx) dt$$

Test Case: true $g(\theta)$ (black) and expected Fourier est $\lambda=0.50$ (red) \pm stdev, $N=4000$



Now for $\lambda=0.333$

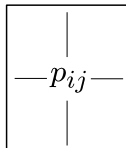


Discretization

- Θ in $(\theta_1, \theta_2, \dots, \theta_m)$ and X in (x_1, x_2, \dots, x_n)
- $\mathbf{g} = (g_1, g_2, \dots, g_m)$ and $\mathbf{f} = (f_1, f_2, \dots, f_n)$
- $\mathbf{P}_{m \times n} = (p_{ij})$ with $p_{ij} = \Pr\{X = x_i | \Theta = \theta_j\}$
- Then $\mathbf{f} = \mathbf{P}\mathbf{g}$
- *Test case*

$$\boldsymbol{\theta} = (-3, -2.8, \dots, 3) \text{ and } \mathbf{x} = (-4.4, -4.2, \dots, 5.2)$$

$$m = 31, n = 49 \text{ and } p_{ij} = \varphi(x_i - \theta_j) \cdot 0.20$$



Inversion and Regularization

- $\mathbf{f} = \mathbf{P}\mathbf{g}$ with $\text{svd } \mathbf{P} = \mathbf{U}\mathbf{d}\mathbf{V}'$
- $\mathbf{g} = \mathbf{P}^{-}\mathbf{f}$ with $\mathbf{P}^{-} = \mathbf{V}\frac{1}{\mathbf{d}}\mathbf{U}'$
- Regularization

$$\boxed{\mathbf{g}_r = \mathbf{P}_r^{-}\mathbf{f}} \quad \text{with } \mathbf{P}_r^{-} = \mathbf{V}_r\frac{1}{\mathbf{d}_r}\mathbf{U}_r'$$

(first r eigens; bias/variance trade-off)

- *Stef-Carroll*:

$$\boxed{\mathbf{g} = \mathbf{k}_\lambda\mathbf{f}} \quad \text{with } \mathbf{k}_\lambda(i, j) = k_\lambda(x_i - \theta_j)$$

($\lambda = 1/2$ like $r = 4$; $\lambda = 1/3$ like $r = 7$)

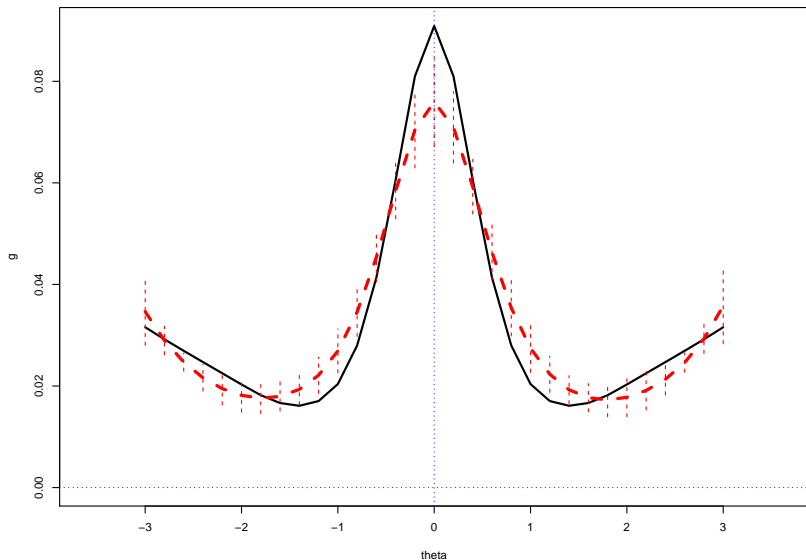
g Modeling

- **Idea** Model $\mathbf{g} = (g_1, g_2, \dots, g_m)$ as exponential family on $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_m)$:

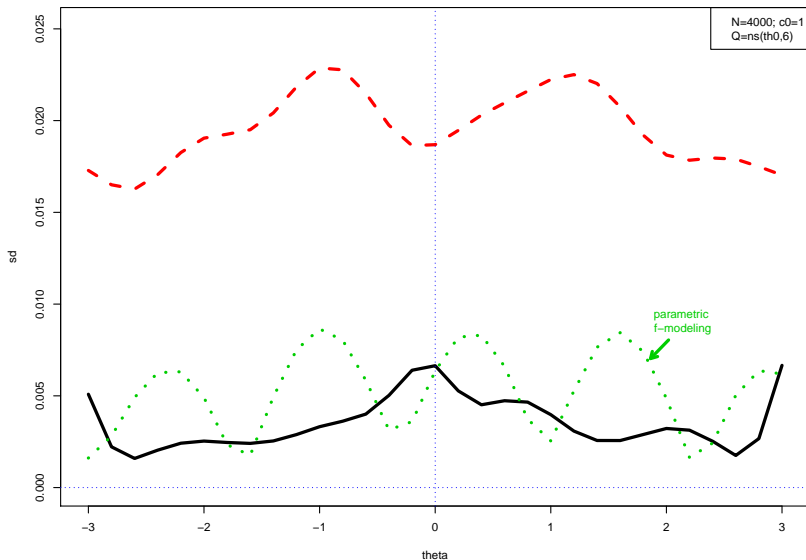
$$\mathbf{g}_\alpha = e^{\mathbf{Q}\alpha} / a_\alpha \quad \left[\begin{array}{l} \mathbf{Q} \text{ structure matrix, } \alpha \text{ natural parameter} \\ m \times p \end{array} \right]$$

- $\mathbf{y} = (y_1, y_2, \dots, y_n)$ vector of counts: $y_i = \#\{X_k = x_i\}$
- **MLE** $\alpha \rightarrow \mathbf{g}_\alpha \rightarrow \mathbf{f}_\alpha = \mathbf{P}\mathbf{g}_\alpha \rightarrow \mathbf{y} \sim \text{Mult}_n(N, \mathbf{f}_\alpha) \xrightarrow{\text{mle}} \hat{\alpha}$
- *Test case* $\mathbf{Q} = ns(\boldsymbol{\theta}, 6)$

g-modeling for Test Case: $Q=ns(\text{theta},6)$, $N=4000$, $c_0=1$;
Solid curve true g; dashed curve mean 400 sims, +1 sd



Test Case Stdevs for ghat: g-modeling (black),
Steff-Car (red), and parametric f-modeling (green)



Fisher Information Calculations

- Define $W_{ij} = g_{\alpha j} \left(\frac{P_{ij}}{f_{\alpha i}} - 1 \right)$
- $\mathbf{W}_{\alpha} = (W_{ij})_{n \times m}$
- *Observed Fisher information* at MLE $\hat{\alpha}$:

$$\hat{\mathbf{I}}_{\hat{\alpha}} = -\ddot{\ell}_{\hat{\alpha}} = \mathbf{Q}' \left\{ \mathbf{W}'_{\hat{\alpha}} \text{diag}(\mathbf{y}) \mathbf{W}_{\hat{\alpha}} - \text{diag}(\mathbf{W}_{\hat{\alpha}} \mathbf{y}) \right\} \mathbf{Q}$$

- *Expected Fisher information* at $\alpha = \hat{\alpha}$:

$$\mathbf{I}_{\hat{\alpha}} = \mathbf{Q}' \left\{ \mathbf{W}'_{\hat{\alpha}} \text{diag}(\mathbf{y}) \mathbf{W}_{\hat{\alpha}} \right\} \mathbf{Q}$$

Regularization and Accuracy for g Models

- $\hat{\alpha} = \arg \max_{\alpha} \{\ell_{\alpha} - \mathbf{s}_{\alpha}\} \quad (\mathbf{s}_{\alpha} = \mathbf{c}_0 \|\alpha\|)$

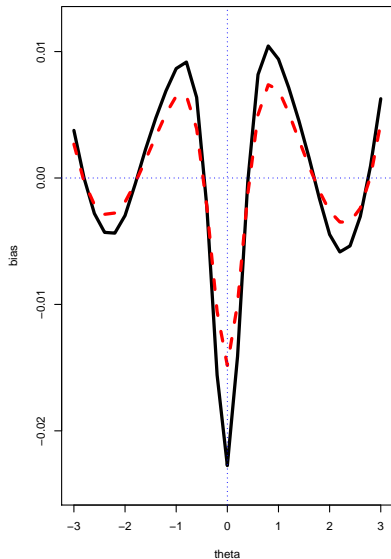
$$\hat{\alpha} - \alpha \sim \left[\underbrace{-\left(\hat{\mathbf{l}} + \ddot{\mathbf{s}}_{\hat{\alpha}}\right)^{-1} \dot{\mathbf{s}}_{\hat{\alpha}}}_{\widehat{\text{Bias}}}, \quad \underbrace{\left(\hat{\mathbf{l}} + \ddot{\mathbf{s}}_{\hat{\alpha}}\right)^{-1} \hat{\mathbf{l}} \left(\hat{\mathbf{l}} + \ddot{\mathbf{s}}_{\hat{\alpha}}\right)^{-1}}_{\widehat{\text{Cov}}} \right]$$

- Letting $\hat{R} = [\text{diag}(\hat{g}) - \hat{g}\hat{g}'] Q$,

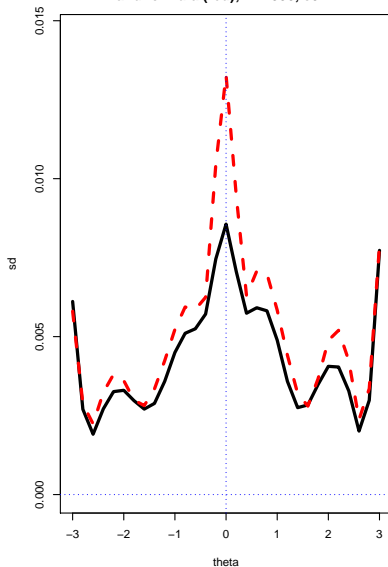
$$g_{\hat{\alpha}} - g \sim \left(\hat{R} \cdot \widehat{\text{Bias}}, \hat{R} \cdot \widehat{\text{Cov}} \cdot \hat{R}' \right)$$

$$\left[\mathbf{c}_0 = 1 \text{ made } \text{tr}(\ddot{\mathbf{s}}_{\hat{\alpha}}) / \text{tr}(\hat{\mathbf{l}}) = 0.03 \right]$$

Test Case: simulation bias (red)
and formula (black)



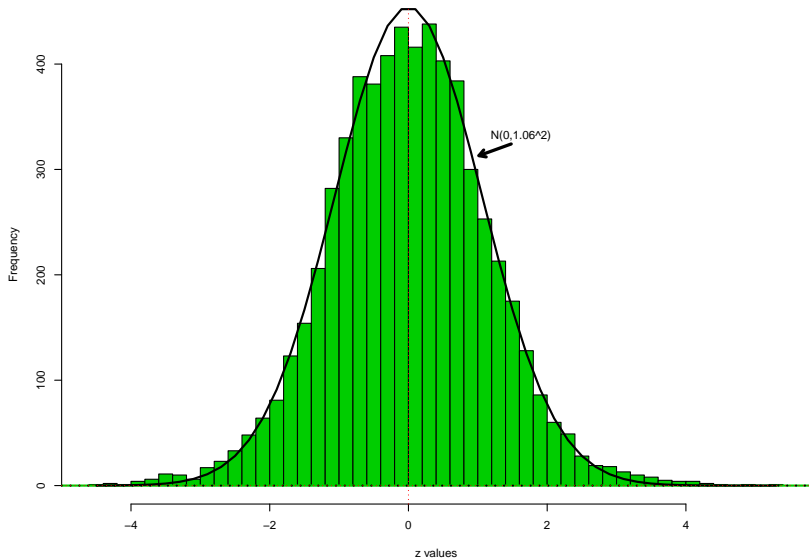
Simulation stdev (black)
and formula (red); N=4000, c0=1



Parametric f Modeling

- *Stef-Carroll*: $\hat{g} = k_\lambda \bar{f}$ where \bar{f} empirical density \mathbf{y}/N
- Instead take $\hat{g} = k_\lambda \hat{f}$, \hat{f} parametric estimate of f
- **Slide 13** $\hat{f} = \text{glm}(\mathbf{y} \sim ns(\mathbf{x}, 6), \text{Poisson})$ est
- Need $X_k = \Theta_k + \epsilon_k$, with ϵ_k iid noise
- Need $g(\theta)$ smooth

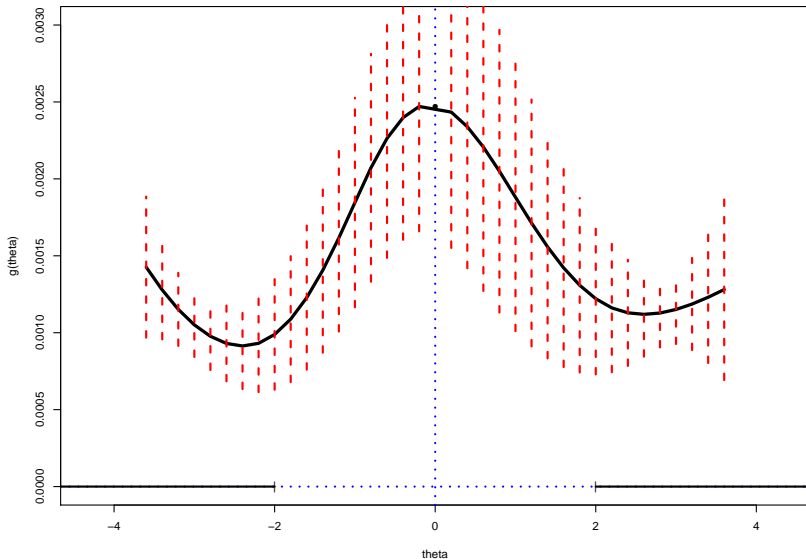
Prostate study data: $N=6033$ z-values for 52 patients vs 50 controls. locfdr: empirical null $\sim N(0, 1.06^2)$, $p_0=0.984$



Estimating $g(\theta)$ for the Prostate Study

- z-value for k th gene $z_k \sim \mathcal{N}(\Theta_k, 1.06^2)$ ($\Theta_k =$ “effect size”)
- **g model** $Q = (\mathbf{1}, ns(\theta, 5))$ (“spike and slab”)
- Tried $c_0 = 0.5, 1, 2$
- *Accuracy formula* gave \widehat{sd} , and $\widehat{Bias} \doteq 0$

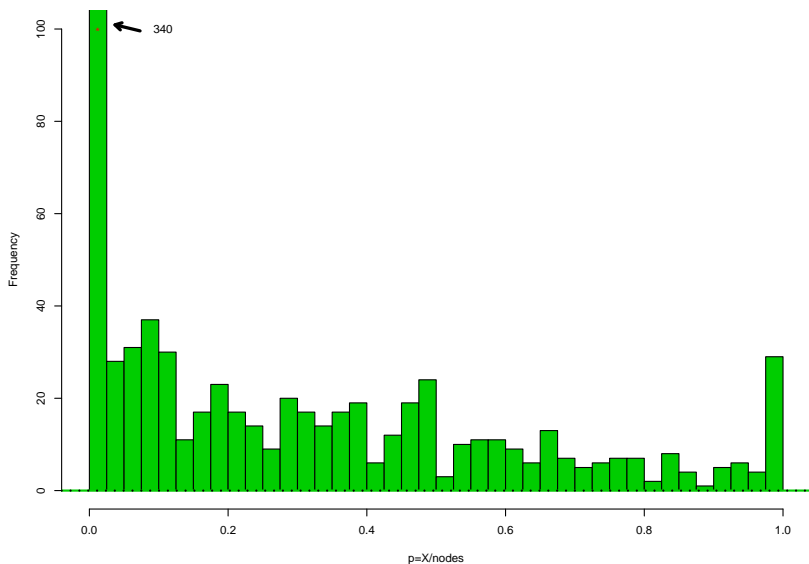
Non-null prior for prostate data, g-model $Q=(1,ns(theta,5),c0=1;$
Null atom .946; Prob{ $|theta| > 2$ }=.02; \pm one stdev



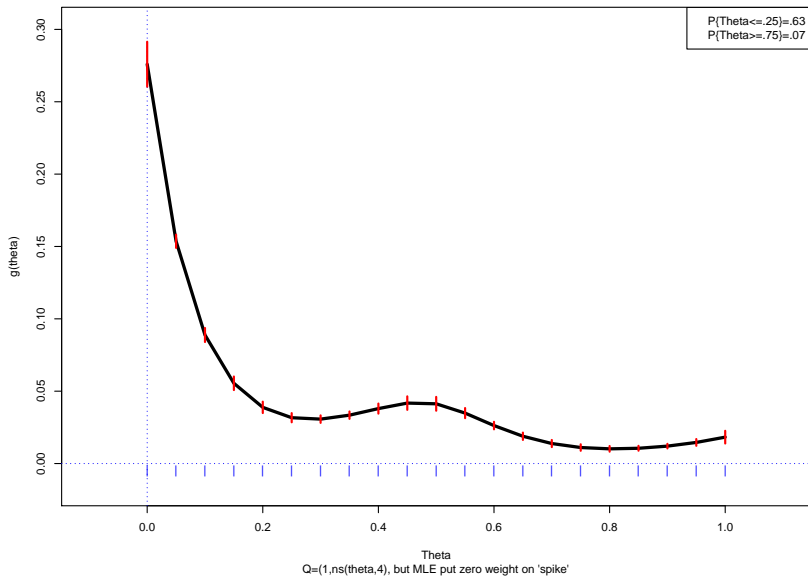
A Binomial Example

- 844 cancer patients: n_k lymph nodes removed; X_k found positive
- *Binomial model* $X_k \sim \text{binom}(n_k, \Theta_k)$ [$\theta = (0, 0.05, 0.10, \dots, 1)$]
- *g modeling* Prior $\mathbf{g} = \mathbf{e}^{\mathbf{Q}\alpha} / \mathbf{a}_\alpha$ with $\mathbf{Q} = (1, ns(\theta, 4))$
- $P_{844 \times 21} : P_{kj} = \binom{n_k}{x_k} \theta_j^{x_k} (1 - \theta_j)^{n_k - x_k}$; $\mathbf{f}_\alpha = \mathbf{P}\mathbf{g}_\alpha$
- *MLE* $\ell_\alpha = \sum_{k=1}^{844} \log(f_{\alpha k})$ [sd's from $-\ddot{\ell}_{\hat{\alpha}}$]
- Fan (1991): Binomial easier than normal

Nodes study: ratio $p=X/n$ for 844 cases;
n ranging from 1 to 69



G-model estimate of prior distribution $g(\theta)$, 844 cases;
Theta the true effect size, nodes study; +- stdev



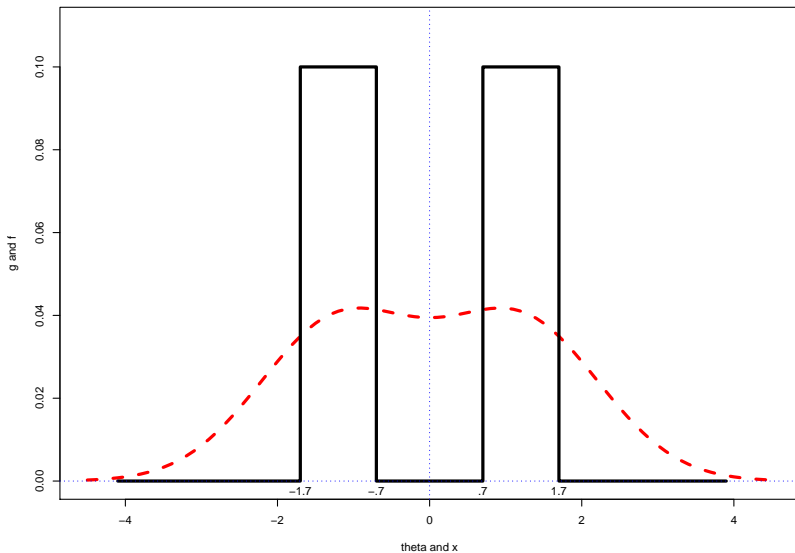
A Difficult Example

- “Two Towers”

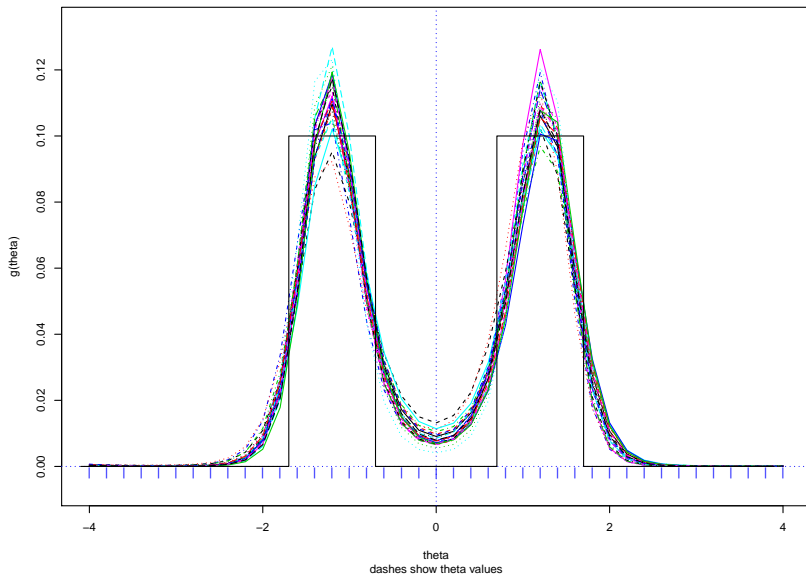
$$g(\theta) = \begin{cases} 1 & \text{on } [-1.7, -0.7] \text{ and } [0.7, 1.7] \\ 0 & \text{otherwise} \end{cases}$$

- $X \sim \mathcal{N}(\Theta, 1)$
- g modeling $\theta = (-4, -3.8, -3.6, \dots, 3.8, 4)$
- $Q = ns(\theta, 6), \quad c_0 = 1$

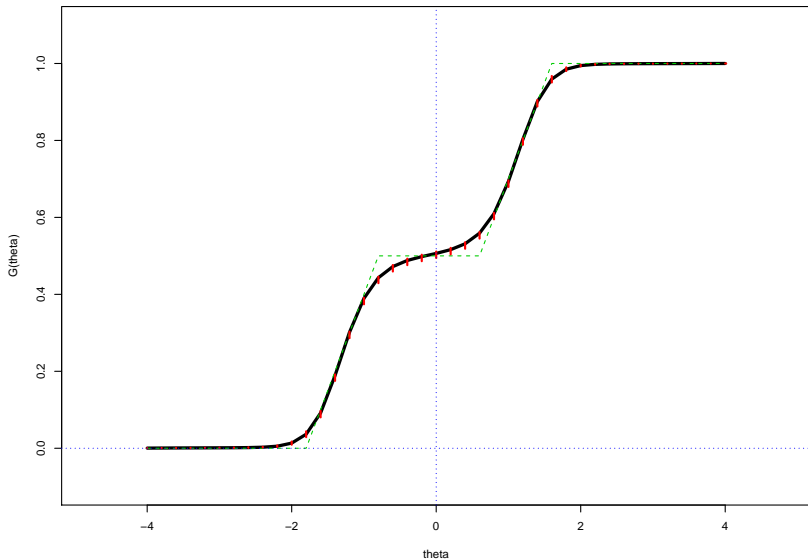
Two towers example: $g(\theta)$ black, $f(x)$ red;
 $X = \Theta + N(0,1)$



G modeling for Two Towers: $Q=ns(\text{theta},6)$, $N=4000$, $c_0=1$;
first 25 simulations



Mean and Sd of 400 gmodel cdf's for Two Towers example;
dashed line is true cdf



Influence Functions

- Derivative matrix of $\hat{\mathbf{g}}$ with respect to \mathbf{y}

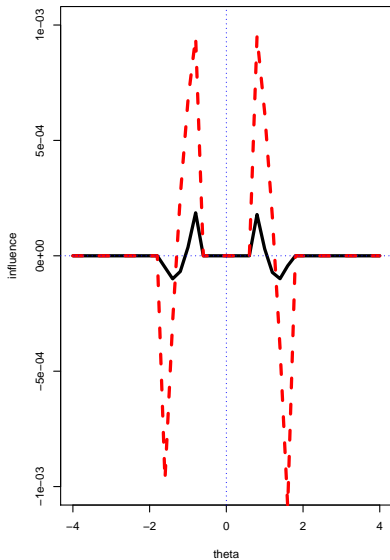
- g model $\frac{d\hat{\mathbf{g}}}{d\mathbf{y}} = \left[\text{diag}(\mathbf{f}_{\hat{\alpha}}) - \mathbf{f}_{\hat{\alpha}}\mathbf{f}'_{\hat{\alpha}} \right] Q (\hat{\eta} + \ddot{s}_{\hat{\alpha}})^{-1} Q' \mathbf{W}'_{\hat{\alpha}}$

- f model $\tilde{\mathbf{g}} = \mathbf{k}_{\lambda} \hat{\mathbf{f}}, \quad \hat{\mathbf{g}} = \tilde{\mathbf{g}}/\tilde{g}_+$,

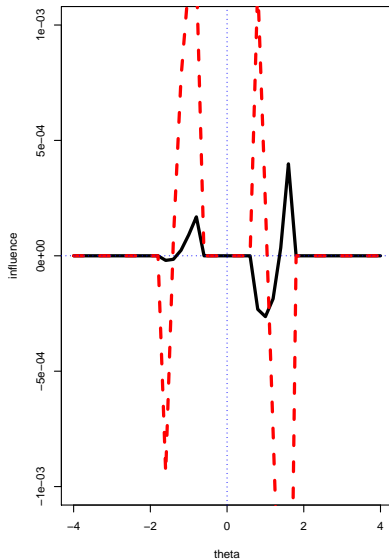
and $\hat{\mathbf{f}}$ from $\text{glm}(\mathbf{y} \sim X, \text{Poisson})$

$$\frac{d\hat{\mathbf{g}}}{d\mathbf{y}} = \frac{1}{N} \text{diag}(\hat{\mathbf{g}}) \left[\text{diag} \left(\frac{1}{\hat{\mathbf{g}}} \right) - \frac{1}{\tilde{g}_+} \right] \text{diag}(\hat{\mathbf{f}}) \left[X (X' \text{diag}(\hat{\mathbf{f}}) X)^{-1} X' \right]$$

Influence function dg/dy at $x=0$, two towers;
gmod black, fmod red



Influence function dg/dy at $x=3$



Empirical Bayes Information

- $\text{CoeffVar}\{\hat{g}(\theta)\} \doteq 1/Ni_\theta$ ($i_\theta =$ “information in single X ”)
- *Test case* (g model)

θ	-1	0	1	2	3
i_θ	.0038	.0117	.0030	.0034	.0040
$N_{0.1}$	26,300	8500	32,800	23,100	25,000

($N_{0.1} = \#X_k$ needed for $cv = 0.1$)

Summary

- **Nonparametric f modeling:** deprecated
- **Parametric f modeling:** OK for smooth g , additive noise
(preferred for Robbins/Tweedie situations)
- **g modeling:** flexible and reasonably efficient for a wide variety of situations

References

- Carroll, R. J. and Hall, P. (1988). Optimal rates of convergence for deconvolving a density. *J. Amer. Statist. Assoc.* 83: 1184–1186.
- Efron, B. (2011). Tweedie's formula and selection bias. *J. Amer. Statist. Assoc.* 106: 1602–1614, doi: 10.1198/jasa.2011.tm11181.
- Efron, B. (2014). Two modeling strategies for empirical Bayes estimation. *Statist. Sci.* To appear.
- Fan, J. (1991). On the optimal rates of convergence for nonparametric deconvolution problems. *Ann. Statist.* 19: 1257–1272, doi: 10.1214/aos/1176348248.
- Stefanski, L. and Carroll, R. J. (1990). Deconvoluting kernel density estimators. *Statistics* 21: 169–184, doi: 10.1080/02331889008802238.