

BAYES, ORACLE BAYES, AND EMPIRICAL BAYES

Bradley Efron

Stanford University

EMPIRICAL BAYES INFERENCE

- Robbins (1951) *Compound Decision Procedures*
- Robbins (1956) *Empirical Bayes*
- **Question** How does *empirical Bayes* relate to Bayesian and frequentist inference?
- Intermediate framework: Oracle Bayes

A FAMILIAR EMPIRICAL BAYES SETUP

- Unknown prior $g(\theta)$ produces unseen parameters

$$g(\theta) \longrightarrow \theta_1, \theta_2, \dots, \theta_N$$

- Each θ_i independently gives observation $x_i \sim \mathcal{N}(\theta_i, 1)$,

$$\mathbf{x} = (x_1, x_2, \dots, x_N)$$

- **Robbins** Use all of \mathbf{x} to estimate each θ_i

ORACLE BAYES

JIANG AND ZHANG (2009)

- Oracle tells us the order statistic of the θ 's,

$$\theta_{\text{ord}} = (\theta_{(1)}, \theta_{(2)}, \dots, \theta_{(N)})$$

(but not which x_i goes with which θ_i)

- We want to estimate $\hat{\theta}_i$ that minimizes

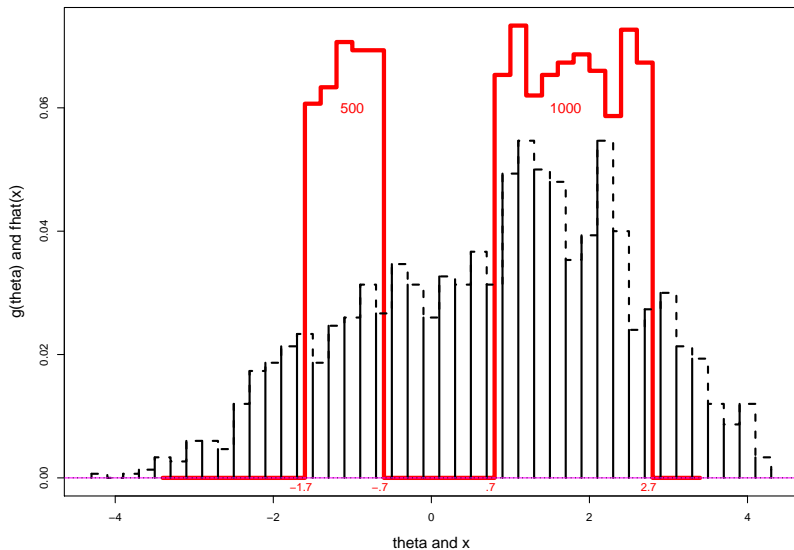
Expected Average Squared Error

$$\text{EASE} = E \left\{ \frac{1}{N} \sum_{i=1}^N (\hat{\theta}_i - \theta_i)^2 \right\}$$

- Example: “Two Towers”



TWO TOWERS EXAMPLE: Oracle (Red), Observations (Black)



USING THE ORACLE

- $\bar{g}(\theta)$ empirical density of the θ 's
(probability $1/N$ on $\theta_{(1)}, \theta_{(2)}, \dots, \theta_{(N)}$)
- *Bayes estimate* (ϕ the standard normal density)

$$\hat{\theta}_{\bar{g}}(x) = \frac{\sum_1^N \theta_{(i)} \phi(x - \theta_{(i)})}{\sum_1^N \phi(x - \theta_{(i)})}$$

- Minimizes EASE among rules $\hat{\theta}(x) = t(x)$
- **Two Towers** $\text{EASE}_{\bar{g}} = 0.563$
- MLE estimates $\hat{\theta}_i = x_i$ has $\text{EASE}_{\text{MLE}} = 1.0$

EMPIRICAL BAYES

- **No Oracle** Use $\mathbf{x} = (x_1, x_2, \dots, x_N)$ to form estimate $\hat{\theta}_{\mathbf{x}}(x)$ of Oracle rule $\hat{\theta}_{\bar{g}}(x)$
- *Empirical Bayes regret*

$$\text{EASE}_{\mathbf{x}} - \text{EASE}_{\bar{g}} = \text{EBregret}$$

- **Two Towers** $\text{EBregret} = 0.008$
- $\text{EASE}_{\mathbf{x}} = 0.563 + 0.008 = 0.571$

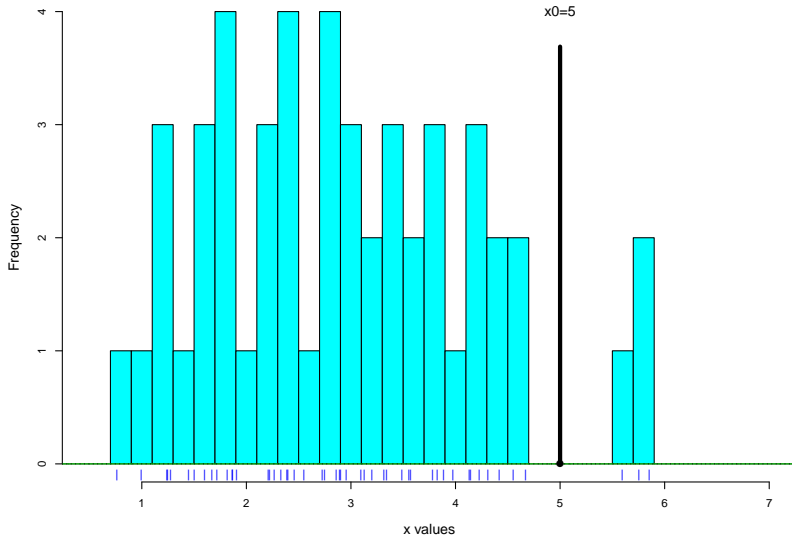
THE FREQUENTIST FACE OF EMPIRICAL BAYES

- **Oracle** application is entirely frequentist:
 - ▶ Minimizing EASE is a frequentist criterion
 - ▶ Assumption $\theta_i \sim g(\theta)$ is irrelevant!
- **Question** Is the 40% EASE reduction meaningful?
- **Answer** Yes for EASE, No for some other applications

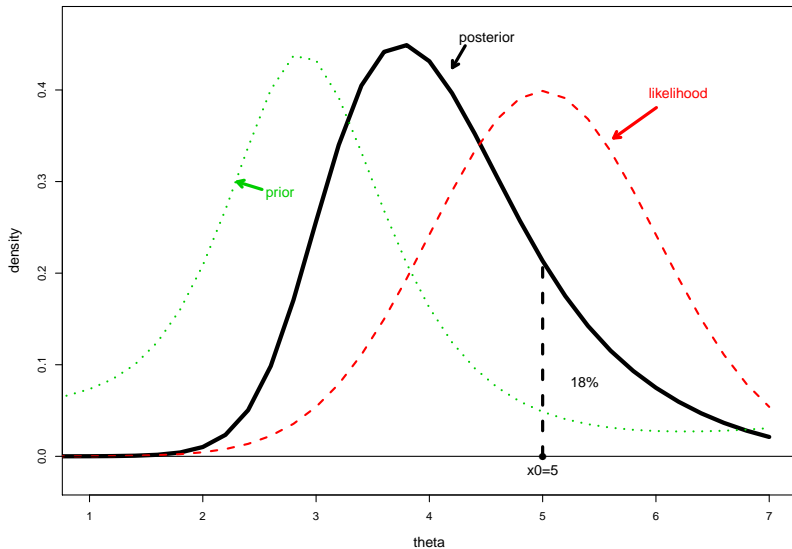
FINITE BAYES INFERENCE

- $\theta_i \stackrel{\text{ind}}{\sim} g(\theta)$ and $x_i | \theta_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\theta_i, 1)$ for $i = 0, 1, 2, \dots, N$
- We want posterior inference for θ_0 given x_0 and the “sibling” observations $\mathbf{x} = (x_1, x_2, \dots, x_N)$
- If $N \rightarrow \infty$ then finite Bayes \rightarrow Bayes
- *Next example:* $x_0 = 5, N = 50$

FINITE BAYES INFERENCE: 50 'sibling' observations to $x_0=5$; what can we say about θ_0 ?



Red: likelihood $N(x_0=5,1)$; Green: Estimated prior from sibs,
Black: Estimated posterior for θ_0



STANDARD BAYES (JUST ONE θ AND x)

- Prior $\theta \sim g(\theta)$ ■ Observe $x \sim p_\theta(x)$ [$\mathcal{N}(\theta, 1)$]
- Marginal density $f(x) = \int_{\mathcal{X}} g(\theta) p_\theta(x) d\theta$
- *Bayesian inference* $x \sim f(x)$ and

$$\theta | x \sim [e_g(x), v_g(x)]$$

conditional
expectation ↗

↖ conditional
variance

- For squared error: $e_g(x)$ is Bayes estimate “ $\hat{\theta}_g(x)$ ”
- *Bayes risk* $\mathcal{R}_g = E \left\{ (\theta - e_g(x))^2 \right\} = \int_{\mathcal{X}} v_g(x) f(x) dx$
- $\mathcal{R}_{\bar{g}} = \text{EASE}_{\bar{g}}$ (= 0.563)

TWEEDIE'S FORMULA (EFRON, 2011)

- If $x \mid \theta \sim \mathcal{N}(\theta, 1)$:

$$e_g(x) = x + l'(x) \quad [l(x) = \log f(x)]$$

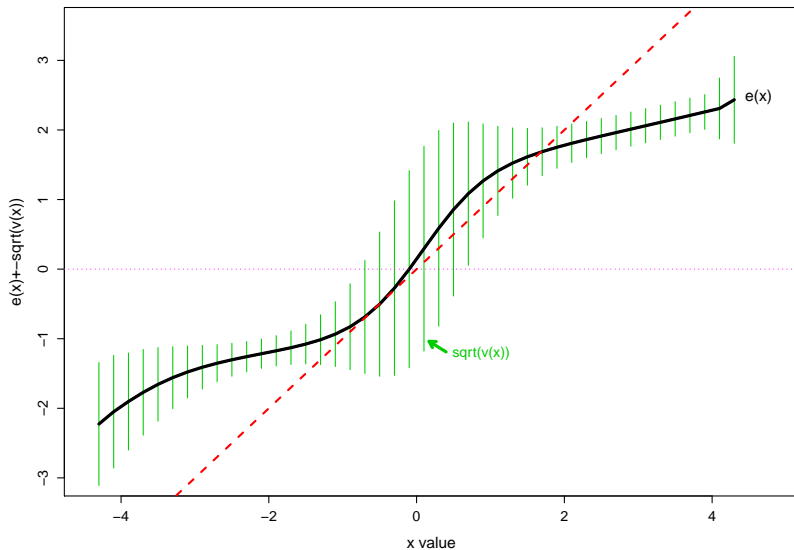
$$v_g(x) = 1 + l''(x)$$

- So $\hat{\theta}_g(x) = x + l'(x)$ and

$$\mathcal{R}_g = \int_{\mathcal{X}} [1 + l''(x)] f(x) = 1 - \int_{\mathcal{X}} [l'(x)]^2 f(x)$$

- Next: $g(\theta) = \bar{g}(\theta)$, Two Towers

theta|x: expectation $e(x) \pm \text{var}(x)^{.5}$,
Two Towers example



EMPIRICAL BAYES RISK AND REGRET

- **Idea** Estimate “ $\hat{f}(x)$ ” from data $\mathbf{x} = (x_1, x_2, \dots, x_N)$
- $\hat{e}(x) = x + \hat{l}'(x)$ and $\hat{v}(x) = 1 + \hat{l}''(x)$ $[\hat{l}(x) = \log \hat{f}(x)]$
- *Empirical Bayes risk* $\mathcal{R}(g, \hat{e}) = E\{(\hat{e}(x) - \theta)^2\}$

LEMMA

$$\mathcal{R}(g, \hat{e}) = \mathcal{R}_g + \int_{\mathcal{X}} [\hat{e}(x) - e_g(x)]^2 f(x) dx$$

Bayes
risk ↗

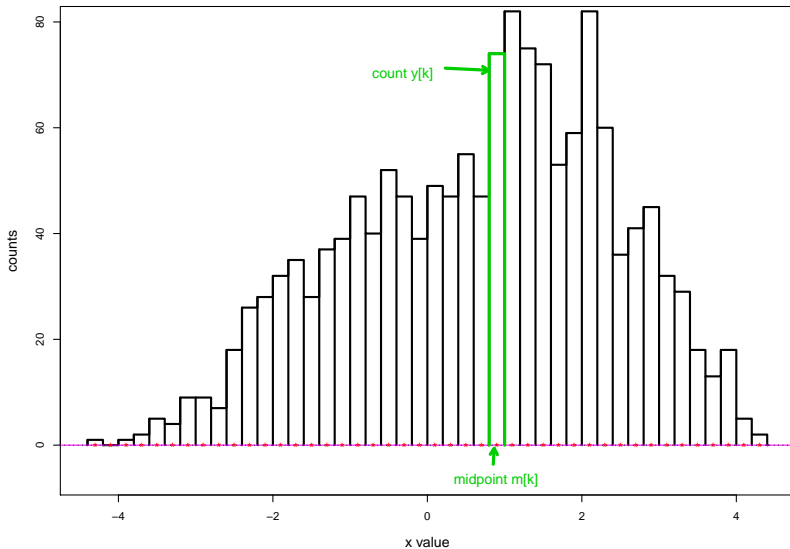
↖ EBregret

- Only need f and its estimate \hat{f}

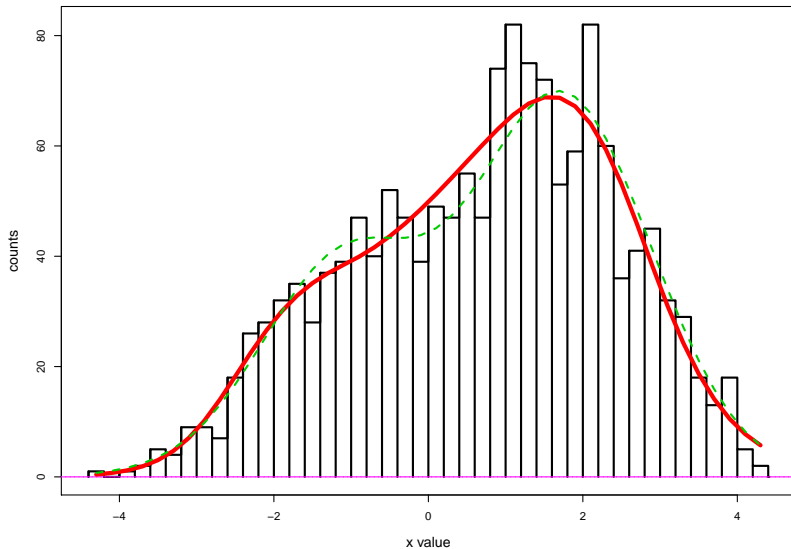
ESTIMATING $f(x)$: “ f -MODELING”

- Bin data x_1, x_2, \dots, x_N
- y_k = number of x 's in bin k
- m_k = midpoint of bin k
- $\hat{f} = \text{glm}(y \sim \text{ns}(m, df), \text{poisson})\est/N
- **Two Towers** $K = 44$ bins ($df = 5$)

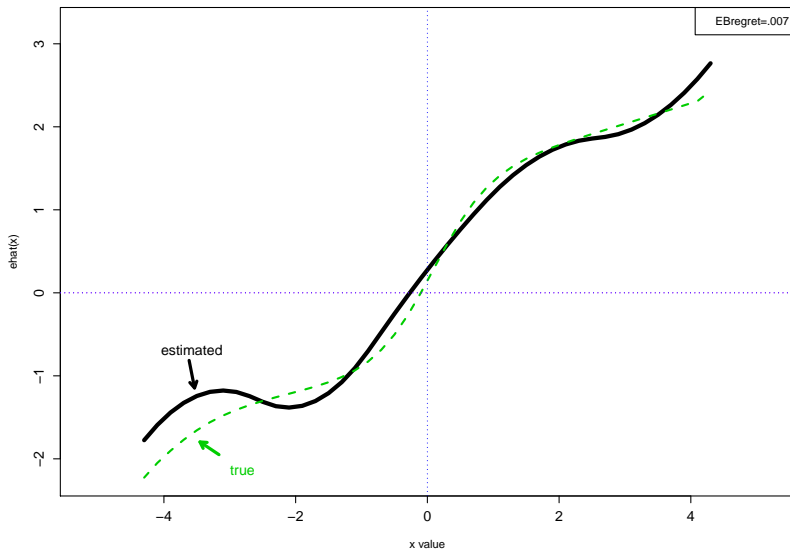
Histogram counts for the Two Towers data $x[i], i=1:1500$;
 $K=44$ bins ('bin[k]' in green)



Fitted marginal density $\hat{f}_{\text{hat}}(x)$ (*1500) from
`glm(counts~ ns(x,5),poisson)`; green = true $f(x)$



Estimated conditional expectation $e(x)=E(\theta|x)$
Two Towers, $\text{glm}(y\sim\text{ns}(m,5),\text{poisson})$; true $e(x)$ in green



DATA-BASED FORMULA FOR EBREGRET

- $\mathbf{M}_{44 \times 6} = (\text{ns}(\mathbf{m}, 5), \mathbf{1})$, k th row \mathbf{M}_k
- $\dot{\mathbf{M}}_{44 \times 6}$, k th row $\frac{d\mathbf{M}_k}{dm_k}$
- $\hat{\mathbf{f}}$ = binned estimate of $f(x)$ from $\text{glm}(\mathbf{y} \sim \mathbf{M}, \text{poisson})$

THEOREM

$$E\{\text{EBregret}\} \doteq \frac{1}{N} \text{trace} \left\{ \left[\mathbf{M}' \text{diag}(\hat{\mathbf{f}}) \mathbf{M} \right]^{-1} \left[\dot{\mathbf{M}}' \text{diag}(\hat{\mathbf{f}}) \dot{\mathbf{M}} \right] \right\}$$

- Decreases as $1/N$
- Doesn't account for bias of $\hat{\mathbf{e}}$

THE POISSON CASE

- $\theta \sim g(\cdot)$ and $x \mid \theta \sim \text{Poi}(\theta)$ $[p_\theta(x) = e^{-\theta}\theta^x/x!]$

- $f(x) = \int g(\theta)p_\theta(x) d\theta$ and $\theta \mid x \sim [e_g(x), v_g(x)]$

- **Robbins**
$$e_g(x) = \frac{(x+1)f(x+1)}{f(x)}$$

and $v_g(x) = e_g(x) [e_g(x+1) - e_g(x)]$

- *Bayes estimate* $\hat{\theta}_g(x) = e_g(x)$ with $\mathcal{R}_g = \sum_x f(x)v_g(x)$

$$\mathcal{R}[g, e(x)] = \mathcal{R}_g + \sum_x f(x) [e(x) - e_g(x)]^2$$

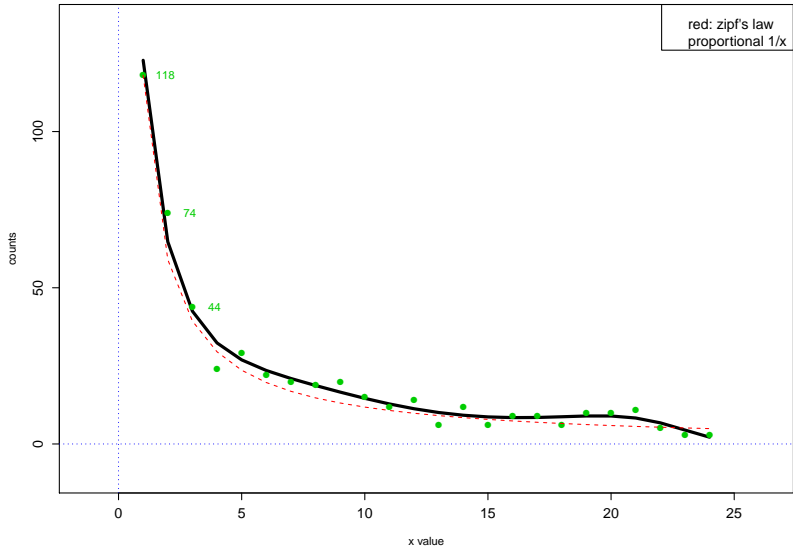
BUTTERFLY DATA

CORBET IN MALAYSIA, 1941–42

- $y_1 = 118$ species trapped just one time each
- $y_2 = 74$ species trapped twice each...

x	1	2	3	4	...	24
y	118	74	44	24	...	3

Butterfly data: observed counts and fitted
model $\text{glm}(y \sim \text{ns}(x, 5), \text{poisson})$



ZIPF'S LAW AND ROBBINS' FORMULA

- If $f(x) = c/x$ then Robbins' formula gives

$$e_g(x) = \frac{(x+1)f(x+1)}{f(x)} = x$$

- So the MLE $\hat{\theta}_i = x_i$ is Bayes!
- **Butterfly estimate** $\mathcal{R}_g = 5.24$ and $\mathcal{R}(g, x) = 5.39$

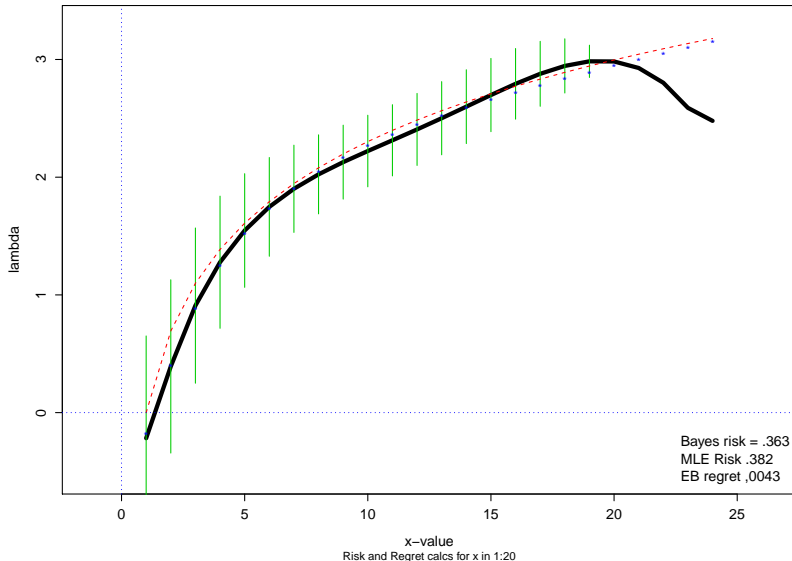
TWEEDIE POISSON ESTIMATION

- $\theta \sim g(\cdot)$ and $x \mid \theta \sim \text{Poi}(\theta)$
- $f(x) = \int g(\theta) (e^{-\theta} \theta^x / x!) d\theta$
- $\lambda = \log \theta$ is “natural parameter”

$$E\{\lambda \mid x\} = e_g(x) = \text{lgamma}(x + 1)' + l'(x)$$
$$\text{Var}\{\lambda \mid x\} = v_g(x) = \text{lgamma}(x + 1)'' + l''(x)$$

- Same for truncated Poisson!
- Next: Butterfly data, \hat{f} from `ns(df=5)`

Tweedie for lambda: post expectation $e(x) \pm v(x)^{.5}$;
Red line is $(x, \log(x))$; Points from g-modeling



EMPIRICAL BAYES: *g*-MODELING

- *f*-modeling: Never need to estimate prior $g(\theta)$
- But only for criteria that depend just on $f(x)$:
 - ▶ EASE
 - ▶ Tweedie
 - ▶ Robbins
 - ▶ false discovery rates...
- *g*-modeling (Efron, 2016): Direct parametric modeling of prior $g(\theta)$

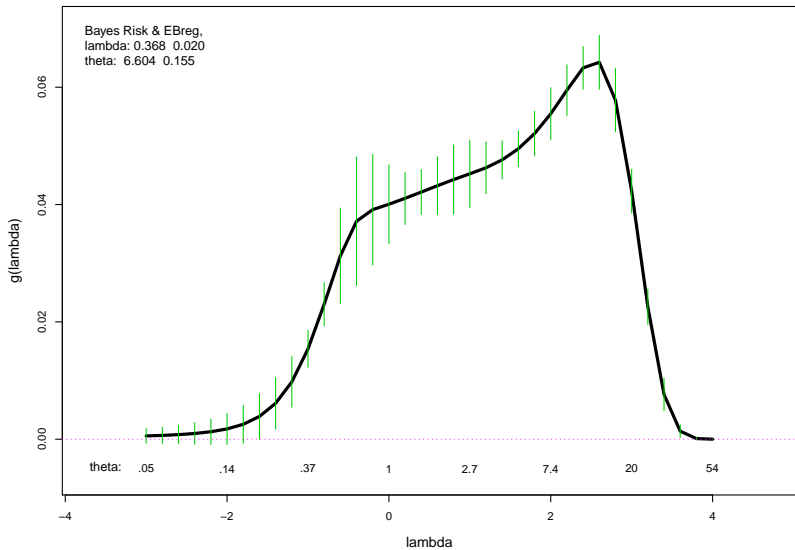
PARAMETRIC MODELS FOR $g(\theta)$

- *Exponential family* $\log g(\theta) = m(\theta)' \beta$ (“hidden GLM”)
- β a p -dimensional parameter; $m(\theta)$ say $\text{ns}(\theta, 5)$
- **MLE** $\hat{\beta}$ found by nonlinear maximization:

$$\beta \rightarrow g_{\beta}(\theta) \rightarrow f_{\beta}(x) \rightarrow y \sim \text{Mult}(N, f_{\beta})$$

- **Trouble** $y \sim \text{Mult}(N, f_{\beta})$ not exponential family
- **Advantage:** Estimate $g(\theta | x)$, $\Pr\{\theta > 2 | x\}$, etc.

Estimated prior $g(\lambda)$ for the Butterfly data,
 $\lambda = \log(\theta)$; prior model $M = ns(df=5)$



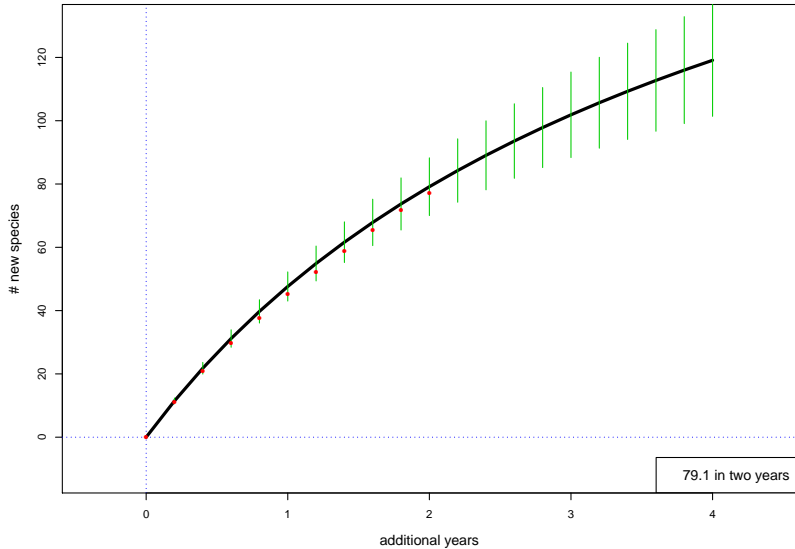
MISSING SPECIES PROBLEM

- **Corbet** How many new species in t more years of observation?
- Fisher, Good (Turing) \rightarrow Poisson process model

$$\frac{E\{\#\text{new}\}}{\#\text{old}} = \int e^{-\theta} \frac{1 - e^{-\theta t/2}}{1 - e^{-\theta}} g(\theta) d\theta$$

- *f-modeling*: Clever formula in terms of \hat{f}
(Good and Toulmin, 1956)
- *g-modeling*: Substitute $\hat{g}(\theta)$ for $g(\theta)$

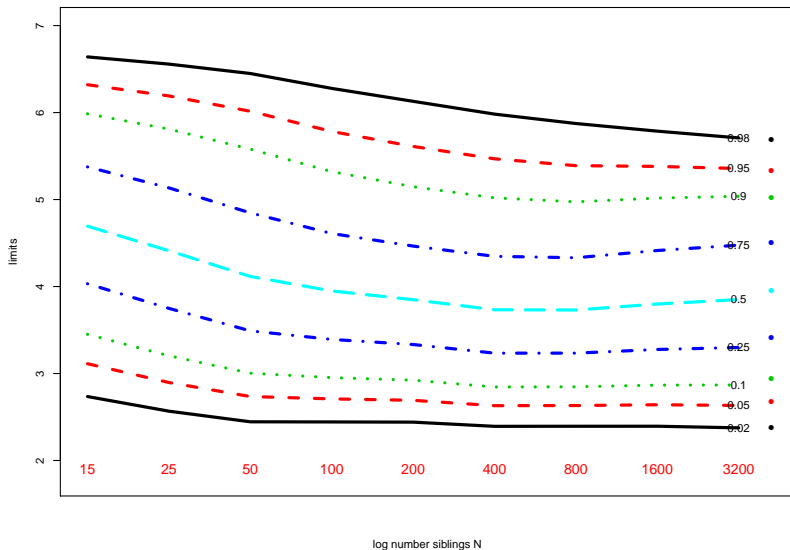
Expected number new species seen in t years additional trapping;
Red dots from Fisher–Good–Gaskins nonparametric formula



A g -MODELING EXAMPLE

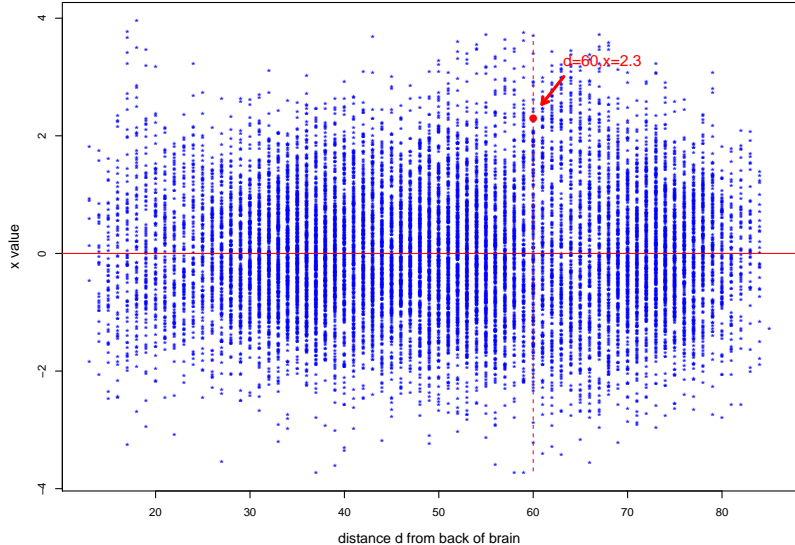
- “gamnorm” $\theta_i \sim \text{Gamma}_9/3$, $x_i \sim \mathcal{N}(\theta_i, 1)$, $i = 1, 2, \dots, 3200$
- $\mathbf{x} = (x_1, x_2, \dots, x_N)$ for $N = 15, 25, 50, \dots, 3200$
- $g(\theta) = \text{ns}(\theta, 5)$
- Next: posterior quantiles of $\theta_0 \mid x_0 = 5$
(finite Bayes inference as N increases)

POSTERIOR PERCENTILES of θ_0 given $x_0=5$ as the number of siblings increases; true gamnorm %iles at right

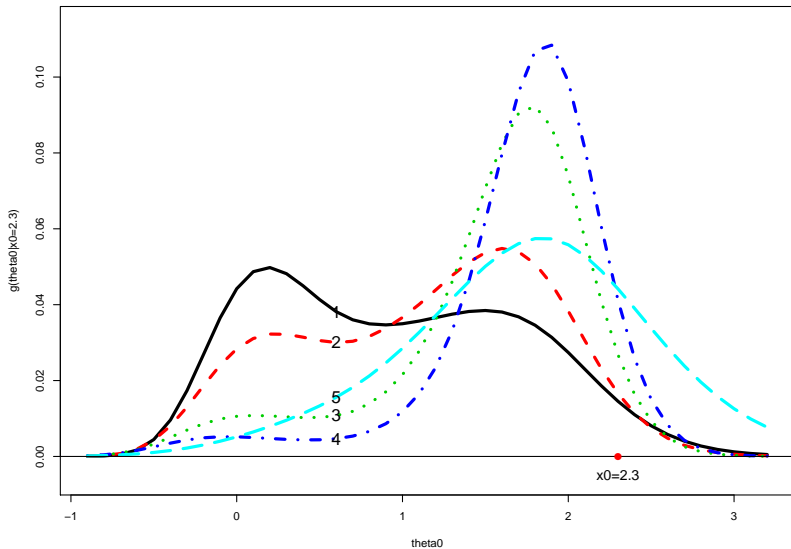


- 12 children, 6 dyslexic vs 6 normal controls
- $N = 15433$ brain voxels
- “ x_i ”, x -value for i th voxel, $x_i \sim \mathcal{N}(\theta_i, 1)$
- $\theta_i =$ true effect size ($= E\{x_i\}$)

DTI STUDY: x-values for 15443 voxels,
versus distance d from back of brain



$g(\theta_0|x_0=2.3)$ for decreasing sibling sets:
(1) All (2) d in 40-80 (3) 50-70 (4) 55-65 (5) only $d=60$



REFERENCES

- Efron, B. (2011). Tweedie's formula and selection bias. *J. Amer. Statist. Assoc.* 106: 1602–1614.
- Efron, B. (2016). Empirical Bayes deconvolution estimates. *Biometrika* 103: 1–20.
- Good, I. and Toulmin, G. (1956). The number of new species, and the increase in population coverage, when a sample is increased. *Biometrika* 43: 45–63.
- Jiang, W. and Zhang, C.-H. (2009). General maximum likelihood empirical Bayes estimation of normal means. *Ann. Statist.* 37: 1647–1684.
- Robbins, H. (1951). Asymptotically subminimax solutions of compound statistical decision problems. In *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability, 1950*. UC Press, 131–148.
- Robbins, H. (1956). An empirical Bayes approach to statistics. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, 1954–1955, Vol. I*. UC Press, 157–163.