

Correlated z -Values and the Accuracy of Large-Scale Statistical Estimates

Bradley Efron

Stanford University

Correlation and Accuracy

- **Modern Scientific Studies** N cases (genes, SNPs, pixels, ...) each with its own summary statistic “ z_i ”, $i = 1, 2, \dots, N$
- *Estimate of interest* $\hat{\theta} = s(\mathbf{z})$ [e.g., $\hat{\theta} = \#\{z_i > 3\}/N$]
- **Question** How accurate is $\hat{\theta}$?
- Easy answer if z_i 's independent (but usually not!)
- Practical answer possible if z_i 's are correlated normal variates

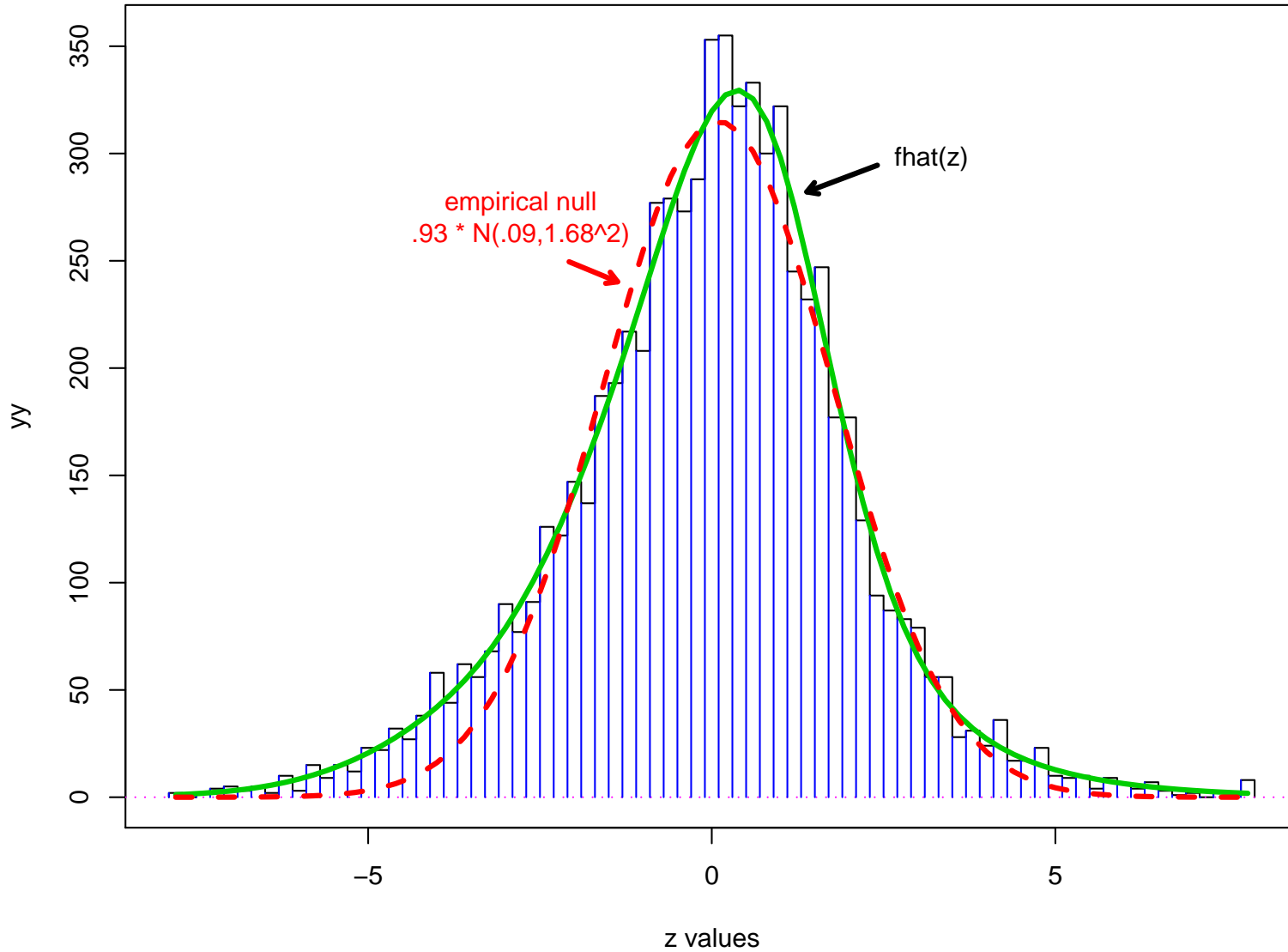
Leukemia Microarray Study

(Golub et al., 1999)

- 72 leukemia patients: $n_1 = 47$ “ALL”, $n_2 = 25$ “AML”
 - Acute Lymphoblastic Leukemia
 - Acute Myeloid Leukemia
- $N = 7128$ genes
- t_i = two-sample z -statistic comparing AML with ALL patient expression levels for gene i
- $z_i = \Phi^{-1}F_{70}(t_i)$ [Φ, F_{70} cdf's $\mathcal{N}(0, 1), t_{70}$ distributions]

$$H_0 : z_i \sim \mathcal{N}(0, 1) \quad (\text{“theoretical null”})$$

z values for 7128 genes, Golub leukemia study; solid curve is fitted density $\hat{f}(z)$; dashed is empirical null $.93 * N(.09, 1.68^2)$




Accuracy of Empirical CDF

- $\hat{F}(x) = \#\{z_i > x\}/N$
- How accurate?

$x:$	1	2	3	4	5
$\hat{F}(x)$.29	.13	.057	.025	.010
$1000 \cdot \widehat{sd}$	17.1	21.5	10.1	4.0	1.9
$1000 \cdot \widehat{sd}_0$	5.4	4.0	2.7	1.8	1.2
$1000 \cdot \widehat{sd}_{perm}$	20.9	9.9	1.4	0.1	.000

CDF Accuracy Formula

$$\text{Var} \{ \hat{F}(x) \} \doteq \underbrace{\left\{ \frac{\hat{F}(x)(1-\hat{F}(x))}{N} \right\}}_{\text{independence}} + \underbrace{\left\{ \frac{\hat{\sigma}_0^2 \hat{\alpha} \hat{f}^{(1)}(x)}{\sqrt{2}} \right\}^2}_{\text{correlation penalty}}$$


- $\hat{\sigma}_0 = 1.68$ from empirical null
- $\hat{\alpha} = .11$ estimated RMS correlation
- $\hat{f}^{(1)}(x)$ first derivative of estimate $\hat{f}(x)$
- Depends on normality: $z_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$

Digression: The Non-Null Distribution of z -Values

- z -value is a test statistic $\sim \mathcal{N}(0, 1)$ under H_0
- *Theorem* Under reasonable conditions the non-null distribution of z is

$$z \sim \mathcal{N}(\mu, \sigma^2) + O_p(1/n)$$

where

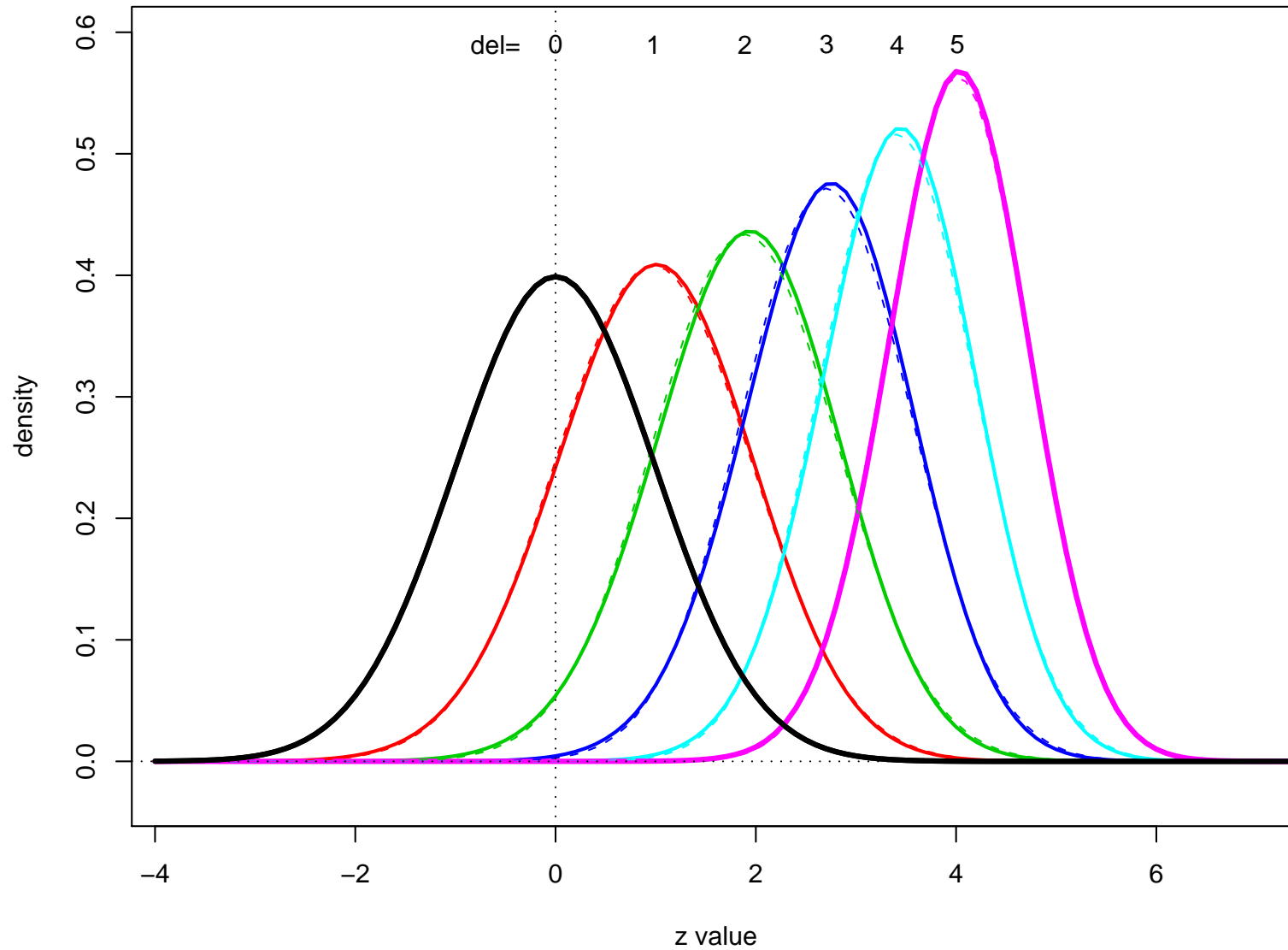
$$\sigma^2 = 1 + O\left(1/n^{\frac{1}{2}}\right)$$

- Normality degrades more slowly than unit standard deviation
- Helps justify model $z_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$

Student- t z -Values

- $t \sim t_\nu(\delta)$ [noncentral- t , noncentrality δ , $df = \nu$]
- $H_0 : \delta = 0$
- $z = \Phi^{-1}F_\nu(t)$ [F_ν central t cdf, $df = \nu$]
so under H_0 , $z \sim \mathcal{N}(0, 1)$
- What if $\delta \neq 0$?

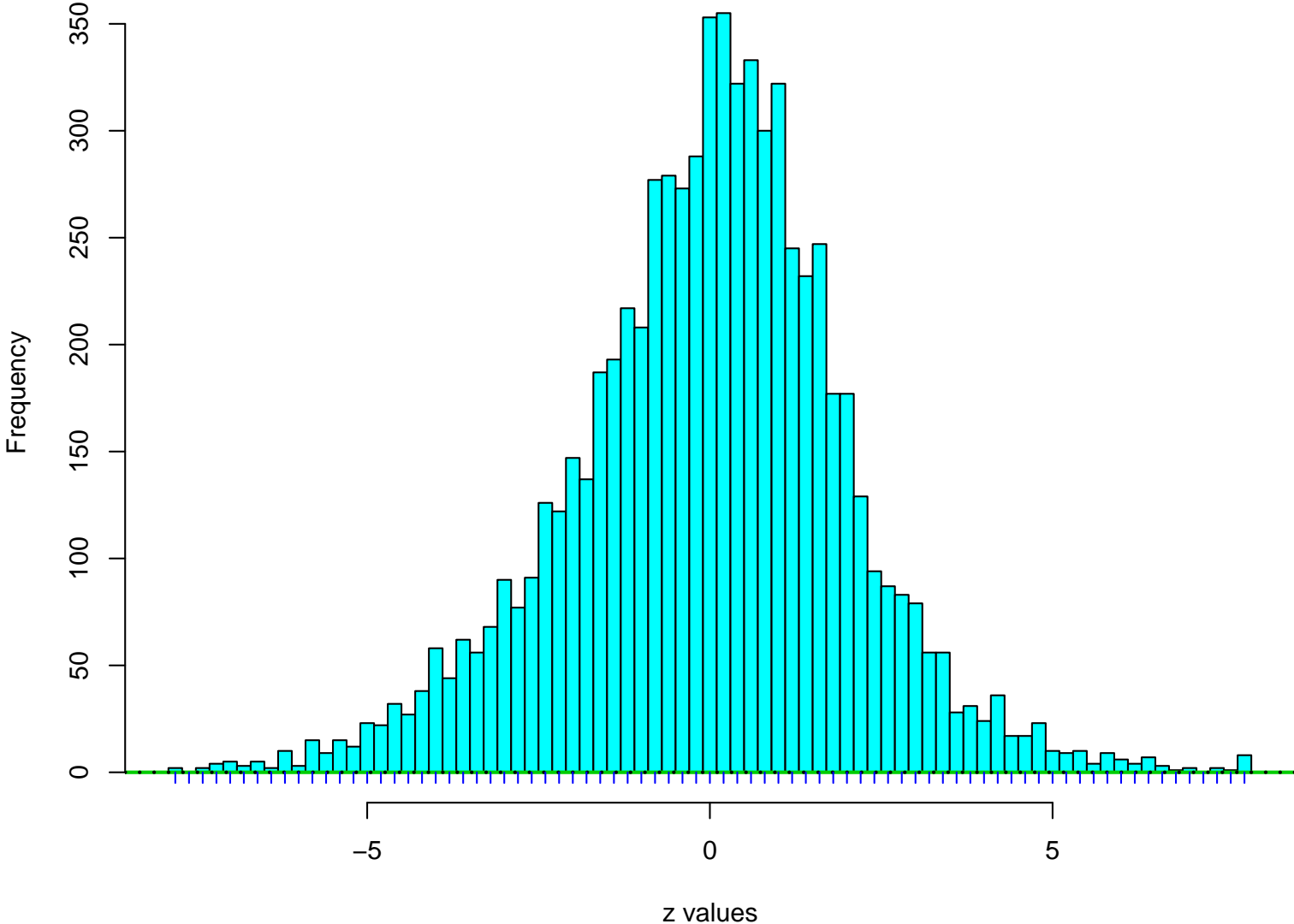
Densities for $z = \Phi^{-1}(F_{\nu}(t))$, $t \sim t(\text{del}, \nu=20)$, for $\text{del}=0,1,2,3,4,5$; Dotted dashed lines are matching $N(M, SD)$



The Count Vector \mathbf{y}

- Partition range \mathcal{Z} of z into K bins: $\mathcal{Z} = \bigcup_{k=1}^K \mathcal{Z}_k$
- Each bin of width “ Δ ”
- Bin centers “ x_k ”, $k = 1, 2, \dots, K$
(Leukemia histogram: $\mathcal{Z} = [-7.9, 7.9]$, $\Delta = .2$, $K = 79$)
- Counts $y_k = \#\{z_i \in \mathcal{Z}_k\}$ • $\mathbf{y} = (y_1, y_2, \dots, y_K)'$
- Count vector \mathbf{y} is discretized order statistic of z
(most statistics of interest of form $\hat{\theta} = m(\mathbf{y})$)

**Leukemia data, K=79 bins of width $\Delta=0.2$, $[-7.9, 7.9]$;
Bar heights are $y[k]$ values; dashes show bin centers $x[k]$**



Multi-Class Model

- Suppose z_i 's are in "classes" C_1, C_2, \dots, C_C , with

$$z_i \sim \mathcal{N}(\mu_c, \sigma_c^2) \quad \text{for } z_i \in C_c$$

- $N_c = \# \{C_c\}$, $p_c = N_c/N$ [so $\sum_c N_c = N$, $\sum_c p_c = 1$]
- *Correlation distribution* $g_{cd}(\rho)$ = empirical density of $N_c \cdot N_d$ correlations between members of C_c, C_d
- Likewise $g_{cc}(\rho)$ for $N_c(N_c - 1)/2$ in C_c
- Assume g_{cd} all equal $g(z)$
[actually only need 2nd moments]

Crucial Definitions

- $x_{kc} = \frac{x_k - \mu_c}{\sigma_c}$ “Adjusted bin center”
- $\pi_{kc} = \Delta \cdot \varphi(x_{kc}) / \sigma_c \doteq \text{Prob}_c \{z_i \in \mathcal{Z}_k\}$ $[\varphi(u) = e^{-u^2/2} / \sqrt{2\pi}]$
- *Expectation of \mathbf{y} :* $E\{\mathbf{y}\} = N \sum_c p_c \boldsymbol{\pi}_c$ $[\boldsymbol{\pi}_c = (\dots, \pi_{kc}, \dots)']$
- $\varphi_\rho(u, v) =$ standard bivariate normal density, correlation ρ
- $\lambda_\rho(u, v) = \frac{\varphi_\rho(u, v)}{\varphi(u)\varphi(v)} - 1$ $[= 0 \text{ if } \rho = 0]$

$$\lambda(u, v) = \int_{-1}^1 \lambda_\rho(u, v) g(\rho) d\rho$$

Exact Covariance of \mathbf{y}

Theorem Under the multi-class model,

$$\text{cov}(\mathbf{y}) = \mathbf{cov}_0 + \mathbf{cov}_1,$$

$$\mathbf{cov}_0 = N \sum_c p_c \{ \text{diag}(\boldsymbol{\pi}_c) - \boldsymbol{\pi}_c \boldsymbol{\pi}_c' \} \quad [\textit{independence}]$$

and

$$\mathbf{cov}_1 = N^2 \sum_c \sum_d p_c p_d \mathbf{B}_{cd} - N \sum_c p_c \mathbf{B}_{cc} \quad \langle 1 \rangle [\textit{corr penalty}]$$

where $B_{cd}(k, l) = \pi_{kc} \pi_{ld} \lambda(x_{kc}, x_{ld})$.

The second term of $\langle 1 \rangle$ is negligible for N large.

Mehler's Identity (Lancaster, 1958)

- $\lambda_\rho(u, v) = \sum_{j \geq 1} \frac{\rho^j}{j!} h_j(u) h_j(v)$ [$h_j = j$ th Hermite polynomial]

- So $\int_{-1}^1 \lambda_\rho(u, v) g(\rho) d\rho = \sum_{j \geq 1} \frac{\alpha_j}{j!} h_j(u) h_j(v)$ [$\alpha_j = \int_{-1}^1 \rho^j g(\rho) d\rho$]

and

$$B_{cd}(x_{kc}, x_{ld}) = \pi_{kc} \pi_{ld} \sum_{j \geq 1} \frac{\alpha_j}{j!} h_j(x_{kc}) h_j(x_{ld})$$

$$= \Delta^2 \sum_{j \geq 1} \frac{\alpha_j}{j!} \varphi^{(j)}(x_{kc}) \varphi^{(j)}(x_{ld}) / \sigma_c \sigma_d \quad \langle 2 \rangle$$

- $\varphi^{(j)}$ the j th derivative of φ [using $\varphi h_j = (-1)^j \varphi^{(j)}$]

Three Simplifications of **cov**

- Drop second term of $\langle 1 \rangle$
- Microarray standardization methods make $\alpha_1 \doteq 0$ in $\langle 2 \rangle$
- Now $\alpha_2 = \int_{-1}^1 \rho^2 g(\rho)$ is the lead term
- Higher terms ignorable if α_2 small

Simplified Formula

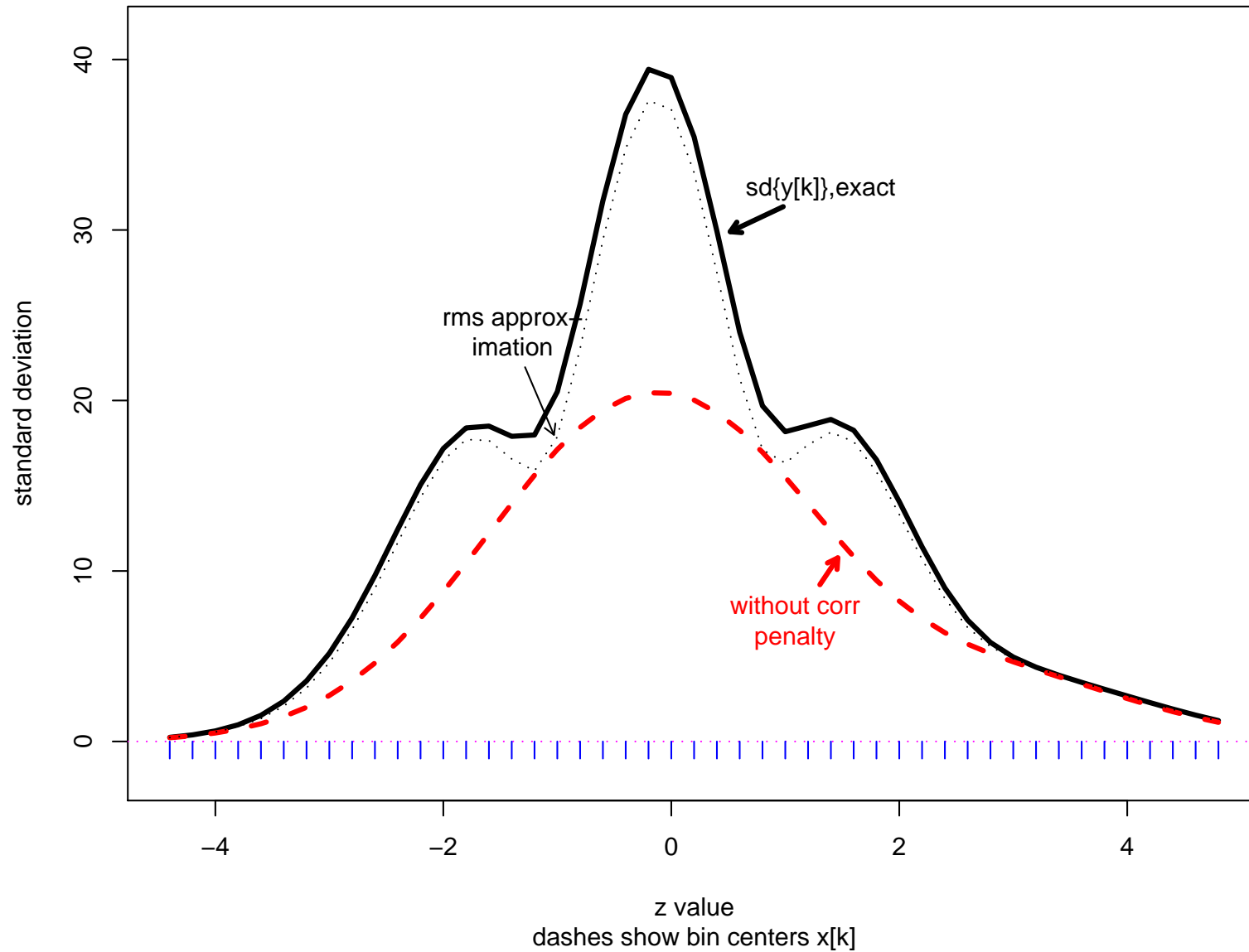
Letting $\alpha = \alpha_2^{\frac{1}{2}}$ and $\bar{\phi}_k^{(2)} = \sum_c p_c \varphi^{(2)} \left(\frac{x_{kc} - \mu_c}{\sigma_c} \right) / \sigma_c$

$$\mathbf{cov}_1 \doteq (N\Delta\alpha)^2 \bar{\phi}^{(2)} \bar{\phi}^{(2)' / 2}$$

Numerical Comparison

- $N = 6000, \alpha = .1$
- Two classes: $(p_c, \mu_c, \sigma_c) = \begin{cases} (.95, 0, 1) \\ (.05, 2.5, 1) \end{cases}$
- Next figure compares standard deviations (square roots diagonal elements) of exact $\mathbf{cov}(\mathbf{y})$ and rms approximation

Compare $\text{sd}\{y[k]\}$ from exact formula (solid) with rms approx (dashed);
 $N=6000$, $\alpha=.1$, $(p_0, \mu_0, \text{sig}_0)=(.95, 0, 1)$ and $(.05, 2.5, 1)$



CDF Accuracy (Owen, 2005)

- Let $\hat{F}_k = \#\{z_i \text{ in bins } \geq k\}/N$ [right-sided empirical cdf]

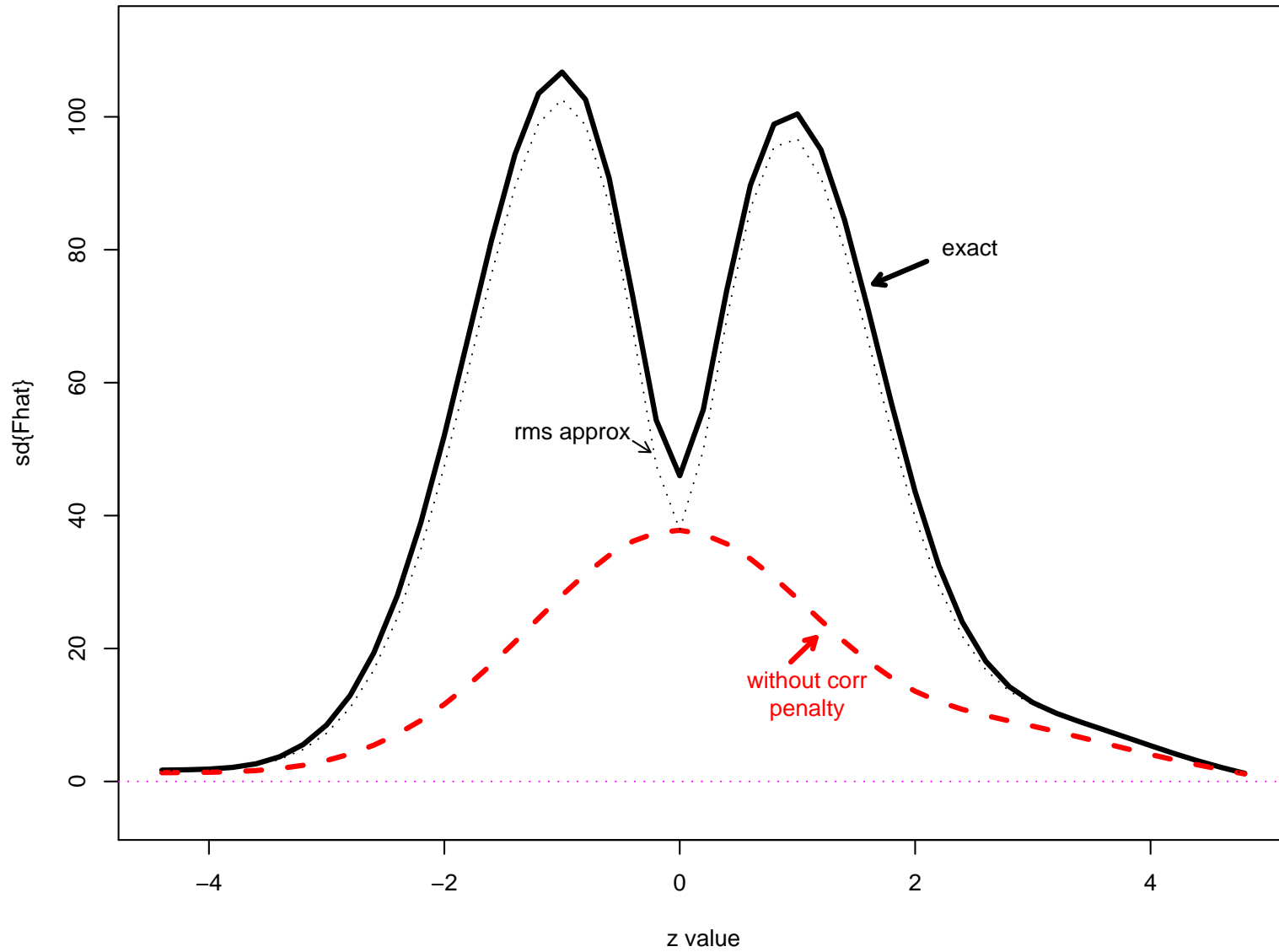
so
$$\hat{\mathbf{F}} = \frac{1}{N} \mathbf{M} \mathbf{y} \quad \text{where} \quad M_{jk} = \begin{cases} 1 & \text{if } j \leq k \\ 0 & \text{if } j > k \end{cases}$$

- Then $\mathbf{Cov}(\hat{\mathbf{F}}) = \mathbf{Cov}_0 + \mathbf{Cov}_1$ where \mathbf{Cov}_0 from usual multinomial, and

$$\mathbf{Cov}_1 \doteq \frac{\alpha^2}{2} \bar{\boldsymbol{\varphi}}^{(1)} \bar{\boldsymbol{\varphi}}^{(1)'}$$

$$\left[\bar{\varphi}_k^{(1)} = \sum_c p_c \varphi^{(1)} \left(\frac{x_k - \mu_c}{\sigma_c} \right), \quad \varphi^{(1)}(u) = -\varphi(u)u \right]$$

Same numerical example, now $\text{sd}\{\hat{F}[k]\}$



Practical Implementation

- **Cov**₁ requires estimating α and $\bar{\varphi}^{(1)}$
- $\alpha = \text{rms correlation}$: not difficult in microarray studies
- $\bar{\varphi}_k^{(1)} = \sum_c p_c \varphi^{(1)} \left(\frac{x_k - \mu_c}{\sigma_c} \right)$: seems to require class parameters (p_c, μ_c, σ_c) , but can be finessed in some situations

Estimation of RMS Correlation α

- $\hat{\rho}_{ii'}$ = empirical correlation, rows i, i' of \mathbf{X} ,
 $N \times n$ expression matrix
- $\{\hat{\rho}_{ii'}\}$ has mean and variance (m, v)
[leukemia = (.00, .19²)]

$$\hat{\alpha}^2 = \frac{n}{n-1} \left(v - \frac{1}{n-1} \right)$$

	ALL	AML	Both
$\hat{\alpha}$:	.121	.109	.114

Estimation of $\bar{\varphi}_k^{(1)}$

- *Multi-Class Model:* $f(x_k) = \sum_c p_c \varphi\left(\frac{x_k - \mu_c}{\sigma_c}\right) / \sigma_c$, so

$$f^{(1)}(x_k) = \sum_c p_c \varphi^{(1)}\left(\frac{x_k - \mu_c}{\sigma_c}\right) / \sigma_c^2$$

- *If $\sigma_c \equiv \sigma_0$* $\sigma_0^2 f^{(1)}(x_k) = \sum_c p_c \varphi'\left(\frac{x_k - \mu_c}{\sigma_c}\right) = \bar{\varphi}_k^{(1)}$
- $y \rightarrow \hat{f} \rightarrow \hat{f}^{(1)} \rightarrow \bar{\varphi}^{(1)}$ [★ formula correlation penalty]

More General Accuracy Estimates

- “ Q ” q -dimensional statistic of interest: $Q = Q(\mathbf{y})$
- *Influence Function*

$$dQ = \hat{D} d\mathbf{y} \quad \left[\hat{D}_{jk} = \partial Q_j / \partial y_k \right]$$

$$\widehat{\text{cov}}(Q) = \hat{D} \text{cov}(\mathbf{y}) \hat{D}'$$

Example: Accuracy of \hat{f}

- $z \rightarrow y \rightarrow \hat{f}$ by Poisson GLM
of counts y_k on polynomial (x_k)
- $Q = \log(\hat{f}) = (\dots \log f(x_k) \dots)'$

- $\hat{D} = M [M' \text{diag}(\hat{f}) M] M' / N\Delta$

with M the GLM structure matrix (Efron, 2007b)

The Local False Discovery Rate

- Bayes: $\left\{ \begin{array}{ll} p_0 = \Pr\{\text{null}\} & f_0(z) = \text{null density} \\ p_1 = \Pr\{\text{non-null}\} & f_1(z) = \text{non-null density} \end{array} \right\}$

$$f(z) = p_0 f_0(z) + p_1 f_1(z)$$

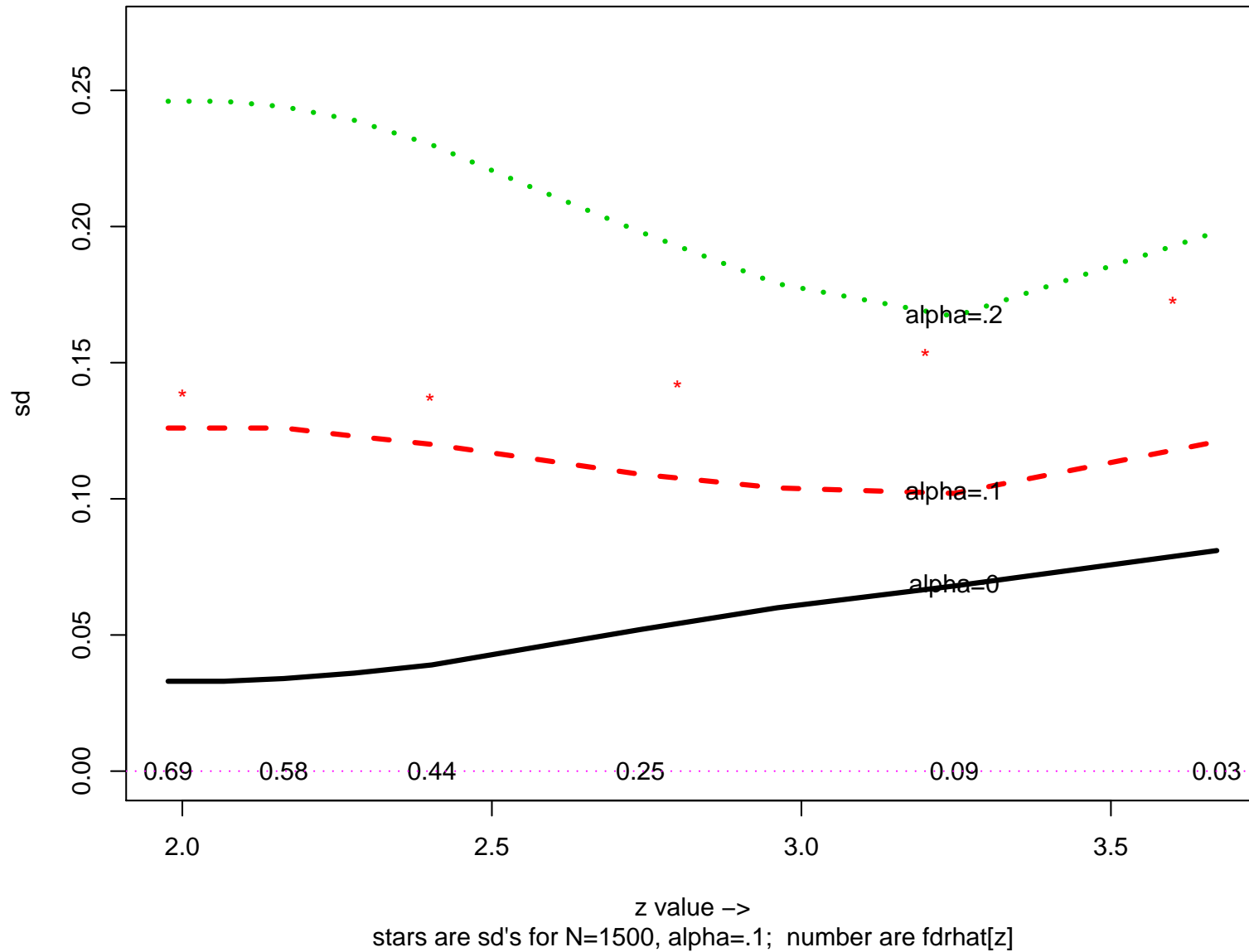
- *Local False Discovery Rate*

$$\text{fdr}(z) = p_0 f_0(z) / f(z) = \Pr\{\text{null}|z\}$$

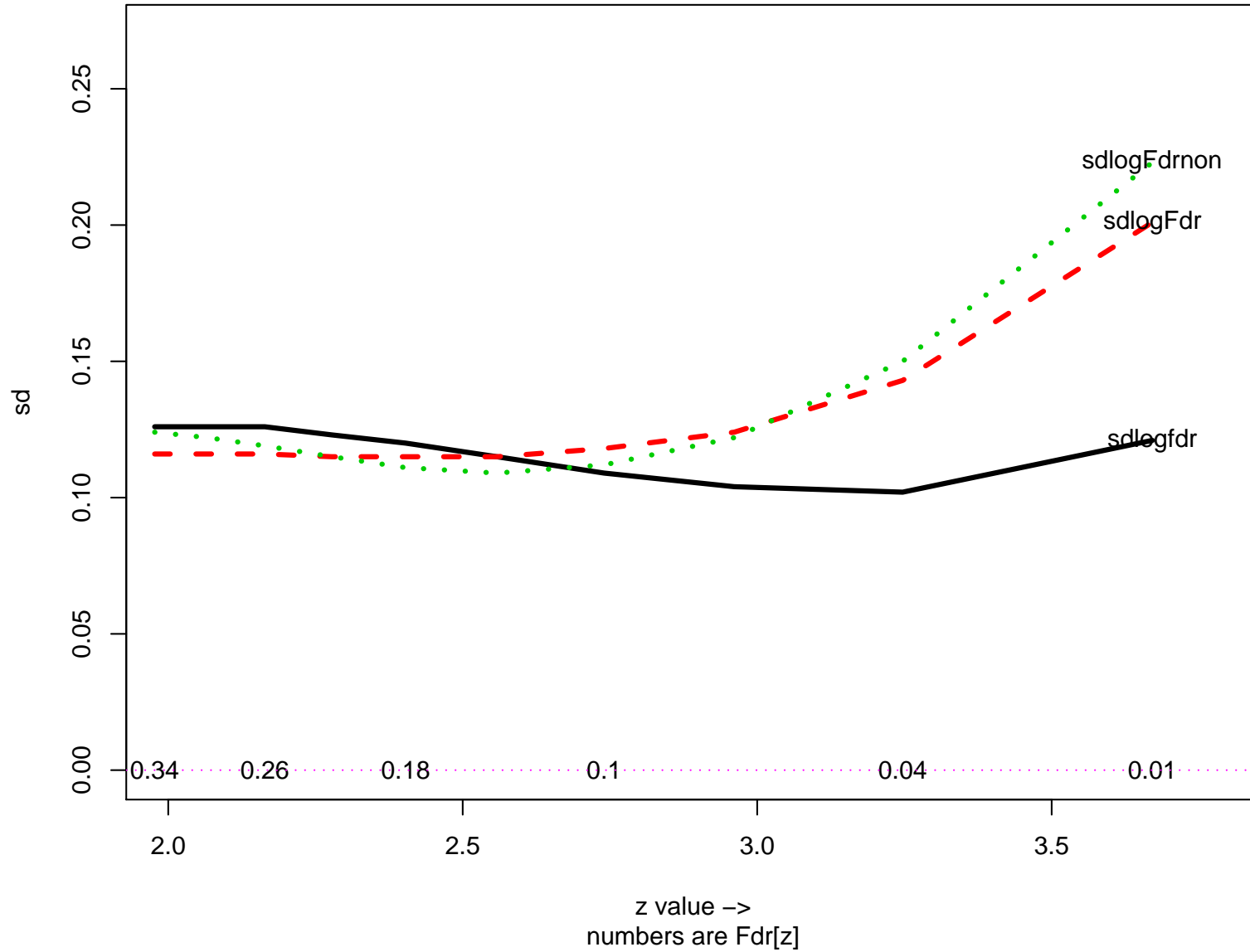
- $\widehat{\text{fdr}} = p_0 \widehat{f}_0 / \widehat{f}$ $\widehat{\text{lfdr}} = \log(\widehat{\text{fdr}}) = \log(p_0 \widehat{f}_0) - \log \widehat{f}$

- If $p_0 \widehat{f}_0$ known then $\widehat{\text{cov}}(\widehat{\text{lfdr}}) = \widehat{\text{cov}}(\log \widehat{f})$

sd{log fdrhat(z)} ; N=6000, alpha=0, .1, and .2,
 (p0,mu,sig) = (.95,0,1) and (.05,2.5,1)



Now compare sd's for $\log\{\hat{fdr}\}$ and $\log\{Fdr\}$,
 $\alpha=.1$



Accuracy When $p_0 f_0(z)$ Must Be Estimated

- Program `locfdr` estimates p_0 and f_0 from central histogram counts, the “empirical null”
- Section 5 of Efron (2007b) gives influence functions \hat{D}
- *Next Table* $\text{sd}\{\log \widehat{\text{fdr}}(z)\}$ using “central matching” method, for $\alpha = .1$ case, as in Figures

$$\text{sd}\{\log \widehat{\text{fdr}}(z)\}$$

z:	1.98	2.16	2.40	2.74	3.25
empirical null:	.18	.26	.36	.54	.83
theoretical null:	.13	.13	.12	.11	.10
fdr:	.69	.58	.44	.25	.09

- Still worse for $\log \widehat{\text{Fdr}}$
- Better using MLE option
- Some of increase is “signal” rather than noise (Efron, 2007a)

References

Csörgő, S. and Mielniczuk, J. (1996). The empirical process of a short-range dependent stationary sequence under Gaussian subordination. *Probab. Theory Related Fields* 104: 15–25.

Efron, B. (2007a). Correlation and large-scale simultaneous significance testing. *J. Amer. Statist. Assoc.* 102: 93–103.

Efron, B. (2007b). Size, power and false discovery rates. *Ann. Statist.* 35: 1351–1377.

Golub, T., Slonim, D. K. and Tamayo, P. et al. (1999). Molecular Classification of Cancer: Class Discovery and

Class Prediction by Gene Expression Monitoring. *Science* 286: 531–537.

Lancaster, H. O. (1958). The structure of bivariate distributions. *Ann. Math. Statist.* 29: 719–736.

Owen, A. B. (2005). Variance of the number of false discoveries. *J. R. Stat. Soc. Ser. B Stat. Methodol.* 67: 411–426.

Qiu, X., Klebanov, L. and Yakovlev, A. (2005). Correlation between gene expression levels and limitations of the empirical Bayes methodology for finding differentially expressed genes. *Stat. Appl. Genet. Mol. Biol.* 4: Art. 34, 32 pp. (electronic).