

False Discovery Rates and Copy Number Variation

Bradley Efron and Nancy Zhang

Stanford University

Three Statistical Centuries

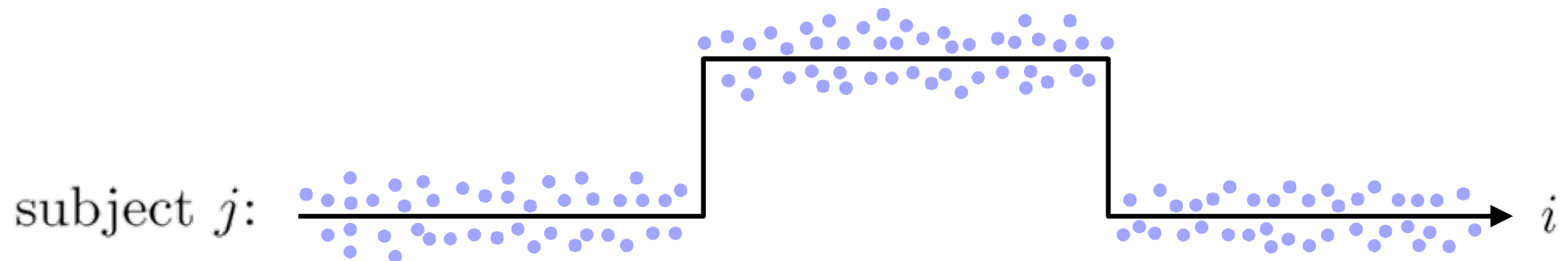
- **19th** (*Quetelet*) Huge data sets, simple questions
- **20th** (*Fisher, Neyman, Hotelling, . . .*) Small data sets, simple questions
- **21st** (*Scientific mass production*) Huge data sets, complicated questions

Example: Copy Number Variation

- *CNV* Gains and losses of chromosome segments (disease association)
- Instead of 2 copies, might have 0, 1, 3, 4,...
- **Data** x_{ij} = noisy msmnt of copy number for subject j at marker position i
- $i = 1, 2, \dots, N$ (5000) and $j = 1, 2, \dots, n$ (150) (< 1% of data!)
- x_{ij} approx. normal with mean 0 if copy number = 2

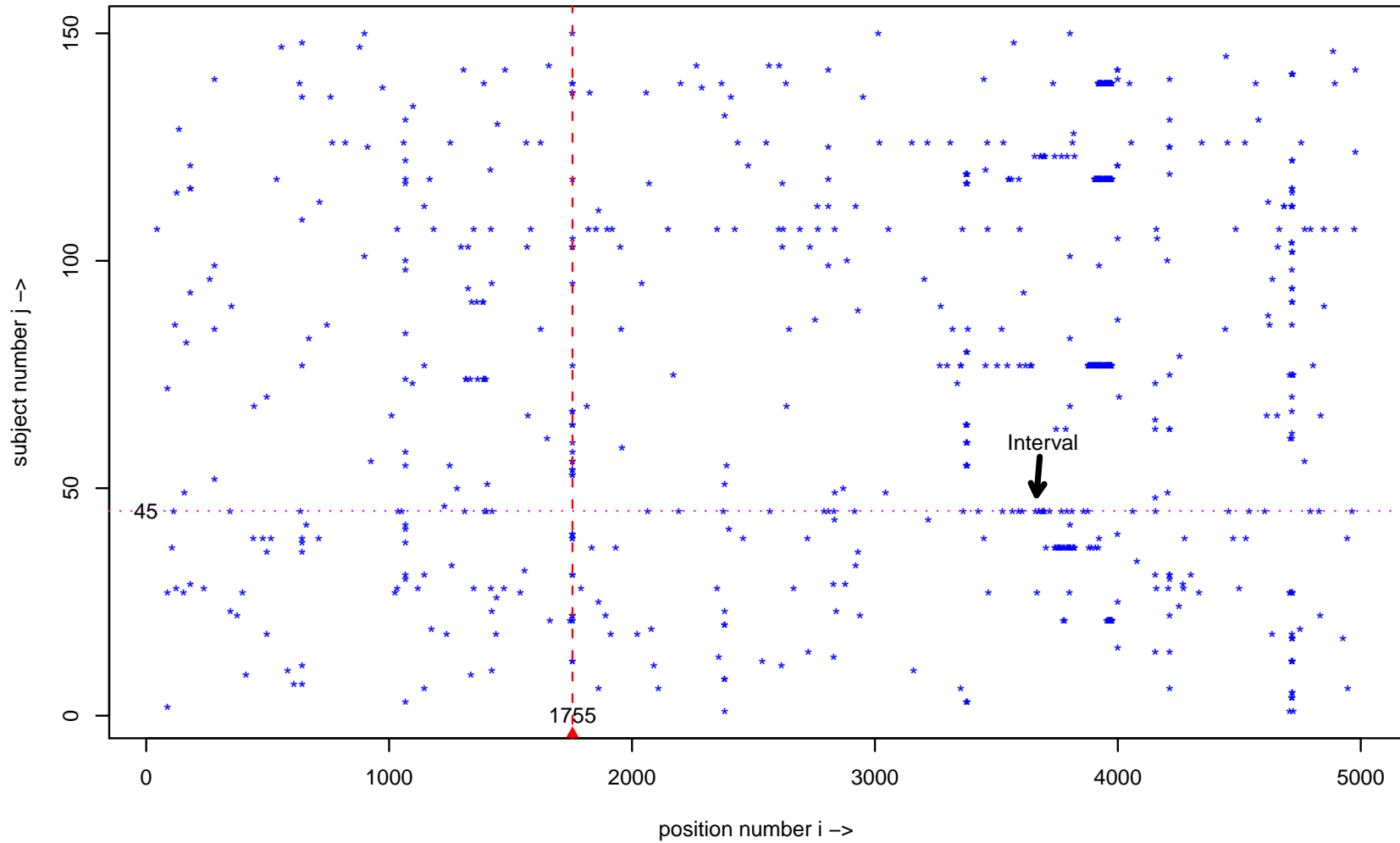
What We Expect To See

- *Hot positions*: i where several of subjects j show unusually high (or low) values x_{ij}
- For some subjects j : intervals of high (or low) values x_{ij}



- *Information* on CNV locations in both directions

Lowest .001 of the 750,000 entries $x[i,j]$;
Subject 45 shows interval of low values around position 3800;
Is position 1755 cnv prone?



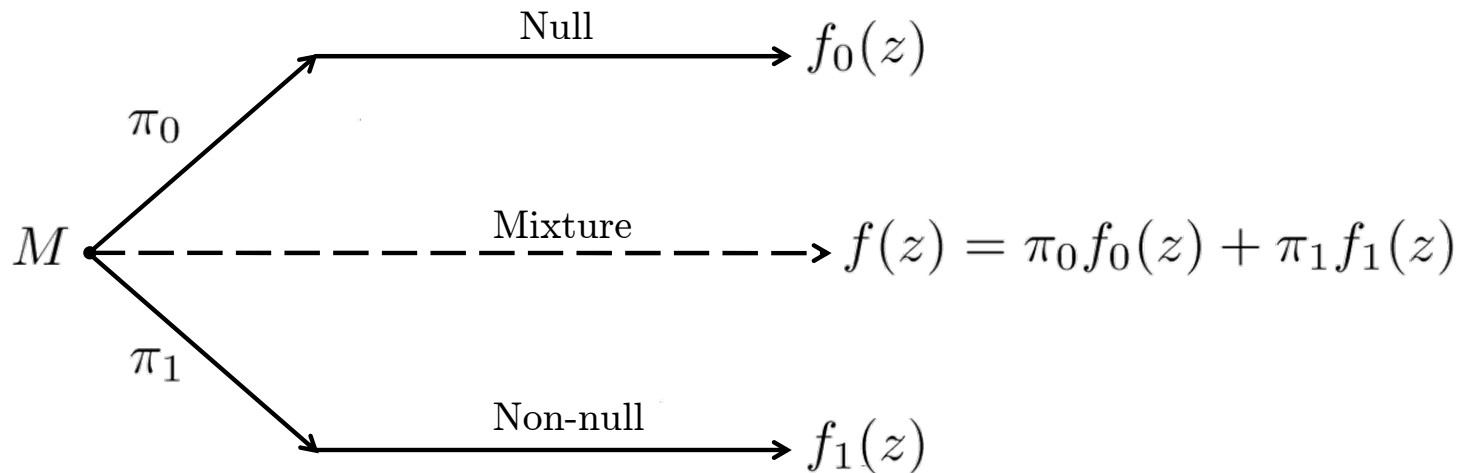
Z-Values

- \mathbf{X} : x_{ij} for $i = 1, 2, \dots, N = 5000$ positions,
 $j = 1, 2, \dots, n = 150$ subjects
- C_i = all subjects at i th position ($n = 150$)
- **Moving averages** Replace x_{ij} with $\bar{x}_{ij} = \sum_{l=i-5}^{i+5} x_{lj} / 11$
 $\Rightarrow \bar{\mathbf{X}}$ (j 's msmts averaged over nearby positions)
- *Standardize rows of $\bar{\mathbf{X}}$* $\left\{ \begin{array}{l} \text{subtract row median} \\ \text{divide by row robust standardization} \end{array} \right.$
- Gives \mathbf{Z} matrix $z_{ij} \Rightarrow$ iterative $\widehat{\text{fdr}}_i$ estimates

Simultaneous Hypothesis Testing

- M null hypotheses $H_{01}, H_{02}, \dots, H_{0M}$
($M = 750,000$ for CNV)
- Case m has test statistic z_m , null density $f_0(z)$
- **The problem** Given $\mathbf{z} = (z_1, z_2, \dots, z_M)$,
simultaneously test all M null hypotheses — and don't
make many mistakes!

The Bayesian Two-Groups Model



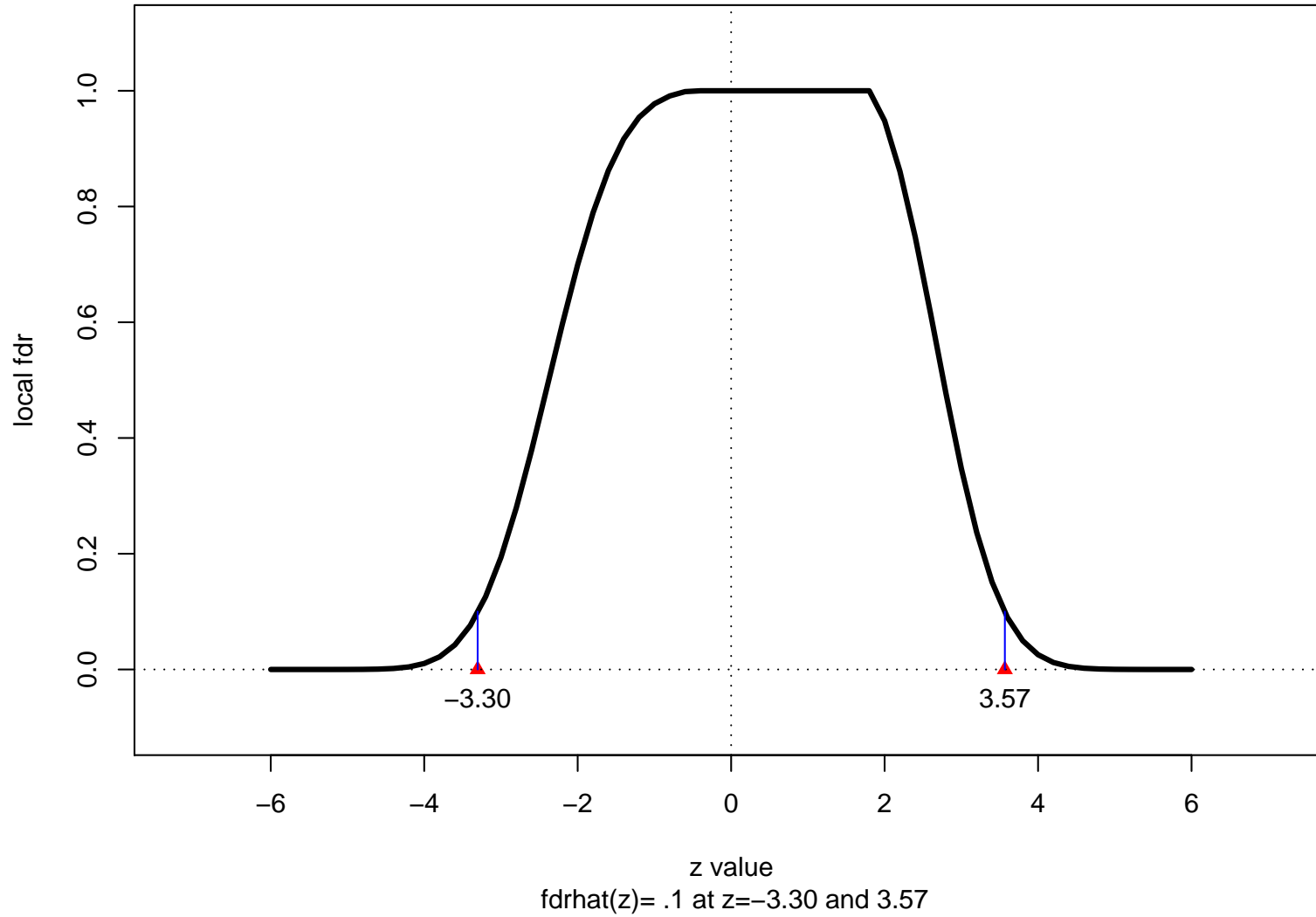
- Local false discovery rate

$$\text{fdr}(z) = \Pr\{\text{null}|z\} = \pi_0 f_0(z) / f(z)$$

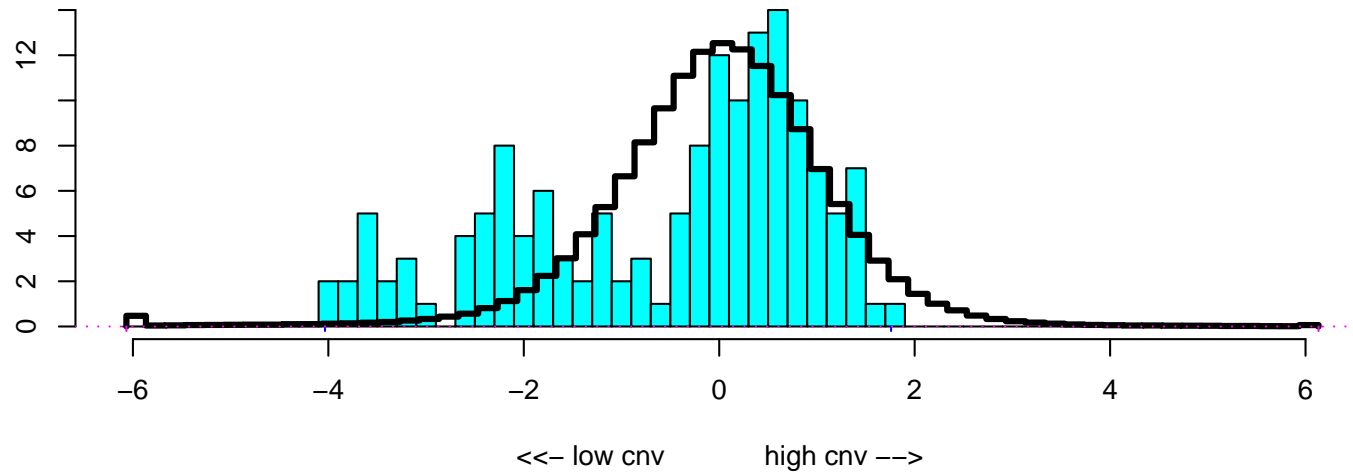
- Empirical Bayes $z \longrightarrow \hat{\pi}_0, \hat{f}_0, \hat{f} \longrightarrow \widehat{\text{fdr}}(z) = \hat{\pi}_0 \hat{f}_0(z) / \hat{f}(z)$

- Reject H_{0m} if $\widehat{\text{fdr}}(z_m)$ small (see Efron, 2008)

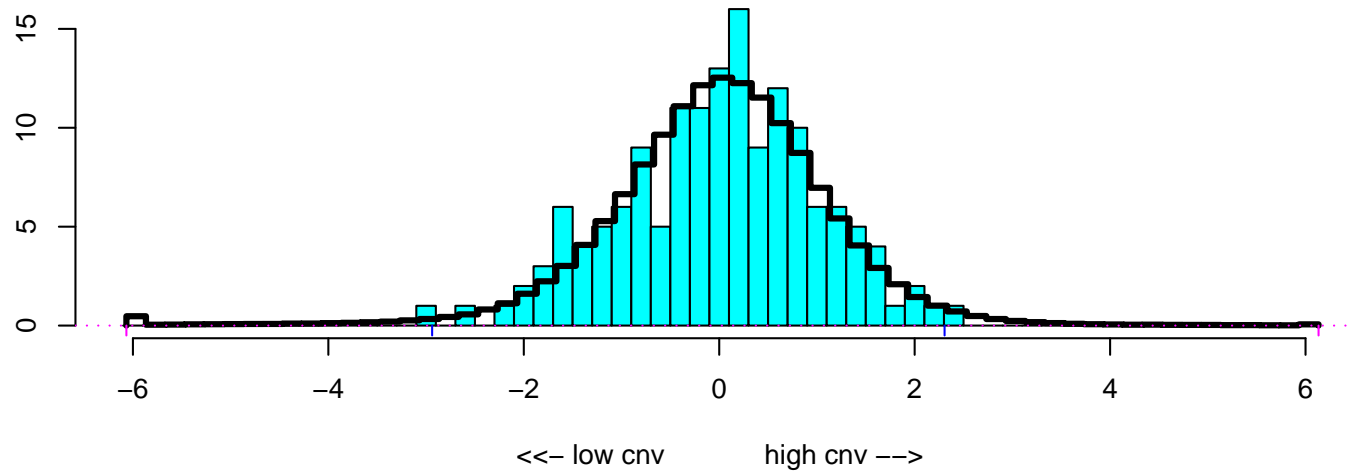
**Estimated local false discovery rate, all 750,000 z-values;
pihat0=.954, estimated null density $N(.04,.93^2)$**



z-values at position i=1755 (solid histogram) compared to all the others (line)

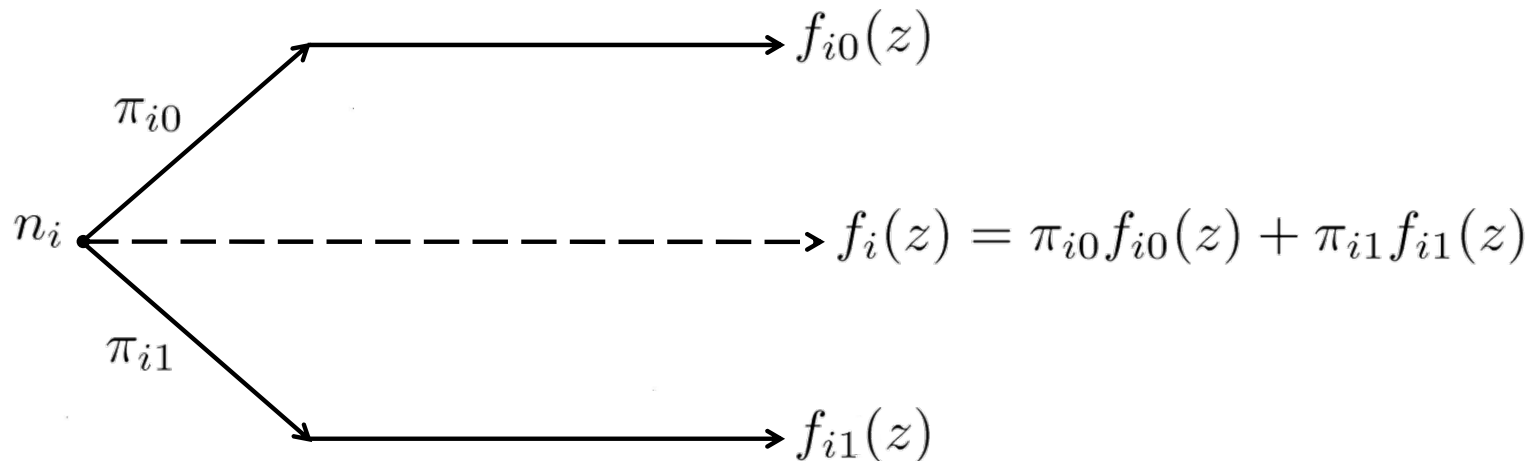


Now for position i=2000



A More General Model

- *Classes:* $C_1, C_2, \dots, C_i, \dots, C_N$ with n_i cases in C_i
- **CNV:** $C_i = i$ th column, $n_i = n = 150$
(the $n = 150$ subjects measured at position i)



- $\text{fdr}_i(z) = \pi_{i0}f_{i0}(z)/f_i(z) = \Pr\{\text{null}|z, C_i\}$

Combined and Separate Fdr's

- *Strategy*: Estimate $\text{fdr}(z) = \pi_0 f_0(z) / f(z)$ from combined data and then modify for C_i
- Assume $f_{i0}(z)$ and $f_{i1}(z)$ do not depend on i , only $\pi_{i1} = \Pr\{\text{non-null} | C_i\}$ varying across classes:

$$\text{fdr}_i(z) = \text{fdr}(z) / [1 + \text{tdr}(z)S_i]$$

- $\text{tdr}(z) \equiv 1 - \text{fdr}(z) =$ “true discovery rate”

$$\text{and } S_i = \frac{\pi_{i1}}{\pi_1} \bigg/ \frac{\pi_{i0}}{\pi_0} - 1$$

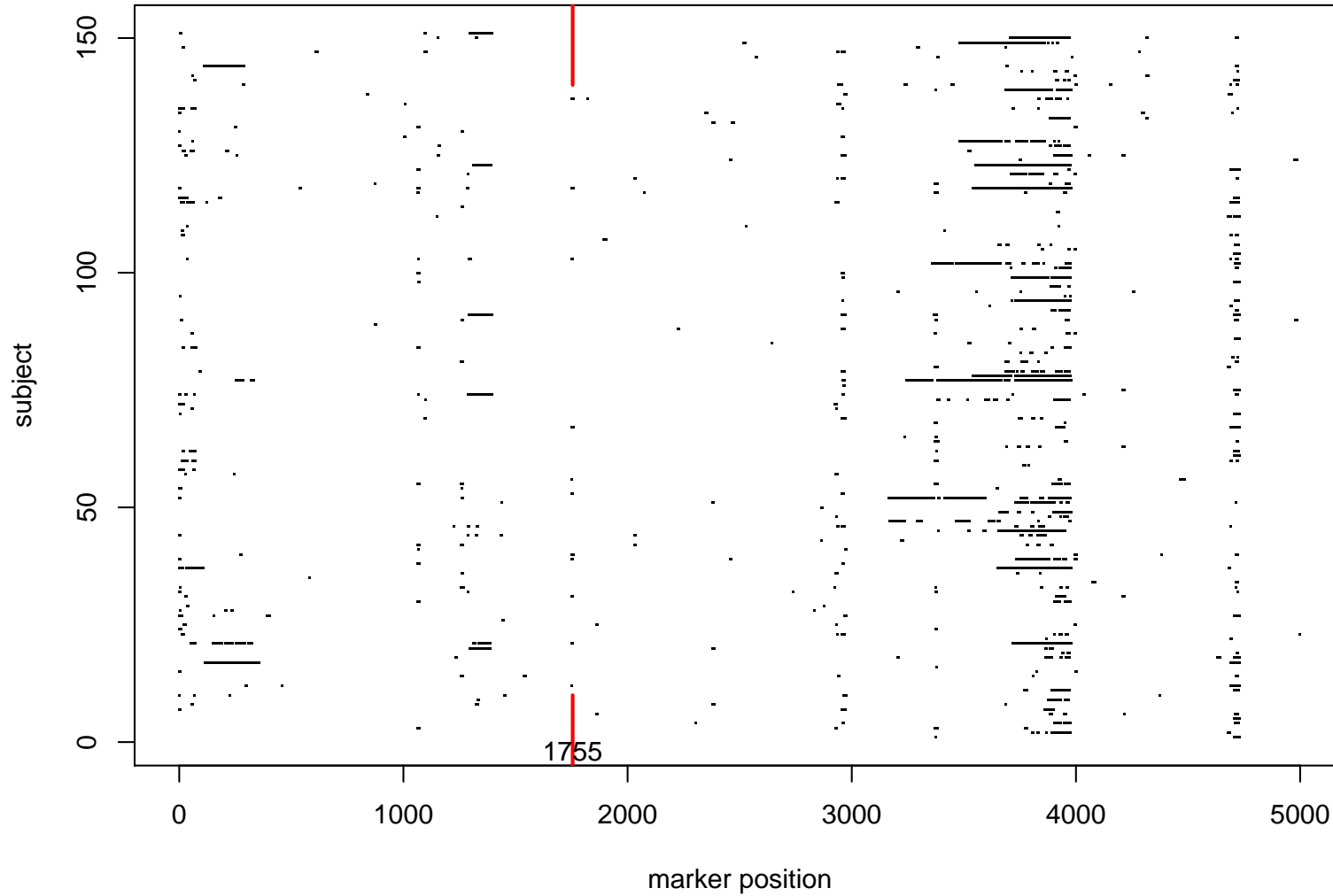
Iterative Estimation of $\text{fdr}_i(z)$ (Model 1)

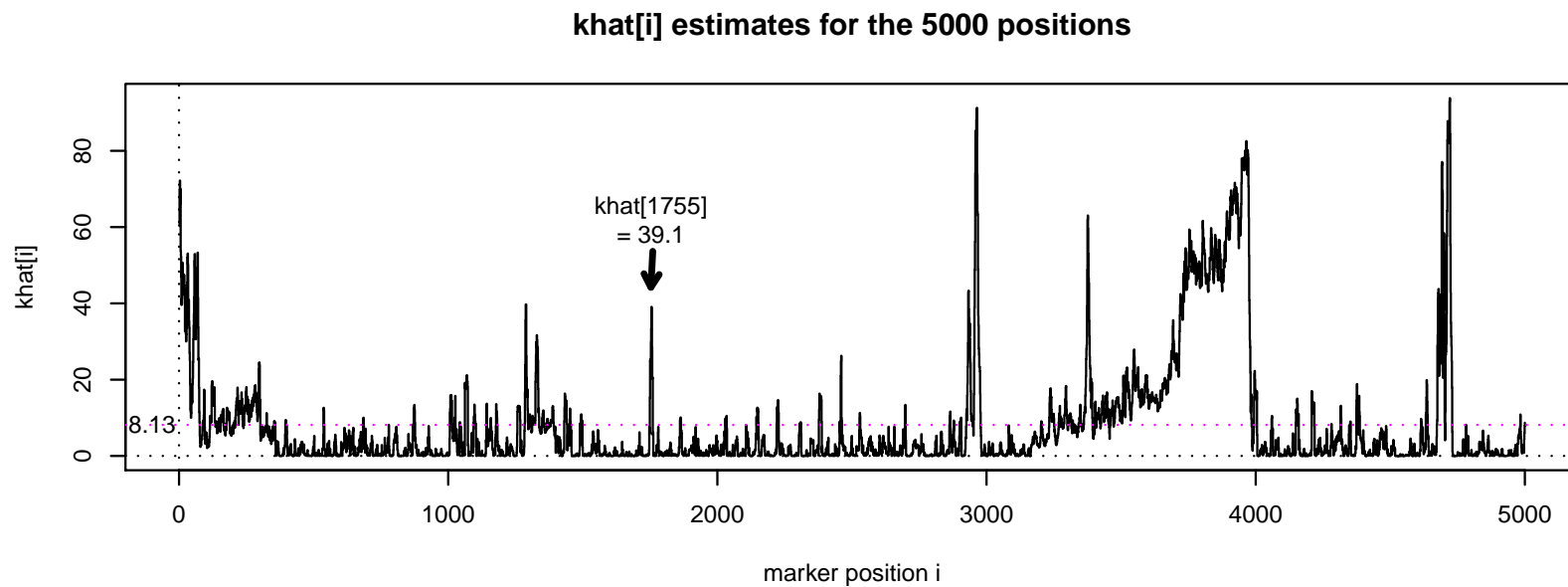
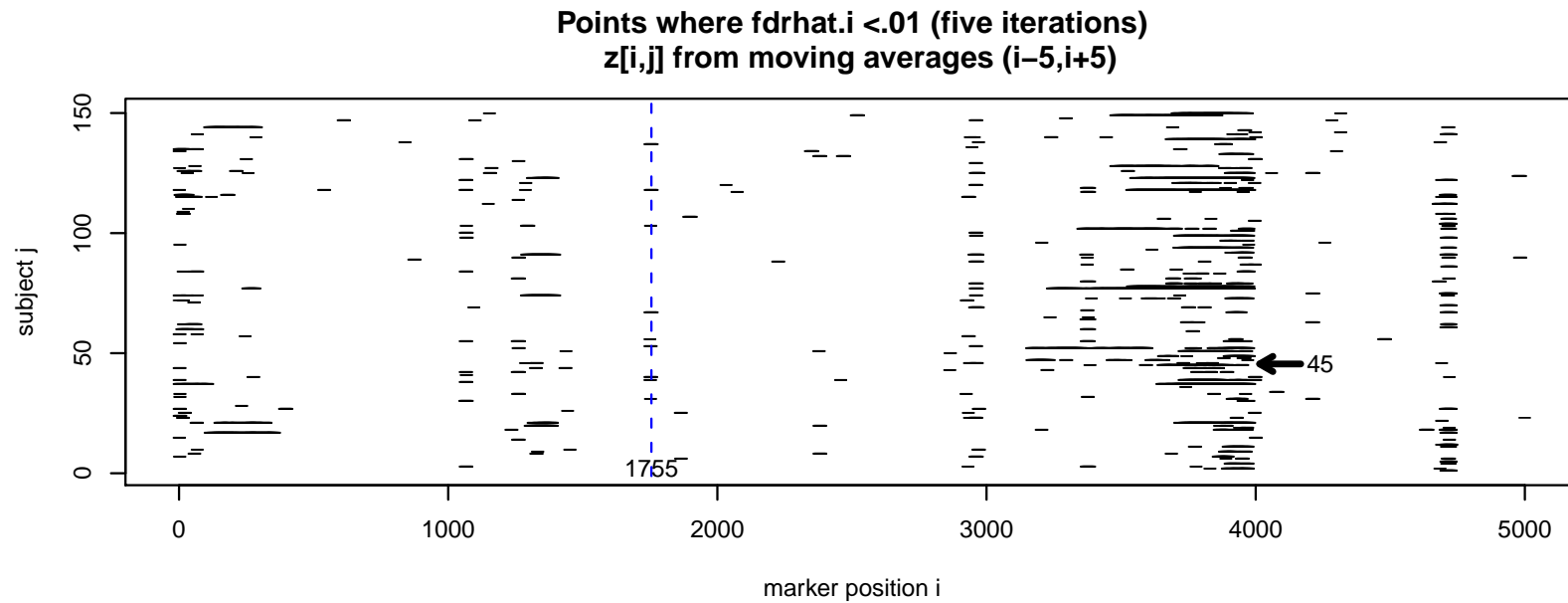
- **First** Estimate $\widehat{\text{fdr}}(z) = \hat{\pi}_0 \hat{f}_0(z) / \hat{f}(z)$ from combined data (z_1, z_2, \dots, z_M)
- If k_i non-nulls in C_i : $\hat{\pi}_{i1} = k_i/n$ gives \hat{S}_i and

$$\widehat{\text{fdr}}_i(z) = \frac{\widehat{\text{fdr}}(z)}{1 + \widehat{\text{tdr}}(z)\hat{S}_i} \quad (\widehat{\text{tdr}}_i = 1 - \widehat{\text{fdr}}_i)$$

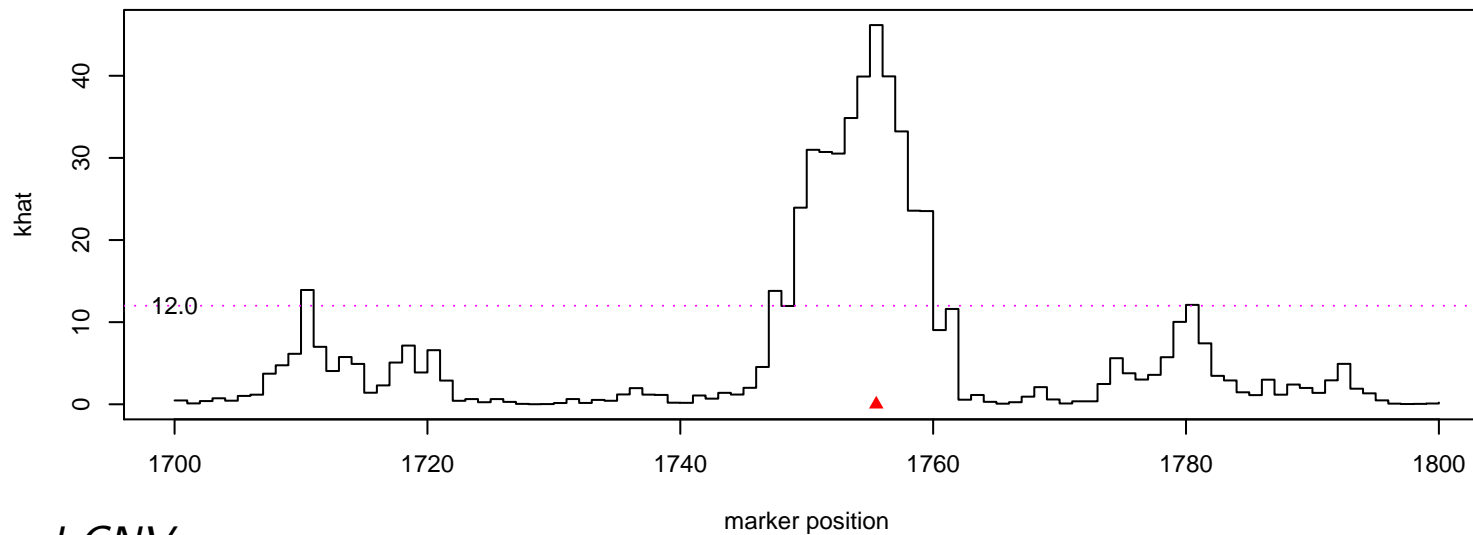
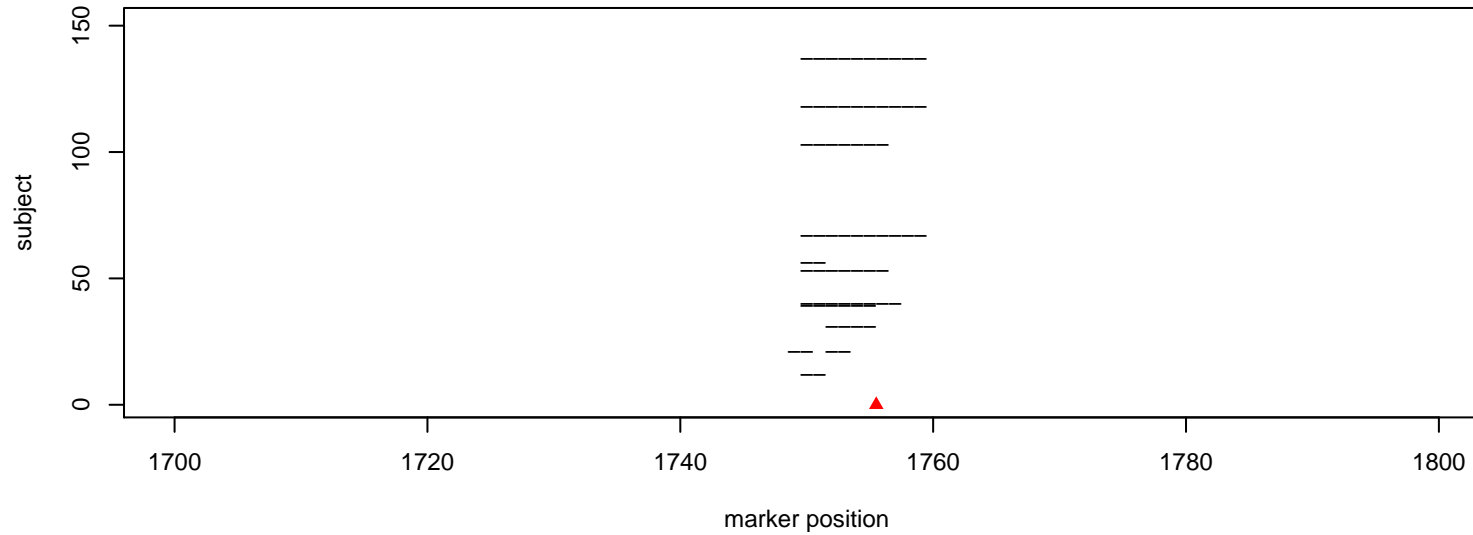
- But $\hat{k}_i = \sum_{C_i} \widehat{\text{tdr}}_i(z_m)$ estimates k_i
- Iterate! (5 cycles plenty in what follows)

Points where $\text{fdrhat} < .01$. Five iterations of Model 1,
 $z[i,j]$ from moving averages $(i-5,i+5)$





points where $fdrk < .01$; close-up positions 1700:1800;
shows possible CNV region at 1750:1759



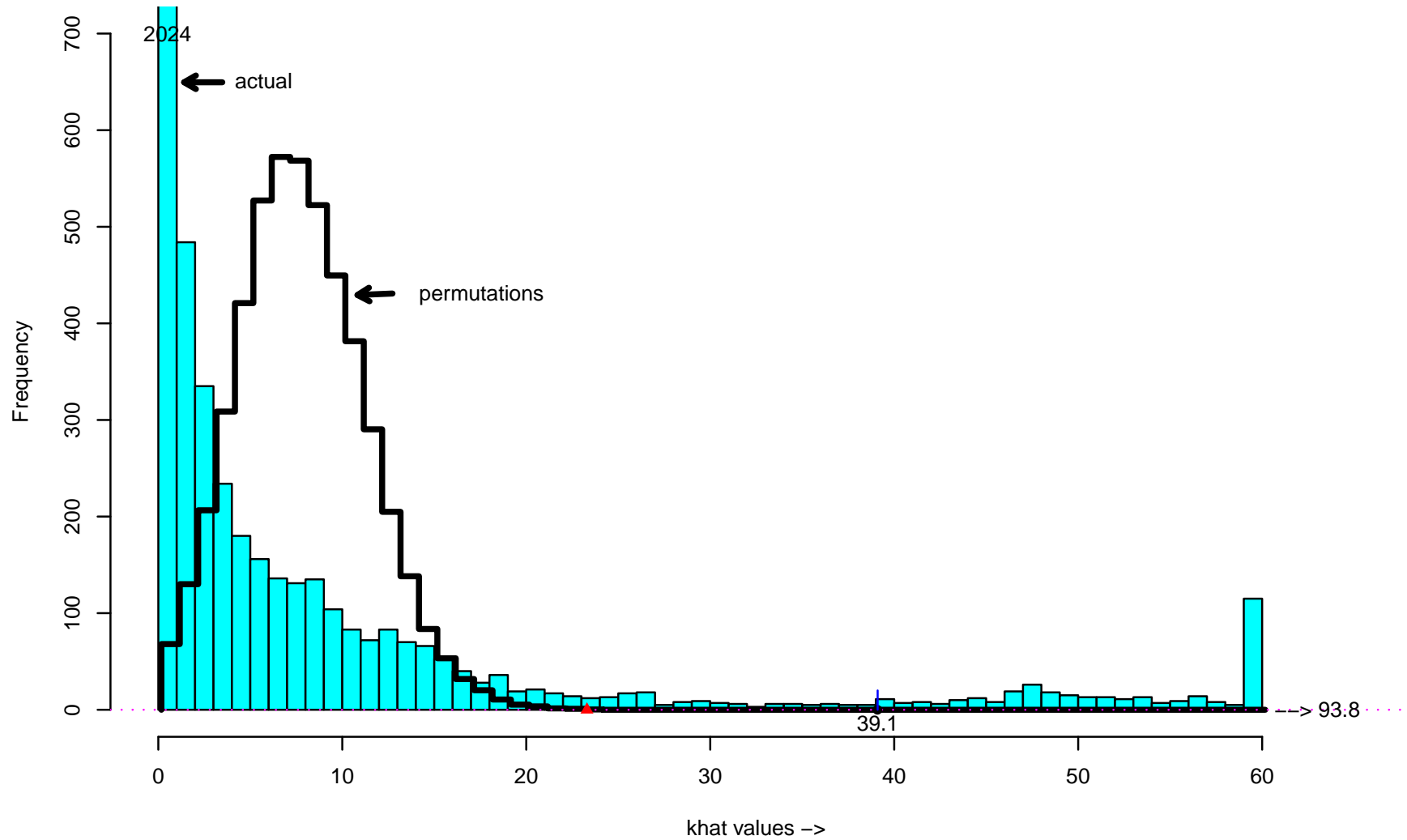
Is Position 1755 “Significant” ?

- $\hat{k}_{1755} = 39.1$ • Believe CNV action at 1755? [$\bar{k} = 8.13$]
- **Permutation test** Randomly shift row j of \mathbf{X} by s_j units left (with wraparound):

$$\mathbf{x}_j^* = (x_{s+1,j}, x_{s+2,j}, \dots, x_{5000,j}, x_{1j}, x_{2j}, \dots, x_{sj})$$

- Do this for all 150 rows
- *Recalculate* \hat{k}_i^* values
- Compare $\hat{k}_{1755} = 39.1$ with $\{\hat{k}_i^*, i = 1, 2, \dots, 5000\}$

Actual khat distribution compared to permutation distribution;
Maximum khat* = 23.3



Locally Most Powerful Tests

- Let $r_i = \pi_{i1}/\pi_1 = \Pr\{\text{non-null}|C_i\} / \Pr\{\text{non-null}\}$.

- $l_i = \sum_1^n \{1 + (r_i - 1)T(z_{ij})\}$ where $T(z) = \frac{\text{tdr}(z) - \pi_1}{\pi_0}$

- \hat{k}_i nearly MLE in this model

- Test

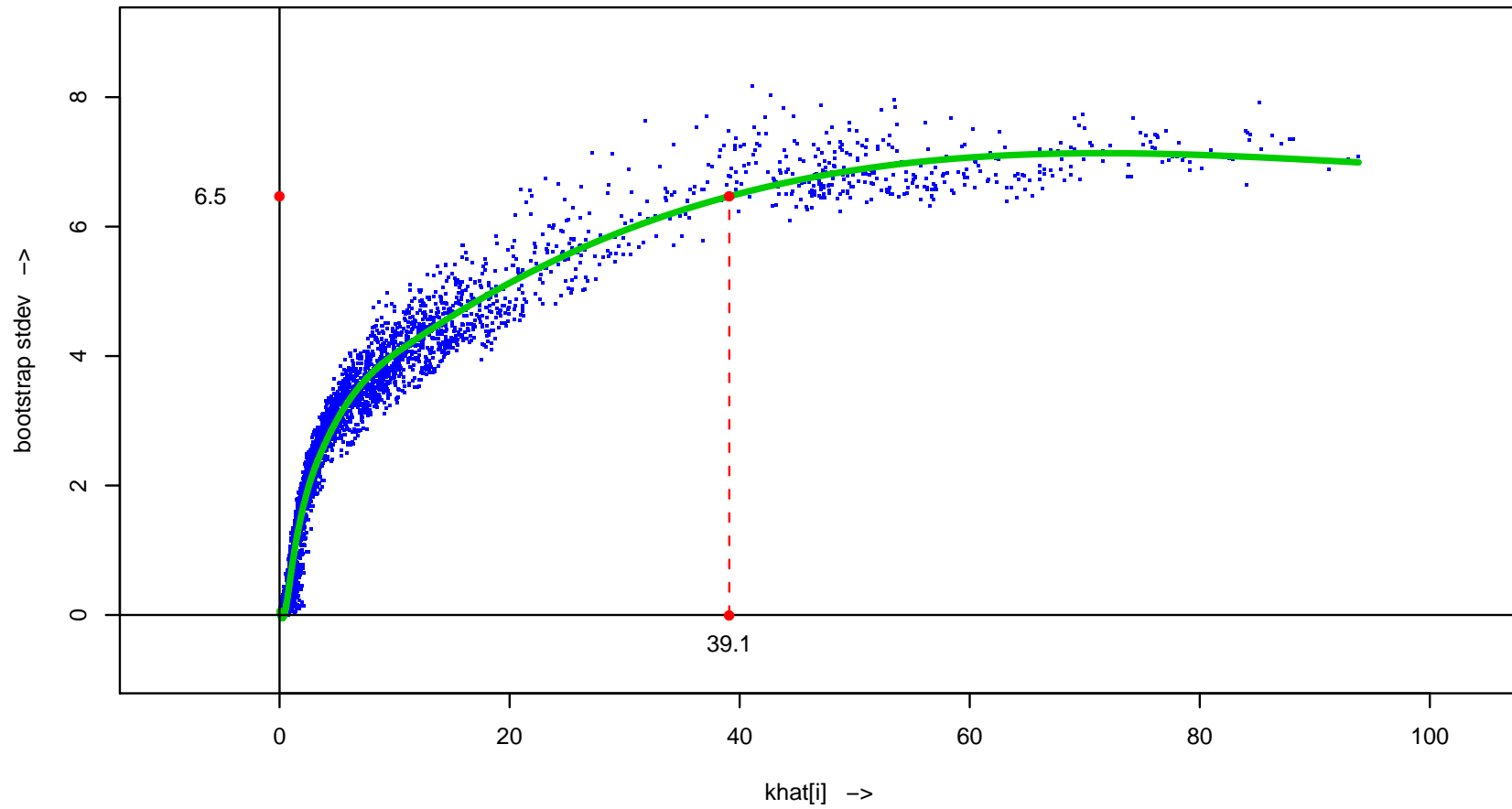
$$H_{0i} : r_i = 1 \quad \text{vs} \quad r_i > 1.$$

- Locally most powerful test rejects for large values of $\hat{k}_i^{(1)}$.
- Use permutations to get p -values.

Bootstrapping \hat{k}_i Estimates

- Resample rows (i.e., subjects)
- Recompute iterative estimate \hat{k}_i^* (5 iterations model)
- $\widehat{\text{sd}}_i$ = boot stdev of \hat{k}_i^* , $B = 100$ resamples
(did not recompute original $\widehat{\text{fdr}}$ curve each time)
- $\hat{k}_i^* \sim \mathcal{N}(\hat{k}_i, \widehat{\text{sd}}_i^2)$ $[6 \leq \widehat{\text{sd}}_i \leq 7 \text{ for } \hat{k}_i > 20]$

Bootstrap estimates of standard deviations for \hat{k}_i values,
(5 iterations) plotted vs \hat{k}_i ; $\text{sdhat}[1755]=6.5$



Brown–Stein–Robbins Estimation

- Suppose $\mu \sim g(\cdot)$ and $x|\mu \sim \mathcal{N}(\mu, \sigma^2)$
- $l(x)$ log marginal density of x

$$\mu|x \sim (x + \sigma^2 l'(x)), \sigma^2 (1 + \sigma^2 l''(x))$$

- Apply with $\mu = k_i$, $x = \hat{k}_i$,
 $\hat{l}(x) = \log$ smoothed density $\{\hat{k}_i\}$
- For $\hat{k}_i = 39.1$, $\hat{\sigma} = 6.5$, gave $k_{1755} \sim (41.3, 7.4^2)$
- **Conclusion** Even taking account of selection effects, k_{1755} is probably much larger than $\bar{k} = 8.13$.

More General Model for $\text{fdr}_i(z)$

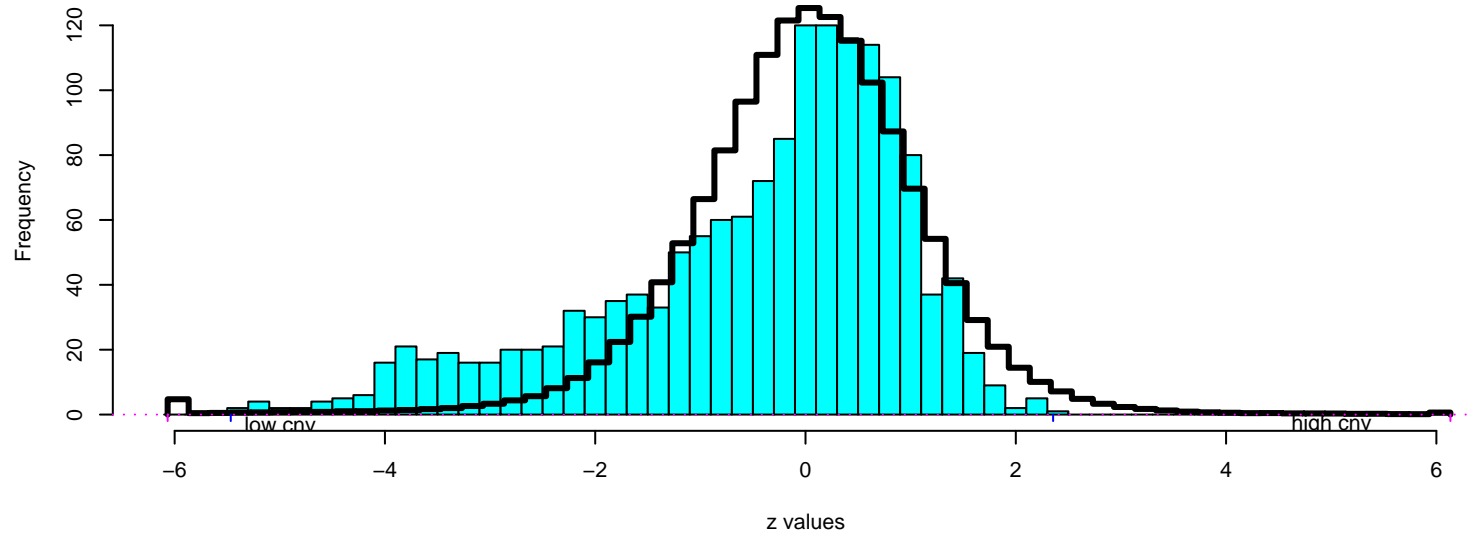
- “Method 2”: Multiclass Bayes model $\pi_{i0}, f_{i0}(z), f_{i1}(z)$ with all $f_{i0} = f_0$, but *drop assumption* that non-null distributions $f_{i1}(z)$ the same.

- Define: $w_i(z) = \Pr\{C_i|z\}$

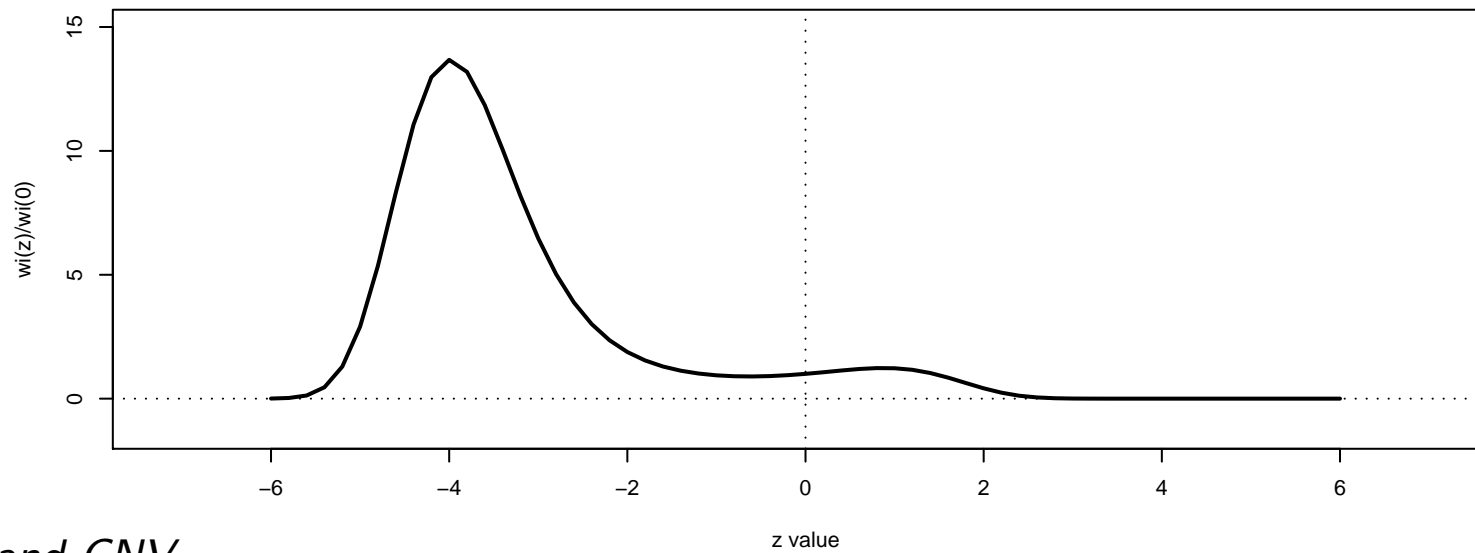
$$\text{fdr}_i(z) \doteq \text{fdr}(z) \cdot \frac{w_i(0)}{w_i(z)}$$

- Empirical Bayes Estimate $w_i(z)$ by logistic regression of C_i indicator on z_m .

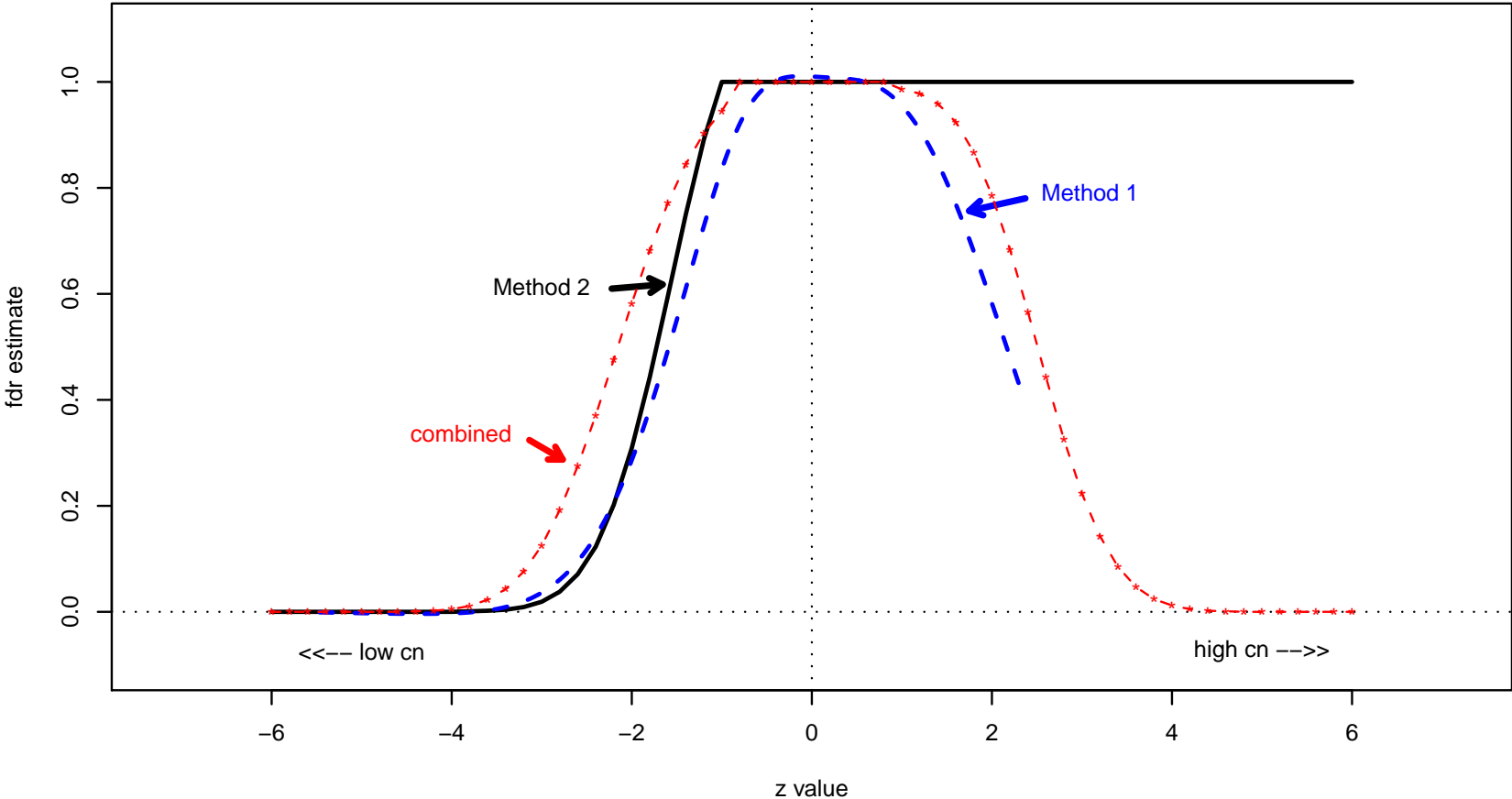
z-values for positions 1750:1759 (solid) compared to all the other positions (line)



logistic regression estimate of $w_i(z) = \text{Prob}\{1750:1759 \mid z\}$



Three estimates of fdrhat for positions 1750:1759



References

EFRON, B. (2008). Simultaneous inference: When should hypothesis testing problems be combined? *Ann. Appl. Statist.*, **2** 197–223.

TIBSHIRANI, R. and WANG, P. (2008). Spatial smoothing and hot spot detection for CGH data using the fused lasso. *Biostatistics*, **9** 18–29.

WALTHER, G. (2009). Optimal and fast detection of spatial clusters with scan statistics. Online, URL <http://stat.stanford.edu/~gwalther/>.

WANG, P., KIM, Y., POLLACK, J., NARASIMHAN, B.

and TIBSHIRANI, R. (2005). A method for calling gains and losses in array CGH data. *Biostatistics*, **6** 45–58.

ZHANG, N., SIEGMUND, D., JI, H. and LI, J. (2009). Detecting simultaneous change-points in multiple sequences. *Biometrika*. Accepted for publication, URL <http://stat.stanford.edu/~nzhang/>.