

Hypothesis Testing

Wolfgang Huber, EMBL



Aims for this lecture

Understand the basic principles of decision making by hypothesis testing, pitfalls, strengths, use cases and limitations

What changes when we go from single to multiple testing?

- false discovery rates
- p-values
- multiple testing ‘adjustments’
- hypothesis filtering and weighting

See also

www.huber.embl.de/msmb Chapter 6

The screenshot shows the Cambridge University Press Academic website. At the top, there is a navigation bar with links to 'Academic', 'Cambridge English', 'Education', 'Bibles', 'Digital Products', 'About Us', 'Careers', and a language selector for 'Italy'. A search bar is located in the top right corner. Below the header, the word 'Academic' is prominently displayed, followed by the tagline 'Unlocking potential with the best learning and research solutions'. A search bar with a magnifying glass icon and a link to 'Include historic titles' are also present. The main content area features a book cover for 'MODERN STATISTICS FOR MODERN BIOLOGY' by Susan Holmes and Wolfgang Huber. The book cover has a green and blue abstract design. To the left of the book image, there is a red button labeled 'LOOK INSIDE'. Below the book image, there is a link to 'I want this title to be available as an eBook'. To the right of the book image, the title 'Modern Statistics for Modern Biology' is listed, along with its classification as a 'TEXTBOOK'. Below the title, the authors 'Susan Holmes' and 'Wolfgang Huber' are mentioned, along with their respective institutions. The publication date is February 2019. The availability is 'In stock', the format is 'Paperback', and the ISBN is 9781108705295. There are two buttons: 'Add to cart' and 'Add to wishlist'. Below these buttons, there is a link to 'Request inspection copy' with a circular icon containing a download symbol. A note states: 'Lecturers may request a copy of this title for inspection' and a 'Request' button. At the bottom of the page, there are links for 'Description', 'Contents', 'Resources', 'Courses', and 'About the Authors'.

How to make rational decisions based on noisy, finite data?

Examples:

- Testing efficacy of a drug on people
 - lack of complete experimental control
 - finite sample size
- Effect of a fertilizer, a genetic variant, ... on phenotype of plants / animals in an outdoors field trial
- Lady testing tea, clairvoyant, telepath, ...
- Toxicology

+: No understanding of mechanism involved / needed / desired

-: Wouldn't we *want* to use any available understanding or 'priors'?

The fundamental tradeoff of statistical decision making

← bias

accuracy-

dispersion ↑



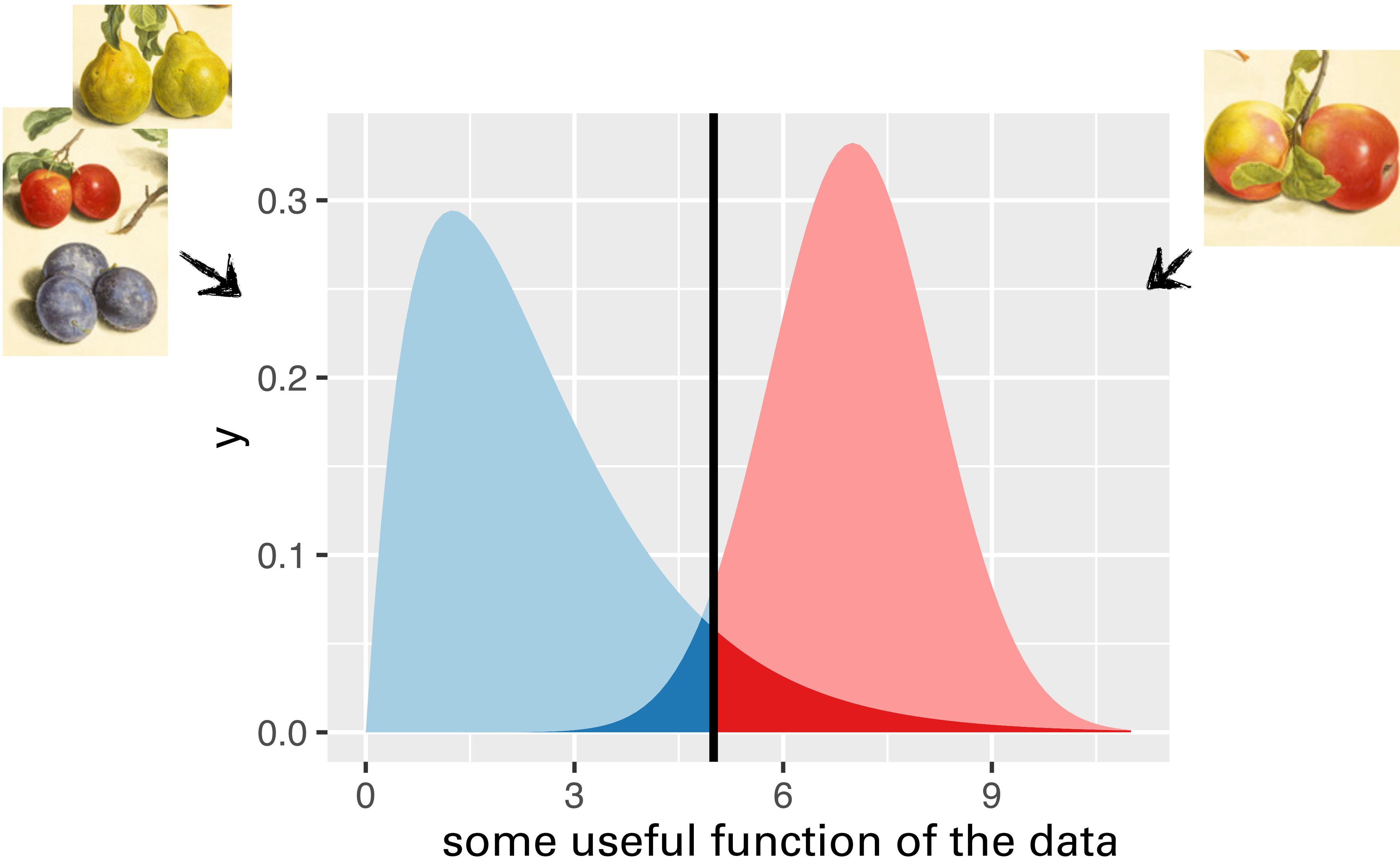
Comes in various
guises

Accuracy vs Precision

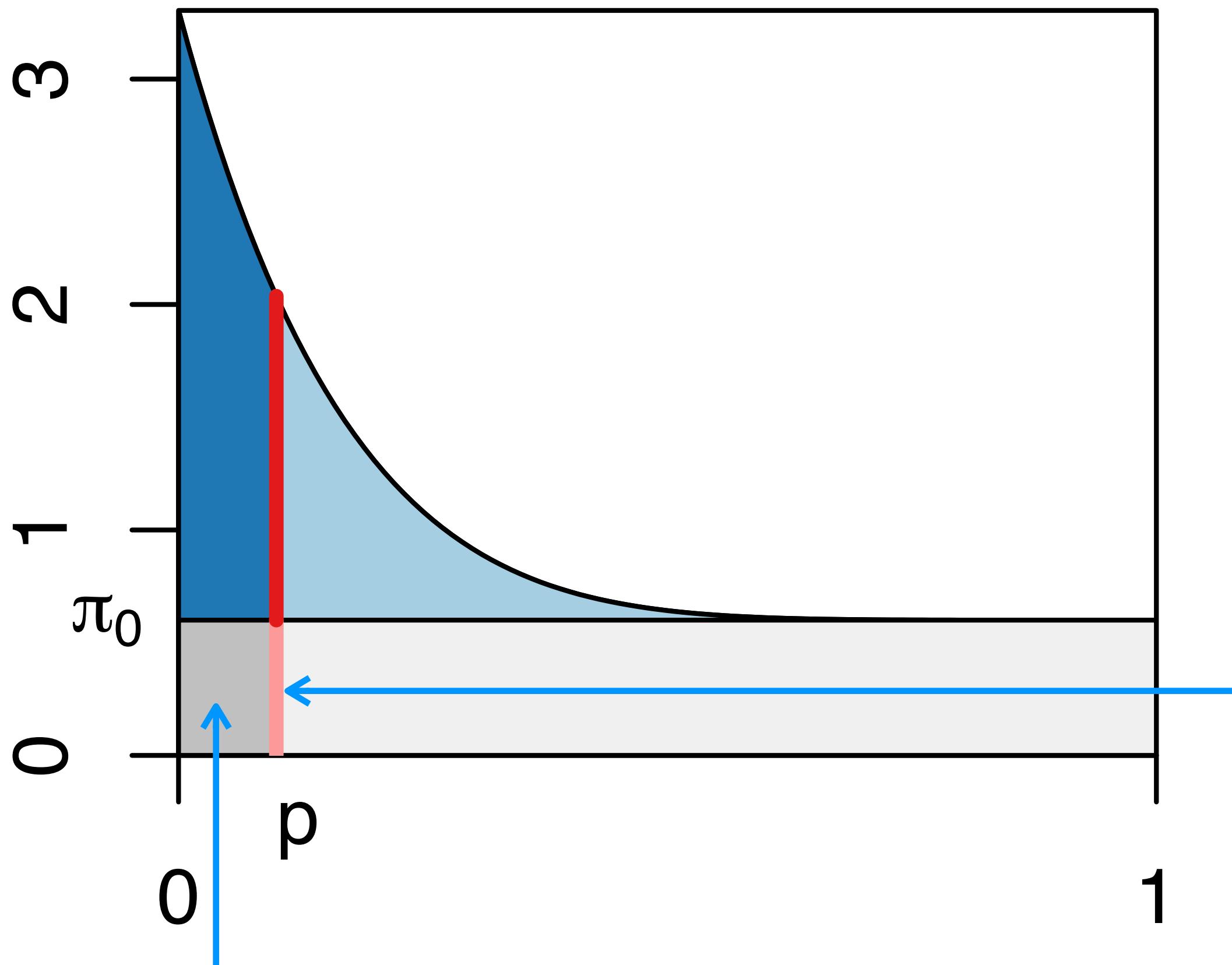
Bias vs Variance

Model complexity vs overfitting

Basic problem: two classes



The two-groups model and the (local) false discovery rate



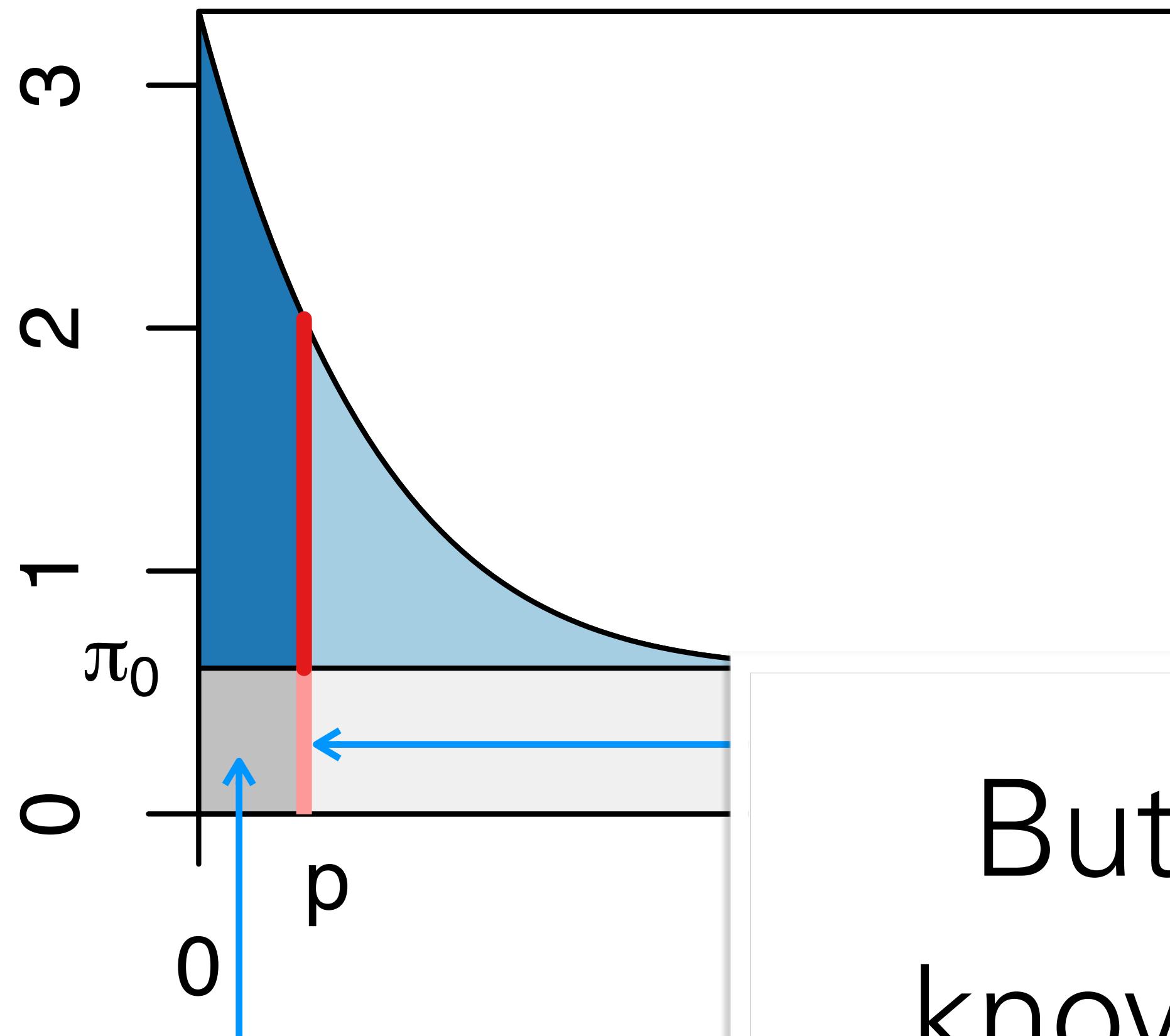
FDR: Ratio between the areas. A set property. It applies to a set of hypotheses (discoveries).

$$f(p) = \pi_0 + (1 - \pi_0)f_{\text{alt}}(p),$$

$$\text{fdr}(p) = \frac{\pi_0}{f(p)}.$$

fdr: Ratio between the line segment lengths. An individual property. It applies to individual discoveries.

The two-groups model and the (local) false discovery rate



FDR: Ratio between the number of false discoveries and the total number of discoveries in a set. It applies to a set of hypotheses (discoveries).

$$f(p) = \pi_0 + (1 - \pi_0)f_{\text{alt}}(p),$$

$$\text{fdr}(p) = \frac{\pi_0}{f(p)}$$

But how do we know π_0 and f_{alt} ?

the line
n
t

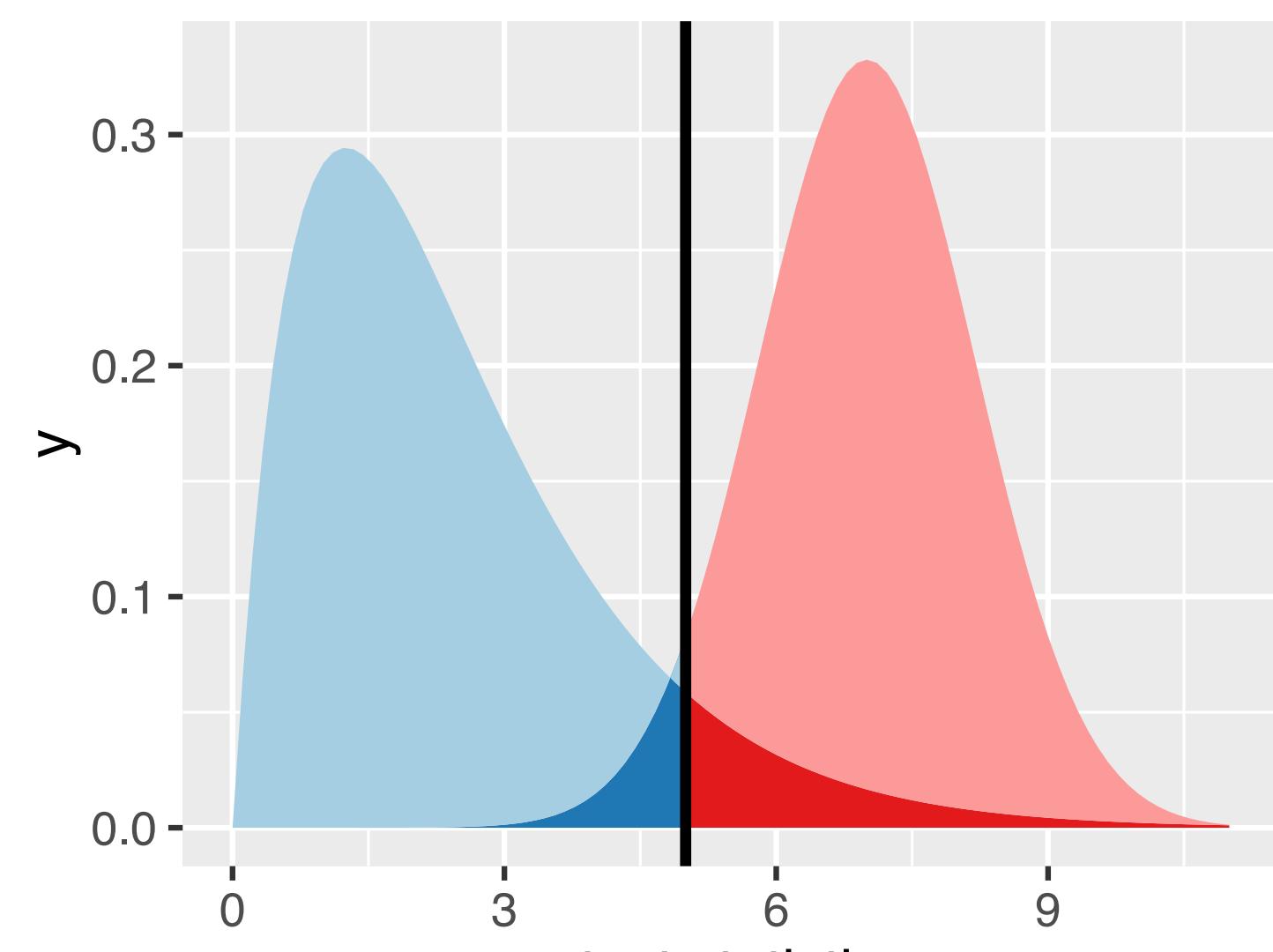
Machine Learning

Data-driven (induction)

Need lots of training data

Using multiple variables

... or objects that are not even
'variables' (e.g. images)



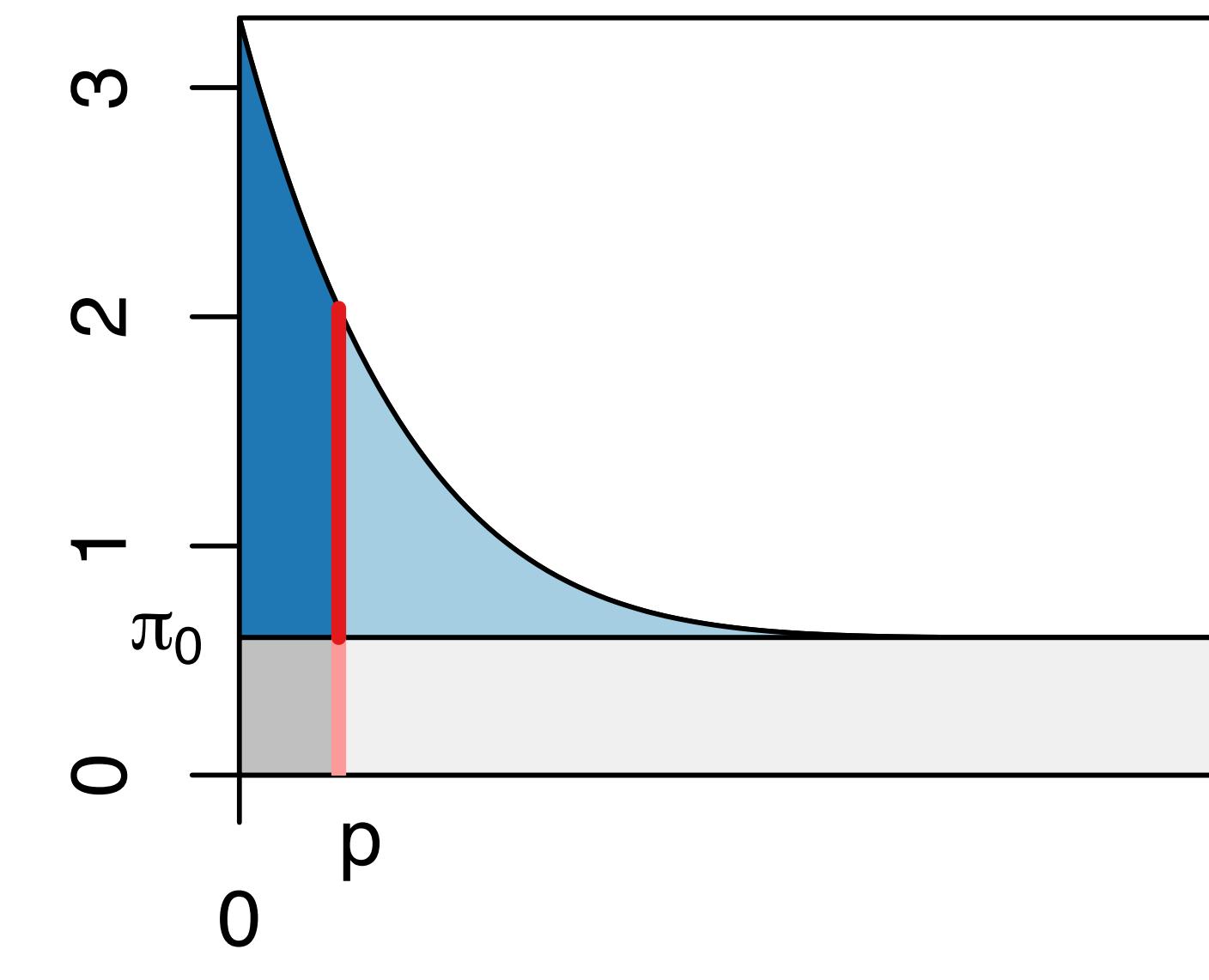
Hypothesis testing

Theory/Model-driven (deduction)

No training data

Choice of model and process of
deduction

Regulatory use



Example

Toss a coin a number of times ⇒

If the coin is fair, then heads should appear half of the time (roughly).



But what is “roughly”? We use combinatorics / probability theory to quantify this.

Suppose we flipped the coin 100 times and got 59 heads. Is this ‘significant’?

Binomial distribution

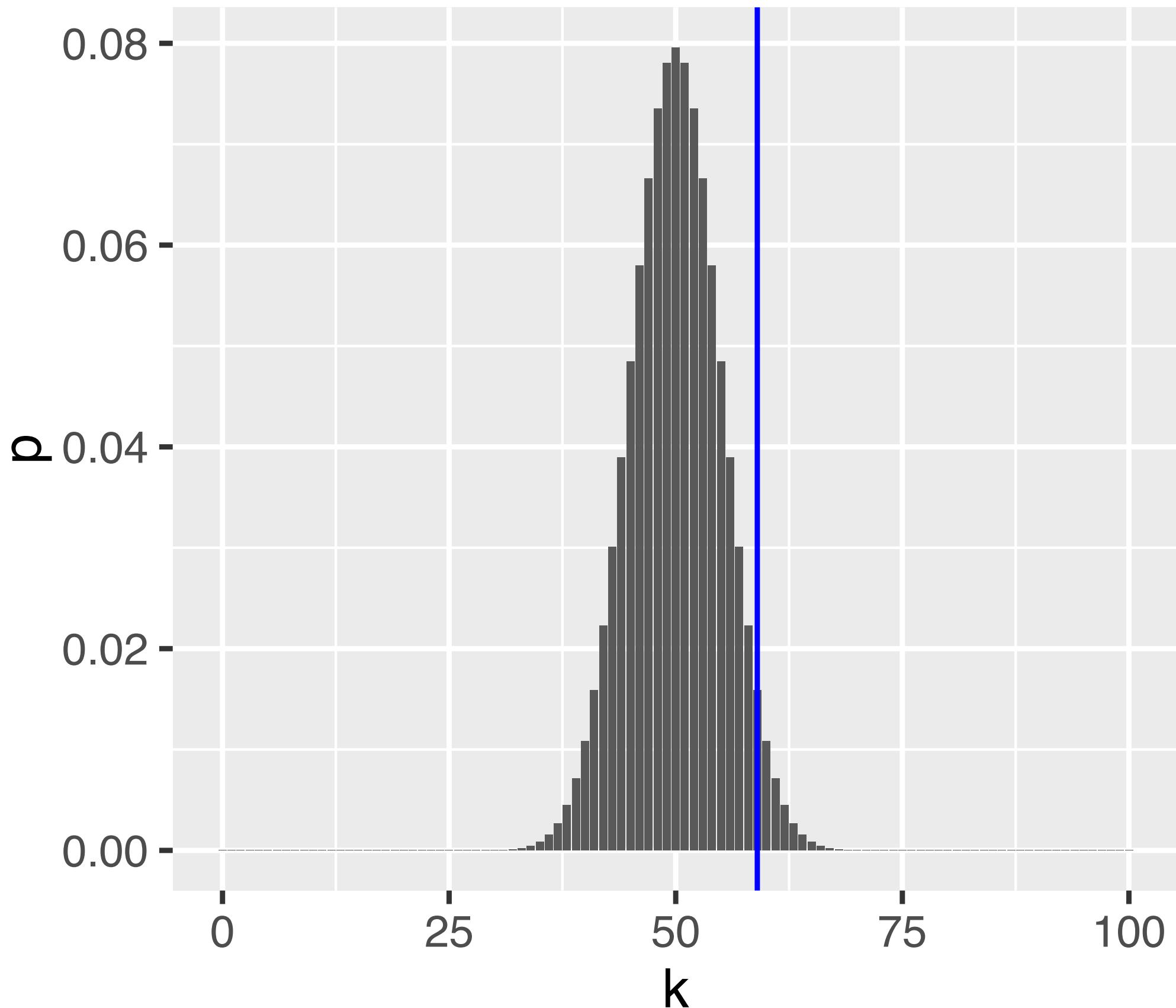


Figure 6.3: The binomial distribution for the parameters $n = 100$ and $p = 0.5$,

$$P(K = k \mid n, p) = \binom{n}{k} p^k (1 - p)^{n-k}$$

Rejection region

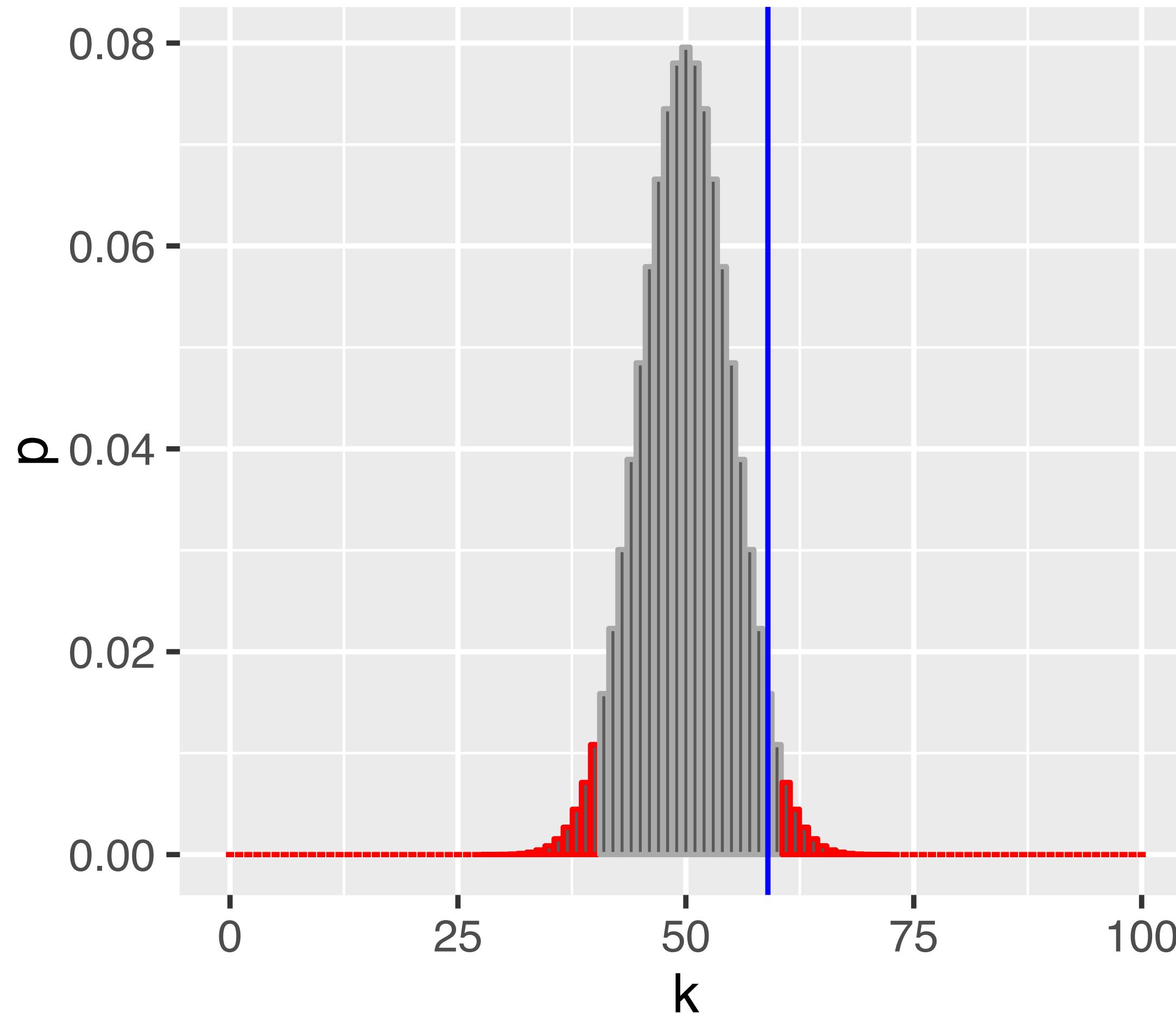


Figure 6.5: As Figure 6.3, with rejection region (red) whose total area is $\alpha = 0.05$.

Questions

- Does the fact that we don't reject the null hypothesis mean that the coin is fair?
- Would we have a better chance of detecting an unfair coin if we did more coin tosses? How many?
- If we repeated the whole procedure and again tossed the coin 100 times, might we then reject the null hypothesis?
- Our rejection region is asymmetric - its left part ends with 40, while its right part starts with 61. Why is that? Which other ways of defining the rejection region might be useful?

The Five Steps of Hypothesis Testing

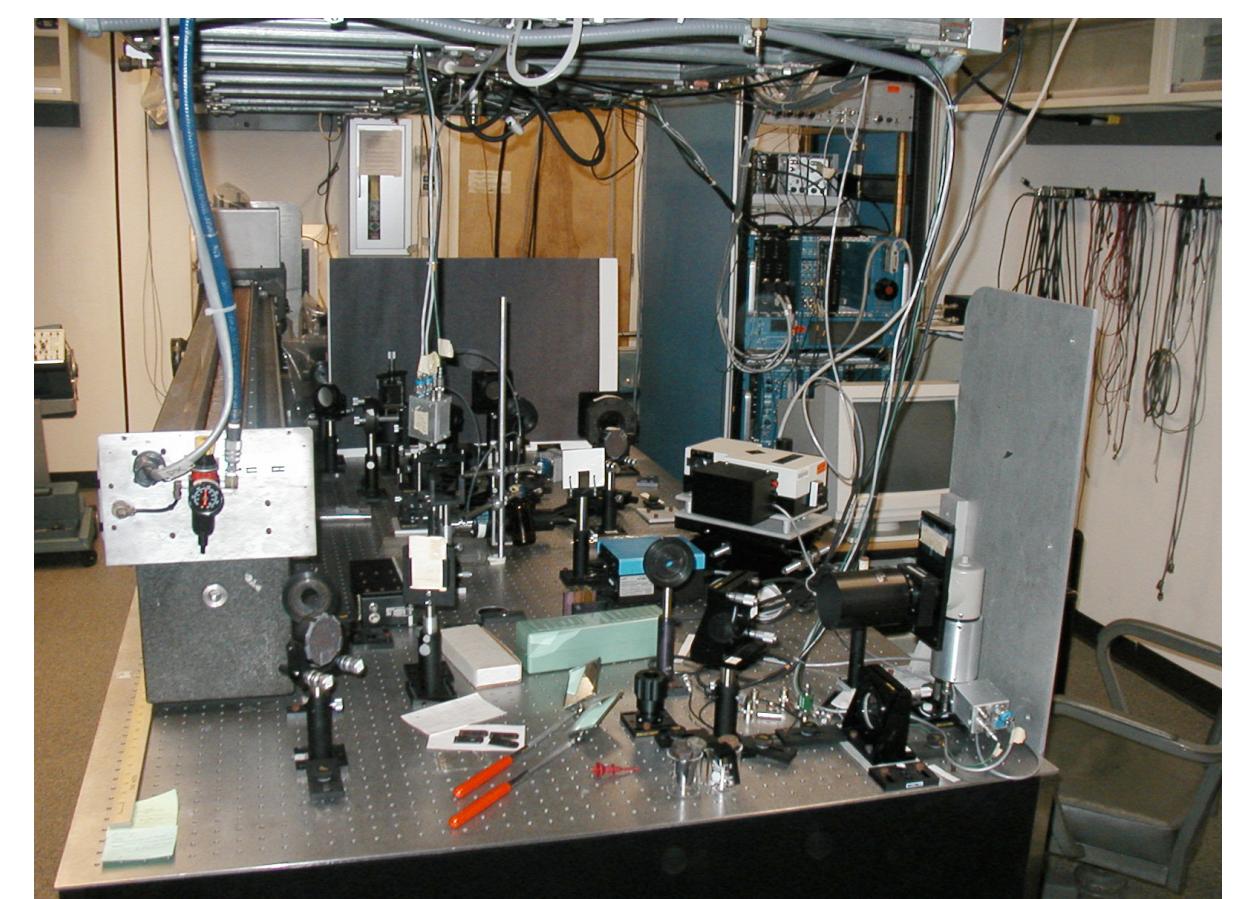
Choose an experimental design and a data summary function for the effect that you are interested in: the **test statistic**

Set up a **null hypothesis**: a simple, computationally tractable model of reality that lets you compute the null distribution of the test statistic, i.e. all its possible outcomes and each of their probabilities.

Decide on the **rejection region**, i.e., a subset of possible outcomes whose total probability is small (**significance level**).

Do the experiment, collect data, compute the test statistic.

Make a **decision**: reject null hypothesis if the test statistic is in the rejection region.



The Five Steps of Hypothesis Testing

Choose an experimental design and a data summary function for the effect that you are interested in:

Set up a null hypothesis:

that lets you compute the possible outcomes and easily

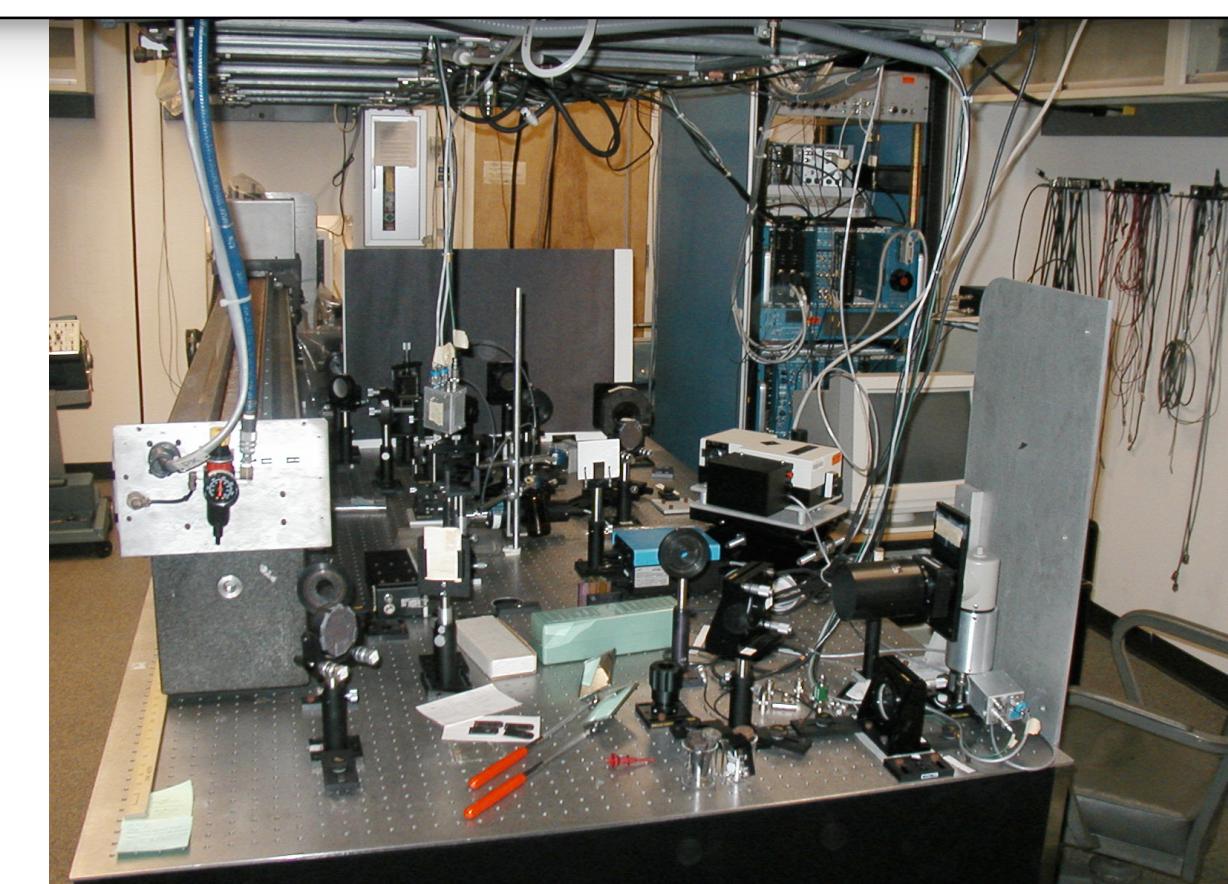
Decide on the rejection rule: if the total probability is small (say,

Do the experiment, collect data, compute the test statistic.

Make a decision: reject null hypothesis if the test statistic is in the rejection region.

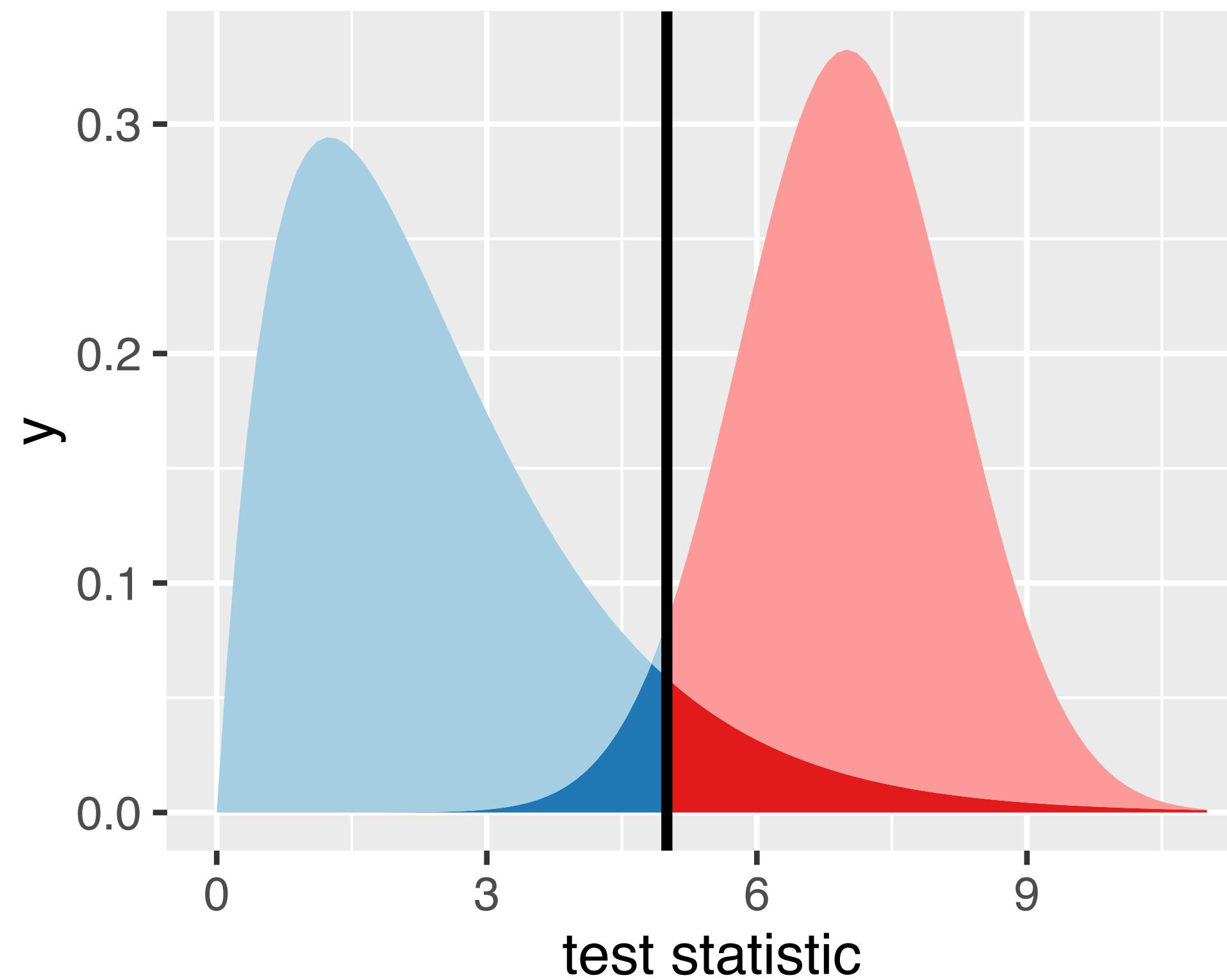
This is the idealised scenario, “orthodoxy”:

Reality, esp. in retrospective ‘data-mining’ can be quite different.



Types of Error in Testing

Test vs reality	Null hypothesis is true	... is false
Reject null hypothesis	Type I error (false positive)	True positive
Do not reject	True negative	Type II error (false negative)



Parametric Theory vs Simulation

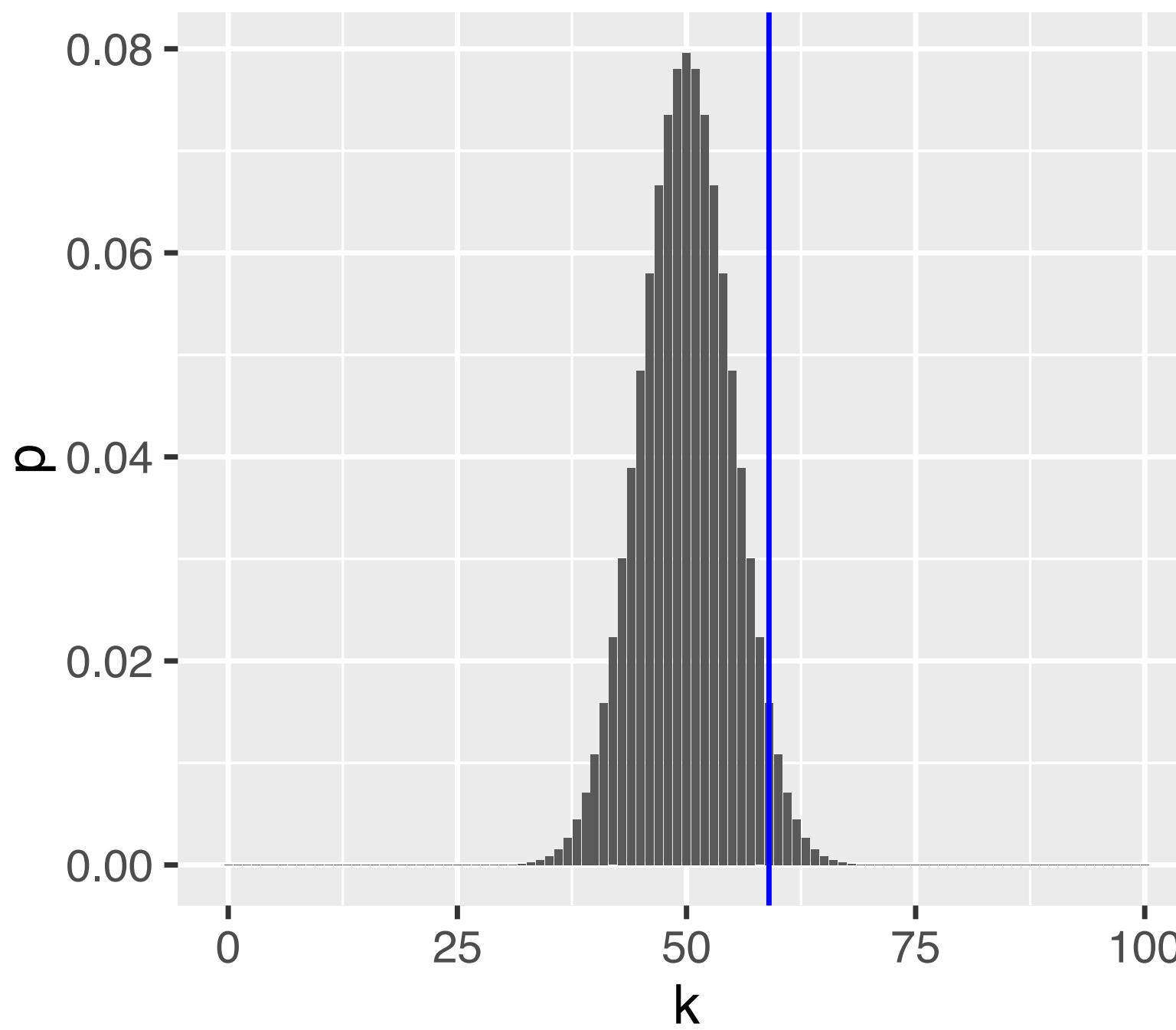


Figure 6.3: The binomial distribution for the parameters $n = 100$ and $p = 0.5$, according to Equation (6.1).

$$P(K = k \mid n, p) = \binom{n}{k} p^k (1 - p)^{n-k}$$

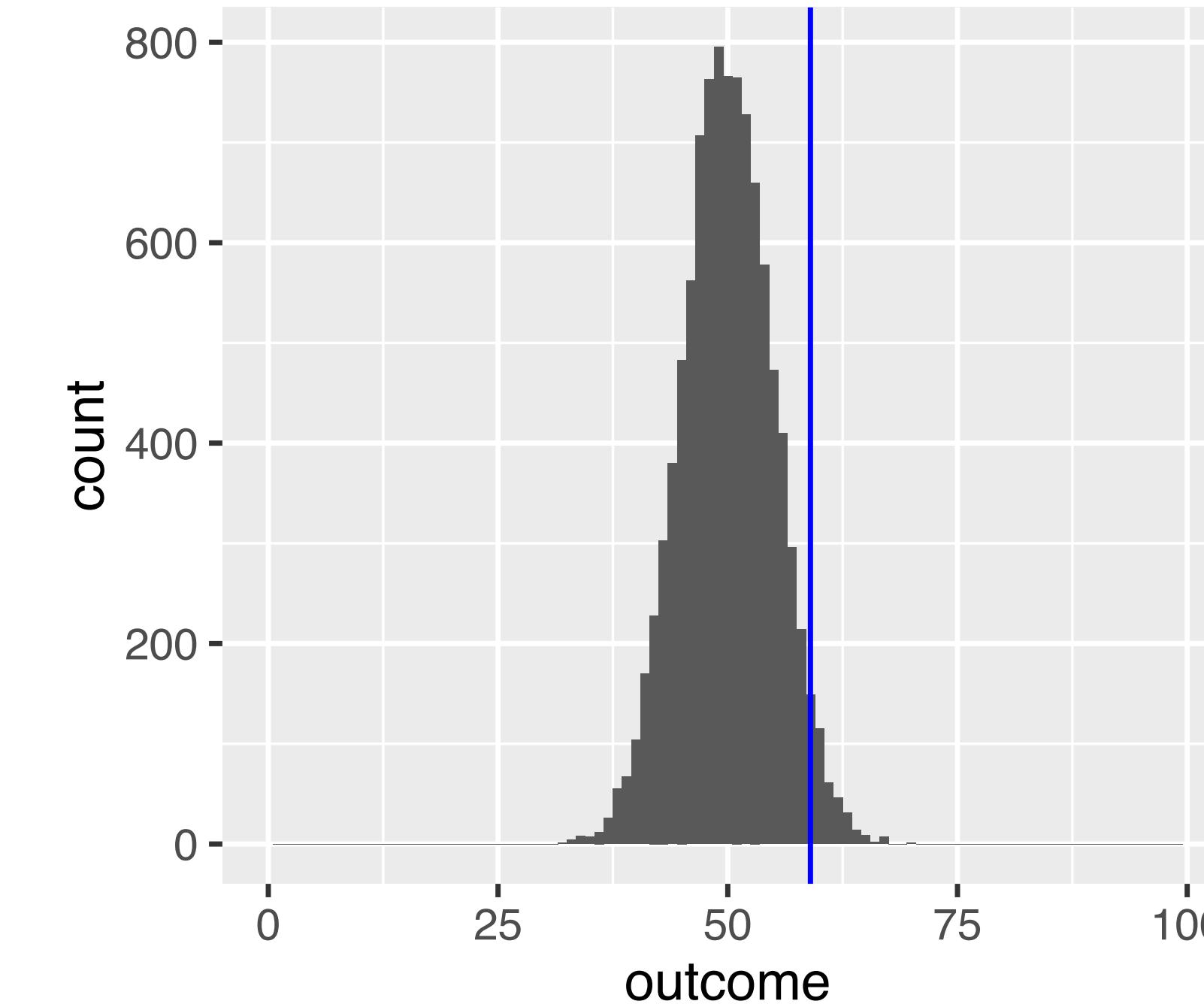
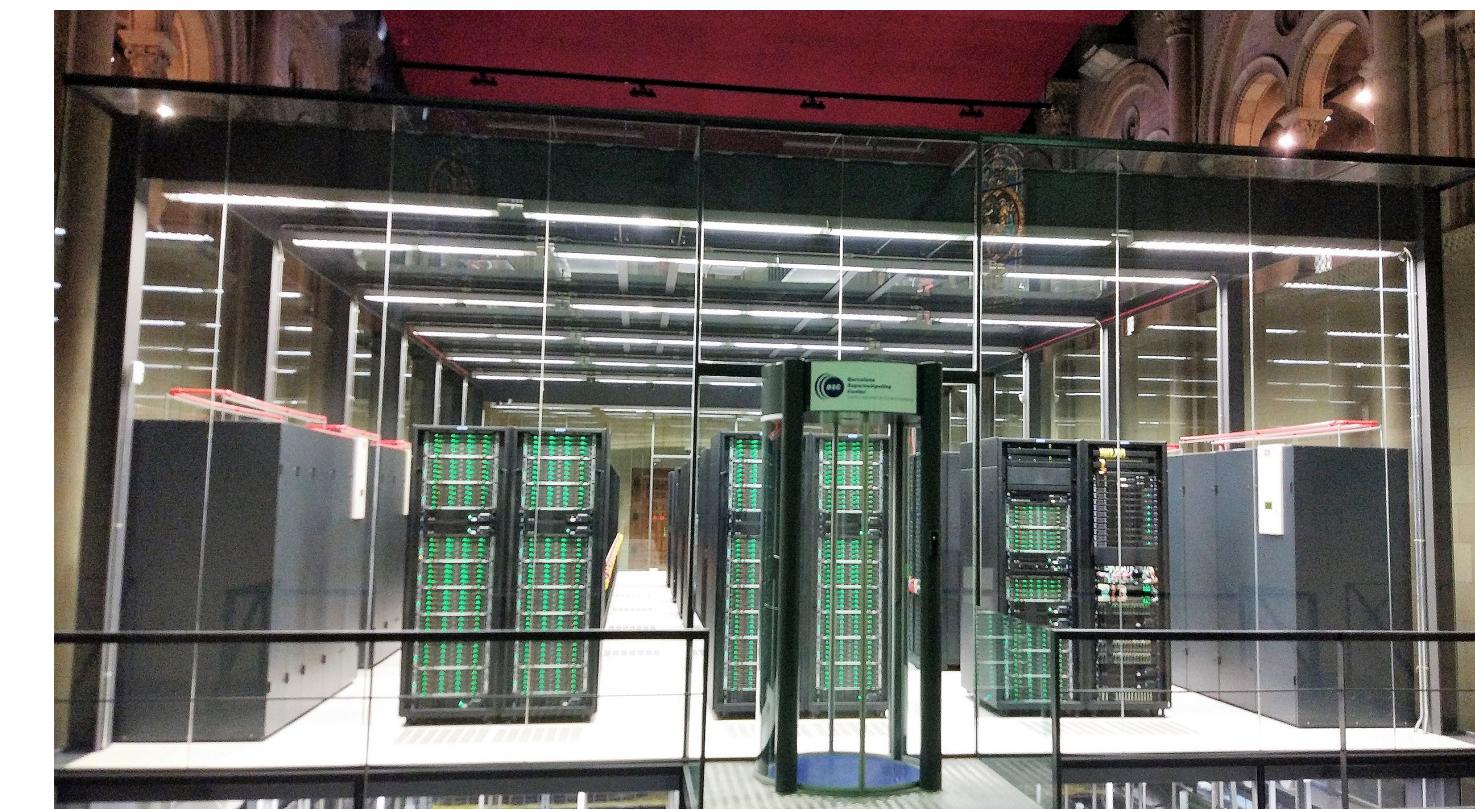


Figure 6.4: An approximation of the binomial distribution from 10^4 simulations (same parameters as Figure 6.3).



Parametric Theory vs Simulation

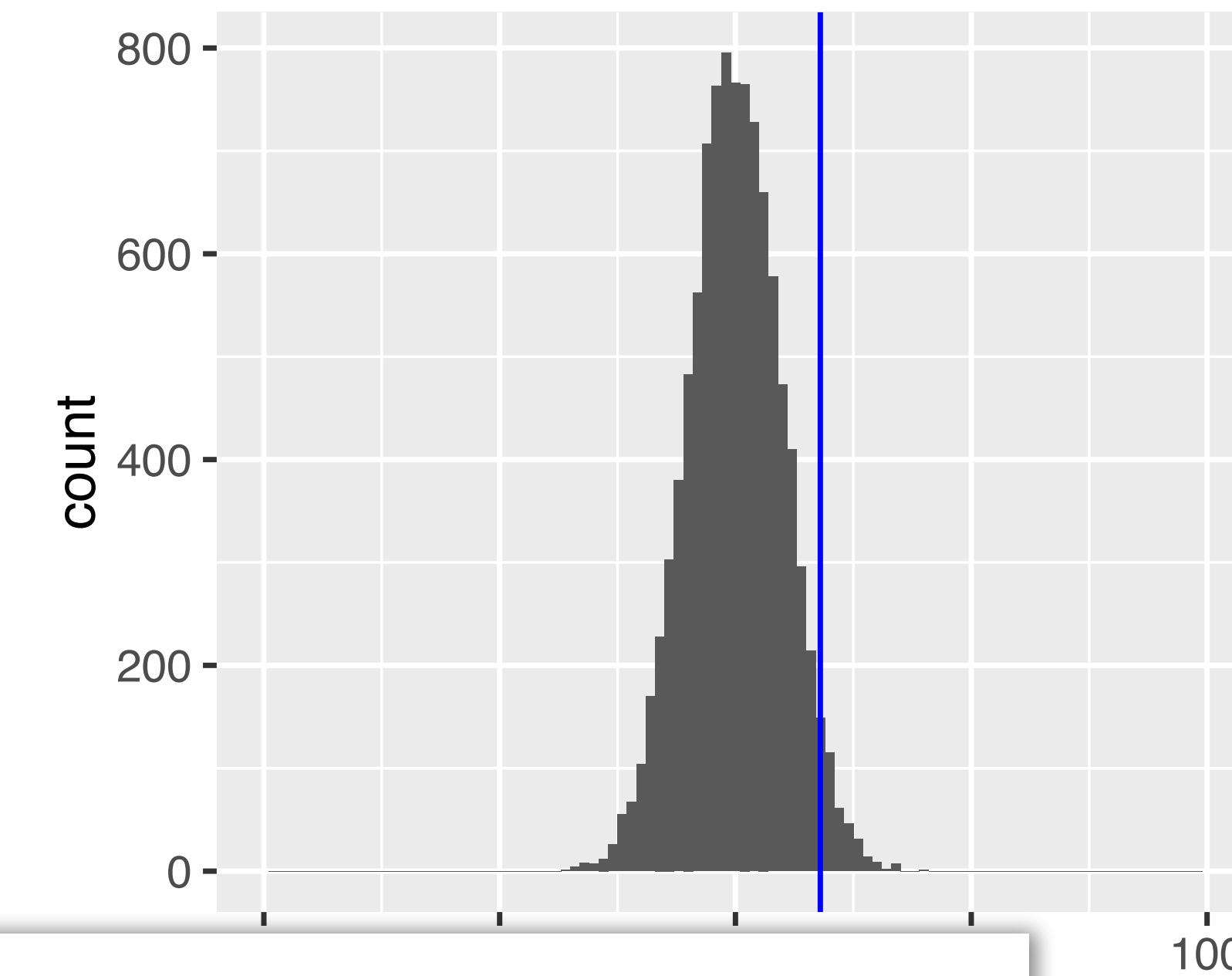
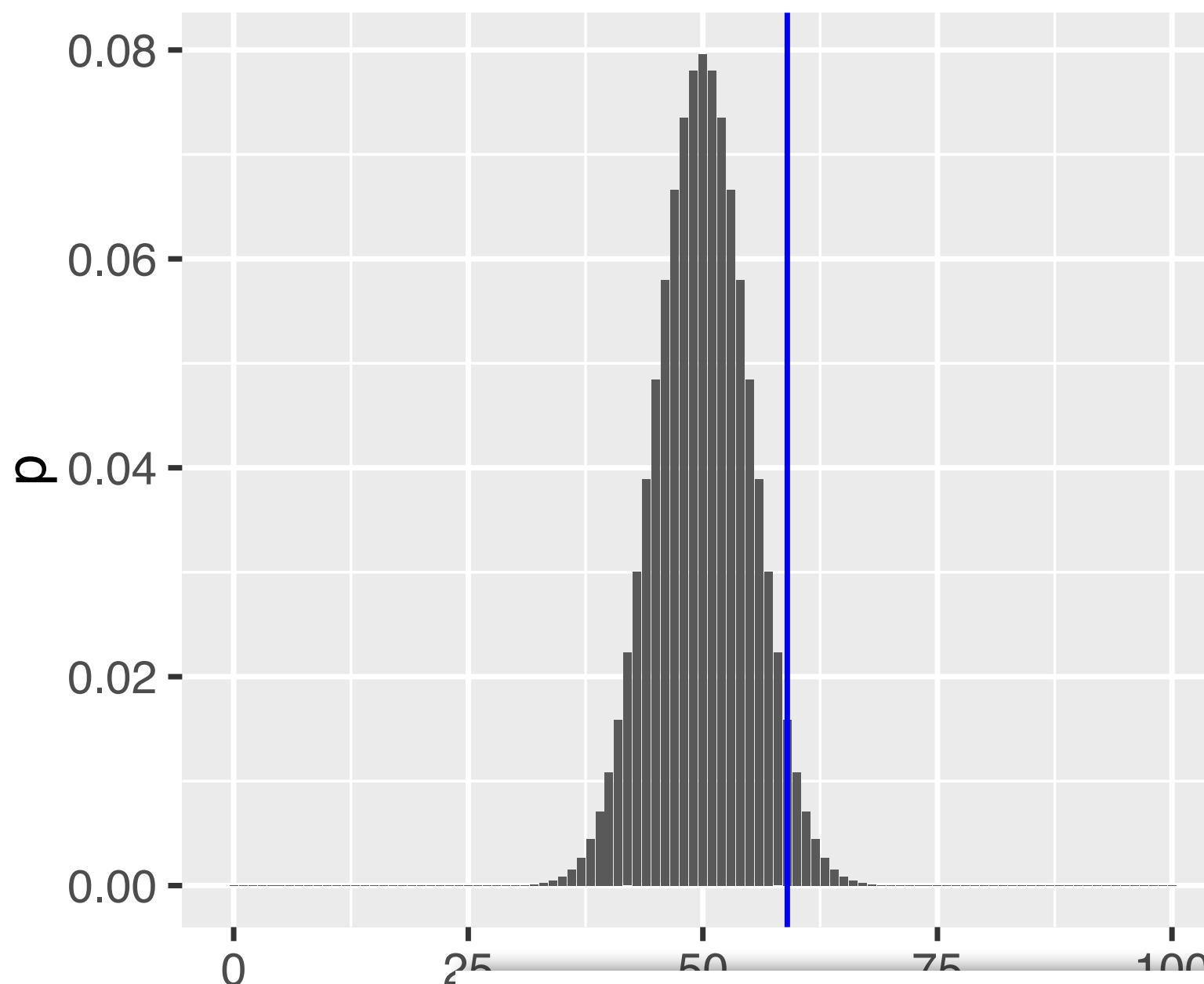
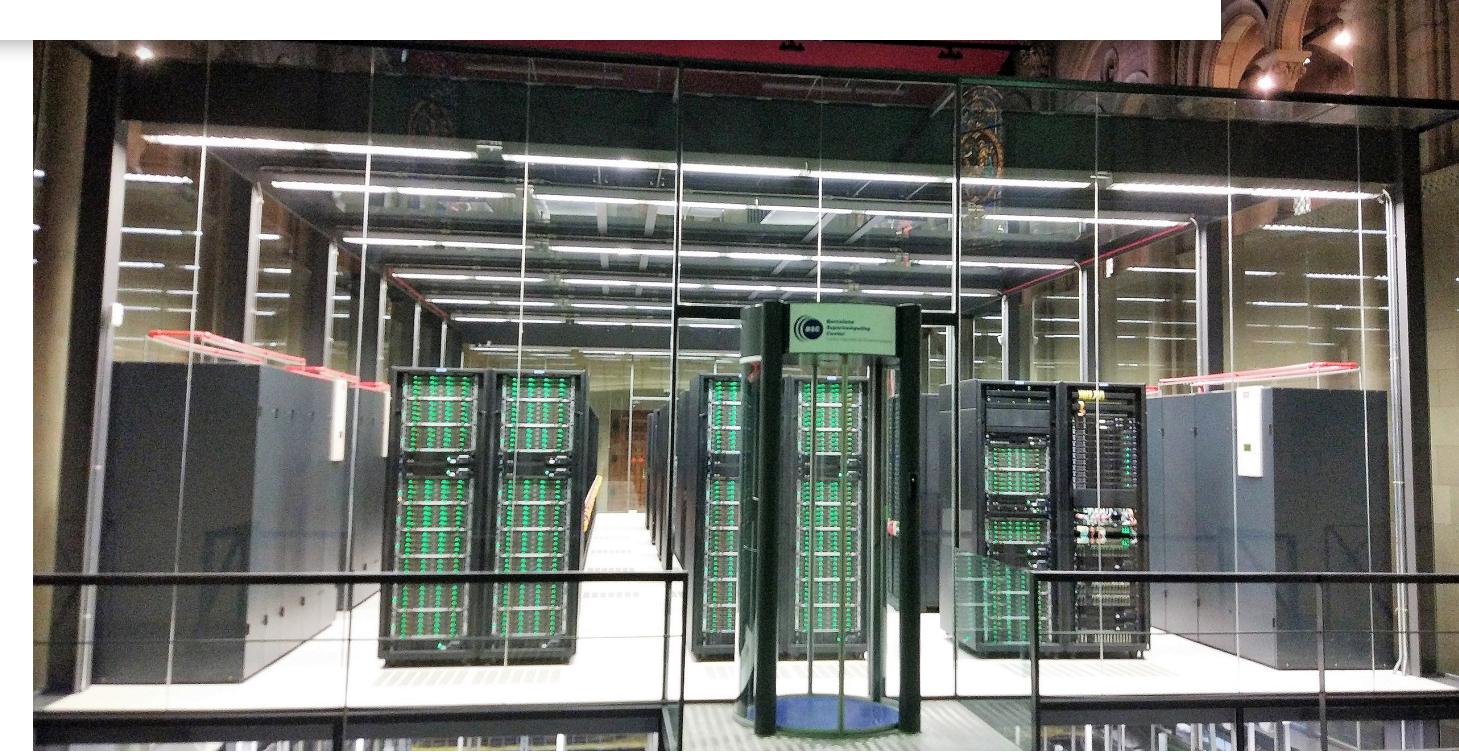


Figure 6.3: The k
the parameters κ
according to Equ

Q:

Discuss pros and contras for each

$$P(K = k \mid n, p) = \binom{n}{k} p^k (1 - p)^{n-k}$$



The choice of the test statistic

Suppose we observed 50 tails in a row, and then 50 heads in a row. Is this a perfectly fair coin?

We could use a different test statistic: number of times we see two tails in a row

Is this statistic generally and always preferable?

Power

There can be several test statistics, with different power, for different types of alternative

Continuous data: the t-statistic

$$t = c \frac{m_1 - m_2}{s}$$

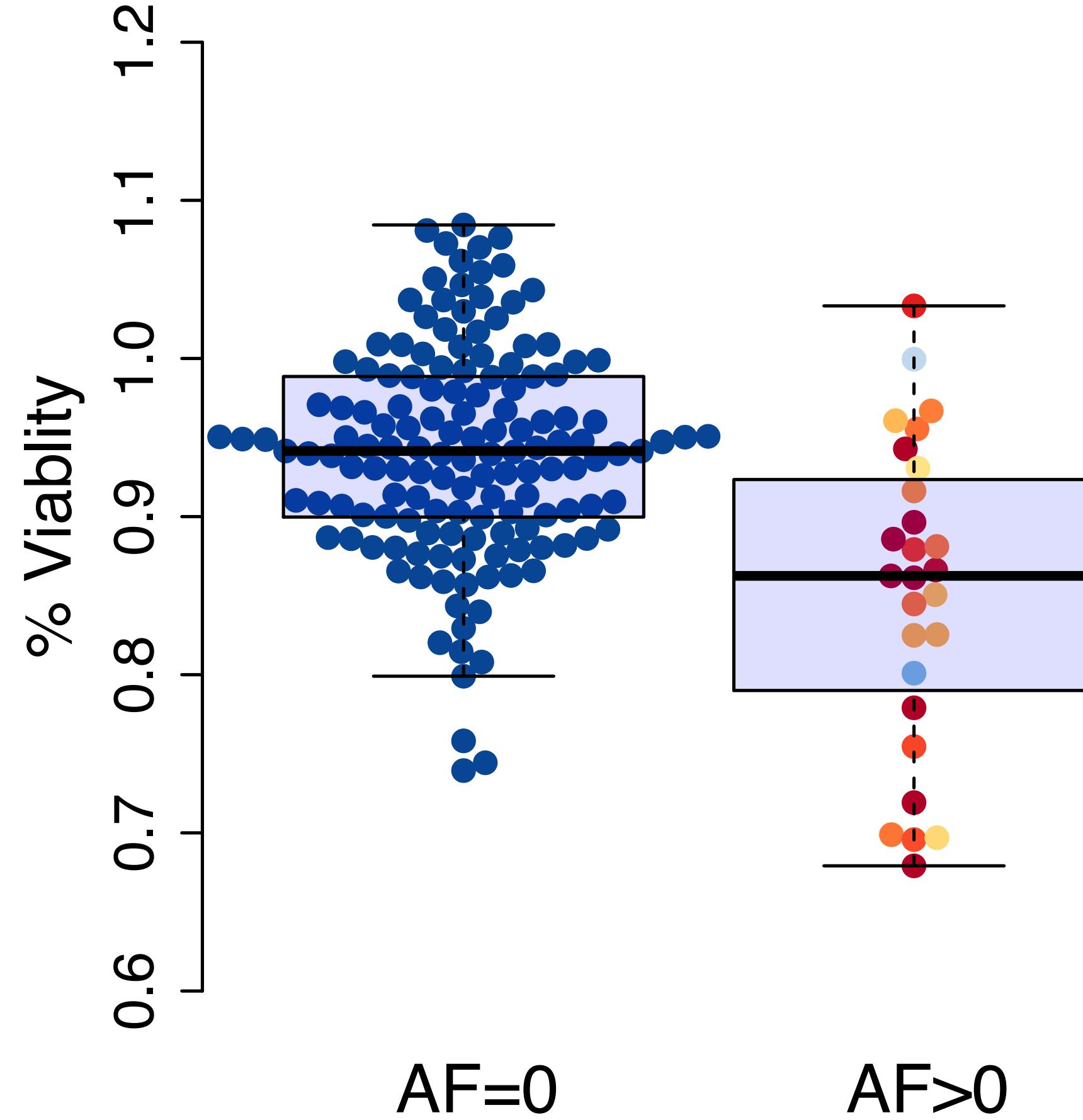
- Can also be adapted to one group only
- Relation to z-score

$$m_g = \frac{1}{n_g} \sum_{i=1}^{n_g} x_{g,i} \quad g = 1, 2$$

$$s^2 = \frac{1}{n_1 + n_2 - 2} \left(\sum_{i=1}^{n_1} (x_{1,i} - m_1)^2 + \sum_{j=1}^{n_2} (x_{2,j} - m_2)^2 \right)$$

$$c = \sqrt{\frac{n_1 n_2}{n_1 + n_2}}.$$

**selumetinib 0.156 μ M ~ trisomy12
($p = 3.02e-08$)**



t-distribution

If the data are identically normal distributed and independent, then under H_0 , t follows a ' t -distribution' with parameter $n_1 + n_2$ (a.k.a. degrees of freedom)

t-distribution

If the data are identically normal distributed and independent, then under H_0 , t follows a 't-distribution' with parameter $n_1 + n_2$ (a.k.a. degrees of freedom)

Q:

How does the distribution of $|t|$ look?

t-distribution

If the data are identically normal distributed and independent, then under H_0 , t follows a ' t -distribution' with parameter $n_1 + n_2$ (a.k.a. degrees of freedom)

Q:

How does the distribution of $|t|$ look?

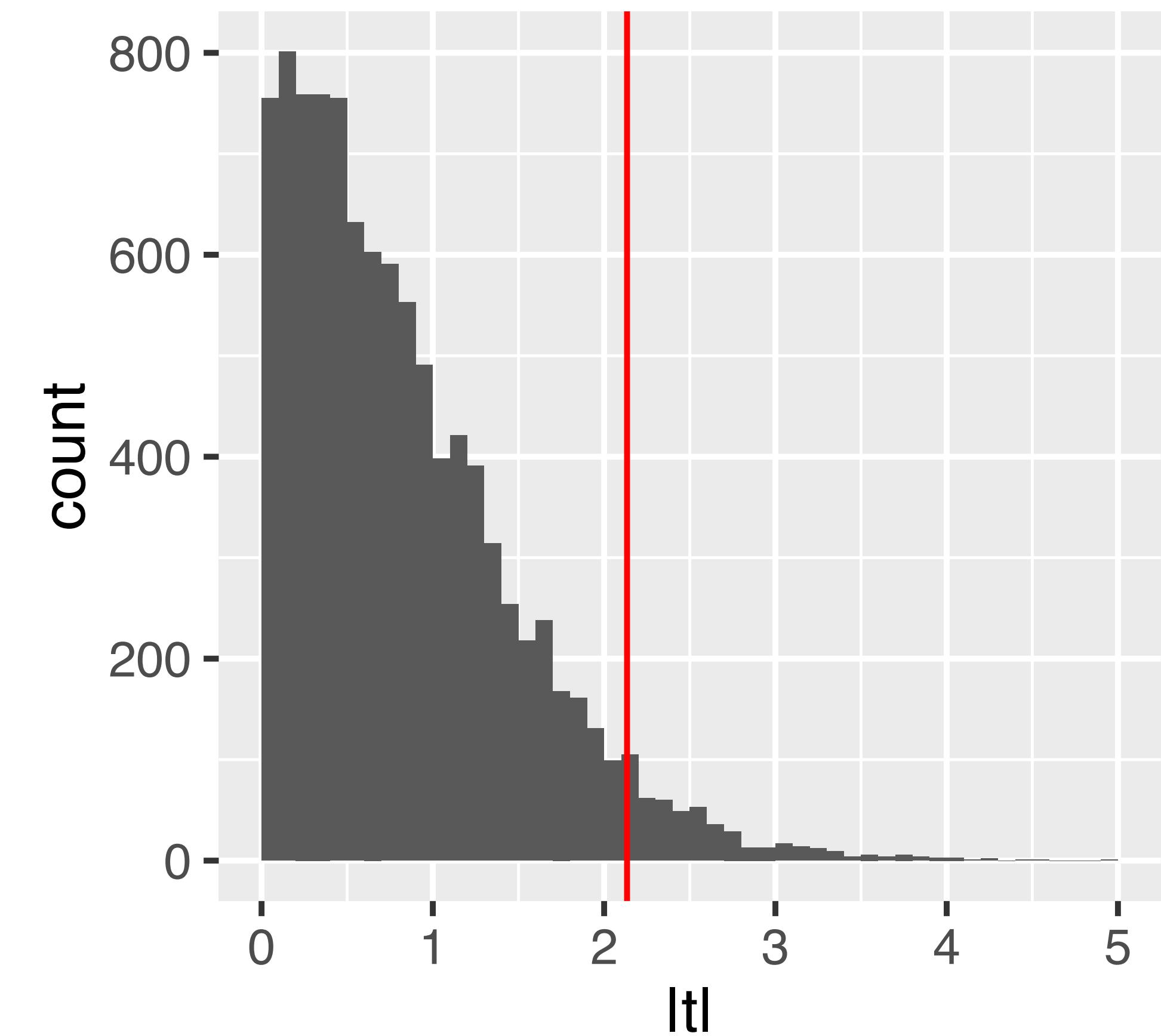


Figure 6.8: The null distribution of the (absolute) t -statistic determined by simulations – namely, by random permutations of the group labels.

Comments and Pitfalls

The proof that the t -statistic follows a t -distribution assumes that observations are independent and follow a normal distribution: this is a sufficient, but not necessary, condition

Comments and Pitfalls

The proof that the t -statistic follows a t -distribution assumes that observations are independent and follow a normal distribution: this is a sufficient, but not necessary, condition

[Deviation from normality](#) (heavier tails): test typically maintains type-I error control, but no longer has optimal power.

Comments and Pitfalls

The proof that the t -statistic follows a t -distribution assumes that observations are independent and follow a normal distribution: this is a sufficient, but not necessary, condition

[Deviation from normality](#) (heavier tails): test typically maintains type-I error control, but no longer has optimal power.

[Options](#): use permutations;
use a different test (e.g., Wilcoxon)

Comments and Pitfalls

The proof that the t -statistic follows a t -distribution assumes that observations are independent and follow a normal distribution: this is a sufficient, but not necessary, condition

Deviation from normality (heavier tails): test typically maintains type-I error control, but no longer has optimal power.

Options: use permutations;
use a different test (e.g., Wilcoxon)

Deviation from independence: type-I error control is lost, p-values will likely be totally wrong (e.g., for positive correlation, too optimistic).

Comments and Pitfalls

The proof that the t -statistic follows a t -distribution assumes that observations are independent and follow a normal distribution: this is a sufficient, but not necessary, condition

Deviation from normality (heavier tails): test typically maintains type-I error control, but no longer has optimal power.

Options: use permutations;
use a different test (e.g., Wilcoxon)

Deviation from independence: type-I error control is lost, p-values will likely be totally wrong (e.g., for positive correlation, too optimistic).

No easy options:

- ... try to model the dependence / remove it ...
- ... empirical null (Efron et al.) ...

Avoid Fallacy

The p-value is the probability that the data could happen, under the condition that the null hypothesis is true.

It is not the probability that the null hypothesis is true.

Absence of evidence ≠ evidence of absence



Recap: Single Hypothesis Testing

p-values are random variables: uniformly distributed if the null hypothesis is true - and should be close to zero if the alternative holds.

Note: We only observe one draw.

We prove something by disproving ('rejecting') the opposite (the null hypothesis). Reject = Discover.

Not rejecting does not prove the null hypothesis

Repeating the experiment (under the null): Around 5% of the times the p-value will be less than 0.05 by chance

All this reasoning is probabilistic. Testing & p-values are for rational decision making in uncertain contexts.

Limitations of p-value based hypothesis testing

Too much power: often, the 'null' is small (point-like), alternative is large (region-like)

Summarizing the data into one single number mushes together effect size and sample size

No place to take into account plausibility or 'prior' knowledge

What is p-value hacking ?

On the same data, try different tests until one is significant

On the same data, try different hypotheses until one is significant (HARKing -
hypothesizing after results are known)

Moreover...:

retrospective data picking

'outlier' removal

the 5% threshold and publication bias

What is p-value hacking ?

On the same data, try different tests until one is significant

On the same data, try different hypotheses until one is significant (HARKing -
hypothesizing after results are known)

Moreover....:

retrospective data picking

'outlier' removal

the 5% threshold and publication bias

The ASA's Statement on p-Values: Context,
Process, and Purpose

Ronald L. Wasserstein & Nicole A. Lazar

DOI: 10.1080/00031305.2016.1154108

What is p-value hacking ?

On the same data, try different tests until one is significant

On the same data, try different hypotheses until one is significant (HARKing -
hypothesizing after results are known)

Moreover....:

retrospective data picking

'outlier' removal

the 5% threshold and publication bias

The ASA's Statement on p-Values: Context,
Process, and Purpose

Ronald L. Wasserstein & Nicole A. Lazar

DOI: 10.1080/00031305.2016.1154108

What can we do about this?

The right answer to the wrong question

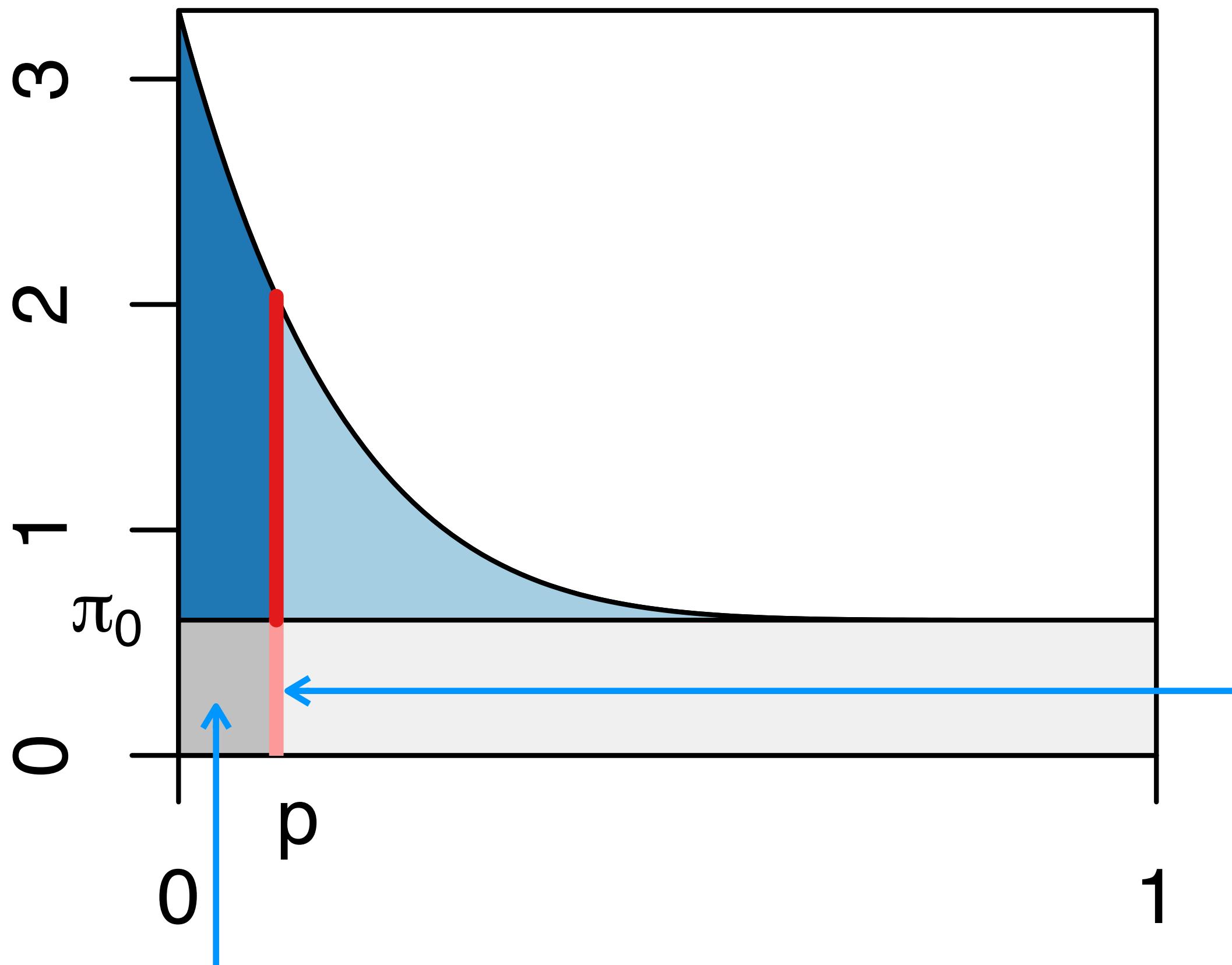
Researchers (regulators, investors, etc.) usually want to know:

If I publish this finding (allow this drug, invest in this product, ...), what is the probability that I'll later be proven wrong (cause harm, lose my money, ...)?

The p-value is the probability of seeing the data if the null hypothesis is true. It has little to do with the probability that my subsequent decision is wrong (a.k.a. "false discovery").

Can we compute a *false discovery probability* instead?

The two-groups model and the (local) false discovery rate



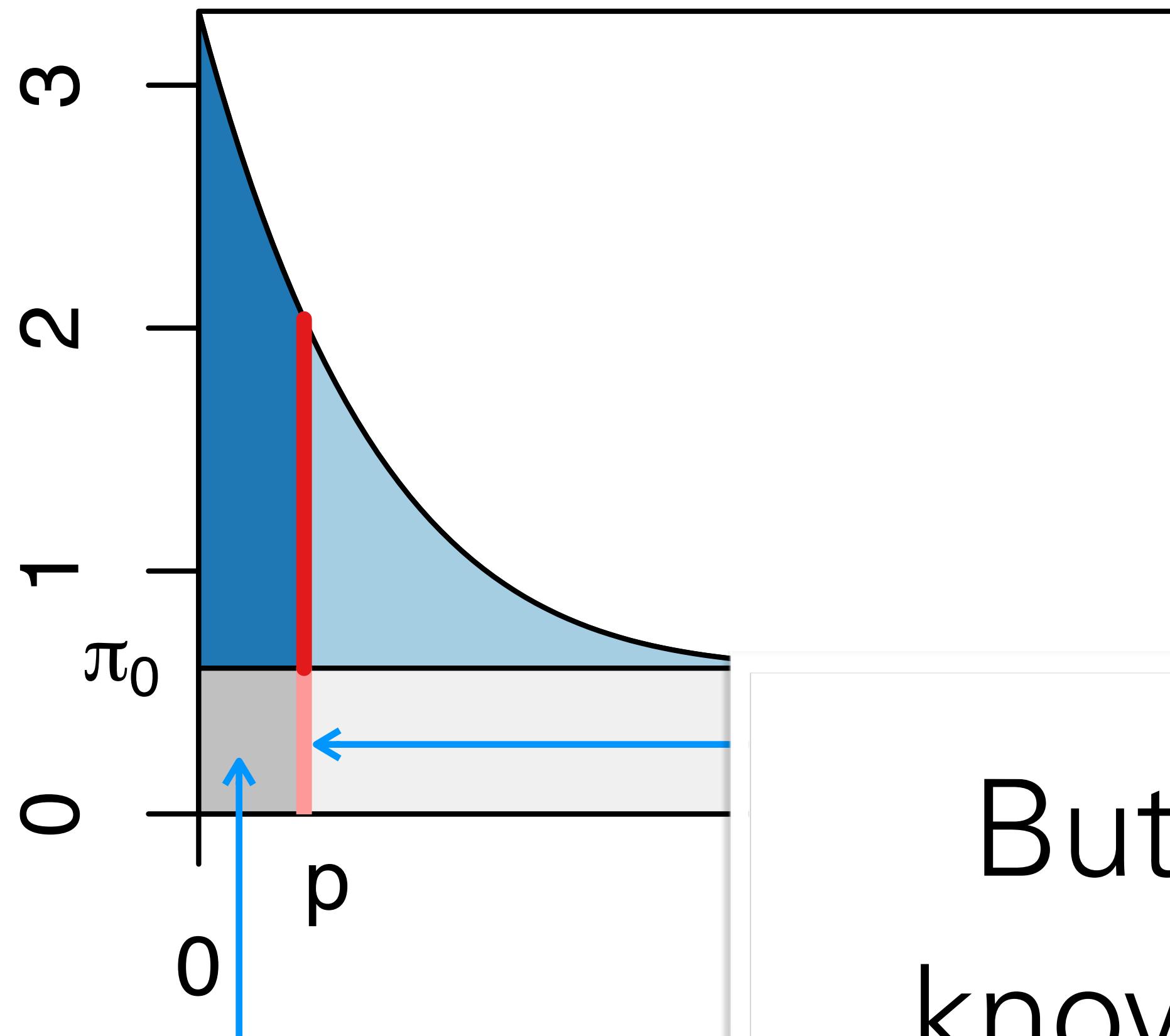
FDR: Ratio between the areas. A set property. It applies to a set of hypotheses (discoveries).

$$f(p) = \pi_0 + (1 - \pi_0)f_{\text{alt}}(p),$$

$$\text{fdr}(p) = \frac{\pi_0}{f(p)}.$$

fdr: Ratio between the line segment lengths. An individual property. It applies to individual discoveries.

The two-groups model and the (local) false discovery rate



FDR: Ratio between the number of false discoveries and the total number of discoveries in a set. It applies to a subset of hypotheses (discoveries).

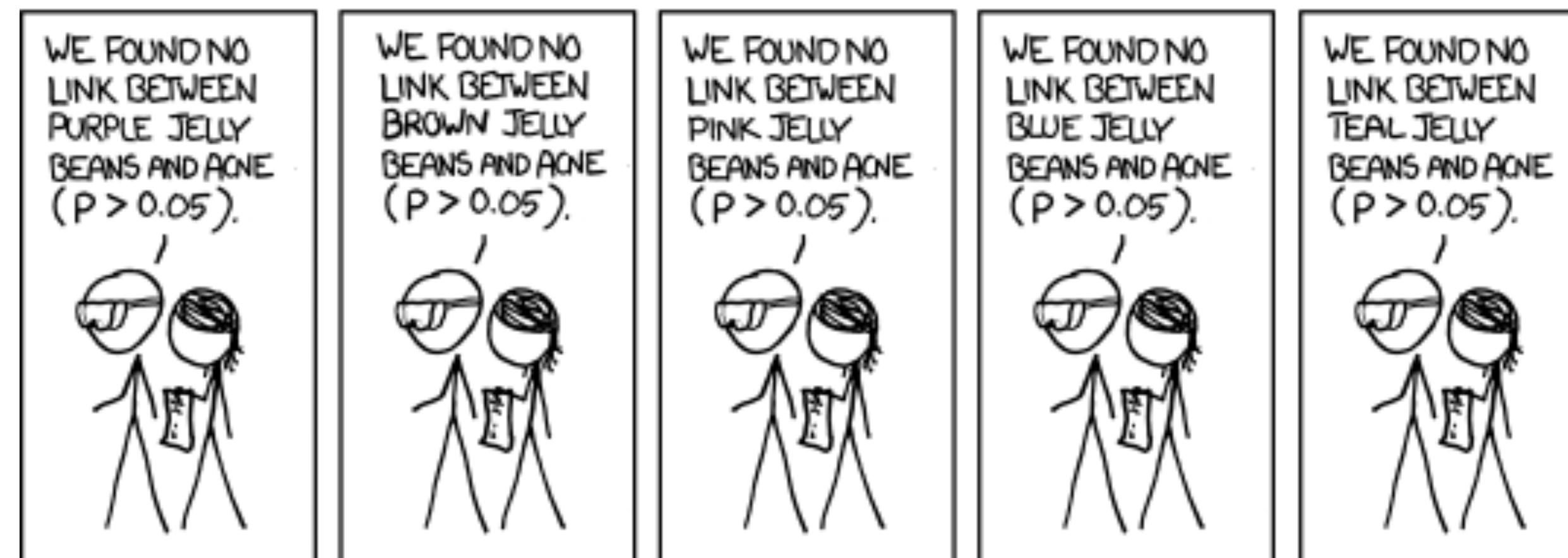
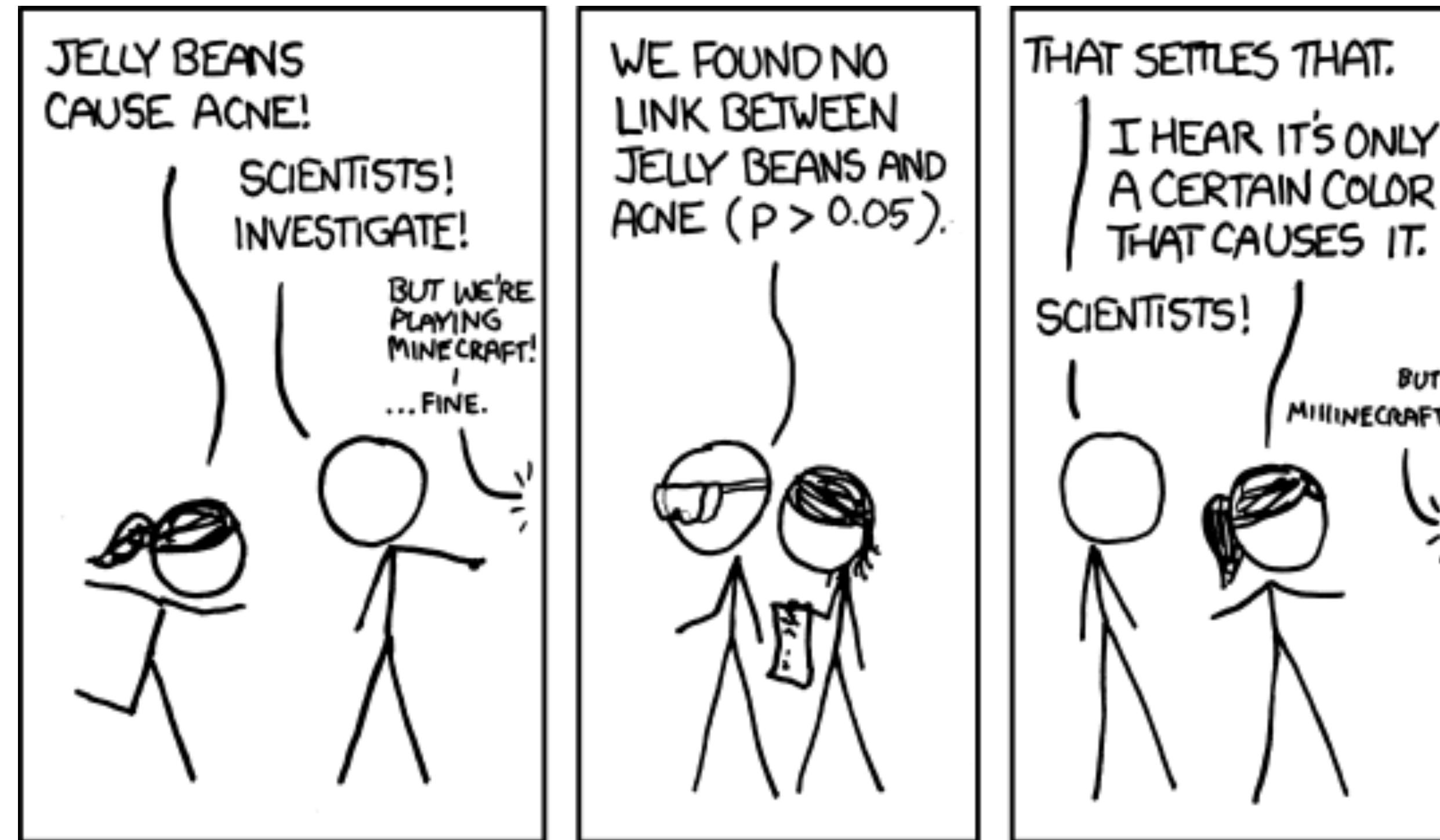
$$f(p) = \pi_0 + (1 - \pi_0)f_{\text{alt}}(p),$$

$$\text{fdr}(p) = \frac{\pi_0}{f(p)}$$

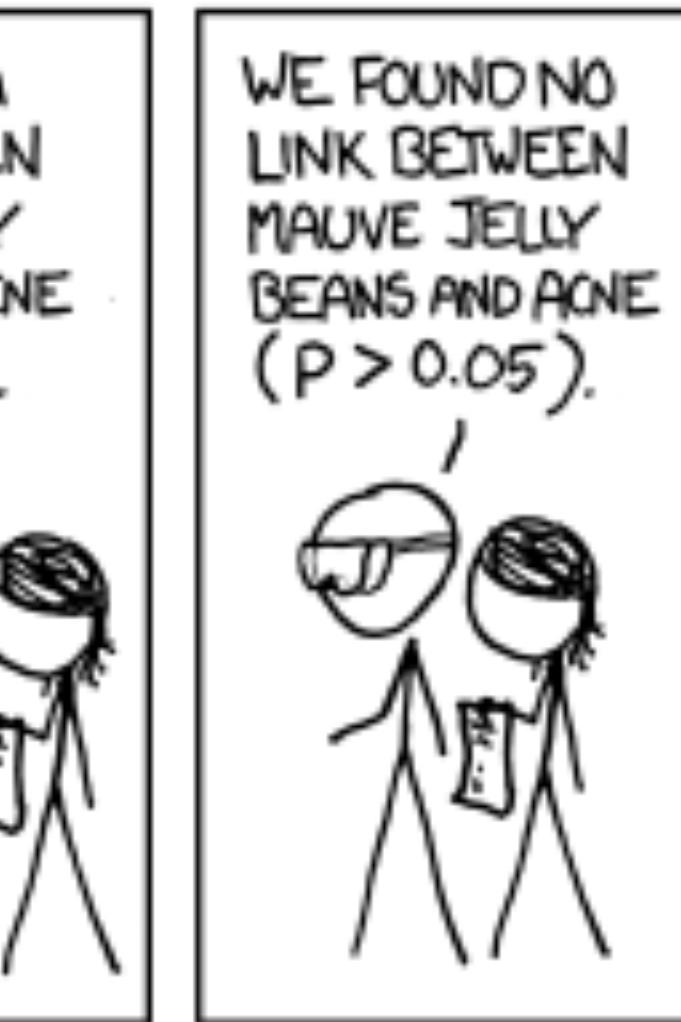
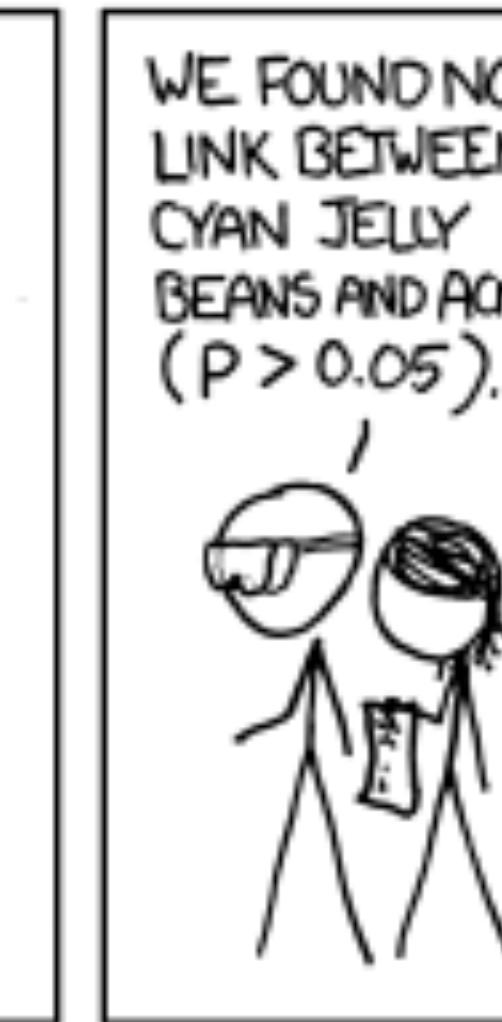
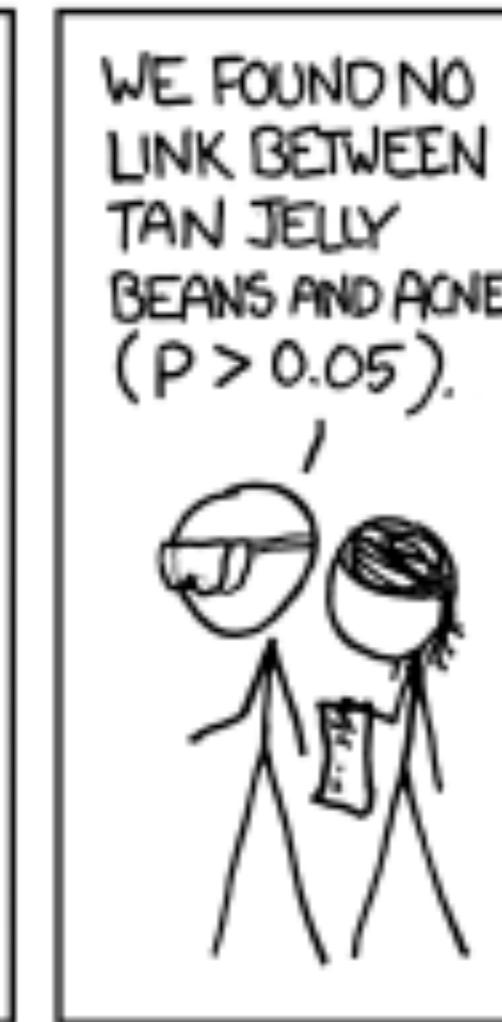
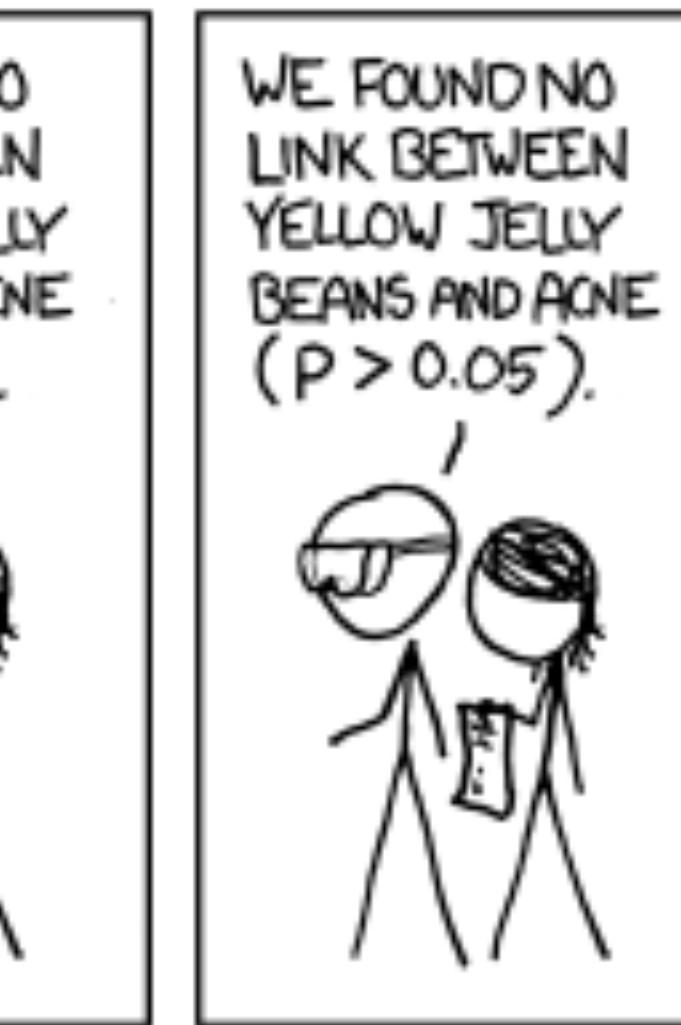
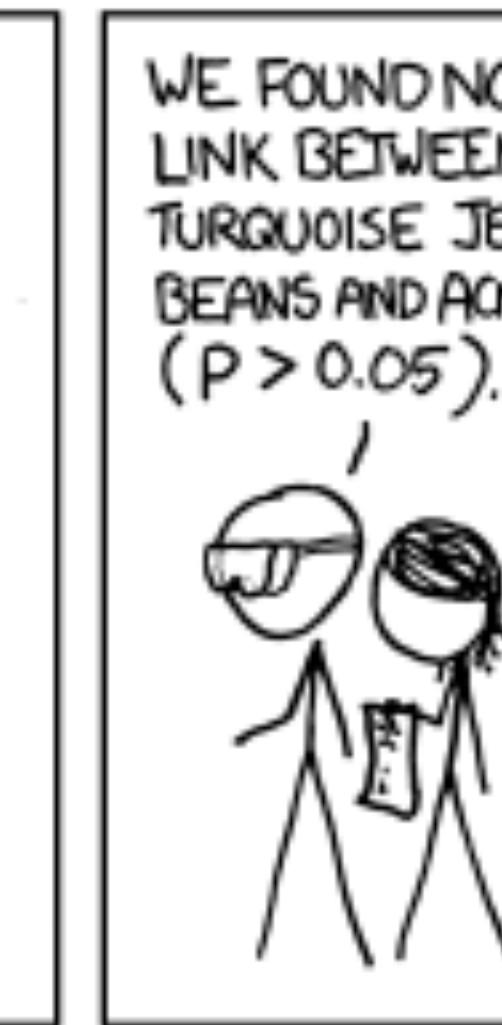
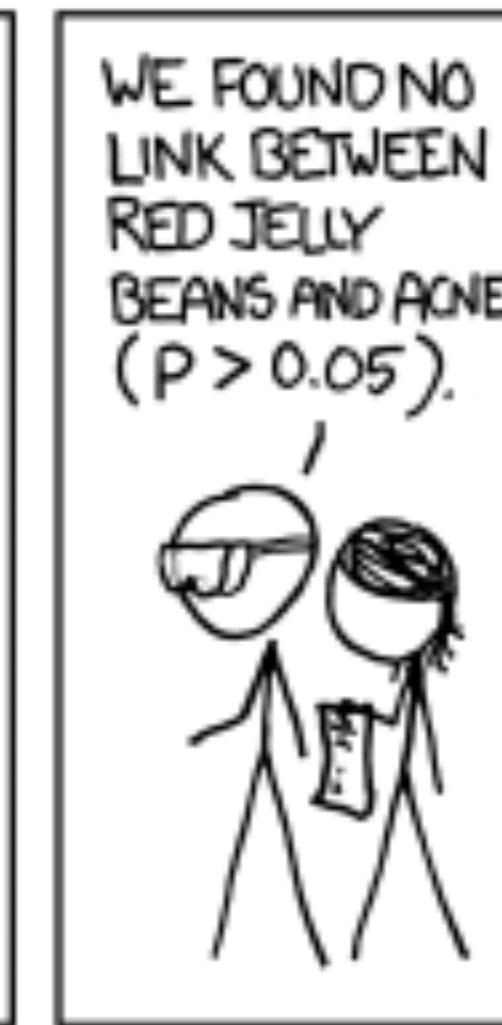
But how do we know π_0 and f_{alt} ?

the line
n
t

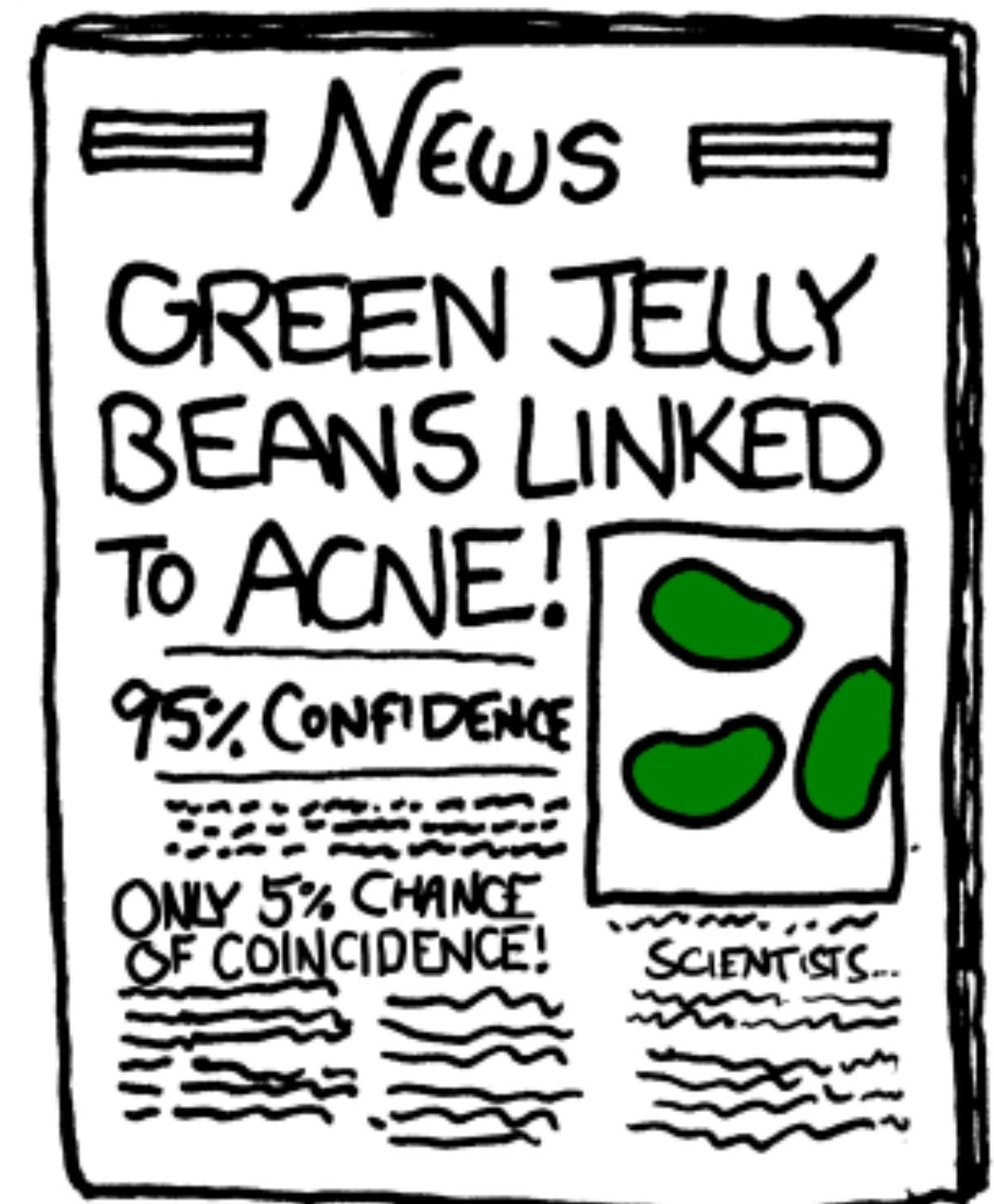
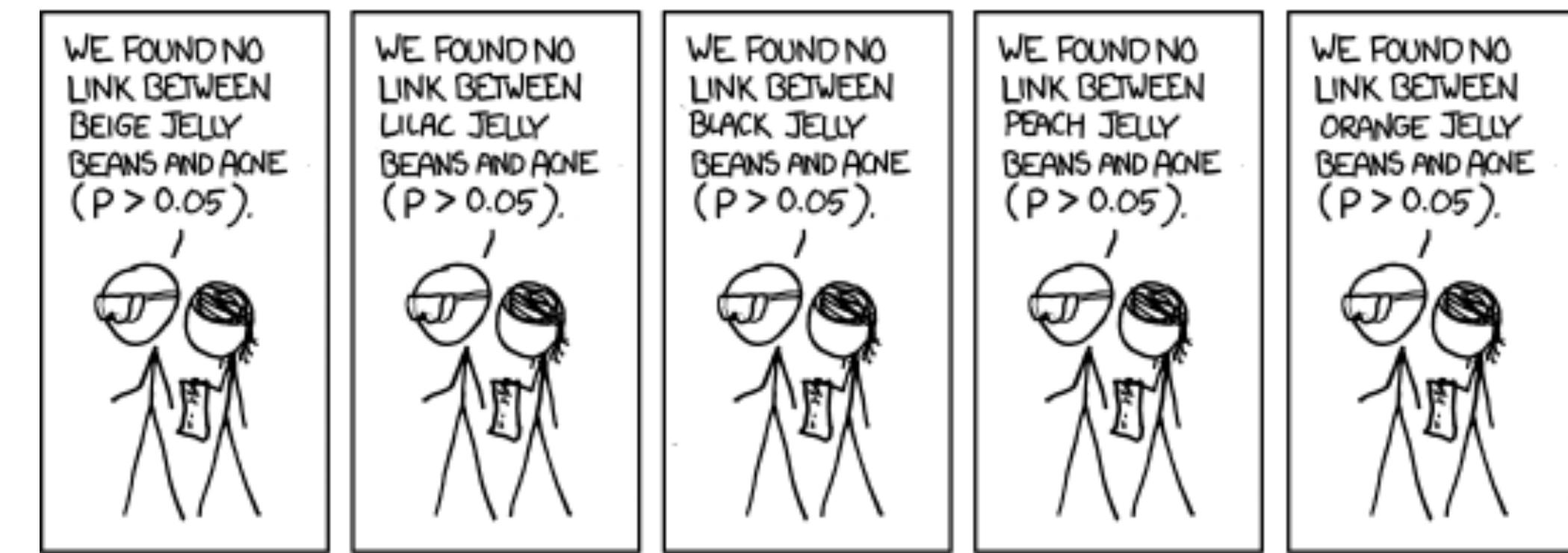
Multiple Testing



Multiple Testing



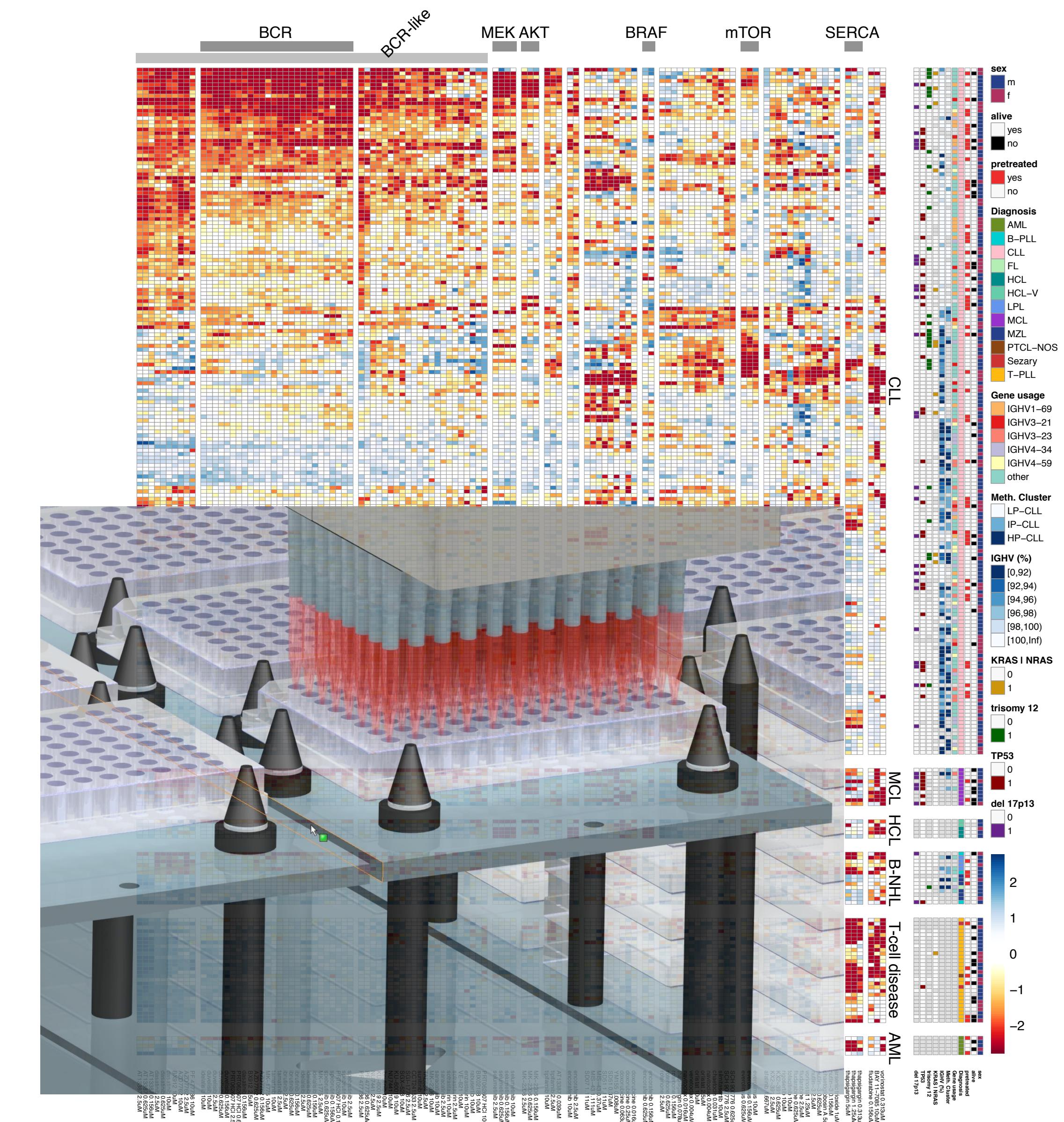
Multiple Testing



Multiple Testing

Many data analysis approaches in genomics employ item-by-item testing:

- Expression profiling
- Differential microbiome analysis
- Genetic or chemical compound screens
- Genome-wide association studies
- Proteomics
- Variant calling
- ...



False Positive Rate and False Discovery Rate

FPR: fraction of FP among all true negatives

FDR: fraction of FP among hits called

Example:
20,000 genes, 500 are d.e., 100 hits called, 10 of them wrong.

FPR: $10/19,500 \approx 0.05\%$

FDR: $10/100 = 10\%$



"Wait a minute! Isn't anyone here a real sheep?"

Experiment-Wide Type I Error Rates

Test vs Reality	Null Hypothesis is true	... is false	Total
Rejected	V	S	R
Not rejected	U	T	$m - R$
Total	m_0	$m - m_0$	m

- m : total number of hypotheses
- m_0 : number of null hypotheses
- V : number of false positives (a measure of type I error)

Family-wise error rate (FWER): The probability of one or more false positives, $P(V > 0)$. For large m_0 , this is difficult to keep small.

False discovery rate (FDR): The expected fraction of false positives among all discoveries, $E[V / \max \{R, 1\}]$.

NB: if $m_0 = m$, then FDR=FWER

The Multiple Testing Burden

When performing several tests, type I error goes up: for $\alpha = 0.05$ and n indep. tests, probability of no false positive result is

$$\underbrace{0.95 \cdot 0.95 \cdot \dots \cdot 0.95}_{n\text{-times}} \lll 0.95$$



Bonferroni Correction



For m tests, multiply each p -value with m .

Then see if anyone still remains below α .

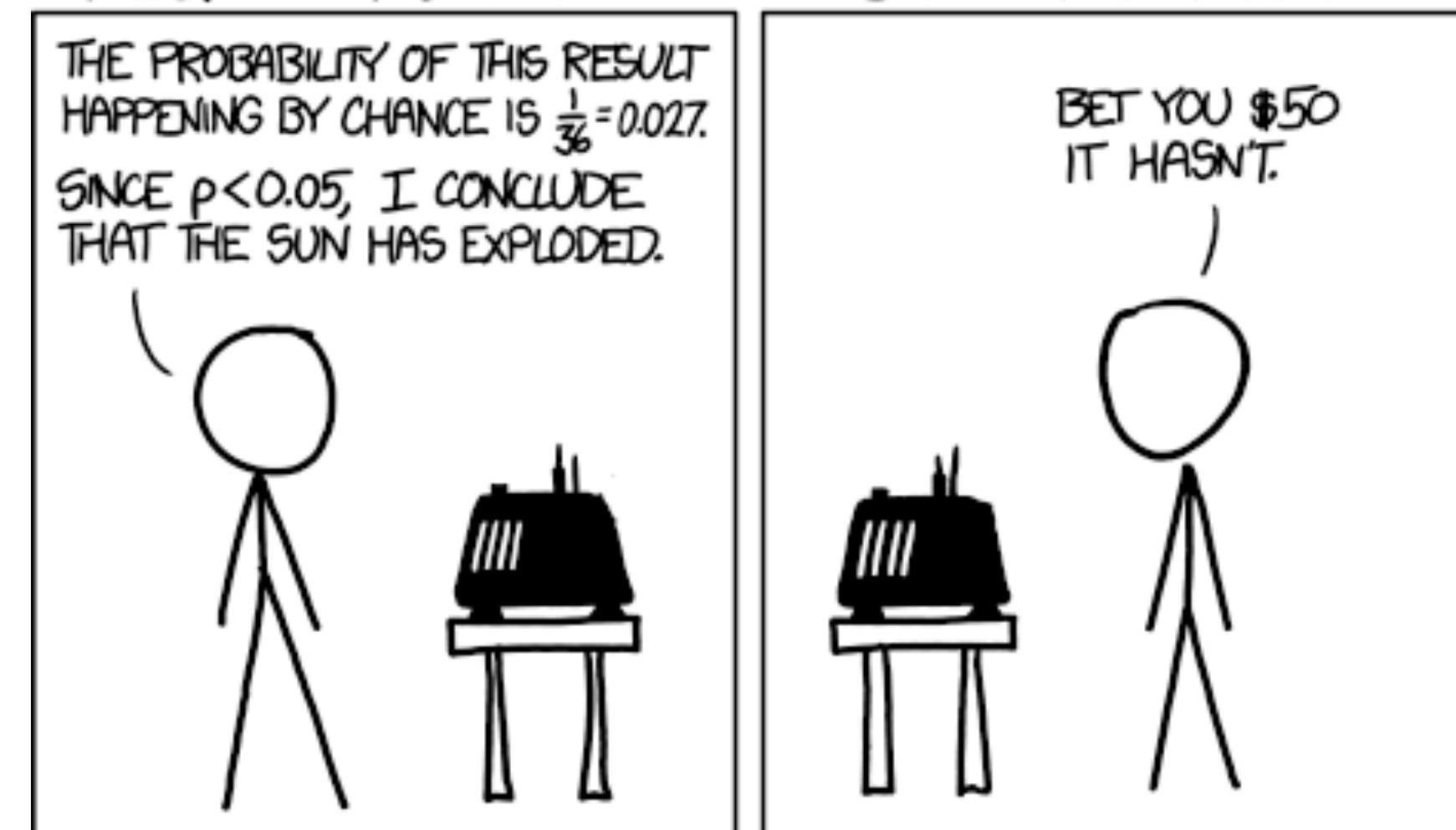
The Multiple Testing Opportunity

DID THE SUN JUST EXPLODE?
(IT'S NIGHT, SO WE'RE NOT SURE.)



FREQUENTIST STATISTICIAN: BAYESIAN STATISTICIAN:

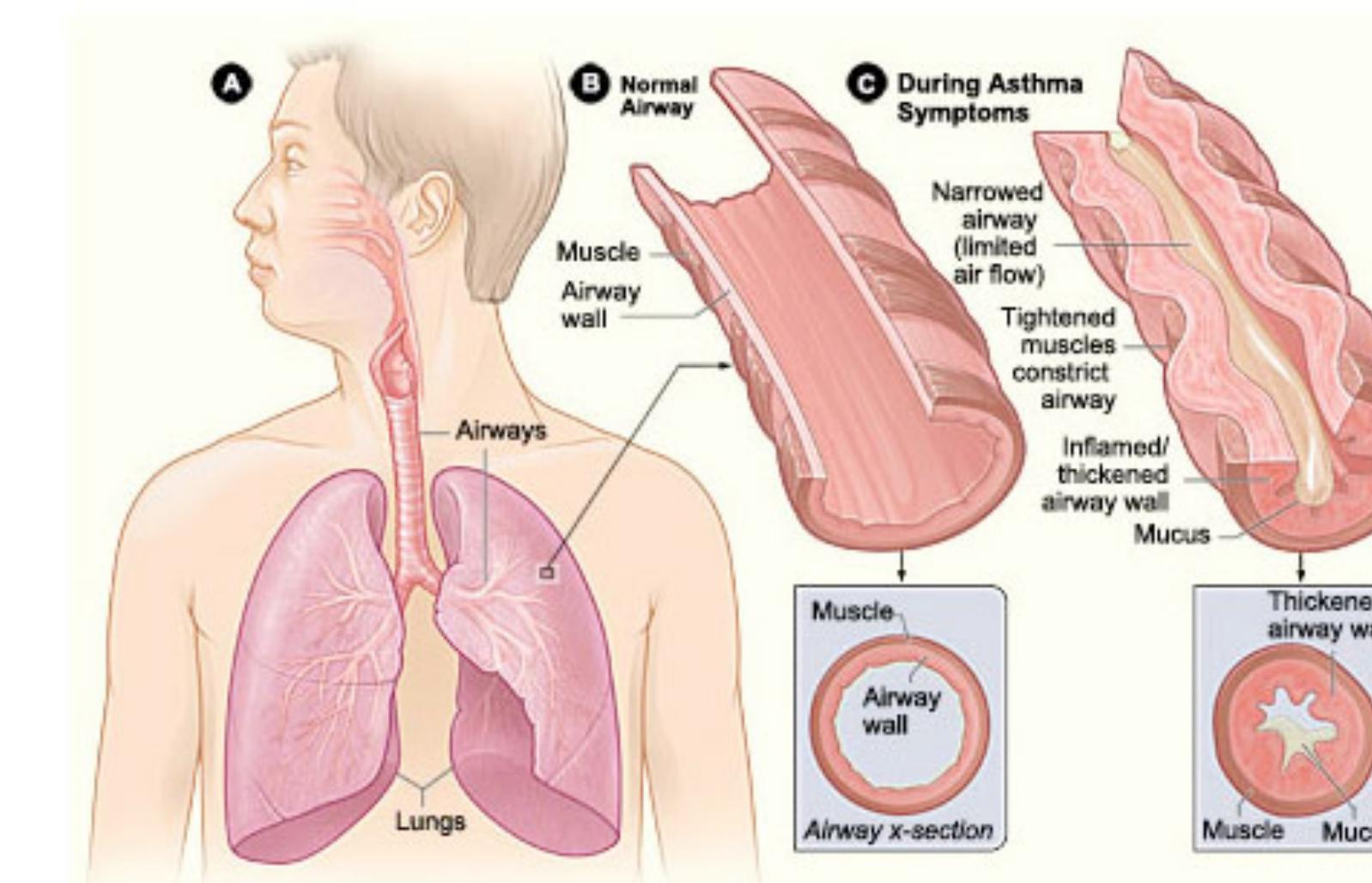
THE PROBABILITY OF THIS RESULT HAPPENING BY CHANCE IS $\frac{1}{36} = 0.027$. SINCE $p < 0.05$, I CONCLUDE THAT THE SUN HAS EXPLODED.



Data set 1: RNA-Seq

Transcriptome changes in four samples of primary human airway smooth muscle cells treated with dexamethasone, a synthetic glucocorticoid. 1 μM for 18 h.

cellline	dexamethasone
N61311	untrt
N61311	trt
N052611	untrt
N052611	trt
N080611	untrt
N080611	trt
N061011	untrt
N061011	trt



DESeq2 differential expression analysis:

gene i , sample j :

$$K_{ij} \sim \text{NB}(\text{mean} = \mu_{ij}, \text{dispersion} = \alpha_j)$$

$$\mu_{ij} = s_j q_{ij}$$

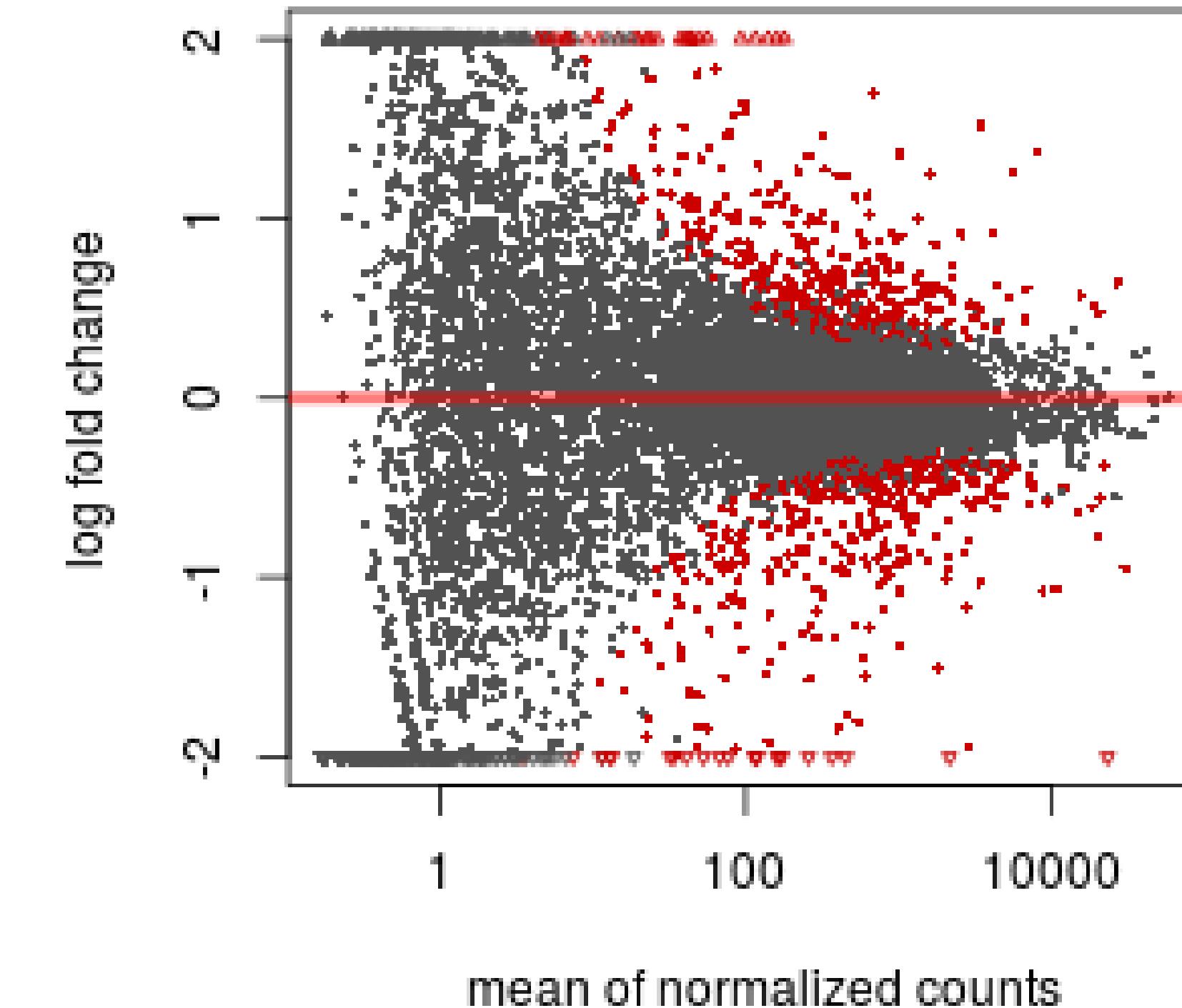
$$\log q_{ij} = \sum_r x_{jr} \beta_{rj}$$

design <- ~ cellline + dexamethasone

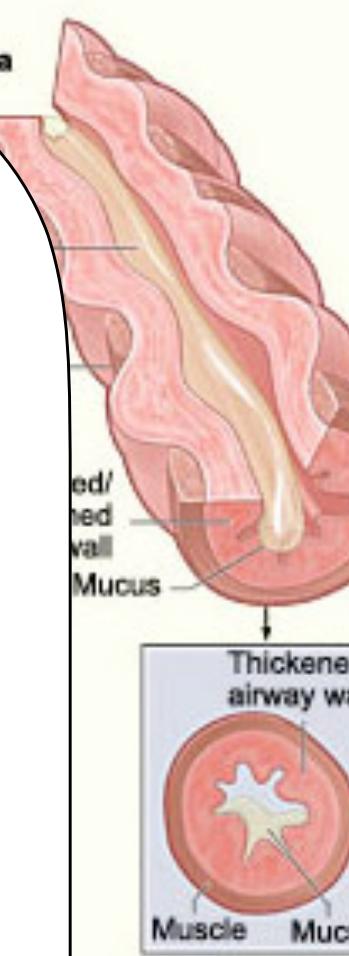
Data set 1: RNA-Seq

Transcription samples
smooth muscle
dexamethasone
glucocorticoids

cellline
N61011
N61011
N05
N05
N08
N08
N061011
N061011 trt



design <- ~ cellline + dexamethasone



analysis:

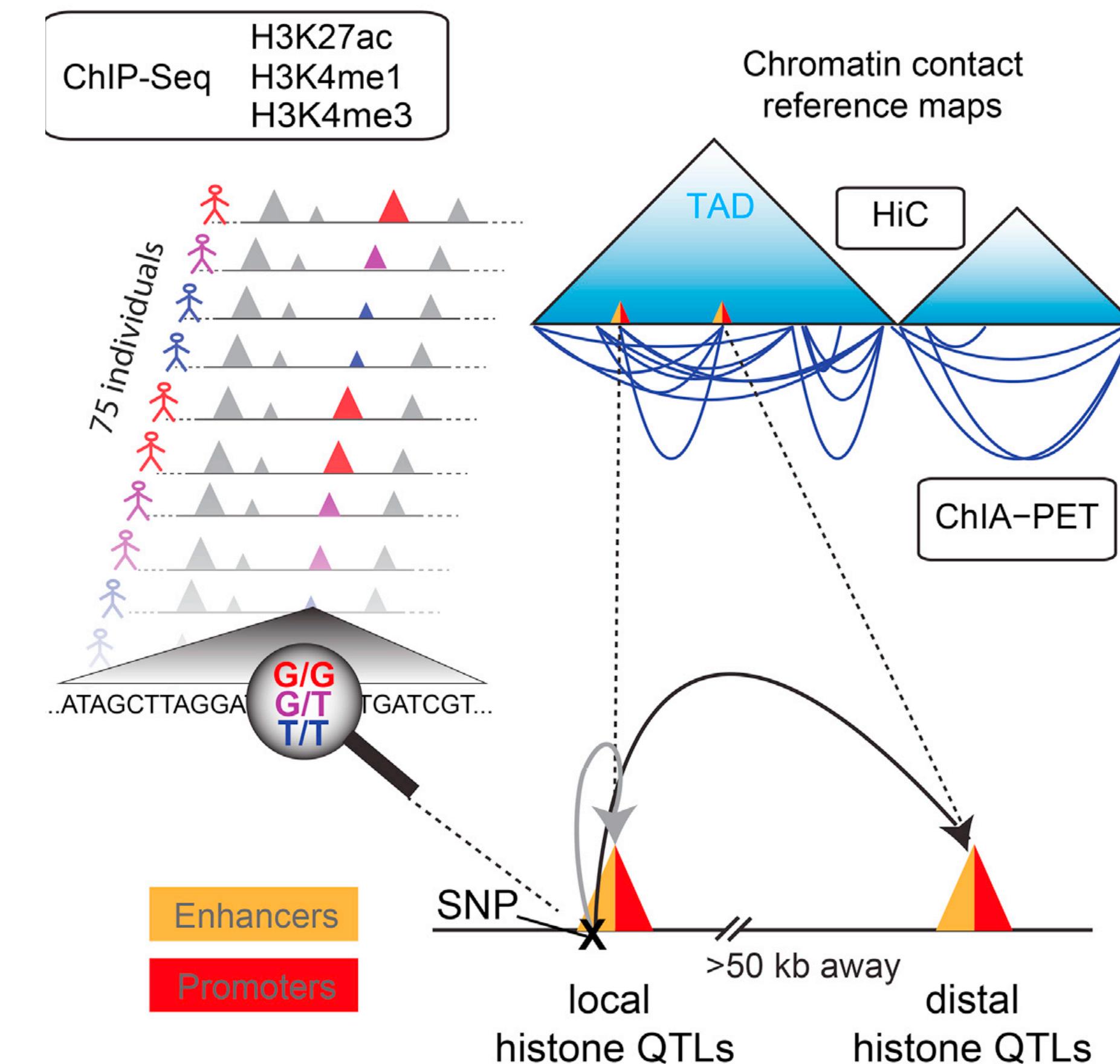
dispersion = α_j)

Data set 2: hQTL

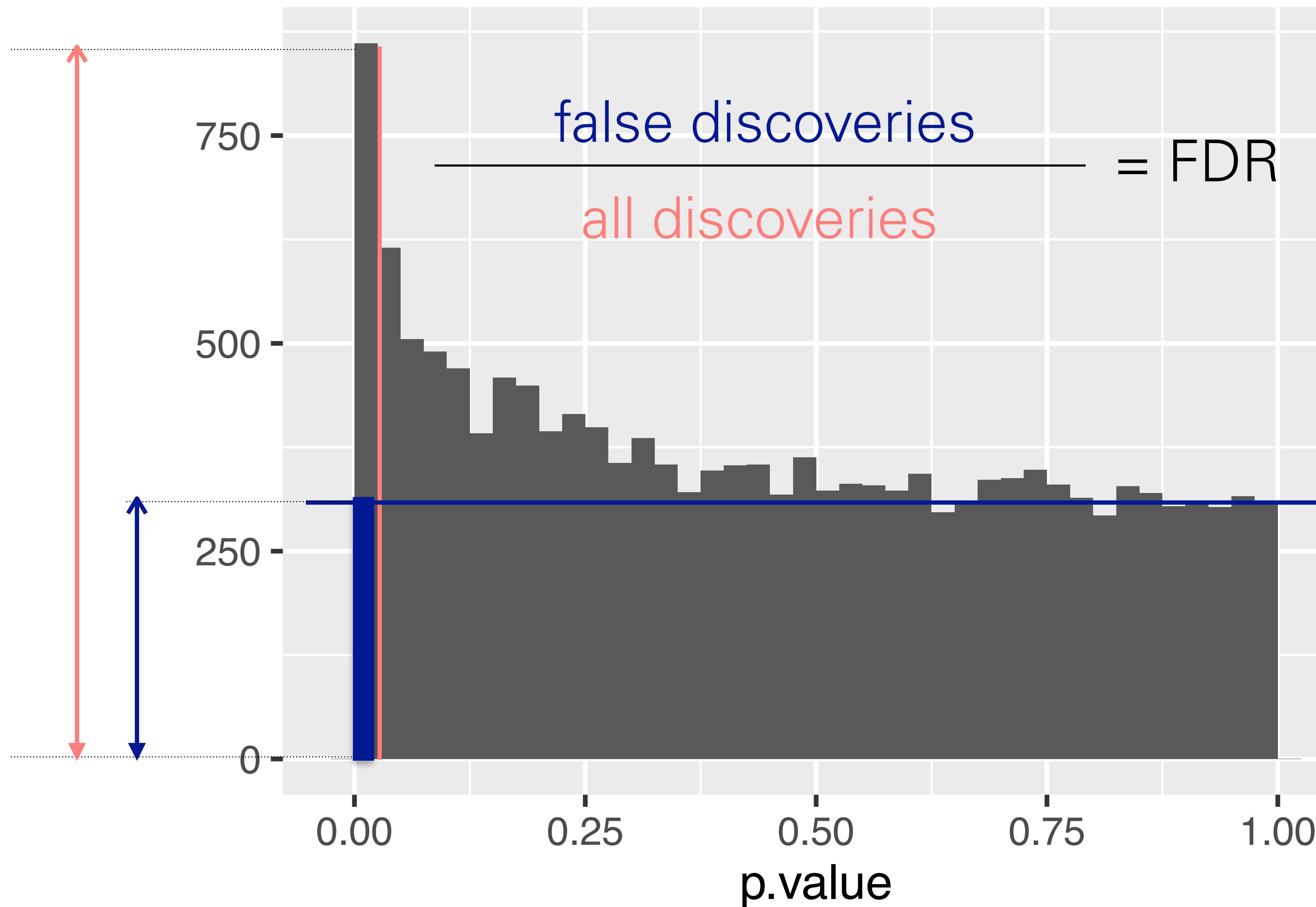
ChIP-seq for histone marks in lymphoblastoid cell lines from 75 sequenced individuals.

Local QTLs: find best-correlated SNP within 2kb of peak boundaries: 14,142 hQTLs, involving ~10% of all H3K27ac peaks (FDR=0.1, permutations)

Distal: distance cutoffs from 50 to 300 kb; also HiC

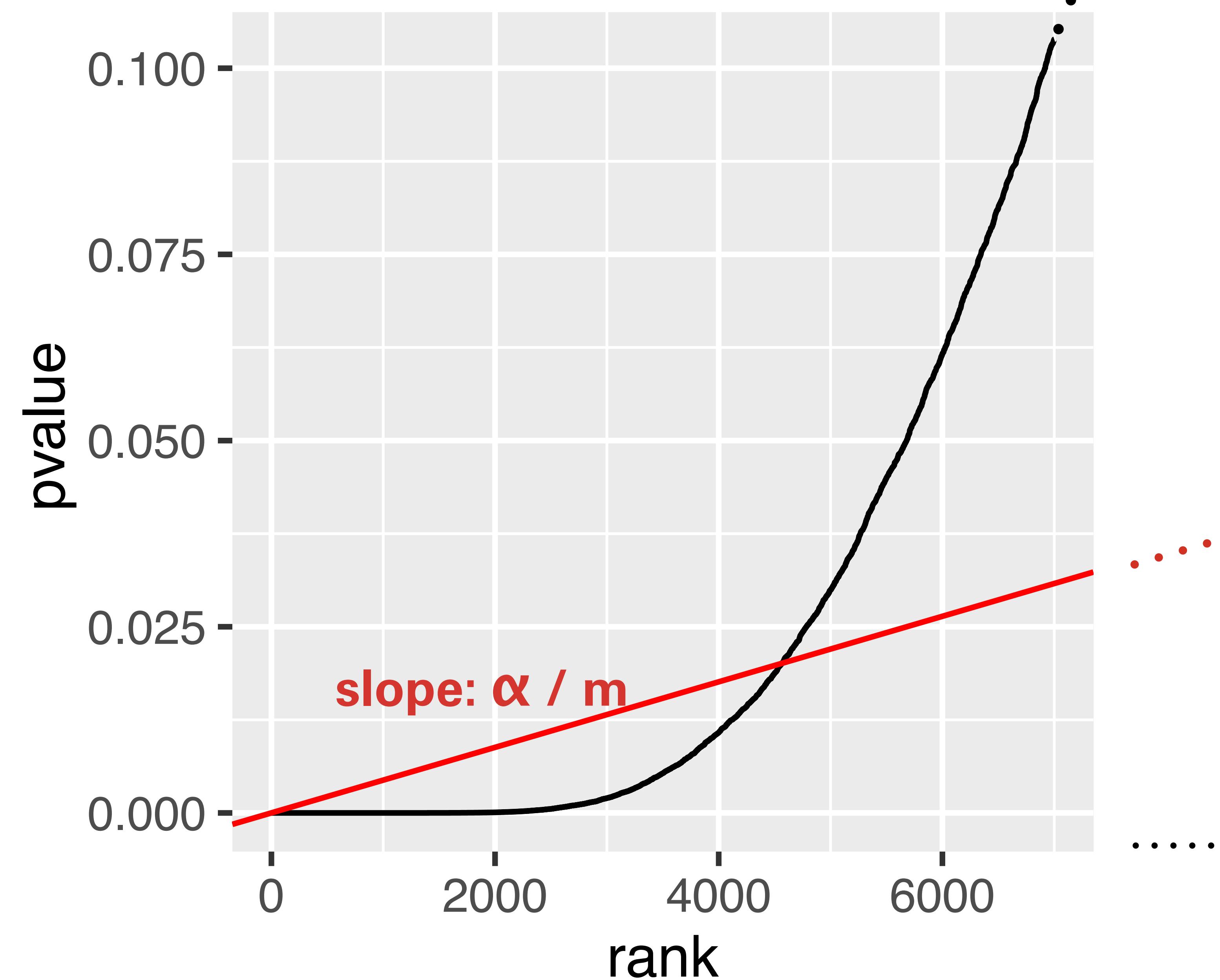


False Discovery Rate



Method of Benjamini & Hochberg (1995)

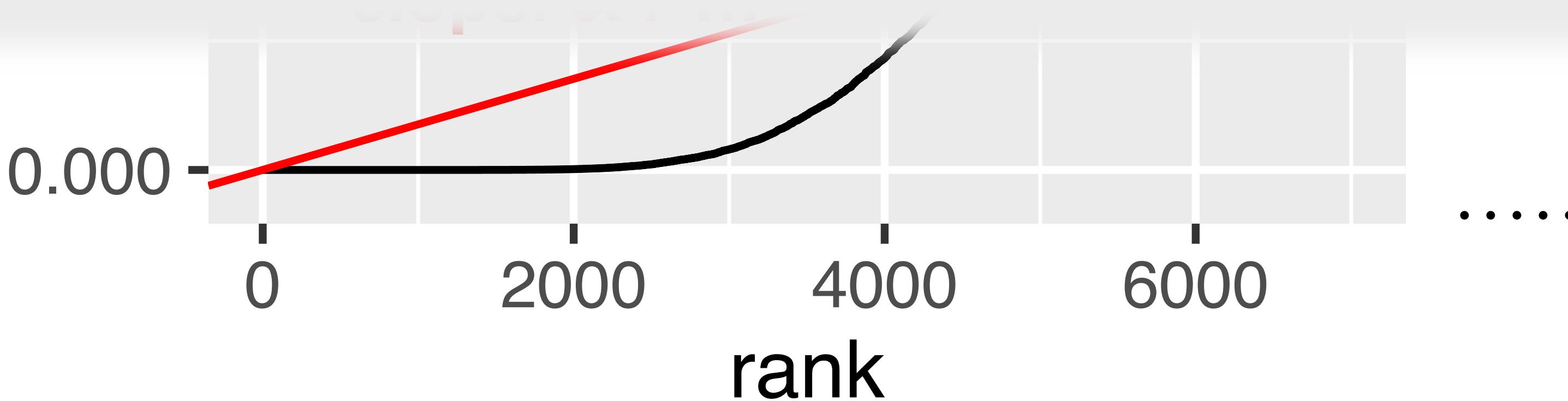
Method of Benjamini & Hochberg



Method of Benjamini & Hochberg

```
BH = {  
    i <- length(p) :1  
    o <- order(p, decreasing = TRUE)  
    ro <- order(o)  
    pmin(1, cummin(n/i * p[o])) [ro]  
}
```

takes a list of p-values as input and returns a matched list of 'adjusted' p-values.



Not all Hypothesis Tests are Created Equal

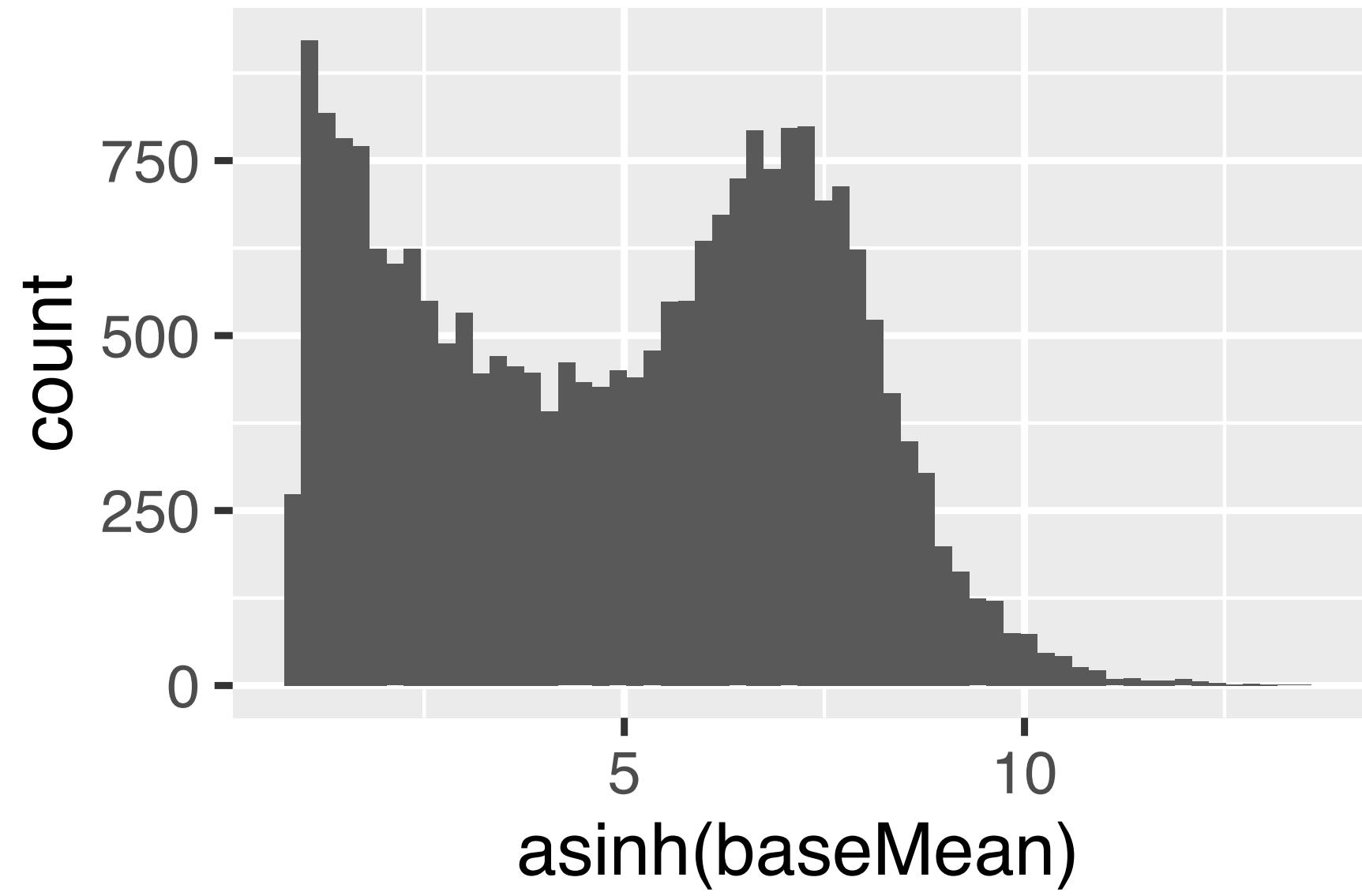
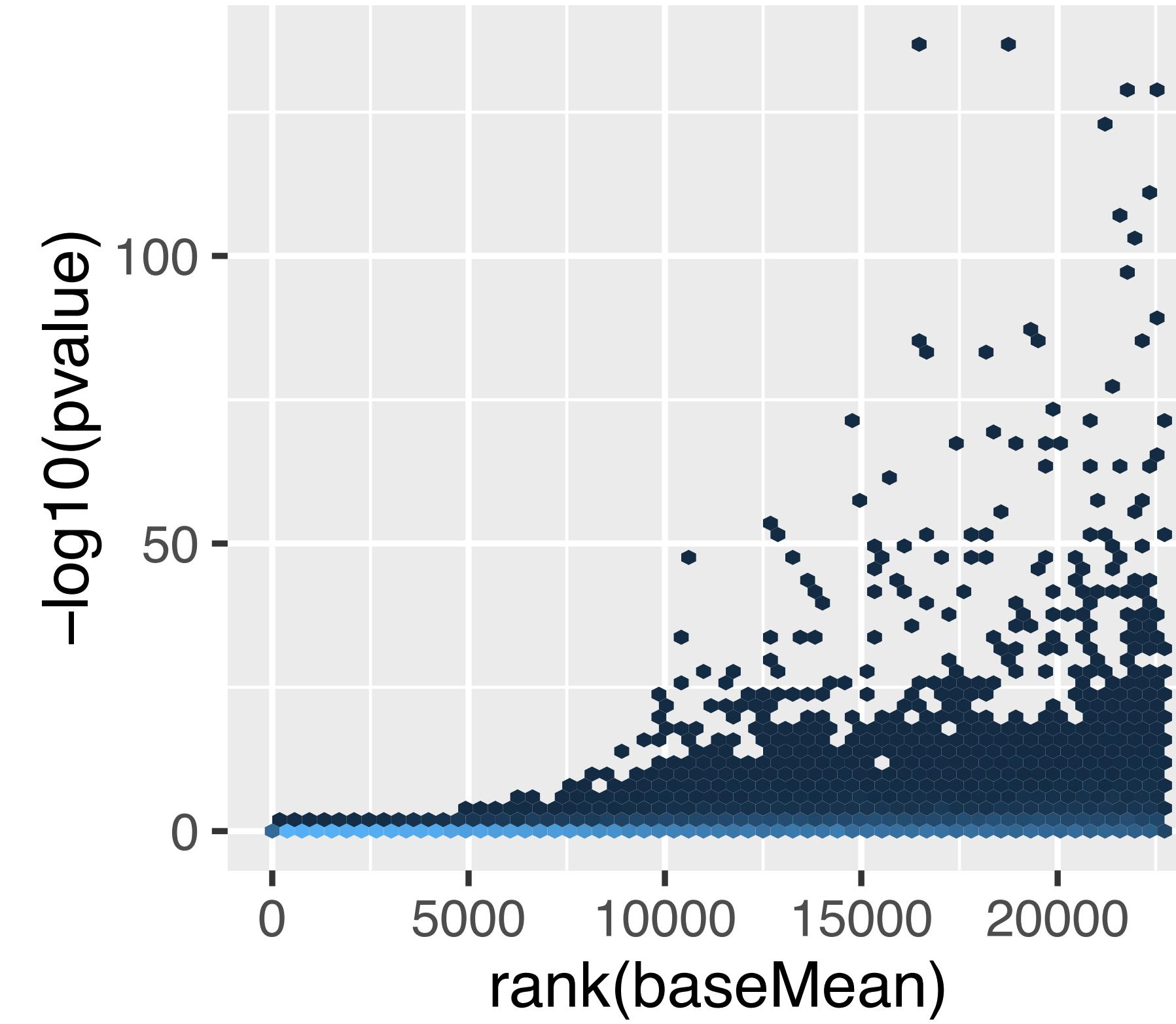


Figure 6.15: Histogram of `baseMean`. We see that it covers a large dynamic range, from close to 0 to around 3.3×10^5 .



Covariates - examples

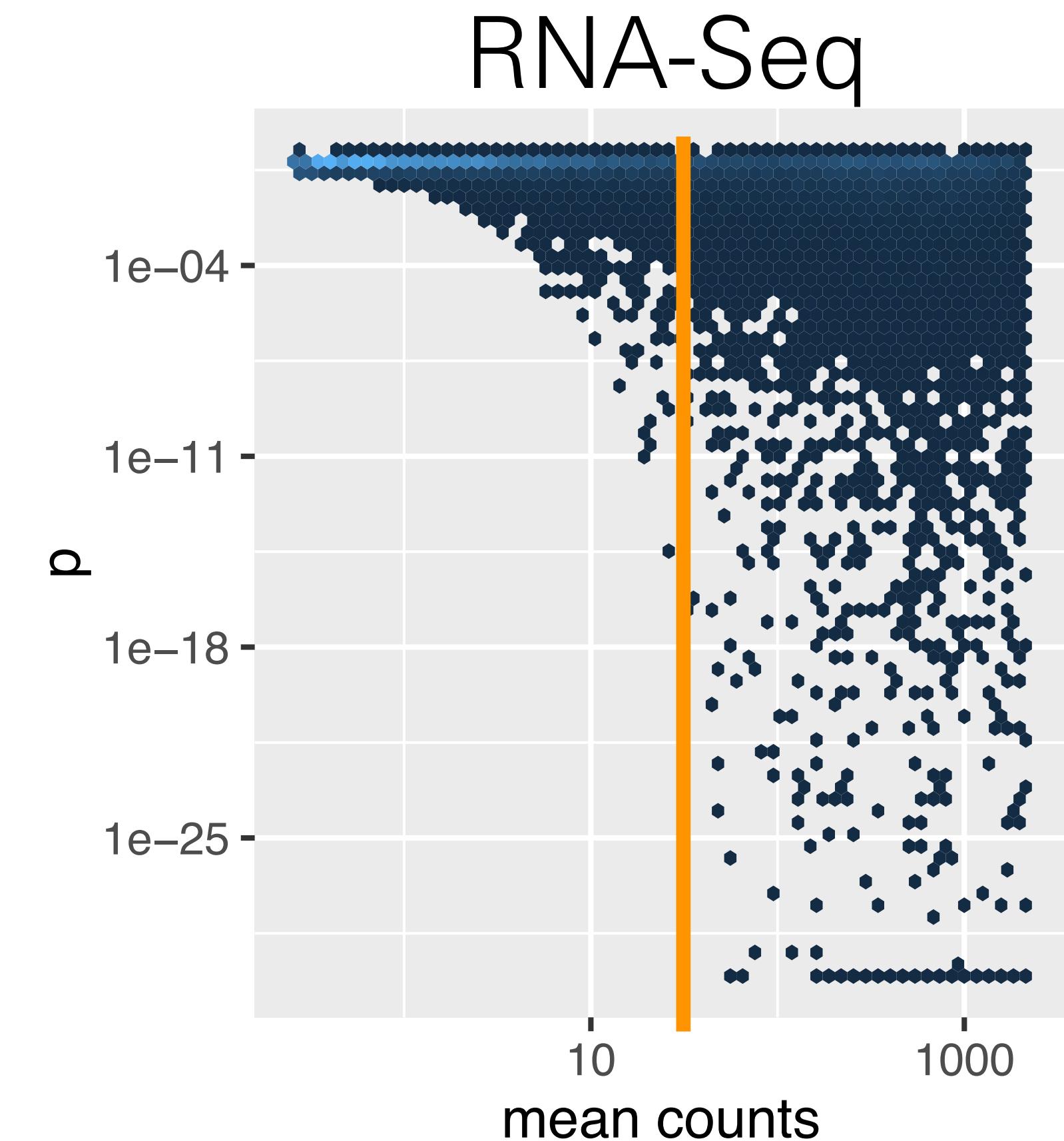
Application	Covariate
Differential RNA-Seq, ChIP-Seq, CLIP-seq, ...	(Normalized) mean of counts for each gene
eQTL analysis	SNP – gene distance
GWAS	Minor allele frequency
<i>t</i> -tests	Overall variance
Two-sided tests	Sign
All applications	Sample size; measures of signal-to-noise ratio

Independent Filtering

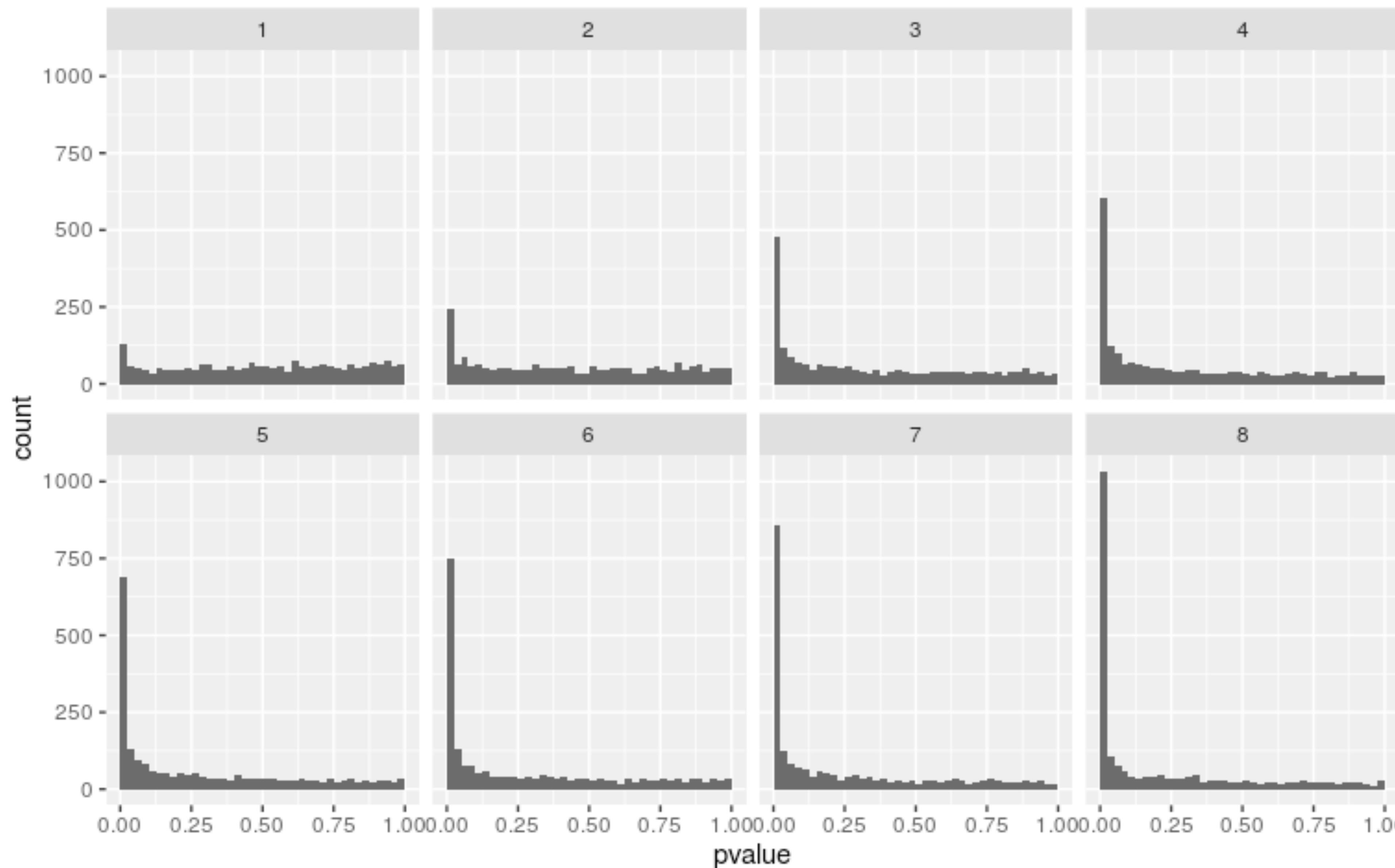
Two steps:

- All hypotheses H_i with $X_i < x$ get filtered.
- Apply BH to remaining hypotheses.

(Bourgon, Gentleman, Huber
PNAS 2010)



RNA-Seq p-value histogram stratified by average read count



Weighted Benjamini-Hochberg method

- Let $w_i \geq 0$ and $\frac{1}{m} \sum_{i=1}^m w_i = 1$ ("weight budget").
- Define $Q_i = P_i/w_i$.
- Apply BH to Q_i instead of P_i .
- Proven Type-I error (FDR) control (Genovese, Roeder, Wasserman *Biometrika* 2006).
- If $w_i > 1$, then H_i is easier to reject.
- $Q_i \leq t \Leftrightarrow P_i \leq w_i t =: t_i$

Weighted Benjamini-Hochberg method

- Let $w_i \geq 0$ and $\frac{1}{m} \sum_{i=1}^m w_i = 1$ ("weight budget").
- Define $Q_i = P_i/w_i$.
- Apply BH to Q_i instead of P_i .
- Proven Type I error control under, Wasserman et al.
- If $w_i > 1$,
 $Q_i \leq t \Leftrightarrow P_i \leq t w_i$
- $Q_i \leq t \Leftrightarrow P_i \leq t w_i$



Weighted Benjamini-Hochberg method

- Let $w_i \geq 0$ and $\frac{1}{m} \sum_{i=1}^m w_i = 1$ ("the budget").
- Define $Q_i = P_{\text{reject}}(H_0)$.
- Apply BH.
- Proven by Wasserman et al. (1996).
- If $w_i > 1$, then $Q_i < t$.
- $Q_i \leq t \Leftrightarrow Q_i / w_i \leq t / w_i$.

Problem: how to know the weights?



Independent hypothesis weighting (IHW): basic idea

- Stratify the tests into G bins, by covariate X
- Choose α
- For each possible weight vector $\mathbf{w} = (w_1, \dots, w_G)$ apply weighted BH procedure. Choose \mathbf{w} that maximizes the number of rejections at level α .
- Report the result with the optimal weight vector \mathbf{w}^* .



Nikos Ignatiadis

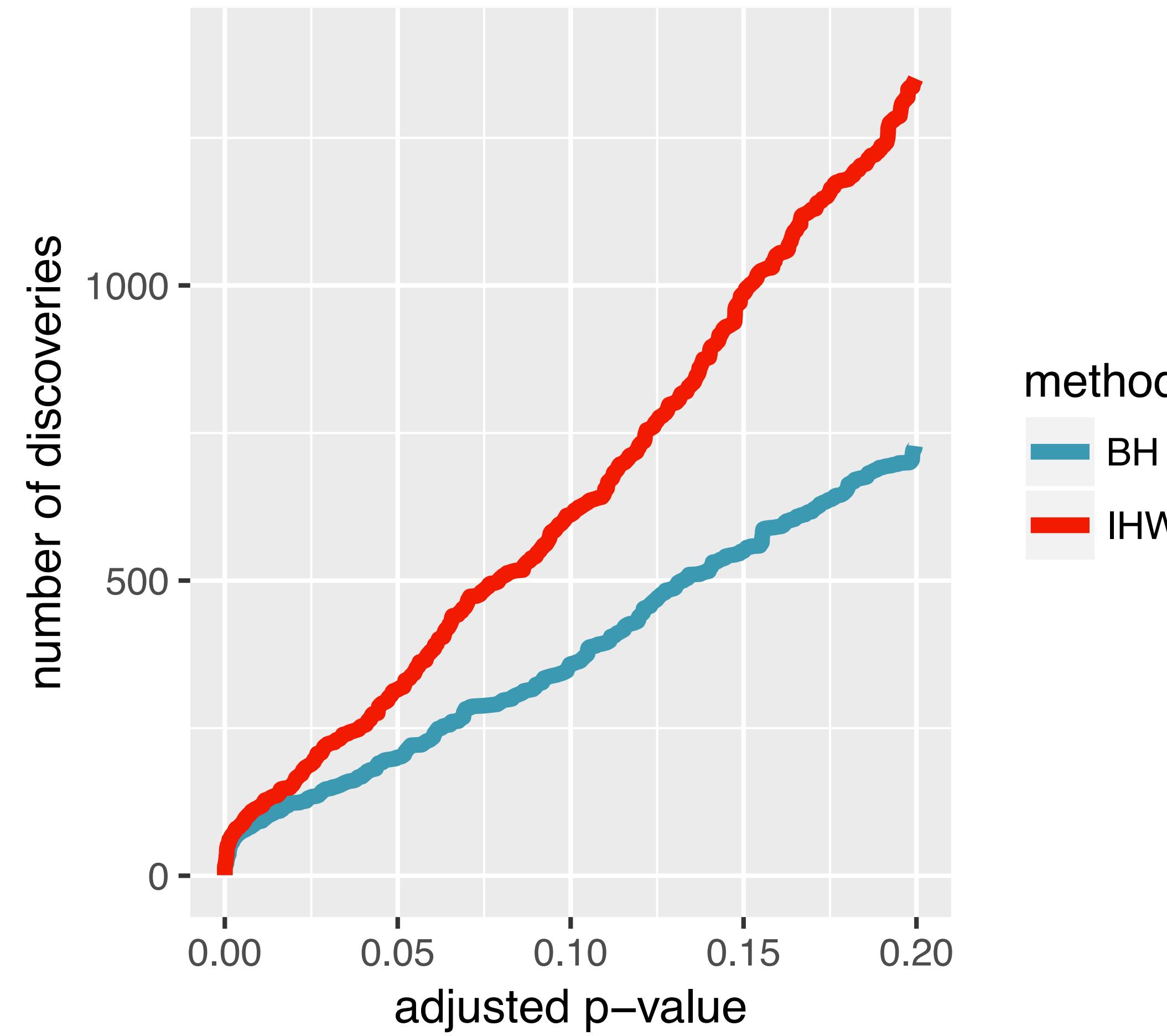
Ignatiadis et al.,

- Nature Methods 2016, DOI10.1038/nmeth.3885
- arXiv:1701.05179

Bioconductor package IHW

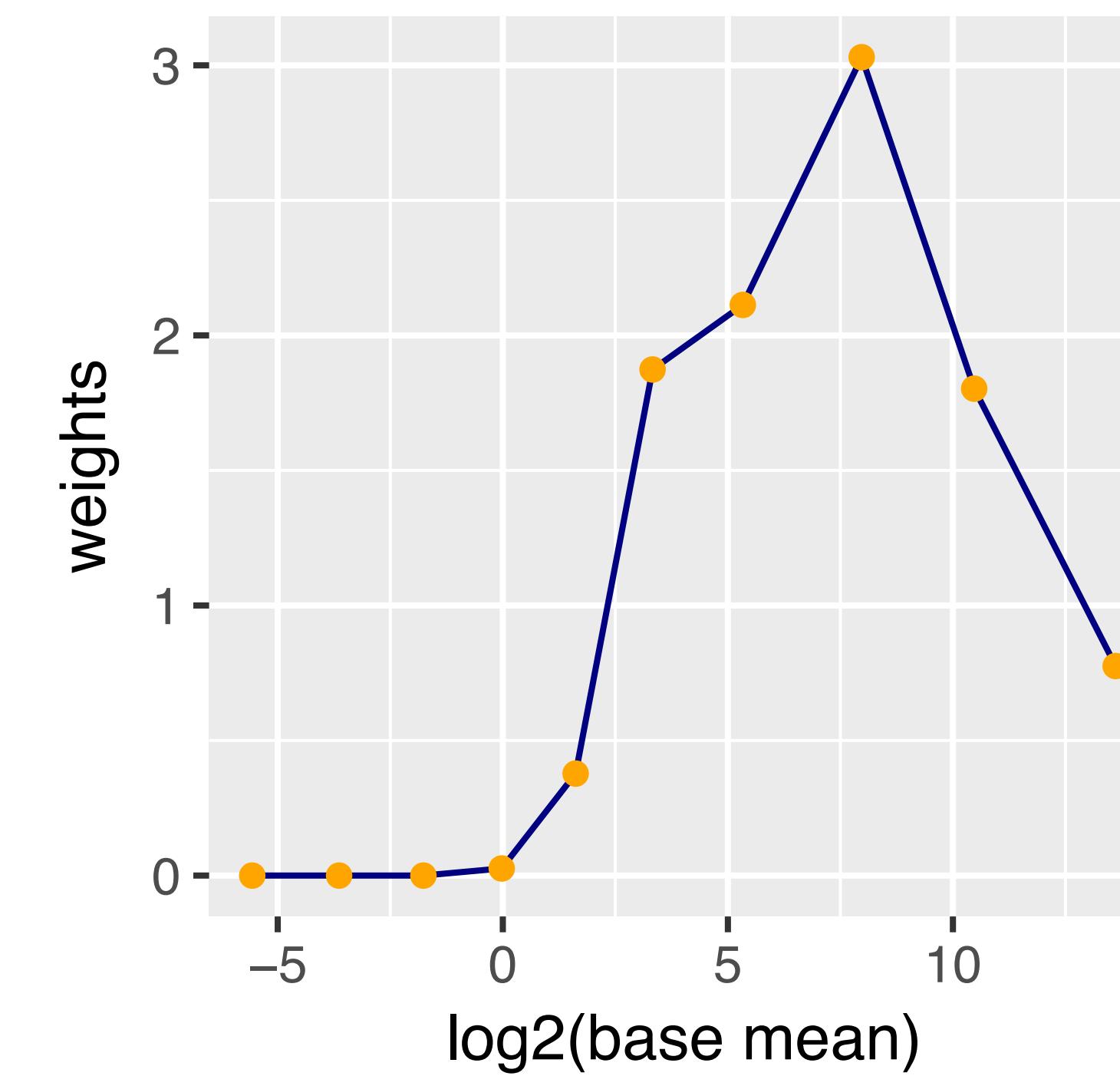
RNA-Seq example (DESeq2)

power

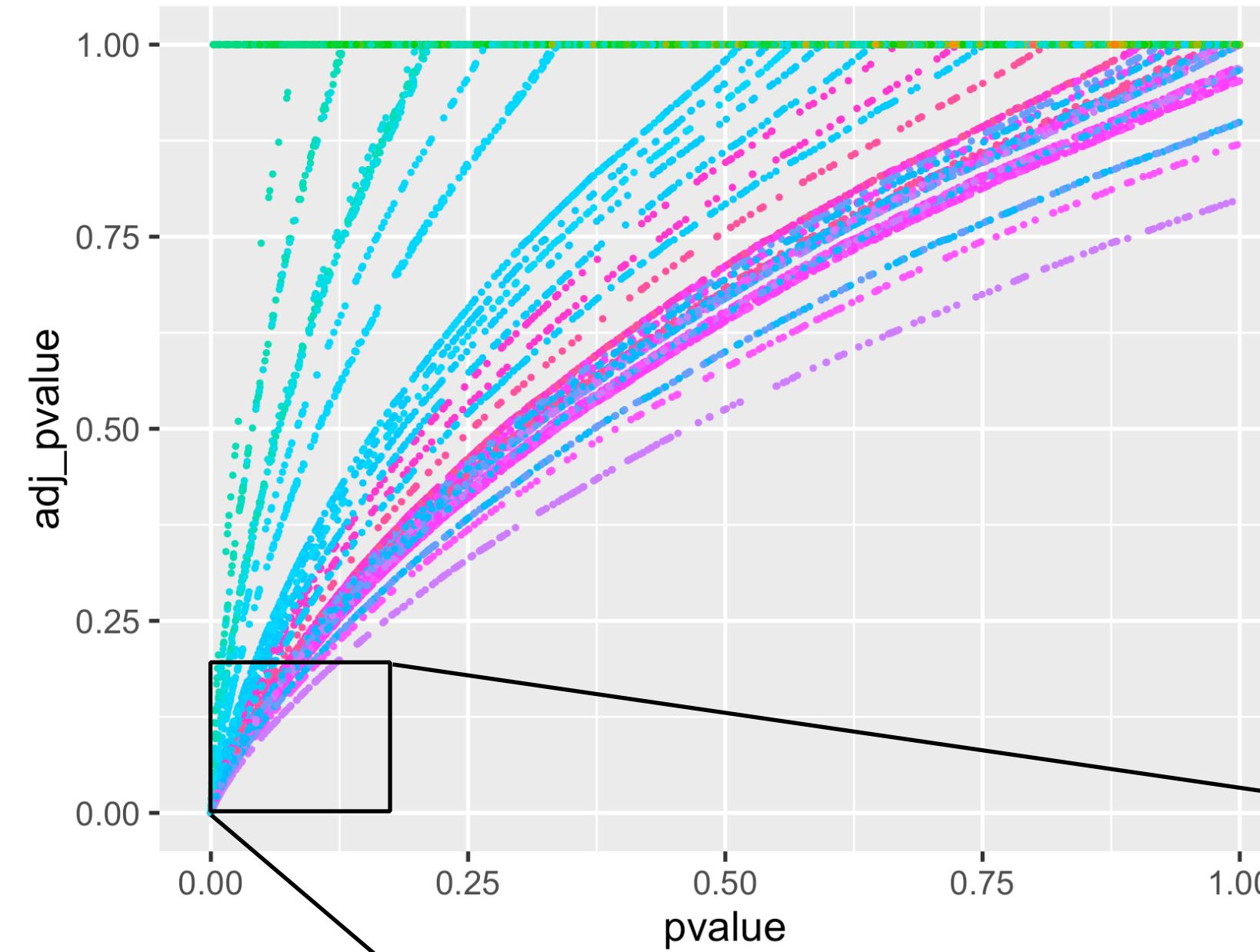


method
BH
IHW

weights

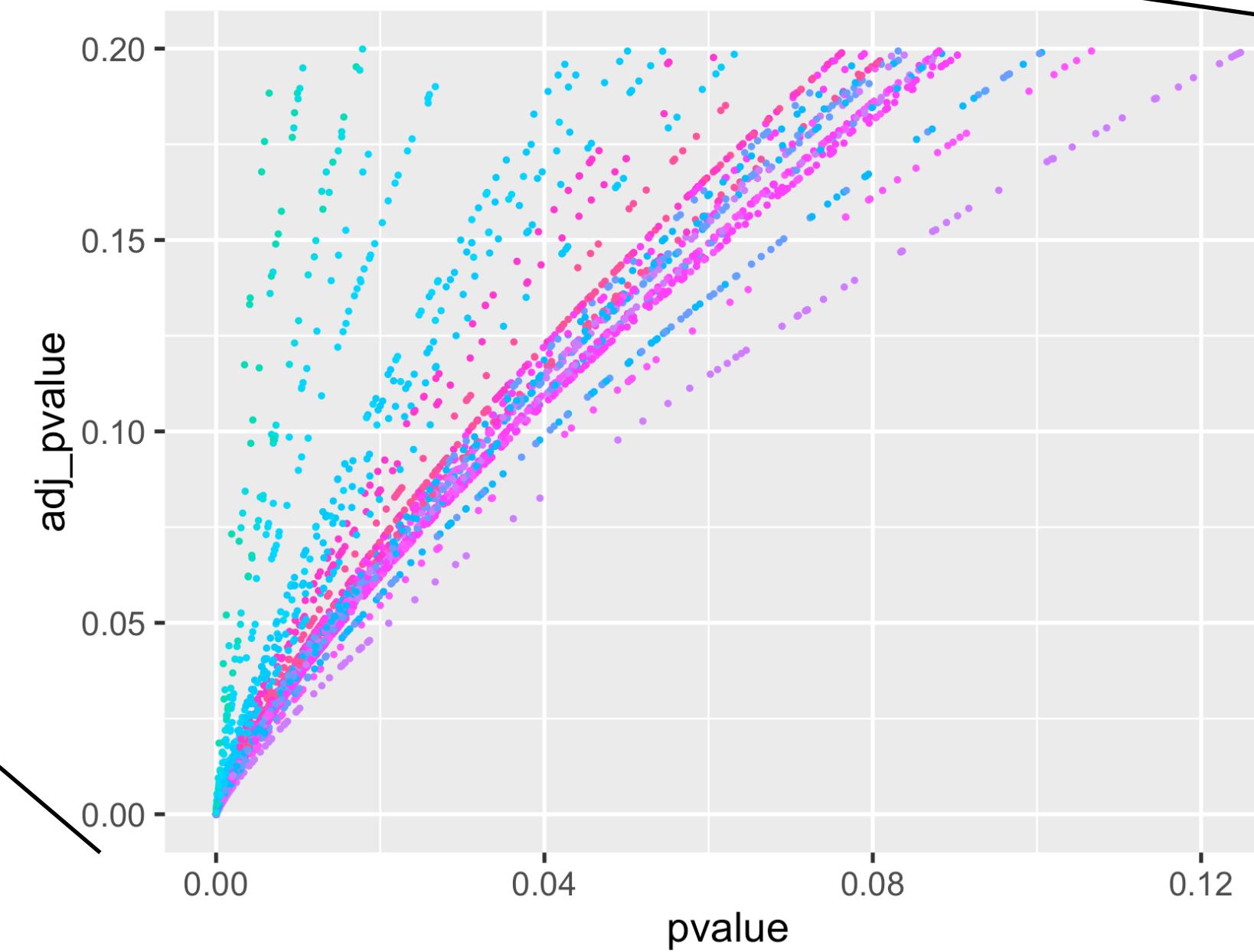


Ranking is not monotonous in raw p-values

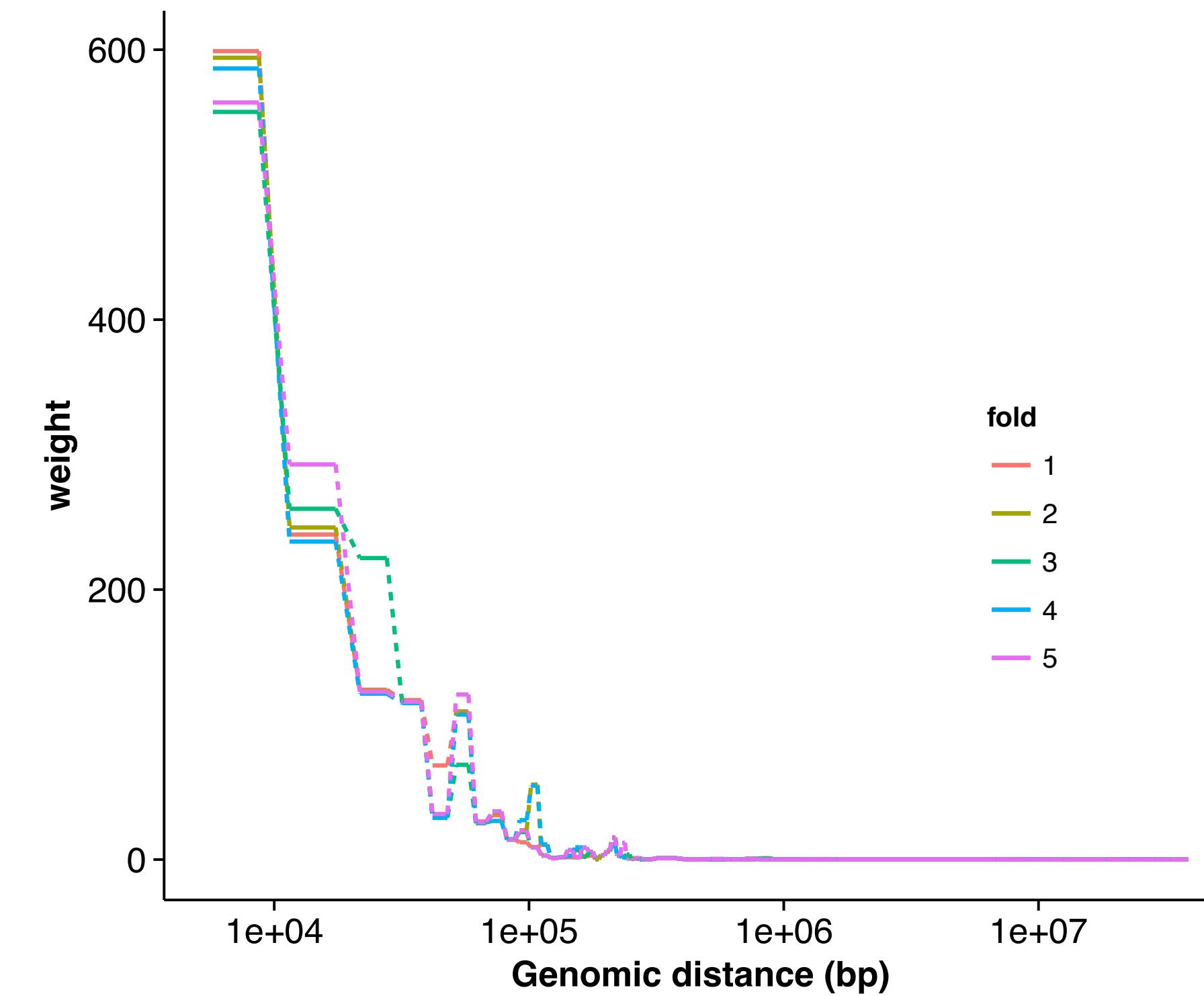
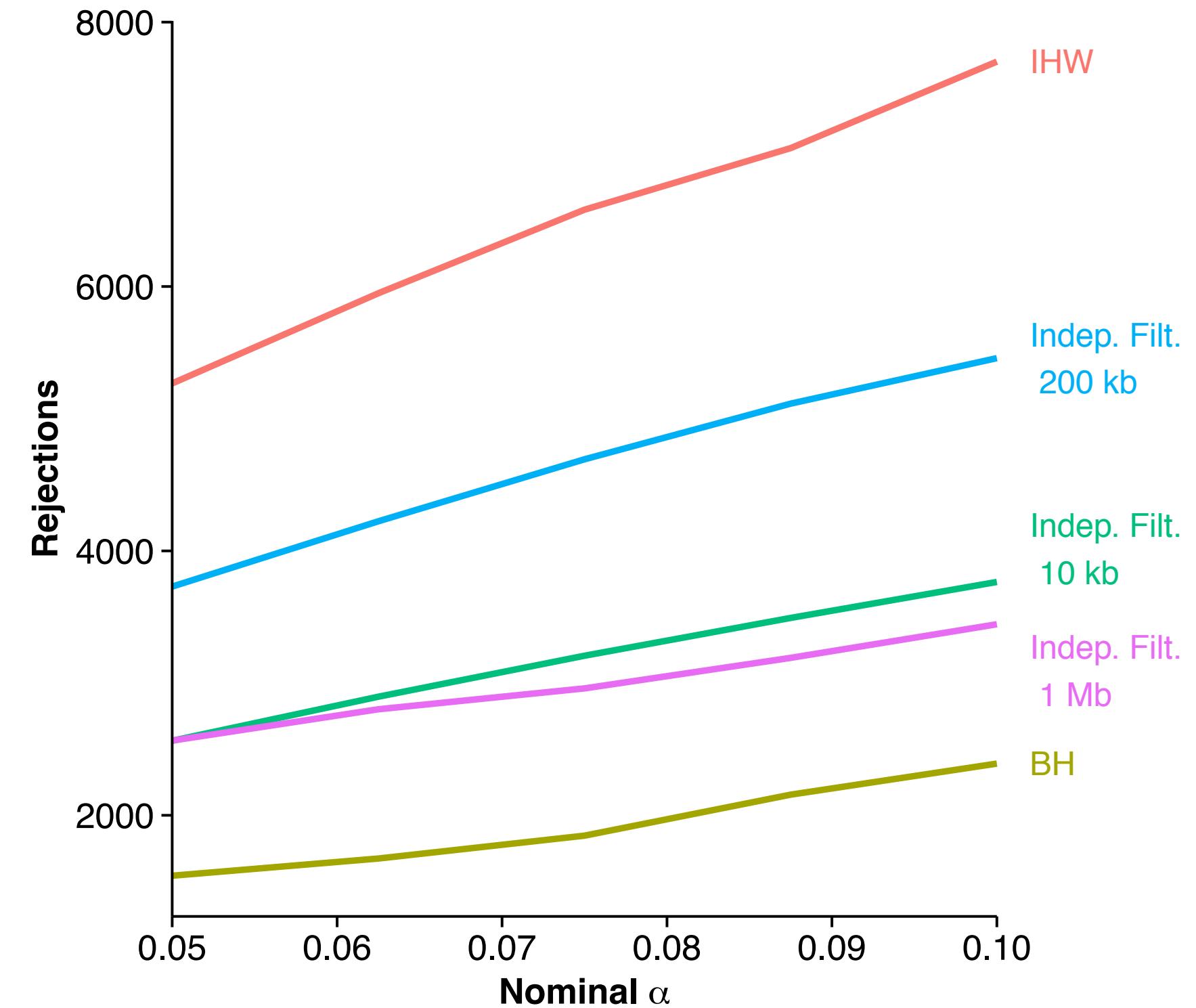


group

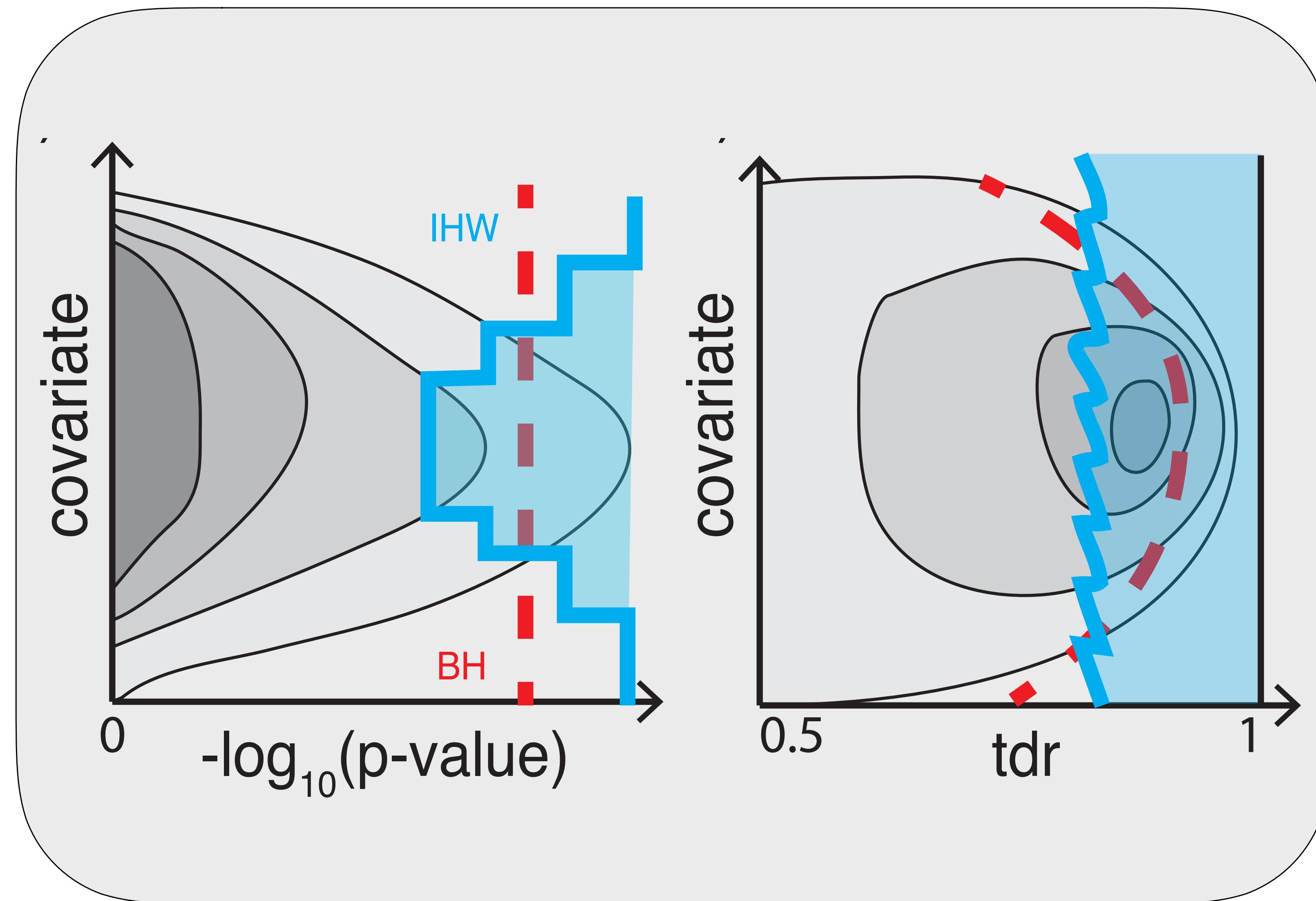
group	color
1	red
2	orange
3	yellow
4	light green
5	green
6	purple
7	magenta
8	pink
9	light blue
10	blue
11	cyan
12	light red
13	light orange
14	light yellow
15	light green
16	light purple
17	light magenta
18	light pink
19	light blue
20	light cyan
21	light red
22	light orange



Histone-QTL example (H3K27ac)

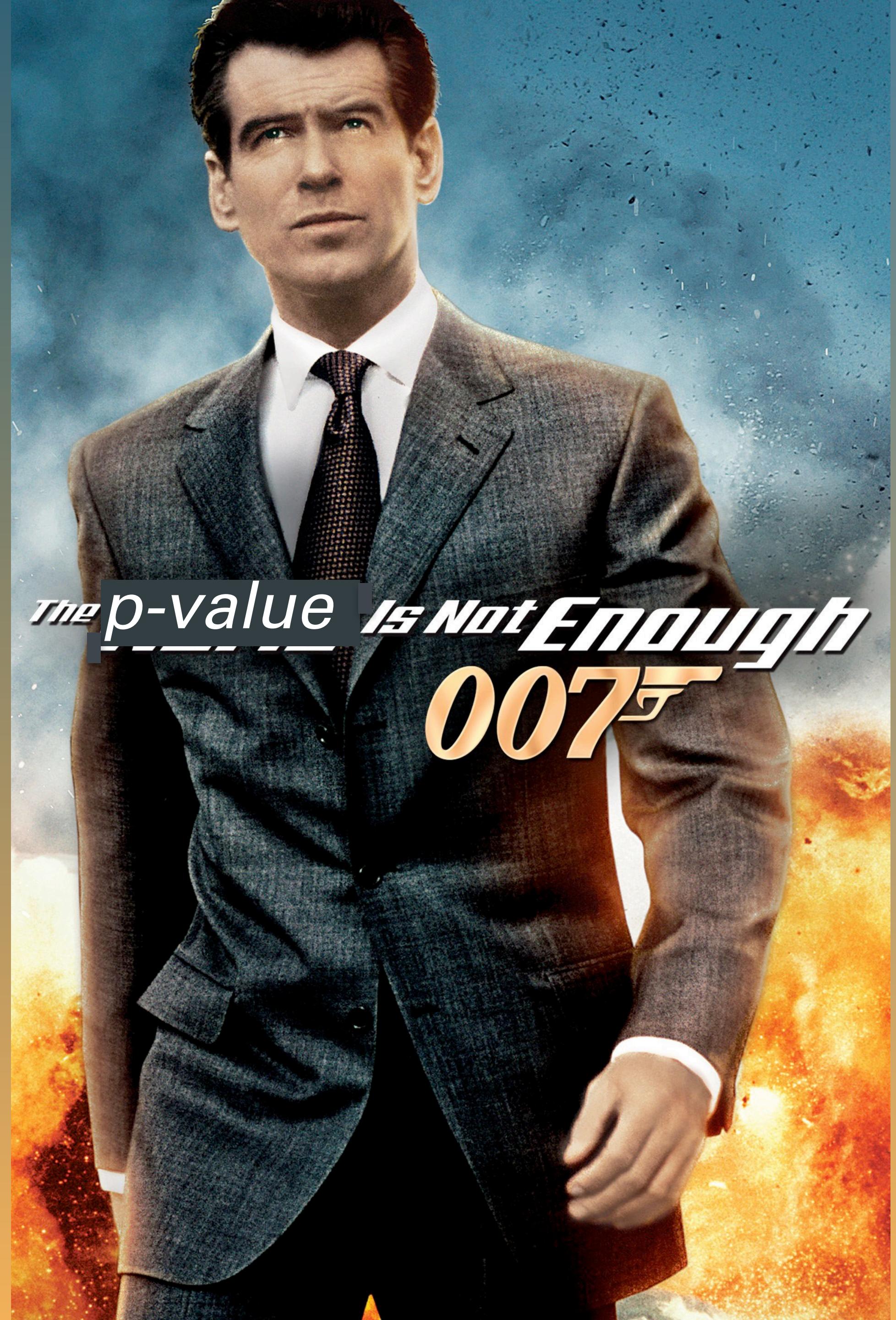


2D decision boundaries



Summary

- Multiple testing is not a problem but an opportunity
- Heterogeneity across tests
- Informative covariates are often apparent to domain scientists
 - independent of test statistic under the null
 - informative on π_1 , F_{alt}
- Data-driven weighting
- Scales well to millions of hypotheses
- Controlling ‘overoptimism’



<u>P-VALUE</u>	<u>INTERPRETATION</u>
0.001	
0.01	
0.02	HIGHLY SIGNIFICANT
0.03	
0.04	
0.049	SIGNIFICANT
0.050	OH CRAP. REDO CALCULATIONS.
0.051	ON THE EDGE OF SIGNIFICANCE
0.06	
0.07	HIGHLY SUGGESTIVE,
0.08	SIGNIFICANT AT THE $p < 0.10$ LEVEL
0.09	
0.099	HEY, LOOK AT
≥ 0.1	THIS INTERESTING SUBGROUP ANALYSIS