

Hypothesis Testing

Wolfgang Huber, EMBL



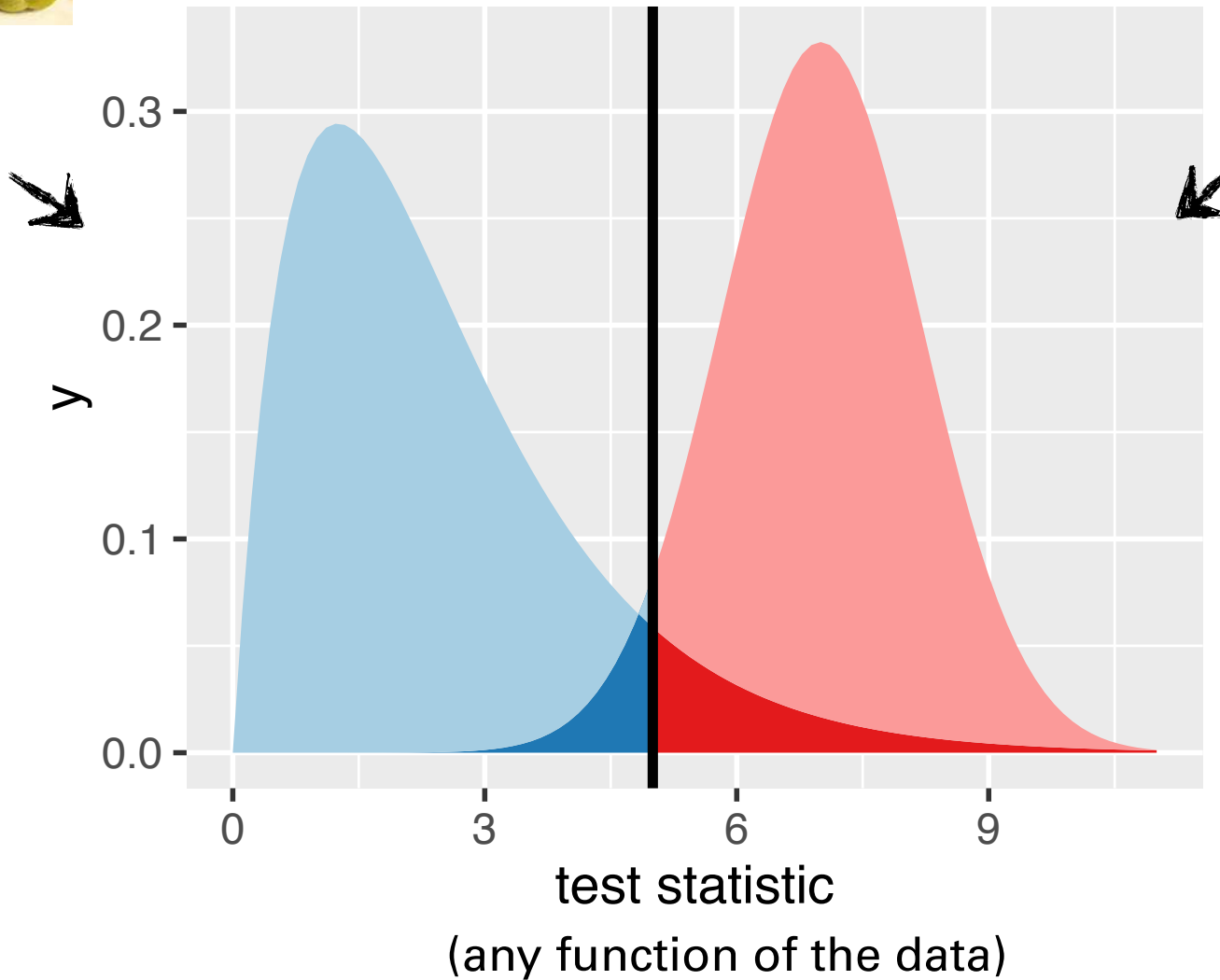
Aims for this Lecture

Understand the basic principles of hypothesis testing, its pitfalls, strengths, use cases and limitations

What changes when we go from single to multiple testing?

False discovery rates, p-value 'adjustments', filtering and weighting

Testing vs Classification



Accuracy vs Precision - Bias vs Variance

← bias

accuracy →

dispersion →



← precision



Karl Popper (1902-1994)



Logical asymmetry between verification and falsifiability.

No number of positive outcomes at the level of experimental testing can confirm a scientific theory, but a single counterexample is logically decisive: it shows the theory is false.

Example



Toss a coin a number of times \Rightarrow

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

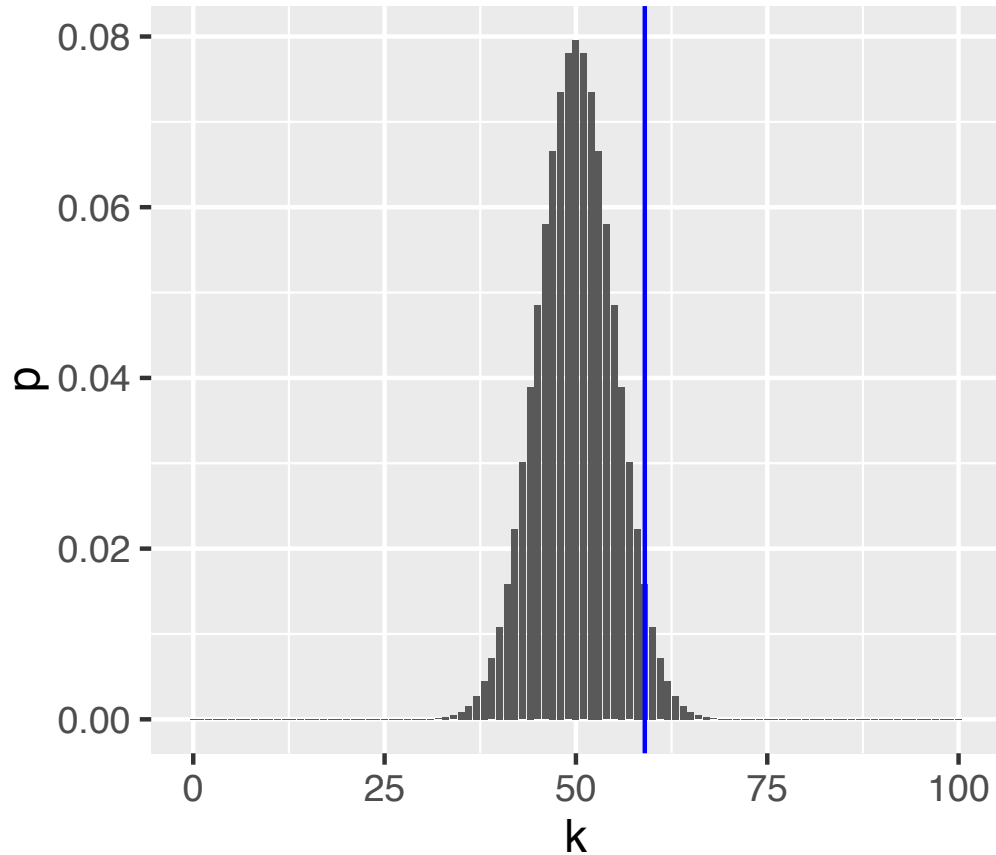


Figure 6.3: The binomial distribution for the parameters $n = 100$ and $p = 0.5$,

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

Rejection Region

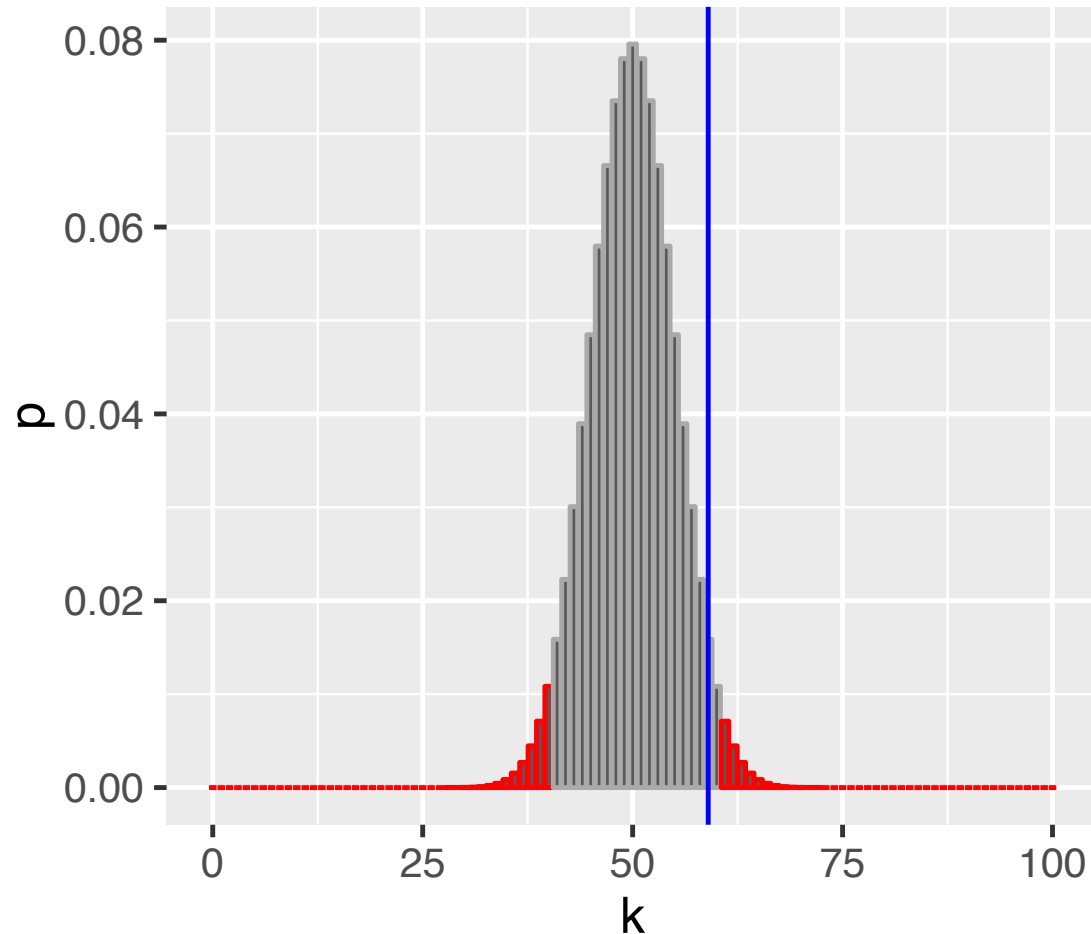


Figure 6.5: As Figure 6.3, with rejection region (red) whose total area is $\alpha = 0.05$.

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 that the coin is not fair 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. the possible outcomes and each of their probabilities.

Decide on the **rejection region**, i.e., a subset of possible outcomes whose total probability is small (\leq **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: the **test statistic**

Set up a **null** model of real world test statistic, probabilities.

Decide on the outcomes with α (\leq **significance**)

This is the idealised scenario, "orthodoxy". Reality, esp. in retrospective 'data-mining' is typically different.

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

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



Examples of Null Hypotheses:

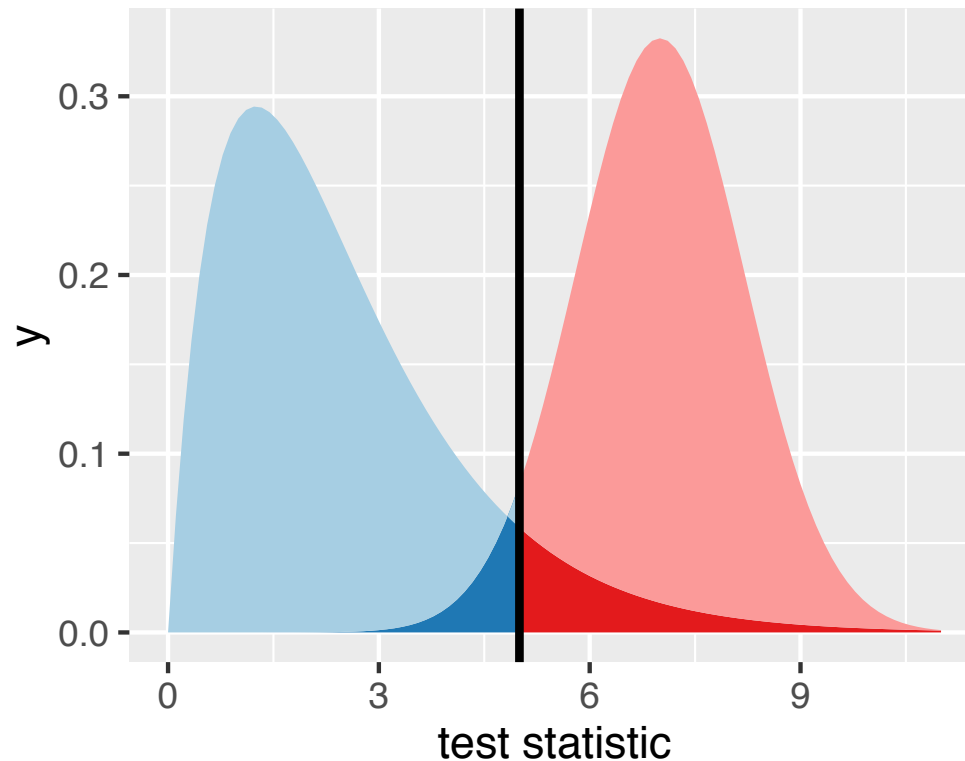
- The coin is fair
- The new drug is no better or worse than a placebo
- The effect of that RNAi-treatment on my cells is no different than that of a negative control treatment

These are not Null Hypotheses:

- The number of heads and tails were the same
- The coin is not fair
- The drug is worth its money

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

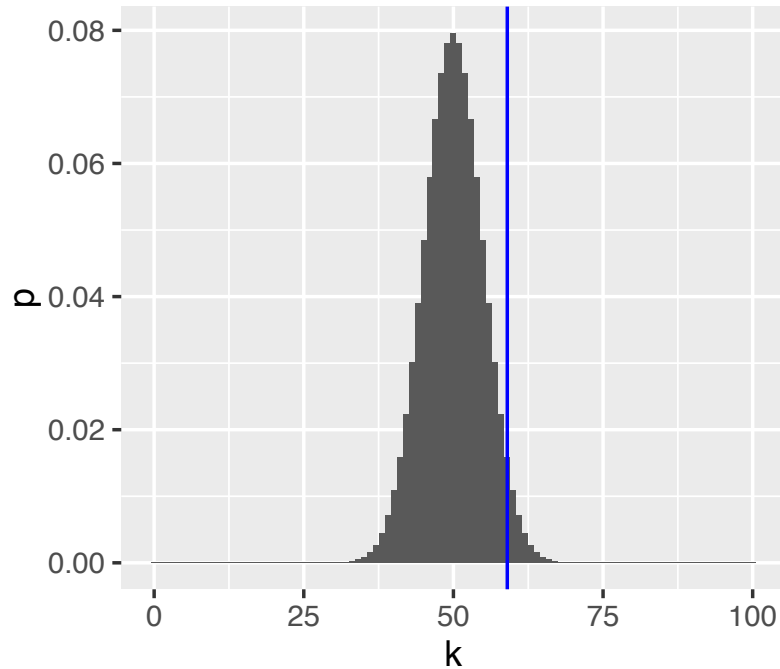


Figure 6.3: The binomial distribution for the parameters $n = 100$ and $p = 0.5$, according to Equation (6.1).

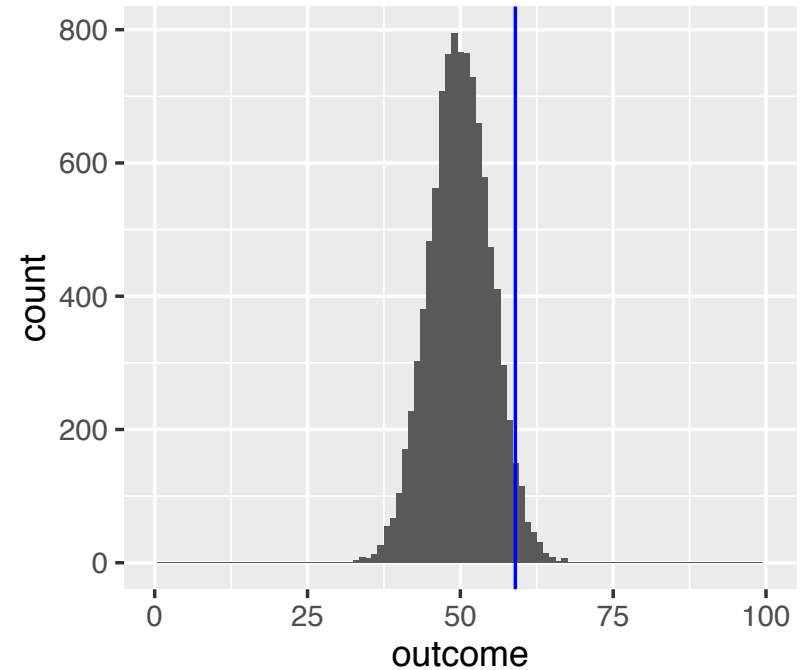


Figure 6.4: An approximation of the binomial distribution from 10^4 simulations (same parameters as Figure 6.3).

Parametric Theory vs Simulation

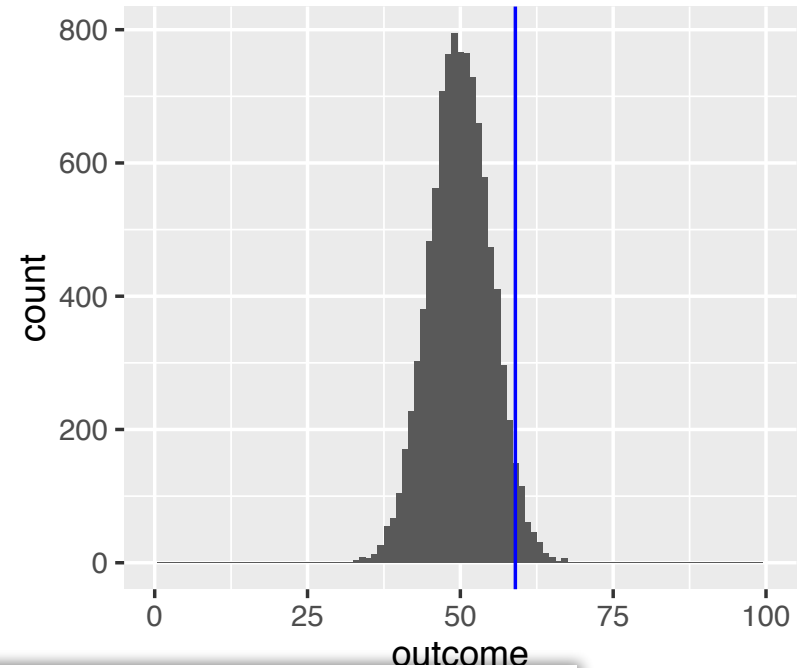
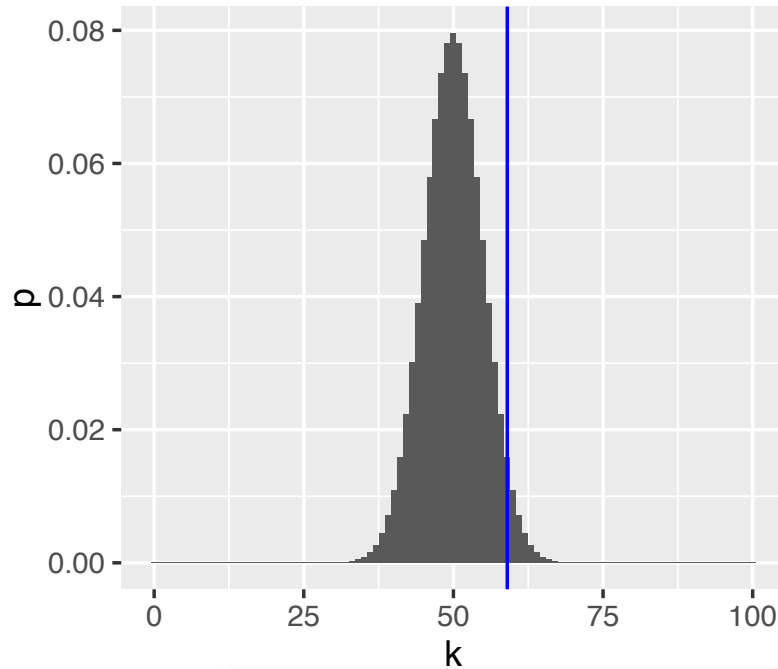


Figure 6.3: The distribution of the parameter k according to

Q:

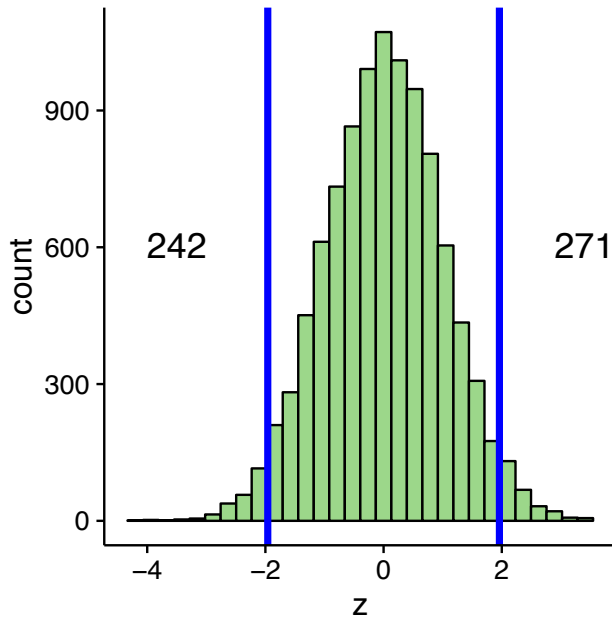
Discuss pros and cons for each

of the simulations).

p-Values as Random Variables

test statistic

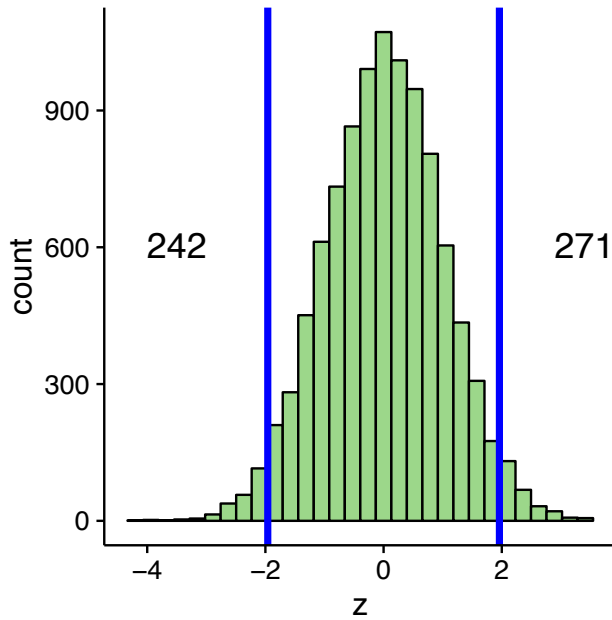
p-value



distribution function

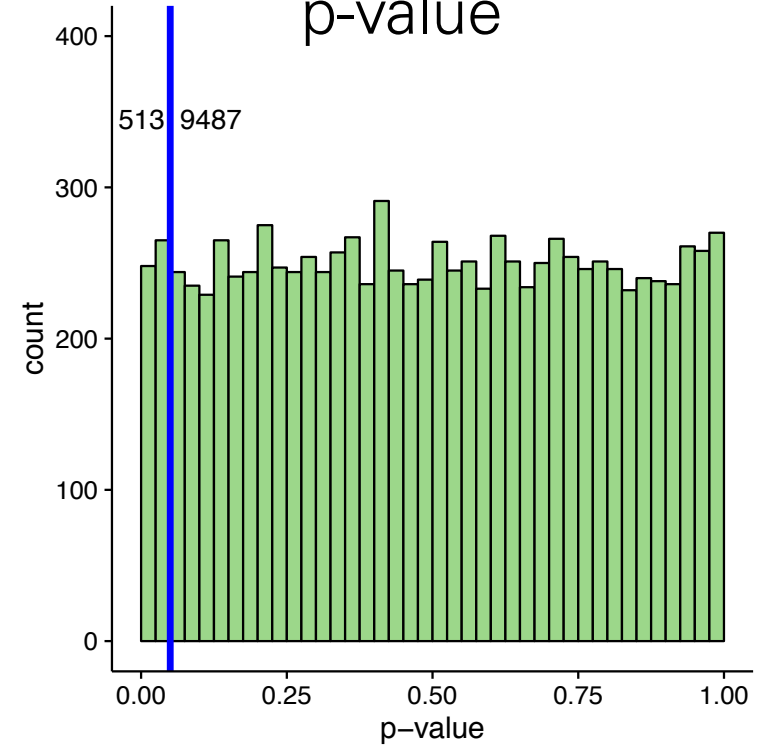
p-Values as Random Variables

test statistic



distribution
function

p-value



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

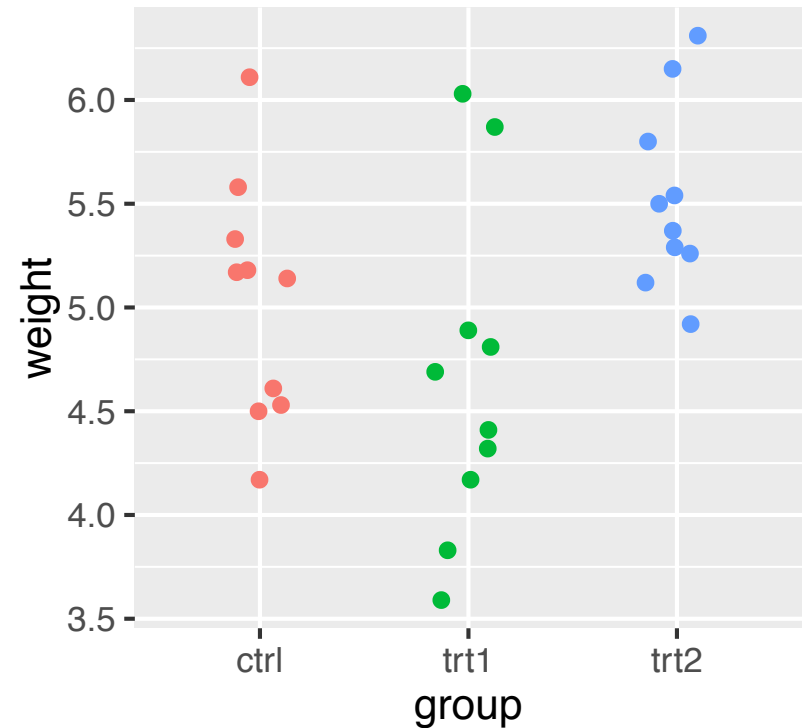


Figure 6.7: The PlantGrowth data.

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

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

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

t- (and |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- (and |t|-) Distribution

If the data are identically normally 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- (and |t|-) Distribution

If the data are identically normally 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?

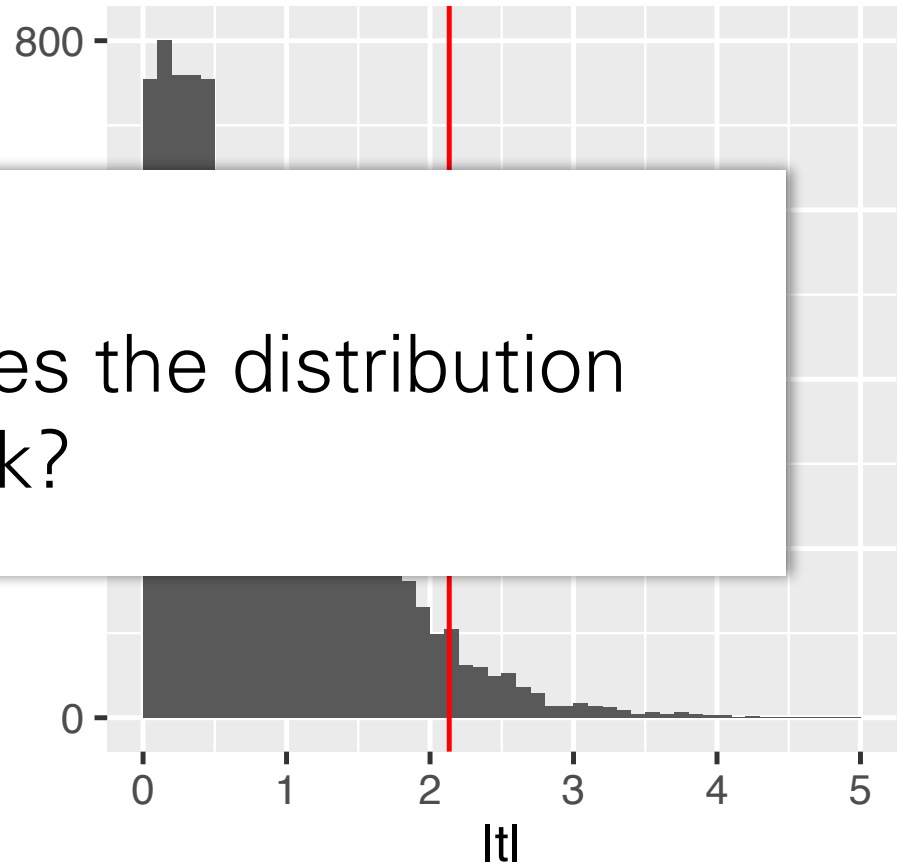


Figure 6.8: The null distribution of the (absolute) t -statistic determined by simulations – namely, by random permutations of the group labels.

Comments and Pitfalls

Comments and Pitfalls

Derivation of the t -distribution assumes that observations are independent and follow a normal distribution: a sufficient, but not necessary, condition

Comments and Pitfalls

Derivation of the t -distribution assumes that observations are independent and follow a normal distribution: a sufficient, but not necessary, condition

Deviation from normality - heavier tails: test still maintains type-I error control, but may no longer have optimal power.

Comments and Pitfalls

Derivation of the t -distribution assumes that observations are independent and follow a normal distribution: a sufficient, but not necessary, condition

Deviation from normality - heavier tails: test still maintains type-I error control, but may no longer have optimal power.

Options: use permutations; transform (e.g. ranks - Wilcoxon test)

Comments and Pitfalls

Derivation of the t -distribution assumes that observations are independent and follow a normal distribution: a sufficient, but not necessary, condition

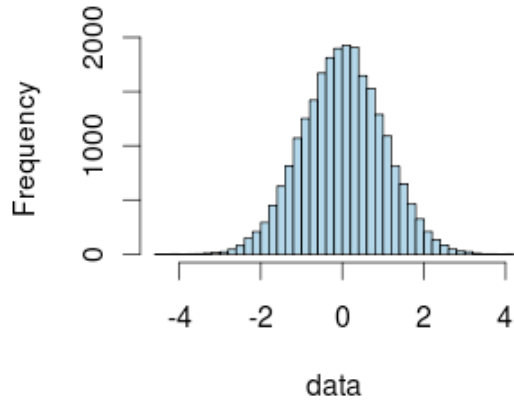
Deviation from normality - heavier tails: test still maintains type-I error control, but may no longer have optimal power.

Options: use permutations; transform (e.g. ranks - Wilcoxon test)

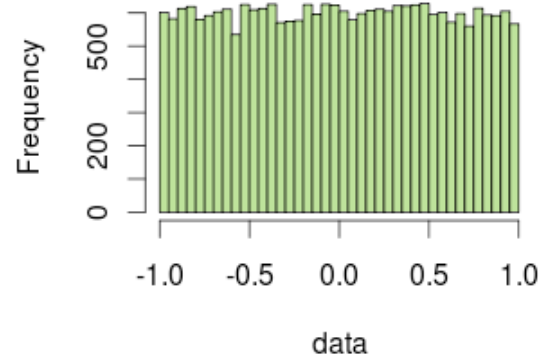
If the data are **dependent**, then p-values will likely be totally wrong (e.g., for positive correlation, too optimistic).

Different Data Distributions – Independent Case

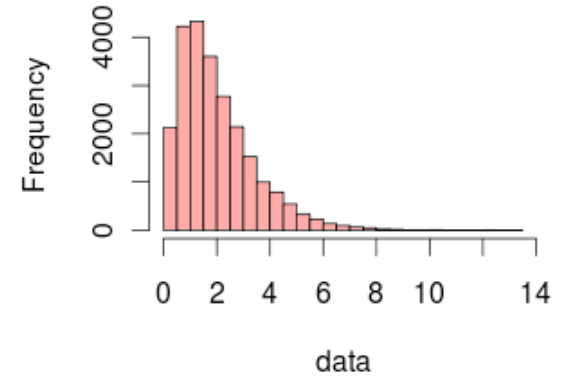
Normal(0,1)



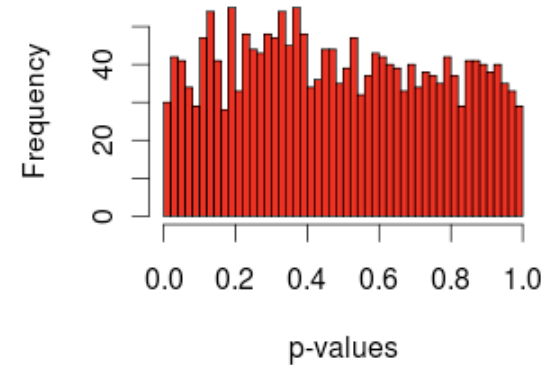
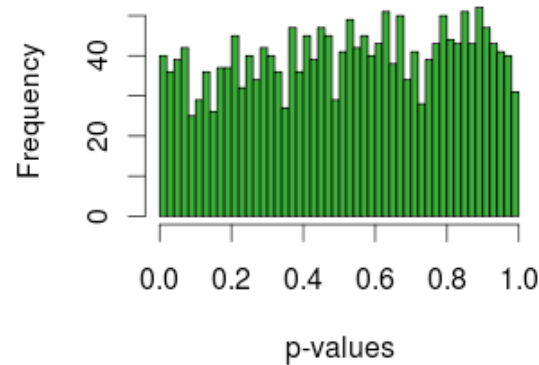
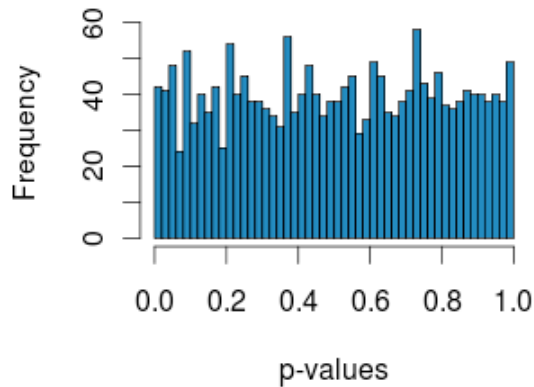
Uniform(-1,1)



Gamma(2, 1)

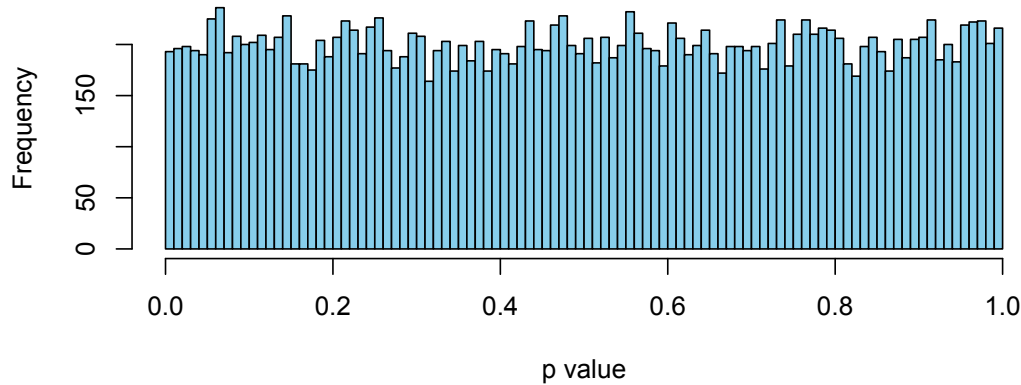


df=10

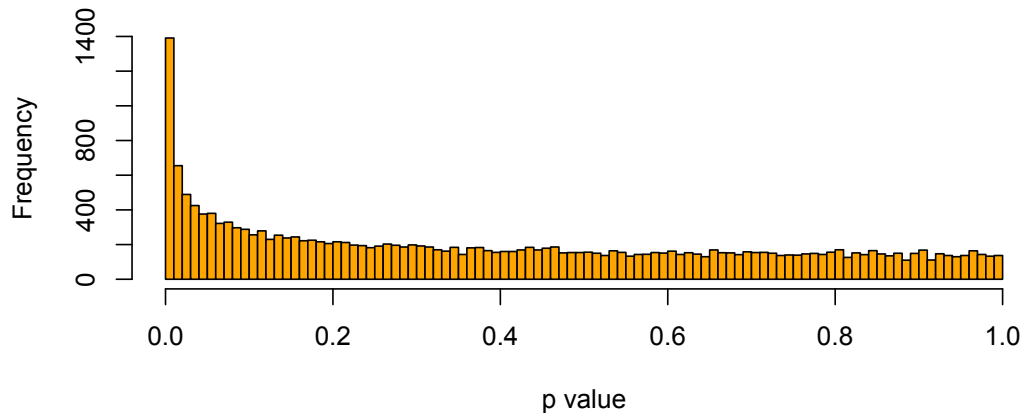


t-Test Loses Error Control if Independence Assumption Does not Hold

uncorrelated



correlated (band-diagonal)



```
library("mvtnorm")
library("genefilter")

p = 30    ## number of samples
n = 20000 ## number of genes

mu = rep(0, p)
dp = diag(p)

sigma = list(
  'uncorrelated' = dp,          ## unity matrix
  'correlated (band-diagonal)' = ## band diagonal
    (row(dp)==col(dp)) + 0.5 * (abs(row(dp)-col(dp))==1))

lapply(sigma, print)

## generate data
x = lapply(sigma, function(s) rmvnorm(n = n, mean = mu, sigma = s))

## tests
tt = lapply(x, rowttests)

par(mfrow=c(length(tt), 1))
for(i in seq(along=tt))
  hist(tt[[i]]$p.value, breaks=100, col=c("skyblue", "orange")[i],
       main=names(tt)[i], xlab="p value")
```

t-Test Loses Error Control if Independence Assumption Does not Hold

uncorrelated

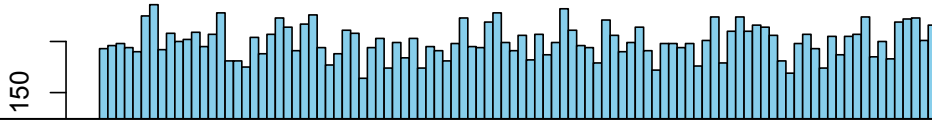
```
library("mvtnorm")  
library("genefilter")  
  
p = 30      ## number of samples  
30000      ## number of genes
```

batch effects!

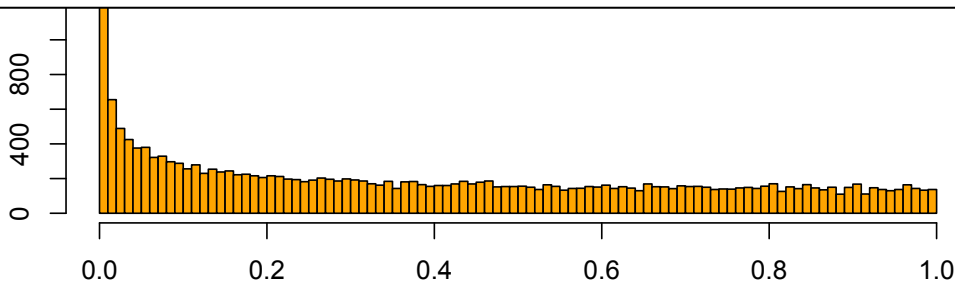
= s))

```
par(mfrow=c(length(tt), 1))  
for(i in seq(along=tt))  
  hist(tt[[i]]$p.value, breaks=100, col=c("skyblue", "orange")[i],  
       main=names(tt)[i], xlab="p value")
```

Frequency



Frequency



p value

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 \neq
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.

P-VALUE

INTERPRETATION

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

What is p-Value Hacking ?

What is p-Value Hacking ?

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

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)

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.
Lazara 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

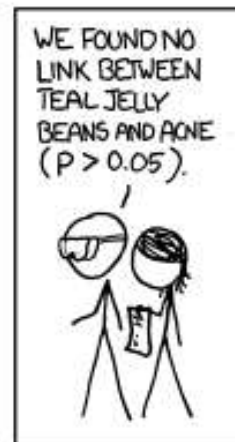
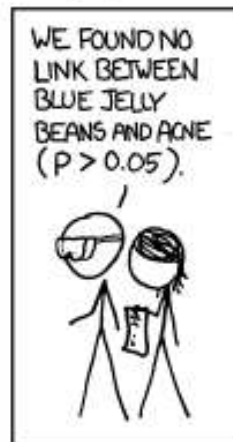
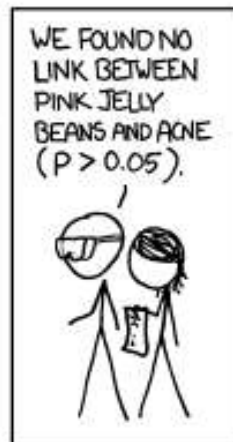
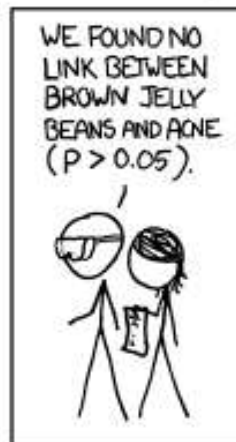
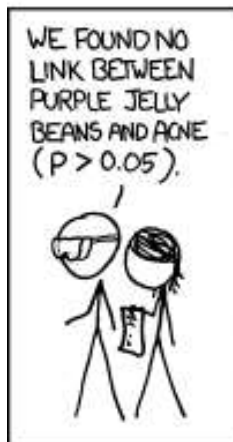
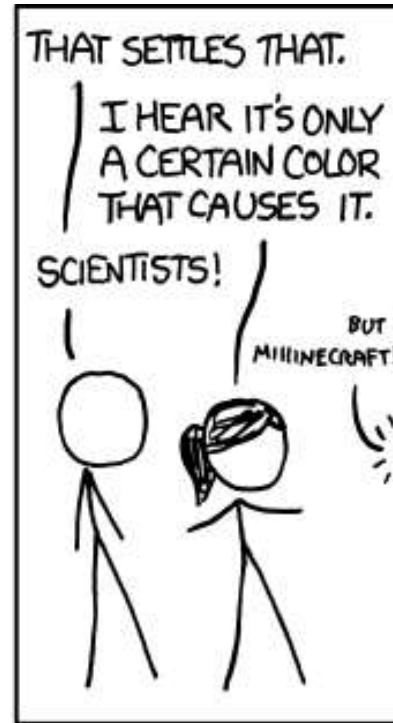
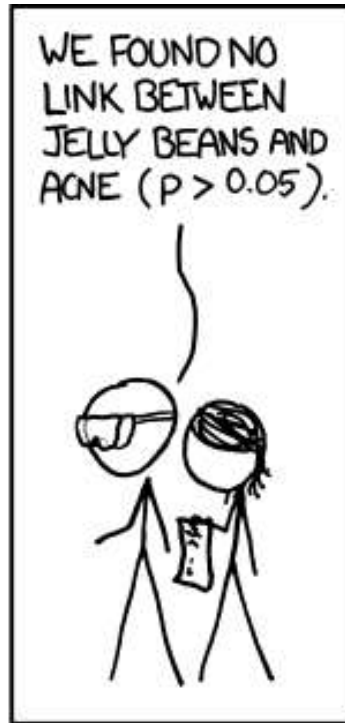
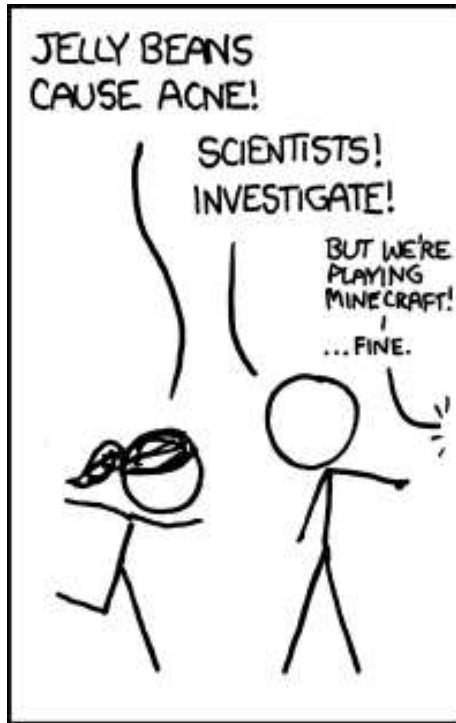
'outlier' removal

the 5% threshold and publication bias

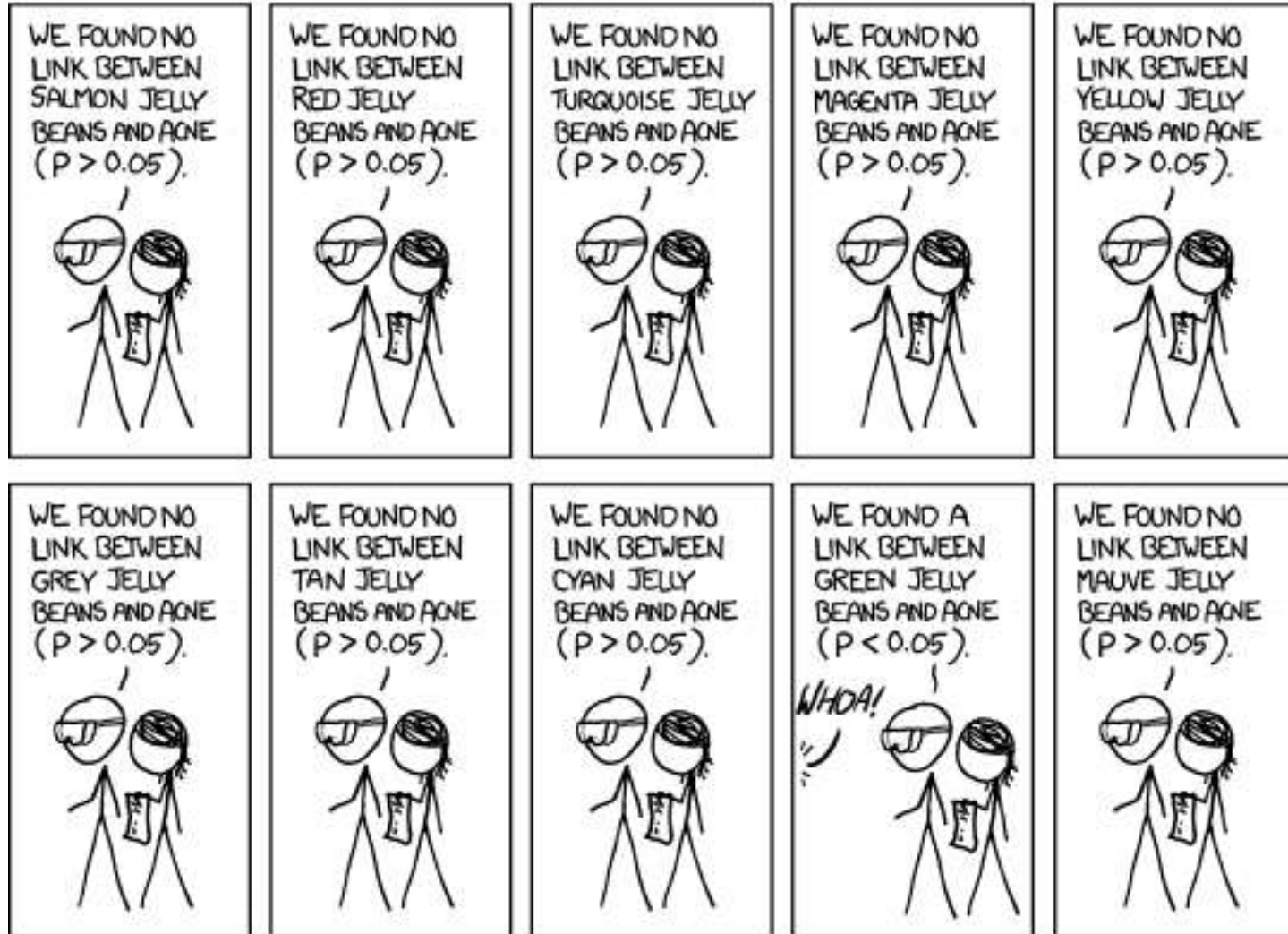
The ASA's Statement on p-Values:
Context, Process, and Purpose
Ronald L. Wasserstein & Nicole A.
Lazara DOI:
10.1080/00031305.2016.1154108

What can we do about this?

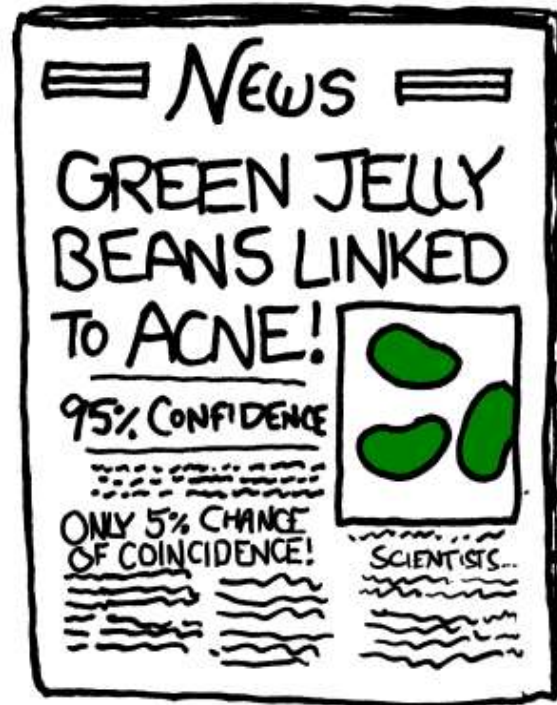
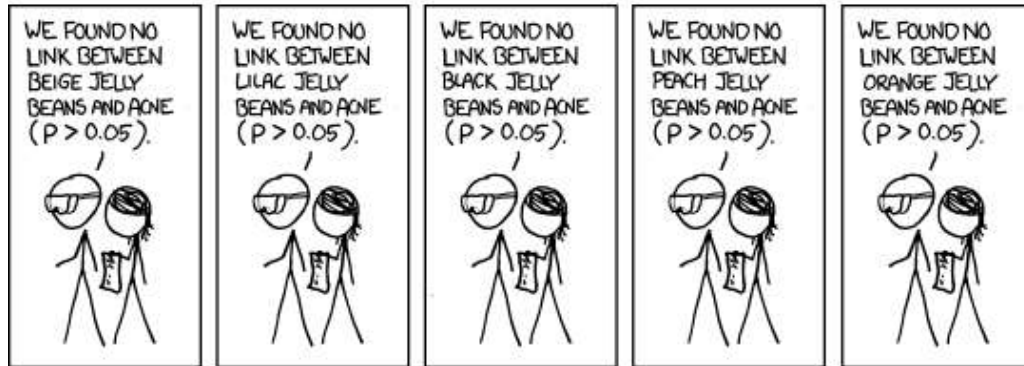
Multiple Testing



Multiple Testing



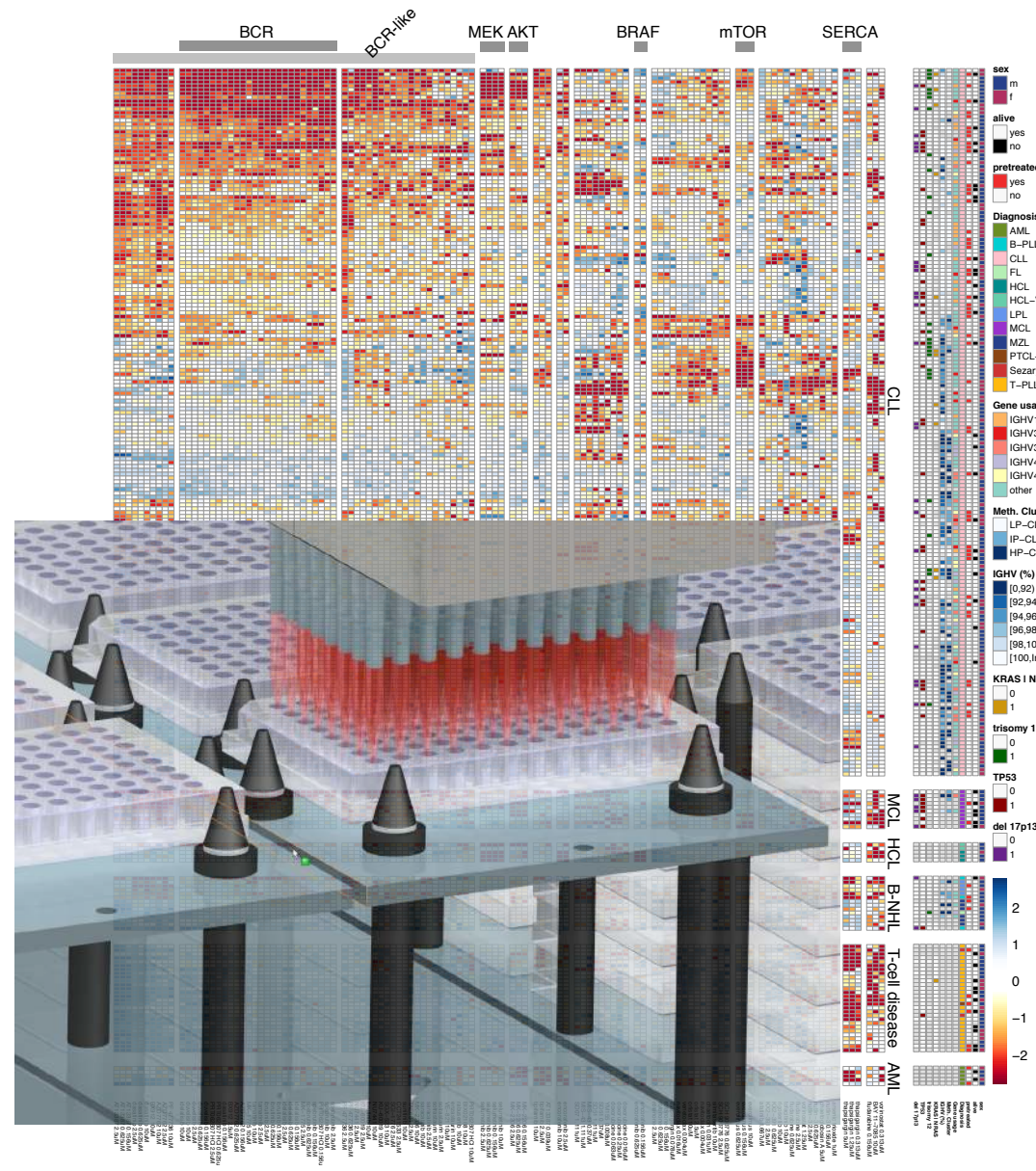
Multiple Testing



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
- ...



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

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



False Positive Rate and False Discovery Rate

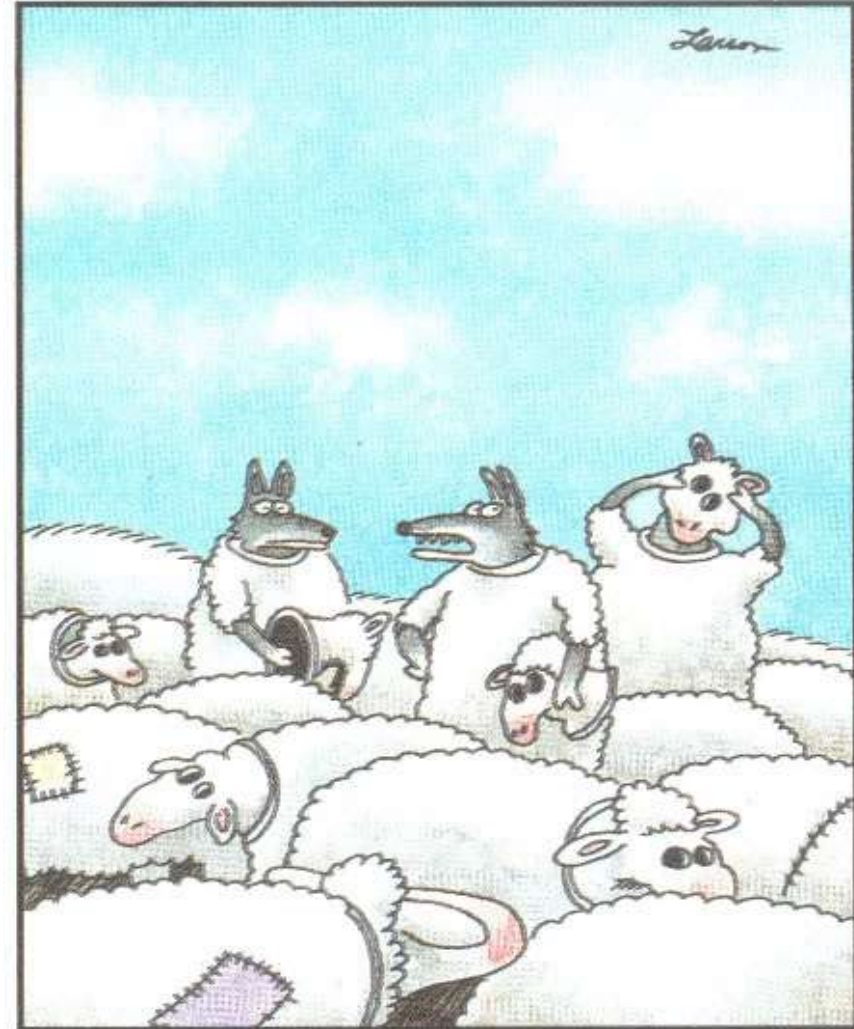
FPR: fraction of FP among all true negatives

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 $FDR = FWER$

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.)

THIS NEUTRINO DETECTOR MEASURES
WHETHER THE SUN HAS GONE NOVA.

THEN, IT ROLLS TWO DICE. IF THEY
BOTH COME UP SIX, IT LIES TO US.
OTHERWISE, IT TELLS THE TRUTH.

LET'S TRY.
DETECTOR! HAS THE
SUN GONE NOVA?



FREQUENTIST 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.



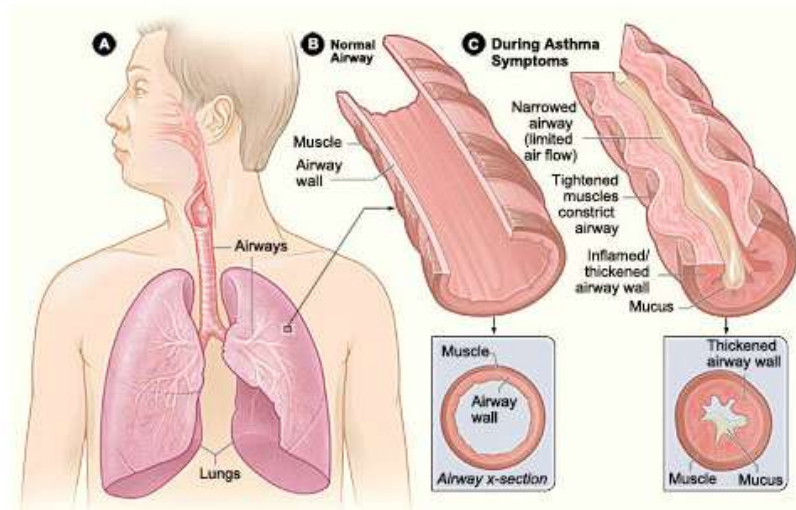
BAYESIAN STATISTICIAN:

BET YOU \$50
IT HASN'T.



Data set 1: RNA-Seq

Transcriptome changes in four samples of primary human airway smooth muscle cells treated with dexamethasone, a synthetic glucocorticoid. 1 μ M 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 :

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

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

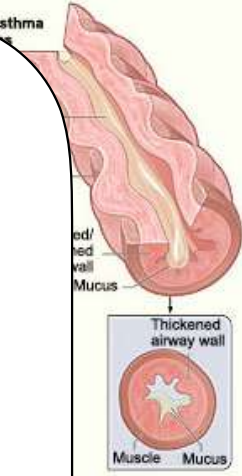
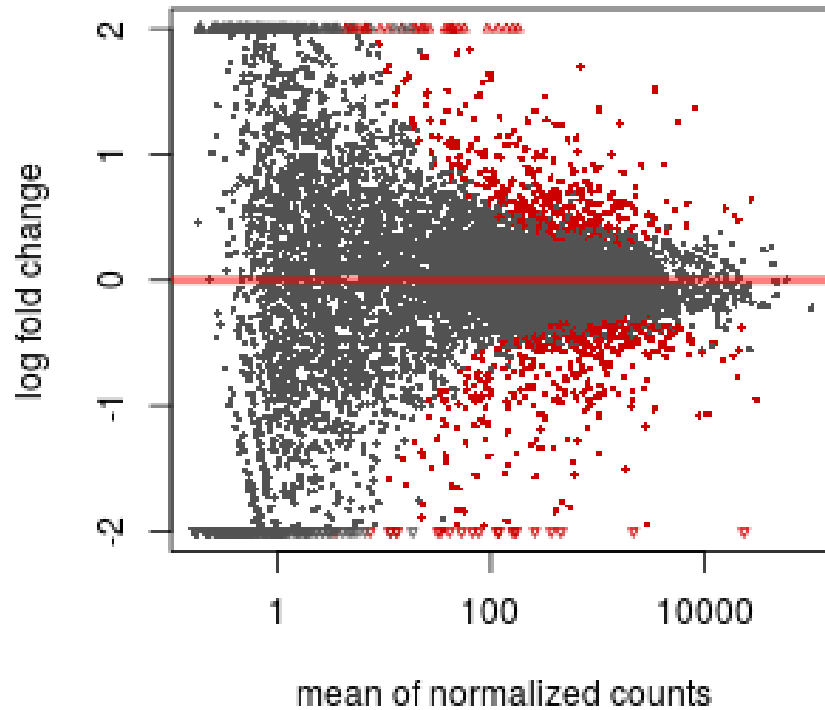
$$\log q_{ij} = \sum_r x_{jr} \beta_{rj}$$

design <- ~ cellline + dexamethasone

Data set 1: RNA-Seq

Transcriptome
samples
smooth muscle
dexamethasone
glucocorticoid

cell
N61
N61
N05
N05
N08
N08
N06101
N061011 trt



analysis:

expression = α_j)

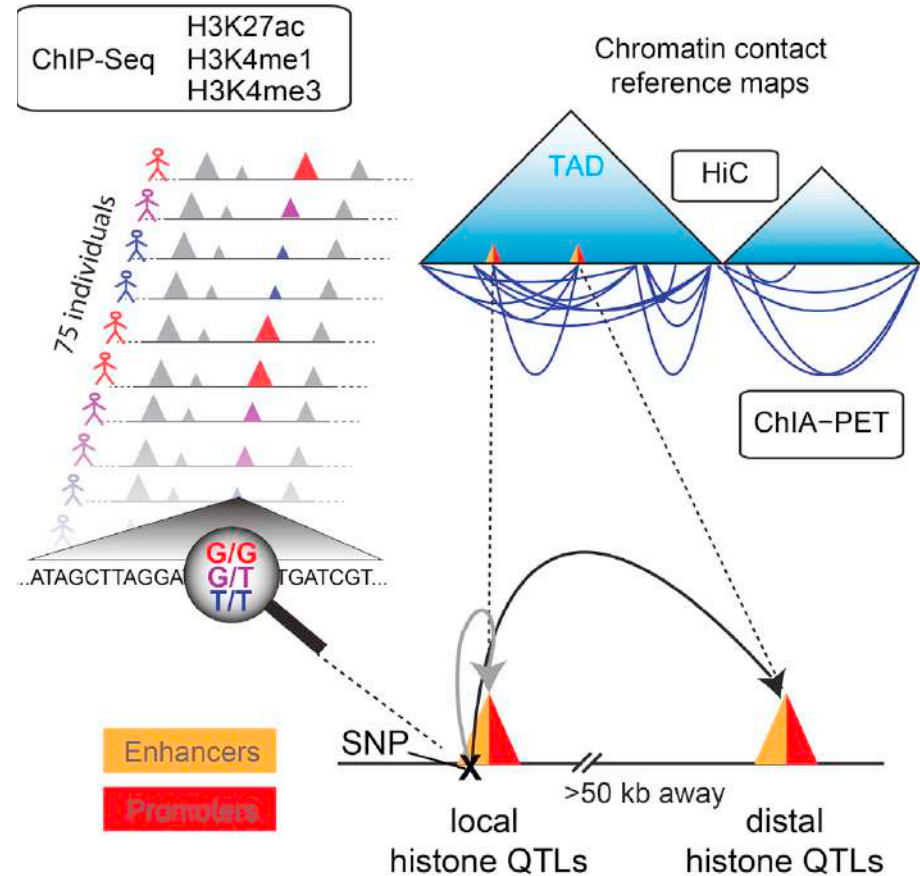
```
design <- ~ cellline + dexamethasone
```

Data set 2: hQTL

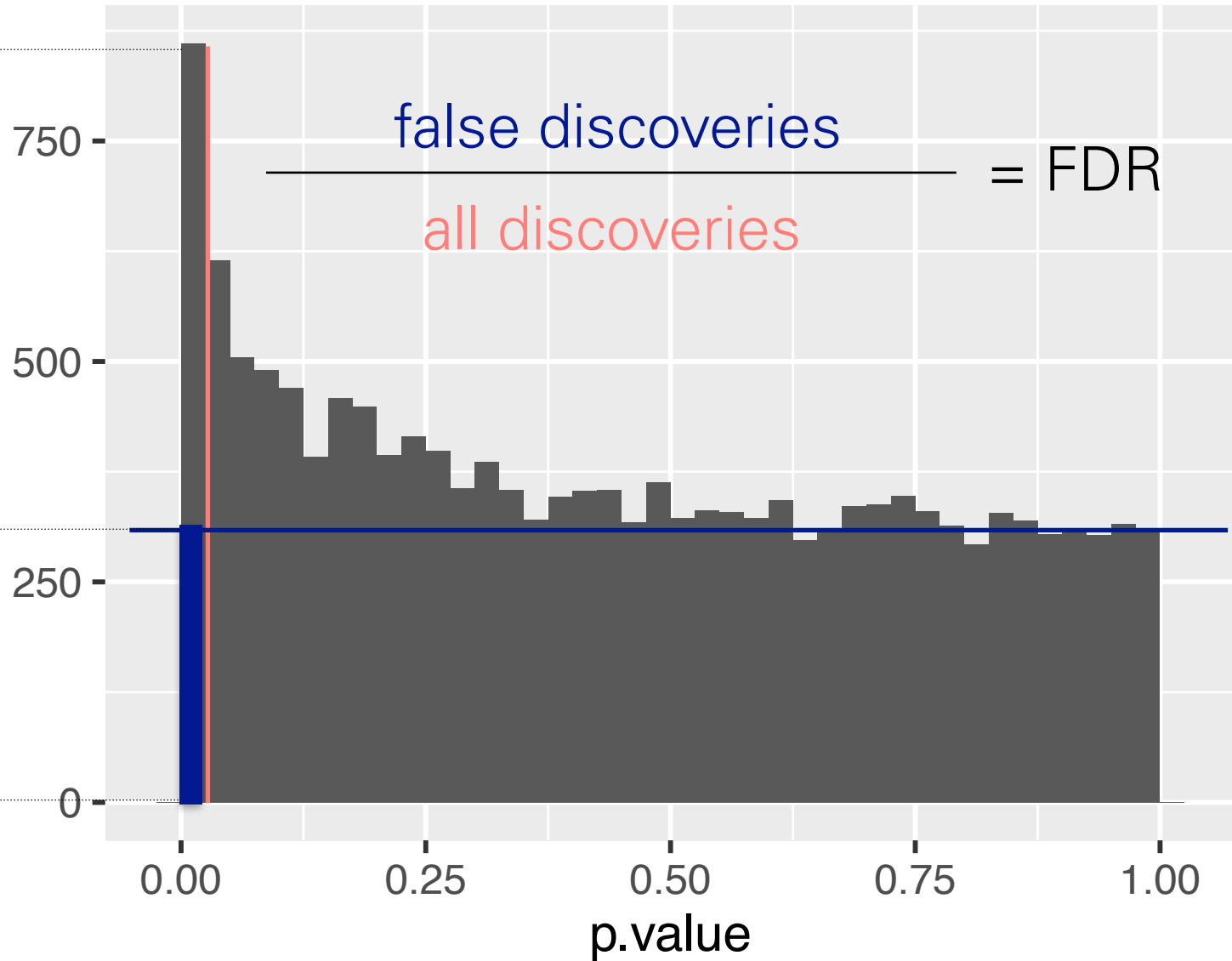
ChIP-seq for histone marks in lymphoblastoid cell lines from 75 sequenced individuals.

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

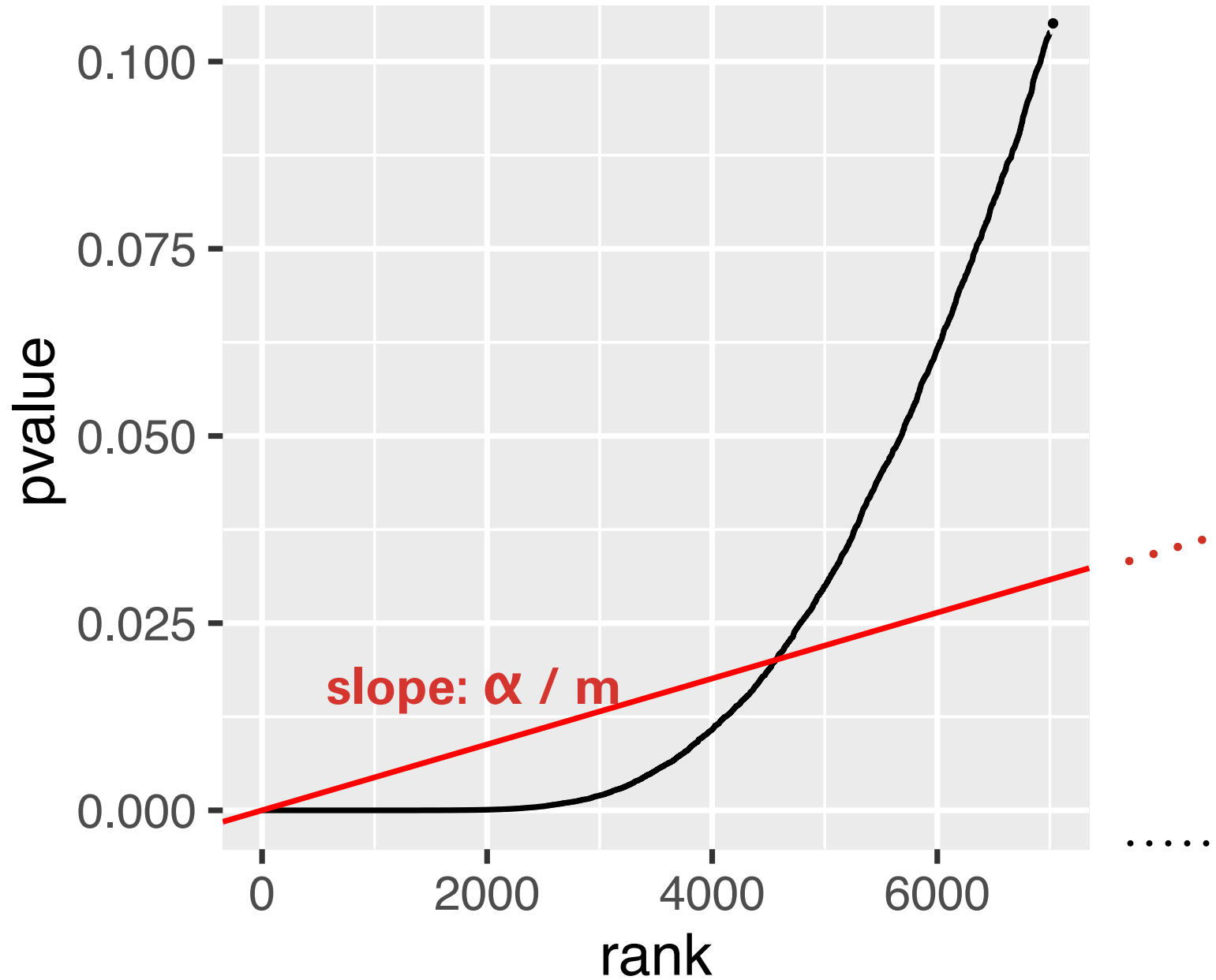


False Discovery Rate

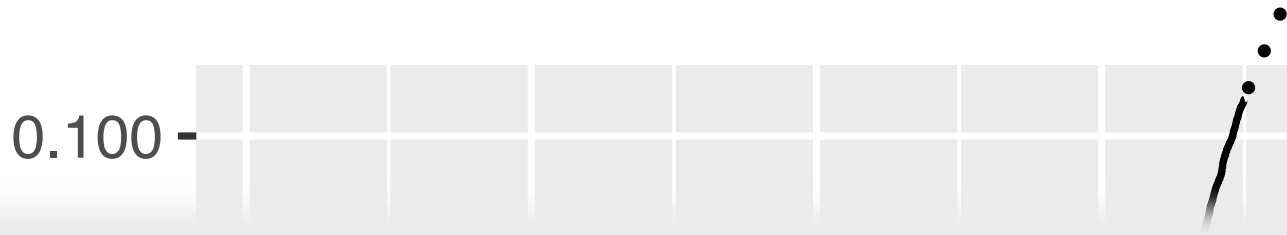


Method of Benjamini & Hochberg (1995)

Method of Benjamini & Hochberg



Method of Benjamini & Hochberg



BH = {

```
i <- length(p) : 1
```

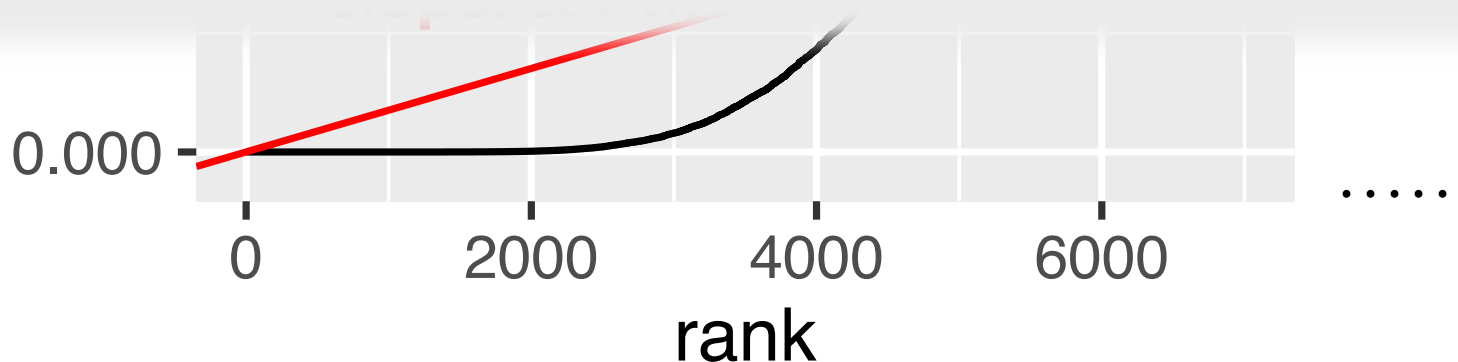
```
o <- order(p, decreasing = TRUE)
```

```
ro <- order(o)
```

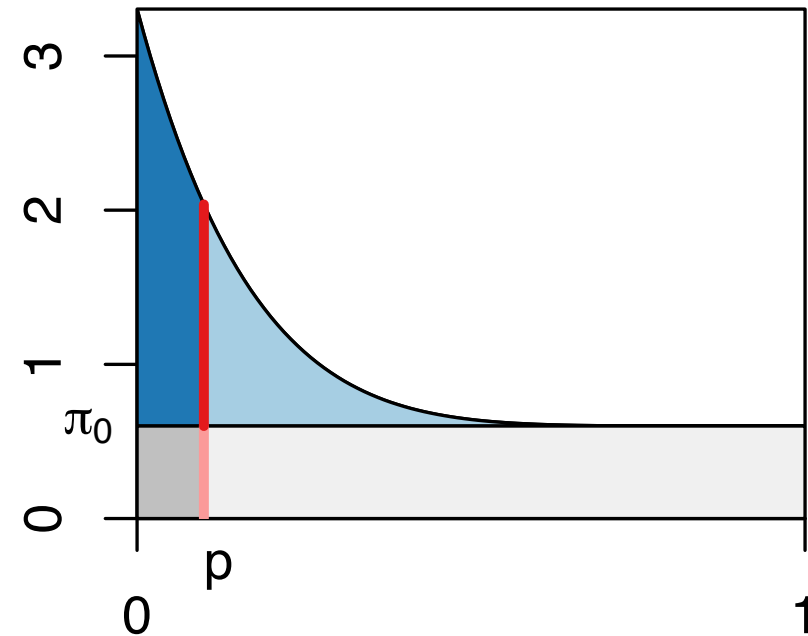
```
pmin(1, cummin(n/i * p[o]))[ro]
```

}

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



The Two-Groups Model and the Local False Discovery Rate



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

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

FDR: a set property. A single number that applies to a whole set of discoveries.

fdr: a local property. It applies to individual hypotheses.

Not all Hypothesis Tests are Created Equal

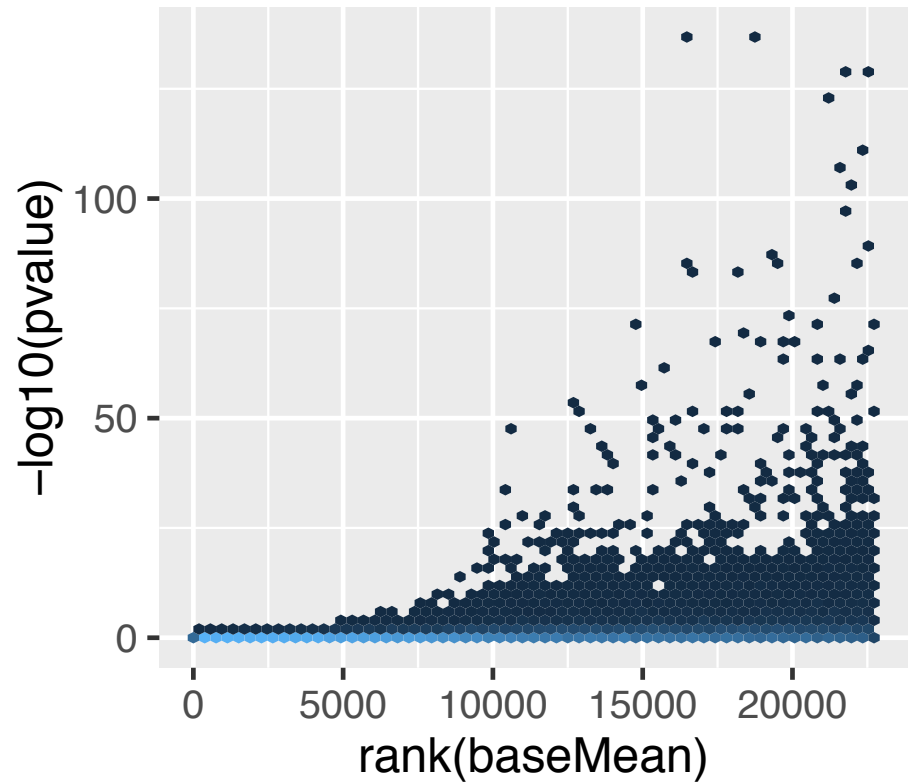
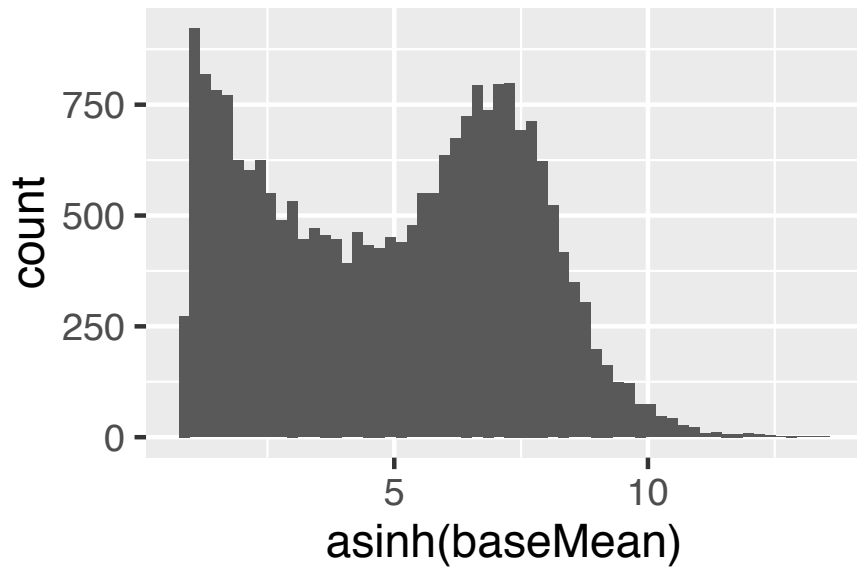


Figure 6.15: Histogram of `baseMean`. We see that it covers a large dynamic range, from close to 0 to around 3.3×10^5 .

Covariates - examples

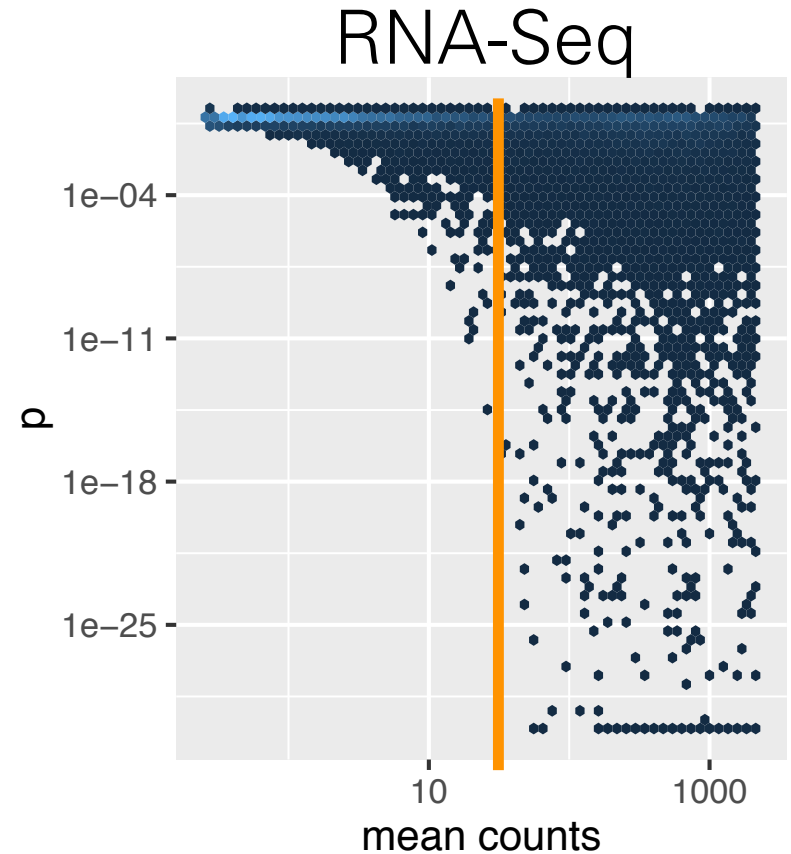
Application	Covariate
Differential RNA-Seq, CHIP-Seq, CLIP-seq, ...	(Normalized) mean of counts for each gene
eQTL analysis	SNP – gene distance
GWAS	Minor allele frequency
<i>t</i> -tests	Overall variance
Two-sided tests	Sign
All applications	Sample size; measures of signal-to-noise ratio

Independent Filtering

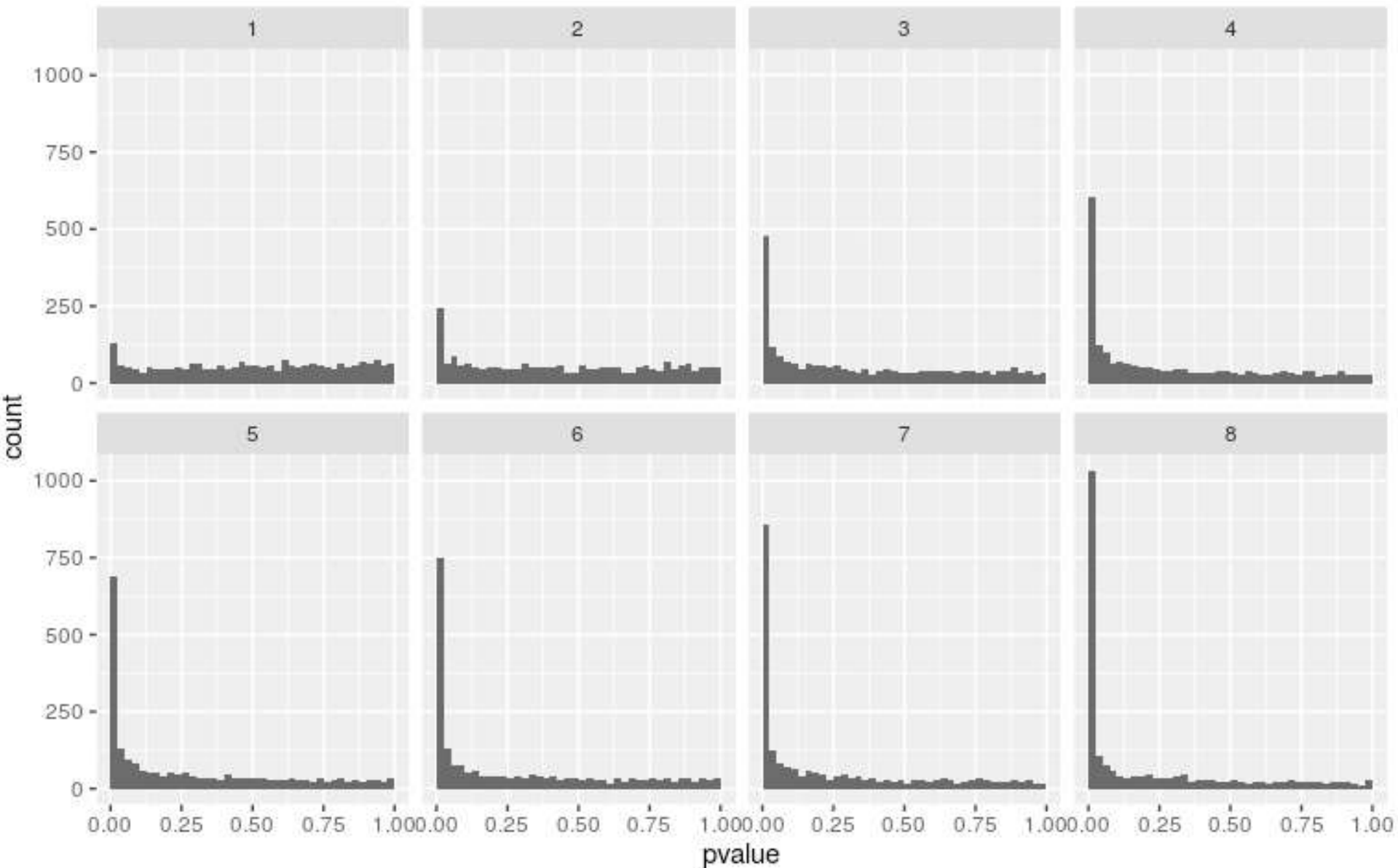
Two steps:

- 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 instead of P .
- Proven Type I error, Wasserma
- If $w_i > 1$,
- $Q_i \leq t \Leftrightarrow$



Weighted Benjamini-Hochberg method

- Let $w_i \geq 0$ and $\frac{1}{m} \sum_{i=1}^m w_i = 1$ ("budget").
- Define $Q_i = P_i / w_i$.
- Apply BH to Q_i .
- Proven to control the FDR at level α .
- Wasserman (2004) shows that $Q_i \leq t \Leftrightarrow P_i \leq t w_i$.
- If $w_i > 1$, $Q_i \leq t \Leftrightarrow P_i \leq t w_i$.
- $Q_i \leq t \Leftrightarrow P_i \leq t w_i$.

Problem: how to know the weights?

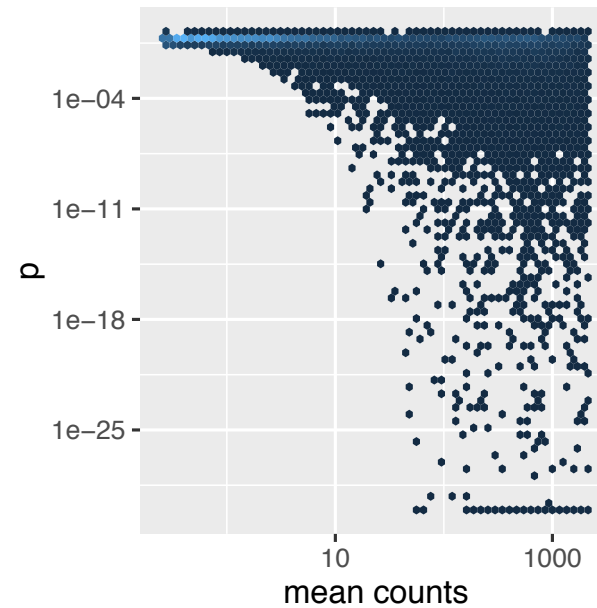
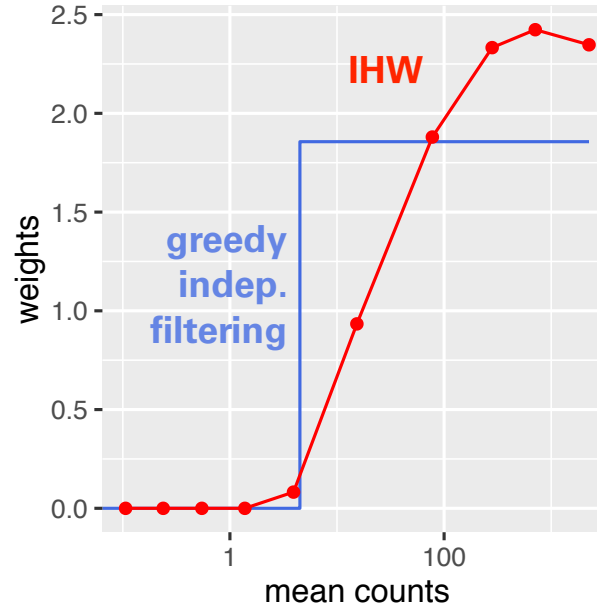
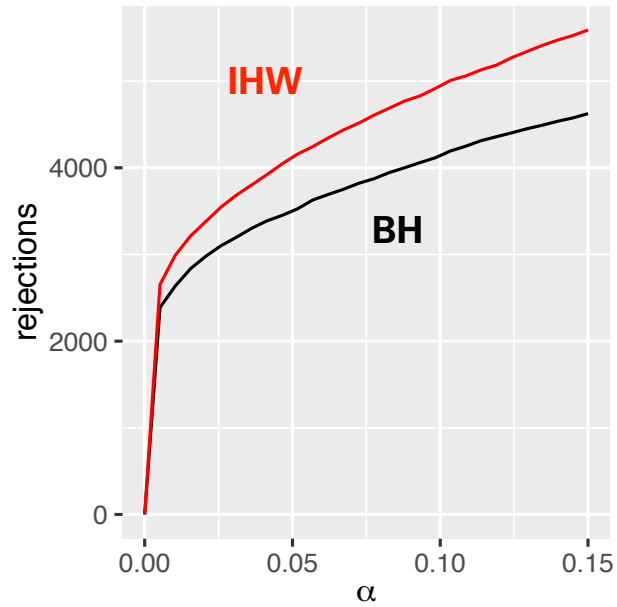


der,

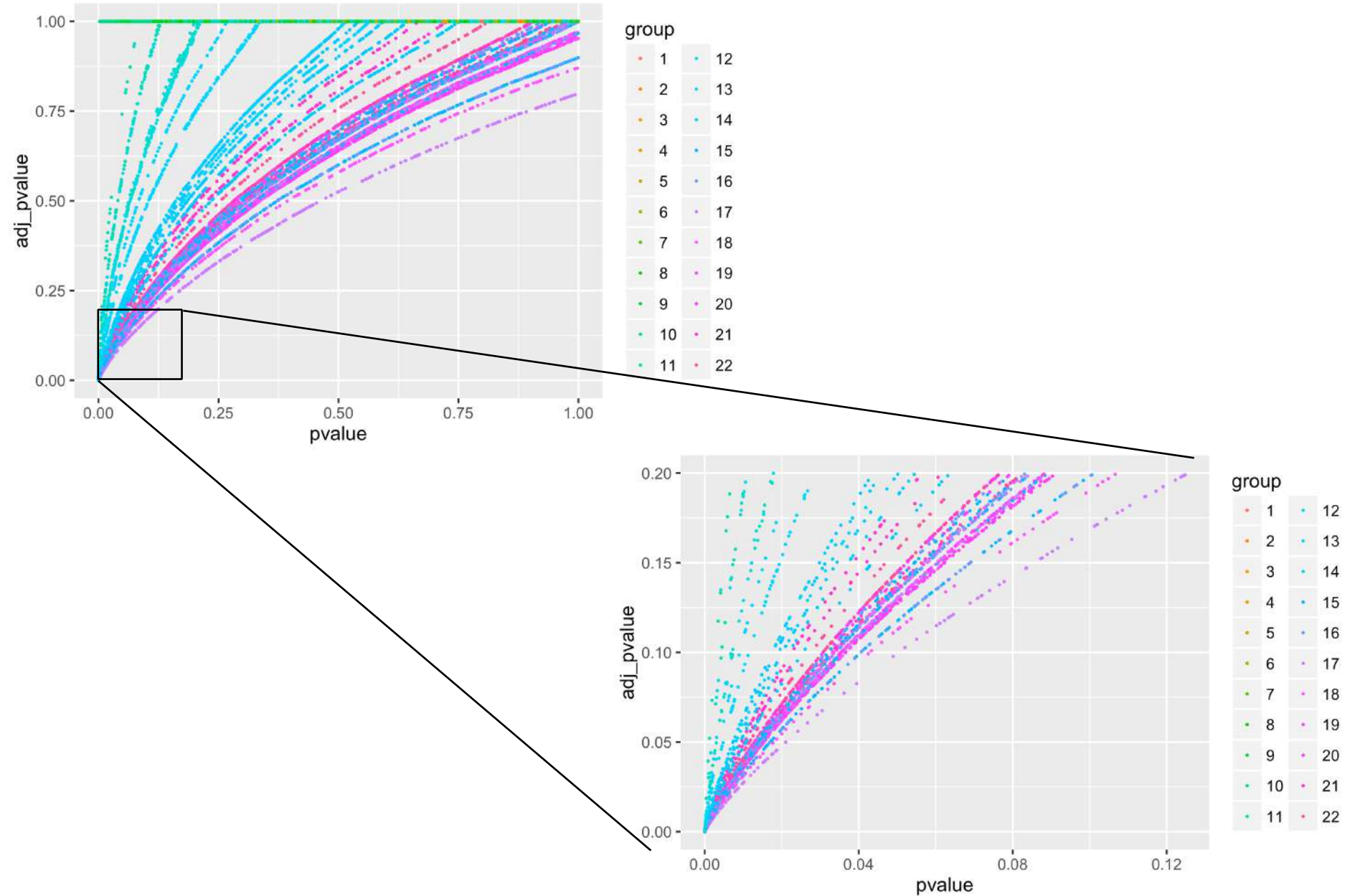
Independent hypothesis weighting (IHW): basic idea

- 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}^* .

RNA-Seq example (DESeq2)



Ranking is not monotonous in raw p-values



Avoiding overfitting



Nikos Ignatiadis

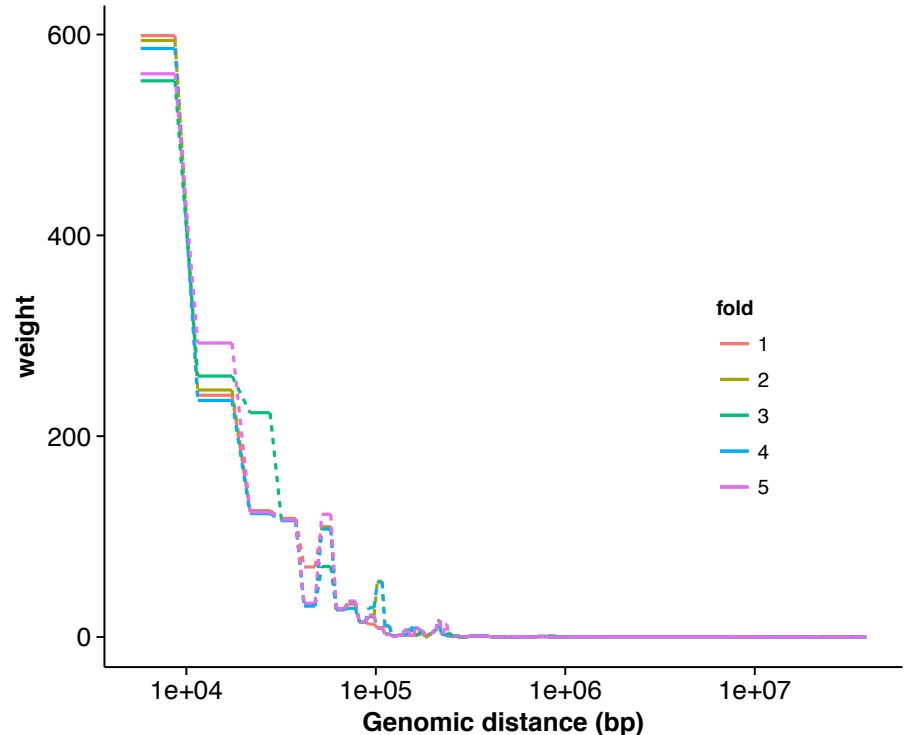
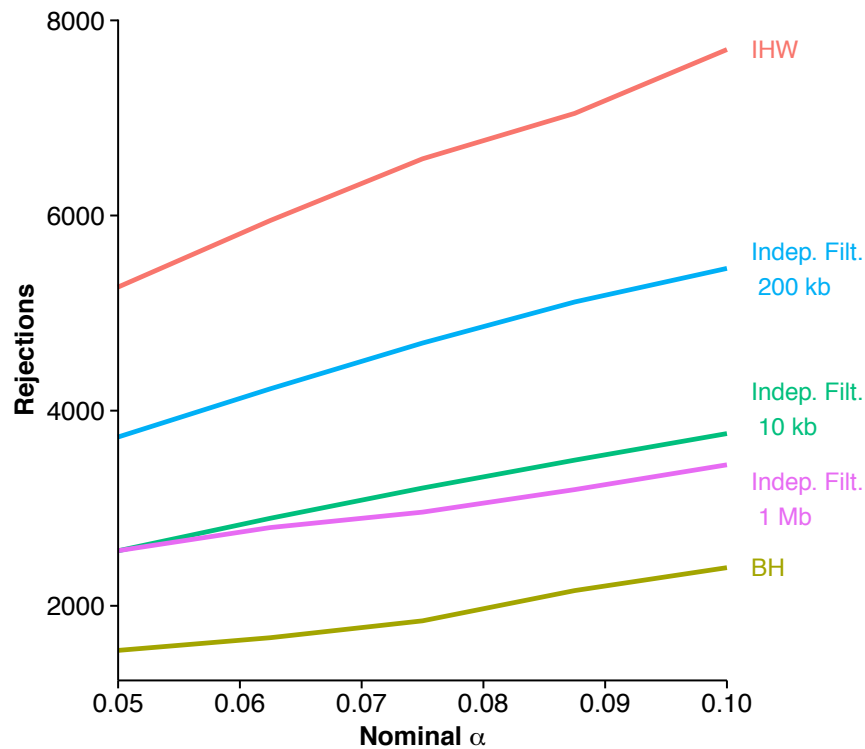
Hypothesis splitting: randomly split hypotheses into k folds. Learn weights for the hypotheses in a fold from the other $k-1$ folds

Regularisation:

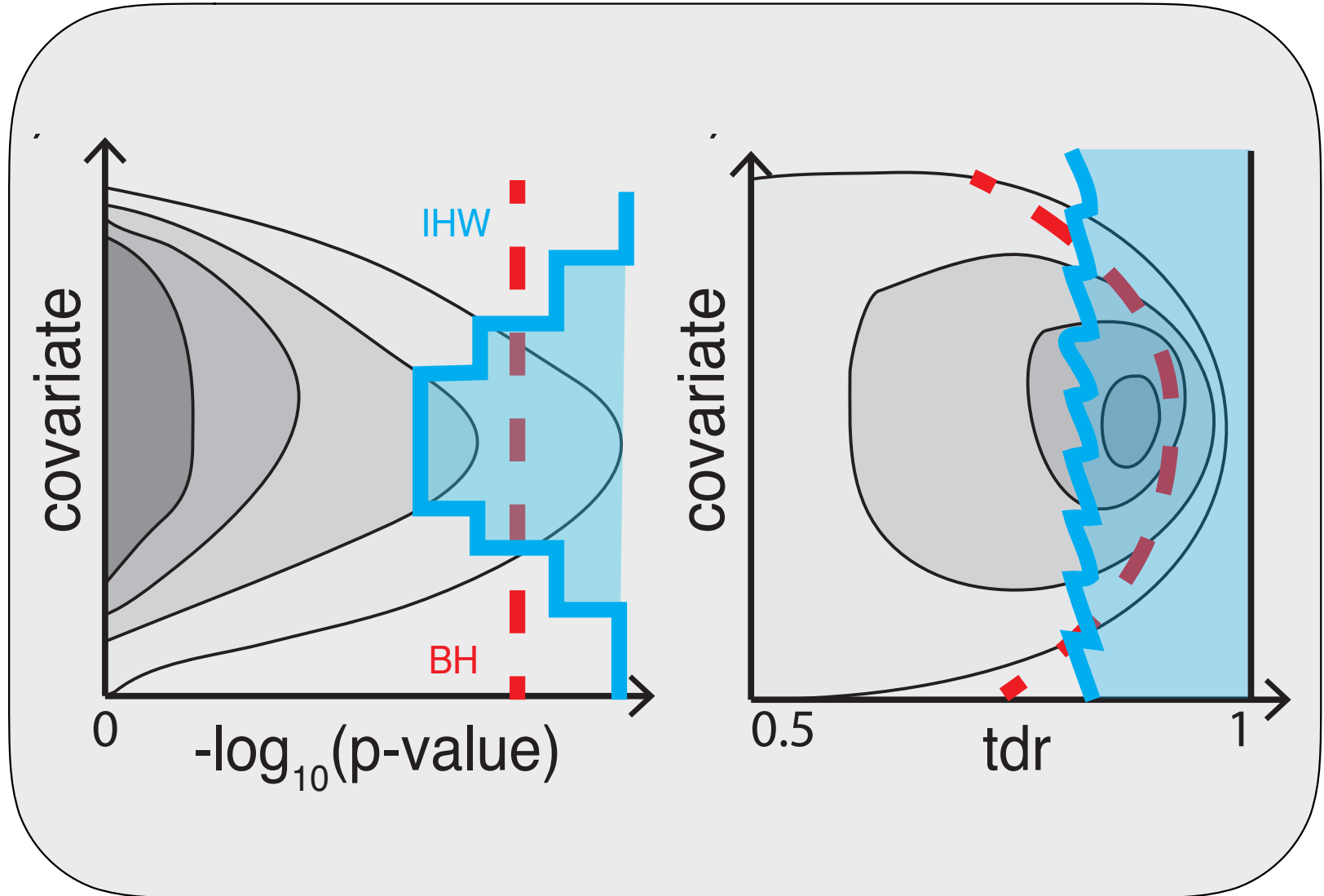
- for ordered covariate: $\sum_g |w_g - w_{g-1}| \leq \lambda$
- for categorical covariate: $\sum_g |w_g - 1| \leq \lambda$

Convex relaxation: for weight optimisation (only), replace ECDFs of the p-values with Grenander estimators (least concave majorant of the ECDF)

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'



Nikos Ignatiadis

A promotional image for the James Bond film 'Casino Royale'. James Bond, played by Daniel Craig, is shown from the waist up, wearing a grey tweed suit, a white shirt, and a dark tie. He is looking slightly to the right with a serious expression. The background is a mix of a blue sky with light clouds and a bright orange and yellow fire or explosion at the bottom. The text 'The p-value Is Not Enough' is overlaid in a white, italicized font, with 'p-value' underlined. Below it, '007' is written in a large, stylized, orange and yellow font.

The p-value Is Not Enough
007

Availability

Bioconductor - IHW (devel... x) +

https://www.bioconductor.org/packages/devel/t Search

Bioconductor
OPEN SOURCE SOFTWARE FOR BIOINFORMATICS

Home Install Help Developers About

Search:

Home » Bioconductor 3.3 » Software Packages » IHW (development version)

IHW

platforms **all** downloads **available** posts **0** in Bioc **devel only**

build **ok** commits **0.17** test coverage **unknown**

This is the **development** version of IHW; to use it, please install the [devel version](#) of Bioconductor.

Independent Hypothesis Weighting

Bioconductor version: Development (3.3)

Independent hypothesis weighting (IHW) is a multiple testing procedure that increases power compared to the method of Benjamini and Hochberg by assigning data-driven weights to each hypothesis. The input to IHW is a two-column table of p-values and covariates. The covariate can be any continuous-valued or categorical variable that is thought to be informative on the statistical properties of each hypothesis test, while it is independent of the p-value under the null hypothesis.

Author: Nikos Ignatiadis [aut, cre]

Maintainer: Nikos Ignatiadis <nikos.ignatiadis01 at gmail.com>

Citation (from within R, enter `citation("IHW")`):

Ignatiadis N, Klaus B, Zaugg J and Huber W (2015). "Data-driven hypothesis weighting increases detection power in big data analytics." *bioRxiv*.

Installation

To install this package, start R and enter:

Documentation »

Bioconductor

- Package [vignettes](#) and manuals.
- [Workflows](#) for learning and use.
- [Course and conference](#) material.
- [Videos](#).
- Community [resources](#) and [tutorials](#).

R / [CRAN](#) packages and [documentation](#)

Support »

Please read the [posting guide](#). Post questions about Bioconductor to one of the following locations:

- [Support site](#) - for questions about Bioconductor packages
- [Bioc-devel](#) mailing list - for package developers

Papers:
Nature Methods, June 2016
arXiv, January 2017