

## 14. hypothesis testing

Programmers using the Eclipse IDE make fewer errors

- (a) Hooey. Errors happen, IDE or not.
- (b) Yes. On average, programmers using Eclipse produce code with fewer errors per thousand lines of code

Black Tie Linux has way better web-server throughput than Red Shirt.

- (a) Ha! Linux is linux, throughput will be the same
- (b) Yes. On average, Black Tie response time is 20% faster.


This coin is biased!

- (a) “Don’t be paranoid, dude. It’s a fair coin, like any other,  $P(\text{Heads}) = 1/2$ ”
- (b) “Wake up, smell coffee:  $P(\text{Heads}) = 2/3$ , totally!”



(a) lbsoff.com sells diet pills. 10 volunteers used them for a month, reporting the net weight changes of:

```
x <- c(-1.5, 0, .1, -0.5, -.25, 0.3, .1, .05, .15, .05)
> mean(x)
[1] -0.15
```



lbsoff proudly announces “Diet Pill Miracle! See data!”

(b) Dr. Gupta says “Bunk!”

Does smoking cause\* lung cancer?

- (a) No; we don't know what causes cancer, but smokers are no more likely to get it than non-smokers
- (b) Yes; a much greater % of smokers get it

\*Notes: (1) even in case (b), “cause” is a stretch, but for simplicity, “causes” and “correlates with” will be loosely interchangeable today. (2) we really don't know, in mechanistic detail, what causes lung cancer, nor how smoking contributes, but the *statistical* evidence strongly points to smoking as a key factor.

Our question: How to do the statistics?

How do we decide?

*Design* an experiment, gather *data*, *evaluate*:

In a sample of N smokers + non-smokers, does % with cancer differ? Age at onset? Severity?

In N programs, some written using IDE, some not, do error rates differ?

Measure response times to N individual web transactions on both.

In N flips, does putatively biased coin show an unusual excess of heads? More runs? Longer runs?

A complex, multi-faceted problem. Here, emphasize evaluation:  
What N? How large of a difference is convincing?

### General framework:

1. Data
2.  $H_0$  – the “null hypothesis”
3.  $H_1$  – the “alternate hypothesis”
4. A decision rule for choosing between  $H_0/H_1$  based on data
5. Analysis: What is the probability that we get the right answer?

### Example:

100 coin flips

$$P(H) = 1/2$$

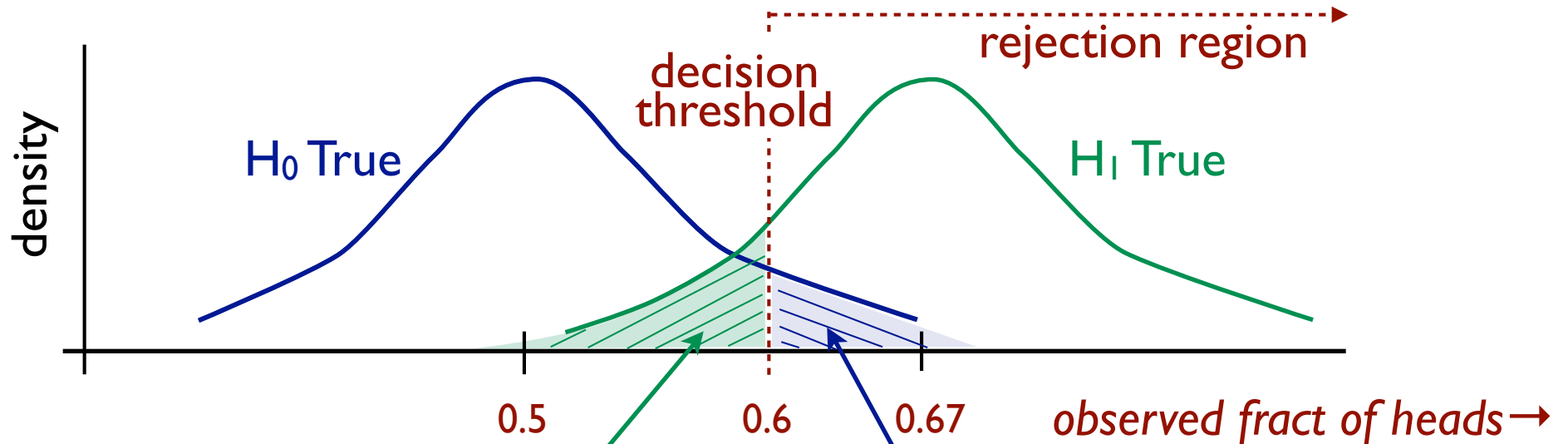
$$P(H) = 2/3$$

“if  $\#H \leq 60$ , accept null, else reject null”

$$P(H \leq 60 \mid 1/2) = ?$$

$$P(H > 60 \mid 2/3) = ?$$

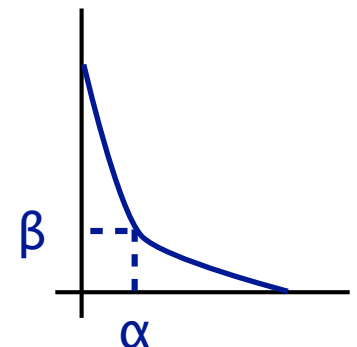
By convention, the null hypothesis is usually the “simpler” hypothesis, or “prevailing wisdom.” E.g., Occam’s Razor says you should prefer that, unless there is *strong* evidence to the contrary.



Type II error: false accept;  
accept  $H_0$  when it is false.  
 $\beta = P(\text{type II error})$

Type I error: false reject;  
reject  $H_0$  when it is true.  
 $\alpha = P(\text{type I error})$

Goal: make both  $\alpha$ ,  $\beta$  small (but it's a tradeoff; they are interdependent).  
 $\alpha \leq 0.05$  common in scientific literature.



Is coin fair (1/2) or biased (2/3)? How to decide? Ideas:

1. Count: Flip 100 times; if number of heads observed is  $\leq 60$ , accept  $H_0$   
or  $\leq 59$ , or  $\leq 61$  ...  $\Rightarrow$  different error rates
2. Runs: Flip 100 times. Did I see a longer run of heads or of tails?
3. Runs: Flip until I see either 10 heads in a row (reject  $H_0$ ) or 10 tails in a row (accept  $H_0$ )
4. Almost-Runs: As above, but 9 of 10 in a row
5. ...

Limited only by your ingenuity and ability to analyze.  
But how will you optimize Type I, II errors?

A generic decision rule: a “Likelihood Ratio Test”

$$\frac{L(x_1, x_2, \dots, x_n \mid H_1)}{L(x_1, x_2, \dots, x_n \mid H_0)} \geq c \quad \begin{cases} < c & \text{accept } H_0 \\ = c & \text{arbitrary} \\ > c & \text{reject } H_0 \end{cases}$$

E.g.:

$c = 1$ : accept  $H_0$  if observed data is *more* likely under that hypothesis than it is under the alternate, but reject  $H_0$  if observed data is more likely under the *alternate*

$c = 5$ : accept  $H_0$  unless there is *strong* evidence that the alternate is more likely (i.e., 5×)

Changing  $c$  shifts balance of Type I vs II errors, of course

Given: A coin, either fair ( $p(H)=1/2$ ) or biased ( $p(H)=2/3$ )

Decide: which

How? Flip it 5 times. Suppose outcome  $D = \text{HHHTH}$

Null Model/Null Hypothesis  $M_0: p(H) = 1/2$

Alternative Model/Alt Hypothesis  $M_1: p(H) = 2/3$

Likelihoods:

$$P(D | M_0) = (1/2) (1/2) (1/2) (1/2) (1/2) = 1/32$$

$$P(D | M_1) = (2/3) (2/3) (2/3) (1/3) (2/3) = 16/243$$

$$\text{Likelihood Ratio: } \frac{p(D | M_1)}{p(D | M_0)} = \frac{16/243}{1/32} = \frac{512}{243} \approx 2.1$$

I.e., alt model is  $\approx 2.1 \times$  more likely than null model, given data



## more jargon: simple vs composite hypotheses

---

A *simple* hypothesis has a single, fixed parameter value

E.g.:  $P(H) = 1/2$

A *composite* hypothesis allows multiple parameter values

E.g.;  $P(H) > 1/2$

Note that LRT is problematic for composite hypotheses; *which* value for the unknown parameter would you use to compute its likelihood?

### The Neyman-Pearson Lemma

If an LRT for a simple hypothesis  $H_0$  versus a simple hypothesis  $H_1$  has error probabilities  $\alpha, \beta$ , then any test with type I error  $\alpha' \leq \alpha$  must have type II error  $\beta' \geq \beta$  (and if  $\alpha' < \alpha$ , then  $\beta' > \beta$ )

In other words, to compare a simple hypothesis to a simple alternative, a likelihood ratio test *is as good as any* for a given error bound.

$H_0: P(H) = 1/2$  | Data: flip 100 times

$H_1: P(H) = 2/3$  | Decision rule: Accept  $H_0$  if  $\#H \leq 60$

$\alpha = P(\text{Type I err}) = P(\#H > 60 | H_0) \approx 0.018$

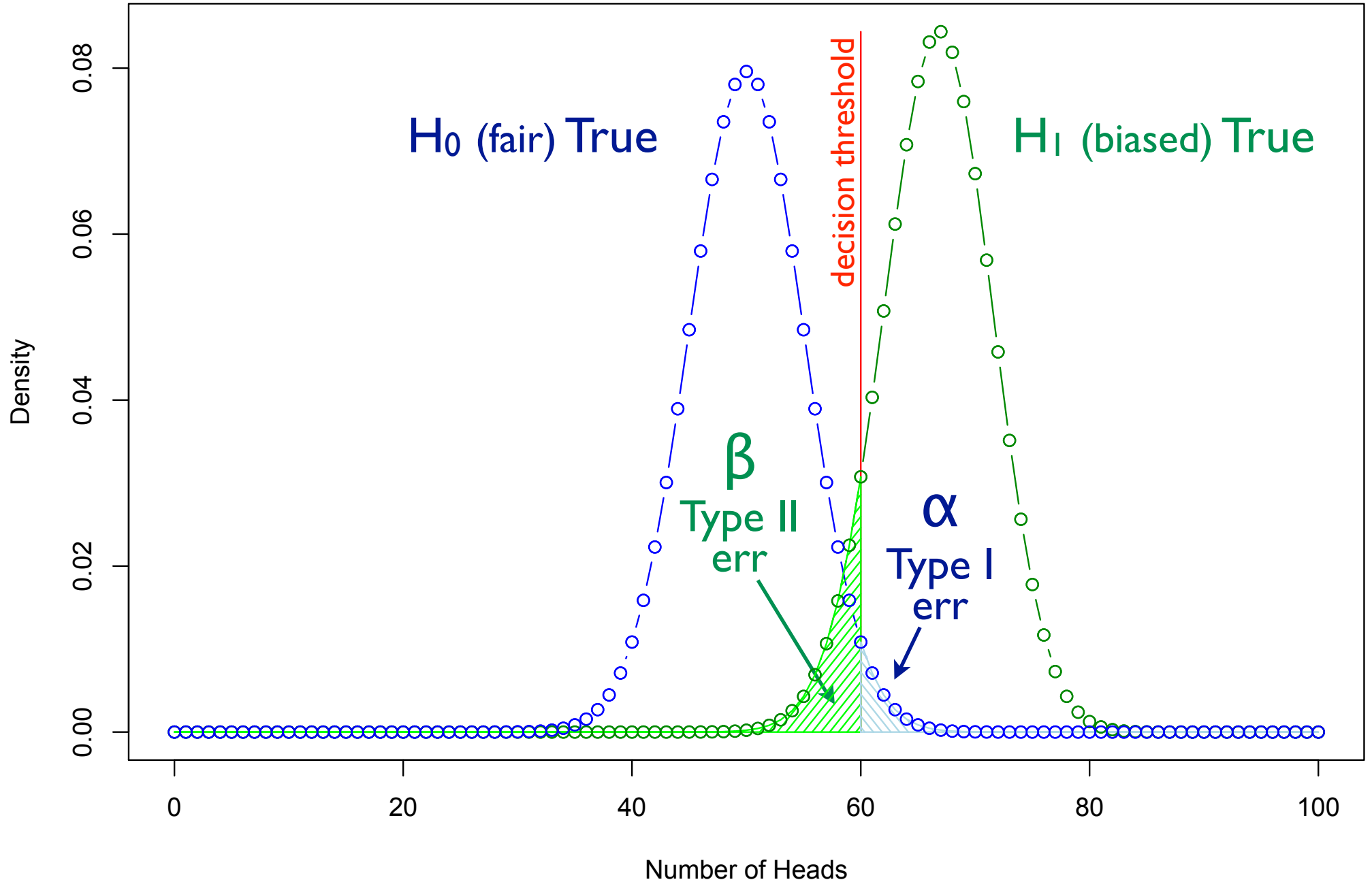
$\beta = P(\text{Type II err}) = P(\#H \leq 60 | H_1) \approx 0.097$

$$\frac{L(59 \text{ heads} | H_1)}{L(59 \text{ heads} | H_0)} \approx 1.4; \frac{L(60 \text{ heads} | H_1)}{L(60 \text{ heads} | H_0)} \approx 2.8; \frac{L(61 \text{ heads} | H_1)}{L(61 \text{ heads} | H_0)} \approx 5.7$$

$$\frac{L(60 \text{ heads} | H_1)}{L(60 \text{ heads} | H_0)} = \frac{\text{dbinom}(60, 100, 2/3)}{\text{dbinom}(60, 100, 1/2)} \approx 2.835788$$

↑ “R” pmf/pdf functions

$$\frac{L(60 \text{ heads} | H_1)}{L(60 \text{ heads} | H_0)} \approx \frac{\text{dnorm}(60, 100 \cdot 2/3, \sqrt{100 \cdot 2/3 \cdot 1/3})}{\text{dnorm}(60, 100 \cdot 1/2, \sqrt{100 \cdot 1/2 \cdot 1/2})} \approx 2.883173$$



Log of likelihood ratio is equivalent, often more convenient

add logs instead of multiplying...

“Likelihood Ratio Tests”: reject null if  $LLR > \text{threshold}$

$LLR > 0$  disfavors null, but higher threshold gives stronger evidence against

Neyman-Pearson Theorem: For a given error rate, LRT is as good a test as any (subject to some fine print).

Null/Alternative hypotheses - specify distributions from which data are assumed to have been sampled

Simple hypothesis - one distribution

E.g., “Normal, mean = 42, variance = 12”

Composite hypothesis - more than one distribution

E.g., “Normal, mean  $\geq 42$ , variance = 12”

Decision rule; “accept/reject null if sample data...”; *many* possible

Type 1 error: false reject/reject null when it is true

Type 2 error: false accept/accept null when it is false

Balance  $\alpha = P(\text{type 1 error})$  vs  $\beta = P(\text{type 2 error})$  based on “cost” of each

Likelihood ratio tests: for simple null vs simple alt, compare ratio of likelihoods under the 2 competing models to a fixed threshold.

Neyman-Pearson: LRT is best possible in this scenario.

# Significance Testing

B & T 9.4

# Recall

## (binary ) hypothesis testing

---

2 competing hypotheses  $H_0$  (the *null*),  $H_1$  (the *alternate*)

E.g.,  $P(\text{Heads}) = 1/2$  vs  $P(\text{Heads}) = 2/3$

Gather data,  $X$

Look at likelihood ratio  $\frac{L(X|H_1)}{L(X|H_0)}$ ; is it  $> c$ ?

Type I error/false reject rate  $\alpha$ ;

Type II error/false non-reject rate  $\beta$

Neyman-Pearson Lemma: no test will do better (for simple hyps)

Often the likelihood ratio formula can be massaged into an equivalent form that's simpler to use, e.g.

“Is #Heads  $> d$ ?”

Other tests, not based on likelihood, are also possible, say

“Is hyperbolic arc sine of #Heads in prime positions  $> 42$ ?”

but Neyman-Pearson still applies...



What about more general problems, e.g. with *composite* hypotheses?

E.g.,  $P(\text{Heads}) = \frac{1}{2}$  vs  $P(\text{Heads}) \text{ not } = \frac{1}{2}$

NB: LRT won't work – can't calculate likelihood for “ $p \neq \frac{1}{2}$ ”

Can I get a more nuanced answer than accept/reject?

General strategy:

Gather data,  $X_1, X_2, \dots, X_n$

Choose a real-valued *summary statistic*,  $S = h(X_1, X_2, \dots, X_n)$

Choose *shape* of the rejection region, e.g.  $R = \{X \mid S > c\}$ ,  $c$  t.b.d.

Choose *significance level*  $\alpha$  (upper bound on false rejection prob)

Find *critical value*  $c$ , so that, assuming  $H_0$ ,  $P(S > c) < \alpha$

No Neyman-Pearson this time, but (assuming you can do or approximate the math for last step) you now know the *significance* of the result – i.e., probability of falsely rejecting the null model.

## example: fair coin or not?

I have a coin. Is  $P(\text{Heads}) = 1/2$  or not?

General strategy:

Gather data,  $X_1, X_2, \dots, X_n$

Choose a real-valued *summary statistic*,  $S = h(X_1, X_2, \dots, X_n)$

Choose *shape* of the rejection region, e.g.  $R = \{X \mid S > c\}$ ,  $c$  t.b.d.

Choose *significance level*  $\alpha$  (upper bound on false rejection prob)

Find *critical value*  $c$ , so that, assuming  $H_0$ ,  $P(S > c) < \alpha$

For this example:

Flip  $n = 1000$  times:  $X_1, \dots, X_n$

*Summary statistic*,  $S = \#$  of heads in  $X_1, X_2, \dots, X_n$

*Shape* of the rejection region:  
 $R = \{X \text{ s.t. } |S - n/2| > c\}$ ,  $c$  t.b.d.

Choose *significance level*  
 $\alpha = 0.05$

Find *critical value*  $c$ , so that, assuming  $H_0$ ,  $P(|S - n/2| > c) < \alpha$

Given  $H_0$ ,  $(S - n/2)/\sqrt{n/4}$  is  $\approx \text{Norm}(0, 1)$ , so  $c = 1.96 * \sqrt{250} \approx 31$  gives the desired 0.05 significance level.

E.g., if you see 532 heads in 1000 flips you can reject  $H_0$  at the 5% significance level

The *p-value* of an experiment is:

$p = \min \{ \alpha \mid H_0 \text{ would be rejected at the } \alpha \text{ significance level} \}$

I.e., observed  $S$  is right at the critical value for  $\alpha = p$

I.e.,  $p = \text{prob of outcome as, or more, unexpected than observed}$

Why?

Shows directly your leeway w.r.t. any desired significance level.

Avoids pre-setting the significance level (pro/con)

Examples:

531 heads in 1000 flips has a p-value of 0.0537,  $> \alpha = 0.05$

532 heads in 1000 flips has a p-value of 0.0463,  $< \alpha = 0.05$

550 heads in 1000 flips has a p-value of 0.00173,  $\ll \alpha = 0.05$

nonrandom;  
it is or it isn't



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

It's the probability of seeing data this extreme, *assuming* null is true

example: is the mean zero or not ( $\sigma^2$  known)?

---

Suppose  $X \sim \text{Normal}(\mu, \sigma^2)$ , and  $\sigma^2$  is *known*.

$$H_0: \mu = 0 \quad \text{vs} \quad H_1: \mu \neq 0$$

Data:  $X_1, X_2, \dots, X_n$

Summary statistic – want something related to mean; how about:

$$S = \frac{X_1 + X_2 + \dots + X_n}{\sigma \sqrt{n}}$$

(*assuming*  $H_0$ ,  $\sum X_i$  has mean = 0, var =  $n \sigma^2$ , so  $S \sim N(0, 1)$  )

If we make rejection region  $R = \{ X \text{ s.t. } |S| > 1.96 \}$ , this will reject the null at the  $\alpha = 0.05$  significance level. I.e., *assuming*  $\mu = 0$ , an extreme sample with  $|S| > 1.96$  will be drawn only 5% of the time.

Similarly, if we observe  $S = 2.5$ , say, then p-value = 0.0124

## example: the t-test: is the mean zero or not ( $\sigma^2$ unknown)?

Suppose  $X \sim \text{Normal}(\mu, \sigma^2)$ , and  $\sigma^2$  is *unknown*.

$$H_0: \mu = 0 \quad \text{vs} \quad H_1: \mu \neq 0$$

Data:  $X_1, X_2, \dots, X_n$

$$\text{Let } \hat{\mu} = \sum_{i=1}^n \frac{x_i}{n} \quad ; \quad \hat{\sigma}^2 = \sum_{i=1}^n \frac{(x_i - \hat{\mu})^2}{n-1}$$

$$S = \frac{X_1 + X_2 + \dots + X_n}{\hat{\sigma} \sqrt{n}}$$

$S$  has a t-distribution with  $n-1$  degrees of freedom

Look up desired values in t-tables (e.g., B&T p 473; *see next slide*). E.g.,

$$\left. \begin{array}{l} \text{for } n = 10, \text{ use } R = \{ x \text{ s.t. } |S| > 2.26 \} \\ \text{for } n = 31, \text{ use } R = \{ x \text{ s.t. } |S| > 2.04 \} \end{array} \right\} \text{not } 1.96$$

to obtain  $\alpha = 0.05$  significance level. E.g.,  $n=10, S=3.25 \Rightarrow \text{p-value} = 0.01$

$\alpha/2$ 


	0.100	0.050	<u>0.025</u>	0.010	0.005	0.001
1	3.078	6.314	12.71	31.82	63.66	318.3
2	1.886	2.920	4.303	6.965	9.925	22.33
3	1.638	2.353	3.182	4.541	5.841	10.21
4	1.533	2.132	2.776	3.747	4.604	7.173
5	1.476	2.015	2.571	3.365	4.032	5.893
6	1.440	1.943	2.447	3.143	3.707	5.208
7	1.415	1.895	2.365	2.998	3.499	4.785
8	1.397	1.860	2.306	2.896	3.355	4.501
9	1.383	1.833	2.262	2.821	3.250	4.297
10	1.372	1.812	2.228	2.764	3.169	4.144
11	1.363	1.796	2.201	2.718	3.106	4.025
12	1.356	1.782	2.179	2.681	3.055	3.930
13	1.350	1.771	2.160	2.650	3.012	3.852
14	1.345	1.761	2.145	2.624	2.977	3.787
15	1.341	1.753	2.131	2.602	2.947	3.733
20	1.325	1.725	2.086	2.528	2.845	3.552
30	1.310	1.697	2.042	2.457	2.750	3.385
60	1.296	1.671	2.000	2.390	2.660	3.232
120	1.289	1.658	1.980	2.358	2.617	3.160
$\infty$	1.282	1.645	1.960	2.326	2.576	3.090

 $n-1$  $n-1$ 

CDF  $\Psi_{n-1}(z)$  of the  $t$ -distribution w/  $n-1$  degrees of freedom


lbsoff.com sells diet pills. 10 volunteers used them for a month, reporting the net weight changes of:

```
x <- c(-1.5, 0, .1, -0.5, -.25, 0.3, .1, .05, .15, .05)
> mean(x)
[1] -0.15
```



lbsoff proudly announces “Diet Pill Miracle!”

```
> cat("stddev=", sd(x), "tstat=", sum(x)/sd(x)/sqrt(10))
stddev= 0.5244044 tstat= -0.904534
> t.test(x)
t = -0.9045, df = 9, p-value = 0.3893
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval: -0.5251363 0.2251363
```



What do you think?

## significance testing – summary

---

Setup much like LRT case: Null  $H_0$  vs Alternate  $H_1$  hypotheses;  
Type I vs Type II errors;  $\alpha$  vs  $\beta$

Especially useful for *composite* hyps (where LRT is problematic)

Formulate a test statistic,  $S = h(X_1, \dots, X_n)$

Choose “rejection region”  $R$ , i.e., values of  $S$  that are too unlikely under  $H_0$  to be credible, typically parameterized by some constant  $c$

Choose “significance level”  $\alpha$  (e.g., 0.05), then calculate threshold  $c$  s.t. rejection probability  $< \alpha$ , and/or calculate p-value of  $S = h(X_1, \dots, X_n)$  i.e., probability of seeing data as extreme as, or more extreme than observed.

Bottom line: data in rejection region, w/ low  $\alpha$  and/or low p-value, is very unlikely assuming  $H_0$  is true; hinting towards  $H_1$

Now that you get p-values: here’s an amusing/depressing story:

<http://io9.com/i-fooled-millions-into-thinking-chocolate-helps-weight-1707251800>



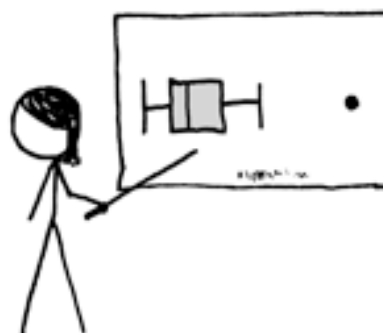
CAN MY BOYFRIEND  
COME ALONG?



I'M NOT YOUR  
BOYFRIEND!  
/ YOU TOTALLY ARE.  
I'M CASUALLY  
DATING A NUMBER  
OF PEOPLE.



BUT YOU SPEND TWICE AS MUCH  
TIME WITH ME AS WITH ANYONE  
ELSE. I'M A CLEAR OUTLIER.



YOUR MATH IS  
IRREFUTABLE.

FACE IT—I'M  
YOUR STATISTICALLY  
SIGNIFICANT OTHER.



Something Completely  
Different

Gene expression

Advance Access publication January 28, 2012

# A new approach to bias correction in RNA-Seq

Daniel C. Jones<sup>1,\*</sup>, Walter L. Ruzzo<sup>1,2,3</sup>, Xinxia Peng<sup>4</sup> and Michael G. Katze<sup>4</sup>

<sup>1</sup>Department of Computer Science and Engineering, University of Washington, Seattle, WA 98195-2350, <sup>2</sup>Department of Genome Sciences, University of Washington, Seattle, WA 98195-5065, <sup>3</sup>Fred Hutchinson Cancer Research Center, Seattle, WA 98109 and <sup>4</sup>Department of Microbiology, University of Washington, Seattle, WA

Associate Editor: Alex Bateman

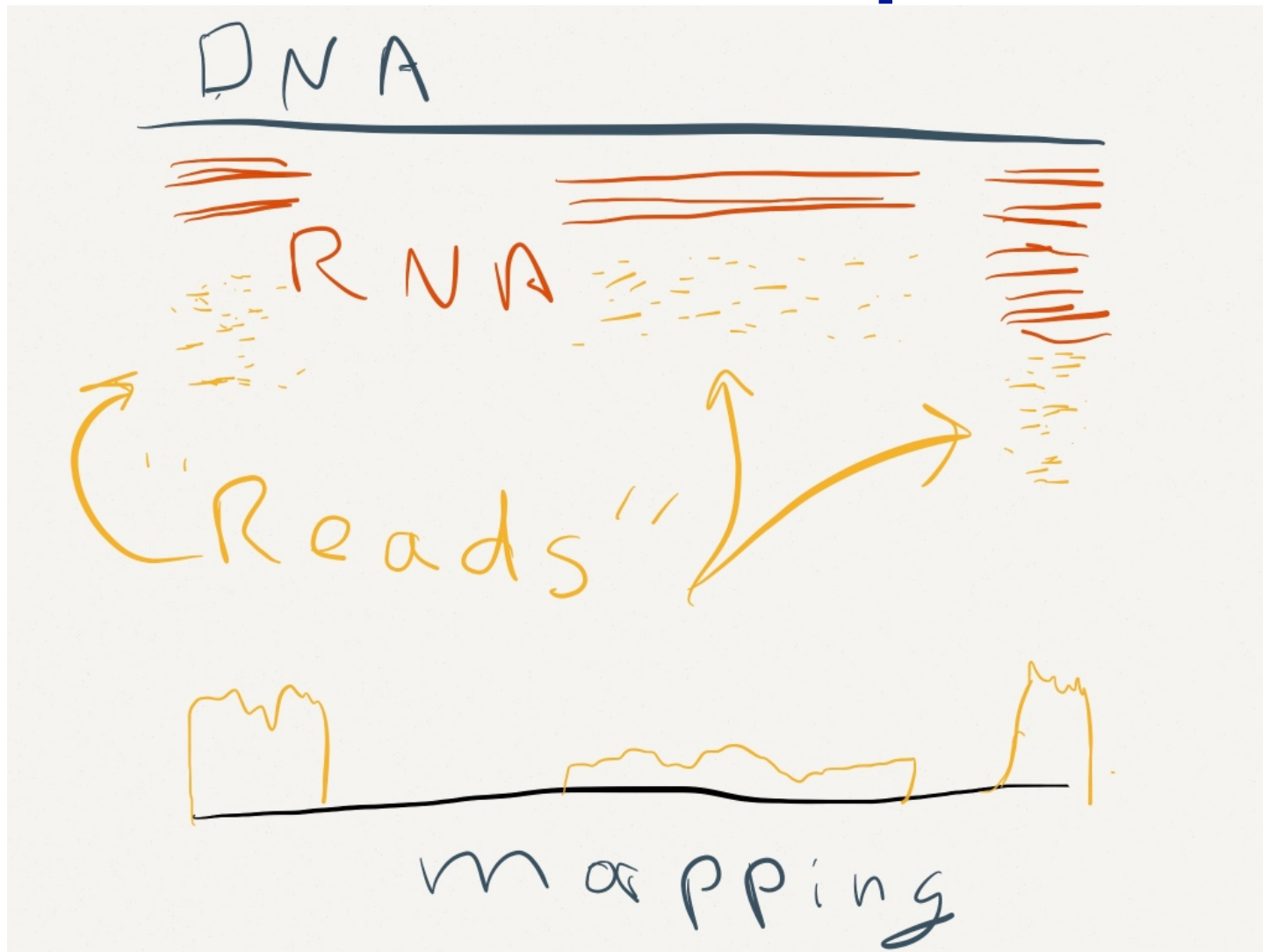
## ABSTRACT

**Motivation:** Quantification of sequence abundance in RNA-Seq experiments is often conflated by protocol-specific sequence bias. The exact sources of the bias are unknown, but may be influenced by... These biases may adversely affect...

Ok, Ok — Not on the final...

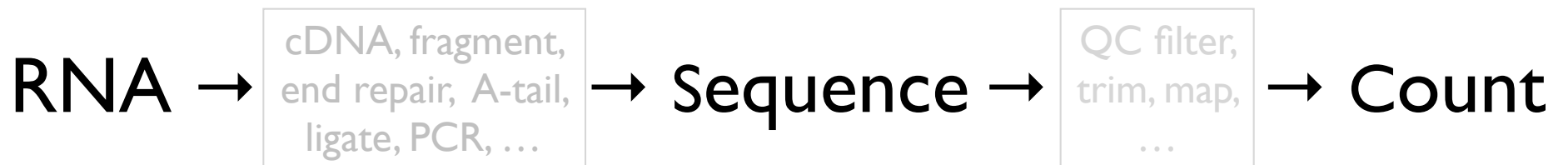
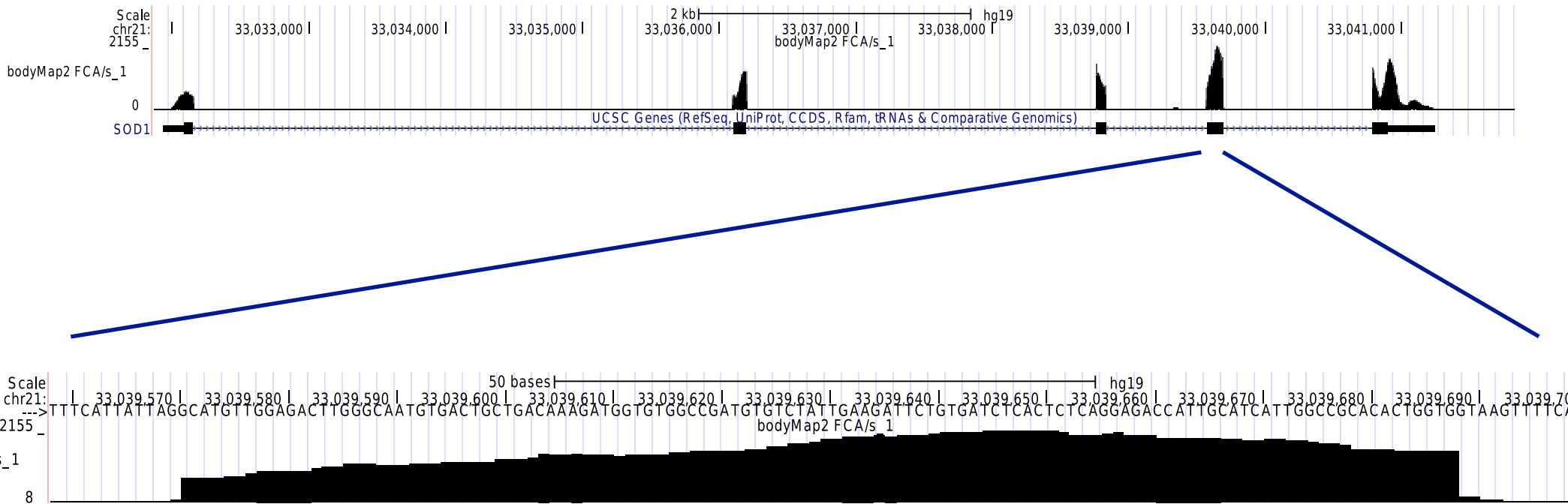


# RNAseq

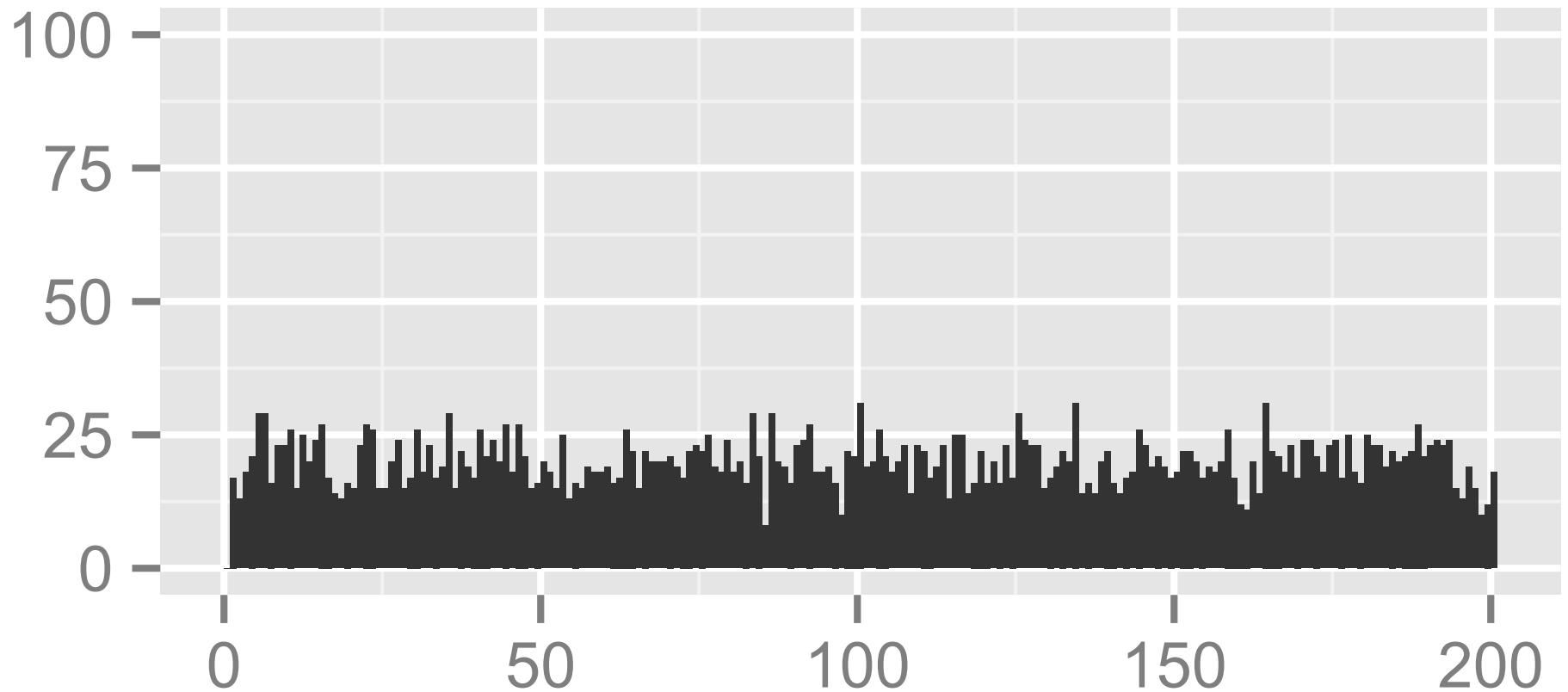


Cells make RNA. Biologists “read” it  
– a (biased!) random sampling process

# RNA seq



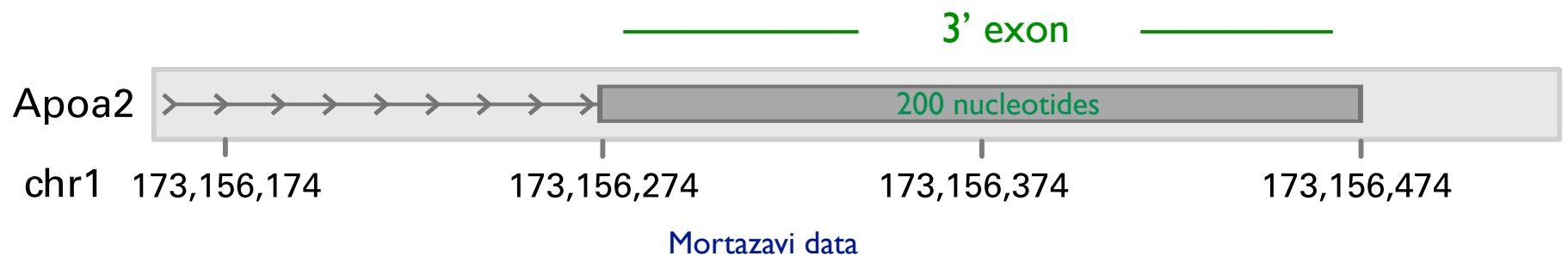
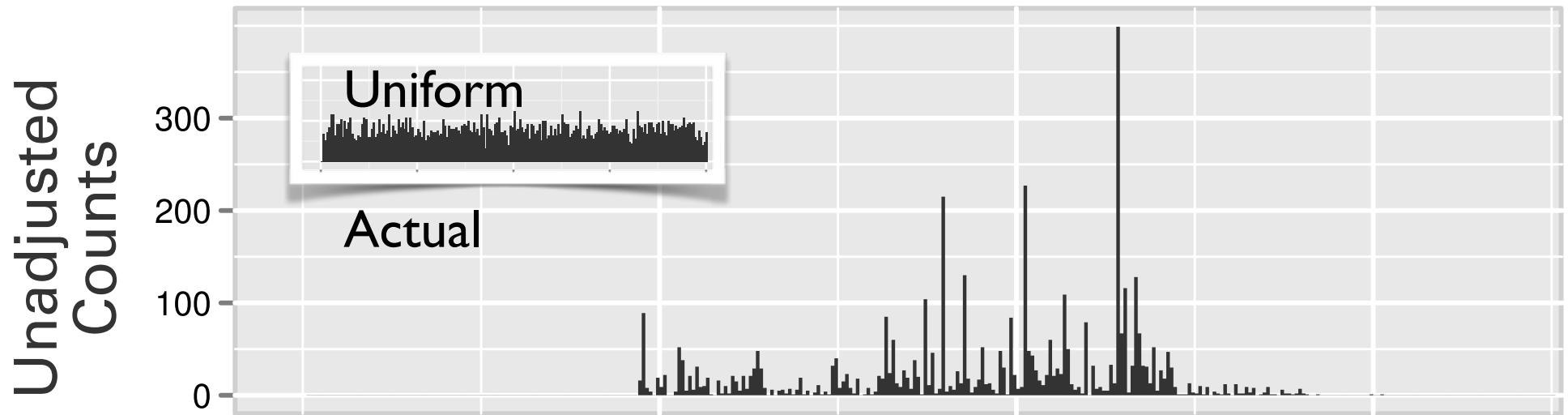
# What we expect: Uniform Sampling



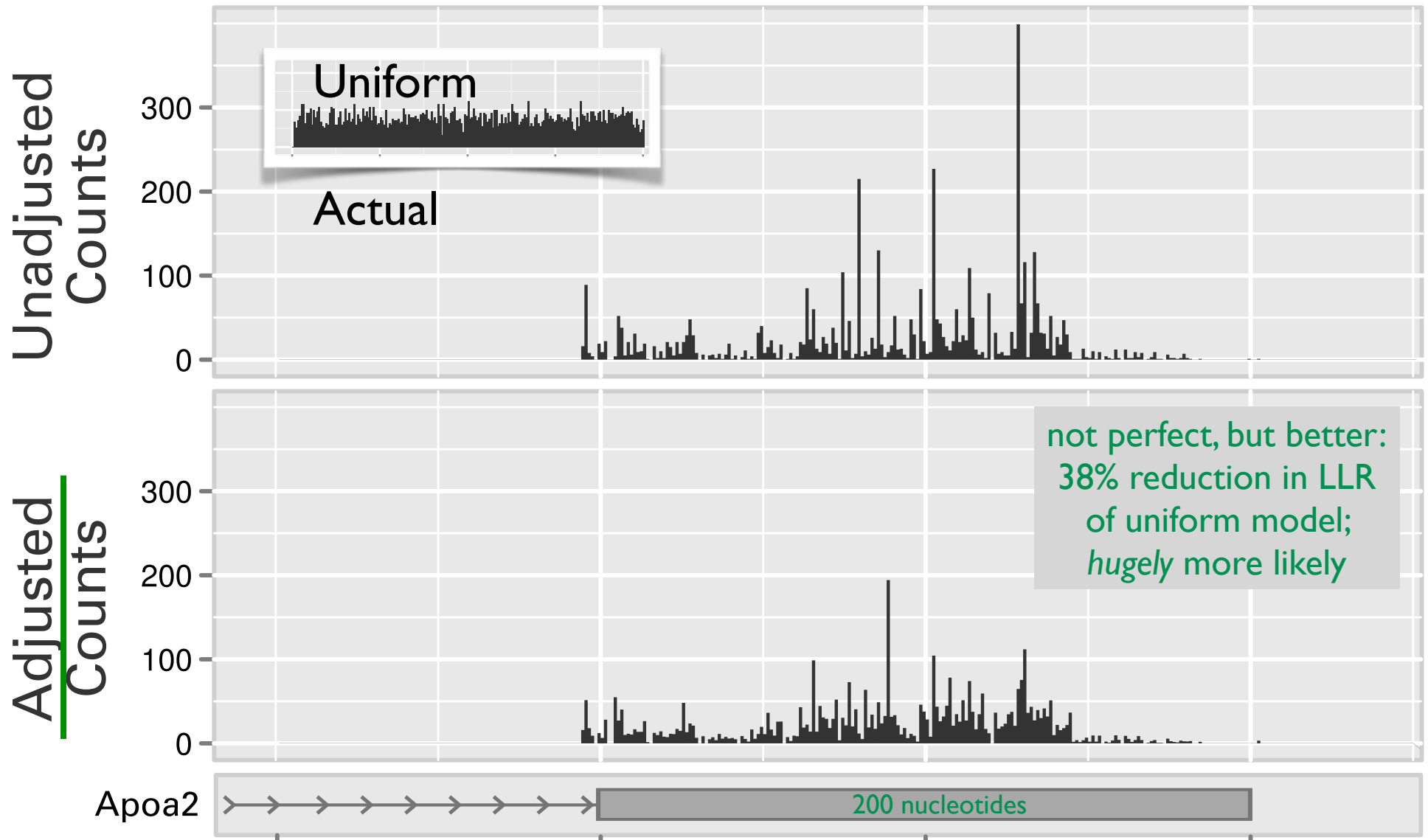
Uniform sampling of 4000 “reads” across a 200 bp “exon.”  
Average  $20 \pm 4.7$  per position, min  $\approx 9$ , max  $\approx 33$   
I.e., as expected, we see  $\approx \mu \pm 3\sigma$  in 200 samples

# What we get: *highly non-uniform coverage*

E.g., assuming uniform, the 8 peaks above 100 are  $\geq +10\sigma$  above mean



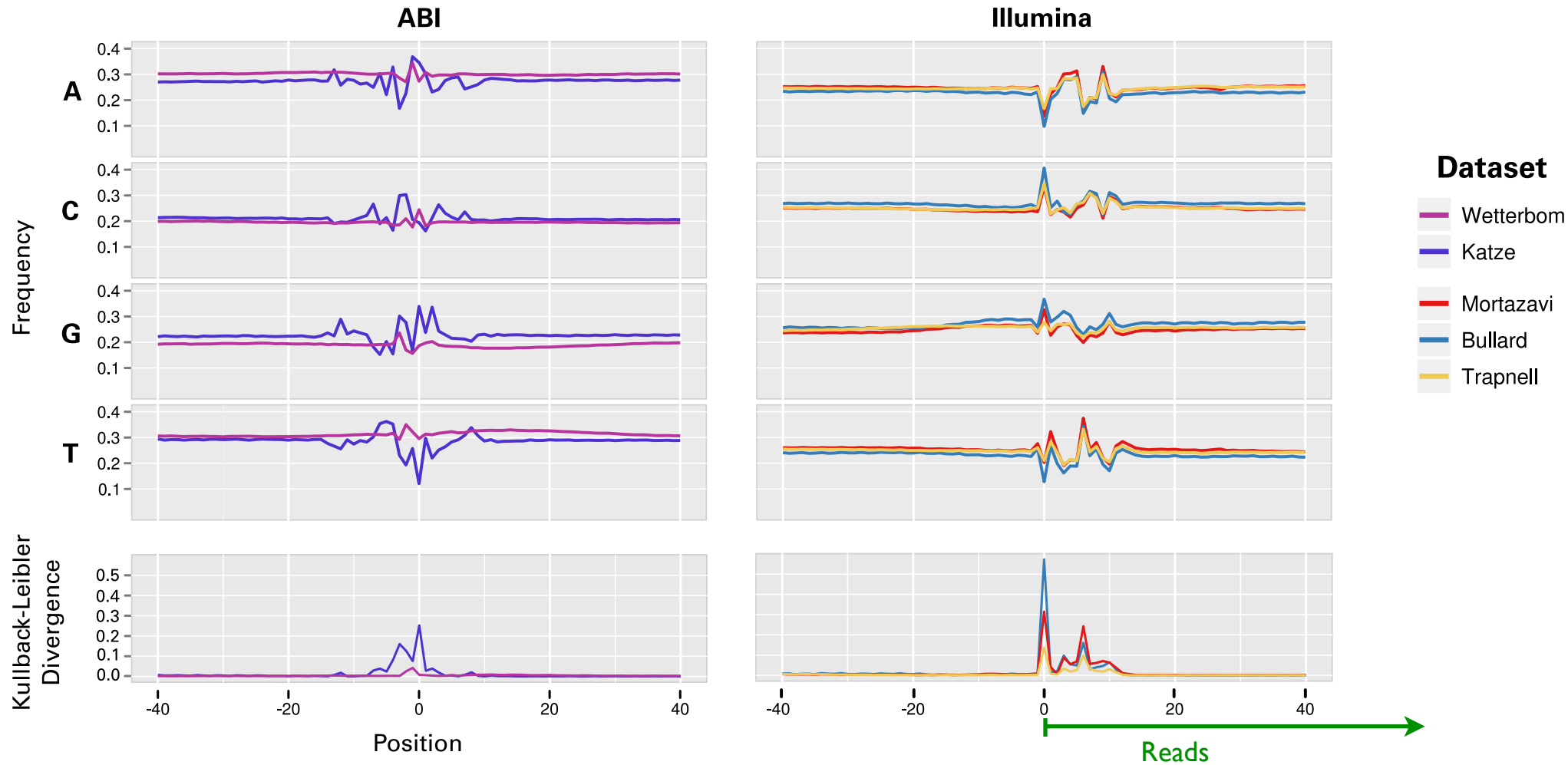
# What we get: *highly non-uniform coverage*



**The Good News:** we can (partially) correct the bias



# (in part) Bias is $\wedge$ sequence-dependent



and platform/sample-dependent

Fitting a model of the sequence surrounding read starts lets us predict which positions have more reads.

you know  
this

$$E[x_i|s_i] = N \Pr[m_i|s_i] = N \Pr[m_i] = E[x_i]$$

From Bayes' rule,

$$\Pr[m_i|s_i] = \frac{\Pr[s_i|m_i]\Pr[m_i]}{\Pr[s_i]}$$

This suggests a natural scheme in which observations may be reweighted to correct for bias. First, define the *sequence bias*  $b_i$  at position  $i$  as  $b_i = \Pr[s_i]/\Pr[s_i|m_i]$ .

Now, if we reweight the read count  $x_i$  at position  $i$  by  $b_i$ , we have,

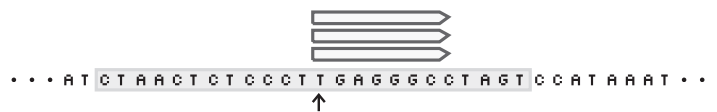
$$\begin{aligned} E[b_i x_i | s_i] &= b_i E[x_i | s_i] \\ &= N b_i \Pr[m_i | s_i] \\ &= N \frac{\Pr[m_i | s_i] \Pr[s_i]}{\Pr[s_i | m_i]} \\ &= N \Pr[m_i] \\ &= E[x_i] \end{aligned}$$

you could  
do this

Thus, the reweighted read counts are made unbiased.

# Method Outline

(a) sample foreground sequences



(b) sample background sequences

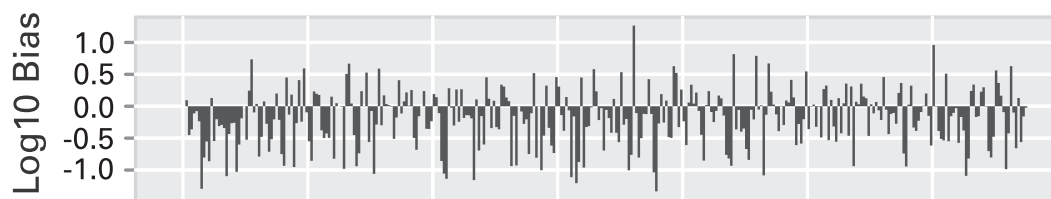


(c) train Bayesian network



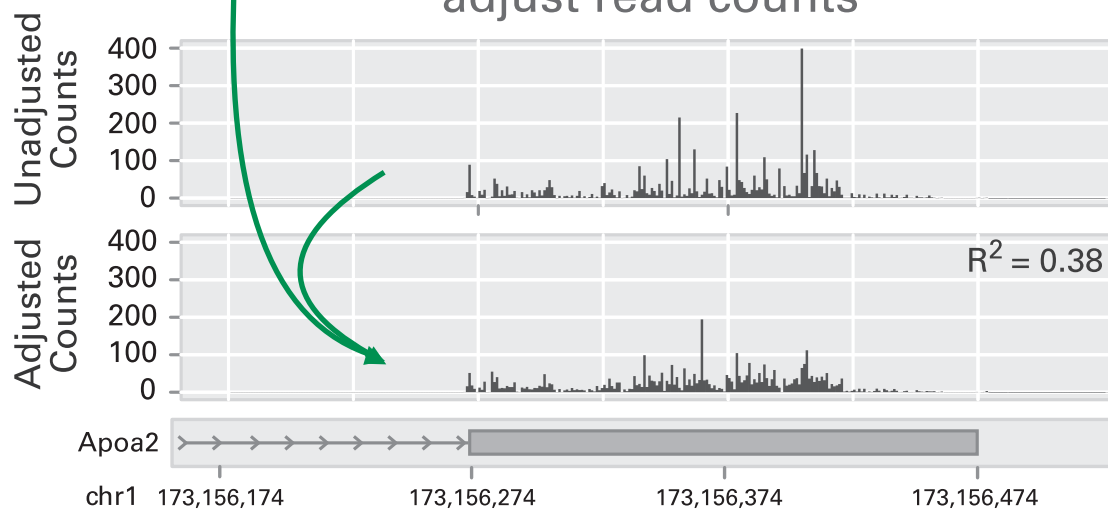
(d)

predict bias



(e)

adjust read counts



Want a probability distribution over k-mers,  $k \approx 40$

Some obvious choices

Full joint distribution:  $4^k - 1$  parameters

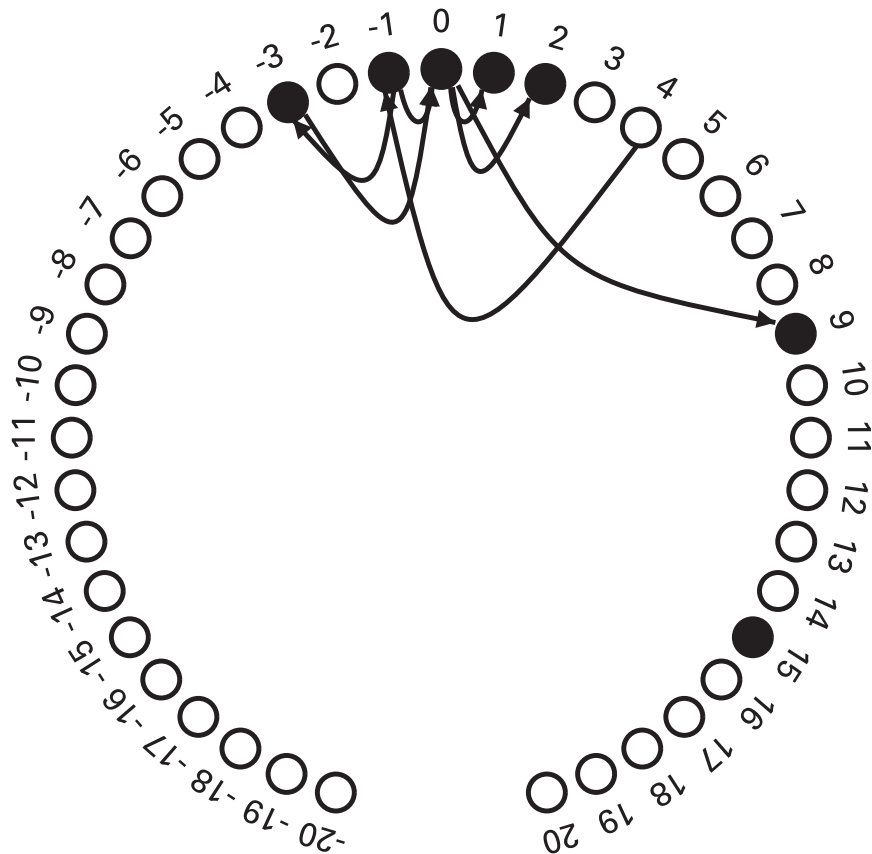
PWM (0-th order Markov):  $(4 - 1) \cdot k$  parameters

Something intermediate

Directed Bayes network

# Form of the models:

## Directed Bayes nets



**Wetterbom  
(282 parameters)**

One “node” per nucleotide,  
 $\pm 20$  bp of read start

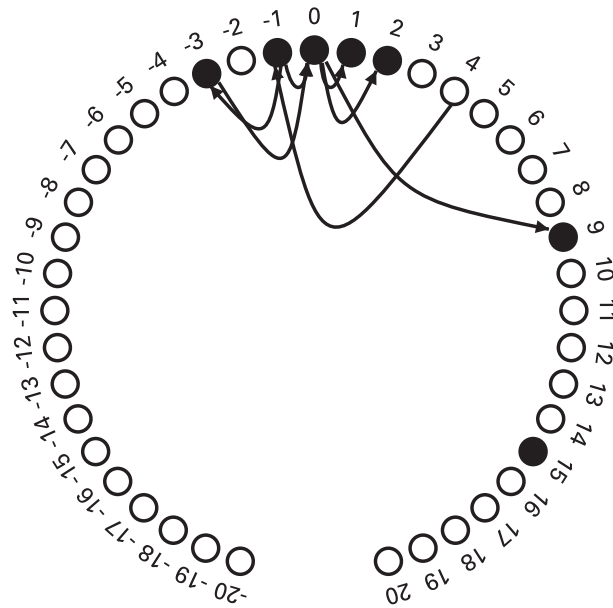
- Filled node means that position is biased
- Arrow  $i \rightarrow j$  means letter at position  $i$  modifies bias at  $j$
- For both, numeric params say how much

How—optimize:

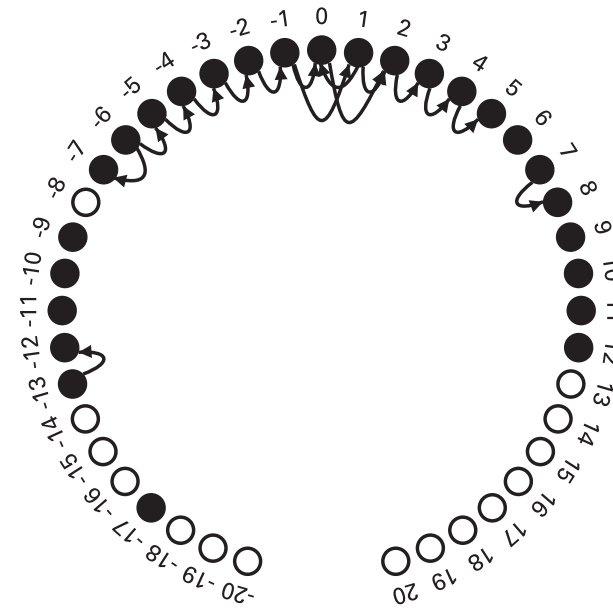
$$\ell = \sum_{i=1}^n \log \Pr[x_i | s_i] = \sum_{i=1}^n \log \frac{\Pr[s_i | x_i] \Pr[x_i]}{\sum_{x \in \{0,1\}} \Pr[s_i | x] \Pr[x]}$$

you could do this: somewhat like EM

ABI



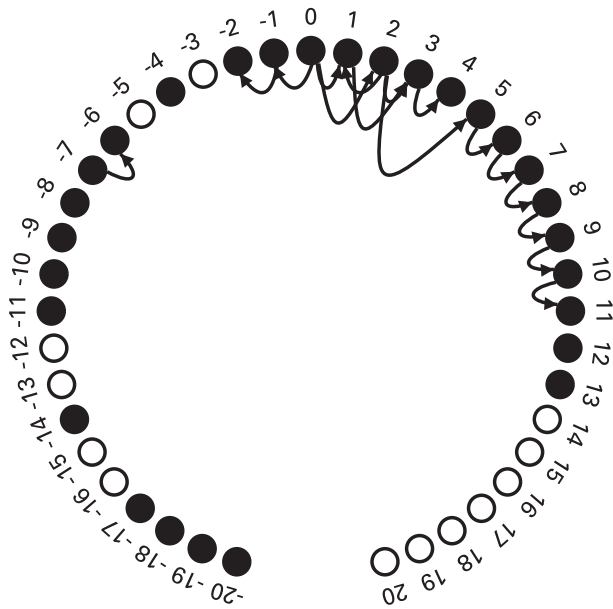
**Wetterbom**  
(282 parameters)



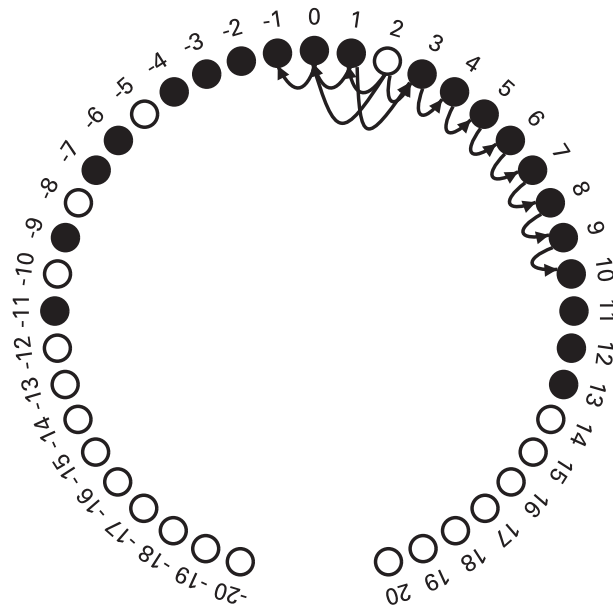
**Katze**  
(684 parameters)

- NB:**
- Not just initial hexamer
  - Span  $\geq 19$
  - All include negative positions
  - All different, even on same platform

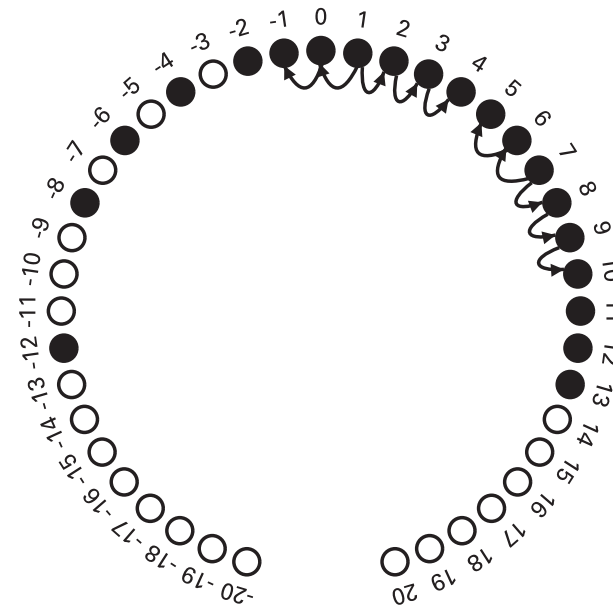
Illumina



**Bullard**  
(696 parameters)



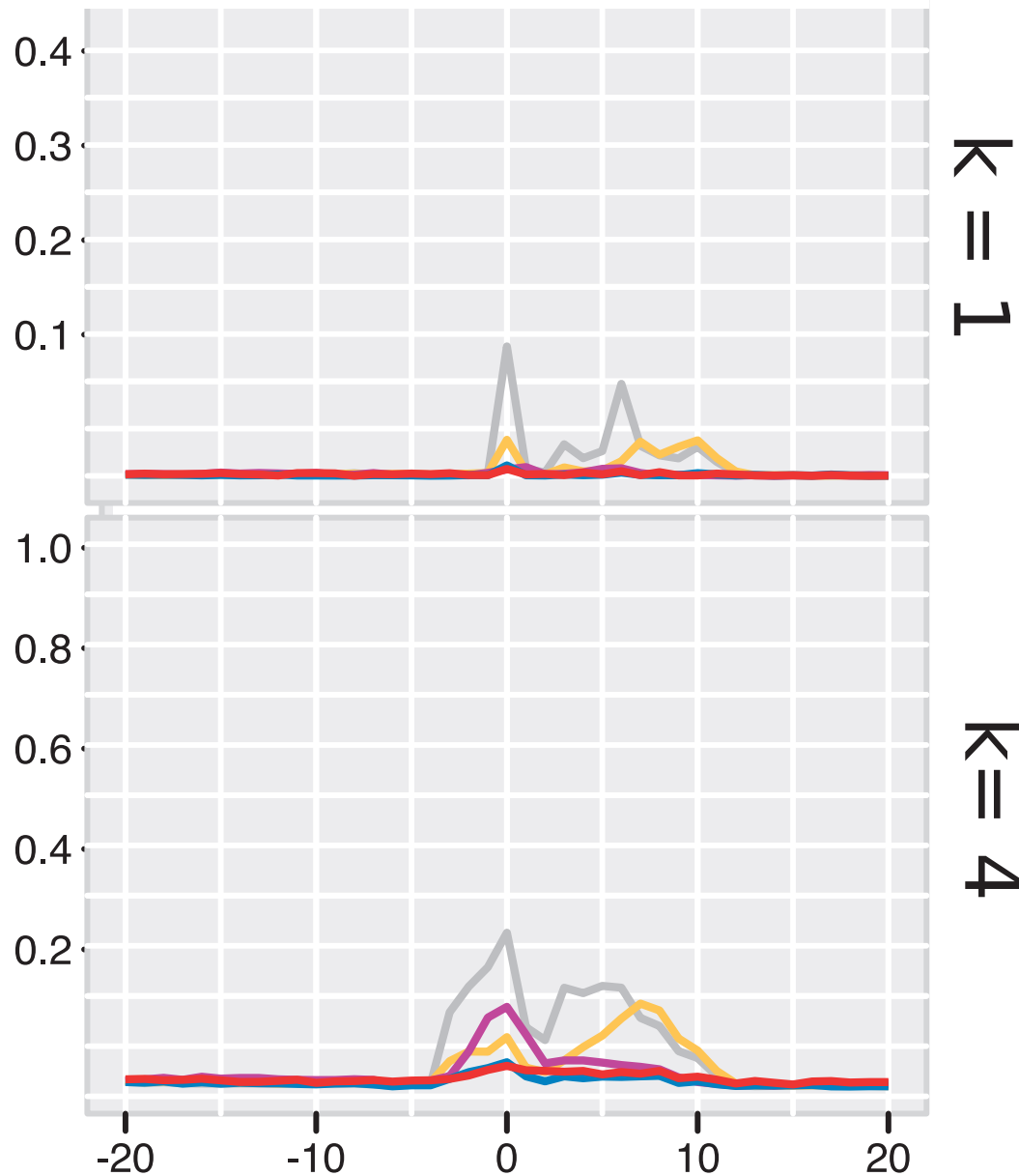
**Mortazavi**  
(582 parameters)



**Trapnell**  
(360 parameters)

# Result – Increased Uniformity

Kullback-Leibler Divergence

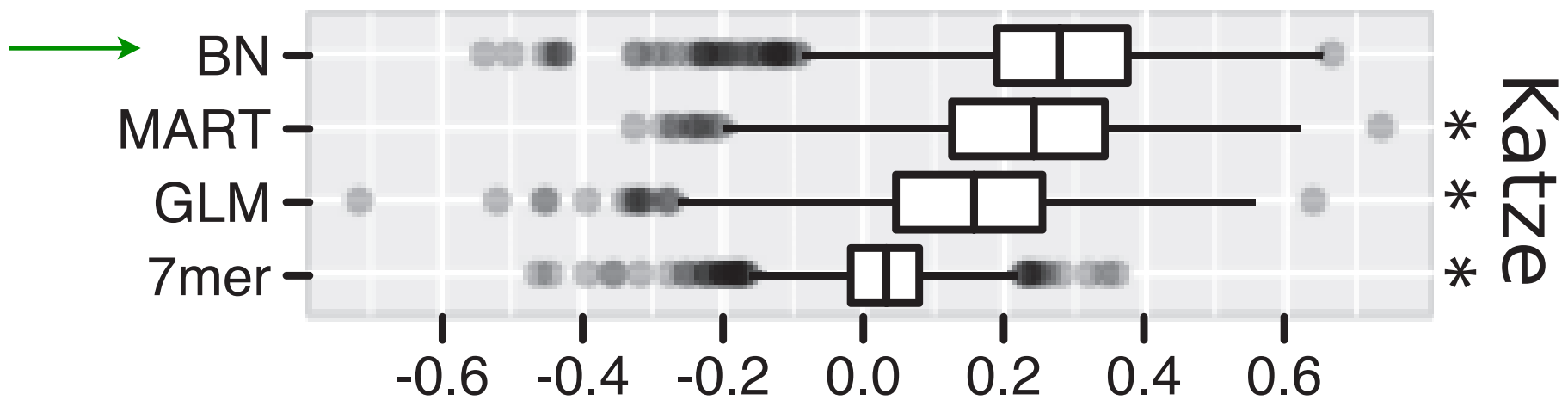
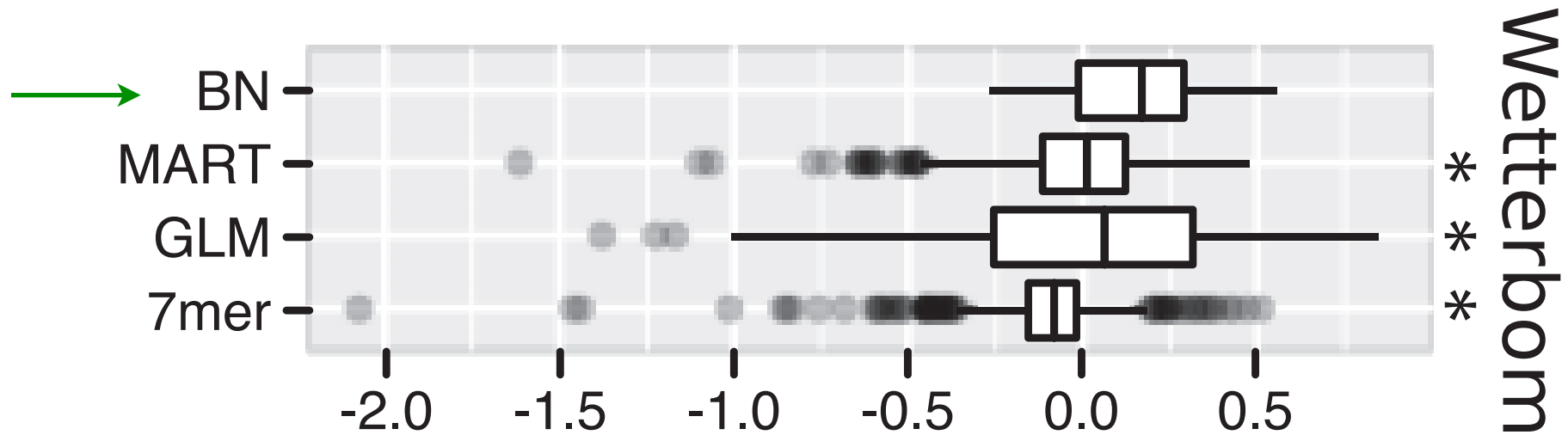


Method

- BN ← Jones
- MART | Li et al
- GLM |
- 7mer Hansen et al
- Unadjusted

Trapnell Data

# Result – Increased Uniformity



Fractional improvement  
in log-likelihood under  
uniform model across  
1000 exons ( $R^2=1-L'/L$ )

→  $R^2$

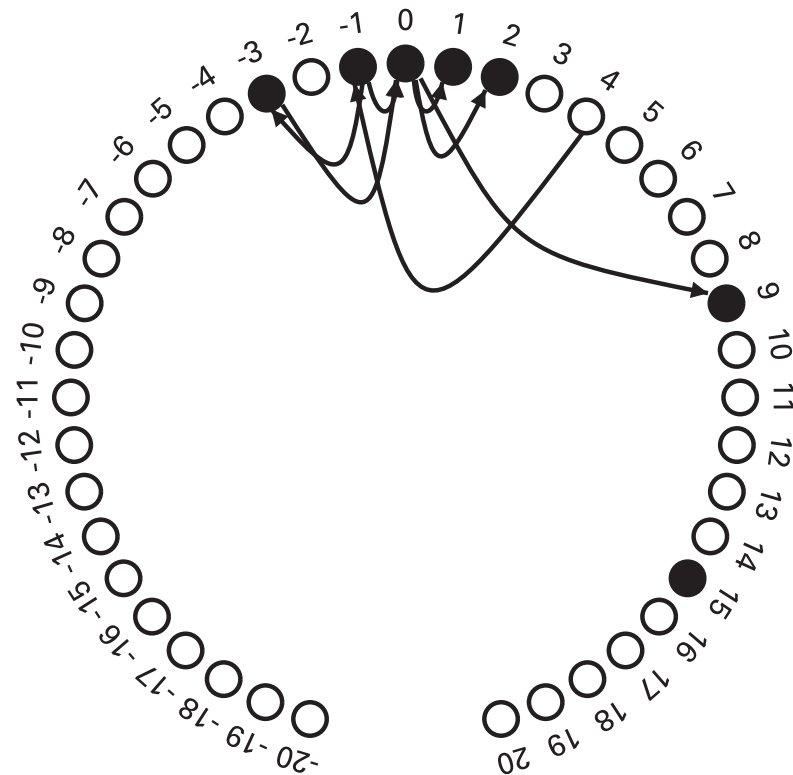
\* = p-value <  $10^{-23}$

you could do this: a  
hypothesis test "Is BN  
better than X?"



1. How does the amount of training data effect accuracy of the resulting model?

2. What is the chance that we will learn an incorrect model? E.g., learn a biased model from unbiased input?



**Wetterbom  
(282 parameters)**

# “First, do no harm”

Theorem:

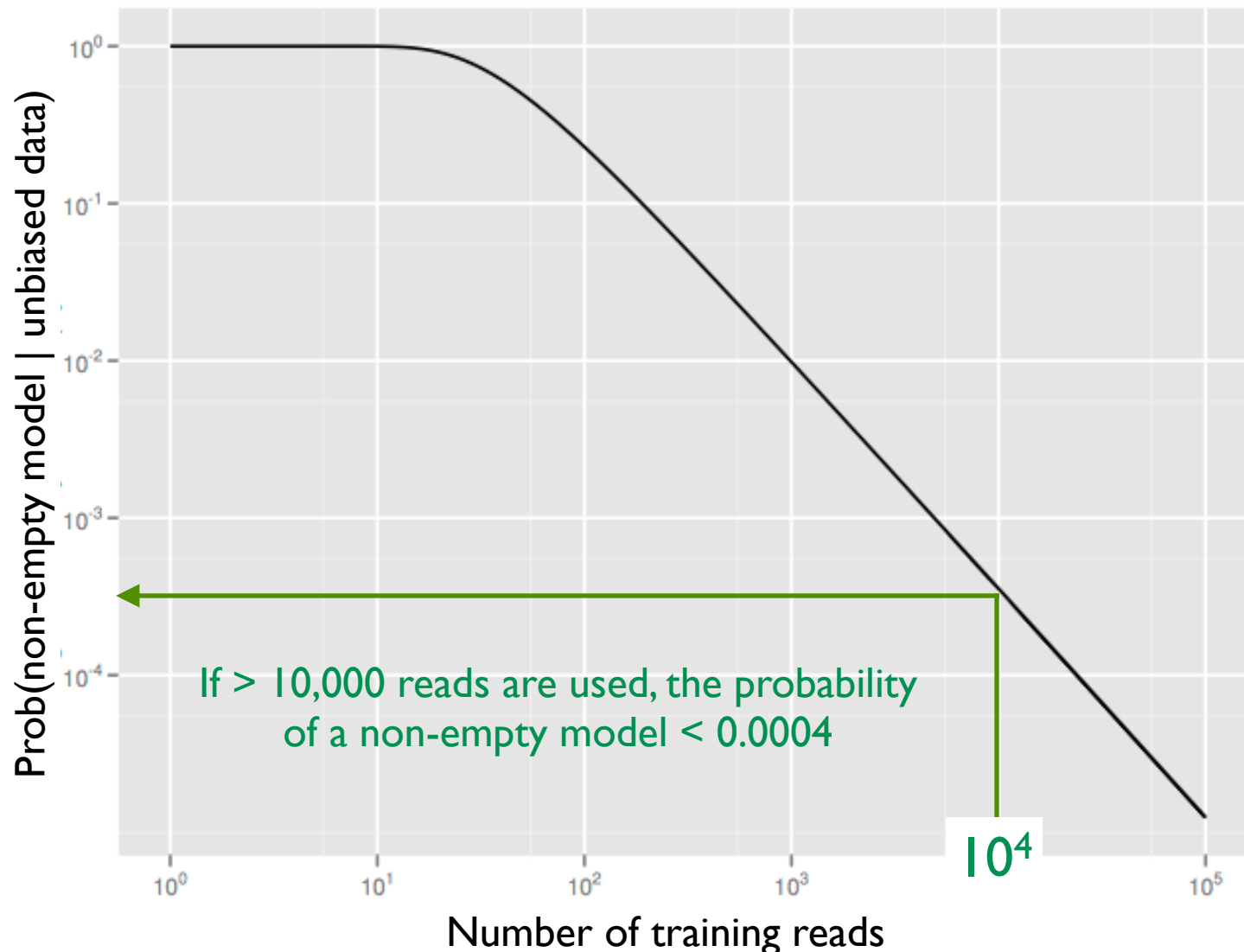
The probability of “false bias discovery,” i.e., of learning a non-empty model from  $n$  reads sampled from *unbiased* data is less than

$$1 - (\Pr(X < 3 \log n))^{2h}$$

where  $h$  = number of nucleotides in the model and  $X$  is a random variable that (asymptotically in  $n$ ) is  $\chi^2$  with 3 degrees of freedom. ( $E[X] = 3$ )

# “First, do no harm”

*Theorem:* The probability of “false bias discovery,” i.e., of learning a non-empty model from  $n$  reads sampled from unbiased data, declines *exponentially* with  $n$ .



## how different are two distributions?

Given:  $r$ -sided die, with probs  $p_1 \dots p_r$  of each face. Roll it  $n=10,000$  times; observed frequencies =  $q_1, \dots, q_r$ , (the MLEs for the unknown  $p_i$ 's). How close is  $p_i$  to  $q_i$ ?

*Kullback-Leibler divergence*, also known as *relative entropy*, of  $Q$  with respect to  $P$  is defined as

$$H(Q||P) = \sum_i q_i \ln \frac{q_i}{p_i}$$

where  $q_i$  ( $p_i$ ) is the probability of observing the  $i^{\text{th}}$  event according to the distribution  $Q$  (resp.,  $P$ ), and the summation is taken over all events in the sample space (e.g., all  $k$ -mers). In some sense, this is a measure of the dissimilarity between the distributions: if  $p_i \approx q_i$  everywhere, their log ratios will be near zero and  $H$  will be small; as  $q_i$  and  $p_i$  diverge, their log ratios will deviate from zero and  $H$  will increase.

Fancy name, simple idea:  $H(Q||P)$  is just the expected per-sample contribution to log-likelihood ratio test for “was  $X$  sampled from  $H_0: P$  vs  $H_1: Q$ ?”

So, assuming the null hypothesis is false, in order for it to be rejected with say, 1000 : 1 odds, one should choose  $m$  to be inversely proportional to  $H(Q||P)$ :

$$mH(Q||P) \geq \ln 1000$$
$$m \geq \frac{\ln 1000}{H(Q||P)}$$

you  
could  
do this

---

Continuing the notation above, suppose  $P$  as an unknown distribution with parameters  $p_1, \dots, p_r$ ,  $\sum p_i = 1$  where  $r$  is the number of points in the sample space (e.g.  $r = 4^k$  in the case of  $k$ -mers). Given a random sample  $X_1, X_2, \dots, X_r$  of size  $n = \sum_i X_i$  from  $P$ , it is well known that the maximum likelihood estimators for the parameters are  $q_i = \frac{X_i}{n} \approx p_i$ . How good an estimate for  $P$  is this distribution  $Q$ ? The estimators are unbiased:

$$E[q_i] = E\left[\frac{X_i}{n}\right] = \frac{E[X_i]}{n} = \frac{np_i}{n} = p_i$$

and the standard deviation of each estimate is proportional to  $1/\sqrt{n}$ , so these estimates are increasingly accurate as the sample size increases. A more quantitative assessment of the accuracy of the estimator is obtained by evaluating the KL divergence:

$$H(Q||P) = \sum_{i=1}^r q_i \ln \frac{q_i}{p_i} = \sum_{i=1}^r q_i \ln \left(1 + \frac{q_i - p_i}{p_i}\right)$$

---

Using the first two terms of the Taylor series for  $\ln(1 + x)$ , this is

$$\begin{aligned} H(Q||P) &\approx \sum_{i=1}^r q_i \left( \frac{q_i - p_i}{p_i} - \frac{1}{2} \left( \frac{q_i - p_i}{p_i} \right)^2 \right) \\ &= \sum_{i=1}^r q_i \frac{q_i - p_i}{p_i} - \frac{q_i}{2p_i} \frac{(q_i - p_i)^2}{p_i} \end{aligned}$$

Since  $\sum_{i=1}^r q_i = \sum_{i=1}^r p_i = 1$ ,  $\sum_{i=1}^r p_i \frac{q_i - p_i}{p_i} = 0$ , so

$$\begin{aligned} H(Q||P) &\approx \sum_{i=1}^r q_i \frac{q_i - p_i}{p_i} - p_i \frac{q_i - p_i}{p_i} - \frac{q_i}{2p_i} \frac{(q_i - p_i)^2}{p_i} \\ &= \sum_{i=1}^r \frac{(q_i - p_i)^2}{p_i} \left( 1 - \frac{q_i}{2p_i} \right) \\ &\approx \frac{1}{2} \sum_{i=1}^r \frac{(q_i - p_i)^2}{p_i} \end{aligned}$$

since  $q_i \approx p_i$ . Multiplying by  $n^2/n^2$  we have,

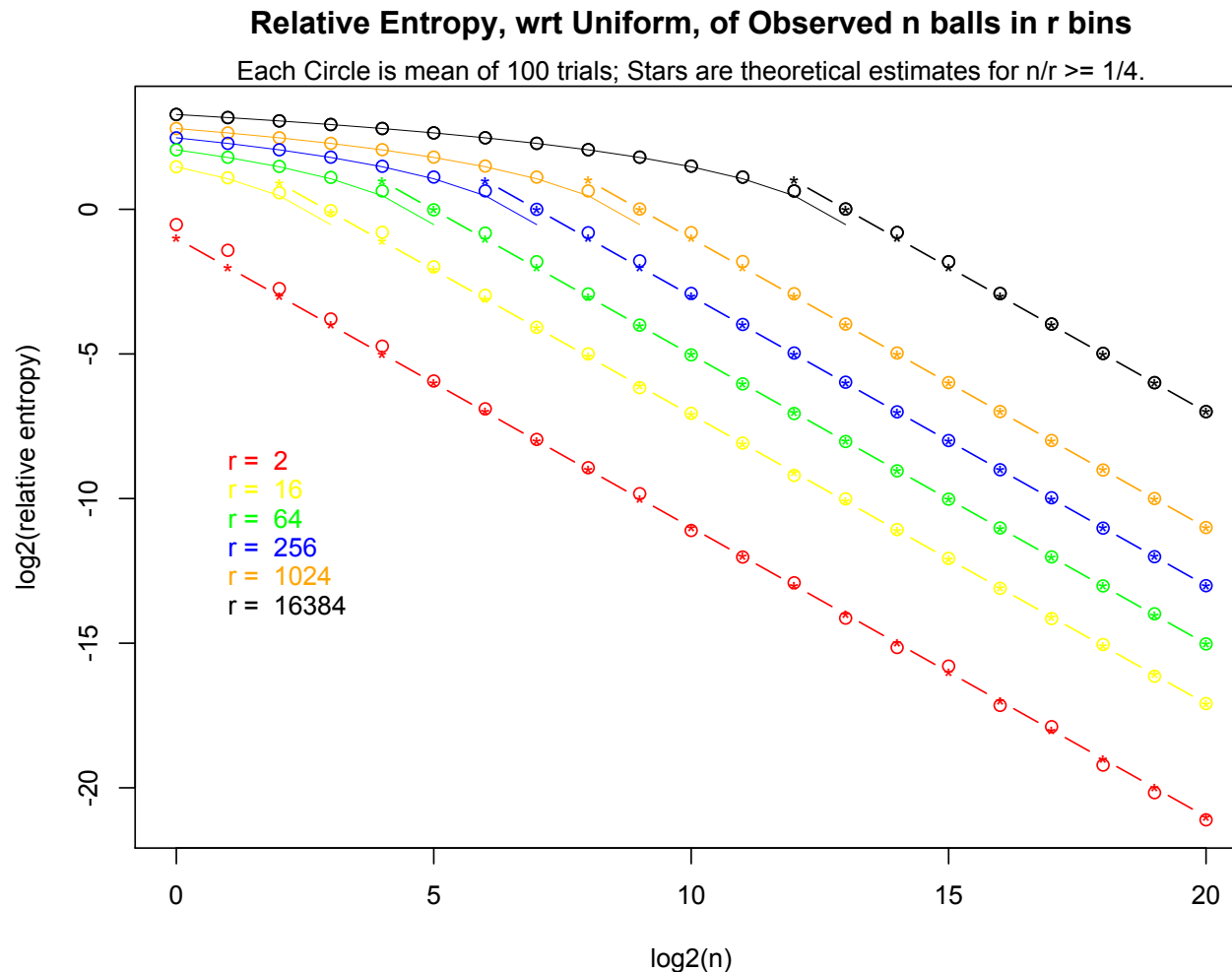
$$\begin{aligned} H(Q||P) &\approx \frac{1}{2n} \sum_{i=1}^r \frac{(nq_i - np_i)^2}{np_i} \\ &= \frac{1}{2n} \sum_{i=1}^r \frac{(X_i - E[X_i])^2}{E[X_i]} \end{aligned}$$

... and after a modicum of algebra:

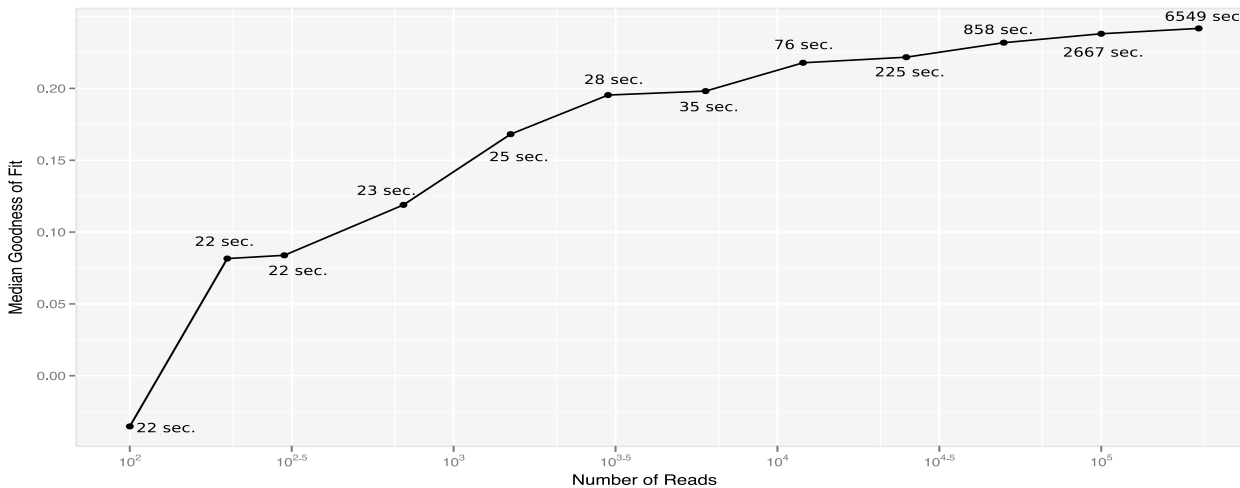
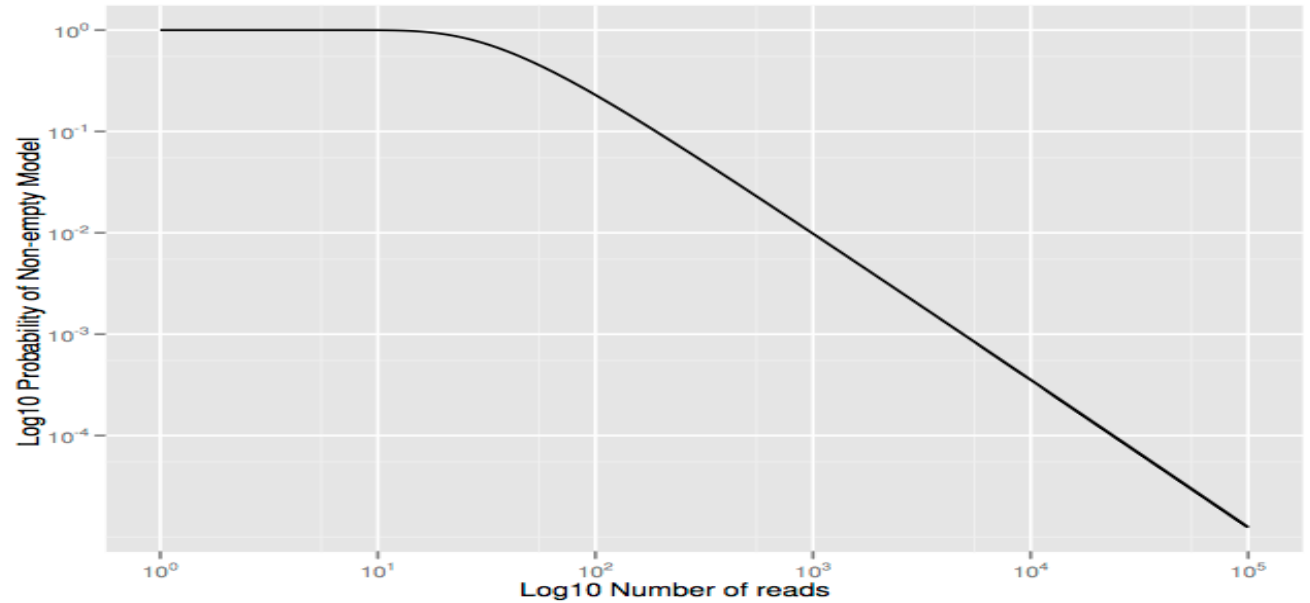
$$E[H(Q||P)] \approx \frac{r-1}{2n}$$

You could do this, too:  
LLR of error rises with  
number of parameters  $r$ ;  
declines with size of  
training set  $n$

... which empirically is a good approximation:



... and so the probability of falsely inferring “bias” from an unbiased sample declines rapidly with size of training set (while runtime rises)



you could do this, too: more algebra (albeit Daniel was a bit clever)

Figure 8: Median  $R^2$  is plotted against training set size. Each point is additionally labeled with the run time of the training procedure.





Home

Install

Help

Home » [Bioconductor 2.12](#) » [Software Packages](#) » seqbias

## seqbias

### Estimation of per-position bias in high-throughput sequencing data

Bioconductor version: Release (2.12)

This package implements a model of per-pos using a simple Bayesian network, the structu reads and a reference genome sequence.

Author: Daniel Jones <dcjones at cs.washing

Maintainer: Daniel Jones <dcjones at cs.wasl

To install this package, start R and enter:

```
source("http://bioconductor.org/
biocLite("seqbias")
```

To cite this package in a public

```
citation(""
```

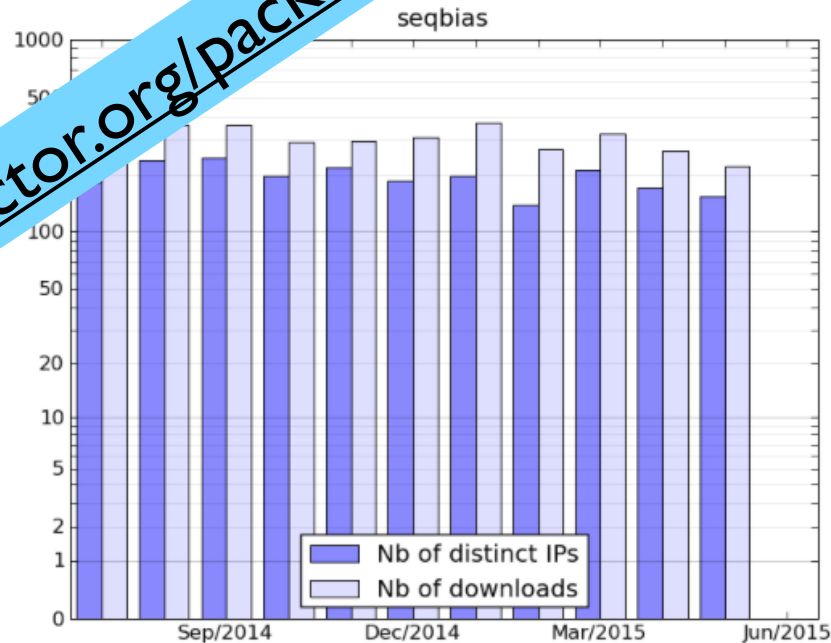
#### Docum

Assessing and Adjusting  
Reference Manual

### Download stats for Software package seqbias

This page was generated on 2015-06-01 06:29:02 -0700 (Mon, 01 Jun 2015).

seqbias home page: [release version](#), [devel version](#).



Month	Nb of distinct IPs	Nb of downloads
Jul/2014	181	252
Aug/2014	236	360
Sep/2014	242	360
Oct/2014	197	292
Nov/2014	217	299
Dec/2014	186	311
Jan/2015	195	371
Feb/2015	138	270
Mar/2015	211	327
Apr/2015	170	264
May/2015	153	220
Jun/2015	0	0
<b>All months</b>	<b>1648</b>	<b>3326</b>

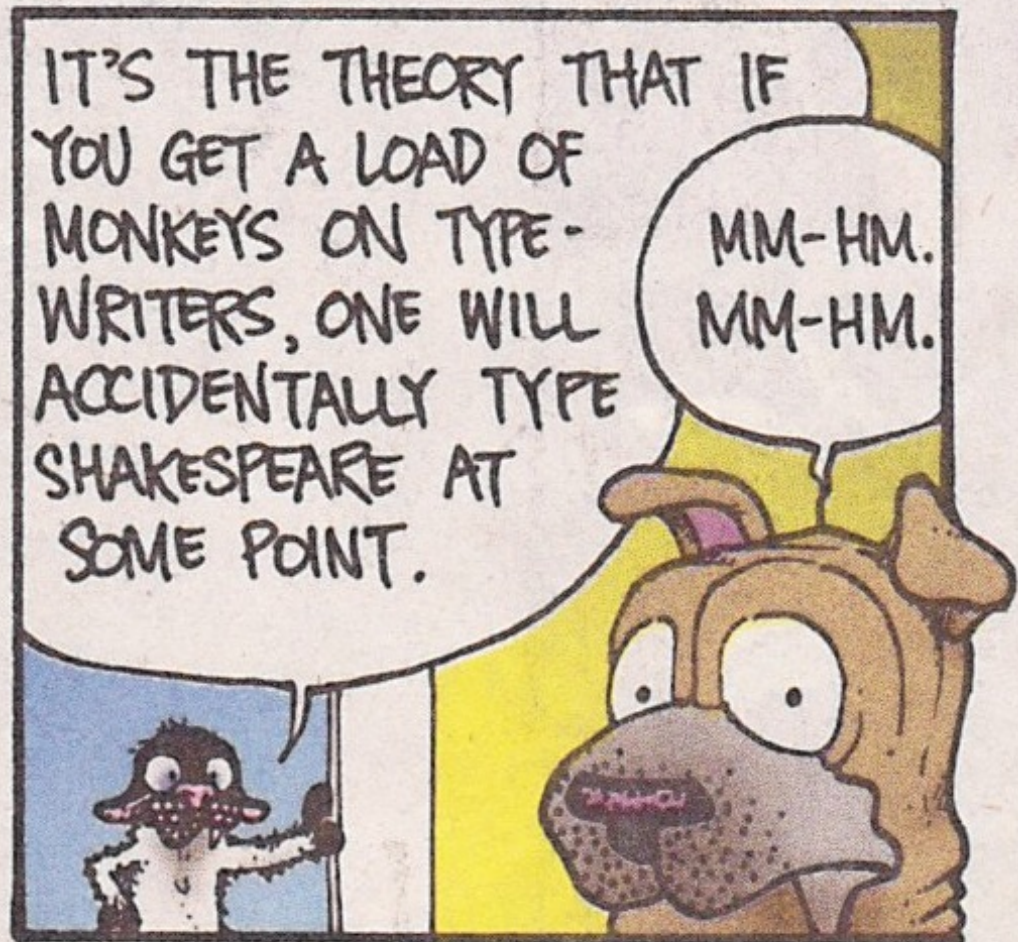
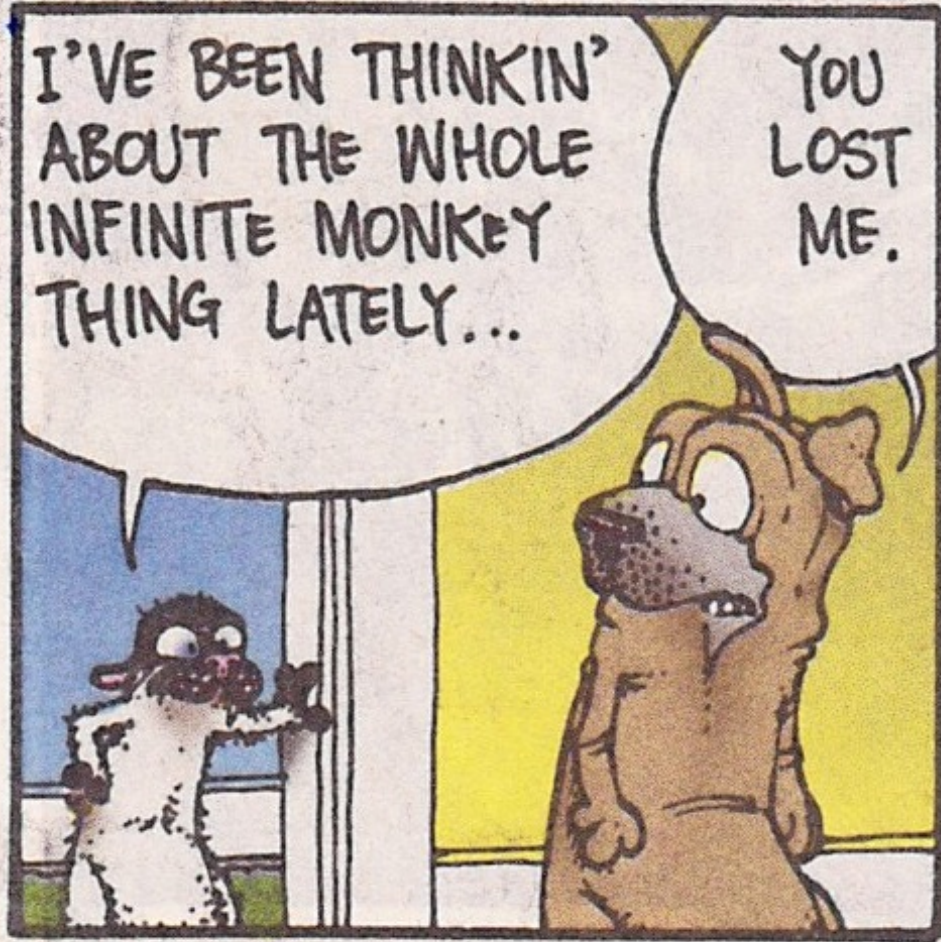
<http://bioconductor.org/packages/release/bioc/html/seqbias.html>

Prob/stats we've looked at is actually useful, giving you tools to understand contemporary research in CSE (and elsewhere).

I hope you enjoyed it!

*And One Last Bit of Probability Theory*

# GET FUZZY



© 2009 Darby Conley Dist. by UFS, Inc.



# by Darby Conley

WELL, THE WHOLE THEORY IS FLAWED. "INFINITE" IS TOO MANY MONKEYS. OVER 8 MONKEYS AND YOU'RE RUNNING INTO DISCIPLINE AND HYGIENE ISSUES.

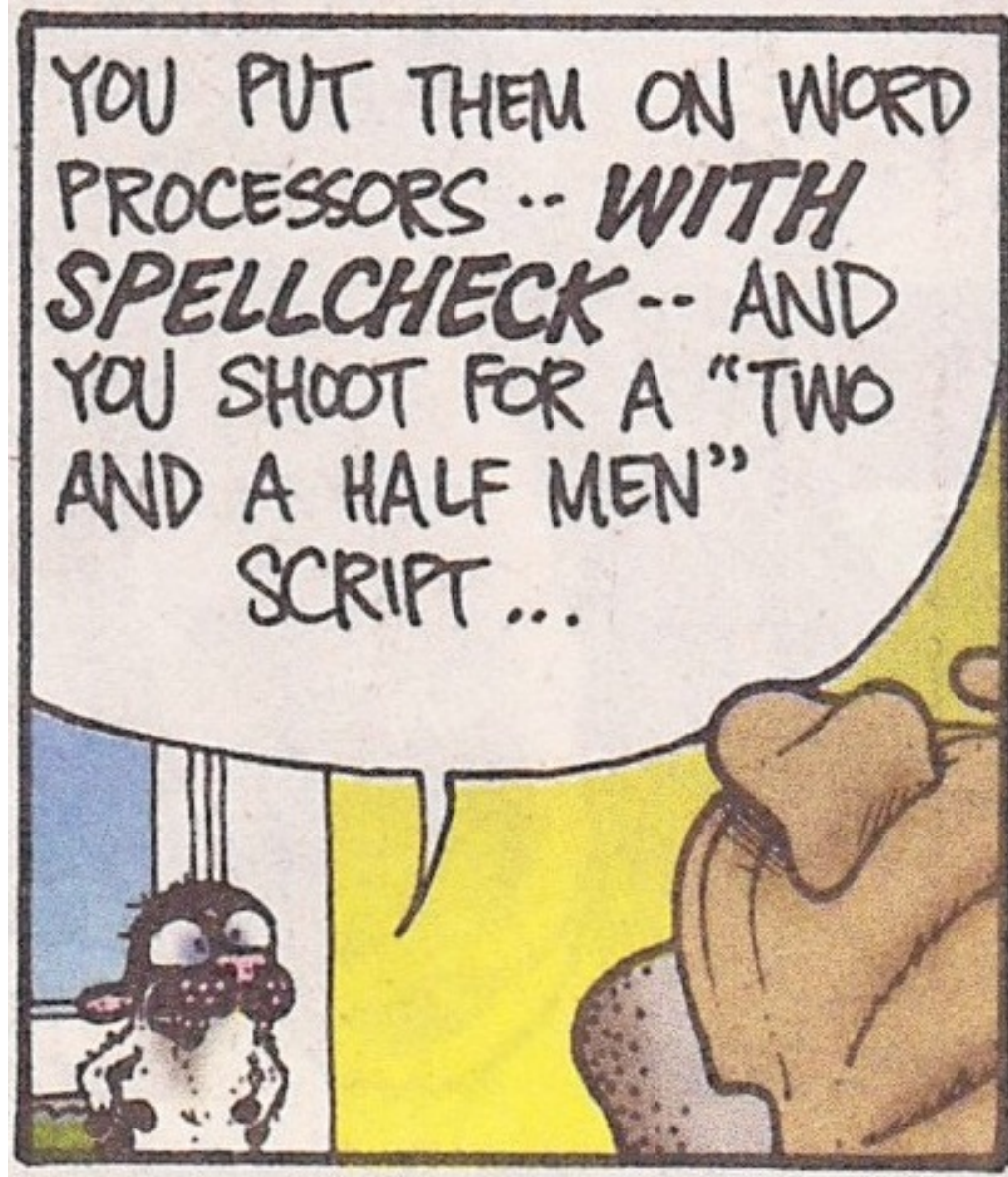


AND WHO'S GONNA READ INFINITE MONKEY SCRIPTS? SOME CHIMP COULD HAVE WRITTEN THE NEXT DA VINCI CODE, BUT *NEWSFLASH*: HE'S EATING THAT SCRIPT BEFORE YOU EVER SEE IT.













See also:

<http://mathforum.org/library/drmath/view/55871.html>

[http://en.wikipedia.org/wiki/Infinite\\_monkey\\_theorem](http://en.wikipedia.org/wiki/Infinite_monkey_theorem)