

Better than Chance: the importance of null models

Kevin Karplus

karplus@soe.ucsc.edu

Biomolecular Engineering Department
Undergraduate and Graduate Director, Bioinformatics
University of California, Santa Cruz



Outline of Talk

- 🐉 What is a null model (or null hypothesis) for?
- 🐉 Example 1: is a conserved ORF a protein?
- 🐉 Example 2: is residue-residue contact prediction better than chance?
- 🐉 Example 3: how should we remove composition biases in HMM searches?



Scoring (Bayesian view)

- 👉 Model M is a computable function that assigns a probability $P(A | M)$ to each sequence A .
- 👉 When given a sequence A , we want to know how likely the model is. That is, we want to compute something like $P(M | A)$.
- 👉 Bayes Rule:

$$P(M | A) = P(A | M) \frac{P(M)}{P(A)} .$$

- 👉 Problem: $P(A)$ and $P(M)$ are inherently unknowable.



Null models

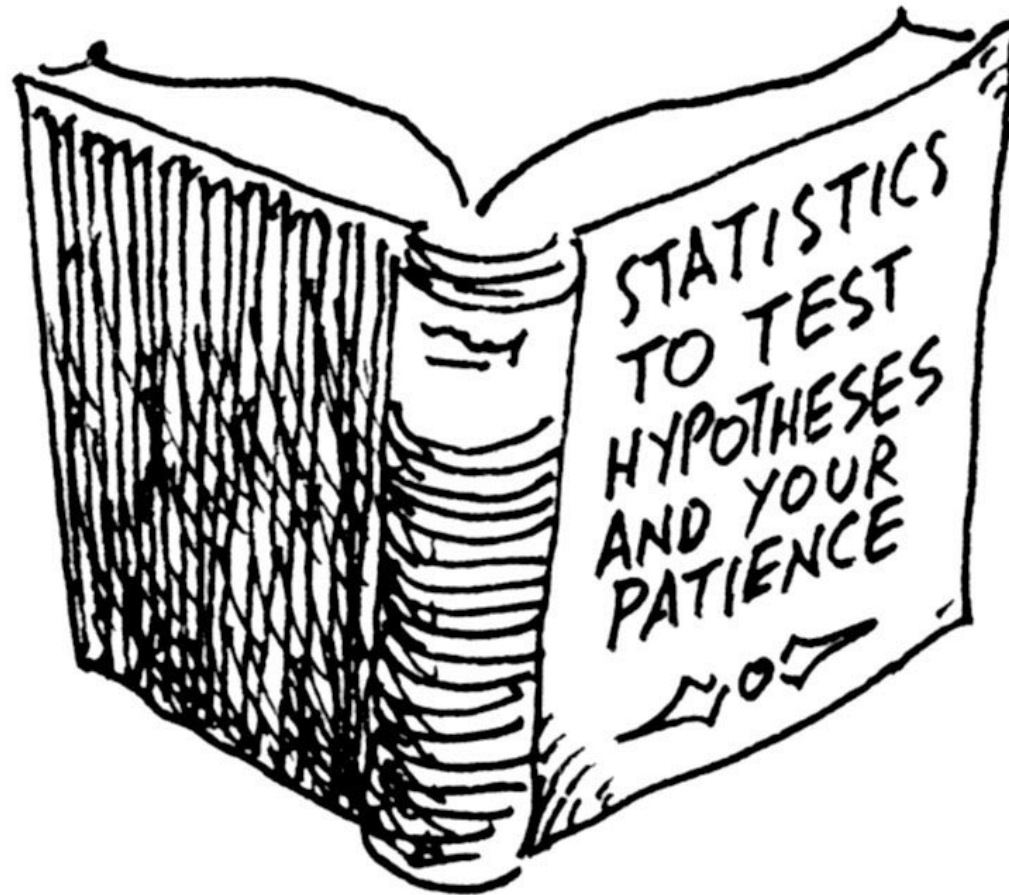
Standard solution: ask how much more likely M is than some *null hypothesis* (represented by a *null model*):

$$\frac{P(M | A)}{P(N | A)} = \frac{P(A|M)}{P(A|N)} \frac{P(M)}{P(N)} .$$

posterior odds likelihood ratio prior odds



Test your hypothesis



Thanks to Larry Gonick *The Cartoon Guide to Statistics*

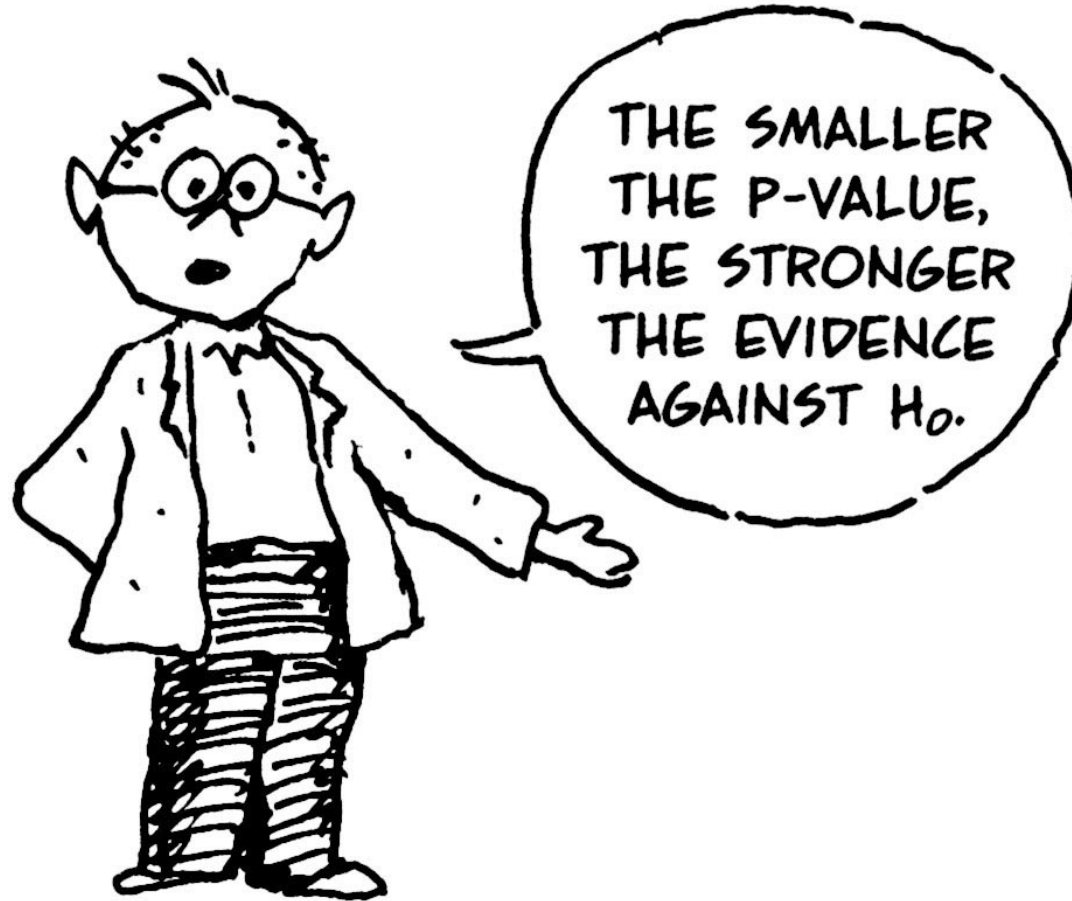


Scoring (frequentist view)

- 🐉 We believe in models when they give a large score to our observed data.
- 🐉 Statistical tests (p-values or E-values) quantify how often we should expect to see such good scores “by chance”.
- 🐉 These tests are based on a null model or null hypothesis.



Small p-value to reject null hypothesis



Thanks to Larry Gonick *The Cartoon Guide to Statistics*



Statistical Significance (2 approaches)

Markov's inequality For any scoring scheme that uses

$$\ln \frac{P(\text{seq} \mid M)}{P(\text{seq} \mid N)}$$

the probability of a score better than T is less than e^{-T} for sequences distributed according to N .

Parameter fitting For “random” sequences drawn from some distribution other than N , we can fit a parameterized family of distributions to scores from a random sample, then compute P and E values.



Null models

- 🐉 P-values (and E-values) often tell us nothing about how good our hypothesis is.
- 🐉 What they tell us is how bad our null model (null hypothesis) is at explaining the data.
- 🐉 A badly chosen null model can make a very wrong hypothesis look good.



Example 1: long ORF

- 🦖 A colleague found an ORF in an archæal genome that was 388 codons long and was wondering if it coded for a protein and what the protein's structure was.
- 🦖 We know that short ORFs can appear “by chance”.
- 🦖 So how likely is this ORF to be a chance event?



Null Model 1

🦖 DNA is undergoing no selection at all

🦖 G+C content bias. (GC is 36.7%, AT is 63.3%.)

🦖 Probability of stop codon

TAG= $0.3165 \times 0.3165 \times 0.1835 = 0.0184$, TGA=0.0184,
TAA=0.0317, so $p(\text{STOP})=0.0685$.

🦖 $P(\text{388 codons without stop}) = (1 - p(\text{STOP}))^{388} = 1.1e-12$

🦖 E-value in a 3 Megabase genome is about $3.3e-6$.

🦖 We can easily reject the null hypothesis!



Null Model 2

- 🐉 I forgot to tell you: this ORF is on the opposite strand of a known 560-codon ribosomal gene.
- 🐉 What is the probability of this long an ORF, on opposite strand of known gene?
- 🐉 Generative model: simulate random codons using the codon bias of the organism, take reverse complement, and see how often ORFs 388-long or longer appear.
- 🐉 Taking 100,000 samples, we get estimates of P-value in the range $3e-05$ to $6e-05$.
- 🐉 There are about 3000 genes, giving us an E-value of 0.09 to 0.18.



Null Model 3

- 🐉 We can do the same sort of simulation, but restrict the codons to ones that would code for exactly the same protein on the forward strand.
- 🐉 Now we get P-value of around 0.01 for long ORFs on the reverse strand of genes coding for this protein.



Protein or chance ORF?



Thanks to Larry Gonick *The Cartoon Guide to Statistics*



Not a protein

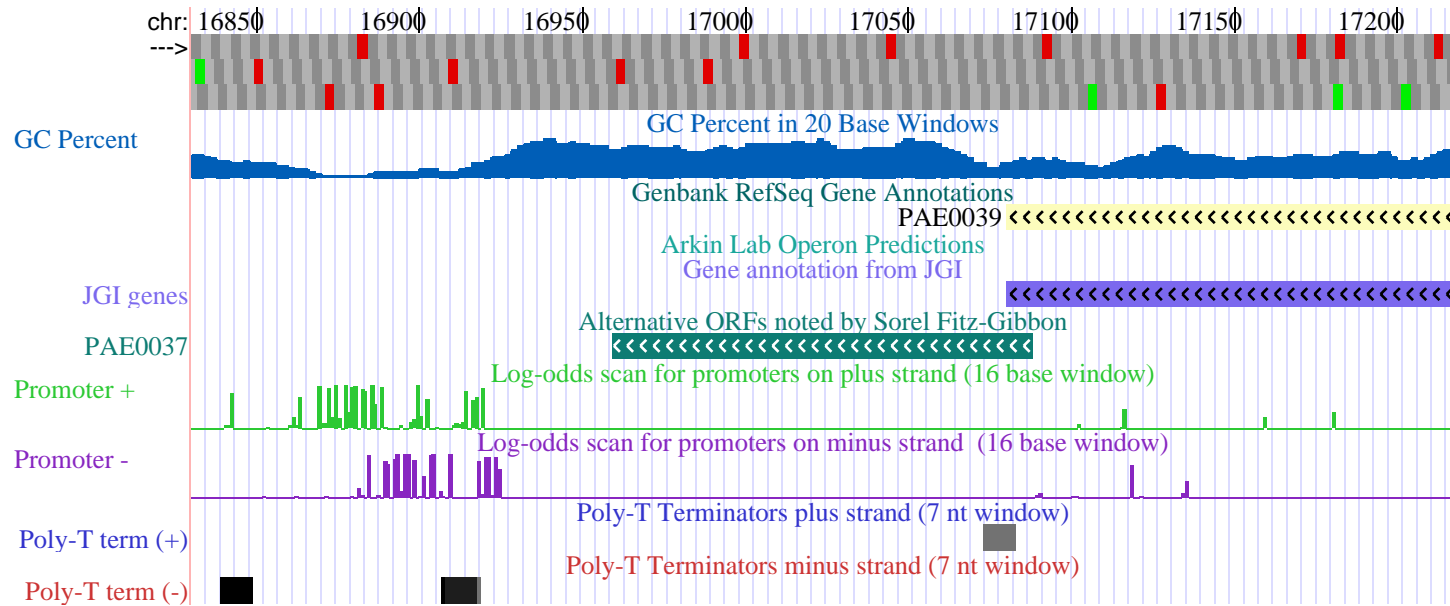
- + A tblastn search with the sequence revealed similar ORFs in many genomes.
- All are on opposite strand of homologs of same gene.
- “Homologs” found by tblastn often include stop codons.
- There is no evidence for a TATA box upstream of the ORF.
- No strong evidence for selection beyond that explained by known gene.

Conclusion: it is rather unlikely that this ORF encodes a protein.



Example 1b: another ORF

 pae0037: ORF, but probably not protein gene in *Pyrobaculum aerophilum*



 Promoter on wrong side of ORF.

 High GC content (need local, not global, null)

 Strong RNA secondary structure.



Example 2: contacts

- 🐉 Is residue-residue contact prediction better than chance?
- 🐉 Early predictors (1994) reported results that were 1.4 to 5.1 times “better than chance” on a sample of 11 proteins.
- 🐉 But they used a uniform null model:

$$P(\text{residue } i \text{ contacts residue } j) = \text{constant} .$$

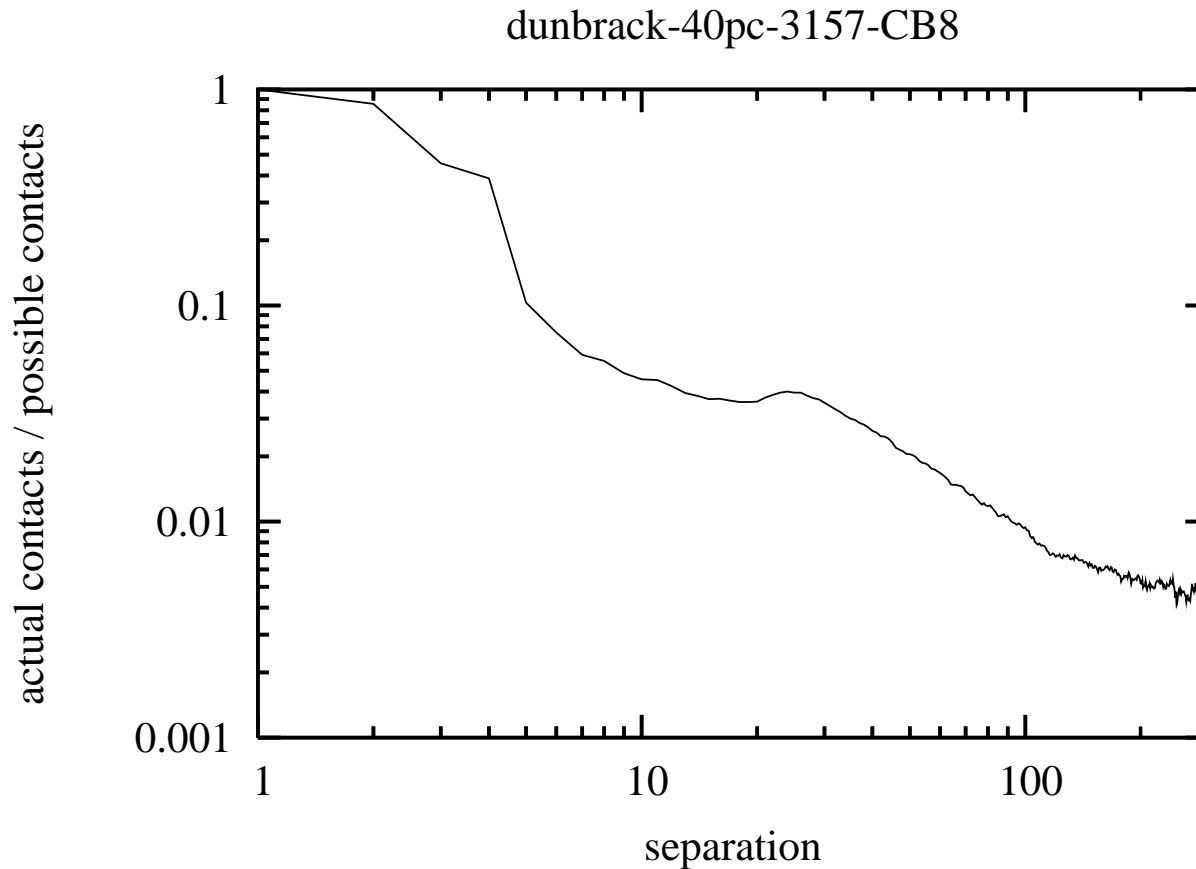
- 🐉 A better null model:

$$P(\text{residue } i \text{ contacts residue } j) = P(\text{contact} \mid \text{separation} = |i - j|) .$$



P(contact|separation)

Using CASP definition of contact, CB within 8 Å, CA for GLY.



Can get accuracy of 100%

- 🐉 By ignoring chain separations, the early predictors got what sounded like good accuracy (0.37–0.68 for L/5 predicted contacts)
- 🐉 But just predicting that i and $i + 1$ are in contact would have gotten accuracy of 1.0 for even more predictions.
- 🐉 More recent work has excluded small-separation pairs, with different authors choosing different thresholds.
- 🐉 CASP uses separation ≥ 6 , ≥ 12 , and ≥ 24 , with most focus on ≥ 24 .



Evaluating contact prediction

Two measures of contact prediction:

🐉 Accuracy:

$$\frac{\sum \chi(i, j)}{\sum 1}$$

🐉 Weighted accuracy:

$$\frac{\sum \frac{\chi(i, j)}{P(\text{contact} | \text{separation} = |i - j|)}}{\sum 1}$$

= 1 if predictions no better than chance, independent of separations for predicted pairs.



Separation as predictor

If we predict all pairs with given separation as in contact, we do much better than uniform model.

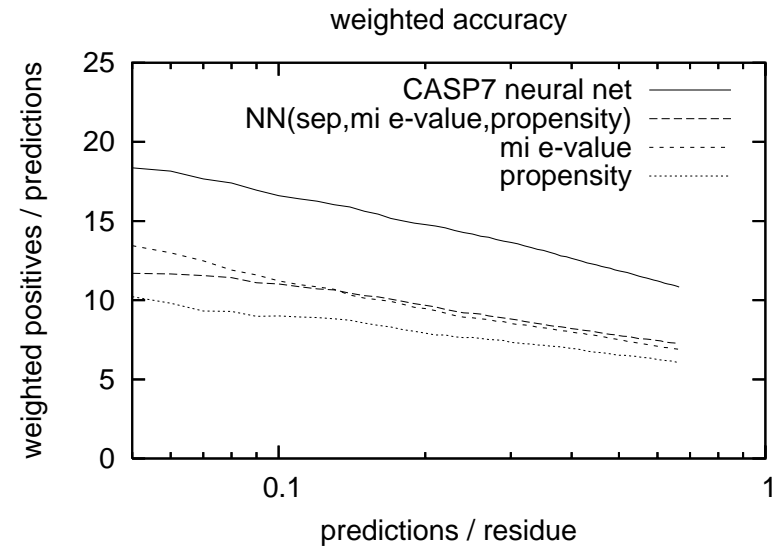
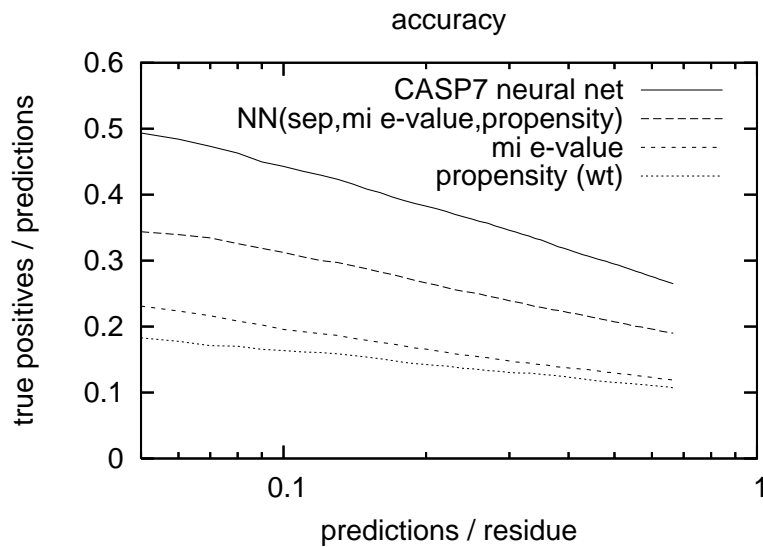
sep	$P(\text{contact} \mid i - j = \text{sep})$	$P(\text{contact} \mid i - j \geq \text{sep})$	“better than chance”
6	0.0751	0.0147	4.96
9	0.0486	0.0142	3.42
12	0.0424	0.0136	3.13
24	0.0400	0.0116	3.46



Now doing better

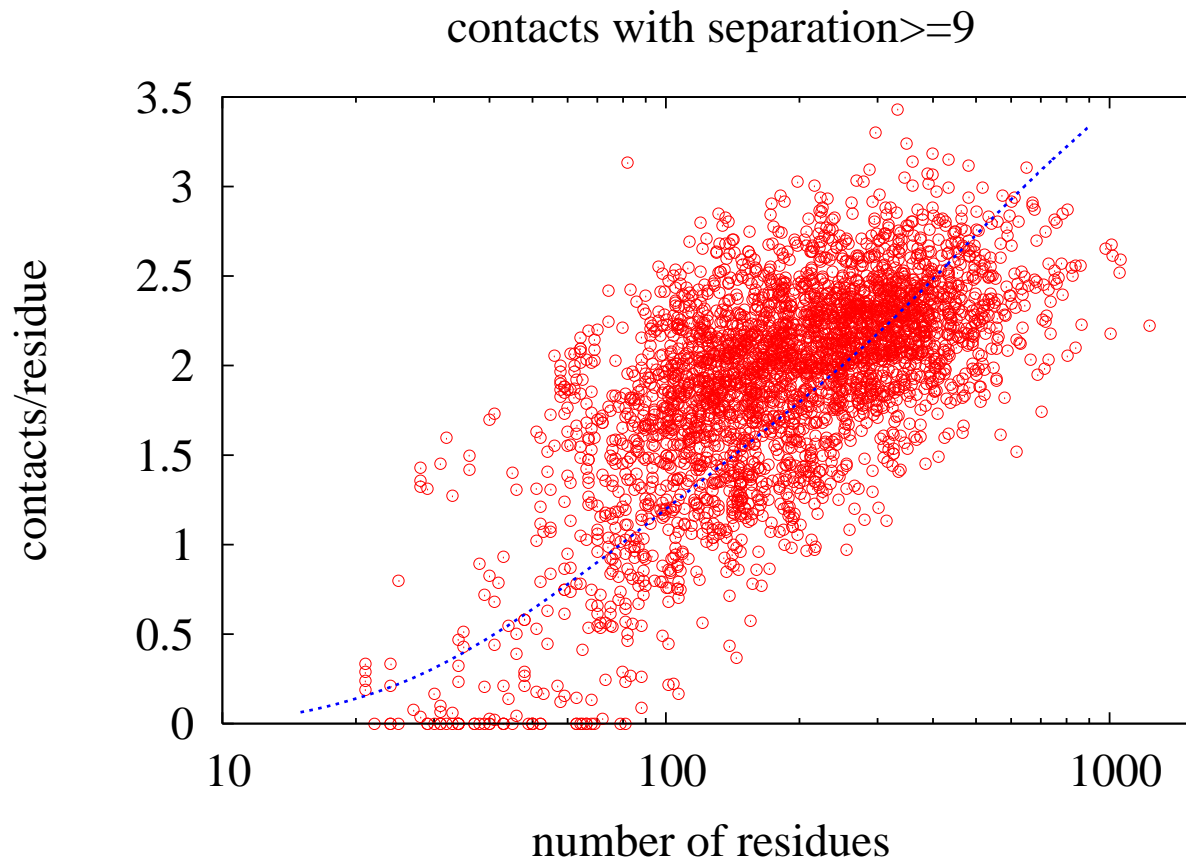
separation ≥ 9

contacts/residue taken separately for each protein



Contacts per residue

We can also use our null model to predict the number of contacts per residue (which is not a constant).



Example 3: HMM

- 🐉 Hidden Markov models assign a probability to each sequence in a protein family.
- 🐉 A common task is to choose which of several protein families (represented by different HMMs) a protein belongs to.



Standard Null Model

👉 Null model is an i.i.d (independent, identically distributed) model.

$$P(A \mid N, \text{len}(A)) = \prod_{i=1}^{\text{len}(A)} P(A_i) .$$

$$P(A \mid N) = P(\text{sequence of length } \text{len}(A)) \\ \prod_{i=1}^{\text{len}(A)} P(A_i) .$$



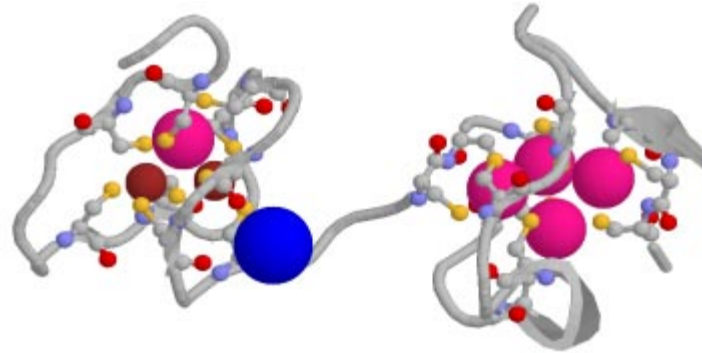
Composition as source of error

- 🦖 When using the standard null model, certain sequences and HMMs have anomalous behavior. Many of the problems are due to unusual composition—a large number of some usually rare amino acid.
- 🦖 For example, metallothionein, with 24 cysteines in only 61 total amino acids, scores well on any model with multiple highly conserved cysteines.

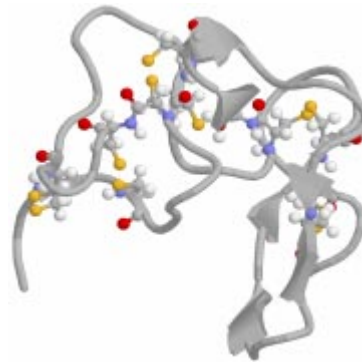


Composition examples

Metallothionein Isoform II (4mt2)

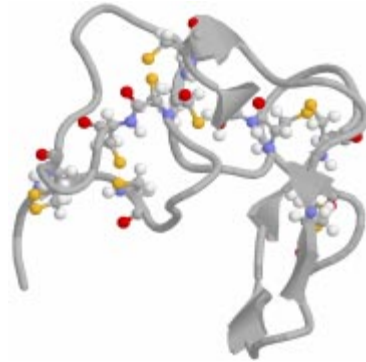


Kistrin (1kst)

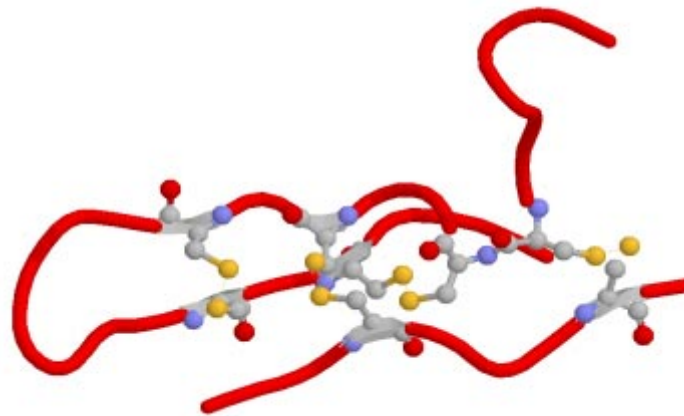


Composition examples

Kistrin (1kst)



Trypsin-binding domain of Bowman-Birk Inhibitor (1tabl)



Reversed model for null

- 👉 We avoid this (and several other problems) by using a reversed model M^r as the null model.
- 👉 The probability of a sequence in M^r is exactly the same as the probability of the reversal of the sequence given M .
- 👉 This method corrects for composition biases, length biases, and several subtler biases.



Helix examples

Tropomyosin (2tmaA)



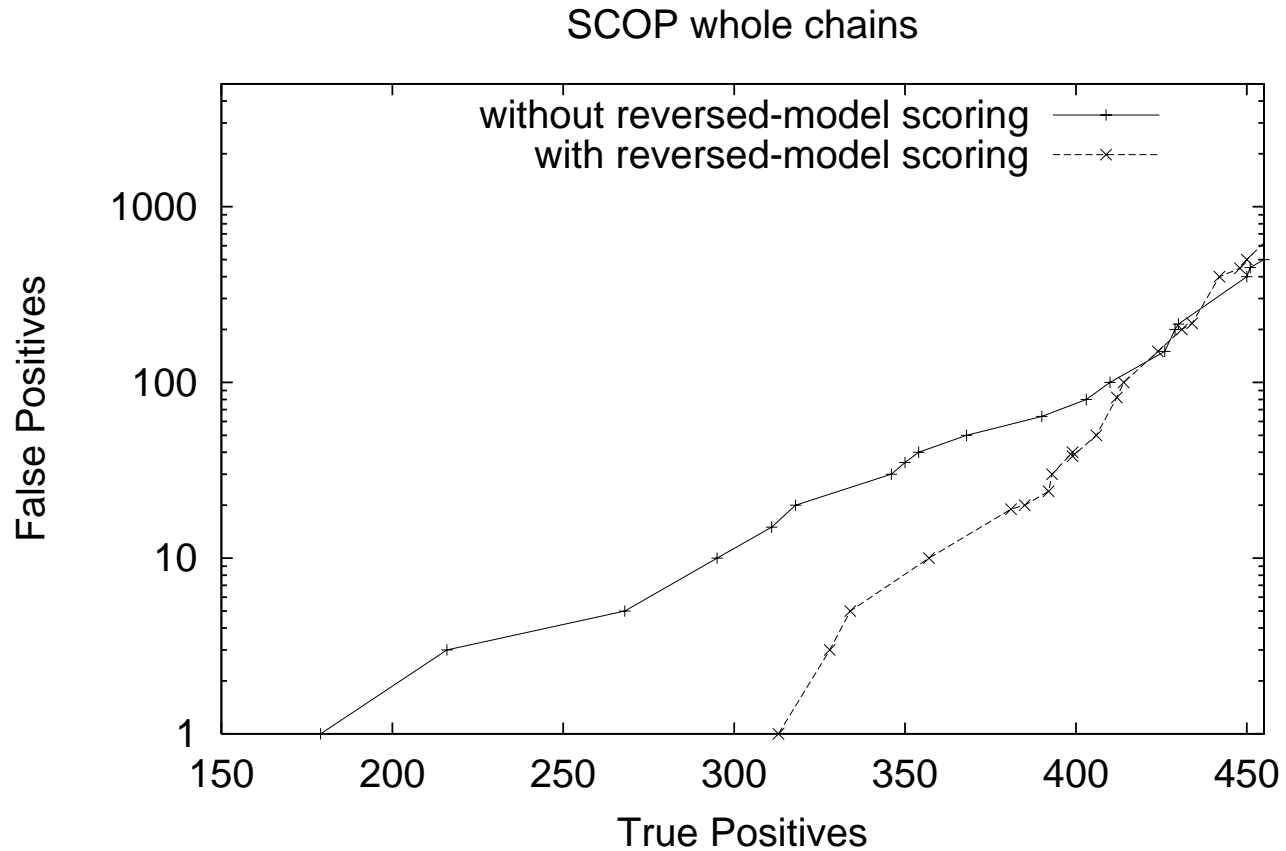
Colicin Ia (1cii)



Flavodoxin mutant (1vsgA)



Improvement from reversed model



Take-home messages

- 🐉 Base your null models on biologically meaningful null hypotheses, not just computationally convenient math.
- 🐉 Generative models and simulation can be useful for more complicated models.
- 🐉 Picking the right model remains more art than science.



Web sites

List of my papers: <http://www.soe.ucsc.edu/~karplus/papers/>

These slides: [http://www.soe.ucsc.edu/~karplus/papers/
better-than-chance-jul-07.pdf](http://www.soe.ucsc.edu/~karplus/papers/better-than-chance-jul-07.pdf)

Reverse-sequence null: Calibrating E-values for hidden Markov models with reverse-sequence null models. *Bioinformatics*, 2005. 21(22):4107–4115;
[doi:10.1093/bioinformatics/bti629](https://doi.org/10.1093/bioinformatics/bti629)

Archæal genome browser: <http://archaea.ucsc.edu>

UCSC bioinformatics (research and degree programs) info:

<http://www.soe.ucsc.edu/research/compbio/>

