Any questions?

- Blast
- Phylogenies
- Likelihood Ratio Tests

| Related Structures | | |
|--|--------|---------|
| | Score | E |
| Sequences producing significant alignments: | (bits) | Value |
| gi 604525 gb AAC37229.1 fertilization protein >gi 1097388 | 275 | 2e-73 |
| gi 604527 gb AAC37230.1 fertilization protein >gi 1097389 | 248 | 2e-65 |
| gi 604533 gb AAC37231.1 fertilization protein >gi 1097390 | 212 | 3e-54 |
| gi 3513683 gb AAC33930.1 fertilization protein precursor [| 94 | 8e-19 |
| gi 604529 gb AAC37232.1 fertilization protein >gi 1097391 | 82 | 5e-15 |
| gi 604531 gb AAC37233.1 fertilization protein >gi 1097392 | 67 | 1e-10 |
| gi 12084519 pdb 1GAK A Chain A, Crystal Structure Of Green | 62 | 3e-09 S |
| gi 538388 gb AAA21518.1 lysin | 40 | 0.022 |
| gi 4704776 gb AAD28265.1 sperm lysin precusor [Tegula fune | 36 | 0.25 |
| gi 602977 gb AAB59216.1 sperm lysin | 36 | 0.25 |
| gi 538406 gb AAB59168.1 lysin | 36 | 0.32 |
| gi 538400 gb AAB59167.1 lysin | 36 | 0.32 |
| di 602973 db 5557303 1 enerm lucin | 3.6 | 0.32 |

The current size of the NR protein database is 680,984,053 and the PDB is 3,816,875.

Which should you use in a Blast search to have the most sensitivity in detecting homologs with known 3D structures? Why? BlastX translates a DNA sequence into the 6 reading frames and searches each one against the protein NR database. BlastP takes a known translation and searches against the protein NR database.

Do you expect both searches to have the identical E-value? Why or why not?

Analysis of variation in the dn/ds ratio between sites or lineages

- Reading for today on website

 Review article of methods discussed today.
- Reading for Thursday on website

 Example of adaptive evolution using variation between lineages.

$\underline{d_N}/\underline{d_S}$ ratio (ω) estimated across all sites is inefficient at detecting positive selection

e.g. 3/44 (6.8%) sites subjected to positive selection

1 MSLAVLTFLVLCGFSFQHQAVGKWLTAAQKHPISGRMIRIRTKE 2 MSLAVLTFLVLCGWSFQHQAVGKWLTLAQKHPISGKMIRIRTKE

 $\texttt{3} \quad \texttt{MSLAVLTFLVLCGYSFQHQAVGKWLTSAQKHPISGHMIRIRTKE}$





This approach can work well if you know, a priori, where the functionally important parts of the gene are located (from functional or structural studies).

BUT, what if you don't have a prior functional or structural information?

| What is you | have | sequ | ienco | es fro | om n | nulti | ple s | pecie | es? |
|-------------|------|------|-------|--------|------|-------|-------|-------|-----|
| | | | | | | | | | |
| | | | | | | | | | |
| PINK | CGC | CAC | CGC | TTC | CGT | TTT | ATT | CCA | CAT |
| BLACK | | CGC | CGC | TAT | CAA | TTT | GTT | CAA | CAT |
| RED | | CGC | AGC | TGG | CAT | TAT | GTT | GAA | CCC |
| WHITE | | CGC | CGC | TGG | CAT | TAT | GTT | CCA | CCC |
| PINTO | | CGC | AGC | TGG | ACT | TAT | GTT | CAA | CCC |
| FLAT | | CGC | CGC | TGG | AAT | TTT | GTT | ACA | CCC |
| GREEN | | CGC | CGA | TGG | ACT | TTT | GTT | CGA | TAT |
| | | | | | | | | | |

Markov models of codon evolution

Goldman & Yang 1994 MBE 11:725-736

Muse & Gaut 1994 MBE 11:715-724



Why use a likelihood model of codon evolution?

- 1. We can take advantage of the phylogeny
- 2. Computation of transition probabilities accomplishes the following in 1 step:
 - i. estimation of parameters (t, κ , ω)
 - ii. correction for multiple hits
 - iii. weight evolutionary pathways between codons

Codon models

Important parameters:

- Transition/transversion rate ratio: κ
- Biased codon usage: π_i for codon j
- Nonsynonymous/synonymous rate ratio: $\omega = d_N/d_S$















Pairwise comparisons

- Calculate dn/ds from data
- Determine if value is significantly > 1.

For pairwise comparisons, we must determine if d_N/d_S is *significantly* greater than 1

- Estimate d_N/d_S ratio for pairwise comparison.
- Estimate d_N/d_S ratio with d_N/d_S fixed at one.
- Compare likelihoods using likelihood ratio test.

Example pairwise estimates

Estimating the dn/ds ratio: H__sorenseni ... 1 H__rufescens lnL = -567.039906 t= 0.1406 S= 108.2 N= 257.8 dN/dS= 6.0570 dN= 0.0622 dS= 0.0103

Fixing the dn/ds ratio equal to one:

 $\begin{array}{ll} H._sorenseni) ... 1 \ H._rufescens \\ lnL = -569.395789 \\ t= 0.1390 \ S= \ 105.2 \ N= \ 260.8 \ dN/dS= 1.0000 \ dN= 0.0463 \ dS= 0.0463 \end{array}$

Likelihood ratio test:

Example pairwise estimates

Estimating the dnlds ratio: H_sorenseni ... 1 H_rufescens InL = -567.039906 t = 0.1406 S= 108.2 N= 257.8 dN/dS=6.0570 dN=0.0622 dS=0.0103

Fixing the dn/ds ratio equal to one: H._sorenseni) ... 1 H._rufescens InL = -569.395789 t= 0.1390 S= 105.2 N= 260.8 dN/dS= 1.0000 dN= 0.0463 dS= 0.0463

Likelihood ratio test: -2[-569.4 - (-567.0)] = 4.8 + selection

Example pairwise estimates

Estimating the dn/ds ratio: H._walallensis ... 2 H._sorenseni lnL = -568.604732

t= 0.1292 S= 95.2 N= 270.8 dN/dS= 1.2521 dN= 0.0454 dS= 0.0363 *Fixing the dn/ds ratio equal to one:* H._walallensis ... 2 H._sorenseni

H._walallensis ... 2 H._sorenseni lnL = -568.661120

 $t{=}\;0.1292\;\;S{=}{-}94.5\;\;N{=}\;\;271.5\;\;dN/dS{=}\;1.0000\;\;dN{=}\;0.0431\;\;dS{=}\;0.0431$

Likelihood ratio test:

Example pairwise estimates

Estimating the dn/ds ratio: H._walallensis ... 2 H._sorenseni lnL = -568.604732

t= 0.1292 S= 95.2 N= 270.8 dN/dS= 1.2521 dN= 0.0454 dS= 0.0363 *Fixing the dn/ds ratio equal to one:* H__walallensis ... 2 H__sorenseni lnL = -568.661120

 $t{=}\;0.1292\;\;S{=}\;\;94.5\;\;N{=}\;\;271.5\;\;dN/dS{=}\;1.0000\;\;dN{=}\;0.0431\;\;dS{=}\;0.0431$

Likelihood ratio test: -2[-568.6 - (-568.7)] = 0.2

Can not reject null model = neutral evolution.

Problem: averaging over a pair

In a pairwise analysis we must average the $\boldsymbol{\omega}$ ratio over:

- 1. all sites
- 2. the entire evolutionary history

In a large-scale pairwise database search, only 17 out of 3,595 genes were found to be under positive selection, at <0.5% (Endo *et al.* 1996 *MBE* **13**: 685-690)

Problem: averaging over a pair has very low power if the questions are about "when" or "where"!

Solution: Phylogenetic estimation of selection pressure

- variable ω over branches (when?)
- variable ω over sites (what fraction?)

Statistical methods (maximum likelihood) have been developed to test models of positive selection using genes. If a selection model fits the data better than a neutral model, one can identify sites subjected to positive selection and infer functional importance.



- Model different classes of codon sites – Purifying, neutral, positive selection.
- Compare selection model with one class that has dn/ds > 1, to neutral model where all classes have dn/ds < 1 using likelihood ratio test.





| Model | Code | NP | Parameters |
|----------------|------|------|---|
| One-ratio | M0 | 1 | ω |
| Neutral | M1 | 1 | p_0 |
| Selection | M2 | 3 | p_{0}, p_{1}, ω_{2} |
| Discrete | M3 | 2K-1 | $p_{0}, p_{1, \dots, p_{K-2}}$ |
| | | | $\omega_0, \omega_1, \ldots, \omega_{K-2}$ |
| Frequency | M4 | 5 | $p_{0}, p_{1},, p_{4}$ |
| Gamma | M5 | 2 | α, β |
| 2Gamma | M6 | 4 | $p_{0}, \alpha_{0}, \beta_{0}, \alpha_{1}$ |
| Beta | M7 | 2 | p, q |
| Beta& <i>w</i> | M8 | 4 | p_{0}, p, q, ω |
| Betaγ | M9 | 5 | $p_{0}, p, q, \alpha, \beta$ |
| Beta&normal+1 | M10 | 5 | p_0, p, q, α, β |
| 3eta&normal>1 | M11 | 5 | p_{0}, p, q, μ, σ |
|)&2normal>1 | M12 | 5 | $p_{0}, p_{1}, \mu_{2}, \sigma_{1}, \sigma_{2}$ |
| 3normal>0 | M13 | 6 | $v_0 v_1, \mu, \sigma_0, \sigma_1, \sigma_2$ |







- 6 alleles of human class I MHC aligned - Known to be subjected to positive selection
- · Calculated dn/ds across all sites
- Test for variation in the dn/ds ratio between sites - Identification of sites under selection







150 – 200 mya

35 mva

Ϋ́́A

ß alobin aene cluster

40 – 80 mya

δ β Chrom, 11







Problems with lineage analysis

- Estimating dn/ds across all sites – There are "branch - site" models
- Many branches increase the degrees of freedom

 Can limit estimates to particular branches, such as
 human lineage or following duplication

Summary

- 1. Pairwise methods have very low power to detect adaptive evolution.
- 2. Branch models allow variation among branches but assume one ω for all sites, and have low power to detect positive selection.
- 3. Site models assume allow variation among sites but assume selection pressure does not change among branches, and will have higher power if positive selection is long term
- 1. Why is averaging the dn/ds ratio across all sites in a gene inefficient at detecting adaptive evolution
- What is one way to detect variation in rates of evolution along lineages.
 Model 7 is a neutral model of DNA evolution estimates for a beta distribution with the parameters p and q. Model 8 add an additional proportion of sites with a dn/ds ratio estimated from the data. What are the degrees of freedom between
- these two models?4. Describe a selective and neutral model that would be a useful for testing adaptive evolution by a likelihood ratio tests. Include what parameters are different between the models.
- 5. How do you determine the degrees of freedom for a likelihood ratio test?6. What are some important factors to consider when calculating dn/ds ratios?

PAML demo