

University of Oxford, Department of Zoology Evolutionary Biology Group Department of Zoology University of Oxford South Parks Road Oxford OX1 3PS, U.K. Fax: +44 1865 271249

Evolutionary Analysis of Viral Genomes

Lecture 4: Natural Selection & Viral Adaptation

Oliver G. Pybus Department of Zoology, University of Oxford

http://evolve.zoo.ox.ac.uk



Natural Selection

"This preservation of favourable variations and the rejection of injurious variations, I call natural selection. Variations neither useful nor injurious would not be affected by natural selection, and would be a left a fluctuating element" (Darwin, On The Origin of Species).

The **selection coefficient** (*s*) of a mutation measures the degree by which it increases or decreases the fitness of an organism.

Mutations are classified as:
advantageous / benficial (s>0)
disadvantageous / deleterious (s<0)
neutral (s=0)

The Fate of Mutations (Simplified)



Fixation = Nucleotide Substitution

The Fate of Mutations (Simplified)



Distribution of Selection Coefficients Selectionist Model Neutralist Model s<0 s>0 Most *fixed* mutations Most *fixed* mutations s≈0 are advantageous are neutral

A Model Of Substitution Rate

 \bigcirc *m* = mutation rate per site per individual

N = population size

f(s) = distribution of selection coefficients

= probability that a given mutation has selection coefficient *s* p(s) = probability of fixation of mutations with selection coefficient *s* ~ (1-exp(-2s)) / (1-exp(-2Ns)) in a constant-size population r(s) = rate at which mutations with selection coefficient *s* are fixed = N m f(s) p(s)

• $\mu = \text{total substitution rate}$ $\mu = \int r(s).ds$ $= Nm \int f(s) p(s).ds$

A Model Of Substitution Rate

•For neutral mutations, *s*=0

p(0) = 1/N , $r(0) \approx mNf(0)/N = mf(0)$

For advantageous mutations, s>0

 $p(s>0) \approx 2s$, $r(s>0) \approx 2Nmsf(s>0)$

For disadvantageous mutations, s<0</p>

 $p(s < 0) \approx 0$ hence $r \approx 0$

dN/dS

Are fixed mutations neutral or advantageous?

Silent mutations are likely to be neutral.

Replacement changes could be neutral, beneficial or deleterious.

 $\bigcirc dN/dS = \omega$ = ratio of replacement and silent substitution rates.

dN/dS = 1 if all replacement changes are neutral dN/dS < 1 if some replacement changes are deleterious dN/dS = 0 if all replacement changes are deleterious dN/dS > 1 only if some replacement changes are beneficial

When are silent changes not neutral?

• Overlapping genes and reading frames. Very common in viruses (e.g. Hepatitis B Virus)

Secondary RNA / DNA structure

e.g. stem-loop structures

Codons for the same amino acid differ in fitness May result from different translational efficiencies. Gives rise to codon bias.

Outranslated or regulatory regions

Detecting Selection Using dN/dS

Pairwise method: Calculates *dN/dS* from each pair of sequences, then averages *dN/dS* across all pairs.

e.g. Nei & Gojobori 1986 (uses Jukes-Cantor model)

Parsimony method: Start with a phylogeny. A parsimony algorithm then estimates the ancestral codons at the internal nodes. dN/dS is then calculated from the tree with specified ancestral codons.

e.g. ADAPTSITE (Suzuki & Gojobori 1999) http://mep.bio.psu.edu

Substitution model method: Starts with a phylogeny. Uses a codon-based substitution model that contains *dN/dS* as a parameter. Allows *dN/dS* to vary among codons and can identify individual codons under selection. Uses realistic models of substitution.

e.g. CODEML (Yang et al. 2000) http://abacus.gene.ucl.ac.uk/

Codon Substitution Model

• Define a 61×61 substitution matrix (**Q**). Each element in the matrix $q_{i,j}$ defines the substitution rate of codon *i* to codon *j*.

| <i>q</i> _{<i>ij</i>} = { | 0, | if <i>i</i> and <i>j</i> differ at more than one position |
|-----------------------------------|-------------------|---|
| | π_j , | forsynonymous transversion |
| | κπ _j , | forsynonymous transition |
| | $\omega \pi_j$, | for nonsynonymous transversion |
| | ωκ $π_j$, | for nonsynonymous transition |

• κ is the transition/transversion rate ratio • ω is the silent/replacement rate ratio • π_i is the frequency of codon *j*

Codon Substitution Models

Codon models are calculated like nucleotide models.

• The genetic distance between codon *i* and codon *j* is given by the distribution $P_{i,j}(t) = \exp(\mathbf{Q}t)$.

This probability is multiplied across all codons and across all branches in the tree. As before, it is then integrated over all possible ancestral codons at internal nodes.

• This gives the final phylogenetic probability: *P*(*seqs*|**T**,**B**,**Q**)

If ω varies among codons, then the final probability is obtained by integrating across the specified distribution of ω among sites.

Disadvantages of dN/dS methods

Assumes all changes are the result of past fixation events. Caution is required if genetic variation represents non-fixed polymorphism (e.g. within-patient viral data sets).

Conservative. Will fail to detect many cases of selection, especially if a site has only undergone a single selective sweep.

Such methods therefore tend to identify sites that are repeatedly positively selected.

Assumes substitutions are independent at each codon (no epistasis or clonal interference).

| Gene | Organism | Refs | Gene | Organism | Refs | | | |
|--|--|--------------------------------|---|--------------------------|------|--|--|--|
| Genes involved in defensive systems of | r immunity | Genes involved in reproduction | | | | | | |
| Class I chitinase gene | Arabis and Arabidopsis | 41 | 18-kDa fertilization protein gene | Abalone (Haliotis) | 61 | | | |
| Colicin genes | Escherichia coli | 45 | Acp26Aa | Drosophila | 62 | | | |
| Defensin genes | Rodents | 46 | Androgen-binding protein gene | Rodents | 63 | | | |
| Fv1 | Mus | 47 | Bindin gene | Echinometra | 64 | | | |
| Immunoglobulin V _H genes | Mammals | 48 | Egg-laying hormone genes | Aplysia californica | 3 | | | |
| MHC genes | Mammals | 49 | Ods homeobox gene | Drosophila | 65 | | | |
| Polygalacturonase inhibitor | Legume and dicots | 50 | Pem homeobox gene | Rodents | 66 | | | |
| genes | | | Protamine P1 gene | Primates | 67 | | | |
| RH blood group and RH50 | Primates and rodents | 51 | Sperm lysin gene | Abalone (Haliotis) | 61 | | | |
| genes | | | S-Rnase gene | Rosaceae | 68 | | | |
| Ribonuclease genes | Primates | 52 | Sry gene | Primates | 69 | | | |
| Transferrin gene | Salmonid fishes | 53 | Course investigation discretion | | | | | |
| Type I interferon- v gene | Mammals | 54 | Genes involvea in aigestion | Devide | 70 | | | |
| a 1-Proteinase inhibitor genes | Rodents | 55 | K-casein gene | Bovids | 70 | | | |
| Concentration and in a defension of | | | Lysozyme gene | Primates | 23 | | | |
| Genes involved in evaling dejensive s | Siems or immunity | Toxin protein genes | | | | | | |
| CED TDAD MEA 2 and DE92 | FIVID VIIUS | 42 56 | Conotoxin genes | Conus gastropods | 71 | | | |
| CSP, IRAP, MSA-2 and PF83 | Plasmodium laiciparum | 50 57 | Phospholipase A $_2$ gene | Crotalinae snakes | 72 | | | |
| Delta-antigen coding region | Hepatitis D virus Dhagaa $C4 \neq V174$ and | 2/ | | 1/ | | | | |
| E gene | Phages $G4$, $IXI/4$, and | 3 | Genes related to electron transport and | a/or AIP synthesis | 2 | | | |
| | 515 | 40 | ATP synthase F $_0$ subunit gene | | 3 | | | |
| Envelope gene | HIV | 40 | COX/A isoform genes | Primates | 13 | | | |
| gH glycoprotein gene | Pseudorables virus | 3 | COX4 gene | Primates | /4 | | | |
| Hemagglutinin gene | Human influenza A virus | 33 | Cytokine genes | | | | | |
| Invasion plasmid antigen genes | Shigella | 3 | Granulocyte-macrophage SF gene | Rodents | 75 | | | |
| Merozoite surface antigen-1 gene | Plasmodium falciparum | 58 | Interleukin-3 gene | Primates | 75 | | | |
| msp 1 a | Anaplasma marginale | 3 | Interleukin-4 gene | Rodents | 75 | | | |
| net | HIV | 38 | 0 | | | | | |
| Outer membrane protein gene Chlamydia | | 3 | Miscellaneous | | | | | |
| Polygalacturonase genes | Fungal pathogens | 50 | CDC6 | Saccharomyces cerevisiae | 3 | | | |
| Porin protein 1 gene | Neisseria | 59 | Growth hormone gene | Vertebrates | 76 | | | |
| S and HE glycoprotein genes | Murine coronavirus | 60 | Hemoglobin <i>b</i> -chain gene | Antarctic fishes | 77 | | | |
| Sigma-1 protein gene | Reovirus | 3 | Jingwei | Drosophila | 78 | | | |
| Virulence determinant gene | Yersinia | 3 | Prostatein peptide C3 gene | Rat | 3 | | | |

Table 1. Selected examples of protein-coding genes in which positive selection was detected by using the d_N/d_S ratio



varies among genes (HIV-1 group O)

dN/dS

Example:

Example: *dN/dS* varies among species



dN/dS varies among codons



High *dN/dS* found in codons that form the antigen-recognition site

(Yang & Swanson. 2002. Mol. Biol. Evol. 19:49-57)

dN/dS varies among codons



(Sheridan, Pybus, Holmes & Klenerman. 2004. J. Virol. 78:3447-3454)

Mutation Site Frequencies



Mutation Site Frequencies

The distribution of site frequencies (the site-frequency spectrum) in a sample of sequences contains information about natural selection.

Can be analyzed using Poisson Random Field (PRF) models.

PRF models assume that polymorphic sites are unlinked (effectively on different chromosomes).

PRF models are more suitable for human genome analyses than for pathogen genome analyses.

McDonald-Kreitman Method

study species (main alignment)

sister species (outgroup alignment) Circles represent *fixation* between species. Diamonds represent *polymorphism* within species.

If there is no selection, then both fixation and polymorphism are driven by genetic drift.

The ratio of silent/replacement fixations should equal the ratio of silent/non-replacement polymorphisms.

A statistical deviation from equality suggests natural selection.

Assumes all polymorphism is neutral.

Estimating Viral Adaptation Rate

later sample (main alignment)

earlier sample (outgroup

alignment)

Can also be applied to pathogen sequences sampled through time.

 $s_{<1} = #$ silent polymorphisms (\diamondsuit)

 $r_{<1} = #$ replacement polymorphisms (\blacklozenge)

 $s_1 = #$ silent fixations (\bigcirc)

 $r_1 = #$ replacement fixations (\bigcirc)

 $a_1 = #$ positively selected fixations

If there is no selection, then, on average: $(r_1 - a_1)/r_{<1} = s_1/s_{<1}$

Rearranging, we can estimate a₁ using:

$$a_1 = r_1 \left[1 - \left(\frac{s_1}{r_1} \right) z \right], \text{ where } z = \left(\frac{r_{<1}}{s_{<1}} \right).$$

Estimating Viral Adaptation Rate

In virus populations, mutations that are highly-frequent but not yet fixed may also be adaptations.

 $s_{<0.5} = #$ silent polymorphisms with a frequency < 0.5

 $r_{<0.5}$ = # replacement polymorphisms with a frequency < 0.5

 $s_{>0.5} = #$ silent polymorphisms with a frequency > 0.5

 $r_{>0.5} = \#$ replacement polymorphisms with a frequency > 0.5

 $a_{>0.5} = #$ positively selected polymorphisms with a frequency > 0.5

• We can now estimate both a_1 and $a >_{0.5}$ using:

$$a_{1} = r_{1} \left[1 - \left(\frac{s_{1}}{r_{1}} \right) z \right], \text{ where } z = \left(\frac{r_{<0.5}}{s_{<0.5}} \right) \qquad a_{>0.5} = r_{>0.5} \left[1 - \left(\frac{s_{>0.5}}{r_{>0.5}} \right) z \right], \text{ where } z = \left(\frac{r_{<0.5}}{s_{<0.5}} \right).$$

This approach assumes that only polymorphisms with a frequency<0.5 are neutral.</p>

Accumulation of adaptations during HIV infection



(Williamson. 2003. Mol. Biol. Evol. 20:1318-1325)