

Sliding Through Phylogenetics

Daisy Shepherd

The University of Auckland

dshe078@aucklanduni.ac.nz

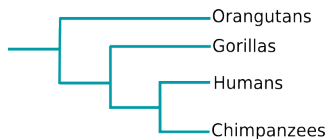
01/12/2015



- 1 Introduction & Motivation
- 2 Study Design
- 3 Results
- 4 Conclusions & Future Work

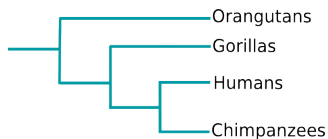
What is Phylogenetics?

- All organisms have DNA.
- Map the differences in DNA.
- How closely related are these groups?
- **Aim:** Derive their evolutionary history.



What is Phylogenetics?

- All organisms have DNA.
- Map the differences in DNA.
- How closely related are these groups?
- **Aim:** Derive their evolutionary history.



The Problem

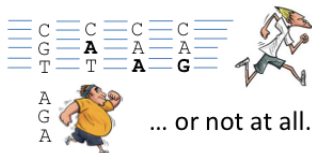
- Sites evolve over time in a number of ways.
- Usually assume each site evolves at a similar rate.



However, this is not always the case...

Some sites change *quickly*...

... whilst others change slowly...



Sites which have the same rate display rate **homogeneity**.



Sites which have differing rates display rate **heterogeneity**.



How do we model these?



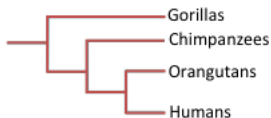
Fixed, constant rate



Use γ distribution
 γ = rate parameter

So why is this a problem?

- Models need to detect and handle this heterogeneity.
- Poorly fitting models lead to poor estimation of evolutionary relationships.



The Current Approach

- When selecting a model (homogeneous vs. heterogeneous), usually look at the whole alignment.
- Heterogeneous model - one estimate for the rate parameter γ .

```
Species 1: ACTACGTACGAGATTAGTGTACGATCA...
Species 2: CTGAGATCGCGTCATAGAGATGTCAGTT...
Species 3: ACACTGACGTACGTAGATCATGTGCACT...
Species 4: AATTCGGTGAACTCATGTTTCGTGCTA...
Species 5: ACGATATGCCCTCGATCGCCGCCTCCGC...
Species 6: GTACACACATGATAATGTGTACACTGTG...
```

Complete alignment analysis

(Traditional approach)
Single analysis

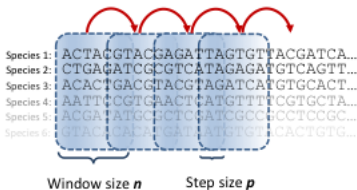
- ✓ Simple and quick
- ✓ Common practice, so widely applied in software
- ✗ **Generalises behaviour across all sites**

The Proposed Solution



How does it work?

- 1 Look at the first n sites.
- 2 Fit a model to only this window of sites.
- 3 Slide window along p sites, to a new group of n sites.
- 4 Fit a model to this new window.
- 5 Repeat until entire alignment has been covered.



Sliding Window analysis

(Our proposed approach)
Multiple analyses

What did we want to do?

Species 1: ACTACGTACGAGATTAGTGTACGATCA...
Species 2: CTGAGATCGCGTCATAGAGATGTCAGTT...
Species 3: ACACTGACGTACGTAGATCATGTGCACT...
Species 4: AATTCCGTGAACACTCATGTTTTCGTGCTA...
Species 5: ACGATATGCGCTCGATCGCCGCTCCGC...
Species 6: GTACACACATGATAATGTGTACACTGTG...

Complete alignment analysis

(Traditional approach)
Single analysis

Species 1: ACTACGTACGAGATTAGTGTACGATCA...
Species 2: CTGAGATCGCGTCATAGAGATGTCAGTT...
Species 3: ACACTGACGTACGTAGATCATGTGCACT...
Species 4: AATTCCGTGAACACTCATGTTTTCGTGCTA...
Species 5: ACGATATGCGCTCGATCGCCGCTCCGC...
Species 6: GTACACACATGATAATGTGTACACTGTG...

Sliding Window analysis

(Our proposed approach)
Multiple analyses

Window size n Step size p

Aim: Test whether the Sliding Window approach improves our ability to detect variation in evolutionary rates.

1. The Data:

Simulated 300 alignments (5000 sites)



- 10, 50 and 100 taxa
- Random γ parameter values used for hetero model
- Generated random topologies
- Random insertion point (for hetero region into alignment)
- Window size $n = 500$ sites
- Step size $p = 50$ sites

2. Testing:

- i. Perform a complete alignment analysis.
- ii. Implement the Sliding Window (SW) approach.
- iii. Compare the results.

3. What are we looking for?

- Are heterogeneous or homogeneous rates detected?
- Accurate estimate of the rate?
- Is the SW approach an improvement?

Results

Complete Alignment Analysis

Based on p -values from the likelihood ratio test:

| Taxa Number | n | Accept H_0 | Reject H_0 |
|---|-----|--------------|--------------|
| 10 | 100 | 25 | 75 |
| 50 | 100 | 20 | 80 |
| 100 | 100 | 0 | 100 |
| H_0 : Homogeneous model is the true model | | | |

Table : The number of simulations which accepted and rejected the null hypothesis under the complete alignment analysis.

Did not detect varying rates in 15% of the alignments.

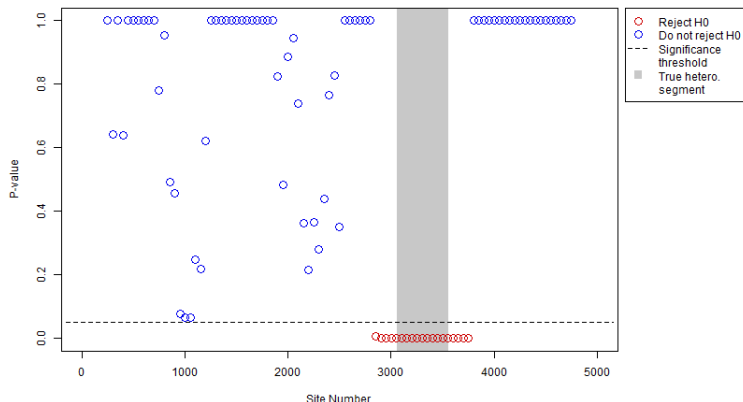
Results

Sliding Window Analysis

1. Were the heterogeneous regions detected?

LRT comparing homogeneous vs. heterogeneous model.

(Complete alignment analysis favoured **heterogeneous** model.)



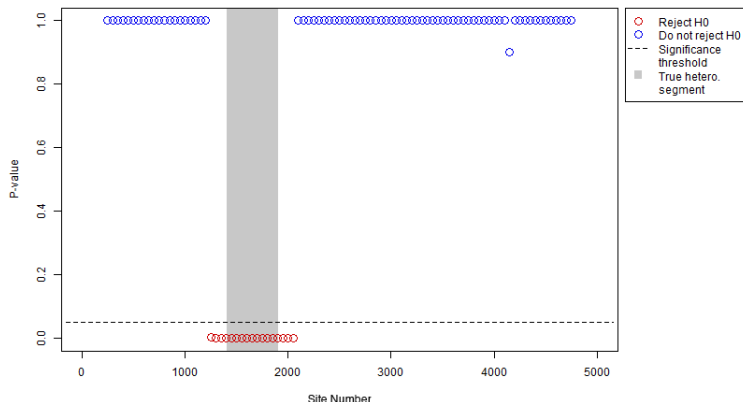
Results

Sliding Window Analysis

1. Were the heterogeneous regions detected?

LRT comparing homogeneous vs. heterogeneous model.

(Complete alignment analysis favoured **homogeneous** model.)

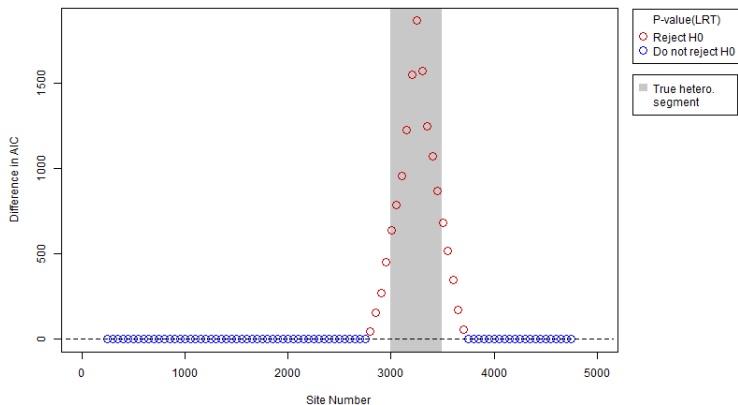


Results

Sliding Window Analysis

1. Were the heterogeneous regions detected?

Difference in AIC, BIC comparing homogeneous vs. heterogeneous model.

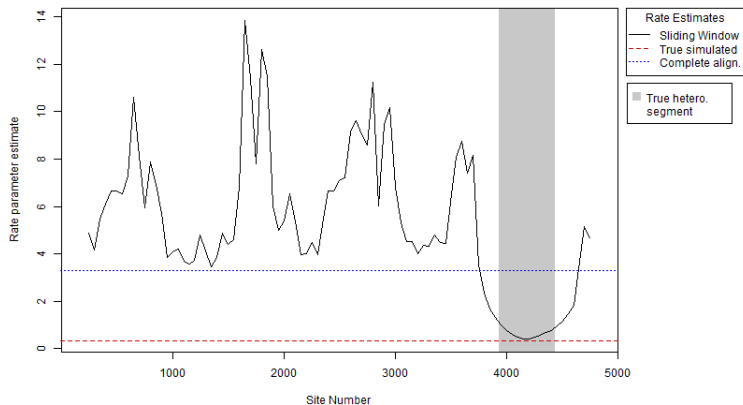


Results

Sliding Window Analysis

2. How accurate was the rate estimate?

Recovered from fitting heterogeneous model.



What did we find?

- ✓ SW approach detected heterogeneous rates consistently.
- ✓ Detected in more situations than under traditional approach.
- ✓ Strong benefits from profiling an alignment.
- ✓ Better overall detection of rate heterogeneity.

Overview: The Sliding Window Approach

Strengths:

- ✓ Gain a deeper insight into the true behaviour of the alignment.
- ✓ Allows multiple analyses which better detect any pattern variation (rates, topology etc.)

Weaknesses:

- ✗ Choosing an appropriate window and step size can be tricky.
- ✗ Computation time can be high.

- SW approach is undoubtedly a useful tool.
- Potential to better detect heterogeneity, and to improve the statistical models we use.

Where to from here?

- ⇒ Continue to test the SW approach in other phylogenetic applications (outlier detection, different forms of heterogeneity).
- ⇒ Create software to make the SW approach more accessible.
- ⇒ Finding optimal window and step sizes?

Acknowledgements

Steffen Klaere



Jessica Leigh

