# Alignment
# Phylogenetic Biology – Week 5

Biology 1425
Professor: Casey Dunn, dunnlab.org
Brown University

# Front matter...

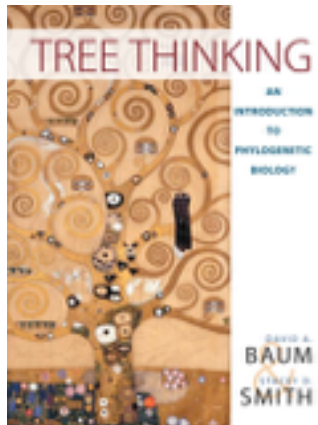All original content in this document is distributed under the following license:
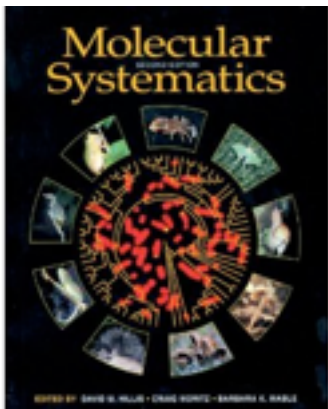
See sources for copyright of non-original content
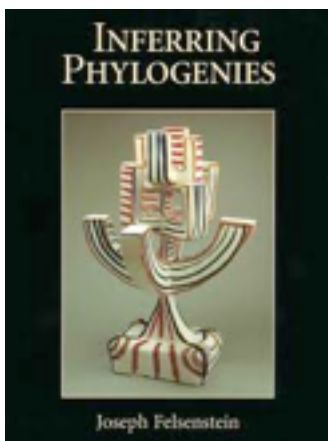
# Sources

Some non-original content is drawn from:

Baum, D and S. Smith (2012) Tree Thinking: and Introduction to Phylogenetic Biology. Roberts and Company Publishers. ISBN 9781936221165

Swofford, D. L., Olsen, G. J., Waddell, P. J., & Hillis, D. M. (1996). Phylogenetic inference. In: Molecular Systematics, Second Edition. eds: D. M. Hillis, C Moritz, & B. K. Mable. Sinauer Associates. ISBN 9780878932825

Felsenstein, J. (2003) Inferring Phylogenies. Sinauer Associates. ISBN 978-0878931774

Other non-original content is referenced by url.

# What is sequence alignment?

The identification of homologous sites in molecular sequence data.

If sequences didn't evolve and could be observed error free, we could just look for identical sequence regions.

Due to evolution, sequence error, and analysis error, we have to ask - How do we know when the same site in two different sequences is homologous?

# Alignment

Reconciles differences that arise from two processes:

- Substitution
- Insertion/deletion (ie, indels)

# Many applications of alignment

Pairwise sequence alignment (eg blast) to find homologous sequences

Alignment of raw sequence reads to a reference sequence to identify variants

Multiple sequence alignment to build character matrices

# Pairwise alignment

# Read alignment

https://vimeo.com/120429438

# Multiple sequence alignment

# Multiple sequence alignment

Common models used in phylogenetic inference (eg GTR) accommodate substation, but not insertion/ deletion

Most phylogenetic programs therefore don't infer homology, they assume that each column is a set of homologous sequences

They treat gaps introduced by indels as missing data

# Multiple sequence alignment

Need an aligner upstream of phylogenetic inference that infers which sequence differences are due to substitution and which are due to insertion/deletion

Partitions sites according to inferred mechanism: changes due to indels are put in separate columns, each column contains sites that are hypothesized to be only due to substitution

# Multiple sequence alignment

Many MSA tools are available:

mafft
clustalw
muscle
t-coffee

# Multiple sequence alignment

In general, they work by:

1. Defining a set of penalties for site differences and the introduction of gaps

2. Heuristically searching for an alignment that minimizes these penalties

# Multiple sequence alignment

The scoring matrix is used to evaluate site differences.

Explains how surprised we should be to see a particular substitution that leads to a difference between homologous sequences.

Related to site substitution models.

# Multiple sequence alignment

BLOSUM62
scoring matrix

|       | Ala | Arg | Asn | Asp | Cys | Gln | Glu | Gly | His | Ile | Leu | Lys | Met | Phe | Pro | Ser | Thr | Trp | Tyr | Val |
|-------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Ala | 4 |
| Arg | -1 | 5 |
| Asn | -2 | 0 | 6 |
| Asp | -2 | -2 | 1 | 6 |
| Cys | 0 | -3 | -3 | -3 | 9 |
| Gln | -1 | 1 | 0 | 0 | -3 | 5 |
| Glu | -1 | 0 | 0 | 2 | -4 | 2 | 5 |
| Gly | 0 | -2 | 0 | -1 | -3 | -2 | -2 | 6 |
| His | -2 | 0 | 1 | -1 | -3 | 0 | 0 | -2 | 8 |
| Ile | -1 | -3 | -3 | -3 | -1 | -3 | -3 | -4 | -3 | 4 |
| Leu | -1 | -2 | -3 | -4 | -1 | -2 | -3 | -4 | -3 | 2 | 4 |
| Lys | -1 | 2 | 0 | -1 | -3 | 1 | 1 | -2 | -1 | -3 | -2 | 5 |
| Met | -1 | -1 | -2 | -3 | -1 | 0 | -2 | -3 | -2 | 1 | 2 | -1 | 5 |
| Phe | -2 | -3 | -3 | -3 | -2 | -3 | -3 | -3 | -1 | 0 | 0 | -3 | 0 | 6 |
| Pro | -1 | -2 | -2 | -1 | -3 | -1 | -1 | -2 | -2 | -3 | -3 | -1 | -2 | -4 | 7 |
| Ser | 1 | -1 | 1 | 0 | -1 | 0 | 0 | 0 | -1 | -2 | -2 | 0 | -1 | -2 | -1 | 4 |
| Thr | 0 | -1 | 0 | -1 | -1 | -1 | -1 | -2 | -2 | -1 | -1 | -1 | -1 | -2 | -1 | 1 | 5 |
| Trp | -3 | -3 | -4 | -4 | -2 | -2 | -3 | -2 | -2 | -3 | -2 | -3 | -1 | 1 | -4 | -3 | -2 | 11 |
| Tyr | -2 | -2 | -2 | -3 | -2 | -1 | -2 | -3 | 2 | -1 | -1 | -2 | -1 | 3 | -3 | -2 | -2 | 2 | 7 |
| Val | 0 | -3 | -3 | -3 | -1 | -2 | -2 | -3 | -3 | 3 | 1 | -2 | 1 | -1 | -2 | -2 | 0 | -3 | -1 | 4 |

https://commons.wikimedia.org/wiki/File:BLOSUM62.gif

# Multiple sequence alignment

Different types of gap penalties:

**Gap opening penalty** - The cost of creating a gap of one site where there was no gap

**Gap extension penalty** - The cost of adding gaps adjacent to an existing gap.