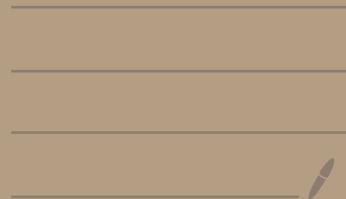


Bioinformatics

WS 19/20



What is Bioinformatics?

- in Bioinfo: intend to develop, optimize, parallelize models, algs, production-level SW for analyzing, storing, extracting knowledge from biological raw data

DNA

- 4 basic nucleotides:

A adenine

C cytosine

G guanine

T thymine / U uracil (RNA)

Ambiguity Code

- char denotes whether it could be A or C, ...

Code	Represents	Complement
A	Adenine	T
C	Cytosine	G
G	Guanine	C
T	Thymine	A
Y	Pirimidine (C or T)	R
R	Purine (A or G)	Y
W	weak (A or T)	S
S	strong (C or G)	W
K	I keto (T or G)	M
M	amino (C or A)	K
D	A, G, T (not C)	H
V	A, C, G (not T)	B
H	A, C, T (not G)	D
B	C, G, T (not A)	V
X/N	any Base	X/N
-	Gap	-

DNA Sequencing

- reading the nucleotide bases in a DNA molecule

> Sanger sequencing:

- accuracy: 99,9%
- sequence length: 300 - 900 nucleotides
- costs: \$2000 per 1 000 000 nucleotides ↗
- few sequences: up to 100 ↗

> Next-generation sequencing

- accuracy: 98 - 99,9%
- sequence length: 100 - 400 nucleotides
- costs: \$1 - \$10 per 1 000 000 nucleotides ↘
- many sequences: 500 - 3 000 000 per sequencer run

Model Organism

- extensively studied / sequenced species → understand particular biological phenomena
- expectation: discoveries provide insight into workings of other organisms

Criteria

- easy experimental manipulation
- ease of genetic manipulation
- easy to grow → short life-cycle (generation times)
- easy to extract DNA data
- economical importance

Examples

- Escherichia coli*: gut bacterium, fast growing, inexpensive to cultivate
- Drosophila melanogaster*: fruit fly, breeds quickly
- Arabidopsis Thaliana*: flower, small genome

Terminology

base pair = pairing of A + T / C + G in double-stranded DNA
ACGGT → 5 nucleotides
↳ 5 base pairs
kB = kilo-bases
Mb = mega-bases
Gb = Giga-bases

Genome = full genetic information of an organism

- contains all chromosomes
- comprises coding + non-coding sequence data
 - encodes proteins
 - does not encode proteins function: partially known, regulation of protein processes

Genome Size

not nec. correlated with organism complexity

· homo sapiens: 3,2 Gb

· marbled lungfish: 130 Gb

plants: often large genomes (redundant info → hybridization)

Shotgun Sequencing

Reading a genome:

- 1) Break up genome randomly into fragments
- 2) Read fragments
- 3) Assemble fragments into a genome

Characteristics

Coverage: how many fragments/reads cover one nucleotide on the genome

Fragment length

paired-end vs single-end reads

de novo vs by reference assembly

by reference assembly = read mapping

- closely related species Y \rightarrow genome available
- map reads of X to genome of Y to assemble

de novo assembly

- no closely related species' genome available
- assemble genome out of read soup

paired-end reads

paired -ends = two ends of the same DNA molecule

\hookrightarrow sequence one \rightarrow turn around \rightarrow sequence the other end

\hookrightarrow 2 sequences = paired-end reads

Gene = coding part of DNA

typically $\approx 1000\text{bp}$ long

encodes either for RNA or protein

Protein, RNA Sequences

T \sim U

proteins: 20 letter alphabet

3 DNA / RNA chars $\stackrel{?}{=}$ 1 protein char
= codon

61 triplets encode 20 proteins

\hookrightarrow redundancy (usually in third codon pos.)

3 triplets for start, stop (start $\stackrel{?}{=}$ ATG $\stackrel{?}{=}$ Met / M acid)

Synonymous substitution / mutation

GCC \rightarrow GCT = Alanine \rightarrow Alanine

Non-synonymous subst. / mut.

GGT \rightarrow GTT = Glycine \rightarrow Valine

DNA \leftrightarrow Protein Translation

DNA → protein: non ambiguous, redundant

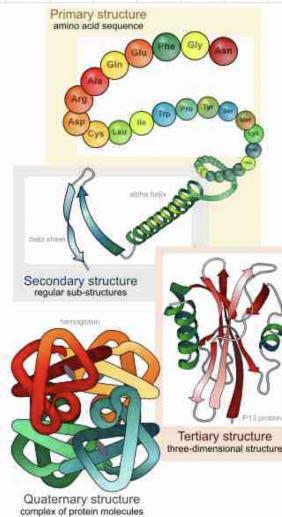
Protein \rightarrow DNA : ambiguous \rightarrow several DNA triplets = same amino acid

Proteins

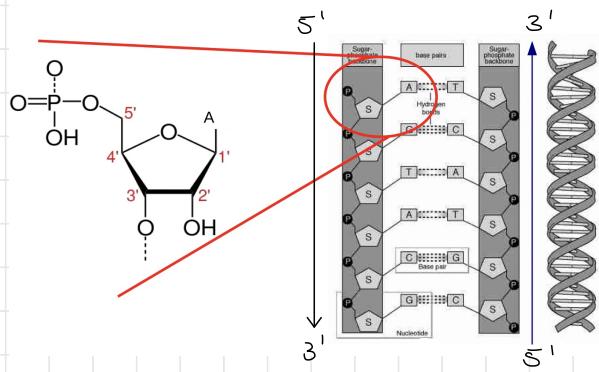
- structural proteins → tissue building blocks
 - enzymatic proteins → catalysts of specific biochemical reactions
 - homo sapiens ≈ 20 000 proteins
 - protein = sequence of amino acid chars
= protein letters = RESIDUES

Structure

- residues = primary structure
 - structure determines function/
effect of a protein



3', 5' end



Convention: DNA sequences
in 5' → 3' direction

→ Genes have a direction: depending on which DNA strand encodes

Prokaryotes / Eukaryotes

Prokaryota: - no cell nucleus, mostly unicellular organisms
· a gene encodes protein / rRNA

Eukaryota: · organism with cell nucleus

· not the entire gene may encode for a protein, just parts of it

Introns: gene not used in protein synthesis

Exons: parts of the gene used for protein synthesis

RNA

· RNA = copy of coding DNA strand (Gene)

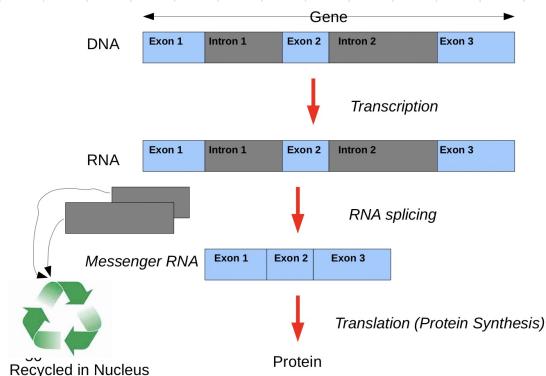
· involved in DNA transcription to construct

- a protein: DNA → RNA → Protein

↳ translation = coding - RNA

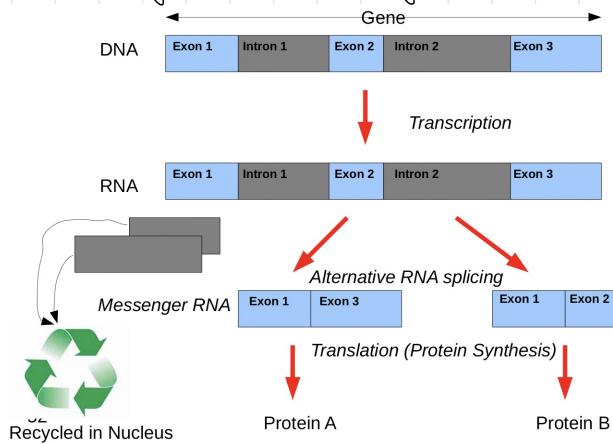
- non-coding RNA: DNA → RNA with other direct cell function

RNA Splicing Eukaryota



Alternative Splicing

Increases coding power of a gene



mRNA = messenger RNA

→ transports RNA data to the ribosome for protein synth.

rRNA = ribosomal RNA

→ carries out translation in ribosome via catalysis

tRNA = transfer RNA

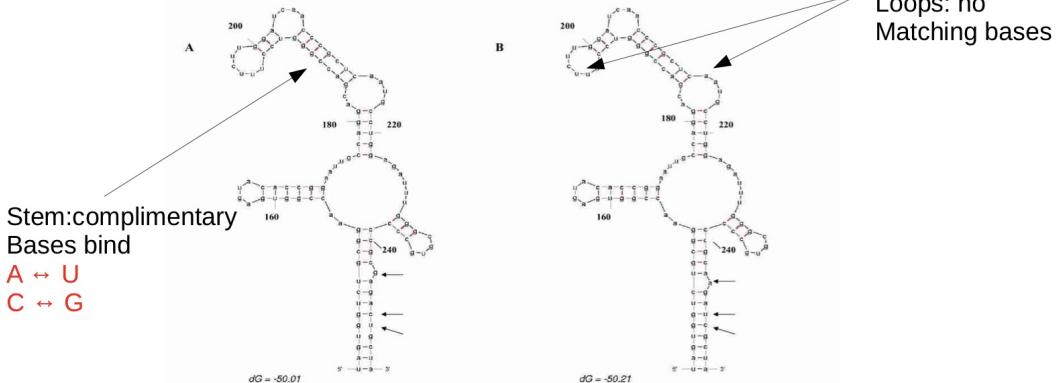
→ brings in the amino acids

only few genes common to all species

e.g. rRNA

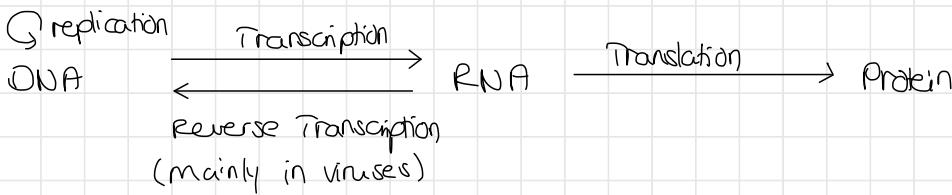
↳ can be used to infer evolutionary relationships among all species

Secondary Structure influences function of the molecule



- matching bases in stem cannot mutate independently from each other

Central Dogma of molecular biology



Transcriptome

- set of all RNA (mRNA, rRNA, tRNA) in a cell
- reflects activity in a cell
- temporal + spatial component: transcriptome varies
 - ↳ different active genes

Meta-Genome

- all genetic material of a community

Chromosomes

- in cell: DNA organized in long molecules = chromosomes
- all chromosomes = genome
- human: 46 (2×23 chr.) prot: 1 chr.
- mouse: 40
- donkey: 62

Eukaryotic Chromosomes

- paired chr. = homologous \rightarrow parental / maternal
- some genes in hom. Chr. exactly identical
some not \rightarrow different genotypes
- genes that appear in different forms = Alleles
- cells containing pairs of chr. = diploid
one of each = haploid (sexual reproduction)

Species

by reproduction:

- 2 species that can reproduce
- bacteria / viruses?

evolutionary species concept:

- via ancestral descent in evolutionary tree

general lineage concept:

- an independently evolving lineage

phylogenetic species concept

by sequence similarity + statistical methods

Taxonomy

- group biological organisms (species) into groups with similar charact.
- define char. of groups at different hierarchy levels
- taxonomic ranks
 - domain \rightarrow 3 domains of life
 - kingdom
 - phylum
 - class
 - Order
 - family
 - genus
 - species

Phylogeny

- unrooted binary tree
- leafs = currently living organisms represented by their DNA/protein sequences
- inner nodes = hypothetical common ancestors
- outgroup = one /more closely related but different species

Taxon

- used to denote clades /subtrees in taxon. (phylog.)
- a group of species that form a biological unit
- in phylogenetics: refer to single leaf as taxon

Pair-wise sequence alignment

DNA, protein sequences = strings

global alignment = align full strings



local alignment = align similar substrings



Approach:

- online algorithms : sequences cannot be preprocessed
→ no index build
- dynamic programming: break down problem into smaller subprob.
→ store results to avoid redundant comp.

Substring

x = substring of y if u, v exist with $y = uxv$

String Distance

- positivity $d(x, y) \geq 0$
- separation $d(x, y) = 0 \Leftrightarrow x = y$
- symmetry $d(x, y) = d(y, x)$
- triangle inequality $d(x, y) \leq d(x, z) + d(z, y) \quad \forall z$

String Operations

- Substitution of a letter
- Deletion of a letter
- Insertion of a letter

Hamming Distance

- for strings with $|x| = |y|$
- counts # differing chars

$$\delta_h(\text{ping}, \text{pong}) = 1$$
$$\delta_h(\text{tata}, \text{ata}) = 4$$

⊕ easy to compute

Edit Distance

$$\delta_e(x, y) = \text{minimum } \# \text{ of operations to transform } x \text{ into } y$$

⊕ more sensful in sequence comparison

String Alignment

- insert gaps into both strings to achieve same length

$$\begin{array}{c} \text{A C G -- A} \\ \text{A T G C T A} \end{array}$$

Cost: hamming dist. btw new sequences

min. cost: edit dist. btw origin sequences

Compute Edit Distance

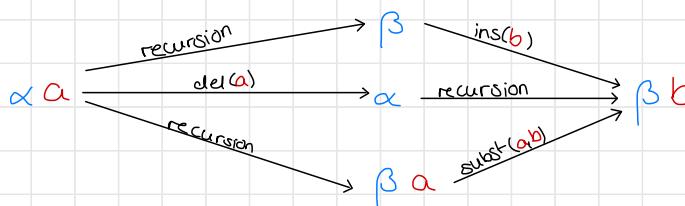
Let $X = \alpha a$ $Y = \beta b$

Compute $\delta(\alpha a, \beta b)$ assuming we know

- $\delta(\alpha, \beta)$
- $\delta(\alpha, \beta b)$
- $\delta(\alpha, \beta)$

→ Recursion until $X = \epsilon$, $Y = \epsilon$

· 3 possible sequences of operations:



Costs:

$$\delta_e(\alpha a, \beta) + 1$$

$$\delta_e(\alpha, \beta b) + 1$$

$$\delta_e(\alpha, \beta) \text{ if } \alpha = b$$

$$\delta_e(\alpha, \beta) + 1 \text{ else}$$

↳ Take shortest path:

$$\delta(x, y) = \min \left(\begin{array}{l} \delta_e(\alpha a, \beta) + 1, \\ \delta_e(\alpha, \beta b) + 1, \\ \delta_e(\alpha, \beta) + 1 * (a == b) \end{array} \right)$$

Needleman - Wunsch algorithm

- computes best global alignment
- complexity $O(|X| \cdot |Y|)$ → filling is most expensive

Algo : 1) Setup table with x, y $\in \Sigma$ $x = \text{monkey}$
 $y = \text{money}$

2) Init first row, first column

3) Fill each element with min of three prev values
store best path

4) backtrace to get chain of operations

	M	O	N	K	E	Y
O	1	2	3	4	5	6
M	1	0	1	2	3	4
O	2	1	0	1	2	3
N	3	2	1	0 ← 1	2	3
E	4	3	2	1	1	2
Y	5	4	3	2	2	1

Arrows indicate the path from one cell to the next, showing the sequence of operations: Match (down), mismatch (left), insertion (up-left), and deletion (up).

Adding weights

penalize less likely operations with weights

e.g. subst. happen 5x more often than ins / del

$$\hookrightarrow \delta_e(x, y) = \min (\delta_e(\alpha a, \beta b) + 5, \delta_e(\alpha, \beta b) + 5, \delta_e(\alpha a, \beta) + 1 * (a \neq b))$$

Similarity: pos. / neg. weights to favor similarities

- ins / del / non-identity subst. neg. weight
- identity subst. pos. weight

! min distance \leftrightarrow max similarity

Local Alignment

more useful for dissimilar sequences with possible similar regions

Approaches:

a) Find substrings with min. distance

\hookrightarrow short strings selected

b) Find substrings with max. similarity

Smith-Waterman algorithm

\approx Needleman-Wunsh

1) Init first row, column with 0

2) Fill table with same recursion

3) Find largest value in table \rightarrow best similarity

4) Traceback until a 0 reached

Substitution Matrices

- DNA/proteins alter over time (DNA mutations)
- subst. matrices describe rate at which a char in a sequence changes to other char states over time
- Transitions $A \leftrightarrow G$ / $C \leftrightarrow T$ more often than transversions

↪ take into account that e.g. $\text{subst}(A, G) < \text{subst}(A, C)$

$$\Rightarrow d(x, y) = \min(\text{del}(\alpha\alpha, \beta) + \text{ins}(b), \\ \text{del}(\alpha, \beta b) + \text{del}(a), \\ \text{del}(\alpha, \beta) + \text{subst}(a, b))$$

BLOSUM matrices

Blocks of Amino Acid Substitution Matrix

· used for sequence alignment of protein sequences

· high numbers (e.g. BLOSUM 80) for closely related sequences

Log-odd scores

$$S_{i,j} = \frac{1}{x} \log \left(\frac{\pi_{i,j}}{\pi_i \times \pi_j} \right)$$

\nwarrow observed frequency of subst.
 \swarrow expected frequency of subst.

π_i = stationary frequency of state i

$S_{i,j} > 0 \rightarrow$ very likely

$< 0 \rightarrow$ unlikely

Hamming Distance with mismatches

$x = \text{long string}, y = \text{short string}$

↪ find all substr. x of X with $|x|=|y|$ and $\delta_h(x,y) \leq k$

Recursion / DP formula

- 1) init first row with 0
- 2) init first col with $|c|+1$
- 3) Fill table: $T(i,j) = T(i-1, j-1) + \text{sub}(i,j)$
(no del/ins with hamming dist.!!)
- 4) Last line: keep all alignments with score $\leq k$

BLAST

Basic Local Alignment Search Tool

fast heuristic to find similar sequences

Assumption: sequence similarity \Leftrightarrow evolutionary / functionally relatedness

Naive Approach: Smith-Waterman algo

- > build pair-wise alignment of query sequence Q with every sequence in DB S_1, \dots, S_d
 - ↪ report best matches
- ① DP matrix for every sequence
 $O(mn)$ compl. for one Sm-W-algo

Algorithm

- 1) Seeding: find seeds = common substrands b/w query and db sequences
- 2) Extension: starting from seed: extend alignment in both directions
→ high-scoring segment pairs (HSP)
- 3) Evaluation: assess statistical significance of each HSP

Seeding

- a) Build L_1 : all subwords (factors) of length ω of query sequence

Query ACTTGCTA , $\omega = 5$
 ACTT G
 CTTGC
 TIGCT
 TGCTA

- b) $\forall w_i \in L_1$: build neighborhood N_i with similar subwords
↳ add only subwords with similarity score $\geq T$

Similarity: via substitution matrix

proteins: BLOSUM

DNA: +2 for match, -3 for mismatch (+5/-4)

$$w_i = TTGCT$$

$$\hookrightarrow TTGAT$$

$$+2 +2 +2 -3 -2 = 5 > 4 \checkmark$$

$$TTGCT$$

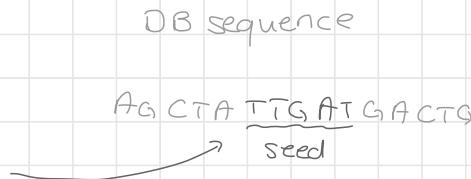
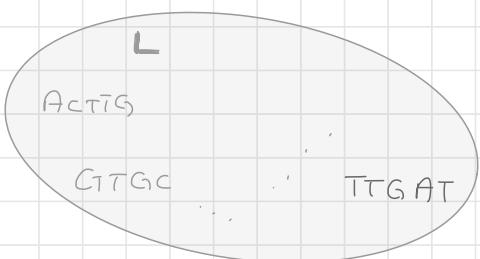
$$TGCCT$$

$$+2 -3 -3 +2 +2 = 0 < 4 X$$

- c) build $L_2 = UN$: combination of all neighborhoods

build final list $L = L_1 \cup L_2$

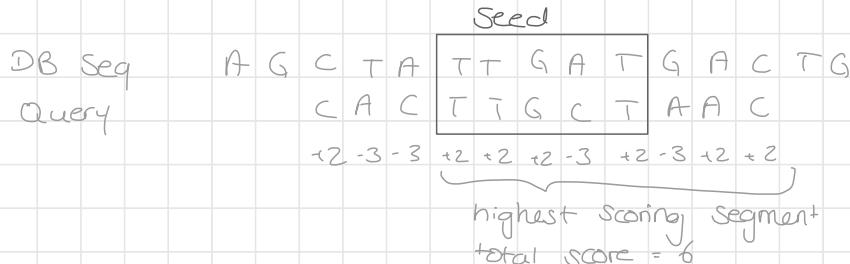
Scan DB for exact matches of subwords in L



Extension

- try to extend alignment to the left + right from the seed
- stop if current total score drops by more than X compared to current max
- trim alignment back to max score

$$X = 3$$



Evaluation

- probability p of observing score $S \geq X$ by chance:

$$p(S \geq x) = 1 - \exp(e^{-\lambda(x-\mu)})$$

- alignment scores follow Gumbel extreme value distribution

↳ E -value: expectation value (takes DB size d into account)

$$E \approx 1 - e^{-p(S > x)d}$$

E -value = expected # of times unrelated DB sequence would obtain score $S \geq x$ by chance

⇒ small E -value for biologically related sequences

Genome Assembly

- reconstructing original DNA from its factors

De novo

- assemble reads into contigs using overlaps between them
- use mate pairs to combine contigs into scaffolds
(contig order + orientation out of paired-end reads)
- manual joining of scaffolds

metrics of assembly quality: length of contigs, scaffolds

1) Overlaps

a) Overlap Graphs

- each read = node
- compute pairwise alignments b/w reads, overlap \rightarrow edge

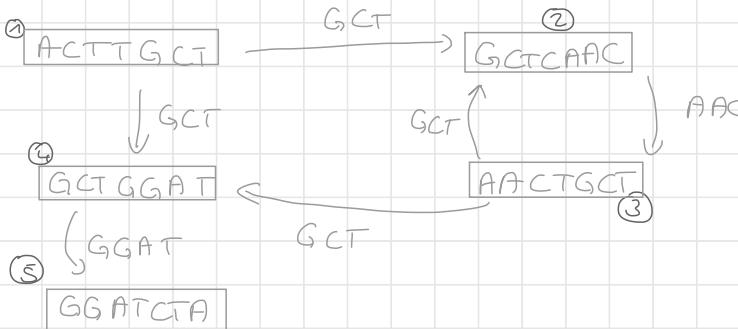
\rightarrow walk along Hamiltonian path (every node exactly once)

\hookrightarrow reconstruct original sequence

- Start:
- circular genome: any node
 - linear genome: node with no inbound edges

Genome ACTT GCT CA ACT G CTG GA TCTA

Overlap Graph.



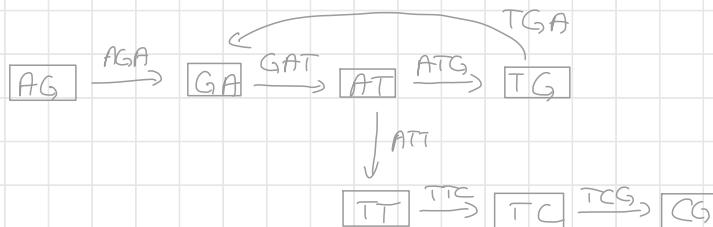
- no efficient algo to find Hamiltonian path
- pair-wise alignments $O(n^2)$ ($n = \# \text{ reads}$)
- only good for small # of reads with significant overlap

b) De Bruijn graphs

- decompose reads into k -mers (us. $k = 20-50$)
 $k < \text{length of shortest read}$
- unique k -mer = edge
- nodes = $(k-1)$ -mers \rightarrow pre/suffix of connecting k -mer edge

\hookrightarrow find original sequence by finding Eulerian path (each edge ex. once)

A G A T G A T T C G
 $k=3$ AGA



- + compact repr. of repeats
- building graph: time $O(N)$ ($N = \text{total length of all reads}$)
space $O(\min(G, N))$ ($G = \text{genome size}$)
- efficient algs to find Eul. path exist

- information loss \rightarrow k -mer extraction
- if reads error-free: Eul. path always exists
- repeats: multiple alternative paths
- disjoint graphs if lack of coverage in some regions

Error correction

- reads contain sequencing errors
- bubbles = errors in middle of the read
- tips = errors at ends of the read
- detection/fixing: use coverage info

By reference

Sliding window approach

- slide each read along reference genome, mark matching pos.
- if gaps allowed: DP algo (eg Smith-Waterman)

- huge complexity

Hashing

- build genome index \rightarrow hash table stores positions of all k-mers
 $k \ll$ read length
- for each read:
 - select proxy k-mer (leftmost/middle part)
 - hash table lookup: find all positions of this k-mer in genome \rightarrow seeds
 - for each seed: try to extend alignment (e.g. Smith-Waterman algo)

Variations

- inexact matches with spaced seeds
↳ binary mask defines pos where mismatches allowed $M1 \underline{00} 1$
- multiple k-mers per read
↳ require $\geq n$ seed matches for mapping location to be considered
- inverted approach: build hash table from reads, search for k-mers present in the reference

Burrows-Wheeler Transformation BWT

- 1) write down all cyclic rotations of source string S
- 2) sort rows lexicographically
- 3) store last column \rightarrow BWT(S)

Rows = sorted list of suffixes
→ efficient substring search

Pattern search:

1st column obtained by sorting last one

2nd analogous

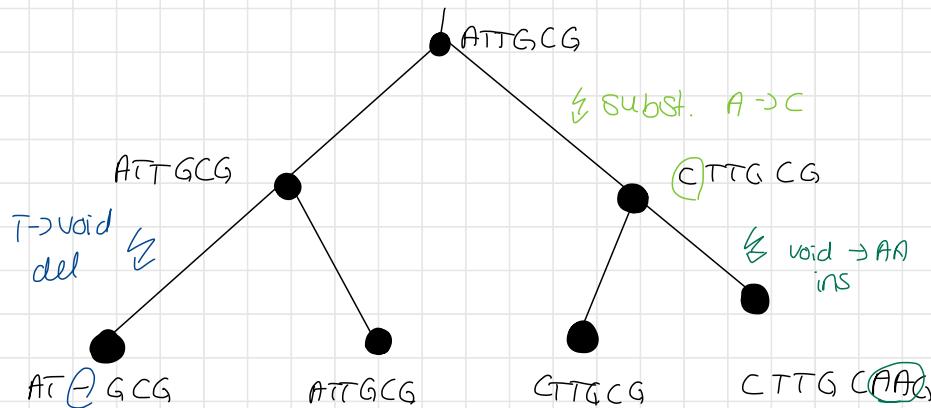
~ proceed and check rows for suffix pattern

Multiple Sequence Alignment

· assess origin sequence for multiple homologous sequences

indel = insertion / Deletion

↳ indel lengths 1-3 frequent (3 = codon length!)



Aligned Data

AT - G C -- G
AT T G C -- G
CT T G C -- G
CT T G C A A G

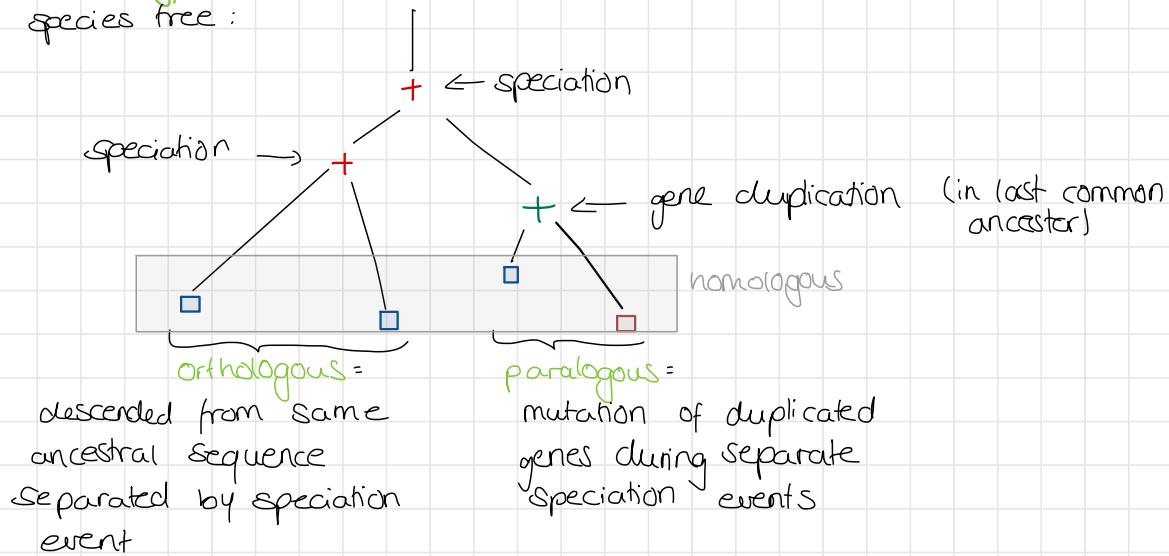
} goal: infer homology
→ which chars share common evolutionary history

Importance of MSA

- input for phylogenetic reconstruction
- discover important parts of protein family (group of evolutionary related genes/proteins in different species with similar function/structure)

Terminology

species tree:



· orthologous sequences = sequences in different species that evolved from same ancestral gene

· homologous characters = chars that share common evolutionary history

- ! sequence similarity \rightarrow homology
 - \hookrightarrow convergent evolution possible (independently evolved twice)
short sequences: by chance

Algorithm Outline

given n orthologous sequences s_1, \dots, s_n of different lengths

Gap Insertion

insert '-' such that:

- all sequences have same length
- some criterion is optimized
- homologous chars s_i, s_j aligned to each other (some alignment col)
- not allowed: columns of entirely gaps

Alignment Criteria

how to define alignment quality

a) Sum of pair measure

- score each site (column), add up scores overall sites
→ penalize mismatches, gaps
favor matches
- per-site score: sum of all pairwise scores btw chars of a site
- example:

S1	A	A	G	A	A	-	A	PC(-,-) = 0
S2	A	T	-	A	A	T	G	
S3	C	T	G	-	G	-	G	
	1	1	1	1	1	1	1	
PC	= edit-dist.	1	1	1	1	1	1	$\sum 1 = 14$



- optimal MSA - computing NP - complete
- complexity: time + space $O(m^n)$ $m =$ length of seq.
- not guaranteed that SP is biologically most plausible criterion
- dependant on arbitrary PC

b) Star Alignment Approximation

- pick S_c = center sequence
 - compute all $\frac{n^2}{2} - n$ optimal pair-wise alignments
 - select seq. with largest similarity to all others
- align remaining seq to S_c (pairwise-seq. alignment algo)
- "once a gap, always a gap"
 - gaps inserted to S_c cannot be removed
- produces NSA with SP score $< 2 \cdot$ optimum

Example:

S_1 : A T T G C C A T T T $\leftarrow S_c$

S_2 : A T G G C C A T T T

S_3 : A T C C A A T T T T

S_4 : A T C T T C T T

S_5 : A C T G A C C

LD S_1 : A T T G C C A T T T - - \nearrow
 S_2 : A T G G C C A T T T -- \nwarrow

S_1 : A T T G C C A T T T - - \leftarrow Gap inserted
 S_3 : A T C - C A A T T T T

S_1 : A T T G C C A T T T - - \swarrow
 S_4 : A T C T T C - T T - -
 S_1 : A T T G C C A T T T - - \searrow
 S_5 : A C T G A C C - - - -

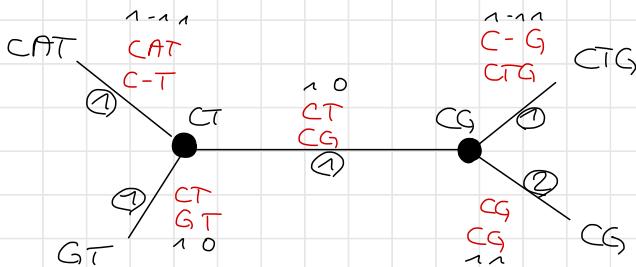
once a gap,
always a gap'

↳ Alignment:
 S_1 : A T T G C C A T T T - -
 S_2 : A T G G C C A T T T - -
 S_3 : A T C C A A T T T T
 S_4 : A T C T T C - T T - -
 S_5 : A C T G A C C - - - -

c) Tree Alignment if evolutionary tree for the seq. available

find assignment of seq. to the inner nodes to maximize sum over similarity scores on all branches

Example:

$$p(a, b) = \begin{cases} 1 & \text{if } a = b \\ 0 & \text{else} \end{cases}$$
$$p(a, -) = -1$$


→ overall score = 6

! NP-hard → ancestral states not given



- hen & egg problem:
 - need MSA to build tree
 - need tree to compute MSA
- ↳ if alignment wrong, tree might be wrong
- ↳ if tree wrong, MSA might be wrong

Idea: simultaneous inference of tree + MSA

↳ solving 2 NP-hard / NP-complete problems at once??

d) Progressive MSA

- Guide tree approach

1) Build guide tree

- compute all $\frac{n^2}{2} - n$ pair-wise distances btw. the n sequences
- hierarchical clustering (e.g. NJ algo)

2) calculate pair-wise sequence-sequence

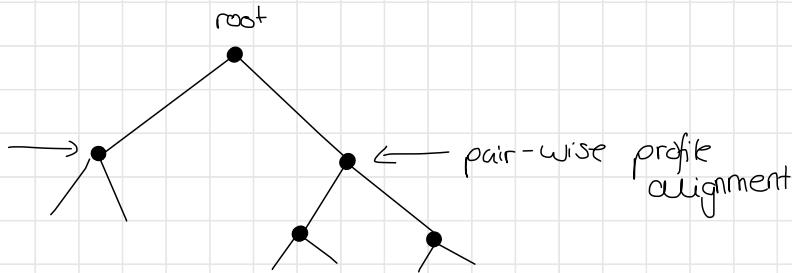
Seq - profile

profile - profile

alignment bottom-up

MSA :

pair-wise
sequence
alignment



- Start with the 2 closest seq. in guide tree
- calc. pair-wi. alignment (Needl. Wunsch)
- calc. profile ("mischsequenz")
- set parent = profile

Loop: take next 2 seq. / profiles that are closest
calc. profile

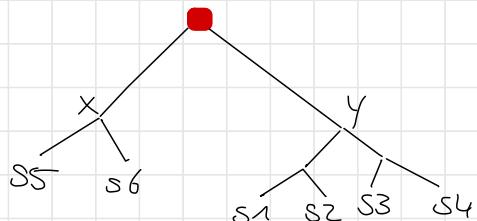
End: one profile left → MSA

Profile Alignment:

· average over all possibilities

example:

	0	1	2	3	4	5	6	7	8	9
S1	P	E	E	K	S	A	V	T	A	L
S2	G	E	E	K	A	A	V	L	A	L
S3	P	A	D	K	T	N	V	K	A	A
S4	A	A	D	K	T	N	V	K	A	A



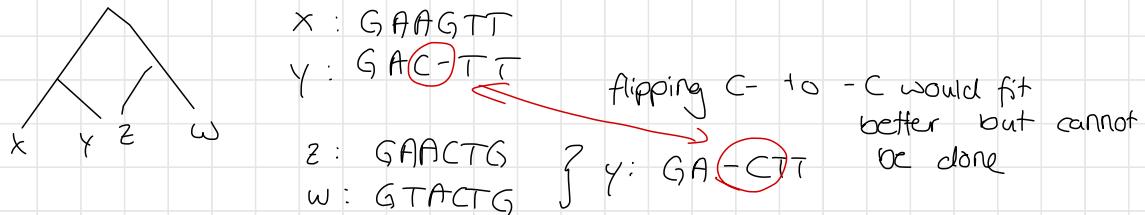
	0	1	2	3	4	5	6	7	8	9
SS	E	G	E	W	G	L	V	L	H	V
S6	A	A	E	K	T	K	I	R	S	A

→ compute score b/w. pos 6 of X and pos 7 of Y

→ weighted average over all 8 possibilities (2×4)

$$\frac{1}{8} * [p(T, V) + p(T, I) + p(L, V) + p(L, I) + p(K, V) + p(K, I) + p(K, I)]$$

(-) Initial pair-wise alignment can't be changed / corrected



e) Iterative progressive MSA

- execute progr. MSA multiple times \rightarrow refine alignment
- e.g. MUSCLE refinement

f) Motif-based approaches

- find small motif (substring) common to all sequences
= anchor / block / region / ...
- if motif found: Align seq. to motif
- align regions around motifs using e.g. progressive alignment

Benchmarking MSAs

- using MSA benchmarks
- using simulation

Simulation

- infer true MSA through simulation from tree
 - \rightarrow disaligns sequences
 - \rightarrow align using MSA algo
 - \rightarrow compare result with true MSA:
 - count correct sites
 - compare SP scores
- \rightarrow infer tree
- \rightarrow compare trees

Phylogenetics

Taxonomy = group of species with similar characteristics

↳ define charact. at different hierarchy levels

→ taxonomic ranks

domain (3 domains of life Bacteria, Archaea, Eukaryota)

Kingdom

phylum

class

order

family

genus

species

Phylogeny

unrooted, binary tree

leaves = extant (currently living) organisms represented by their DNA/
protein sequences

inner nodes = hypothetical common ancestors

outgroup = closely related, but different species (?!)
↳ allows rooting of tree

Taxon

denote clades / subtrees in phylog. / taxonom.

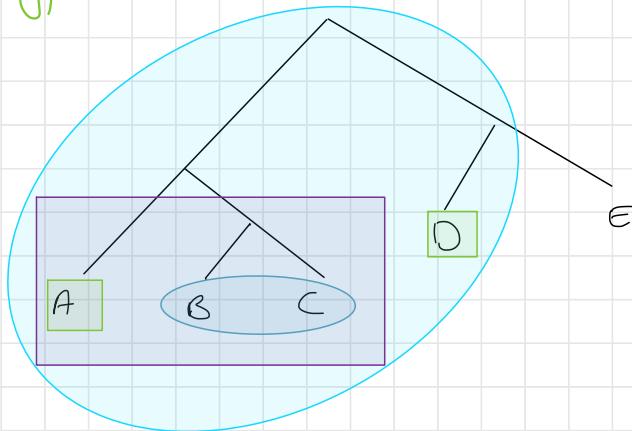
= group of ≥ 1 species forming a biological unit

defined by taxonomists → controversial debates!

phylogen.: single leaf = taxon (pl. taxa)

Terminology

! tree is rooted

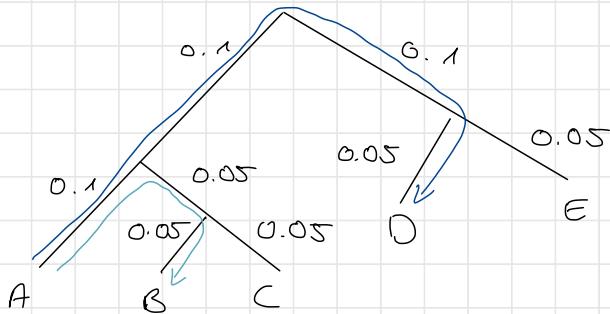


B, C = monophyletic group = sister species \rightarrow sisters to each other

(A, B, C) = monophyletic group = sister species \rightarrow sister to (D, E)

(A, B, C, D) = paraphyletic g: includes ancestor but not all descendants

(A, D) = polyphyletic group = most recent common ancestor excluded



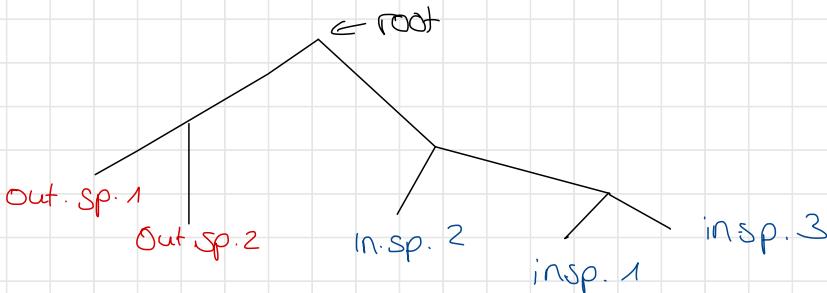
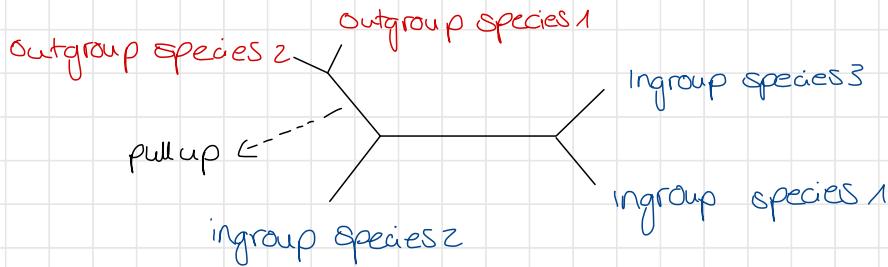
pathistic distance between 2 taxa: sum over branch length along path in tree

$$\text{e.g. } A \leftrightarrow B = 0.1 + 0.05 + 0.05 = 0.2$$

$$A \leftrightarrow D = 0.1 + 0.1 + 0.1 + 0.05 = 0.35$$

Tree Rooting

- root tree using outgroups
 - root tree on connecting edge between outgroup and ingroup



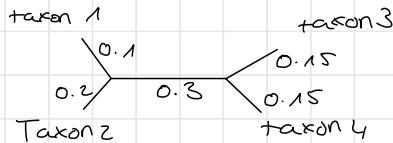
→ choice of outgroup influences tree topology!

Tree Inference

- 1) obtain homologous sequences from the same gene of different species
- 2) MSA program calculates MSA → alignment-free: no MSA calc.
 - ↪ less accurate
- 3) Tree inference program calculates tree

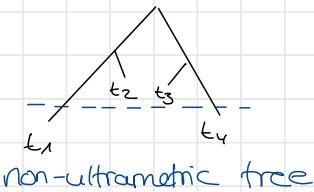
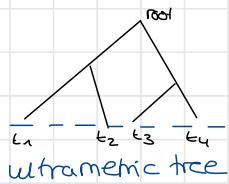
Newick Tree Format

(Taxon1: 0.1, Taxon2: 0.2,
(Taxon3: 0.15, Taxon4: 0.15): 0.3);



Problem: meta-data (→ branch length association) → different interpretations possible

Tree Shapes



↓
evolutionary time
↓ → relative

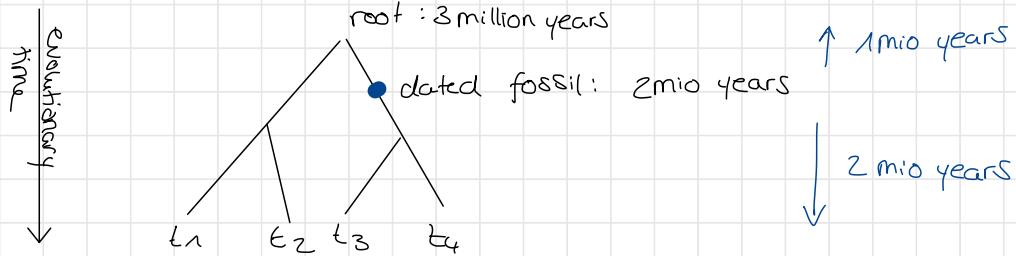
Dating Trees

relative evolutionary time → real time

Requirement: rooted + ultrametric tree

→ rooting: outgroups

ultrametricity: programs for divergence time estimation



→ how to place fossil? no DNA data available

↳ ad hoc empirical knowledge /

computationally using morphological data

structural form

Tree Counts (binary trees)

Rooted trees: tree with n taxa

→ $n-1$ inner nodes

$2n-2$ branches

= # of unrooted trees
for $n+1$ taxa

$$(2n-3)!! = \frac{(2n-3)!}{2^{n-2}(n-2)!} \quad \text{trees} \quad (n \geq 2)$$

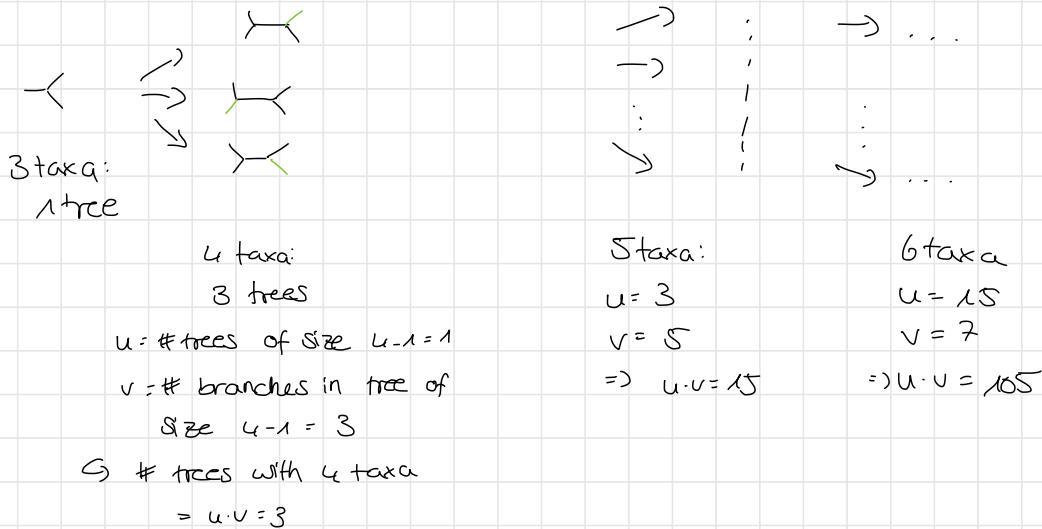
$$\text{e.g. } n=4 \quad \rightarrow \quad (2n-3)!! = \frac{(8-3)!}{2^2(2)!} = \frac{5!}{4 \cdot 2!} = \frac{5 \cdot 4 \cdot 3 \cdot 2 \cdot 1}{4 \cdot 2 \cdot 1} = 5 \cdot 3 = 15$$

Unrooted trees: tree with n taxa:

→ $n-2$ inner nodes

$2n-3$ branches

$$(2n-5)!! = \frac{(2n-5)!}{2^{n-3}(n-3)!}$$



Use of phylogenetic trees

- identifying unknown species
- divergence time estimates
- diversification rates
- viral outbreaks
- forensics
- ...

Building Phylogenetic Trees

distance-based methods

- compute pair-wise distance matrix using MSA
- build tree using these distances
- heuristics (\approx hierarch. clustering methods)
 - e.g. Neighbour Joining / UPGMA
- least-squares method: explicit optimality criterion

distance matrix: memory problem

- ↪ out-of-core / external memory algos for data that is too large for main memory
 - explicit data transfer RAM \leftrightarrow disk in application
 - ↪ circumvent OS paging (use of knowledge \rightarrow faster)

(+) faster (-) less accurate

character-based methods

- optimality criteria f() operate directly on MSA
 - e.g. parsimony
 - max likelihood
 - bayesian inference

? memory intensive!
- calculate score based on current tree topology + MSA
 - ↪ score = how well does MSA data fit tree?

(+) more accurate (-) slow

NP-Problem

- least squares, parsimony, likelihood, bayesian all NP-hard
 - ↪ super-exponential increase in possible trees

Neighbor Joining Algorithm

distance-based

- Start with the distance matrix $D_{i,j}$ ($i, j = 1, \dots, n$)
- For each tip: compute $u_i = \frac{\sum_{j \neq i} D_{i,j}}{n-2}$
- Find pair of tips (i, j) with min. $D_{ij} - u_i - u_j$
↳ connect (i, j) to new ancestral node X
- compute new branch lengths from X to i and j :
 $b_i = 0.5 D_{ij} + 0.5 (u_i - u_j)$
 $b_j = 0.5 D_{ij} + 0.5 (u_j - u_i)$
- Update distance matrix: compute dist. btw X and remaining tips
- Replace tips i, j with $X \rightarrow$ treat X as tip now
- Loop until 2 nodes remain
↳ connect them

Complexity:
Space $O(n^2)$
Time $O(n^3)$

UPGMA

distance-based

Unweighted Pair Group Method with Arithmetic mean

→ produces ultrametric, rooted trees

can be used if known that we have ultrametric tree → usually not the case

- start with distance matrix \mathbb{D}
- find min value $D_{i,j}$
- merge $(i,j) \rightarrow$ new group has $n_{(i,j)} = n_i + n_j$ members
- assign branch lengths $D_{i,j}/2 \rightarrow i \rightarrow (i,j), j \rightarrow (i,j)$
- update distances b/w (i,j) and all $k \neq i, j$:

$$D_{j,k} = \frac{n_i}{n_i + n_j} \cdot D_{i,k} + \frac{n_j}{n_i + n_j} D_{j,k}$$

Complexity:

· naive impl. = $O(n^3)$

~ maintain list of per-col / per-row minima

↳ update-list $O(n)$ } $O(n^2)$
look for min $O(n)$

· time $O(n^2)$

non-ultrametric trees:

can yield misleading results

↳ most trees non-ultrametric! (unequal evolut. rates among all lineages)

Least squares Optimization

distance-based

- given fixed, fully binary tree T with n taxa and pair-wise distance matrix D

- Find branch lengths t_1, \dots, t_{n-3} such that:

sum of squared differences between pair-wise patristic dist. d_{ij} and plain pair-wise dist. D_{ij} is minimized

$$\approx \text{minimize } Q = \sum_{i,j} (D_{ij} - d_{ij})^2$$

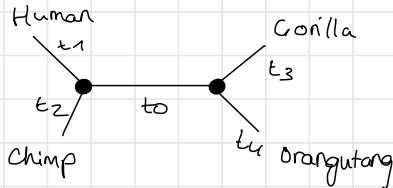
complexity: $\mathcal{O}(n^3)$ if derivative + system of linear equations

$\mathcal{O}(n^2) \dots \mathcal{O}(n)$ if tree-like structure taken into account

- Find tree topology for that minim. Q

\rightarrow NP-hard!

(! single tree in $\mathcal{O}(n)$ - $\mathcal{O}(n^3)$!)



patristic distances:

$$d[H][C] = t_1 + t_2$$

$$d[H][G] = t_1 + t_2 + t_3$$

$$d[H][O] = t_1 + t_2 + t_3 + t_4$$

⋮

Minimum Evolution Method

\approx least squares

explicit criterion

$$\min Q = t_0 + t_1 + t_2 + \dots$$

Parsimony

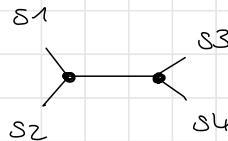
character-based

given a MSA \rightarrow find tree that explains data with least amount of mutations

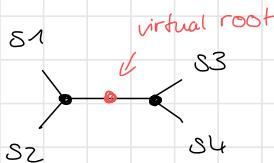
Count Parsimony

MSA

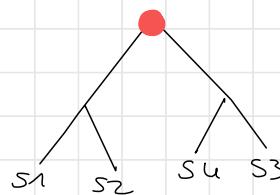
S1	A	A	G	G		
S2	A	T	A	C	indel	~ can be A/C/T/G
S3	A	G	A	G		
S4	T	T	A	T		



\rightarrow Find sequences of inner nodes to min. # of mutations



--->



! Score independent of root position

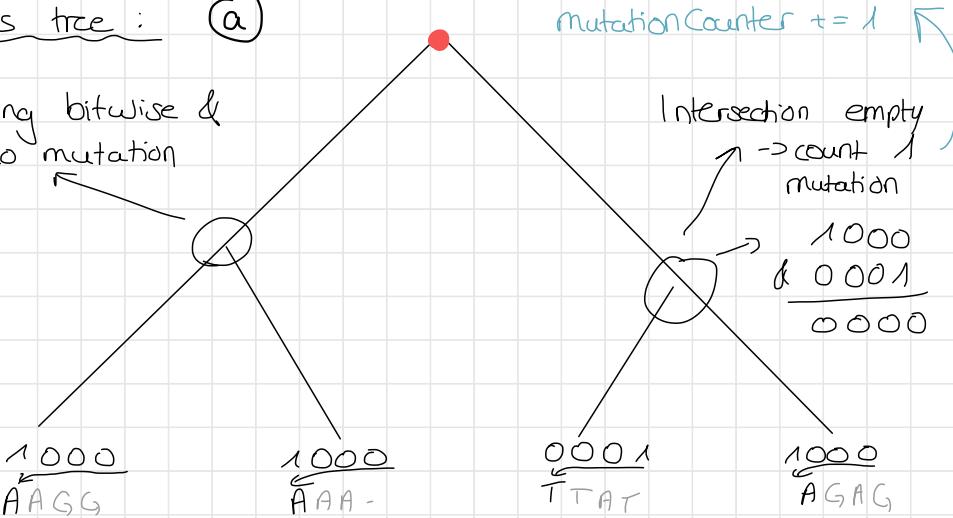
\rightarrow Compute scores by comparing sites (columns) individually
(Assumption: independent evolution of sites)

Step 1 for this tree: (a)

mutationCounter += 1

Intersect using bitwise &
 \rightarrow not empty: no mutation

$$\begin{array}{r} 1000 \\ \& 1000 \\ \hline 1000 \end{array}$$



1000
AAGG

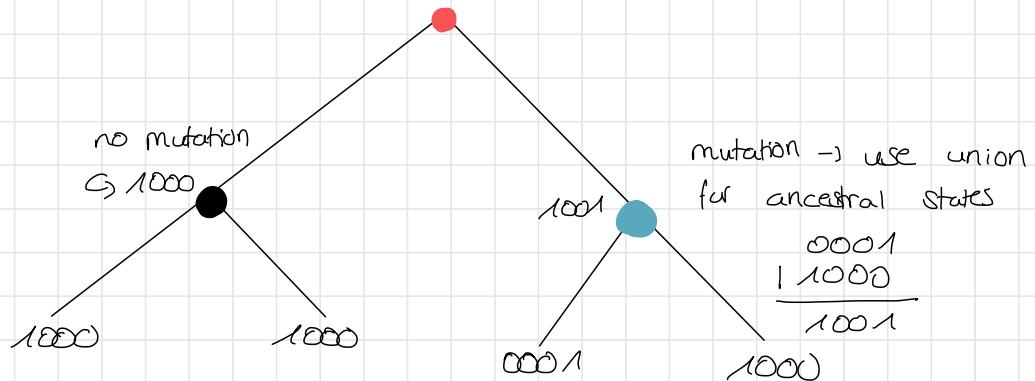
1000
AAA-

0001
TTAT

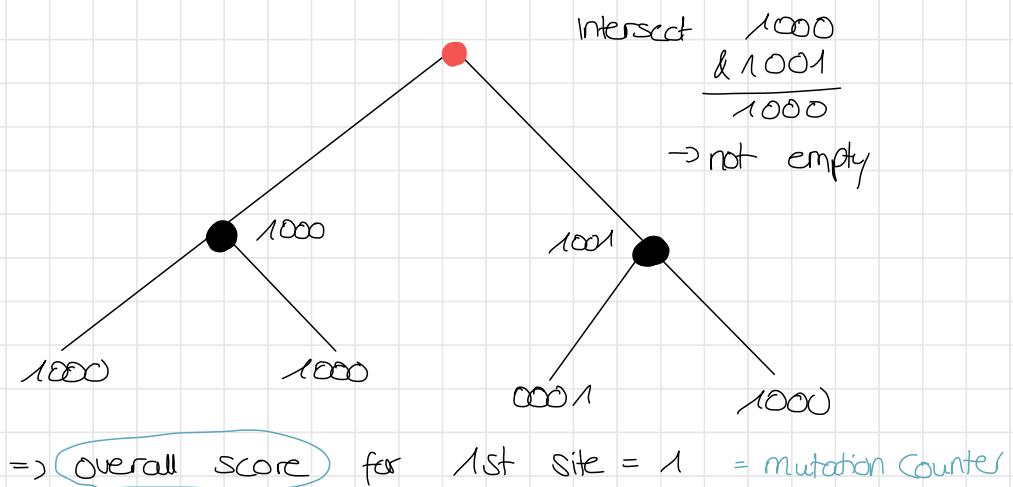
1000
AGAG

(Encode A = 1000 C = 0100 G = 0010 T = 0001)

(b)



(c)



Repeat for all sites, add mutationCounter over all sites

Complexity

Time: to score one tree given MSF with n taxa, m sites
 → $(n-2) * m$ calculations ($n-2 = \#$ inner nodes)
 → $O(nm)$, small constant in $O(1)$

- Space :
- alignment : $n+m+4$ bits
 - ancestral nodes $(n-2)+m+4$ bits
 - score counter : $(n-2) + 82$ bits
 - > $O(nm)$, small constant in $O(nm)$

Tree Search

how to obtain initial starting tree?

↳ comprehensive tree

-> change tree to improve score

- NJ / UPGMA
- random tree
- stepwise addition algo

Random Tree

· insert sequences at random positions

Randomized Stepwise addition algo

- build tree using first 3 randomly picked sequences
- insert following sequences according to best parsimony insertion score

Alter Topology to improve score

alteration mechanisms:

- Hill Climbing
- Simulated annealing
- ...

Basic moves:

- NNI nearest neighbor interchange
- SPR subtree pruning and re-grafting
- TBR tree bisection and reconnection

NNI

- exchange connectivity of 4 subtrees within main tree

SPR

- select + remove subtree from main tree
- reinsert it elsewhere on main tree as new node

TBR

- detach subtree from main tree at inner node
- attempt all possible connections btw. the 2 trees

Long branch attraction

- error where distantly related lineages are incorrectly inferred to be closely related
- happens if both lineages underwent large amount of change

Felsenstein zone = settings where parsimony recovers wrong tree

↳ Solution: max likelihood

Maximum Likelihood

character-based

- calculate likelihood of tree given MSA

$$L(T|D) = P(D|T)$$

T = tree

D = MSA data

assuming sites evolved independently

$$L(T|D) = \prod_{i=1}^n P(S_i|T) = \text{probability that site } 1, \dots, i \text{ evolved to tree}$$

↪ use log (small values cause underflow!)

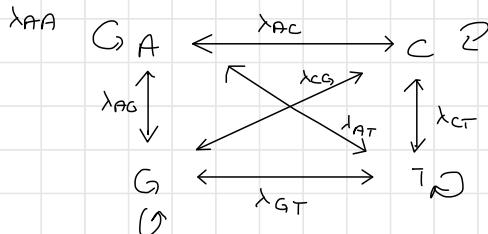
$$\log(L(T|D)) = \sum_{i=1}^n \log(P(S_i|T))$$

Equilibrium Frequencies

$$\vec{\pi} = (\pi_A, \pi_C, \pi_G, \pi_T) \rightarrow \text{must sum to 1}$$

Assumption: time reversibility $\Rightarrow \pi_i P_{i,j}(t) = \pi_j P_{j,i}(t)$

Find substitution model



$$\lambda_{AA} = -\sum_{j \neq A} \lambda_{jA}$$

A C G T

$$\left(\begin{array}{cccc} 0 & & & \\ \lambda_{AC} & 0 & & \\ \lambda_{AG} & \lambda_{CG} & 0 & \\ \lambda_{AT} & \lambda_{CT} & \lambda_{GT} & 0 \end{array} \right) \xrightarrow{\text{Symm}} \text{rate matrix}$$

I Jukes Cantor Model

everything equally likely

$$\begin{matrix} A \\ C \\ G \\ T \end{matrix} \left(\begin{array}{c} \alpha \\ \alpha & \alpha \\ \alpha & \alpha & \alpha \\ \alpha & \alpha & \alpha & \alpha \end{array} \right)$$

$$\vec{\pi} = (1/4, 1/4, 1/4, 1/4)$$

II Felsenstein

$$\begin{matrix} A \\ C \\ G \\ T \end{matrix} \left(\begin{array}{c} \alpha \\ \alpha & \alpha \\ \alpha & \alpha & \alpha \\ \alpha & \alpha & \alpha & \alpha \end{array} \right)$$

$$\pi_i \neq \pi_j$$

III Kimura 2-parameter Model

$$\begin{matrix} A \\ C \\ G \\ T \end{matrix} \left(\begin{array}{c} \cdot \\ \beta & \cdot \\ \alpha & \beta & \cdot \\ \beta & \alpha & \beta & \cdot \end{array} \right)$$

$$\vec{\pi} = (1/4, 1/4, 1/4, 1/4)$$

IV HKY

$$\begin{matrix} A \\ C \\ G \\ T \end{matrix} \left(\begin{array}{c} \cdot \\ \beta & \cdot \\ \alpha & \beta & \cdot \\ \beta & \alpha & \beta & \cdot \end{array} \right)$$

$$\pi_i \neq \pi_j$$

IV GTR General time reversible model

$$\left(\begin{array}{c} \alpha & \cdot \\ \beta & \delta & \cdot \\ \gamma & \varepsilon & ? & \cdot \end{array} \right)$$

$$\pi_i \neq \pi_j$$

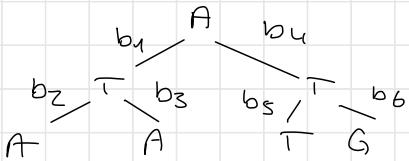
\rightarrow all rates relative to each other

STR/GTR

Computing likelihood

$P_{i,j}(t) \Rightarrow$ probability of starting in state $i \in \{A, C, T, G\}$ and after time t being in state j

1) Assume inner states of tree known:

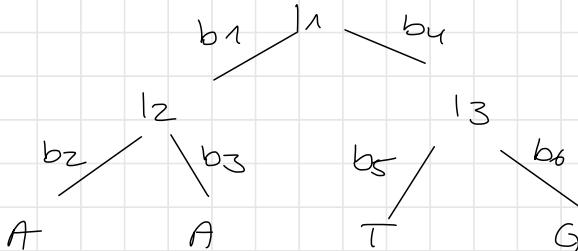


$b_1 \Rightarrow$ from A to T in time b_1

$$\sim \prod_A P_{A,T}(b_1) P_{T,A}(b_2) \cdot P_{T,A}(b_3) \cdot P_{A,T}(b_4) \cdot P_{T,T}(b_5) \cdot P_{T,G}(b_6)$$

↑
likelihood of A as root

2) Now: inner states unknown



\rightarrow Likelihood of all possible combinations:

$$L \left(\begin{array}{c} A \\ \diagup \quad \diagdown \\ A \quad A \end{array} \right) + L \left(\begin{array}{c} A \\ \diagup \quad \diagdown \\ A \quad C \end{array} \right) + \dots$$

$$\Rightarrow L = \underbrace{\sum_{l_1=A}^T \sum_{l_2=A}^T \sum_{l_3=A}^T \prod_{n=1}^N P_{l_1, l_2}(b_1) P_{l_2, l_3}(b_4) \dots}_{\text{AND}} \quad \text{OR}$$

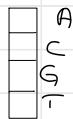
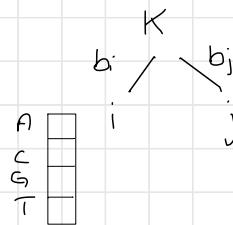
Combinatorial problem: n unknown nodes $\Rightarrow 4^n$ possibilities

\Rightarrow dynamic programming: Conditional likelihood vector

Calculate CLV of K with respect to

i, j

$$CLV_K(A) = \left(\sum_{S_i=A}^T P_{A \rightarrow S_i} (b_i) \cdot CLV_i(S_i) \right) + \left(\sum_{S_j=A}^T P_{A \rightarrow S_j} (b_j) \cdot CLV_j(S_j) \right)$$



CLV of leaves is e.g. leaf = A \rightarrow

1.0
00
00
00

\Rightarrow Likelihood of the root given CLV_{root}

$$L = \sum_{S=A}^T \pi_S \cdot CLV_{\text{root}}(S)$$

Computing transition matrix P

Introduce Q matrix

$$Q = \text{diag}(\vec{\pi}) R \rightarrow \text{time reversability for } R : \pi_i \cdot d_{ij} = \pi_j \cdot d_{ji}$$

$$d = \begin{pmatrix} -3\lambda & \lambda & \lambda & \lambda \\ \lambda & -3\lambda & \lambda & \lambda \\ \lambda & \lambda & -3\lambda & \lambda \\ \lambda & \lambda & \lambda & -3\lambda \end{pmatrix} \quad (\det = 0 !) \quad \text{mit } * = \lambda, \alpha = 1 \text{ for } J \subset N.$$

$$\Rightarrow P(t) = e^{Qt}$$

$Q = UVU^{-1} \rightarrow EV\text{-decomposition}, Q \text{ not symmetric}$
 \rightarrow non real values

$$\sim \text{define symmetric } Q' = \text{diag}(\sqrt{\vec{\pi}})^{-1} \cdot Q \cdot \text{diag}(\sqrt{\vec{\pi}}) \\ = U' V (U')^T$$

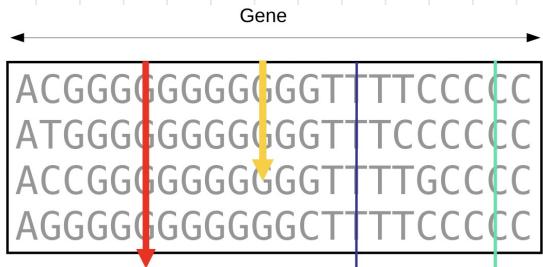
$$\Leftrightarrow Q = \underbrace{\left(\text{diag}(\Gamma^T \cdot U) \right)}_U \cdot V \cdot \underbrace{\left(U^{T^{-1}} \text{diag}(\Gamma^T)^{-1} \right)}_{U^{-1}}$$

$$\begin{aligned} \Rightarrow P(t) &= e^{Qt} = \sum_{k=0}^{\infty} \frac{1}{k!} (Qt)^k \\ &= \sum_{k=0}^{\infty} \frac{1}{k!} (UVU^{-1})^k \cdot t^k \\ &= U \left(\sum_{k=0}^{\infty} \frac{1}{k!} V^k t^k \right) U^{-1} \\ &= U e^{\text{diag}(\lambda) t} U^{-1} \end{aligned}$$

$\text{diag}(\lambda_i)$ = diag(Eigenvalues)

Among Site Rate Heterogeneity

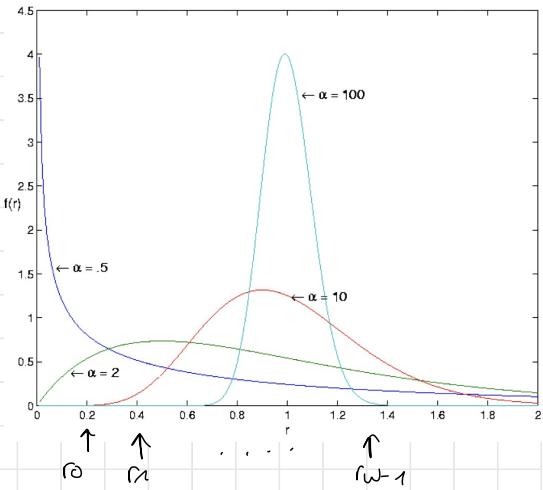
different sites evolve at different speeds



→ model using Γ -distribution

small α : rates differ significantly

→ high rate heterogeneity



e.g.

Discrete Γ -model

discret. Γ distribution at ω rates e.g. $\omega=4$

\rightarrow log likelihood at site i is then

$$\ln(L(i)) = \log\left(\frac{1}{\omega} * (L_0 + L_1 + \dots + L_{\omega-1})\right)$$

$$\text{rate } 0 \\ P(t) = e^{-Q_r t}$$

...

$$\text{rate } \omega-1 \\ P(t) = e^{-Q_r \omega-1 t}$$

Mixture Models

Γ model: 4 discrete evolutionary rates

\hookrightarrow consider all at once without assigning specific rate to a specific site

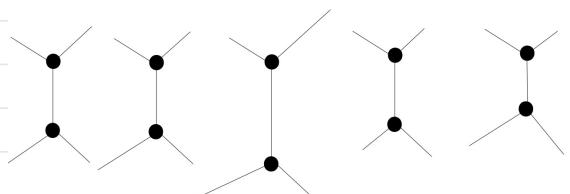
partitioned Dataset

use GTR for each Gene

\hookrightarrow results in same tree topology per gene

Gene 0	Gene 1	Gene 2	Gene 3	Gene 4
a_0 GTR ₀	a_1 GTR ₁	a_2 GTR ₂	a_3 GTR ₃	a_4 GTR ₄

\rightarrow increases # of params in likelihood model by
 $\hookrightarrow + (2n-3)$



Model Testing

more params \Rightarrow improved likelihood

\rightarrow but overparametrized?

Likelihood ratio test if models nested

Models A nested in model B if A's params subset of B's params

$$\rightarrow LR = P(D|A) / P(D|B) = L(A) / L(B) \quad \text{e.g.: can only compare likelihoods for same data}$$
$$\Delta = \ln(LR^2) = 2(\ln(L(A)) - \ln(L(B)))$$

\rightarrow compare Δ to χ^2 distribution with deg. of fr. $k_A - k_B$ to determine if A is significant

Bootstrapping

Goal: find best tree \rightarrow to what degree does Data support resulting tree?

Creating Replicas

given alignment A with length n

- resample n columns/sites (ziehen mit Zurücklegen)
- create tree (via Max. Likl. / Max. Pars.)
- do this n times
- => n trees > bootstrap trees

Probabilities

- prob. of site being sampled per draw : $\frac{1}{n}$
- prob. of site not - " - : $1 - \frac{1}{n}$
- prob. of site not being sampled in any of the n draws : $(1 - \frac{1}{n})^n$
- prob. of site being sampled at least once in n draws
 $\rightarrow \Pr(X \geq 1) = 1 - (1 - \frac{1}{n})^n \approx 1 - \frac{1}{e} \approx 63,2\%$.

\hookrightarrow expectation: searches on alignment and replicates are somehow different

Bipartitions: Inner Branches

· inner branch = smaller unit of evolutionary relationship

AB|CDEF \rightarrow A, B more closely related to each other than taxa C, D, E, F

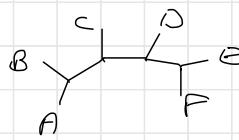
· trivial bipartition: outer taxa A|BCDEF

possible bipartitions for

AB | CDEF

ABC | D E F

ABCD | E F



Algorithm

- 1) Infer ML tree on full alignment
- 2) Create bootstrap tree set
- 3) Extract bipartitions from ML tree and bootstrap tree set
- 4) Annotate ML tree with relative frequency of its bipartitions in the bootstrap tree set

Consensus of Trees with consensus threshold t

- 1) extract bipartitions from tree set
- 2) determine consensus bipartitions \rightarrow occur in $\geq t\%$ of trees
- 3) transform consensus biparts into consensus tree

Consensi Flavors

- majority rule consensus : freq. $> 50\%$.
- strictly consensus : freq. $> 95\%$.
- extended majority rule consensus

Extended majority rule consensus

Bipartitions b_1, b_2 cannot be in same tree if incompatible

$b_1 = B \bar{B}$, $b_2 = C \bar{C}$ compatible if
 $(B \cap C = \emptyset) \vee (B \cap \bar{C} = \emptyset) \vee (\bar{B} \cap C = \emptyset)$

\rightarrow Algo : given consensus Bips B_C , non-consensus Bips B_n

- 1) remove most frequent bip $b \in B_n$ from B_n
- 2) if b compatible to all $c_i \in B_C$: $B_C = B_C \cup \{b\}$
- 3) loop until $|B_C| = (n-3)$ or $B_n = \emptyset$

Hashing

- Store biparts. in hash-table
 - compute biparts. from adjacent biparts
$$a = 001100 = CD | AB EF$$

$$b = 110000 = AB | CDEF$$

$$? = a \cup b = 111100 = ABCD | EF$$

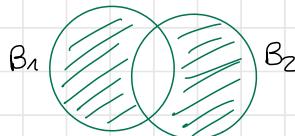
- $r[J]$ one random number per taxon
 - ↳ hash: xor of all corresp. random numbers
 - hash(a) = $r[C] \text{ xor } r[D]$
 - hash(b) = $r[A] \text{ xor } r[B]$
 - hash(?) = hash(a) xor hash(b)

Distances between trees

Robinson - Foulds Distance (metric)

Given unrooted trees T_1, T_2 with biparts B_1, B_2

$$RF(T_1, T_2) = |B_1 \cup B_2| - |B_1 \cap B_2|$$

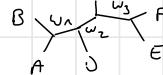


$T_1:$ ~~ABICD ϵ F~~
~~A β BO γ C δ E ϵ F~~
ABOE β FC
 $\underbrace{\hspace{1cm}}$
 $\beta\lambda$

Flavors

$\Delta = \text{Symmetric set difference } \subset X \Delta Y = (X \setminus Y) \cup (Y \setminus X)$

$w_i = \text{Support value}$



	unweighted	weighted
absolute	$ B_1 \cup B_2 - B_1 \cap B_2 $	$\sum_{b \in B_1 \Delta B_2} w(b)$
relative	$\frac{ B_1 \cup B_2 - B_1 \cap B_2 }{2(n-3)}$	$\frac{\sum_{b \in B_1 \Delta B_2} w(b)}{2(n-3)}$

How many replicates?

- 1) Randomly split tree in subsets t' , t''
- 2) Compute consensus trees $c(t')$, $c(t'')$
- 3) Compute weighted RF distance $wRF(c(t'), c(t''))$
- 4) Repeat 1) - 3) with cutoff values l, m
if less than l permutations have $wRF < m$:
stop \rightarrow stable support values

Triplets / Quartets

alternative to biparts.

\rightarrow triplets for rooted
quartets for unrooted trees

\sim for 4 taxa: only 3 possible trees

$\sim O(n^4)$ to extract all quartets

Bayesian Methods

- Sample posterior probability distribution
- ML: searches best peak value

Probabilities

		A		
		a_1	a_2	Σ
B	b_1	•	.	•
	b_2	.	.	.
Σ				

Joint probability $Pr(a_1, b_1)$: prob. of observing a_1 and b_1

Marginal probability: $Pr(b_1) = Pr(b_1, a_1) + Pr(b_1, a_2)$

Conditional probability: Prob. of observing a_1 given b_1 has occurred

$$Pr(a_1 | b_1) = \frac{Pr(a_1, b_1)}{Pr(b_1)}$$

Statistical independence: two events A, B independent if $Pr(A, B) = Pr(A) \cdot Pr(B)$

Bayes Theorem

$$Pr(A, B) - \frac{Pr(B|A)}{Pr(A)} \Rightarrow Pr(B|A) = \frac{Pr(A, B)}{Pr(A)}$$

↑ ↗
unobserved observed
outcome outcome

likelihood prior probability

$$Pr(B|A) = \frac{Pr(A|B) \cdot Pr(B)}{Pr(A)}$$

↑ ↘
posterior marginal probability
prob.

→ Phylogenetics

does alignment fit tree + model?

$$\Pr(\text{Alignment} | \text{Tree}, \text{Params}) \cdot \Pr(\text{Tree}, \text{Params})$$

—————
pr(Alignment)

$$\Pr(\text{Tree}, \text{Params} | \text{Alignment}) =$$

—————

distr. over all possible trees and all model parameter values

- * introduces prior knowledge/assumptions about prob. distn. of trees + params (GTR rates/...)

$n+1 =$ # of possible trees

$$\Pr(\text{Alignment}) = \Pr(\text{Alignment}, t_0) + \dots + \Pr(\text{Alignment}, t_n)$$

$$= \Pr(\text{Alignment}(t_0)) \cdot \Pr(t_0) + \dots + \Pr(\text{Alignment}(t_n)) \Pr(t_n)$$

→ General Bayes Theorem

$$f(\theta | A) = \frac{f(A | \theta) f(\theta)}{\int f(\theta) f(A | \theta) d\theta}$$

$\theta =$ tree topology, model params, branch lengths, ...

$f =$ density function

$$f(x) > 0 \quad \forall x$$

$$\int f(x) = 1$$

$$\Pr(0.2 \leq x \leq 0.3) = \int_{0.2}^{0.3} f(x)$$

Monte-Carlo Methods

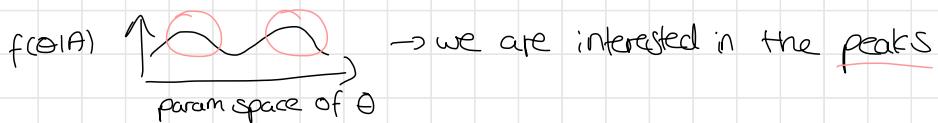
how to compute $\int f(\theta) f(A|\theta) d\theta$?

Idea: approximate with samples

$$\int f(\theta) d\theta = \frac{1}{N} \sum_i f(\theta_i) \quad \text{with } \theta \in \mathbb{R}^m \text{ (domain)}$$

$$\hookrightarrow \text{domain volume } V \cdot \frac{1}{N} \cdot \sum_i f(\theta_i)$$

→ problem: pure MC samples completely random



Markov Chain Monte Carlo Methods

higher sample density in interesting regions

→ biased random walks: probability to find sample in area with high posterior prob. ($f(\theta|A)$) is proportional to posterior distribution

Robot metaphor

- robot has to explore unknown landscape
- Altitude = probability
- generate step to take → evaluate whether to take it

accept/reject step from $P_1 \rightarrow P_2$ based on ratio R of their posterior densities

$$R = \frac{\Pr(P_2 | \text{data})}{\Pr(P_1 | \text{data})} = \frac{\Pr(P_2) \Pr(\text{data} | P_2)}{\Pr(P_1) \Pr(\text{data} | P_1)}$$

$$= \frac{\Pr(P_2) \Pr(\text{data} | P_2)}{\Pr(P_1) \Pr(\text{data} | P_1)}$$

$$\frac{Pr(P_2)}{Pr(P_1)} = \text{prior ratio} \rightarrow \text{for uniform priors} = 1$$

$$\frac{Pr(\text{data} | P_2)}{Pr(\text{data} | P_1)} = \text{likelihood ratio}$$

Target distribution = posterior distr. \rightarrow the one to sample

proposal distribution: decides which point to randomly go to next

- small variance pro: seldom refuses step
con: smaller steps, more steps required

\hookrightarrow difficult to find!

dataset dependent, trial & error

Metropolis Algorithm (for Phylogenetics)

Init: choose random tree with random branch lengths as first sample

Loop:

- propose either
 - new tree topology
 - new branch length

- calculate acceptance ratio of proposal
- accept/reject change
- print current tree + branch lengths for every k iters
 - = thinning

Finish: Summarize Sample using means, histograms, ...

Hastings Correction

asymmetric proposal distribution \Rightarrow bias

$$\rightarrow \text{Hastings ratio } HR = \frac{Q(P_1|P_2)}{Q(P_2|P_1)}$$

if Q symmetric: $Q(P_1|P_2) = Q(P_2|P_1) \Rightarrow HR=1$

$$\Rightarrow R = \frac{\Pr(P_2)}{\Pr(P_1)} \cdot \frac{\Pr(\text{data}|P_2)}{\Pr(\text{data}|P_1)} \cdot HR$$

Bipartitions

posterior prob. of bipart. AB/CD/E

\rightarrow count relative occurrence in generated samples

\hookrightarrow approx. true proportion if MC converged!

Convergence

- can't say that MCNC -chain has converged
 \hookrightarrow can only say it hasn't

\Rightarrow run multiple chains

Metropolis-coupled MCNC Methods

Cold chain: sees landscape as is

Hot chain: sees flatter landscape version \rightarrow moves more quickly between peaks

\hookrightarrow Run several chains simultaneously

1 cold chain

several heated chains

'flattened' acceptance Ratio: $R^{1/H}$ with $H = \text{temperature}$

cold chain: $H=0$

Prior Probabilities

- convey scientist's beliefs before having seen the data
- can bias analysis

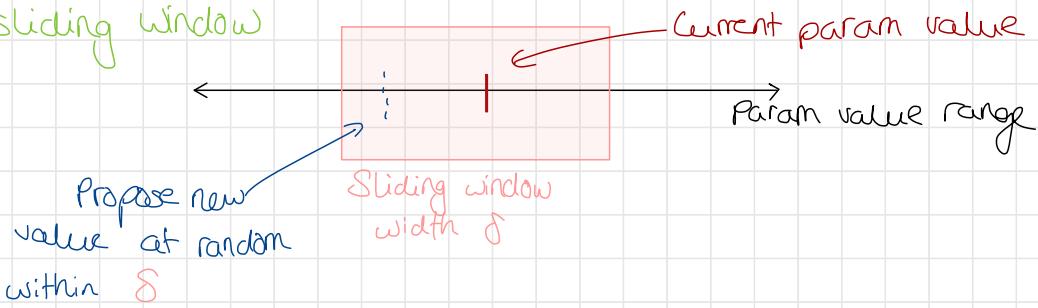
Proposal Mechanisms

- Univariate params + branch lengths:
 - sliding window
- branch lengths:
 - node slider
- topologies:
 - local proposal

Proposal Requirements:

- either no Hastings Ratio needed
- or it can be calculated
- acceptance rate $\approx 25\%$.

Sliding Window



- Hastings Ratio = 1
- δ can be auto-tuned \rightarrow acceptance rate $1/4$ can be obtained
- used for α -shape param in Γ function for rate heterogeneity

Node Slider

- 1) Pick 2 contiguous branches randomly b_1, b_2
- 2) Multiply both by same random number
- 3) Propose new branch ratio b_1/b_2 at random

! HR ≠ 1

Tree proposal

- 1) Pick 3 contiguous branches at random that define 2 subtrees X, Y
- 2) Shrink / grow selected 3 branch segments by random amount
- 3) Prune X or Y (selected rand.)
- 4) Reinsert pruned tree at random position in the 3 branch segment

Different Models NCMC

goal: integrate over different models

- encode rate configurations for models as strings:
 $111111 = \text{JC Model}$
 $123456 = \text{CTR Model}$

→ Bell number: given n models there are B_n combinations possible

$$B_{n+1} = \sum_{k=0}^n \binom{n}{k} B_k$$

1) Model proposal

Split move: choose set of substitution rates with >1 member

1 1 1 2 2 2 (2-param model)

→ split it randomly into two rates

1 1 1 2 2 3 (3-param model)

merge move: choose 2 subst. rate sets

111 223

→ merge into one subst. rate set

111 222

2) Sampling Different Models

→ cannot compare likelihoods for different models

↪ use reversible jump MCMC to jump btw. models with different # of params

⇒ proposal ratio calculation changed

→ acceptance ratio:

$$r = \text{likelihood ratio} * \text{prior ratio} * \text{proposal ratio} * \text{Jacobian}$$

Jacobian defines linear map $\mathbb{R}^n \rightarrow \mathbb{R}^m$ at point x if $f(x)$ at x differentiable

→ posterior probability for one spec. model:

fraction of time (samples) the MCMC chain visited within that model

Evolution = change over time through inheritance

Model param: population size constant or time variant?
→ lecture: constant

Units of Evolution

Population

evolves due to changing frequency of features of individuals

4 main forces:

Genetic Drift (chance/random events)



Gene

inherited from generation to generation

inheritance via **Alleles**

Allele = specific form a or A of gene **A**

diploid: 2 sets of corresponding chromosomes (homologous)

if allele sequences in the 2 chrom. identical: **homozygous**
- " -

different: **heterozygous**

↳ cause phenotyp. differences

humans inherit one allele from father, one from mother

Genotype = set of corresponding alleles in diploid organism

Phenotype = observation for the trait/property the gene controls

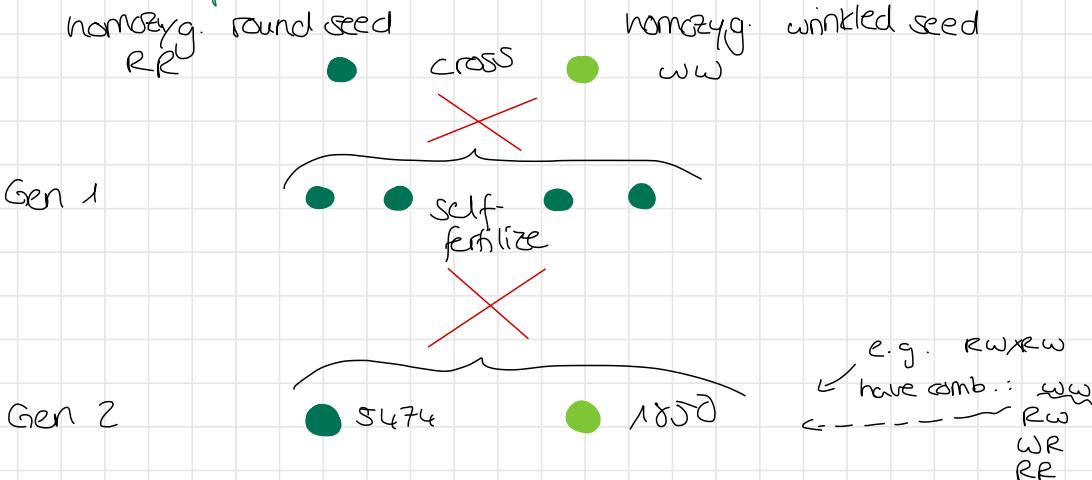
Mendelian inheritance

Dominance

- genotype $RW \rightarrow R$
 $\Rightarrow R$ is dominant
 w is recessive

if no dominance-recession relationship: intermediate phenotype

Mendel Experiment



\Rightarrow traits of 2 parent plants don't blend

\rightarrow both traits passed to Gen 2 without showing in Gen 1

Allele Inheritance

Gene = **A**

corresp. alleles = a/A if only 2 else A_1, A_2, \dots

Polymorphic gene

Gene **A** is polymorphic when there exist multiple alleles

Polymorphic site = SNP ('Snip')

after NGS: observe sites in certain genes with more than one state
 - called Snip SNP = Single Nucleotide Polymorphism

In population genetics: few mutations \Rightarrow MSA has to be correct!

Does dominance affect allele frequency?

assumptions:

- infinite population size
- random mating (\hookrightarrow independent of phenotypes)
- gene A has 2 alleles A, a

given current frequencies at generation 0 of alleles defining genotype

$$f_0(A) = p$$

$$f_0(a) = q$$

$$\text{with } p+q=1$$

$$\text{Generation } t: \quad f_t(A) = f_t(AA) + \frac{1}{2} f_t(Aa)$$

$$f_t(a) = f_t(aa) + \frac{1}{2} f_t(Aa)$$

\rightarrow how to compute $f_0 \rightarrow f_1$?

Punnett Square

		females		\leftarrow occurrence of a with frequency q
		A(p)	a(q)	
males	A(p)	AA(p ²)	Aa(p,q)	
	a(q)	Aa(p,q)	aa(q ²)	

Σ

random mating

$$\rightarrow \underline{f_1(AA) = p^2}$$

$$\underline{f_1(aa) = q^2}$$

$$\underline{f_1(Aa) = 2pq}$$

$$\Rightarrow f_0(AA) \neq f_1(AA)$$

! $f_0(AA) \neq f_0(AA)$

for $t > 0$:

$$f_t(AA) = f_1(AA)$$

$$f_t(aa) = f_1(aa)$$

$$f_t(Aa) = f_1(Aa)$$

$$\Rightarrow f_t(AA) = f_1(AA) + \frac{1}{2}f_1(Aa) = p^2 + pq = p(p+q) = p = f_0(AA)$$

$q=1-p$
↓

→ frequency of allele occurrence does not change
 proportion of genotypes AA:Aa:aa constant after gen 1

Hardy-Weinberg Equilibrium

- null hypothesis in evol. biol.
- What happens to allel freq. without influence of the 4 evolutionary forces

Finite Population Sizes

- N individuals in diploid population $\Rightarrow 2N$ chrom.
- Frequency of allel A = p

Random Genetic Drift

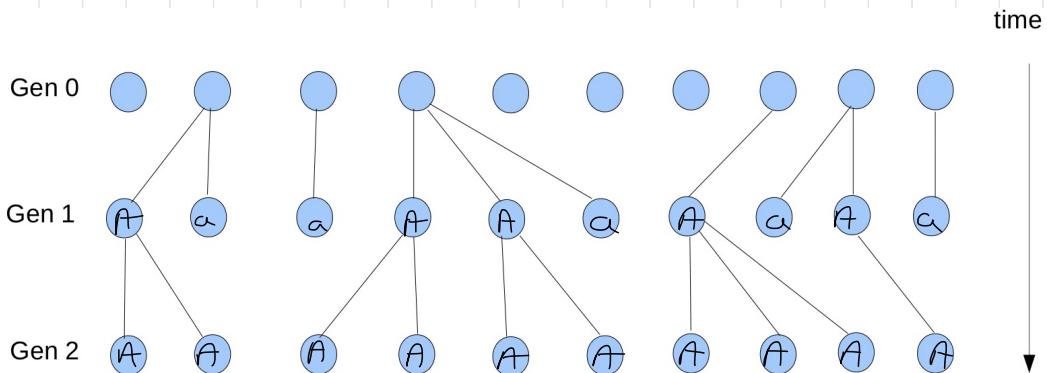
- Some alleles passed to next generation disproportionately without being advantageous/harmful
- happens esp. often in small populations (sampling errors)

Wright-Fisher Model

- population size N ($2N$ chromosomes) constant per generation
- random mating
- non-overlapping generations
- no natural selection
- equal distribution of sexes

Simulation

- each individual from offspring picks random parent
→ all parents equally likely to be picked
parent can be picked > 1 times
- each offspring inherits genetic info of the parent



→ 'a' vanished due to random parent picking

probability to pick allele 'A' as parent is

$$p = \frac{\#A}{2N}$$

if population remains constant: sample $2N$ times from the current generation to construct next gen. with $2N$ offspring

→ p remains constant

→ probability to pick 'A' as ancestor k times

$$P(X=k) = \binom{n}{k} p^k (1-p)^{n-k} \quad P = \frac{\#A}{2N}; n = 2N$$

$$\mathbb{E}(\#A) = 2N * p$$

$$\text{Var}(\#A) = 2N * p * (1-p)$$

! next state depends only on current state \Rightarrow model as Markov Chain

Transition Probabilities

Probability of changing from : alleles in gen t to j alleles in gen $t+1$:

$$\text{Prob } \{ X(t+1) = j | X(t) = i \}$$

$$= p_{ij} = \binom{2N}{j} \left(\frac{i}{2N}\right)^j \left(1 - \frac{i}{2N}\right)^{2N-j}$$

$$i, j = 0, 1, \dots, 2N$$

($2N \rightarrow$ either N diploid organisms or $2N$ haploid)

Absorbing State in Markov Chain

· probability to exit this state is 0

· Wright-Fisher model: only one allele left

-> given frequency of allele A $\neq A/2N$

probability that A will become fixed: $\neq A/2N$

Heterozygosity and Genetic Drift

· Reduction of polymorphism quantified by degree of homozygosity

ς = probability that 2 alleles are identical

Het_t = heterozyg. at gen. t

$$= Het_{t-1} (1 - 1/2N)$$

$$= Het_0 (1 - 1/2N)^t$$

$1 - \frac{1}{2N}$ = prob. that 2 randomly chosen alleles differ

Mutation - Drift Balance

- genetic drift removes polymorphism from population
- mutations introduce polymorphism

Balance

- Loss of heterozygosity per generation due to genetic drift:

$$-\frac{1}{2N} \cdot \text{Het}$$
- μ = mutation rate per gene (! 2 alleles per gene) per generation
- gain of heteroz. due to mutation:

$$2\mu(1 - \text{Het})$$

- Pick two alleles

- Consider transition from generation $t \rightarrow t+1$
- Probability that they are identical: $1 - \text{Het}_t$
- If identical: probability that one of them will mutate:

$$2\mu$$

$$\Rightarrow \text{Het}_{t+1} = \text{Het}_t - \frac{1}{2N} \cdot \text{Het}_t + 2\mu(1 - \text{Het}_t)$$

$$\Delta \text{Het} = -\frac{1}{2N} \cdot \text{Het}_t + 2\mu(1 - \text{Het}_t)$$

$$\Delta \text{Het} = 0 \Rightarrow \text{Het} = \frac{4\mu N}{1 + 4\mu N}$$

Mutation Rate

- μ = mutation prob. per gen. per individual
- $2N$ indiv. $\rightarrow 2N\mu$ mutations per gen.
- probability that mutation will be fixed = $\frac{1}{2N}$

\Rightarrow rate of mutation arising and being fixed in population

$$\frac{1}{2N} \cdot 2N\mu = \mu$$

Natural Selection

- fitness = ability of an individual to survive + reproduce
 - depends on genotype
- selection = one genotype reproduces more freq. than others
- if genotype has better fitness : (e.g. AA)
 - will fix in the population after several gen.
 - ⇒ allele A will also fix
 - ~ natural selection has favored allele A
 - ~ natural selection on A = Positive Selection

Models of Selection

