

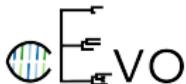
Computational Biology

Lecturers:
Tanja Stadler, Carsten Magnus & Tim Vaughan

Teaching Assistants:
Jūlija Pečerska, Jérémie Sciré,
Sarah Nadeau & Marc Manceau

Computational Evolution
Department of Biosystems Science and Engineering

HS 2019



04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

GWAS, molecular evolution: Questions

- ② Does the odds ratio in a GWAS prove that a minor variant at a SNP position causes a genetic disease?
- ② In a GWAS, why can you not reject your null hypothesis if the p-value is $< \alpha$?

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

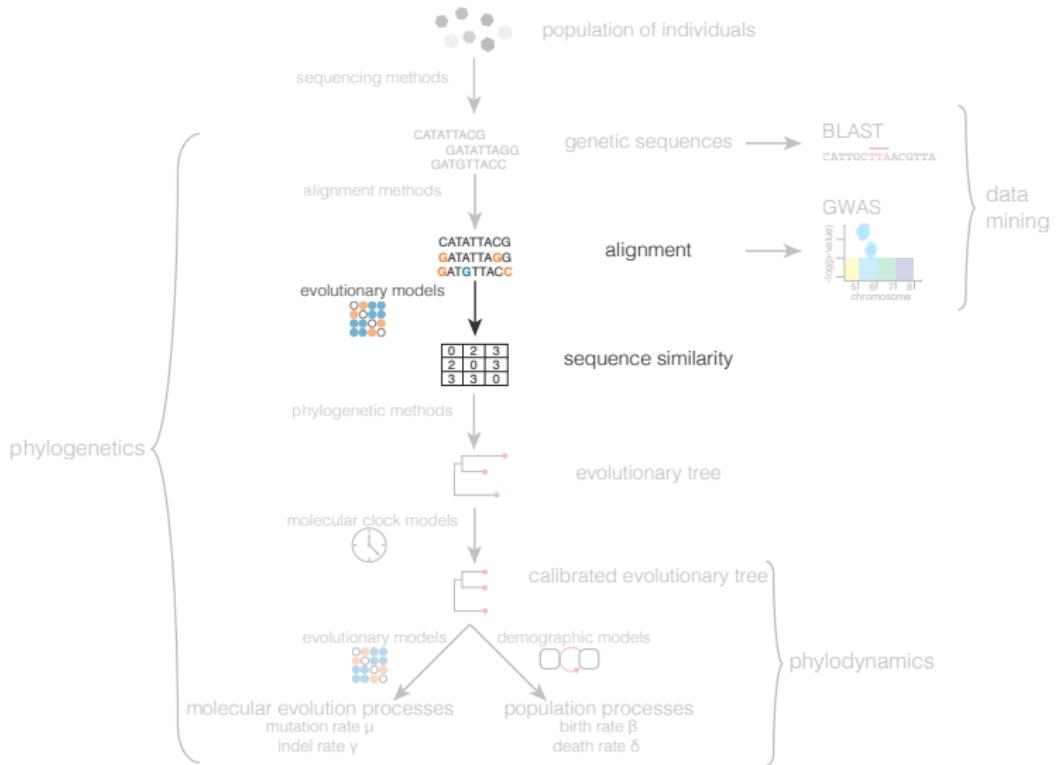
The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Overview



04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

What we will learn today

- ▶ nucleotide substitution models
- ▶ distance estimation under JC69
 - ▶ maximum likelihood approach
- ▶ modelling variable substitution rates across the genome
- ▶ amino acid substitution models
- ▶ codon substitution models

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Summary Markov chain model for nucleotide substitutions.

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices
Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance
Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Why Markov chains are a great model for nucleotide substitutions

- ▶ memorylessness: a nucleotide substitution happens independently from the substitution history at this site
- ▶ substitution rate matrix defines the transition probabilities
- ▶ the transition probabilities take into account every possible substitution path

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Why Markov chains are a great model for nucleotide substitutions

- ▶ memorylessness: a nucleotide substitution happens independently from the substitution history at this site
- ▶ substitution rate matrix defines the transition probabilities
- ▶ the transition probabilities take into account every possible substitution path
- ▶ applying theories of linear algebra we can calculate the transition probability matrix according to:

$$P(t) = e^{Qt} = U \text{diag}(e^{\epsilon_1 t}, e^{\epsilon_2 t}, e^{\epsilon_3 t}, e^{\epsilon_4 t}) U^{-1}$$

- ▶ cookbook recipe provided on moodle
- ▶ for further information on how to derive this formula please consult a textbook on linear algebra

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

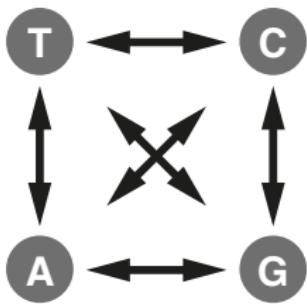
References

Substitution rate matrices and their properties.

The easiest substitution model: JC69

JC69:

- ▶ named after TH Jukes, CR Cantor: Evolution of protein molecules. 1969 [Jukes and Cantor, 1969].
- ▶ all substitution have the same rate, λ



Substitution rate matrix:

$$\begin{matrix} & \text{T} & \text{C} & \text{A} & \text{G} \\ \text{T} & \cdot & \lambda & \lambda & \lambda \\ \text{C} & \lambda & \cdot & \lambda & \lambda \\ \text{A} & \lambda & \lambda & \cdot & \lambda \\ \text{G} & \lambda & \lambda & \lambda & \cdot \end{matrix}$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

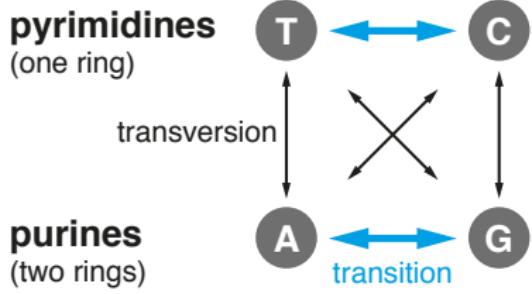
Synonymous and non-synonymous substitutions

References

Accounting for transition/transversion: K80

K80:

- ▶ named after M Kimura: A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide sequences. 1980. [Kimura, 1980]
- ▶ transitions happen at rate α , transversions at rate β



Substitution rate matrix:

	T	C	A	G
T	.	α	β	β
C	α	.	β	β
A	β	β	.	α
G	β	β	α	.

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Including equilibrium frequencies: TN93 and HKY

TN93:

- ▶ named after [Tamura and Nei, 1993]
- ▶ transitions between T \leftrightarrow C happen at rate $\alpha_1 \times$ (nucleotide equilibrium frequency)
- ▶ transitions from A \leftrightarrow G happen at rate $\alpha_2 \times$ (nucleotide equilibrium frequency)
- ▶ transversions happen at rate $\beta \times$ (nucleotide equilibrium frequency)

Substitution rate matrix:

$$\begin{array}{cccc} & T & C & A & G \\ T & \cdot & \alpha_1\pi_C & \beta\pi_A & \beta\pi_G \\ C & \alpha_1\pi_T & \cdot & \beta\pi_A & \beta\pi_G \\ A & \beta\pi_T & \beta\pi_C & \cdot & \alpha_2\pi_G \\ G & \beta\pi_T & \beta\pi_C & \alpha_2\pi_A & \cdot \end{array}$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Including equilibrium frequencies: TN93 and HKY

TN93:

- ▶ named after [Tamura and Nei, 1993]
- ▶ transitions between T \leftrightarrow C happen at rate $\alpha_1 \times$ (nucleotide equilibrium frequency)
- ▶ transitions from A \leftrightarrow G happen at rate $\alpha_2 \times$ (nucleotide equilibrium frequency)
- ▶ transversions happen at rate $\beta \times$ (nucleotide equilibrium frequency)

Substitution rate matrix:

$$\begin{matrix} & \text{T} & \text{C} & \text{A} & \text{G} \\ \text{T} & \cdot & \alpha_1\pi_C & \beta\pi_A & \beta\pi_G \\ \text{C} & \alpha_1\pi_T & \cdot & \beta\pi_A & \beta\pi_G \\ \text{A} & \beta\pi_T & \beta\pi_C & \cdot & \alpha_2\pi_G \\ \text{G} & \beta\pi_T & \beta\pi_C & \alpha_2\pi_A & \cdot \end{matrix}$$

▶ Note that the special case of $\alpha_1 = \alpha_2$ was described earlier [Hasegawa et al., 1984] and is named HKY

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

A more general substitution model: GTR

GTR (REV):

- ▶ generalised time-reversible model
- ▶ based on three papers:
[Tavaré, 1986, Yang, 1994, Zharkikh, 1994]

Substitution rate matrix:

$$\begin{matrix} & \text{T} & \text{C} & \text{A} & \text{G} \\ \text{T} & \cdot & a\pi_C & b\pi_A & c\pi_G \\ \text{C} & a\pi_T & \cdot & d\pi_A & e\pi_G \\ \text{A} & b\pi_T & d\pi_C & \cdot & f\pi_G \\ \text{G} & c\pi_T & e\pi_C & f\pi_A & \cdot \end{matrix}$$

- + quite flexible
- + time-reversible
- not completely general

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Time reversibility

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

$$\begin{aligned}
 & \begin{pmatrix} \cdot & a\pi_C & b\pi_A & c\pi_G \\ a\pi_T & \cdot & d\pi_A & e\pi_G \\ b\pi_T & d\pi_C & \cdot & f\pi_G \\ c\pi_T & e\pi_C & f\pi_A & \cdot \end{pmatrix} \\
 = & \begin{pmatrix} \cdot & a & b & c \\ a & \cdot & d & e \\ b & d & \cdot & f \\ c & e & f & \cdot \end{pmatrix} \cdot \begin{pmatrix} \pi_T & 0 & 0 & 0 \\ 0 & \pi_C & 0 & 0 \\ 0 & 0 & \pi_A & 0 \\ 0 & 0 & 0 & \pi_G \end{pmatrix}
 \end{aligned}$$

The most general substitution model

UNREST:

- ▶ unrestricted model first described in [Yang, 1994]
- ▶ each substitution has a (different) rate

Substitution rate

matrix:

$$\begin{array}{ccccc} & T & C & A & G \\ T & \cdot & a & b & c \\ C & d & \cdot & e & f \\ A & g & h & \cdot & i \\ G & j & k & l & \cdot \end{array}$$

- + most general case
- + all other models are special cases of UNREST
- mathematically very complicated and not handy to use
- not time-reversible

- ▶ Further models described in [Yang, 1994]

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Substitution models

model	parameters	description
JC69	1	all substitutions have the same rate
K80	2	accounts for transition and transversions
HKY	2+3*	distinction between transition and transversions, including equilibrium frequencies
TN93	3+3*	different rates for transitions
GTR	6+3*	general, but still time-reversible
UNREST	12	most general, not time-reversible

* equilibrium frequencies of nucleotides

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Calculating transition probabilities and sequence distance.

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices
Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance
Variable substitution rates across sites

Amino acid substitution models

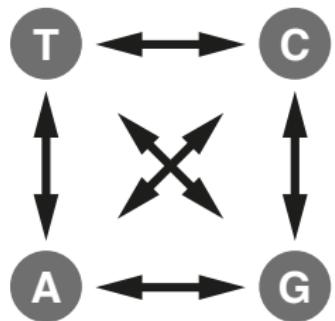
Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References



Substitution rate matrix:

$$Q = \begin{pmatrix} -3\lambda & \lambda & \lambda & \lambda \\ \lambda & -3\lambda & \lambda & \lambda \\ \lambda & \lambda & -3\lambda & \lambda \\ \lambda & \lambda & \lambda & -3\lambda \end{pmatrix}$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

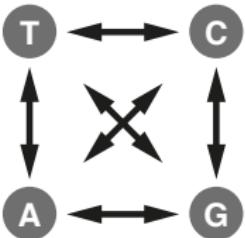
References

Example of transition probabilities: JC69

Substitution rate matrix:

$$Q = \begin{pmatrix} -3\lambda & \lambda & \lambda & \lambda \\ \lambda & -3\lambda & \lambda & \lambda \\ \lambda & \lambda & -3\lambda & \lambda \\ \lambda & \lambda & \lambda & -3\lambda \end{pmatrix}$$

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



04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

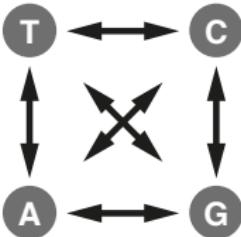
References

Example of transition probabilities: JC69

Substitution rate matrix:

$$Q = \begin{pmatrix} -3\lambda & \lambda & \lambda & \lambda \\ \lambda & -3\lambda & \lambda & \lambda \\ \lambda & \lambda & -3\lambda & \lambda \\ \lambda & \lambda & \lambda & -3\lambda \end{pmatrix}$$

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



transition probability matrix*:

$$P(t) = \begin{pmatrix} p_0(t) & p_1(t) & p_1(t) & p_1(t) \\ p_1(t) & p_0(t) & p_1(t) & p_1(t) \\ p_1(t) & p_1(t) & p_0(t) & p_1(t) \\ p_1(t) & p_1(t) & p_1(t) & p_0(t) \end{pmatrix}$$

$$\text{with } p_0(t) = \frac{1}{4} + \frac{3}{4}e^{-4\lambda t}$$

$$\text{and } p_1(t) = \frac{1}{4} - \frac{1}{4}e^{-4\lambda t}$$

*cookbook recipe provided on moodle (for a proof of the matrix exponentiation please consult a textbook on linear algebra)

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Example of transition probabilities: JC69

Substitution rate matrix:

$$Q = \begin{pmatrix} -3\lambda & \lambda & \lambda & \lambda \\ \lambda & -3\lambda & \lambda & \lambda \\ \lambda & \lambda & -3\lambda & \lambda \\ \lambda & \lambda & \lambda & -3\lambda \end{pmatrix}$$

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

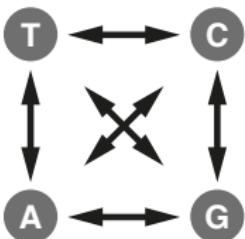
transition probability matrix*:

$$P(t) = \begin{pmatrix} p_0(t) & p_1(t) & p_1(t) & p_1(t) \\ p_1(t) & p_0(t) & p_1(t) & p_1(t) \\ p_1(t) & p_1(t) & p_0(t) & p_1(t) \\ p_1(t) & p_1(t) & p_1(t) & p_0(t) \end{pmatrix}$$

$$\text{with } p_0(t) = \frac{1}{4} + \frac{3}{4}e^{-4\lambda t}$$

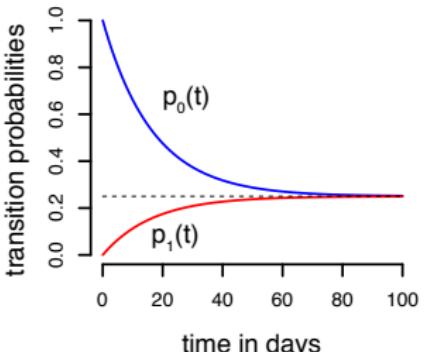
$$\text{and } p_1(t) = \frac{1}{4} - \frac{1}{4}e^{-4\lambda t}$$

*cookbook recipe provided on moodle (for a proof of the matrix exponentiation please consult a textbook on linear algebra)



Example:

$$\lambda = 0.015 \frac{\text{substitutions per site}}{\text{day}}$$



04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

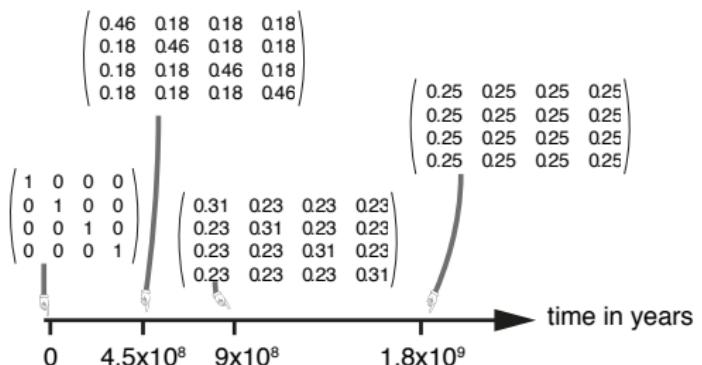
Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69: Stationary distribution

- ▶ sequence TCAG evolves according to JC69
- ▶ $\lambda = 2.2/3 \times 10^{-9}$ substitutions per site
year
- ▶ What is the probability to observe a certain nucleotide $a_i \in \{T, C, A, G\}$ at position $i \in \{1, 2, 3, 4\}$ after time t ?



- ▶ when $t \rightarrow \infty$ stationary distribution is reached
- ▶ Any long sequence (e.g. TTTTTT...) at time 0, will be composed of equal amounts of T,C,A,G after time $t \rightarrow \infty$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Maximum likelihood estimators of sequence distance.

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The idea behind maximum likelihood estimators

Example: We want to estimate the probability with which a six sided die shows 6. (Is the die loaded (unfair)?)



[loaded.dice, 2015]

- ▶ we throw the same die $n = 100$ times and note how many times we obtain a 6 (example: $x = 40$ times)

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The idea behind maximum likelihood estimators

Example: We want to estimate the probability with which a six sided die shows 6. (Is the die loaded (unfair)?)



[loaded.dice, 2015]

- ▶ we throw the same die $n = 100$ times and note how many times we obtain a 6 (example: $x = 40$ times)

Strategy:

- ▶ we use this data to obtain an estimate for the probability of throwing a 6
 - ▶ we denote the probability of throwing a 6 with p ; the probability of throwing x times 6 out of n tries is Binomial-distributed, i.e.
- $$P(x \text{ times } 6) = \binom{n}{x} p^x (1-p)^{n-x}$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The idea behind maximum likelihood estimators

Strategy:

- ▶ We use this data to obtain an estimate for the probability of throwing a 6
- ▶ We denote the probability of throwing a 6 with p . The probability of throwing x times 6 out of n tries is Binomial-distributed, i.e.

$$P(x \text{ times } 6) = \binom{n}{x} p^x (1-p)^{n-x}$$

- ▶ As we know x and n (this is our data), we can interpret this probability as a function of the unknown p :

$$L(p; x) = P(x|p) = \binom{100}{40} p^{40} (1-p)^{60}$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The idea behind maximum likelihood estimators

Strategy:

- ▶ We use this data to obtain an estimate for the probability of throwing a 6
- ▶ We denote the probability of throwing a 6 with p . The probability of throwing x times 6 out of n tries is Binomial-distributed, i.e.

$$P(x \text{ times } 6) = \binom{n}{x} p^x (1-p)^{n-x}$$

- ▶ As we know x and n (this is our data), we can interpret this probability as a function of the unknown p :

$$L(p; x) = P(x|p) = \binom{100}{40} p^{40} (1-p)^{60}$$

- ▶ **trick:** find the maximum of $L(p; x)$ in p (necessary condition: first derivative set to 0)

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The idea behind maximum likelihood estimators

Strategy:

- ▶ We use this data to obtain an estimate for the probability of throwing a 6
- ▶ We denote the probability of throwing a 6 with p . The probability of throwing x times 6 out of n tries is Binomial-distributed, i.e.

$$P(x \text{ times } 6) = \binom{n}{x} p^x (1-p)^{n-x}$$

- ▶ As we know x and n (this is our data), we can interpret this probability as a function of the unknown p :

$$L(p; x) = P(x|p) = \binom{100}{40} p^{40} (1-p)^{60}$$

- ▶ **trick:** find the maximum of $L(p; x)$ in p (necessary condition: first derivative set to 0)

Maximum likelihood estimator: is an estimator of a model parameter that maximises the probability to obtain the observed results

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

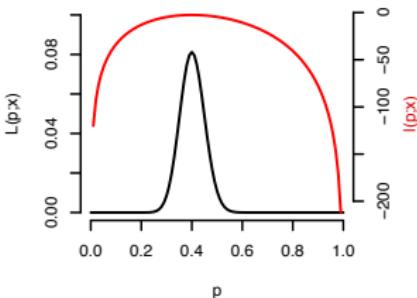
Synonymous and non-synonymous substitutions

References

Mathematical tricks to calculate the MLE

sometimes it is handy to apply some transformation to $L(p; x)$ that does not change the position of the maximum (but maybe the height):

- ▶ looking for the maximum of the log-likelihood function:
 $l(p; x) = \log(L(p; x))$
- ▶ multiplying L with a constant



04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Throwing dice

coming back to our example:

- ▶ throw a die $n = 100$ -times. $x = 40$ times the result is 6

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Throwing dice

coming back to our example:

- ▶ throw a die $n = 100$ -times. $x = 40$ times the result is 6

The likelihood function is:

$$L(p; x) = \binom{n}{x} p^x (1-p)^{n-x}$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance
Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Throwing dice

coming back to our example:

- ▶ throw a die $n = 100$ -times. $x = 40$ times the result is 6

The likelihood function is:

$$L(p; x) = \binom{n}{x} p^x (1-p)^{n-x}$$

Thus

$$l(p; x) = \log \binom{n}{x} + x \log p + (n - x) \log(1 - p)$$

Now we derive l with respect to p and set the derivative to 0:

$$\frac{dl}{dp}(\hat{p}; x, n) = \frac{x}{\hat{p}} - \frac{n-x}{1-\hat{p}} \stackrel{!}{=} 0 \quad \Leftrightarrow \hat{p} = \frac{x}{n}$$

This means that we estimate the probability to throw a 6 to be $\frac{40}{100} = 0.4$ which is much higher than $\frac{1}{6} = 0.1667$.

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Confidence intervals

- interval that tries to capture the uncertainty of a parameter estimate

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Confidence intervals

- interval that tries to capture the uncertainty of a parameter estimate

Definition: Confidence interval: If we repeatedly estimated the parameter from realisations of the random experiment and the interval estimate for each realisation of the random experiment, we could expect 95% of these intervals to contain the true parameter.

definition adapted from [Sokal and Rohlf, 2012]

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Confidence intervals

- interval that tries to capture the uncertainty of a parameter estimate

Definition: Confidence interval: If we repeatedly estimated the parameter from realisations of the random experiment and the interval estimate for each realisation of the random experiment, we could expect 95% of these intervals to contain the true parameter.

definition adapted from [Sokal and Rohlf, 2012]

- note that this interval is an estimate itself based on a realisation of a random experiment

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Confidence intervals

Confidence intervals for the parameter estimate can be calculated based on **likelihood intervals**:

- ▶ technique based on likelihood ratio test

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Confidence intervals

Confidence intervals for the parameter estimate can be calculated based on **likelihood intervals**:

- ▶ technique based on likelihood ratio test
- ▶ Let X be a random variable with a distribution parameterised in θ . Based on collected data x of a huge sample, the MLE for the parameter is $\hat{\theta}$. Then, one can show (under certain conditions) that, $2(l(\hat{\theta}) - l(\theta))$ has a χ_k^2 -distribution, where k are the degrees of freedom (the vector length of θ)

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Confidence intervals

Confidence intervals for the parameter estimate can be calculated based on **likelihood intervals**:

- ▶ technique based on likelihood ratio test
- ▶ Let X be a random variable with a distribution parameterised in θ . Based on collected data x of a huge sample, the MLE for the parameter is $\hat{\theta}$. Then, one can show (under certain conditions) that, $2(l(\hat{\theta}) - l(\theta))$ has a χ_k^2 -distribution, where k are the degrees of freedom (the vector length of θ)
- ▶ Idea: the 95% CI includes all θ s which are not in the 0.05 tail of the χ^2 distribution, i.e. which would not be rejected in a statistical test with $\alpha = 0.05$. This means that when performing the experiment e.g. 100 times, then 5 times the estimated parameter is expected not to be contained in the 95% CI.

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

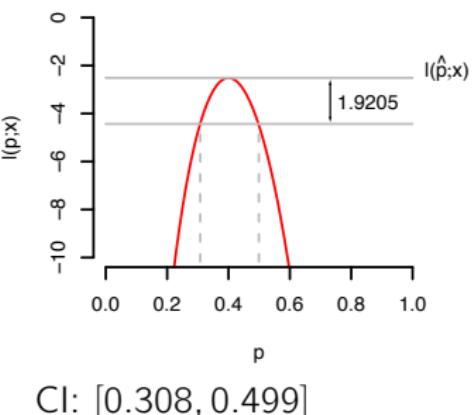
References

Confidence intervals

- Let X be a random variable with a distribution parameterised in θ . Based on collected data x of a huge sample, the MLE for the parameter is $\hat{\theta}$. Then, $2(l(\hat{\theta}) - l(\theta)) \sim \chi^2_k$.

Strategy to obtain a 95% confidence interval (CI):

- determine the value of the log-likelihood function in $\hat{\theta}$, $l(\hat{\theta}; x)$
- subtract $0.5\chi^2_{k,5\%}*$, i.e. calculate $l(\hat{\theta}; x) - 0.5\chi^2_{k,5\%}$
- determine those θ values for which $l(\theta; x) = l(\hat{\theta}; x) - 0.5\chi^2_{k,5\%}$



*A χ^2 table can be found on the moodle page.

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

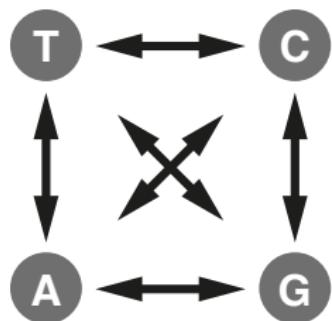
Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69: MLE for sequence distance

We will apply the MLE framework to derive an estimator for sequence distance under JC69.



Transition probability matrix $P(t) =:$

$$\begin{pmatrix} p_0(t) & p_1(t) & p_1(t) & p_1(t) \\ p_1(t) & p_0(t) & p_1(t) & p_1(t) \\ p_1(t) & p_1(t) & p_0(t) & p_1(t) \\ p_1(t) & p_1(t) & p_1(t) & p_0(t) \end{pmatrix}$$

$$\text{with } p_0(t) = \frac{1}{4} + \frac{3}{4} e^{-4\lambda t}$$

$$p_1(t) = \frac{1}{4} - \frac{1}{4} e^{-4\lambda t}$$

Imagine we have two sequences of length n with x differences. The probability that a position is different is $p = 3p_1(t)$. We define $d = 3\lambda t$ (the expected distance in time t).

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

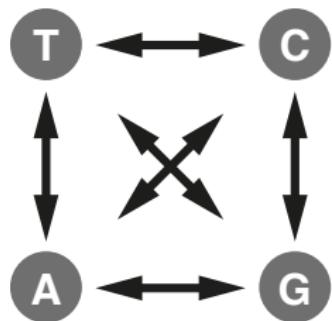
Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69: MLE for sequence distance

We will apply the MLE framework to derive an estimator for sequence distance under JC69.



Transition probability matrix $P(t) =:$

$$\begin{pmatrix} p_0(t) & p_1(t) & p_1(t) & p_1(t) \\ p_1(t) & p_0(t) & p_1(t) & p_1(t) \\ p_1(t) & p_1(t) & p_0(t) & p_1(t) \\ p_1(t) & p_1(t) & p_1(t) & p_0(t) \end{pmatrix}$$

$$\text{with } p_0(t) = \frac{1}{4} + \frac{3}{4} e^{-4\lambda t}$$

$$p_1(t) = \frac{1}{4} - \frac{1}{4} e^{-4\lambda t}$$

Imagine we have two sequences of length n with x differences. The probability that a position is different is $p = 3p_1(t)$. We define $d = 3\lambda t$ (the expected distance in time t).

Thus the probability that x positions out of n are different is:

$${n \choose x} p^x (1-p)^{n-x} = {n \choose x} \left(\frac{3}{4} - \frac{3}{4} e^{-\frac{4}{3}d} \right)^x \left(\frac{1}{4} + \frac{3}{4} e^{-\frac{4}{3}d} \right)^{n-x}$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

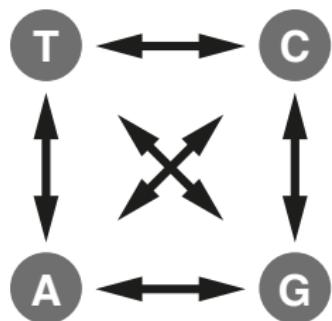
Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69: MLE for sequence distance

We will apply the MLE framework to derive an estimator for sequence distance under JC69.



Transition probability matrix $P(t) =:$

$$\begin{pmatrix} p_0(t) & p_1(t) & p_1(t) & p_1(t) \\ p_1(t) & p_0(t) & p_1(t) & p_1(t) \\ p_1(t) & p_1(t) & p_0(t) & p_1(t) \\ p_1(t) & p_1(t) & p_1(t) & p_0(t) \end{pmatrix}$$

$$\text{with } p_0(t) = \frac{1}{4} + \frac{3}{4} e^{-4\lambda t}$$

$$p_1(t) = \frac{1}{4} - \frac{1}{4} e^{-4\lambda t}$$

Imagine we have two sequences of length n with x differences. The probability that a position is different is $p = 3p_1(t)$. We define $d = 3\lambda t$ (the expected distance in time t).

Thus the probability that x positions out of n are different is:

$${n \choose x} p^x (1-p)^{n-x} = {n \choose x} \left(\frac{3}{4} - \frac{3}{4} e^{-\frac{4}{3}d} \right)^x \left(\frac{1}{4} + \frac{3}{4} e^{-\frac{4}{3}d} \right)^{n-x} = L(d; x)$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69: MLE for sequence distance

$${n \choose x} p^x (1-p)^{n-x} = {n \choose x} \left(\frac{3}{4} - \frac{3}{4} e^{-\frac{4}{3}d} \right)^x \left(\frac{1}{4} + \frac{3}{4} e^{-\frac{4}{3}d} \right)^{n-x} = L(d; x)$$

From this equation we obtain the log-likelihood function:

$$l(d; x) = \log {n \choose x} + x \log \left(\frac{3}{4} - \frac{3}{4} e^{-\frac{4}{3}d} \right) + (n - x) \log \left(\frac{1}{4} + \frac{3}{4} e^{-\frac{4}{3}d} \right)$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69: MLE for sequence distance

$${n \choose x} p^x (1-p)^{n-x} = {n \choose x} \left(\frac{3}{4} - \frac{3}{4} e^{-\frac{4}{3}d} \right)^x \left(\frac{1}{4} + \frac{3}{4} e^{-\frac{4}{3}d} \right)^{n-x} = L(d; x)$$

From this equation we obtain the log-likelihood function:

$$l(d; x) = \log {n \choose x} + x \log \left(\frac{3}{4} - \frac{3}{4} e^{-\frac{4}{3}d} \right) + (n-x) \log \left(\frac{1}{4} + \frac{3}{4} e^{-\frac{4}{3}d} \right)$$

Setting the first derivative to 0 and solving this equation with respect to d leads to the MLE of the JC69 distance:

$$\hat{d} = -\frac{3}{4} \log \left(1 - \frac{4x}{3n} \right)$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

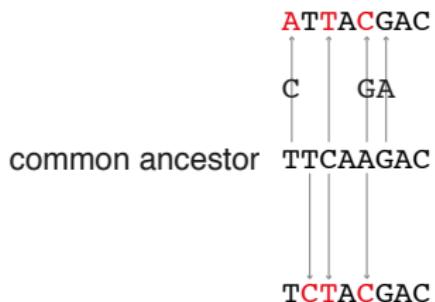
The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69: Example



- ▶ length of gene: $n = 8$
 - ▶ differences between the two sequences: $x = 2$
- $$\Rightarrow \hat{d} = -\frac{3}{4} \log \left(1 - \frac{4 \times 2}{3 \times 8} \right) = 0.3$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

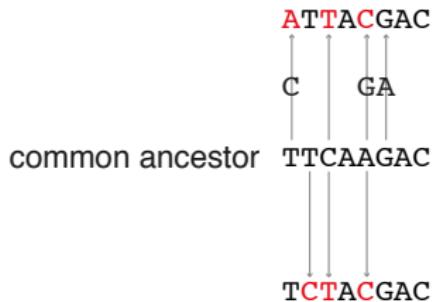
The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

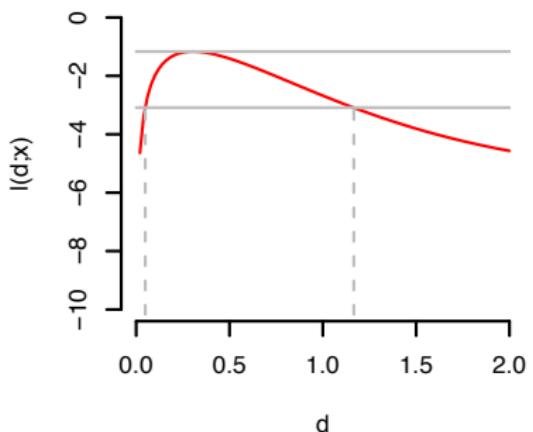
JC69: Example



- ▶ length of gene: $n = 8$
 - ▶ differences between the two sequences: $x = 2$
- $$\Rightarrow \hat{d} = -\frac{3}{4} \log \left(1 - \frac{4 \times 2}{3 \times 8} \right) = 0.3$$

Determination of the 95% confidence interval:

$$x = 2, n = 8: [0.05, 1.17]$$



04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

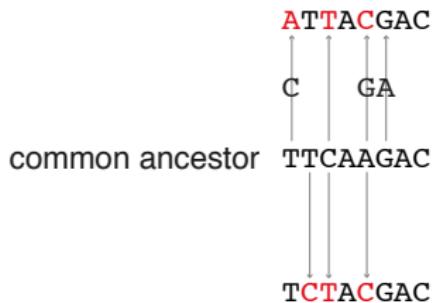
The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

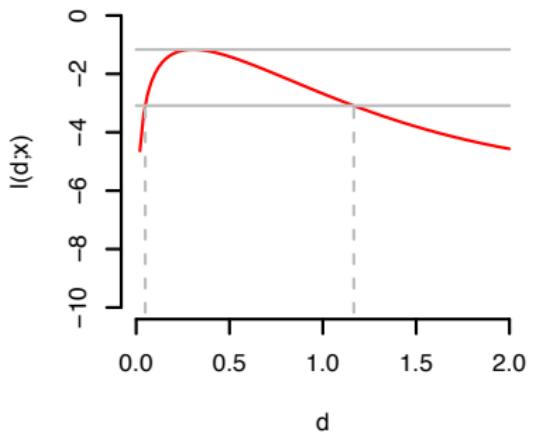
JC69: Example



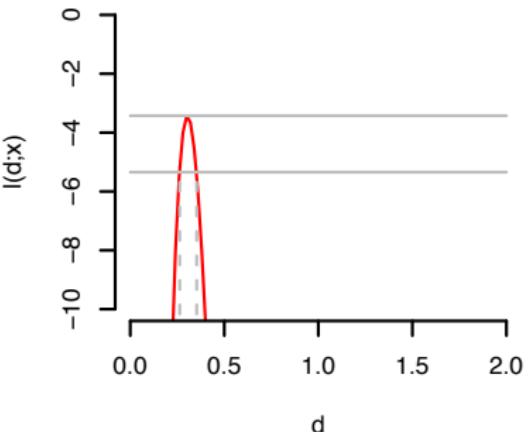
- ▶ length of gene: $n = 8$
 - ▶ differences between the two sequences: $x = 2$
- $$\Rightarrow \hat{d} = -\frac{3}{4} \log \left(1 - \frac{4 \times 2}{3 \times 8} \right) = 0.3$$

Determination of the 95% confidence interval:

$$x = 2, n = 8: [0.05, 1.17]$$



$$x = 200, n = 800: [0.26, 0.35]$$



04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Variable substitution rates across the genome.

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Variable rates

- ▶ so far: all sites in the sequence evolve at the same rate
- ▶ but: substitution rates might differ across the genome
 - ▶ mutation rates might differ across sites
 - ▶ selective pressure might be different on the phenotypic level

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Variable rates

- ▶ so far: all sites in the sequence evolve at the same rate
- ▶ but: substitution rates might differ across the genome
 - ▶ mutation rates might differ across sites
 - ▶ selective pressure might be different on the phenotypic level

We extend the existing models, by replacing the constant rates by Γ -distributed random variables (notation: JC69+ Γ , K80+ Γ , ...)

- ▶ sequence distance: expected number of substitutions per site, averaged over all sites

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

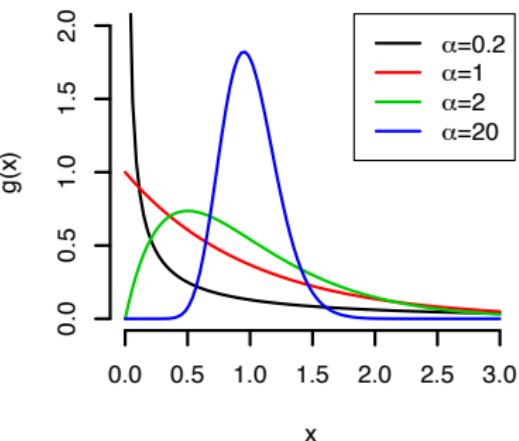
References

Toolbox: The Γ -distribution

- ▶ probability distribution: $g(x; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} e^{-\beta x} x^{\alpha-1}, x \geq 0$
- ▶ Γ -function: $\Gamma(\alpha) = \int_0^\infty e^{-t} t^{\alpha-1} dt; \Gamma(n) = (n-1)!$
- ▶ a $\Gamma(\alpha, \beta)$ -distributed random variable X has mean $E[X] = \frac{\alpha}{\beta}$ and variance $\text{var}[X] = \frac{1}{\alpha}$

Here, we fix the mean of the distribution to be 1, this is the case if and only if $\alpha = \beta$

- ▶ Γ -distribution very flexible



04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69+ Γ

- ▶ we replace the substitution rate λ by λR , where R is a $\Gamma(\alpha, \alpha)$ -distributed random variable (! mean 1)
- ▶ JC 69: $p = 3p_1(t) = \frac{3}{4} - \frac{3}{4}e^{-4\lambda t} = \frac{3}{4} - \frac{3}{4}e^{-\frac{4}{3}d}$
JC69+ Γ : $p = \frac{3}{4} - \frac{3}{4}e^{-4\lambda Rt} = \frac{3}{4} - \frac{3}{4}e^{-\frac{4}{3}dR}$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69+ Γ

- ▶ we replace the substitution rate λ by λR , where R is a $\Gamma(\alpha, \alpha)$ -distributed random variable (! mean 1)
- ▶ JC 69: $p = 3p_1(t) = \frac{3}{4} - \frac{3}{4}e^{-4\lambda t} = \frac{3}{4} - \frac{3}{4}e^{-\frac{4}{3}d}$
 JC69+ Γ : $p = \frac{3}{4} - \frac{3}{4}e^{-4\lambda Rt} = \frac{3}{4} - \frac{3}{4}e^{-\frac{4}{3}dR}$

To obtain the probability of a substitution at one site, we calculate $E[p]$:

$$E[p] = \int_0^\infty \left(\frac{3}{4} - \frac{3}{4}e^{-\frac{4}{3}dr} \right) g(r) dr = \frac{3}{4} - \frac{3}{4} \left(1 + \frac{4d}{\alpha} \right)^{-\alpha}$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69+ Γ

- ▶ we replace the substitution rate λ by λR , where R is a $\Gamma(\alpha, \alpha)$ -distributed random variable (! mean 1)
- ▶ JC 69: $p = 3p_1(t) = \frac{3}{4} - \frac{3}{4}e^{-4\lambda t} = \frac{3}{4} - \frac{3}{4}e^{-\frac{4}{3}d}$
 JC69+ Γ : $p = \frac{3}{4} - \frac{3}{4}e^{-4\lambda Rt} = \frac{3}{4} - \frac{3}{4}e^{-\frac{4}{3}dR}$

To obtain the probability of a substitution at one site, we calculate $\mathbb{E}[p]$:

$$\mathbb{E}[p] = \int_0^\infty (\frac{3}{4} - \frac{3}{4}e^{-\frac{4}{3}dr})g(r)dr = \frac{3}{4} - \frac{3}{4}\left(1 + \frac{4d}{\alpha}\right)^{-\alpha}$$

- ▶ to estimate the distance d , we use the MLE as described before with $l(d; x) = \binom{n}{x}(\mathbb{E}[p])^x(1 - \mathbb{E}[p])^{n-x}$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69+ Γ

- ▶ we replace the substitution rate λ by λR , where R is a $\Gamma(\alpha, \alpha)$ -distributed random variable (! mean 1)
- ▶ JC 69: $p = 3p_1(t) = \frac{3}{4} - \frac{3}{4}e^{-4\lambda t} = \frac{3}{4} - \frac{3}{4}e^{-\frac{4}{3}d}$
 JC69+ Γ : $p = \frac{3}{4} - \frac{3}{4}e^{-4\lambda Rt} = \frac{3}{4} - \frac{3}{4}e^{-\frac{4}{3}dR}$

To obtain the probability of a substitution at one site, we calculate $\mathbb{E}[p]$:

$$\mathbb{E}[p] = \int_0^\infty \left(\frac{3}{4} - \frac{3}{4}e^{-\frac{4}{3}dr} \right) g(r) dr = \frac{3}{4} - \frac{3}{4} \left(1 + \frac{4d}{\alpha} \right)^{-\alpha}$$

- ▶ to estimate the distance d , we use the MLE as described before with $l(d; x) = \binom{n}{x} (\mathbb{E}[p])^x (1 - \mathbb{E}[p])^{n-x}$

$$\hat{d} = \frac{3}{4}\alpha \left(\left(1 - \frac{4}{3}\hat{p} \right)^{-1/\alpha} - 1 \right)$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

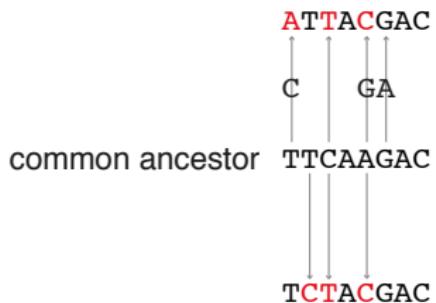
The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69+Γ: Example



- ▶ length of gene: $n = 8$
- ▶ differences between the two sequences: $x = 2$
- ⇒ $\hat{p} = 2/8 = 0.25$
- ▶ site variation:
 $\Gamma(2, 2)$ -distributed, i.e.
 $\alpha = 2$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

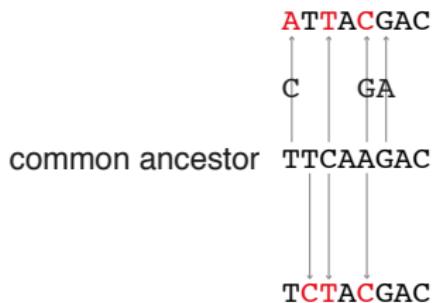
The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69+Γ: Example



- ▶ length of gene: $n = 8$
- ▶ differences between the two sequences: $x = 2$
- $\Rightarrow \hat{p} = 2/8 = 0.25$
- ▶ site variation:
 $\Gamma(2, 2)$ -distributed, i.e.
 $\alpha = 2$

$$\text{Thus: } \frac{3}{4}\alpha \left(\left(1 - \frac{4}{3}\hat{p}\right)^{-1/\alpha} - 1 \right) = 0.34$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

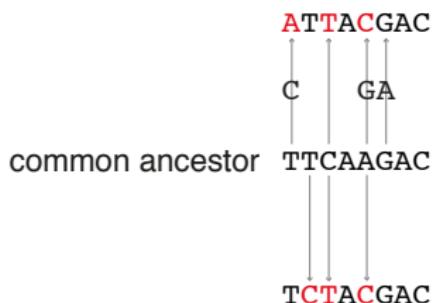
The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

JC69+Γ: Example



- ▶ length of gene: $n = 8$
- ▶ differences between the two sequences: $x = 2$
- $\Rightarrow \hat{p} = 2/8 = 0.25$
- ▶ site variation:
 $\Gamma(2, 2)$ -distributed, i.e.
 $\alpha = 2$

$$\text{Thus: } \frac{3}{4}\alpha \left(\left(1 - \frac{4}{3}\hat{p}\right)^{-1/\alpha} - 1 \right) = 0.34$$

This is bigger than $\hat{d} = 0.3$ in JC69 \Rightarrow

Ignoring site variation leads to underestimation of the sequence distance

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Distance usage: distance based phylogenetic reconstruction

We can now replace the Hamming-distance with the evolutionary distance for distance based phylogenetic reconstruction:



04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

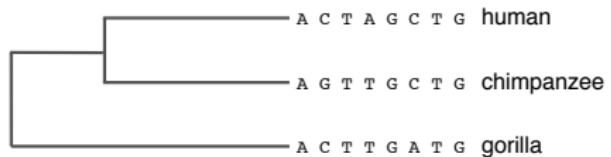
Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Distance usage: distance based phylogenetic reconstruction

We can now replace the Hamming-distance with the evolutionary distance for distance based phylogenetic reconstruction:



- here we only derived distance estimators for the most simple nucleotide substitution model (JC69), other substitution rate models would lead to other distances estimated
- not all phylogenetic reconstruction methods are based on distances (other methods will be discussed from lecture 5 on)

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Amino acid substitution models.

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices
Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance
Variable substitution rates across sites

Amino acid substitution models

Codon substitution models
The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

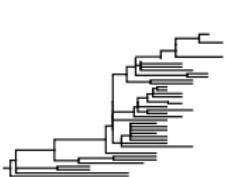
References

Levels of evolution

genotype

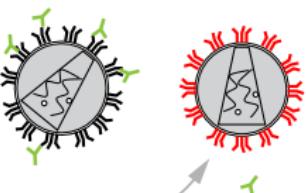
sequence level

ACUGAACGUGACACUCUG
ACUGAACGUAACUACUG



phenotype

e.g. antigenic level: Antibody binding to HIV



codon: three nucleotides encode for one amino acid

one nucleotide change can already change the phenotype

alphabet:

4 nucleotides: DNA: TCAG
RNA: UCAG

$64=4^3$ codons

20 amino acids

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Method for estimating distance between two strains

In general, the same Markov model is used for amino acid substitutions as for nucleotide substitutions. The transition probability matrix can be calculated according to the equation

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

Of course, in case of AA substitutions, $P(t)$ as well as the substitution rate matrix have dimension 20×20 . To determine $P(t)$, we need to derive the Q-matrix:

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Method for estimating distance between two strains

In general, the same Markov model is used for amino acid substitutions as for nucleotide substitutions. The transition probability matrix can be calculated according to the equation

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

Of course, in case of AA substitutions, $P(t)$ as well as the substitution rate matrix have dimension 20×20 . To determine $P(t)$, we need to derive the Q-matrix:

- ▶ more difficult than for nucleotide substitutions
- ▶ two approaches
 - ▶ empiric
 - ▶ mechanistic
- ▶ desired: time-reversibility

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Distance estimation of AA sequences

Let us consider the easiest case of AA substitutions: All substitutions happen at the same rate λ .

- ▶ as we have 20 AA, the substitution rate for any substitution is then 19λ
- ▶ the expected time until a substitution happens is then $\frac{1}{19\lambda}$
- ▶ with $t = \frac{d}{19\lambda}$ we can also derive the distance with the maximum likelihood approach

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Distance estimation of AA sequences

Let us consider the easiest case of AA substitutions: All substitutions happen at the same rate λ .

- ▶ as we have 20 AA, the substitution rate for any substitution is then 19λ
- ▶ the expected time until a substitution happens is then $\frac{1}{19\lambda}$
- ▶ with $t = \frac{d}{19\lambda}$ we can also derive the distance with the maximum likelihood approach

$$\hat{d} = -\frac{19}{20} \log \left(1 - \frac{20x}{19n} \right)$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Distance estimation of AA sequences

Let us consider the easiest case of AA substitutions: All substitutions happen at the same rate λ .

- ▶ as we have 20 AA, the substitution rate for any substitution is then 19λ
- ▶ the expected time until a substitution happens is then $\frac{1}{19\lambda}$
- ▶ with $t = \frac{d}{19\lambda}$ we can also derive the distance with the maximum likelihood approach

$$\hat{d} = -\frac{19}{20} \log \left(1 - \frac{20x}{19n} \right)$$

- ▶ to estimate the distance with an empirical or mechanistic Q-matrix, we also use a MLE

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Codon substitution models.

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices
Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance
Variable substitution rates across sites

Amino acid substitution models

Codon substitution models
The universal genetic code

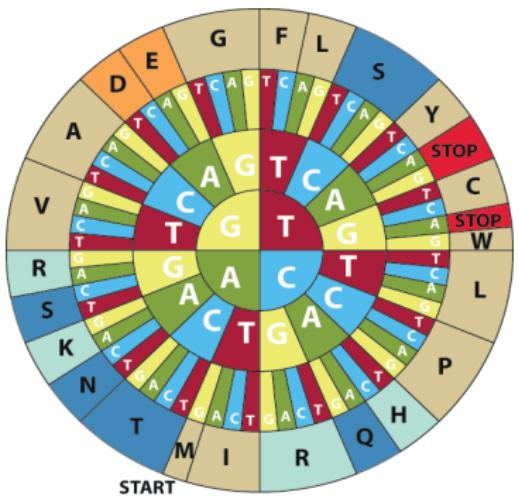
Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The codon sun

A codon consists of three nucleotides, translating to one of the 20 amino acids:



[Sanger, 2015]

Amino Acid	Three-Letter Abbreviation	One-Letter Symbol	Molecular Weight
Alanine	Ala	A	89Da
Arginine	Arg	R	174Da
Asparagine	Asn	N	132Da
Asparticacid	Asp	D	133Da
Asparaginer			
asparticacid	Asx	B	133Da
Cysteine	Cys	C	121Da
Glutamine	Gln	Q	146Da
Glutamiacid	Glu	E	147Da
Glutaminer			
glutamiacid	Glx	Z	147Da
Glycine	Gly	G	75Da
Histidine	His	H	155Da
Isoleucine	Ile	I	131Da
Leucine	Leu	L	131Da
Lysine	Lys	K	146Da
Methionine	Met	M	149Da
Phenylalanine	Phe	F	165Da
Proline	Pro	P	115Da
Serine	Ser	S	105Da
Threonine	Thr	T	119Da
Tryptophan	Trp	W	204Da
Tyrosine	Tyr	Y	181Da
Valine	Val	V	117Da

[Promega, 2015]

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

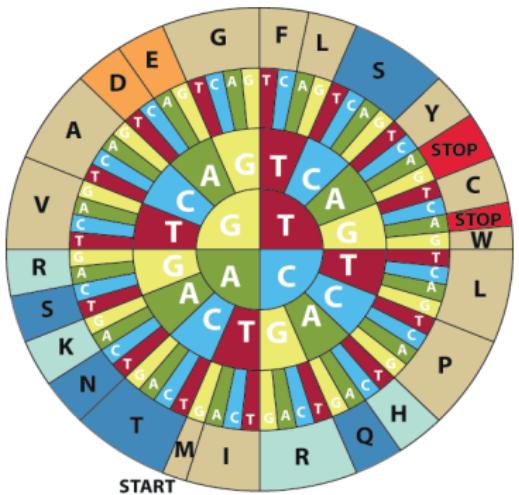
Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Example: Codon CTA

Which one-point mutations lead to the codon CTA for leucine?



[Sanger, 2015]

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

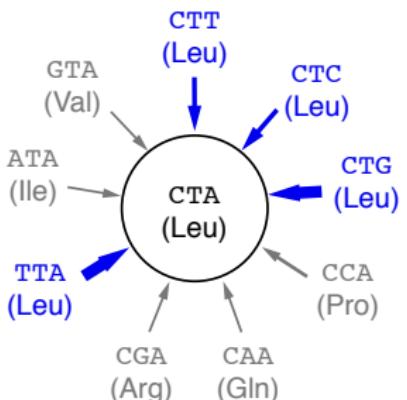
substitution
Synonymous and

non-synonymous substitutions

References

Example: Codon CTA

Overview over substitution rates to the same codon CTA, the thickness of arrows represent different rates:



- ▶ **synonymous substitutions:**
AA does not change
- ▶ **nonsynonymous substitutions:**
AA does change
- ▶ bigger arrows: transition
- ▶ smaller arrows: transversion

adapted from [Yang, 2014]

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The basic model

We also model codon substitution with a Markov chain model, where 61 states (no stop codons) are allowed:

- The substitution rate matrix has dimension 61×61

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The basic model

We also model codon substitution with a Markov chain model, where 61 states (no stop codons) are allowed:

- ▶ The substitution rate matrix has dimension 61×61
- ▶ We denote codons with capital letters, e.g. I, J, but nucleotides with small letters, e.g. i, j

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The basic model

We also model codon substitution with a Markov chain model, where 61 states (no stop codons) are allowed:

- ▶ The substitution rate matrix has dimension 61×61
- ▶ We denote codons with capital letters, e.g. I, J, but nucleotides with small letters, e.g. i, j
- ▶ We incorporate:
 - ▶ κ : transition/transversion rate ratio
 - ▶ ω : nonsynonymous/synonymous rate ratio
 - ▶ π_I : equilibrium frequency of codon I consisting of nucleotides $i_1 i_2 i_3$, with equilibrium frequencies $\pi_{i_1}, \pi_{i_2}, \pi_{i_3}$

$$\pi_I = \frac{1}{C} \pi_{i_1}^* \pi_{i_2}^* \pi_{i_3}^*$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The basic model

We also model codon substitution with a Markov chain model, where 61 states (no stop codons) are allowed:

- ▶ The substitution rate matrix has dimension 61×61
 - ▶ We denote codons with capital letters, e.g. I, J, but nucleotides with small letters, e.g. i, j
 - ▶ We incorporate:
 - ▶ κ : transition/transversion rate ratio
 - ▶ ω : nonsynonymous/synonymous rate ratio
 - ▶ π_I : equilibrium frequency of codon I consisting of nucleotides $i_1 i_2 i_3$, with equilibrium frequencies $\pi_{i_1}, \pi_{i_2}, \pi_{i_3}$
- $$\pi_I = \frac{1}{C} \pi_{i_1}^* \pi_{i_2}^* \pi_{i_3}^*$$

Model:

$$q_{IJ} = \begin{cases} 0 & \text{if I and J differ at more than 1 positions} \\ \pi_J & \text{if I and J differ by a synonymous transversion} \\ \kappa \pi_J & \text{if I and J differ by a synonymous transition} \\ \omega \pi_J & \text{if I and J differ by a nonsynonymous transversion} \\ \omega \kappa \pi_J & \text{if I and J differ by a nonsynonymous transition} \end{cases}$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The substitution rate matrix

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The substitution rate matrix

$$Q = \gamma \begin{pmatrix} & (P) & (*) & (H) & (Q) & (Q) \\ \dots & CCG & CAT & CAC & CAA & CAG \\ \vdots & -\sum_{row} & 0 & 0 & 0 & 0 \\ CCG(P) & 0 & -\sum_{row} & 0 & 0 & 0 & 0 \\ CAT(H) & 0 & 0 & -\sum_{row} w_{CAC} w_{CAA} w_{CAG} & 0 \\ CAC(H) & 0 & 0 & w_{CAC} -\sum_{row} w_{CAA} w_{CAG} & 0 \\ CAA(Q) & 0 & 0 & w_{CAA} w_{CAG} -\sum_{row} w_{CAG} & 0 \\ CAG(Q) & 0 & 0 & w_{CAG} w_{CAG} -\sum_{row} w_{CAG} & 0 \\ \vdots & & & & & & -\sum_{row} \end{pmatrix}$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Evidence for selection

- ▶ synonymous substitutions do not change the protein, thus, these substitutions are seen as neutral
- ▶ nonsynonymous substitutions do change the protein and selective processes can act on the new protein

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Evidence for selection

- ▶ synonymous substitutions do not change the protein, thus, these substitutions are seen as neutral
- ▶ nonsynonymous substitutions do change the protein and selective processes can act on the new protein
- ▶ to discover selection, one compares amounts of nonsynonymous and synonymous substitutions

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Evidence for selection

- ▶ synonymous substitutions do not change the protein, thus, these substitutions are seen as neutral
- ▶ nonsynonymous substitutions do change the protein and selective processes can act on the new protein
- ▶ to discover selection, one compares amounts of nonsynonymous and synonymous substitutions

d_N : distance at nonsynonymous codon positions

d_S : distance at synonymous codon positions

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

d_N/d_S ratio: Counting method

In the following we will look at a method based on counting synonymous and non-synonymous substitutions. We will illustrate this method with the two sequences TTTCCTCCTCCT and TTCCAGCCTCCT:

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

d_N/d_S ratio: Counting method

In the following we will look at a method based on counting synonymous and non-synonymous substitutions. We will illustrate this method with the two sequences TTTCCCTCCTCCT and TTCCAGCCTCCT:

	codon 1	codon 2	codon 3	codon 4
sequence 1	TTT	CCT	CCT	CCT
sequence 2	TTC	CAG	CCT	CCT

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

d_N/d_S ratio: Counting method

In the following we will look at a method based on counting synonymous and non-synonymous substitutions. We will illustrate this method with the two sequences TTTCCTCCTCCT and TTCCAGCCTCCT:

	codon 1	codon 2	codon 3	codon 4
sequence 1	TTT	CCT	CCT	CCT
	F	P	P	P
	F	Q	P	P
sequence 2	TTC	CAG	CCT	CCT

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

d_N/d_S ratio: Counting method

In the following we will look at a method based on counting synonymous and non-synonymous substitutions. We will illustrate this method with the two sequences TTTCCCTCCTCCT and TTCCAGCCTCCT:

	codon 1	codon 2	codon 3	codon 4
sequence 1	TTT	CCT	CCT	CCT
	F	P	P	P
	F	Q	P	P
sequence 2	TTC	CAG	CCT	CCT

Algorithm:

1. count the (non-) synonymous differences
2. count the (non-) synonymous sites
3. account for the possible evolutionary history

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Step 1: number of (non-) synonymous differences

N_d : number of nonsynonymous differences

S_d : number of synonymous differences

... taking into account all possible ways from seq1 to seq2

	codon 1		codon 2		codon 3		codon 4		
seq1	TTT		CCT		CCT		CCT		
seq2	TTC		CAG		CCT		CCT		
$N_d =$	0	+		+	0	+	0	=	
$S_d =$	1	+		+	0	+	0	=	

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Step 1: number of (non-) synonymous differences

N_d : number of nonsynonymous differences

S_d : number of synonymous differences

... taking into account all possible ways from seq1 to seq2

	codon 1	codon 2	codon 3	codon 4	
seq1	TTT	CCT	CCT	CCT	
seq2	TTC	CAG	CCT	CCT	
$N_d =$	0	+	+	0	=
$S_d =$	1	+	+	0	=

pathway	S_d at codon 2	N_d at codon 2
CCT (P) → CAT (H) → CAG (Q)		
CCT (P) → CCG (P) → CAG (Q)		
average		

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Step 1: number of (non-) synonymous differences

N_d : number of nonsynonymous differences

S_d : number of synonymous differences

... taking into account all possible ways from seq1 to seq2

	codon 1	codon 2	codon 3	codon 4	
seq1	TTT	CCT	CCT	CCT	
seq2	TTC	CAG	CCT	CCT	
$N_d =$	0	+	+	0	=
$S_d =$	1	+	+	0	=

pathway	S_d at codon 2	N_d at codon 2
CCT (P) → CAT (H) → CAG (Q)	0	2
CCT (P) → CCG (P) → CAG (Q)	1	1
average	0.5	1.5

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Step 1: number of (non-) synonymous differences

N_d : number of nonsynonymous differences

S_d : number of synonymous differences

... taking into account all possible ways from seq1 to seq2

	codon 1	codon 2	codon 3	codon 4	
seq1	TTT	CCT	CCT	CCT	
seq2	TTC	CAG	CCT	CCT	
$N_d =$	0	+	1.5	+	0 + 0 = 1.5
$S_d =$	1	+	0.5	+	0 + 0 = 1.5

pathway	S_d at codon 2	N_d at codon 2
CCT (P) → CAT (H) → CAG (Q)	0	2
CCT (P) → CCG (P) → CAG (Q)	1	1
average	0.5	1.5

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Step 2: number of (non-) synonymous sites

N : number of nonsynonymous sites

S : number of synonymous sites

... averaging over all potential one point mutations

important: the sum of (non-) synonymous mutations per codon must sum up to 3

	cdn1	cdn2	cdn3	cdn4	
seq1	TTT	CCT	CCT	CCT	
seq2	TTC	CAG	CCT	CCT	
nonsyn					
seq1	+	+	+		=
seq2	+	+	+		=
average					$N =$
syn					
seq1	+	+	+		=
seq2	+	+	+		=
average					$S =$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Step 2: number of (non-) synonymous sites

	cdn1	cdn2	cdn3	cdn4
seq1	TTT	CCT	CCT	CCT
seq2	TTC	CAG	CCT	CCT

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

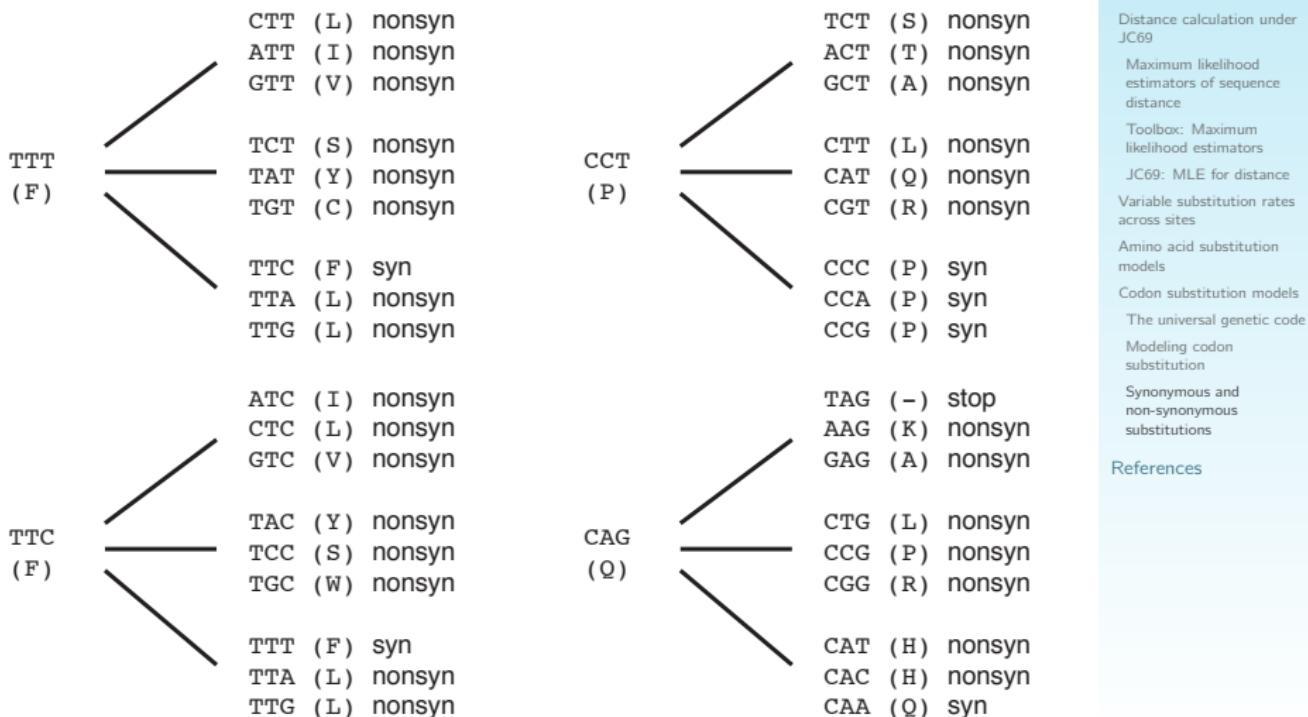
Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Step 2: number of (non-) synonymous sites

	cdn1	cdn2	cdn3	cdn4
seq1	TTT	CCT	CCT	CCT
seq2	TTC	CAG	CCT	CCT



04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Step 2: number of (non-) synonymous sites

	cdn1	cdn2	cdn3	cdn4			
seq1	TTT	CCT	CCT	CCT			
seq2	TTC	CAG	CCT	CCT			
nonsyn							
seq1	8/3	+	2	+	2	+	2 = 8.67
seq2	8/3	+	21/8	+	2	+	2 = 9.29
average							N = 8.98
syn							
seq1	1/3	+	1	+	1	+	1 = 3.33
seq2	1/3	+	3/8	+	1	+	1 = 2.71
average							S = 3.02

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Step 3: accounting for evolution

Under the JC69 substitution model, we can now calculate the distances at (non-)synonymous codon positions, d_N and d_S :

$$d_N = -\frac{3}{4} \log \left(1 - \frac{4 N_d}{3 N} \right)$$

$$d_S = -\frac{3}{4} \log \left(1 - \frac{4 S_d}{3 S} \right)$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Step 3: accounting for evolution

Under the JC69 substitution model, we can now calculate the distances at (non-)synonymous codon positions, d_N and d_S :

$$d_N = -\frac{3}{4} \log \left(1 - \frac{4}{3} \frac{N_d}{N} \right)$$

$$d_S = -\frac{3}{4} \log \left(1 - \frac{4}{3} \frac{S_d}{S} \right)$$

In our example, the d_N/d_S -ratio is:

$$\frac{d_N}{d_S} = \frac{-\frac{3}{4} \log \left(1 - \frac{4}{3} \frac{N_d}{N} \right)}{-\frac{3}{4} \log \left(1 - \frac{4}{3} \frac{S_d}{S} \right)} = \frac{\log \left(1 - \frac{4}{3} \frac{1.5}{8.98} \right)}{\log \left(1 - \frac{4}{3} \frac{1.5}{3.02} \right)} = 0.23$$

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The interpretation of the d_N/d_S -ratio

$d_N/d_S < 1$: nonsynonymous mutations occur less frequently than synonymous mutations (purifying selection)

$d_N/d_S > 1$: nonsynonymous mutations occur more frequently than synonymous mutations (positive selection)

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

The interpretation of the d_N/d_S -ratio

$d_N/d_S < 1$: nonsynonymous mutations occur less frequently than synonymous mutations (purifying selection)

$d_N/d_S > 1$: nonsynonymous mutations occur more frequently than synonymous mutations (positive selection)

☞ In addition, ML methods have been developed. However:

"It is worth noting that the estimation of d_S and d_N is a complicated exercise, as manifested by the fact that even the simple transition/transversion rate difference is nontrivial to deal with. Unfortunately, different methods can often produce very different estimates." [Yang, 2014]

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

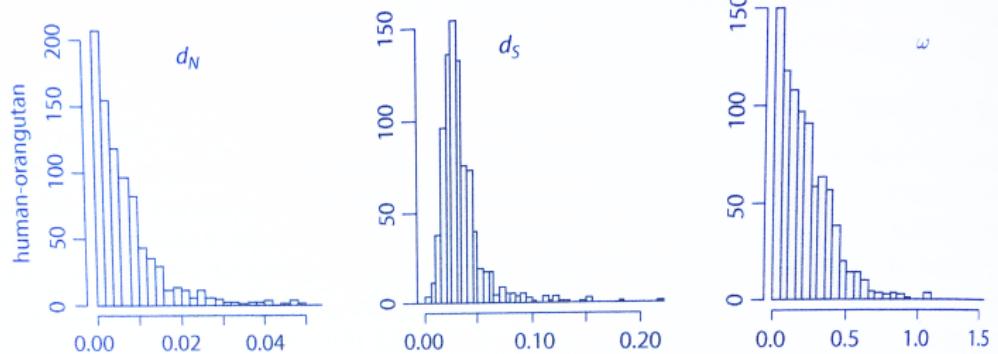
Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Example: d_N , d_S between humans and orangutans

Histograms of d_N , d_S and $\omega = \frac{d_N}{d_S}$ for 857 orthologous genes between humans and orangutans:



adapted from Figure 2.6 [Yang, 2014]

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

What we learnt today

- ▶ nucleotide substitution models
- ▶ maximum likelihood estimators for distances
- ▶ modeling site variation with Γ -distribution
- ▶ amino acid substitution models
- ▶ codon substitution models
 - ▶ d_N and d_S as a tool to study selection

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance
Variable substitution rates across sites

Amino acid substitution models

Codon substitution models
The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

Molecular evolution, AA and codon substitution models: Questions

- ② What are the differences between the JC69 and the TN93 models?
- ② Which of the two models would you chose if you were to perform a phylogenetic analysis based on sequence distances?
- ② In lecture 3 we tried to naively reconstruct a phylogeny based on three sequences (the same tree and sequences appear on slide 33 of this lecture). The pairwise Hamming distances were all the same. Would you expect that any of the presented nucleotide sequence models would result in different trees?

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References

References |

- Hasegawa, M., Yano, T., and Kishino, H. (1984). A New Molecular Clock of Mitochondrial-Dna and the Evolution of Hominoids. *Proceedings of the Japan Academy Series B-Physical and Biological Sciences*, 60(4):95–98.
- Jukes, T. and Cantor, C. (1969). Evolution of protein molecules. *Mammalian Protein Metabolism.*, pages 21–123.
- Kimura, M. (1980). A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide sequences. *Journal of molecular evolution*, 16(2):111–120.
- loaded.dice (2015). How to load dice: <http://rss2.com/feeds/How-to-of-the-Day/page51>.
- Promega (2015). The amino acids: <https://www.promega.com/ /media/files/resources/technical references/amino acid abbreviations and molecular weights.pdf>.
- Sanger (2015). The codon sun:
<ftp://ftp.sanger.ac.uk/pub/yourgenome/downloads/activities/kras-cancer-mutation/krascodonwheel.pdf>.
- Sokal, R. and Rohlf, F. (2012). *Biometry. Fourth Edition.* W.H. Freeman and Company.
- Tamura, K. and Nei, M. (1993). Estimation of the number of nucleotide substitutions in the control region of mitochondrial DNA in humans and chimpanzees. *Molecular Biology and Evolution*, 10(3):512–526.
- Tavaré, S. (1986). Some probabilistic and statistical problems in the analysis of DNA sequences. In *Some mathematical questions in biology—DNA sequence analysis (New York, 1984)*, pages 57–86. Amer. Math. Soc., Providence, RI.
- Yang, Z. (1994). Estimating the pattern of nucleotide substitution. *Journal of molecular evolution*, 39(1):105–111.
- Yang, Z. (2014). *Molecular Evolution – A Statistical Approach.* Oxford University Press.
- Zharkikh, A. (1994). Estimation of evolutionary distances between nucleotide sequences. *Journal of molecular evolution*, 39(3):315–329.

04: Nucleotide, amino acid and codon substitution models

Important concepts of last lecture

Substitution rate matrices

Distance calculation under JC69

Maximum likelihood estimators of sequence distance

Toolbox: Maximum likelihood estimators

JC69: MLE for distance

Variable substitution rates across sites

Amino acid substitution models

Codon substitution models

The universal genetic code

Modeling codon substitution

Synonymous and non-synonymous substitutions

References