Korea-China-Japan Bioinformatics Training Course

June 19, A.S. 0014 Jeju Island, Korea

Introduction to Evolutionary Genomics

SAITOU Naruya

National Institute of Genetics, Mishima Department of Genetics, Graduate University for Advanced Studies (adjunct) Department of Biolo



Textbook of my lecture

"Introduction to Evolutionary Genomics"

by Naruya Saitou

Springer-Verlag 2014 Naruya Saitou Introduction to Evolutionary Genomics

Computational Biology



Preface

Necessity of Evolutionary Studies What Is Evolution? What Is a Genome? Vitalism Versus Mechanism Everything Is History Genome as a Republic of Genes Structure of This Book Acknowledgments

The History of the Earth is recorded in the Layers of its Crust: The History of all Organisms is inscribed in the Chromesomes

March 188

The history of the earth is recorded in the layers of its crus The history of all organisms is inscribed in the chromosomes

Kihara Hitoshi (1946)

Introduction to Evolutionary Genomics

Part 1 Basic Processes of Genome Evolution

Part 2 Evolving Genomes

Part 3 Methods for Evolutionary Genomics

Part 1 Basic Processes of Genome Evolution

Chapter 1 Basic Metabolism Surrounding DNAs

- Chapter 2 Mutation
- Chapter 3 Phylogeny
- Chapter 4 Neutral Evolution
- Chapter 5 Natural Selection

Chapter 3 Phylogeny

3.1 DNA replications generate phylogenies 3.2 Genealogy of individuals 3.3 Gene genealogy 3.4 Species phylogeny 3.5 Basic concepts of trees and networks 3.6 Biological nature of trees and networks

Figure 3-2: A schematic representation of the phylogenetic tree of 10bp DNA sequences



Figure 3-4: Cell genealogy of *C. elegans* (based on ref 3-27)



Fig. 3.5 Genealogy of one diploid individual



Individual in question

Figure 3-7: Gene genealogy for haploids.(A) Animal mitochondrial DNA. (B) Mammalian Y chromosomes(C) Avean W chromosomes.





Fig. 3.8 Contributions of maternal and paternal lineages of mitochondrial DNA and Y chromosomes. M.G.M., M.G.F., P.G.M., and P.G.F. designate maternal grand mother, maternal grand father, paternal grand mother, and paternal grand father, respectively



Fig. 3.10 Alternative gene genealogies for four autosomal genes of two diploid individuals

Figure 3–11: Gene phylogeny with recombinations for gibbon ABO blood group genes (from ref 3–9)



From Kitano et al. (2010)

Fig. 3 (Kitano et al.)

Figure 3-12: Three kinds of gene phylogenies. (A) Temporal and mutational gene phylogeny. (B) Mutational gene phylogeny. (C) Estimated gene phylogeny (from ref 3-13).







Fig. 3.20 Two possible gene genealogies for three species



Figure 3-21: Trees for 6 OTUs. (A) A rooted tree. (B) An



Figure 3-22: Example of nontree networks for 6 OTUs





Fig. 3.23 Three possible unrooted trees for 4 OTUs and their corresponding rooted trees





For n OTUs, [Nr(n)]: possible number of rooted trees [Nu(n)]: possible number of unrooted trees

$$Nr(n) = 1 \times 3 \times 5 \times \dots \times (2n-3)$$

Nu(n) = $1 \times 3 \times 5 \times ... \times (2n-5)$ = $(2n - 5)! / [2^{n-3} (n - 3)!]$

Table 3.1 Number ofpossible unrooted treetopologies for up to 20 OTUs

Numbe	r of OTUs	Possible number of unrooted trees
3		1
4		3
5		15
6		105
7		945
8		10, 395
9		135, 135
10		2, 027, 025
11		34, 459, 425
12		654, 729, 705
13		13, 749, 310, 575
14		316, 234, 143, 225
15		7, 905, 853, 580, 625
16		213, 458, 046, 676, 875
17	6	6, 190, 283, 353, 629, 375
18	191	, 898, 783, 962, 510, 625
19	6, 332	2, 659, 870, 762, 850, 625
20	221, 643	, 095, 476, 699, 771, 875

Fig. 3.25 An example of unrooted tree for nine OTUs



$$\text{Tree}_\text{description} = \left[(1,2), (10,3), (4,5), (11,12), (8,9), (13,6) \right] \quad (3.8)$$

$$\begin{aligned} & \text{Free _description} = \left[(1,10), (2,10), (3,11), (4,12), (5,12), \\ & (6,14), (7,15), (8,16), (9,16), (10,11), \\ & (11,13), (12,13), (13,14), (14,15), (15,16) \right] \end{aligned} (3.6)$$

Tree description = [(1,2), (1,3), (4,5), (3,4), (8,9), (3,6)]

$$\frac{2}{1} + \frac{2}{10} + \frac{6}{14} + \frac{9}{15} + \frac{8}{7} + \frac{1}{12} +$$

(3.7)

How to describe one tree?



Tree description =
$$[(1,2), (10,3), (4,5), (11,12), (6,14), (7,15), (8,9,15)]$$
 (3.9)

Tree_description =
$$\left[\left(\left(((1,2),3), (4,5) \right), 6 \right), 7 \right), 8, 9 \right]$$
 (3.10)

Table 3.2 A splits matrix that describes the tree shown in Fig





	1	2	3	4	5	6	7	8	9
A	+	+	_	_	<u> </u>	_	_	_	_
в	+	+	+	_	_	_	_	_	_
С	+	+	+	_	_	+	+	+	+
D	+	+	+	+	_	_	_	_	_
E	+	+	+	+	+	+	_	_	_
F	+	+	+	+	+	+	+	_	_
G	+	_	+	_		_	-	_	

Table 3.3 A splits matrix for nine OTUs that are not mutually compatible or nested

Fig. 3.26 A nontree network corresponding to the split matrix shown in Table 3.3



Fig. 3.27 A tree with one trifurcation



Fig. 3.30 Two possible unlabeled rooted trees for four spec



Part 2 Evolving Genomes

Chapter 6 Brief History of Life Chapter 7 Prokaryote Genomes Chapter 8 Eukaryote Genomes Chapter 9 Vertebrate Genomes Chapter 10 Human Genomes

Part 3 Methods for Evolutionary Genomics Chapter 11 Genome Sequencing Chapter 12 Omic Data Collection Chapter 13 Databases Chapter 14 Sequence Homology Handling Chapter 15 Evolutionary Distances Chapter 16 Tree and Network Building Chapter 17 Population Genomics

Chapter 15 Evolutionary Distances

15.1 Overview of Evolutionary Distances
15.2 Nucleotide Substitutions
15.3 Synonymous and Nonsynonymous
Substitutions
15.4 Amino Acid Substitutions
15.5 Evolutionary Distances Not Based on
Substitutions

Table 15-1: Example of a distance matrix (data from Table 3 of Ishida et al. [1995; ref 15-54])

	1	2	3	4	5	6	====== 7
1	0	9	11	6	42	38	35
2	9	0	6	5	45	41	38
3	11	6	0	7	47	43	40
4	6	5	7	0	42	38	35
5	42	45	47	5	0	46	43
6	38	41	43	5	46	0	29
7	35	38	40	5	43	29	0

1: Thoroughbred horse (*Equus caballus*), 2: Przewalskii's wild horse (*E. caballus*), 3: Mongolian native horse (*E. caballus*), 4: Japanese native horse (*E. caballus*), 5: mountain zebra (*E. zebra*), 6: donkey (*E. asinus*), and 7: Grevy's zebra (*E. grevyi*).

Table 15-2: Example of a distance matrix showing only lower-triangle values (data from Table 3 of Ishida et al. [1995; ref 15-54])

	1	2	3	4	5	6
2 3	8.0 11.7	5.4				
4	6.5	5.6	5.4			
5	41.8	49.0	46.5	43.7		
6	35.1	41.3	45.8	35.5	47.9	
7	36.6	34.8	34.8	38.2	41.5	29.1

1-7: Same as those of Table 15-1

Figure 15-1: An example of nucleotide sequence evolution only through substitutions


Figure 15-2: Another example of nucleotide sequence evolution only through substitutions



Table 15-3: One-parameter and two-parameter models of nucleotide substitution matrix

(A) One-parameter model

===	=====		======= N I	======================================		=
		А	C	z w T	G	
O L D	A C T G	1–3α α α α	α 1–3α α α	α α 1–3α α	α α α 1–3α	-

(B) Two-parameter model



Figure 15-3: An evolutionary scheme between two presentday sequences





Figure 15-4: Relationship between *p* and *d* under the one-parameter method

Ρ

Figure 15-5: Temporal changes of *P* and *Q* under the two-parameter model with $\alpha = 10\beta$



αt

Table 15-5: Pattern of nucleotide substitutions for the human mitochondrial genome (from Kawai, Kikuchi, and Saitou, unpublished)

===	=====		======= N E	 W	
		Α	С	Т	G
0	Α	-	0.0031	0.0030	0.0901
L	C	0.0070	-	0.0593	0.0028
D	Т	0.0042	0.2566	-	0.0048
	G	0.5574	0.0080	0.0037	-

Table 15-6: Models of nucleotide substitution matrices incorporating nucleotide frequencies

		A	N E W C	Т	G
O L D	A C T G	1-Σλ _Α . π _A α π _A α π _A α	π _C α 1-Σλ _C . π _C α π _C α	$\pi_{T} lpha \ \pi_{T} lpha \ 1-\Sigma \lambda_{T} llowbreak \ \pi_{T} lpha$	π _G α π _G α π _G α 1-Σλ _{G•}

(A) Equal-input model

(-)-						
===:			======		=====	
			NEW			
		Α	С	Т	G	
0	Α	$1 - \Sigma \lambda_{A^{\bullet}}$	$\pi_A \alpha$	$\pi_A \alpha$	$\pi_A \alpha$	
L	С	$\pi_{C}\alpha$	$1-\Sigma\lambda_{C}$	$\pi_{C}\alpha$	$\pi_{C}\alpha$	
D	Т	$\pi_T \alpha$	$\pi_T \alpha$	$1 - \Sigma \lambda_{T}$.	$\pi_T \alpha$	
	G	$\pi_G \alpha$	$\pi_G \alpha$	$\pi_G \alpha$	$1-\Sigma\lambda_{G}$	

(C) Hasegawa-Kishino-Yano model

			NEW		
		Α	С	Т	G
0	Α	$1-\Sigma\lambda_{A\bullet}$	$\pi_{C}\beta$	$\pi_T\beta$	$\pi_G \alpha$
L	С	$\pi_A \beta$	$1-\Sigma\lambda_{C}$.	$\pi_T \alpha$	$\pi_G \beta$
D	Т	$\pi_A \beta$	$\pi_{C}\alpha$	$1-\Sigma\lambda_1$	 π_Gβ
	G	$\pi_A \alpha$	$\pi_C \beta$	$\pi_T\beta$	$1 - \Sigma \lambda_{G^{\bullet}}$

(D) Tamura-Nei model

(B) Equal-output model



Table 15-7: Various models of nucleotide substitutions



(B) Takahata-Kimura 4P (ref 15-18)

Figure 15-6: Relationship of various nucleotide substitution models



Ala Arg Asn Asp Cys Gln Glu Gly His Ile Leu Lys Met Phe Pro Ser Thr Trp Tvr Val *** Ala Ara *** *** Asn Asp *** *** Cys Gln *** Glu *** Gly *** His *** Ile *** *** Leu *** Lys *** Met *** Phe *** Pro Ser *** *** Thr *** Trp *** Tyr Val ***

Table 15-9: Amino acid substitution matrix for vertebrate mitochondrial DNA coded proteins

Chapter 16 Tree and Network Building

- 16.1 Classification of Tree-Building Methods
- 16.2 Distance Matrix Methods
- 16.3 Neighbor-Joining Method
- 16.4 Phylogenetic Network Construction from Distance Matrix Data
- 16.5 Maximum Parsimony Method
- 16.6 The Maximum Likelihood Method
- 16.7 Phylogenetic Network Construction from Character- State Data
- 16.8 Tree-Searching Algorithms
- 16.9 Comparison of Phylogenetic Tree-Making Methods 16.10 Phylogeny Construction Without Multiple Alignment

Table 16-1: Examples of distance matices

(A) When the evolutionary rate is constant

A B C D E F

A	0	2	6	4	16	16
B	2	0	6	4	16	16
С	6	6	0	6	16	16
D	4	4	6	0	16	16
E	16	16	16	16	0	10
F	16	16	16	16	10	0

(B) When all distances are purely additive (from ref 16-1)

	1	2	3	4	5	6	7	8
1	0 7	7	8	11 8	13 10	16 13	13 10	17 14
3	8	5	0	5	7	10	7	11
4 5	11 13	8 10	5 7	0 8	8 0	6	8 6	10
6 7	16 13	13 10	10 7	11 8	5 6	0 9	9 0	13 8
8	17	14	11	12	10	13	8	0

(C) Dis	stance mat	rix estima	ted from r	eal nulced	otide seque	ence data (from ref 1	(6-2)	
	1	2	3	4	5	6	7	8	9
2 3 4 5 6 7 8 9	0.0516 0.0550 0.0483 0.0582 0.0094 0.0125 0.0284 0.0925	0.0031 0.0221 0.0651 0.0416 0.0584 0.0687 0.1221	0.0253 0.0685 0.0450 0.0619 0.0722 0.1259	0.0549 0.0384 0.0551 0.0654 0.1185	0.0549 0.0651 0.0754 0.1370	0.0157 0.0317 0.0820	0.0285 0.0786	0.0927	0 1900
10	0.1921	0.2183	0.2228	0.2054	0.2309	0.1798	0.1795	0.1833	0.1860

OTU ID; 1 = M. *m. domesticus* functional gene, 2 = M. *m. domesticus* pseudogene, 3 = M. *m. castaneus* pseudogene, 4 = M. spicilegus pseudogene, 5 = M. *leggada* pseudogene, 6 = M. *m. domesticus* functional gene, 7 = M. *leggada* functional gene, 8 = M. *platythrix* functional gene, 9 = Rattus norvegicus functional gene, 10 = Homo sapiens functional gene.

Table 16-2: Operation of UPGMA for distance matrix of Table 16-1A

(A) After OTUs A and B are clustered						
	AB	С	D	E	F	
AB C D E F	0 6 4 16 16	6 0 6 16 16	4 6 0 16 16	16 16 16 0 10	16 16 16 10 0	

(B) After OTUs AB and D are clustered

====				====
	ABD	С	Е	F
ABD	0	6	16	16
С	6	0	16	16
Е	16	10	0	10
F	16	16	10	0

(C) After OTUs ABD and C are clustered						
ABDC E F						
ABDC	0	16	16			
Е	16	0	10			
F	16	10	0			

(D) After OTUs E and F are clustered

	ABDC	EF
ABDC EF	0 16	16 0

Figure 16-1: Phylogenetic trees corresponding to distance matices shown in Table 16-1



Figure 16-2: A phylogenetic tree constructed by using UPGMA from distance matix of Table 16-1B



Fig. 16.3 An unrooted tree for three OTUs



$$\begin{split} D_{12} &= BL_{10} + BL_{20}, \\ D_{13} &= BL_{10} + BL_{30}, \\ D_{23} &= BL_{20} + BL_{30}. \end{split}$$

$$BL_{10} = \frac{\left[D_{12} + D_{13} - D_{23}\right]}{2},$$
$$BL_{20} = \frac{\left[D_{12} + D_{23} - D_{13}\right]}{2},$$
$$BL_{30} = \frac{\left[D_{13} + D_{23} - D_{12}\right]}{2}.$$

Fig. 16.4 Explanation of the distance Wagner method



Figure 16-6: Two types of multifurcating trees

(A) Star phylogeny with no interior branch. (B) Multifurcating tree with only one interior branch.



Figure 16-7: Construction of the neighbor-joining tree for 8 OTUs (from ref 16-1)







(F)



Figure 16-8: Expected tree for 10 sequences which produced the distance matrix of Table 16-1C



Figure 16-9: A phylogenetic network with two splits for four OTUs



Figure 16-10: Two phylogenetic networks constructed from a distance matrix data of 2? gibbon ABO blood group gene partial sequences (from ref 16-43)

(A) When the Split Decomposition method (ref 16-42) was used. (B) When the Neighbor-Net method (ref 16-44; ref 16-45) was used



Figure 16-11: The maximum parsimony tree (left) for the sequence data (right)



Table 16-7: Ten informative configurations for five nucleotide sequences

	Α	В	С	D	Е	
	 v	 v	*	*	 *	
1 2	^ X	^ *	X	*	*	
3	Х	*	*	Х	*	
4	Х	*	*	*	Х	
5	*	Х	Х	*	*	
6	*	Х	*	Х	*	
7	*	Х	*	*	Х	
8	*	*	Х	Х	*	
9	*	*	Х	*	Х	
10	*	*	*	Х	Х	

Note: * and X are diffent nucleotides with each ohter.

Table 16-8: Application of the SSJ algorithm to an ideal sequence data

Step 0: Initial sequences after elimination of invarianct sites and sequence identity check		
	00000000011111111112222 12345678901234567890123	
Α	aacgtttcattgagatacgtgca	
В	gc	
С	cg.g	
D	ctagtt	
Е	ctacgagcg	
F	ctacgagtga	
G	ctacgagtgca	
н	ctacgag.c.gt.	
Ι	ctacaaa.caaa.	
J	ctacgag.cggc	

Step 4:		Step 5:		Step 6	Step 6:	
====	2345679	====	4567		4567	
ABC D E FG HIJ	acgttta ta tacga tacgag. tacgagc	ABC D E FG HIJ	gttt cga. cgag cgag	ABCD E FGHIJ	gttt cga. cgag	
====		====				

Ste	p 1: After
eli	mination of non
inf	ormative sites
==	
	0000000001
	1234567890
Α	aacgtttcat
В	
С	c
D	cta
Е	ctacga
F	ctacgagt
G	ctacgagt
Н	ctacgag.c.
Ι	ctacgag.cg
J	ctacgag.cg

	Step 2: After	Step 3: After	
on-	joining identical	elimination o	
5	sequences	informative s	
	0000000001	0000000	
	1234567890	2345679	
	AB aacgtttcat	AB acgttta	
	С с	с	
	D cta	D ta	
	E ctacga	E tacga	
	FG ctacgagt	FG tacgag.	
	H ctacgag.c.	H tacgagc	
	IJ ctacgag.cg	IJ tacgagc	

Step 3: After		
elimination of non-		
informative sites		
	0000000	
	2345679	
AB	acgttta	
С		
D	ta	
Е	tacga	
FG	tacgag.	
Н	tacgagc	
IJ	tacgagc	

Figure 16-12: The maximum parsimony tree for sequence data shown in Table 16-8



Fig. 16.13 An unrooted tree of six sequences



$$\mathbf{L}[i] = \sum_{\mathbf{D}} \left[g P_{\mathbf{D}5} P_{\mathbf{D}6} \left\{ \sum_{\mathbf{C}} P_{DC} P_{\mathbf{C}4} \left\{ \sum_{\mathbf{B}} P_{CB} P_{\mathbf{B}3} \left\{ \sum_{\mathbf{A}} P_{\mathbf{B}A} P_{\mathbf{A}1} P_{\mathbf{A}2} \right\} \right\} \right\} \right],$$

Figure 16-14: Likelihood values for three possible trees with four sequences (from ref 16-71)



Figure 16-15: Application of Saitou's NJ-like stepwise clustering search using ML method (from ref 16-74)



Figure 16-16: A tetrahedron for four nucleotides



Figure 16-17: A phylogenetic network with three splits for four sequences



a 16S rRNA gene



b Concatenated gene

Fig. 16.18 Phylogenetic trees of bacterial species with GC content 32–38 %. (a) 16S rRNA genebased tree. (b) Concatenated gene-based tree. (c) Tri-nucleotide frequency-based tree

Chapter 17 Population Genomics

- 17.1 Evolutionary Distances Between Populations
- 17.2 Mitochondrial DNA Population Genomics17.3 Population Genomics of Prokaryotes17.4 Population Genomics of Nuclear Genomes

Figure 17-1: A schematic gene genealogy of two populations which differenciated long time ago



Figure 17-2: A schematic gene genealogy of two populations which differenciated recently


Figure 17-3: Overlayed gene genealogies of genes 1-4 sampled from two species A and B



Figure 17-4: Dynamics of allele frequency changes during the population differenciation (from ref 17-6)





Fig. 17.5 A phylogenetic tree of N haplogroup mtDNA sequences of Malaysians (From Jinam et al. 2012; [16])



Fig. 17.6 Phylogenetic tree of major haplogroups of human mtDNAs (From http://www. phylotree.org/)

Figure 17-7: distributions of sequence differences between all possible pairs of individuals for non-African populations (from ref 17-19)





Fig. 17.8 An example of the Bayesian Skyline Plot (From Jinam et al. 2012; [16])



Fig. 17.10 Various kinship relationships (kindly provided by Sarah Voisin)



Fig. 17.11 PCA analysis of genome-wide SNP data of Europeans (From Ref. [37])

Figure 17-12: Phylogenetic tree of 9 human populations in East Asia (from ref 17-40)





Fig. 17.13 PCA plots of genome-wide SNP data of human individuals in Southeast Asia (From Ref. [38]). (a) Eight populations. (b) Thirteen populations



Fig. 17.14 STRUCTURE result of 13 human populations (From Jinam et al. 2013; [46])



Fig. 17.15 Correlation of STRUCTURE and PCA results for two Malaysian Negrito populations. Three individuals with arrows are possible recent hybrids with Malays (From Jinam et al. 2013; [46])

Figure 17-16: Phylogenetic networks for four genes (from ref 17-53)





Figure 17-17: P values of microsatellite loci and SNP loci near the ABCC11 gene (from ref 17-61)