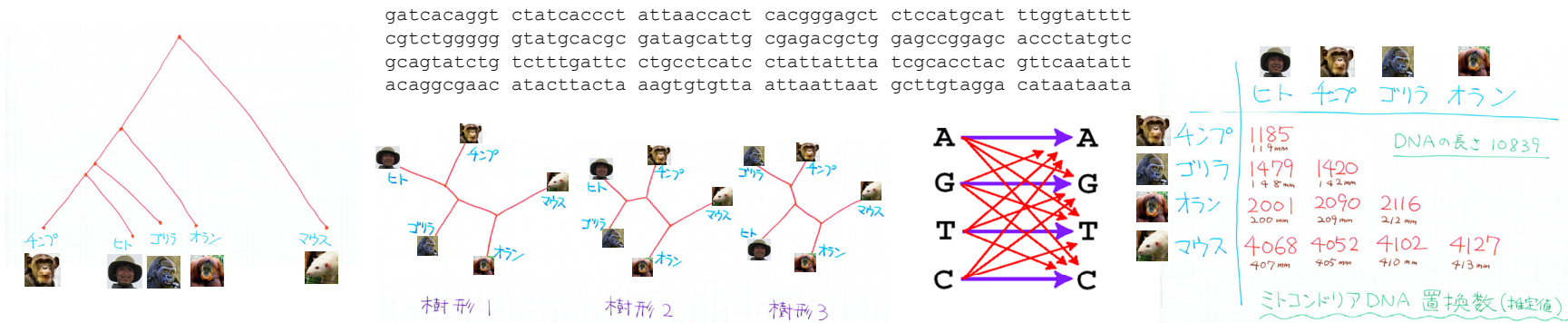


DNA情報からよみとる生物進化 とランダムネス

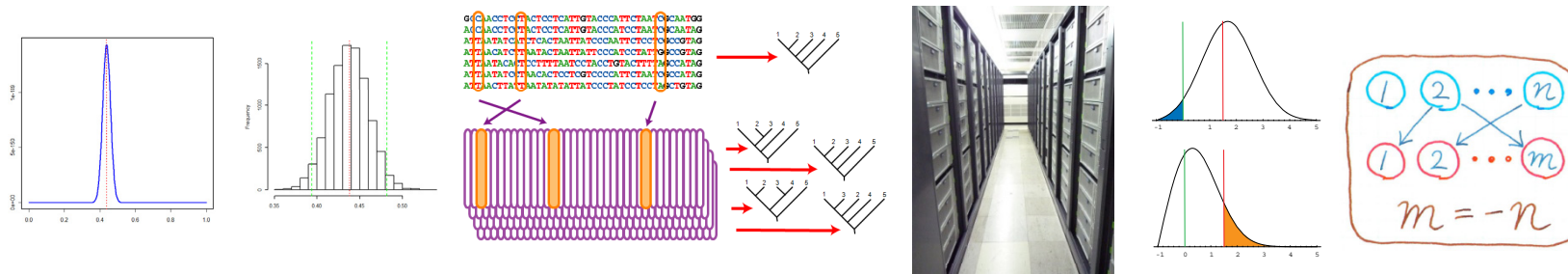
東京工業大学 情報理工学研究科
数理・計算科学専攻
下平英寿

概要

• 系統樹: DNA配列 分子進化 工作



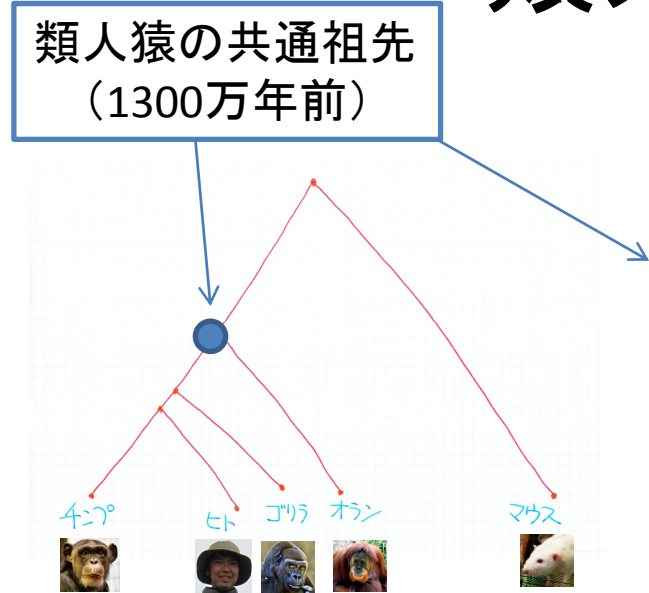
• 統計学: 最尤法 (未知量の推定方法) ブートストラップ法 (ランダムネスを調べる方法)



下平研の成果: Shimodaira-Hasegawa test, マルチスケール・ブートストラップ法

系統樹

類人猿の系統樹



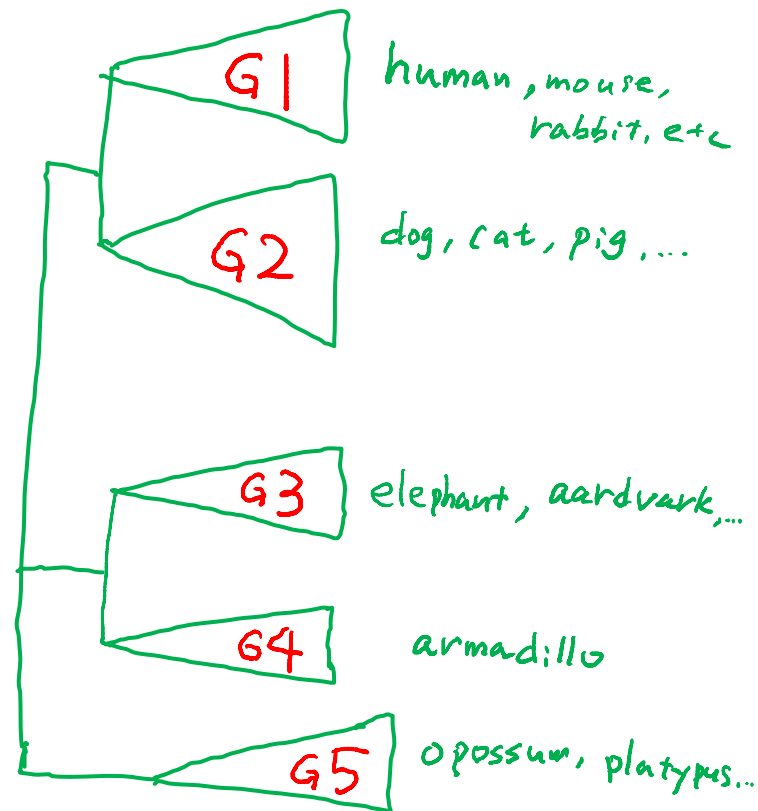
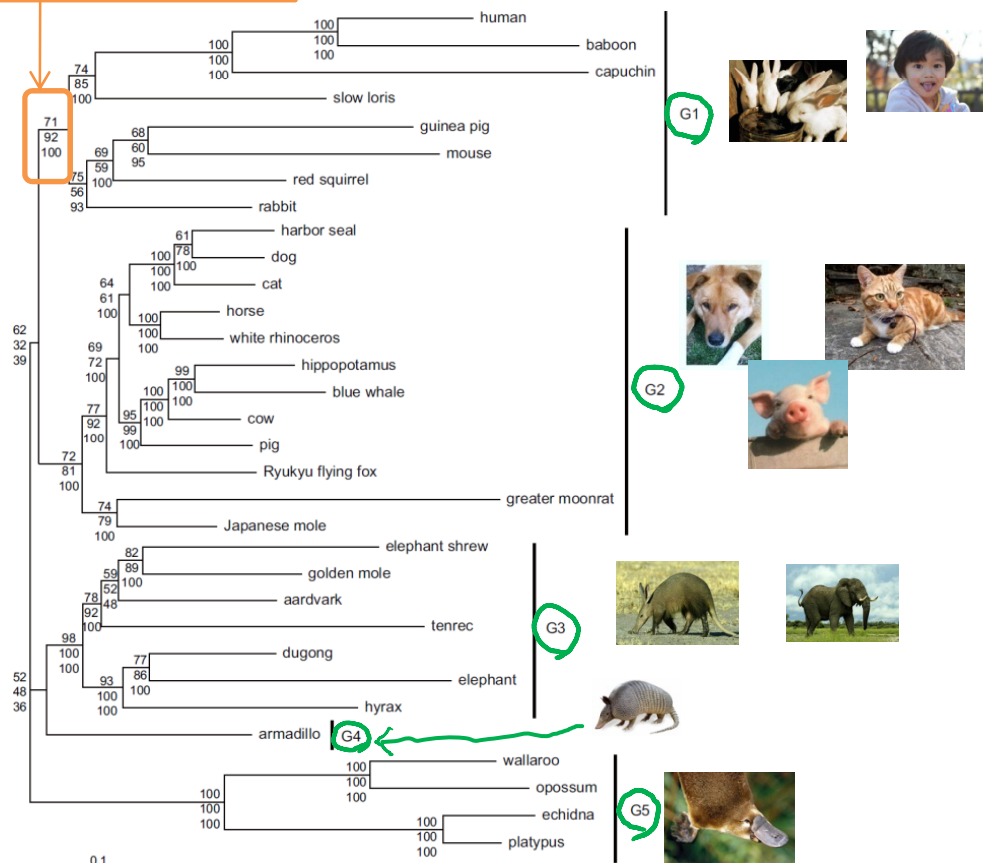
あとで工作します



NC_001807: Homo sapiens: human (old RefSeq)
 NC_012920: Homo sapiens: human
 NC_001643: Pan troglodytes: chimpanzee
 NC_001645: Gorilla gorilla: Western Gorilla
 NC_002083: Pongo abelii: Sumatran orangutan
 NC_010339: Mus musculus musculus: eastern European house mouse

哺乳類の系統樹 (32種)

枝の小さい数字は
推定結果の信頼度

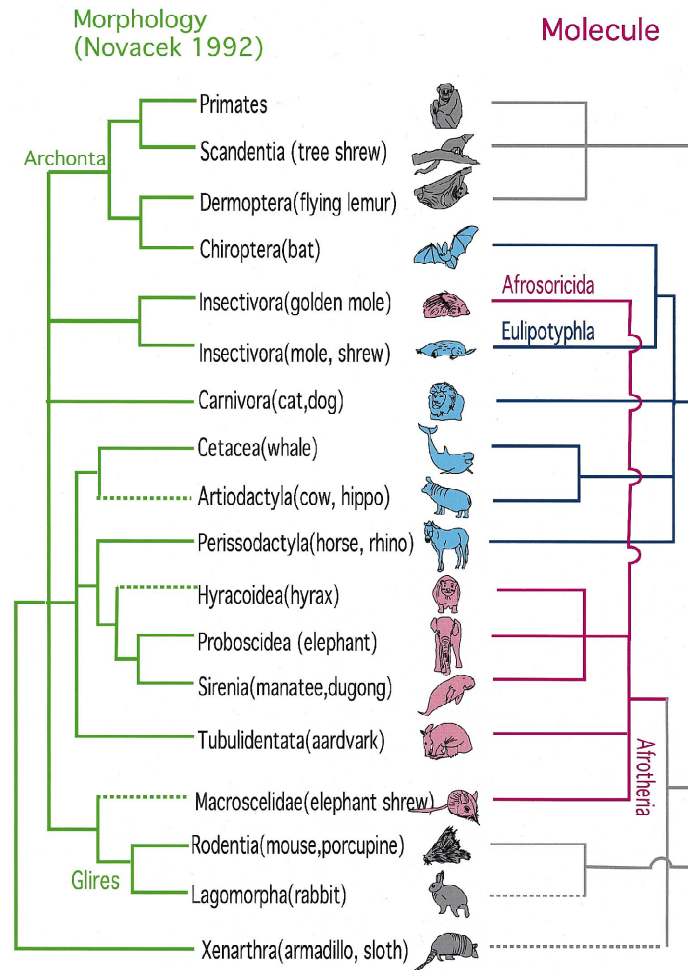


ML tree topology: ((G1,G2), (G3,G4), G5)

Fig 1 of Shimodaira and Hasegawa (2005) from the book (ed. Nielsen)
Data: mt protein sequences of n=3392 amino acids for s=32 species

見た目 vs 分子データ

形態学による
古典的な系統樹



ミトコンドリアDNAから
統計手法(最尤法)
によって推定した系統樹

Cao et al. (2000) *Gene*

ゲノムの情報量

- ある生物種の個体全体を完全な状態に保つために必要な遺伝情報の1セット
- 多数の塩基(A,T,G,C)が並んだDNA分子で遺伝子が記述される
- 1塩基(対)=2ビット, つまり4塩基=1バイト
- ヒトゲノムは3G塩基 = 0.75Gバイト



(730円)

2Gバイトのメモリに
余裕で入る



(10709円)

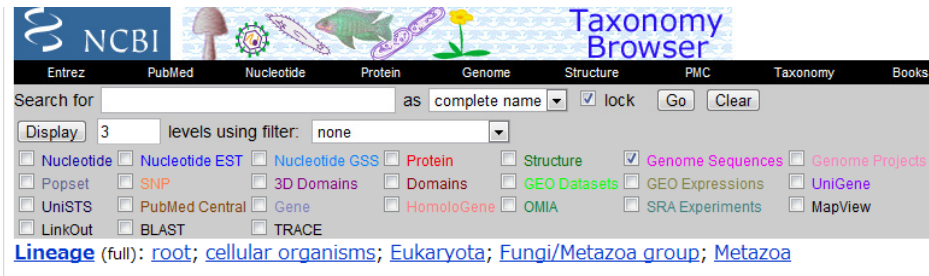
2000Gバイトのハードディスクに2500人の
ゲノムが圧縮無しでそっくり入る

DNAは小さな細胞核に有ることを考えれば, 驚くことではない?

分類	生物種	塩基数	総遺伝子数
哺乳類	Homo sapiens (ヒト)	3Gb	30000
昆虫	Drosophila melanogaster (ショウジョウバエ)	120Mb	13000
植物	Arabidopsis thaliana (シロイヌナズナ)	125Mb	25498
植物	Oryza sativa (イネ)	430Mb	50000?
線虫	Caenorhabditis elegans (線虫)	97Mb	18000
酵母	Saccharomyces cerevisiae (出芽酵母)	12Mb	6286
細菌	Bacillus subtilis (枯草菌)	4.21Mb	4100
細菌	Escherichia coli K12 MG1655 (大腸菌)	4.64Mb	4289
細菌	Escherichia coli O157 Sakai (大腸菌 O157)	5.5Mb	5361
細菌	Haemophilus influenzae Rd (インフルエンザ菌)	1.83Mb	1709
細菌	Helicobacter pylori J99 (ピロリ菌)	1.64Mb	1491
その他	ミトコンドリア	16.6Kb	13

DNAデータの入手 (NCBI)

<http://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?mode=Tree&id=6072&lvl=3&p=genome&lin=f&keep=1&srchmode=1&unlock>



National Center for Biotechnology Information (NCBI)のデータベース

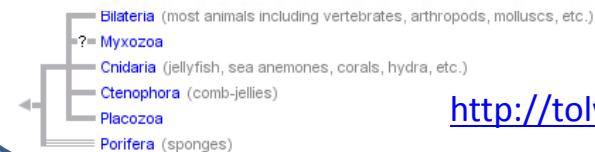
- o [Eumetazoa](#) 2,593 *Click on organism name to get more information.*
 - o [Bilateria](#) 2,555
 - o [Acoelomata](#) 33
 - o [Platyhelminthes](#) (flatworms) 33
 - o [Coelomata](#) 2,461
 - o [Deuterostomia](#) 1,972
 - o [Protostomia](#) 489
 - o [Pseudocoelomata](#) 61
 - o [Acanthocephala](#) (thorny-headed worms) 1
 - o [Cycliophora](#)
 - o [Gastrotricha](#) (gastrotrichs)
 - o [Kinorhyncha](#) (mud dragons)
 - o [Loricifera](#) (loriciferans)
 - o [Micrognathozoa](#) (micrognathozoans)
 - o [Nematoda](#) (roundworms) 52
 - o [Nematomorpha](#) (horsehair worms)
 - o [Rotifera](#) (rotifers) 3
 - o [Bilateria incertae sedis](#)
 - o [Gnathostomulida](#)
 - o [environmental samples](#)
 - o [Bilateria environmental sample](#)
 - o [Cnidaria](#) (cnidarians) 38

真正後生動物

こちらは参考ページです

Animals

Metazoa



<http://tolweb.org/Animals/2374>

リンクをたどっていく (脊椎動物)

http://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?mode=Tree&id=7742&lvl=3&p=genome&p=mapview&p=has_linkout&p=blast_url&p=genome_blast&lin=f&keep=1&srchmode=1&unlock

NCBI Taxonomy Browser

Entrez PubMed Nucleotide Protein Genome Structure PMC Taxonomy Books

Search for _____ as complete name [lock] [Go] [Clear]

Display 3 levels using filter: none

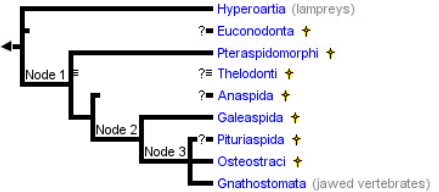
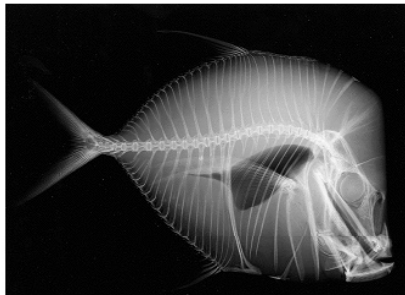
Nucleotide Nucleotide EST Nucleotide GSS Protein Structure Genome Sequences Genome Projects
 Popset SNP 3D Domains Domains GEO Datasets GEO Expressions UniGene
 UniSTS PubMed Central Gene HomoloGene OMA SRA Experiments MapView
 LinkOut BLAST TRACE

Lineage (full): [root](#); [cellular organisms](#); [Eukaryota](#); [Fungi/Metazoa group](#); [Metazoa](#); [Eumetazoa](#); [Bilateria](#); [Coelomata](#); [Deuterostomia](#); [Chordata](#); [Craniata](#)

- **Vertebrata** (vertebrates) [1,911 LinkOut](#) *Click on organism name to get more information.*
 - **Gnathostomata** (jawed vertebrates) [1,909 LinkOut](#)
 - **Chondrichthyes** (cartilaginous fishes) [10 LinkOut](#)
 - **Elasmobranchii** (elasmobranchs) [9 LinkOut](#)
 - **Holocephali** [1 LinkOut](#)
 - **Teleostomi** [1,899](#)
 - **Euteleostomi** (bony vertebrates) [1,899](#)
 - **Hyperoartia** (fish) [2 LinkOut](#)
 - **Petromyzontiformes** [2 LinkOut](#)
 - **Petromyzontidae** (lampreys) [2 LinkOut](#)

脊椎動物

Vertebrata
Animals with backbones
Philippe Janvier



10

さらにたどって (真獣下綱)

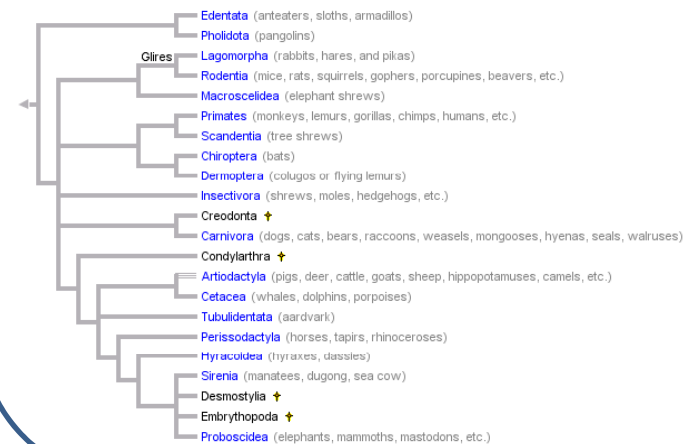
http://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?mode=Tree&id=9347&lvl=3&p=genome&p=mapview&p=has_linkout&p=blast_url&p=genome_blast&lin=f&keep=1&srchmode=1&unlock

The screenshot shows the NCBI Taxonomy Browser search results for 'Eutheria'. The search bar contains 'Eutheria' and the results are displayed as a tree structure. The 'Eutheria' node is highlighted, and its lineage is shown as: root; cellular organisms; Eukaryota; Fungi/Metazoa group; Metazoa; Eumetazoa; Bilateria; Coelomata; Deuterostomia; Chordata; Craniata; Vertebrata; Gnathostomata; Teleostomi; Euteleostomi; Sarcopterygii; Tetrapoda; Amniota; Mammalia; Theria. Below the lineage, the 'Eutheria' node is listed with 682 linkouts. The tree structure below 'Eutheria' includes: Afrotheria (14), Chrysochloridae (Golden moles) (2 linkouts), Hyracoidea (rock rabbits/dassies) (2 linkouts), Macroscelidea (elephant shrews) (2 linkouts), Proboscidea (elephants) (4 linkouts), and Sirenia (manatees and dugongs (seacows)) (2 linkouts).

真獣下綱: 胎盤を持つ哺乳類

Eutheria

Placental Mammals



やっとつきました (ヒト上科)

http://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?mode=Tree&id=9604&lvl=3&p=genome&p=mapview&p=has_linkout&p=blast_url&p=genome_blast&lin=f&keep=1&srchmode=1&unlock

NCBI Taxonomy Browser

Search for as complete name lock

Display 3 levels using filter: none

Nucleotide Nucleotide EST Nucleotide GSS Protein Structure Genome Sequences Genome Projects
 Popset SNP 3D Domains Domains GEO Datasets GEO Expressions UniGene
 UniSTS PubMed Central Gene HomoloGene OMA SRA Experiments MapView
 LinkOut BLAST TRACE

Lineage (full): [root](#); [cellular organisms](#); [Eukaryota](#); [Fungi/Metazoa group](#); [Metazoa](#); [Eumetazoa](#); [Bilateria](#); [Coelomata](#); [Deuterostomia](#); [Chordata](#); [Craniata](#); [Vertebrata](#); [Gnathostomata](#); [Teleostomi](#); [Euteleostomi](#); [Sarcopterygii](#); [Tetrapoda](#); [Amniota](#); [Mammalia](#); [Theria](#); [Eutheria](#); [Euarchontoglires](#); [Primates](#); [Haplorrhini](#); [Simiiformes](#); [Catarrhini](#); [Hominoidea](#)

o [Hominidae](#) (great apes) [132 LinkOut](#) *Click on organism name to get more information.*

o [Homininae](#) [106 LinkOut](#)

o [Gorilla](#) [2 LinkOut](#)

- o [Gorilla beringei](#) (Eastern Gorilla) [LinkOut](#)
- o [Gorilla gorilla](#) (Western Gorilla) [2 LinkOut](#)

ゴリラ

o [Homo](#) [76 LinkOut](#)

- o [Homo sapiens](#) (human) [75 MapView LinkOut BLAST page](#)
- o [Homo sp. Altai](#) (Denisova hominin) [1](#)

ヒト

o [Pan](#) (chimpanzees) [28 LinkOut](#)

- o [Pan paniscus](#) (pygmy chimpanzee) [1 LinkOut](#)
- o [Pan troglodytes](#) (chimpanzee) [27 MapView LinkOut BLAST page](#)
- o [Pan sp. WO2007026120](#)

チンパンジー

o [Ponginae](#) [26 LinkOut](#)

o [Pongo](#) [26 LinkOut](#)

- o [Pongo abelii](#) (Sumatran orangutan) [25 LinkOut BLAST page](#)
- o [Pongo abelii x pygmaeus](#) [LinkOut](#)
- o [Pongo pygmaeus](#) (Bornean orangutan) [1 LinkOut](#)
- o [Pongo sp.](#) [LinkOut](#)

オランウータン

ヒト上科

Hominidae

Humans, great apes, and their extinct relatives



Homo sapiens (ヒト)

NCBI Taxonomy Browser

Search for: as complete name lock

Display: 3 levels using filter: none

Homo sapiens

Taxonomy ID: 9606
 Genbank common name: human
 Inherited blast name: primates
 Rank: species
 Genetic code: [Translation table 1 \(Standard\)](#)
 Mitochondrial genetic code: [Translation table 2 \(Vertebrate Mitochondrial\)](#)
 Other names:
 common name: **man**
 authority: **Homo sapiens Linnaeus, 1758**

[Lineage \(full\)](#)
[cellular organisms](#); [Eukaryota](#); [Fungi/Metazoa group](#); [Metazoa](#); [Eumetazoa](#); [Bilateria](#); [Coelomata](#); [Deuterostomia](#); [Chordata](#); [Craniata](#); [Vertebrata](#); [Gnathostomata](#); [Teleostomi](#); [Euteleostomi](#); [Sarcopterygii](#); [Tetrapoda](#); [Amniota](#); [Mammalia](#); [Theria](#); [Eutheria](#); [Euarchontoglires](#); [Primates](#); [Haplorrhini](#); [Simiiformes](#); [Catarrhini](#); [Hominoidea](#); [Hominidae](#); [Homininae](#); [Homo](#)

Entrez records		
Database name	Subtree links	Direct links
Nucleotide	16,985,984	16,985,959
Nucleotide EST	8,301,515	8,301,515
Nucleotide GSS	1,293,831	1,292,505
Protein	531,420	531,324
Structure	15,730	15,730
Genome Sequences	75	74
Genome Projects	32	32
Popset	18,170	18,170
SNP	31,020,464	31,020,464
3D Domains	61,175	61,175
Domains	8	8
GEO Datasets	9,564	9,564
GEO Expressions	17,689,684	17,689,684
UniGene	123,200	123,200
UniSTS	327,522	327,522
PubMed Central	7,918	7,915
Gene	42,645	42,608
HomoloGene	18,876	18,876
SRA Experiments	8,821	8,821
Taxonomy	2	1

Genome Information

[See the NCBI Genome homepage](#)
[Go to NCBI genomic BLAST page for Homo sapiens](#)

Genome view: 24 chromosomes

Names:

[See the Mitochondrion Genome](#)
[See the TRACE Assembly](#)

Trace records (raw single-pass reads of DNA sequence)

Sequencing Center Name
Record counts per type

ヒトゲノムの選択ページ

NCBI Genome

Search: Genome for

Display: Summary Show 20 Send to

All: 75

Items 1 - 20 of 75 Page 1 of 4 Next

Recent activity

- 1: [NC_012920](#)
 Homo sapiens mitochondrion, complete genome
dsDNA; circular; Length: **16,569 nt**
 Organelle: **mitochondrion**
 Created: **2009/07/08**
- 2: [AC_000156](#)
 Homo sapiens chromosome Y, alternate assembly HuRef, whole genome shotgun sequence
DNA; linear; Length: **19,317,006 nt**
 Replicon Type: **chromosome**
 Replicon Name: **Y**
 Created: **2007/09/27**
- 3: [AC_000155](#)
 Homo sapiens chromosome X, alternate assembly HuRef, whole genome shotgun sequence
DNA; linear; Length: **143,733,266 nt**
 Replicon Type: **chromosome**
 Replicon Name: **X**
 Created: **2007/09/27**
- 4: [AC_000154](#)
 Homo sapiens chromosome 22, alternate assembly HuRef, whole genome shotgun sequence
DNA; linear; Length: **34,107,095 nt**
 Replicon Type: **chromosome**
 Created: **2007/09/27**

ミトコンドリアを選択

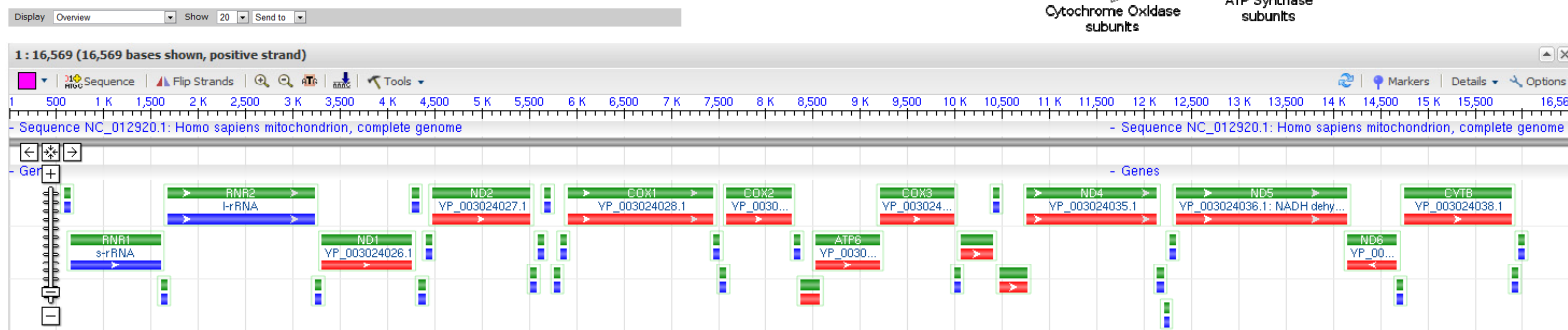
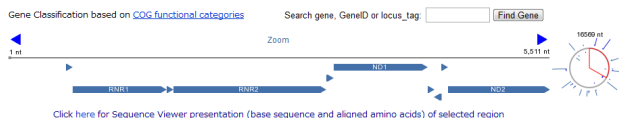
ヒト ミトンドリア ゲノム

<http://www.ncbi.nlm.nih.gov/sites/entrez?Db=genome&Cmd=ShowDetailView&TermToSearch=24730>

NCBI Genome search results for *Homo sapiens mitochondrion, complete genome*. The search shows 1 result with 1 link. The genome length is 16,569 nt.

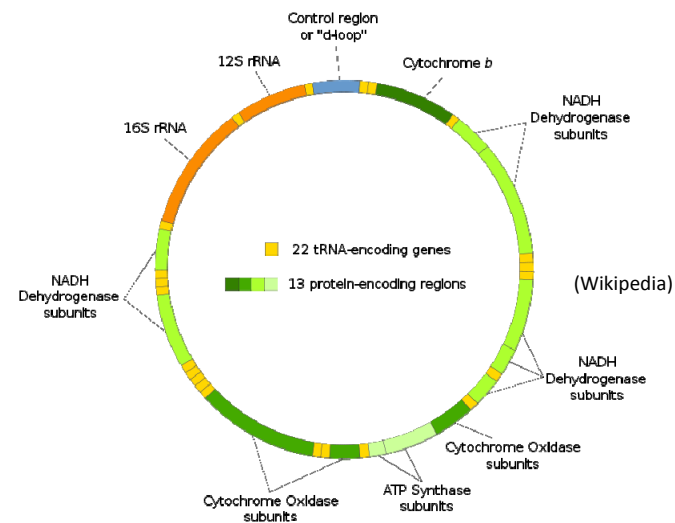
Lineage: Eukaryota; Fungi/Metazoa group; Metazoa; Fumetazoa; Bilateria; Coelomata; Deuterostomia; Chordata; Craniata; Vertebrata; Gnathostomata; Teleostomi; Euteleostomi; Sarcopterygii; Tetrapoda; Amniota; Mammalia; Theria; Eutheria; Euarchontoglires; Primates; Haplorhini; Simiiformes; Catarrhini; Hominoidea; Hominoidea; Hominidae; homininae; Homo; *Homo sapiens*

Genome Info:	Features:	BLAST homologs:	Links:	Review Info:
Refseq: NC_012920	Genes: 37	COG	Genome Project	Publications: [2]
GenBank: J01415	Protein coding: 13	TaxMap	Refseq FTP	Refseq Status: PROVISIONAL
Length: 16,569 nt	Structural RNAs: 24	TaxPlot	GenBank FTP	Seq Status: Completed
GC Content: 44%	Pseudo genes: None	GenePlot	BLAST	Sequencing center: Center for Molecular and Mitochondrial Medicine and Genetics (MAMMAG) University of California, University of California, Irvine, Mitomap.org, USA, Irvine
% Coding: 68%	Others: 705	glmap	TraceAssembly	Completed: 2009/07/08
Topology: circular	Contigs: None		COG	Organism Group
Molecule: dsDNA			Other genomes for species: 3756	



環状ゲノムで長さは16000文字程度
このデータは16569 bp

遺伝子の名前: ND1, ND2, COX1, COX2, ...



(Wikipedia)

NC_012920

NCBI Resources How To My NCBI Sign In

Nucleotide Search: Nucleotide Limits Advanced search Help

Alphabet of Life Search Clear

Display Settings GenBank **Send**

Homo sapiens mitochondrion, complete genome

NCBI Reference Sequence: NC_012920.1

Comment Features Sequence

LOCUS NC_012920 18589 bp DNA circular PRI 30-APR-2010

DEFINITION Homo sapiens mitochondrion, complete genome.

ACCESSION NC_012920 AC_000001

VERSION NC_012920.1 GI:251831108

DBLINK Project:30353

KEYWORDS mitochondrion Homo sapiens (human)

SOURCE

ORGANISM [Homo sapiens](#)

Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominidae; Homo.

REFERENCE 1 (bases 1 to 18589)

AUTHORS Andrews,R.M., Kubacka,I., Chinnery,P.F., Lightowlers,R.N., Turnbull,D.M. and Howell,N.

TITLE Reanalysis and revision of the Cambridge reference sequence for human mitochondrial DNA

JOURNAL Nat. Genet. 23 (2), 147 (1999)

PUBMED [10508508](#)

REFERENCE 2 (bases 324 to 743)

AUTHORS Andrews,R.M., Kubacka,I., Chinnery,P.F., Lightowlers,R.N., Turnbull,D.M. and Howell,N.

TITLE Reanalysis and revision of the Cambridge reference sequence for human mitochondrial DNA

JOURNAL Nat. Genet. 23 (2), 147 (1999)

PUBMED [10508508](#)

REFERENCE 3 (bases 1 to 18589)

AUTHORS Anderson,S., Bankier,A.T., Barrell,B.G., de Bruijn,M.H., Coulson,A.R., Drouin,J., Eperon,I.C., Nierlich,D.P., Roe,D.A., Sanger,F., Schreier,P.H., Smith,A.J., Staden,R. and Young,I.G.

TITLE Sequence and organization of the human mitochondrial genome

JOURNAL Nature 290 (5806), 457-465 (1981)

PUBMED [7219534](#)

REFERENCE 4 (bases 15888 to 15954)

AUTHORS Anderson,S., Bankier,A.T., Barrell,B.G., de Bruijn,M.H., Coulson,A.R., Drouin,J., Eperon,I.C., Nierlich,D.P., Roe,B.A., Sanger,F., Schreier,P.H., Smith,A.J., Staden,R. and Young,I.G.

TITLE Sequence and organization of the human mitochondrial genome

JOURNAL Nature 290 (5806), 457-465 (1981)

PUBMED [7219534](#)

REFERENCE 5 (bases 1 to 18589)

CONSTRM NCBI Genome Project

TITLE Direct Submission

JOURNAL Submitted (08-JUL-2009) National Center for Biotechnology Information, NIH, Bethesda, MD 20894, USA

REFERENCE 6 (bases 1 to 18589)

AUTHORS Kogelnik,A.M. and Lott,W.T.

TITLE Direct Submission

JOURNAL Submitted (24-AUG-2006) Mitomap.org, Center for Molecular and Mitochondrial Medicine and Genetics (MAMMAG) University of California, University of California, Irvine, Irvine, CA 92697-3940, USA

REMARK Sequence update by submitter

REFERENCE 7 (bases 1 to 18589)

AUTHORS Kogelnik,A.M. and Lott,W.T.

TITLE Direct Submission

JOURNAL Submitted (18-APR-1997) Center for Molecular Medicine, Emory University School of Medicine, 1462 Clifton Road, Suite 420, Atlanta, GA 30322, USA

REMARK sequence updated

COMMENT PROVISIONAL REFSEQ: This record has not yet been subject to final

Send

Complete Record

Coding Sequences

Choose Destination

File Clipboard

Collections

Download 1 items.

Format

GenBank

Create File

Recent activity

Turn Off Clear

- Homo sapiens mitochondrion, complete genome Nucleotide
- IGKV10RY-1 immunoglobulin kappa variable 1/ORY-1 pseudogene [Homo gene
- NC_012920 Homo sapiens mitochondrion, complete genome Genome
- gi|17981852|ref|NC_001807.4| Genome
- bid9606[Organism:exp] (75) Genome

See more...



NC_012920.txt

今回は6本のDNA配列を用意

NC_001807: Homo sapiens: human ヒトA (古いRefSeq)
NC_012920: Homo sapiens: human ヒトB (現在のRefSeq)
NC_001643: Pan troglodytes: chimpanzee チンパンジー
NC_001645: Gorilla gorilla: Western Gorilla ゴリラ
NC_002083: Pongo abelii: Sumatran orangutan オランウータン
NC_010339: Mus musculus musculus: eastern European house mouse マウス



work-conc-nuc.txt

16000 bp 程度あるが、そのうち 10839 bpを利用する。
遺伝子コーディングリージョンを使う。

12個の遺伝子 ND1 ND2 COX1 COX2 ATP8 ATP6 COX3 ND3 ND4L ND4 ND5 CYTB

```
CONSENSUS ATACCCATG. CCAACCTCCT ACTCCTCATT GTACCCAT.C TAATCGC.AT AGCATTTCCTA
NC_001807 .....G ..... . . . . .T. ....A.. G.....
NC_012920 .....G ..... . . . . .T. ....A.. G.....
NC_001643 .C.....A ..... . . . . .C. ....A.. .....
NC_001645 ...T....G .T....T. .... . . . .T..C. ....C.. .....
NC_002083 ..G....AA T..... . . . . .A..T...C. ....C.. .....T...
NC_010339 G.GTT.T.TA TT..TA... .ACA...C.C ..C..T..T. ....T.. ..C.....
                10                20                30                40                50                60
CONSENSUS ATGCTAACCG AACGAAAAAT TCTAGGCTAT ATACAACCTAC GCAAAGGCCCA CAAC.TTGTA
NC_001807 .....T.... ..... . . . . .G.....
NC_012920 .....T.... ..... . . . . .G.....
NC_001643 ..... . . . . .C ..... . . . .T.. ..A....
NC_001645 ..... . . . . .T..... .T..... .G.C...
NC_002083 ..... . . . . .C.....C.C .C..... .G.. ..A...G
NC_010339 .CA...GTA. ....C..... CT....G... ..A..... T...A....
```


塩基配列の解読法



(Wikipedia)

情報

塩基3文字
(コドン)



アミノ酸1個



アミノ酸を一系列につなげたものがタンパク質

物質

TTT



フェニルアラニン (F または Phe)

GAA



グルタミン酸 (E または Glu)

TTTGAACCTGGGGAT...



FEPGD...



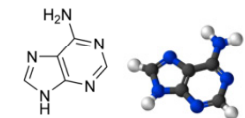
work-conc-nuc.txt



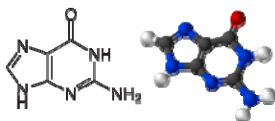
work-conc-ptn.txt

塩基は4種類 (2グループ)

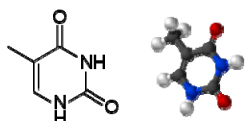
アデニン (A)



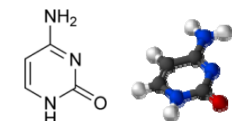
グアシン (G)



チミン (T)



シトシン (C)



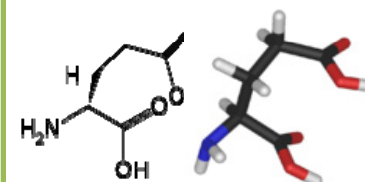
遺伝暗号 変換表2 (脊椎動物ミトコンドリア)

TTT F Phe	TCT S Ser	TAT Y Tyr	TGT C Cys
TTC F Phe	TCC S Ser	TAC Y Tyr	TGC C Cys
TTA L Leu	TCA S Ser	TAA * Ter	TGA W Trp
TTG L Leu	TGG S Ser	TAG * Ter	TGG W Trp
CTT L Leu	CCT P Pro	CAT H His	CGT R Arg
CTC L Leu	CCC P Pro	CAC H His	CGC R Arg
CTA L Leu	CCA P Pro	CAA Q Gln	CGA R Arg
CTG L Leu	CCG P Pro	CAG Q Gln	CGG R Arg
ATT I Ile i	ACT T Thr	AAT N Asn	AGT S Ser
ATC I Ile i	ACC T Thr	AAC N Asn	AGC S Ser
ATA M Met i	ACA T Thr	AAA K Lys	AGA * Ter
ATG M Met i	ACG T Thr	AAG K Lys	AGG * Ter
GTT V Val	GCT A Ala	GAT D Asp	GGT G Gly
GTC V Val	GCC A Ala	GAC D Asp	GGC G Gly
GTA V Val	GCA A Ala	GAA E Glu	GGA G Gly
GTG V Val i	GCG A Ala	GAG E Glu	GGG G Gly

アミノ酸20種類

アミノ酸の例

グルタミン酸



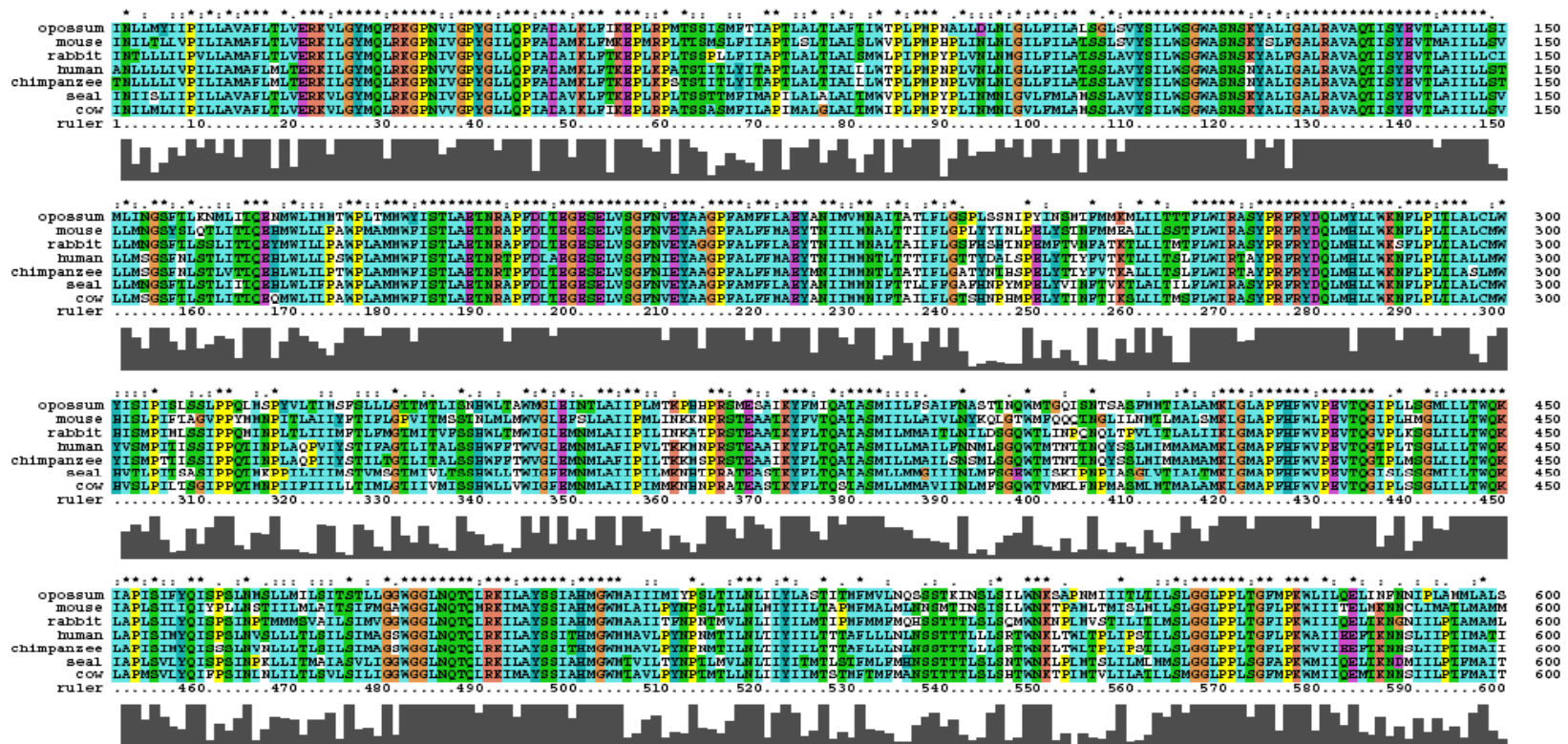
(wikipediaより)

アミノ酸配列データ

CLUSTAL X (1.81) MULTIPLE SEQUENCE ALIGNMENT

File: D: mpmam7bpbn.ps
Page 1 of 6

Date: Wed Sep 25 11:13:46 2002








7種の哺乳類のミトコンドリア・アミノ酸配列の一部

本講演中の哺乳類の系統樹分析は、このアミノ酸配列データを使用した

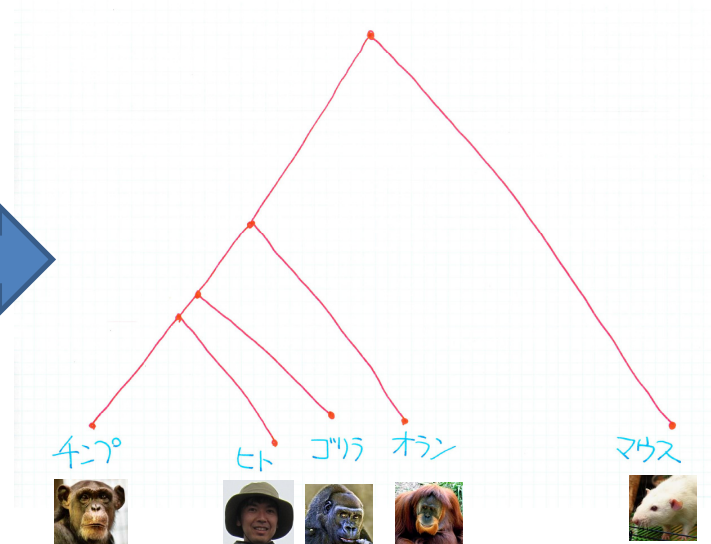
塩基の変異数 (塩基の違いを数えたもの)

DNA配列の長さ: 10839



	NC_001807	NC_012920	NC_001643	NC_001645	NC_002083
 NC_012920		18			
 NC_001643	1063	1061			
 NC_001645	1296	1294	1251		
 NC_002083	1698	1703	1763	1780	
 NC_010339	3148	3147	3134	3158	3180

計算



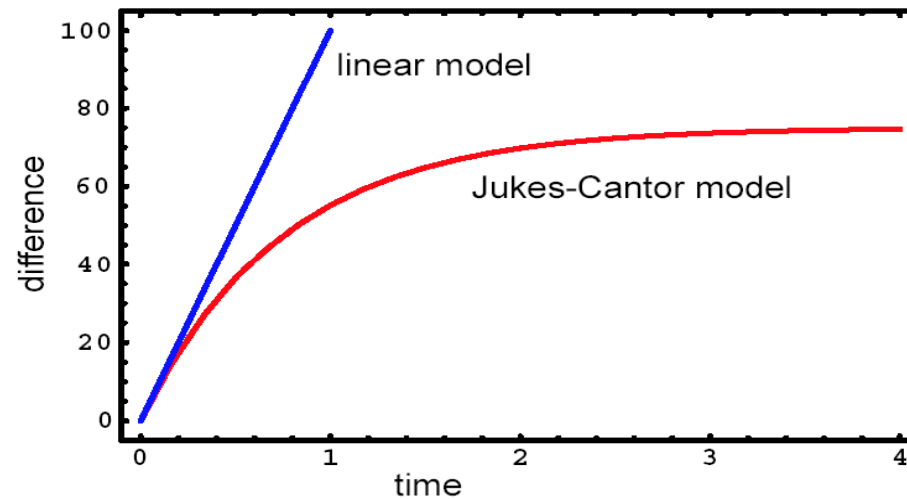
分子時計

- 100万年でDNAの1%が変異すると仮定する
- 200万年でDNAの2%が変異
- 1000万年でDNAの9%が変異
- 1億年でDNAの55%が変異

time	dif	
0	0%	AAAAAAAAAAAAA.....AAAAAAAAAAAA
0.01	1%	AAAA T AAAAA.....AAAAAAAAAAAA
0.1	9%	AAAA T AAAAA.....AAAAA C AAAA
0.2	18%	AAG A T AAAAA.....AA T AA C AAAA
0.3	25%	AAG A AAAAAA.....AA T AA C AAAG
1.0	55%	AAG T TA C AA.....AA T AG C AT G C
3.0	74%	T A C T C G A C T T..... C T A T A G C A G C

分岐時間の推定

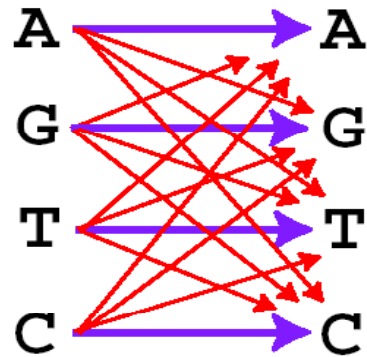
- DNAはランダムに変化
- 確率モデルの一種
- 75%までしか変わらない



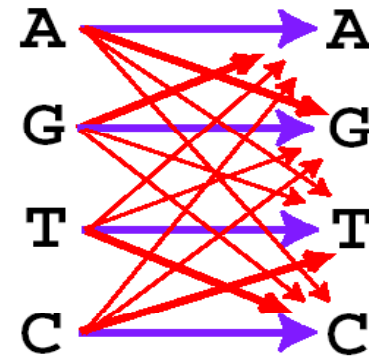
$$P(t) = \frac{3}{4} \left(1 - \exp \left(-\frac{4}{3} \lambda t \right) \right)$$

マルコフ過程

- 塩基置換 (A, T, G, C の 4 種類)



JCモデル



HKYモデル

(トランジションとトランスバージョン, 組成比)

- アミノ酸置換 (20種類)

Dayhoffモデル, **Jones**モデルなど

条件付確率

- Probability of future state- b given current state- a after time- t

$$P_{ba}(t)$$

note: $a, b \in \{A, T, G, C\}$ for nucleotide sequences.

$$P(t) = \begin{pmatrix} P(A|A) & P(G|A) & P(C|A) & P(T|A) \\ P(A|G) & P(G|G) & P(C|G) & P(T|G) \\ P(A|C) & P(G|C) & P(C|C) & P(T|C) \\ P(A|T) & P(G|T) & P(C|T) & P(T|T) \end{pmatrix}$$

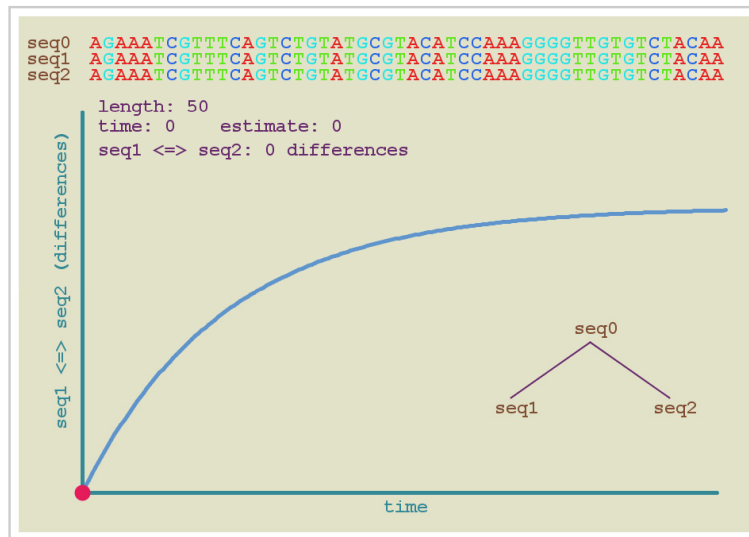
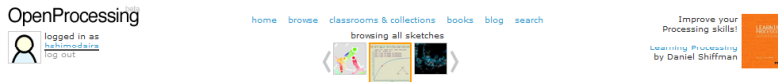
- Transition probability matrix

$$P(t) = \exp(tQ) = \sum_{k \geq 0} \frac{t^k}{k!} Q^k$$

note: $Q = 4 \times 4$ matrix for nucleotide sequences, and $Q = 20 \times 20$ matrix for those of amino-acid.

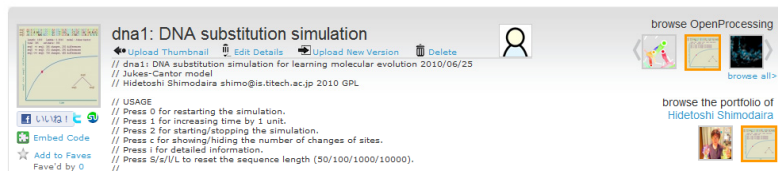
DNA置換のシミュレーション

<http://www.openprocessing.org/visuals/?visualID=10609>



キー操作

- 0: 時間を0にする.
- 1: 時間を1ステップ進める
- 2: 進行のオン, オフ
- c: 置換数のオン, オフ
- l (エル): シーケンス長 = 1000に変更




javaのインストールが必要になります






変異数から置換数へ

DNA配列の長さ: 10839

変異数
(観測値)

計算



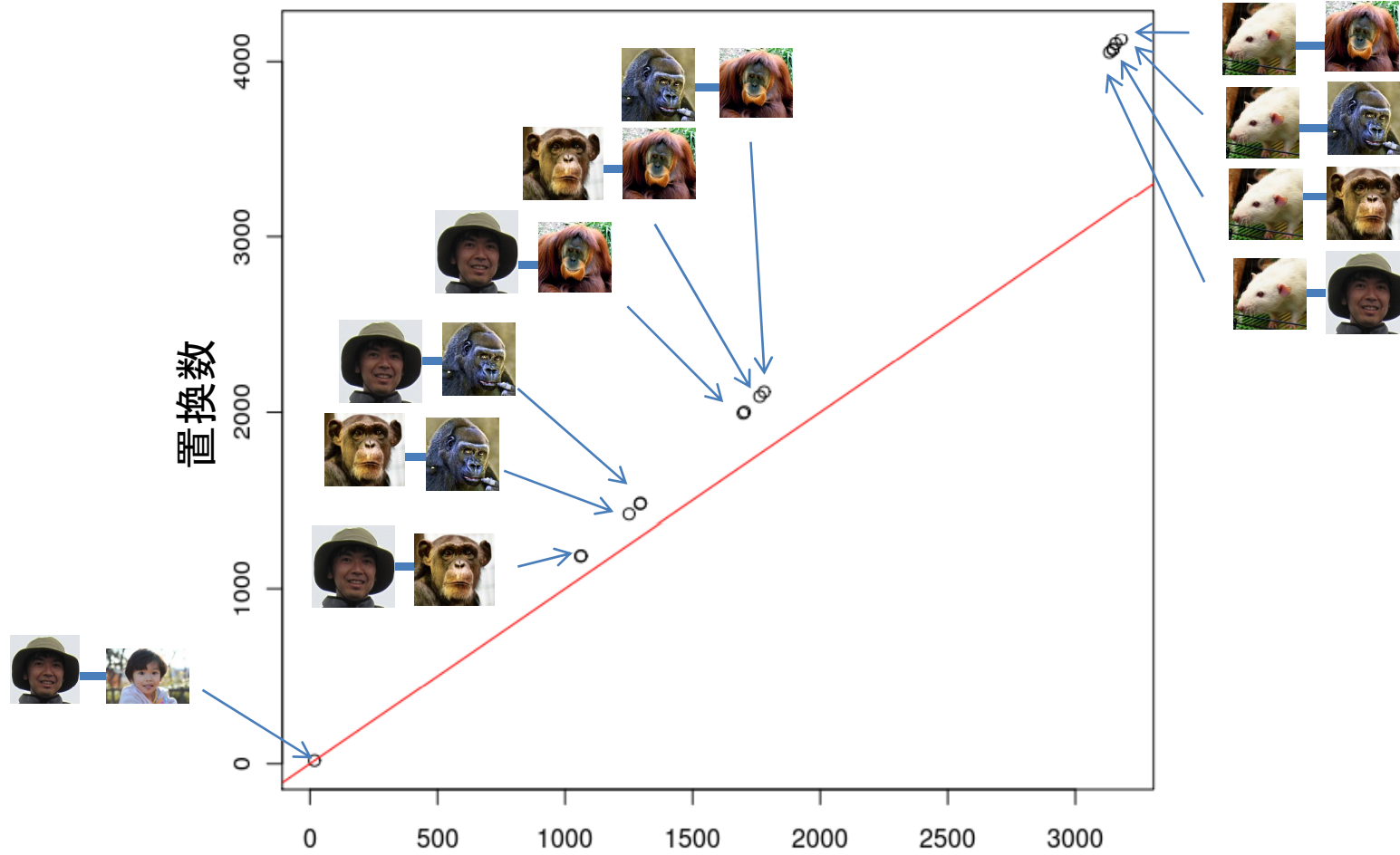
	NC_001807	NC_012920	NC_001643	NC_001645	NC_002083
 NC_012920		18			
 NC_001643	1063		1061		
 NC_001645	1296	1294		1251	
 NC_002083	1698	1703	1763		1780
 NC_010339	3148	3147	3134	3158	3180

	NC_001807	NC_012920	NC_001643	NC_001645	NC_002083
NC_012920		18			
NC_001643	1188		1185		
NC_001645	1481	1479		1420	
NC_002083	1994	2001	2090		2116
NC_010339	4069	4068	4052	4102	4127

注意: 方法の概略を示すために説明を簡略化してある. 実際には変異数を変換したのではなく, DNA配列から最尤法で置換数を推定した. R言語のapeパッケージにてdist.dna(dat, "TN93")を実行. しかし後でJukes-Cantorモデル("JC")を試したところ, 結果は"TN93"とほとんど同じだった. JCならば変異数を変換して置換数を直接計算できる.

NC_001807: Homo sapiens: human
 NC_012920: Homo sapiens: human
 NC_001643: Pan troglodytes: chimpanzee
 NC_001645: Gorilla gorilla: Western Gorilla
 NC_002083: Pongo abelii: Sumatran orangutan
 NC_010339: Mus musculus musculus: eastern European house mouse

変異数と置換数の関係



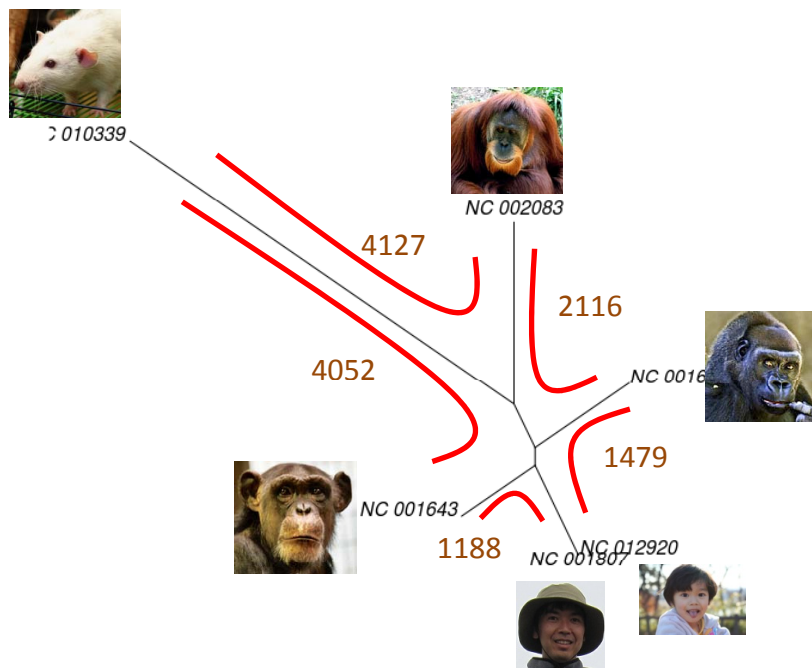
DNA配列の長さ: 10839

変異数

置換数から系統樹を推定

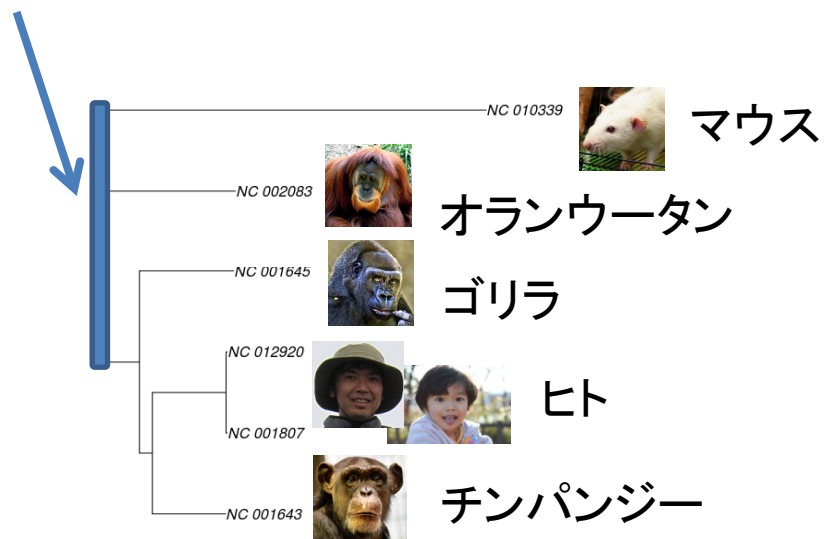
置換数
(推定値)

	NC_001807	NC_012920	NC_001643	NC_001645	NC_002083
NC_012920		18			
NC_001643	1188	1185			
NC_001645	1481	1479	1420		
NC_002083	1994	2001	2090	2116	
NC_010339	4069	4068	4052	4102	4127



推定された系統樹 (無根系統樹)

ココに根があると考える(マウス以外)



推定された系統樹 (有根系統樹)

紙テープで系統樹を工作

2010/06/27

系統樹の工作

東京工業大学 下平英寿 <http://www.is.titech.ac.jp/~shimo/>



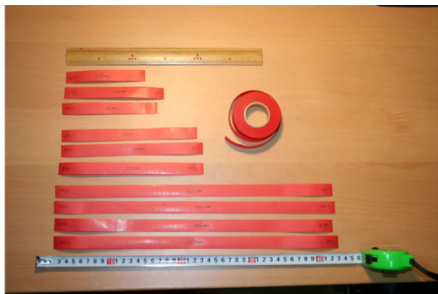
紙テープ(糸でもよい), はさみ, ものさし, のり, ペン, セロテープなど

置換数の長さに紙テープを切る

	ヒト	チンパンジー	ゴリラ	オラン
4ンゴ	1185 11.9mm			
ゴリラ	1479 14.8mm	1420 14.2mm		
オラン	2001 20.0mm	2090 20.9mm	2116 21.2mm	
マウス	4068 40.7mm	4052 40.5mm	4102 41.0mm	4127 41.3mm

DNAの長さ 10839
ミトコンドリアDNA 置換数(推定値)

ここでは置換数10個を1mmにした
1センチがだいたい100万年です



参考: 置換数は最尤法で推定した

- データは類人猿4種(ヒト, チンパンジー, ゴリラ, オランウータン)とマウスのミトコンドリアDNA配列をNCBIのウェブサイトから取得。
- 12個の遺伝子コーディングリージョンをclustalWでアライメントして位置合わせ。長さ10839の配列を得る。
- R言語のapeパッケージにてDNA配列から最尤法で置換数を推定した。dist.dna(dat, "TN93")を実行後、結果を10839倍して整数に丸めた。単純な変異数(塩基が何個異なるか)はdist.dna(dat, "raw")で計算できる。
- 後で確かめたらJukes-Cantor (JC)モデルをつかってもほぼ同じ置換数が得られたので、TN93モデルより解説が容易なJCモデルを工作に使うべきだった。
- 置換数>変異数である。今回のデータでは置換数があまり大きくないので、置換数の代わりに変異数を用いても同じ系統関係が得られた。
- (工作しないで)系統樹まで求めるには、置換数の行列をn関数に与えれば、近隣結合法の無根系統樹が得られる。root()で根を指定する。

YouTubeにも置きました

<http://www.youtube.com/watch?v=jqqJODQytk>

統計学：最尤法

コインが表になる確率 θ の推定

状況設定

コイン投げの実験 500回くりかえす 表=1, 裏=0と書く

実験結果($n=500$) 00111111100011001011011000...

i 回目のコイン投げを x_i と書くと $x_1 = 0, x_2 = 0, x_3 = 1, x_4 = 1, x_5 = 1, \dots, x_{500} = 0$

表になった回数は219回 $C = x_1 + \dots + x_n$ $C = 219$

推定結果

θ の推定値は表になった頻度 $\hat{\theta} = \frac{C}{n} = \frac{219}{500} = 0.438$

ほとんど自明に見えるが、どうしてこの推定量でよいのか？

最尤法 (コイン)

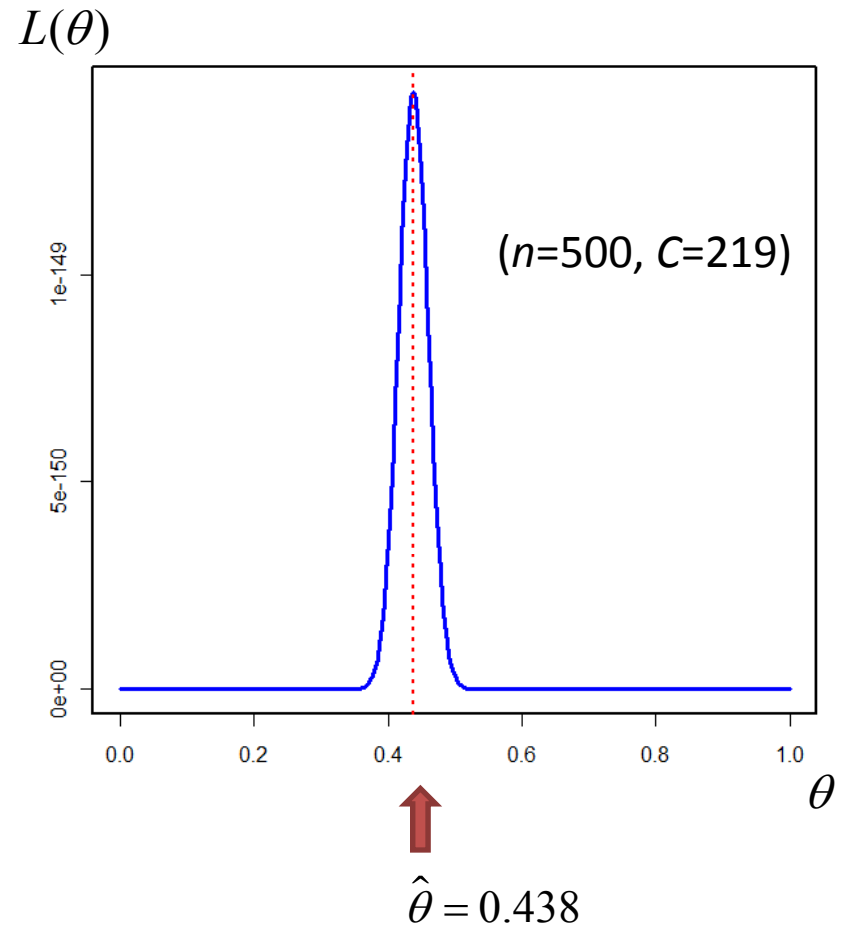
データを観測する確率(尤度)は

$$L(\theta) = \theta^C (1 - \theta)^{n-C}$$

確率を最大にするパラメータ値を選ぶ

$$\max_{\theta} L(\theta) \rightarrow \hat{\theta} \quad \text{最尤推定量}$$

$$\text{尤度を最大にする } \theta \text{ は } \hat{\theta} = \frac{C}{n}$$



最尤法の詳細(コイン)

1回目のコイン投げの結果が x_1 である確率は

$$f(x_1; \theta) = \theta^{x_1} (1 - \theta)^{1 - x_1} = \begin{cases} \theta, & x_1 = 1 \\ 1 - \theta, & x_1 = 0 \end{cases}$$

データ (x_1, \dots, x_n) を観測する確率(尤度)は

$$\begin{aligned} L(\theta) &= f(x_1; \theta) \cdots f(x_n; \theta) \\ &= \theta^{x_1} (1 - \theta)^{1 - x_1} \cdots \theta^{x_n} (1 - \theta)^{1 - x_n} \\ &= \theta^{x_1 + \cdots + x_n} (1 - \theta)^{(1 - x_1) + \cdots + (1 - x_n)} = \theta^C (1 - \theta)^{n - C} \end{aligned}$$

$\log L(\theta) = C \log \theta + (n - C) \log(1 - \theta)$ を微分すると

$$\frac{d \log L(\theta)}{d\theta} = \frac{C}{\theta} - \frac{n - C}{1 - \theta} = \frac{C - n\theta}{\theta(1 - \theta)} \quad \text{なので} \quad \frac{d \log L(\theta)}{d\theta} = 0 \quad \text{を解くと} \quad \hat{\theta} = \frac{C}{n}$$

最尤法 (Jukes-Cantorモデル)

コイン投げと同じところ

塩基が異なる確率と時間の関係

塩基が異なる場合=1, 同じ場合=0 と書く

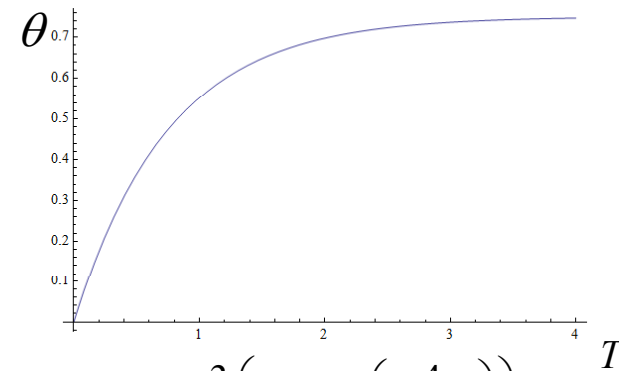
...ATGCCATG...
...ATACCTG...
...00100100...

DNA配列の長さ = n

変異数 = C

塩基が異なる確率 = θ

$$L(\theta) = \theta^C (1 - \theta)^{n-C}$$



$$\theta = \frac{3}{4} \left(1 - \exp\left(-\frac{4}{3}T\right) \right)$$

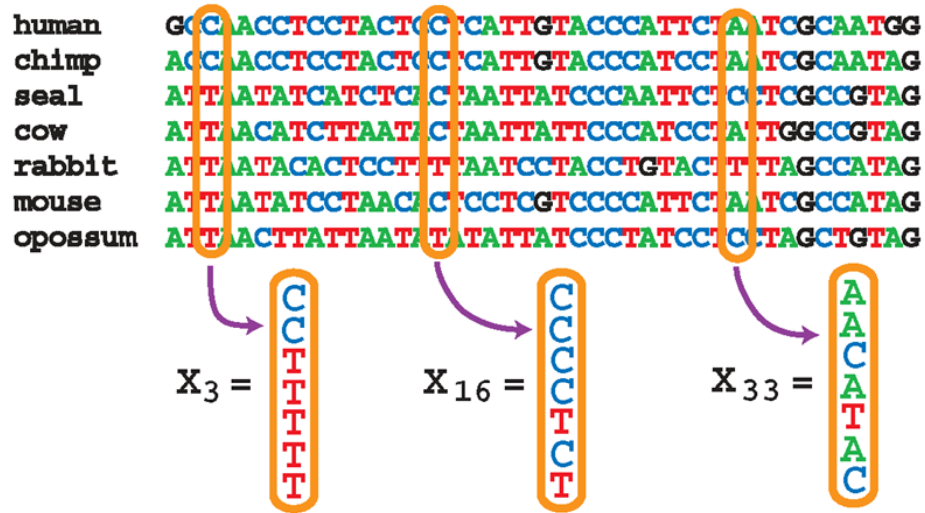
塩基毎の置換数に換
算した時間 = T

$$L_{JC}(T) = L\left(\frac{3}{4} \left(1 - \exp\left(-\frac{4}{3}T\right) \right)\right) \text{ を最大にする } T \rightarrow \hat{T} = -\frac{3}{4} \log\left(1 - \frac{4C}{3n}\right)$$

DNA配列間の置換数の推定値: $n\hat{T}$

($C \geq 0.75n$ なら $\hat{T} = \infty$)

最尤法 (系統樹)



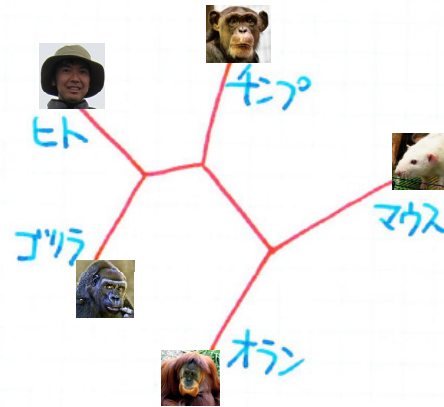
$$L(\text{樹形}, \text{枝長}) = f(x_1; \text{樹形}, \text{枝長}) \cdots f(x_n; \text{樹形}, \text{枝長})$$

$$\max_{\text{樹形}} \max_{\text{枝長}} L(\text{樹形}, \text{枝長}) \longrightarrow (\widehat{\text{樹形}}, \widehat{\text{枝長}}) \quad \text{最尤推定量}$$

樹形



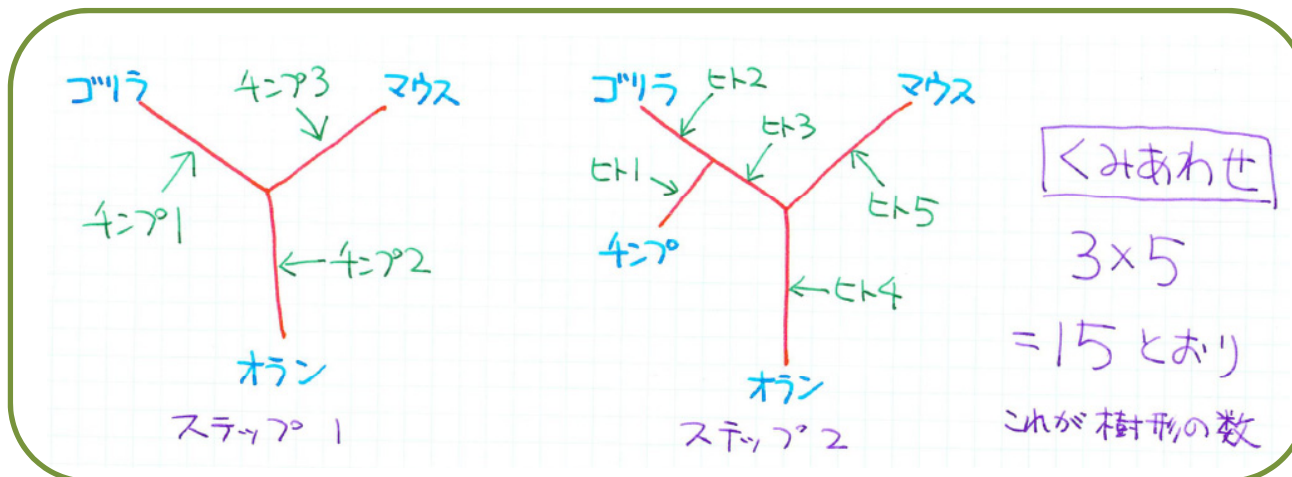
樹形 1



樹形 2



樹形 3



系統樹を数える

```
Table[{n,  $\frac{(2n-5)!}{(n-3)!2^{n-3}}$ }, {n, 3, 20}] // TableForm
```

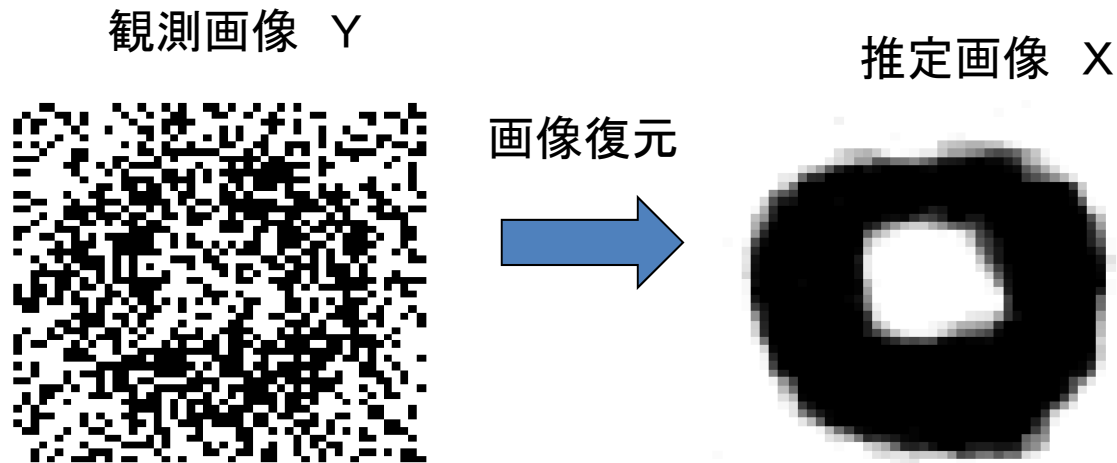
```
3 1
4 3
5 15
6 105
7 945
8 10395
9 135135
10 2027025
11 34459425
12 654729075
13 13749310575
14 316234143225
15 7905853580625
16 213458046676875
17 6190283353629375
18 191898783962510625
19 6332659870762850625
20 221643095476699771875
```

生物種= n のとき, 無根系統樹の樹形の数は


$$1 \times 3 \times 5 \times 7 \times \cdots \times (2n-5) = \frac{(2n-5)!}{(n-3)!2^{n-3}}$$

参考： Markov chain Monte Carlo法

情報科学科3年「データ解析」の課題： ランダムネス計算のアルゴリズムをプログラミング



ベイズ事後確率を計算 $P(X|Y) = \frac{P(Y|X) P(X)}{P(Y)}$



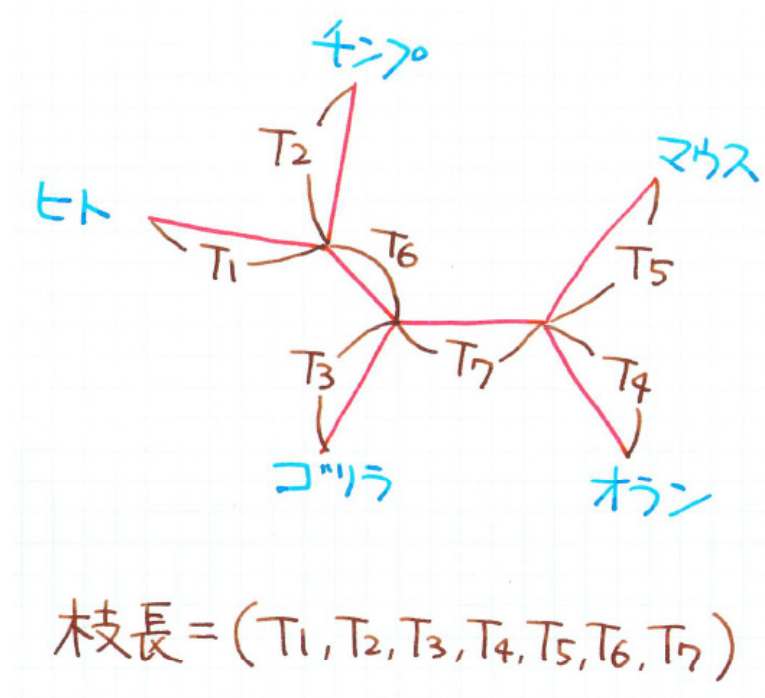
ベイズ
(1702-1761)

50x50=2500ピクセル

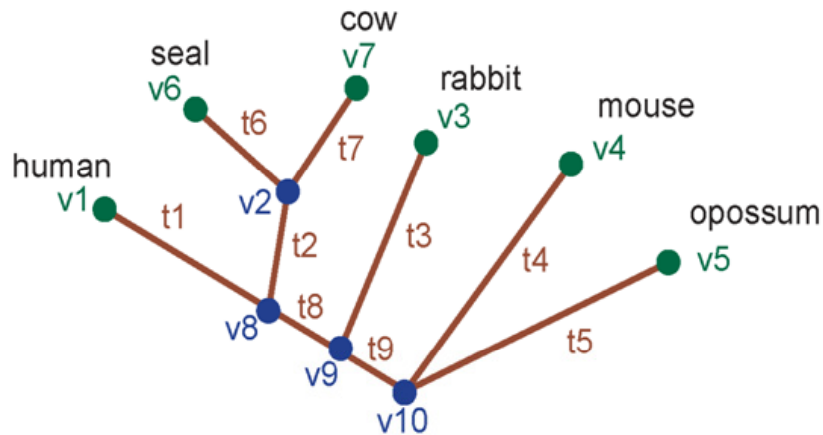
探索する画像の数 = $2^{2500} = 3 \times 10^{752}$

MCMC法(ギブス・サンプラー)

枝長



系統樹の確率モデル



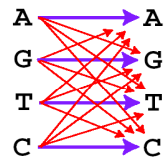
サイト h における配列パターン x_h の確率

$f(x_h; \text{樹形, 枝長}) =$

$$\sum_{v_2, v_8, v_9, v_{10}} P(v_1|v_8; t_1)P(v_2|v_8; t_2)P(v_3|v_9; t_3) \\ \times P(v_4|v_{10}; t_4)P(v_5|v_{10}; t_5)P(v_6|v_2; t_6) \\ \times P(v_7|v_2; t_7)P(v_8|v_9; t_8)P(v_9|v_{10}; t_9)P(v_{10})$$

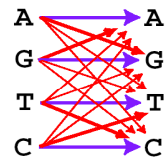
各枝でのマルコフ過程(つまり進化)は独立という仮定

- 塩基置換 (A, T, G, C の 4 種類)



JCモデル

(トランジションとトランスページョン, 組成比)



HKYモデル

- アミノ酸置換 (20 種類)

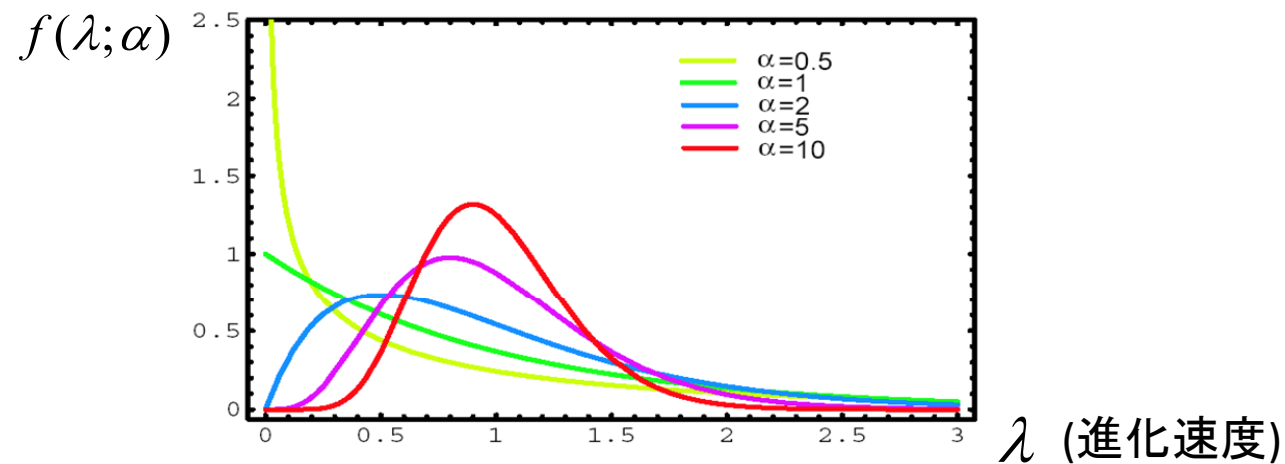
Dayhoffモデル, Jonesモデルなど

$$P(t) = \exp(tQ) = \sum_{k \geq 0} \frac{t^k}{k!} Q^k$$

進化速度の変動

DNAの座位毎に進化速度が異なる

➡ ガンマ分布の確率密度関数(平均=1)でモデル化



$$f(x_h; \text{樹形}, \text{枝長}) = \int_0^{\infty} f(x_h; \text{樹形}, \text{枝長} \times \lambda) f(\lambda; \alpha) d\lambda$$

樹形の尤度

$$L(\text{樹形}) = \max_{\text{枝長}} L(\text{樹形}, \text{枝長})$$

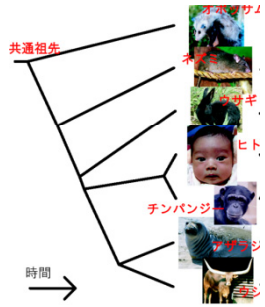
計算には対数尤度を使うことが多い

$$\log L(\text{樹形}) = \max_{\text{枝長}} \log L(\text{樹形}, \text{枝長})$$

$$\max_{\text{樹形}} L(\text{樹形}) \longrightarrow \widehat{\text{樹形}} \quad \text{最尤推定量}$$

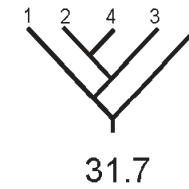
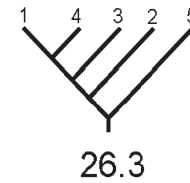
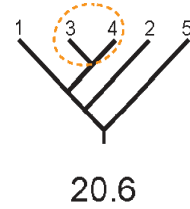
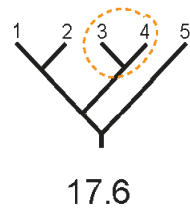
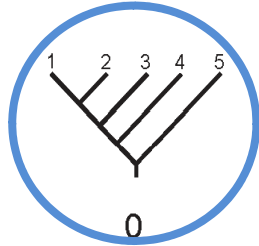
⇕ どちらで考えても同じこと

$$\max_{\text{樹形}} \max_{\text{枝長}} L(\text{樹形}, \text{枝長}) \longrightarrow (\widehat{\text{樹形}}, \widehat{\text{枝長}})$$

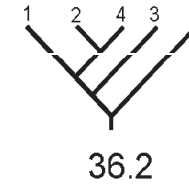
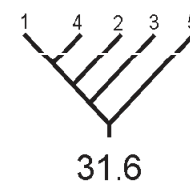
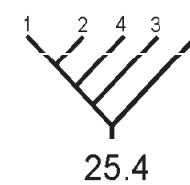
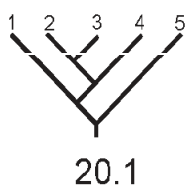
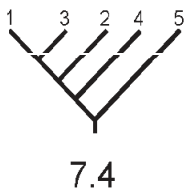
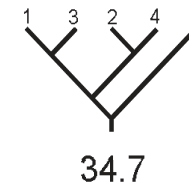
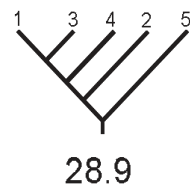
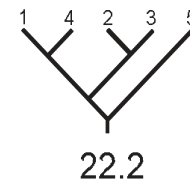
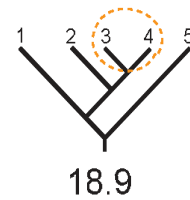
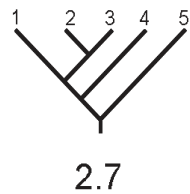


系統樹を尤度の順に並べる

1=human, 2=(seal, cow), 3=rabbit, 4=mouse, 5=opossum



樹形の最尤推定



数値は対数尤度の差 $\log L(\widehat{\text{樹形}}) - \log L(\text{樹形})$

最尤法の一般化：情報量規準

Akaike Information Criterion

$$AIC = -2 \log L(\theta) + 2 \dim \theta$$

第1項：対数尤度 = 確率の対数 = エントロピーに相当

第2項：パラメタ数 = 確率モデルの複雑さ



2006年 京都賞 基礎科学部門 / 数理科学
赤池 弘次 (Hirotugu Akaike)

統計数理研究所 名誉教授

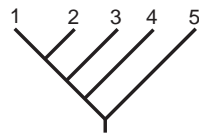


情報数理の基礎概念に基づく、実用性と汎用性の両方を兼ね備えた、統計モデル選択のための規準 Akaike Information Criterion (AIC) の提唱により、データの世界とモデルの世界を結びつける新しいパラダイムを打ち立て、情報・統計科学への多大な貢献をした

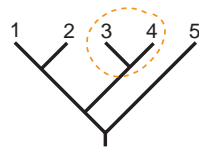
次の話題

系統樹のブートストラップ確率

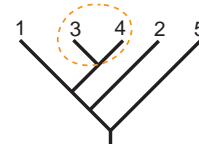
1=ヒト, 2=(アザラシ, ウシ), 3=ウサギ, 4=マウス, 5=オポッサム



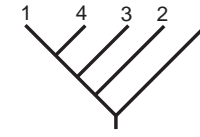
0.58



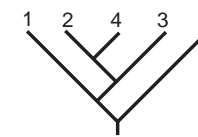
0.01



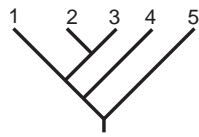
0.02



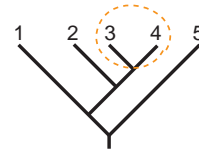
0.00



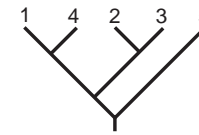
0.00



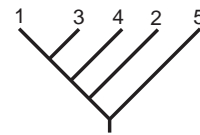
0.31



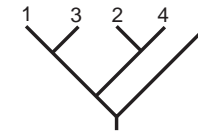
0.04



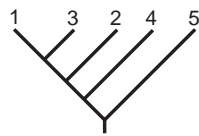
0.00



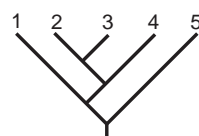
0.00



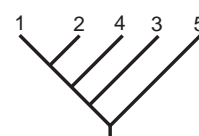
0.00



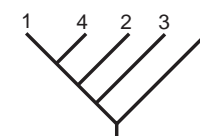
0.04



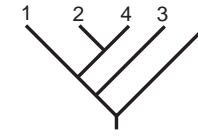
0.01



0.00



0.00

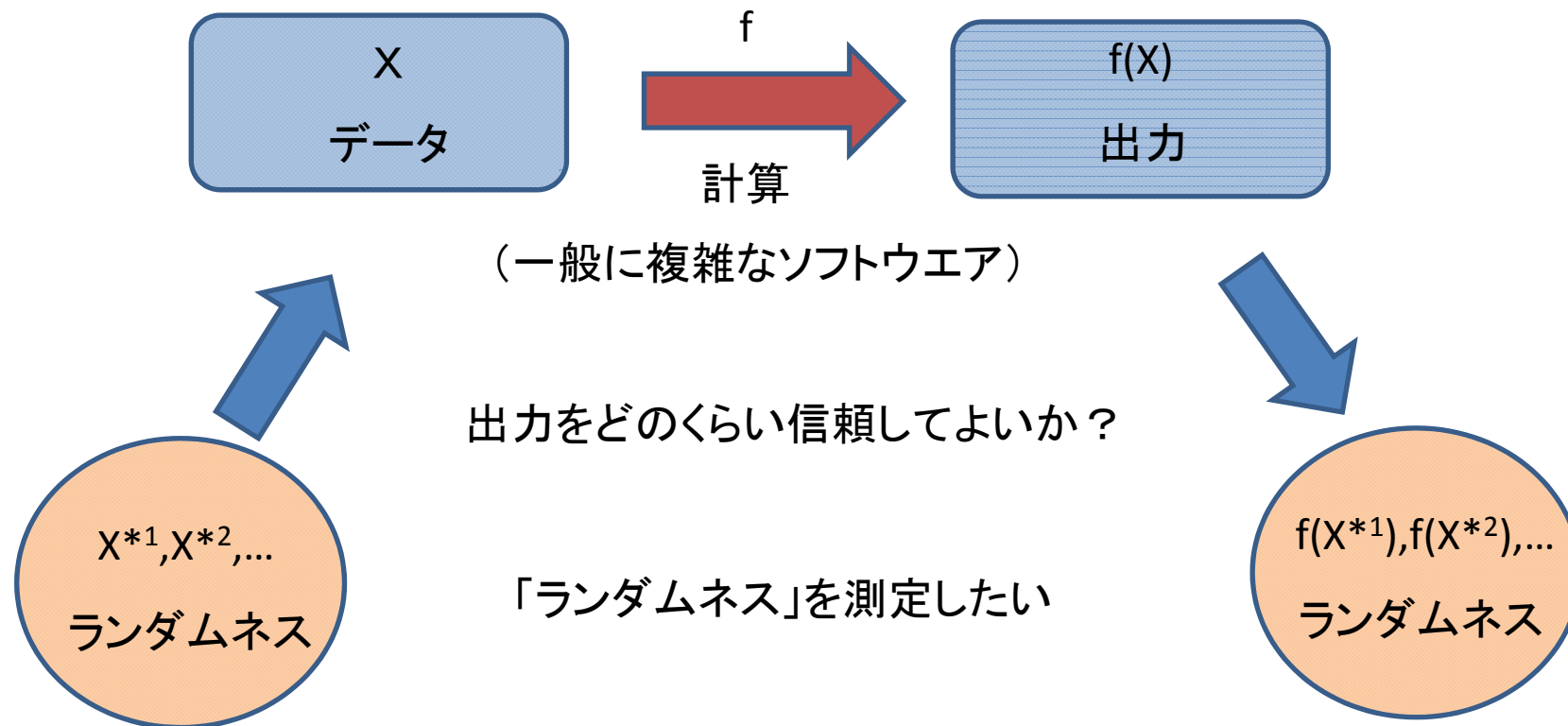


0.00

統計学： ブートストラップ法

ランダムネスの測定

Measuring the Randomness

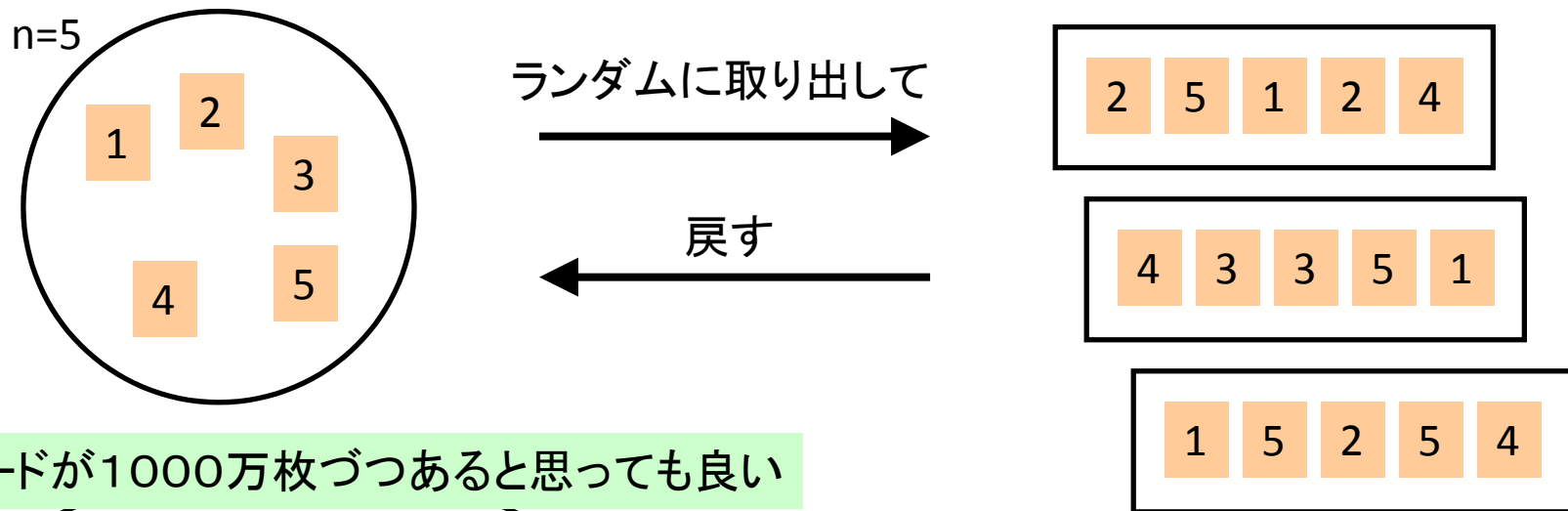


ランダムネスのものさし： 統計学の「確率値 (p-値)」または「信頼度」

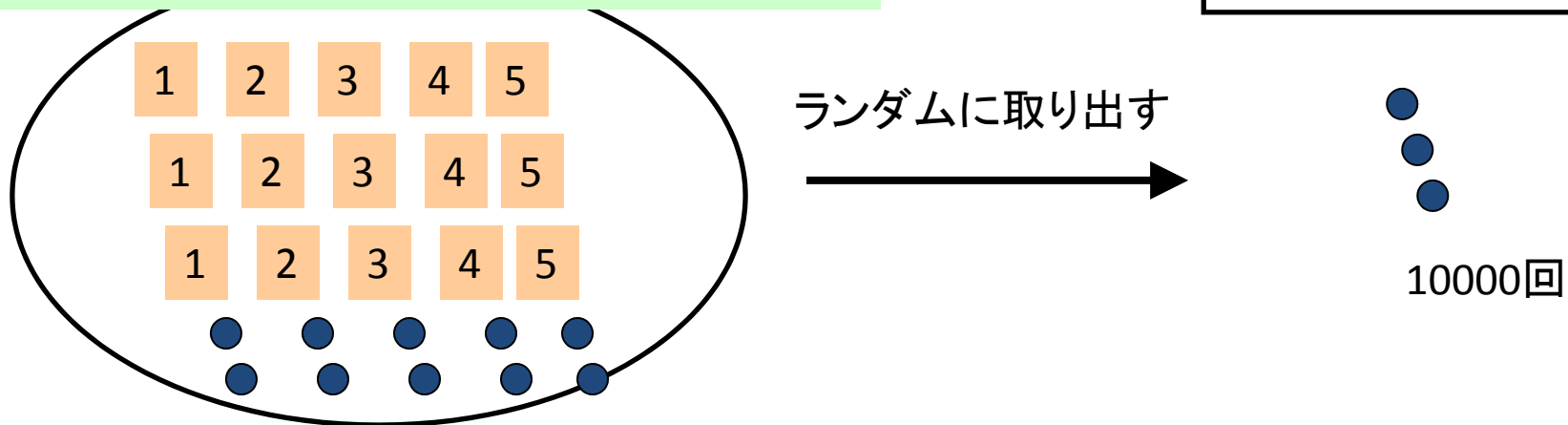
きれいに解が求まらない問題でも使いたい(面倒な式はイヤ)

復元抽出

ブートストラップ法 = データからの復元抽出 (resampling with replacement)



各カードが1000万枚ずつあると思っても良い



ブートストラップ法の応用場面

- コイン投げ
- 系統樹のブートストラップ確率
- 系統樹の枝に関するブートストラップ確率
- ネットワーク推定における枝のブートストラップ確率（例：遺伝子制御ネットワーク推定）
- 機械学習のバギング（例：手書き文字認識）
- 空間統計（例：酸性雨の分布）

ブートストラップ法 (コイン投げ)

シミュレーション

1回目
110000000000111000100010100111...
→ $\hat{\theta}^* = 0.418$

2回目
100100010000101111101010100101...
→ $\hat{\theta}^* = 0.458$



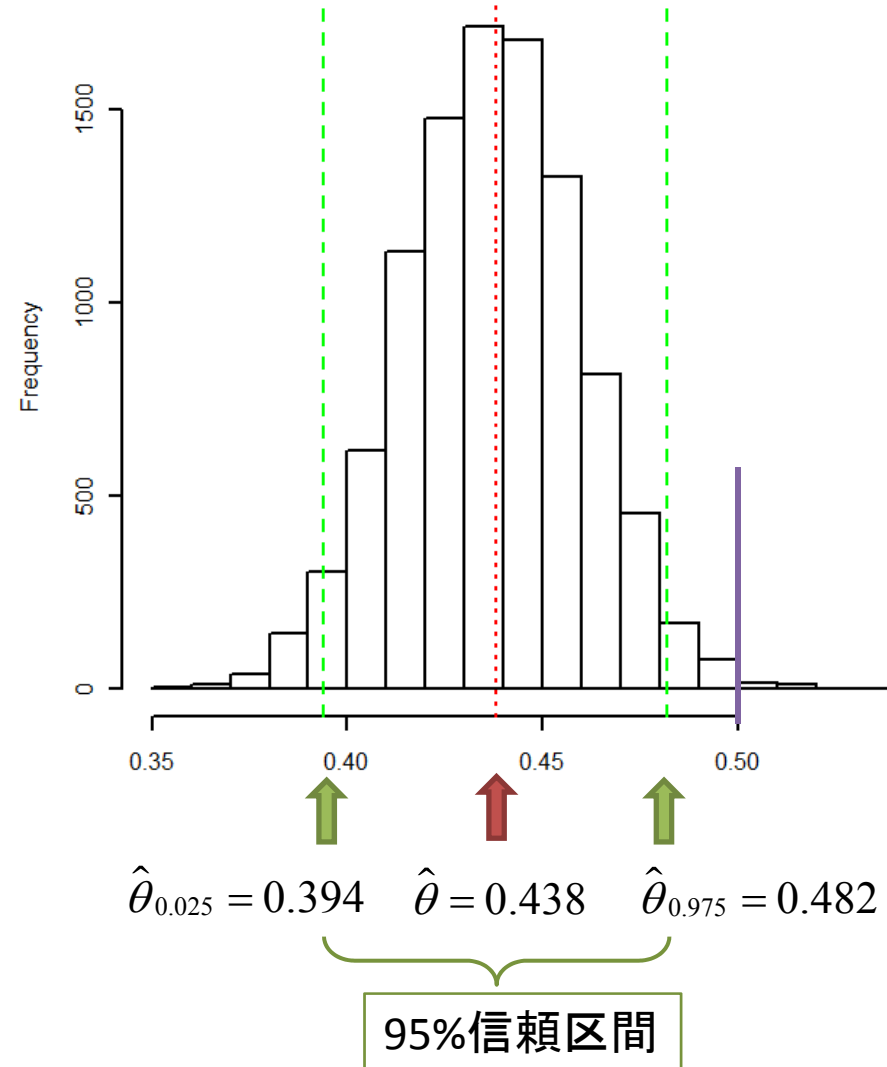
10000 回目
001011011011100011101011000111...
→ $\hat{\theta}^* = 0.420$

信頼区間

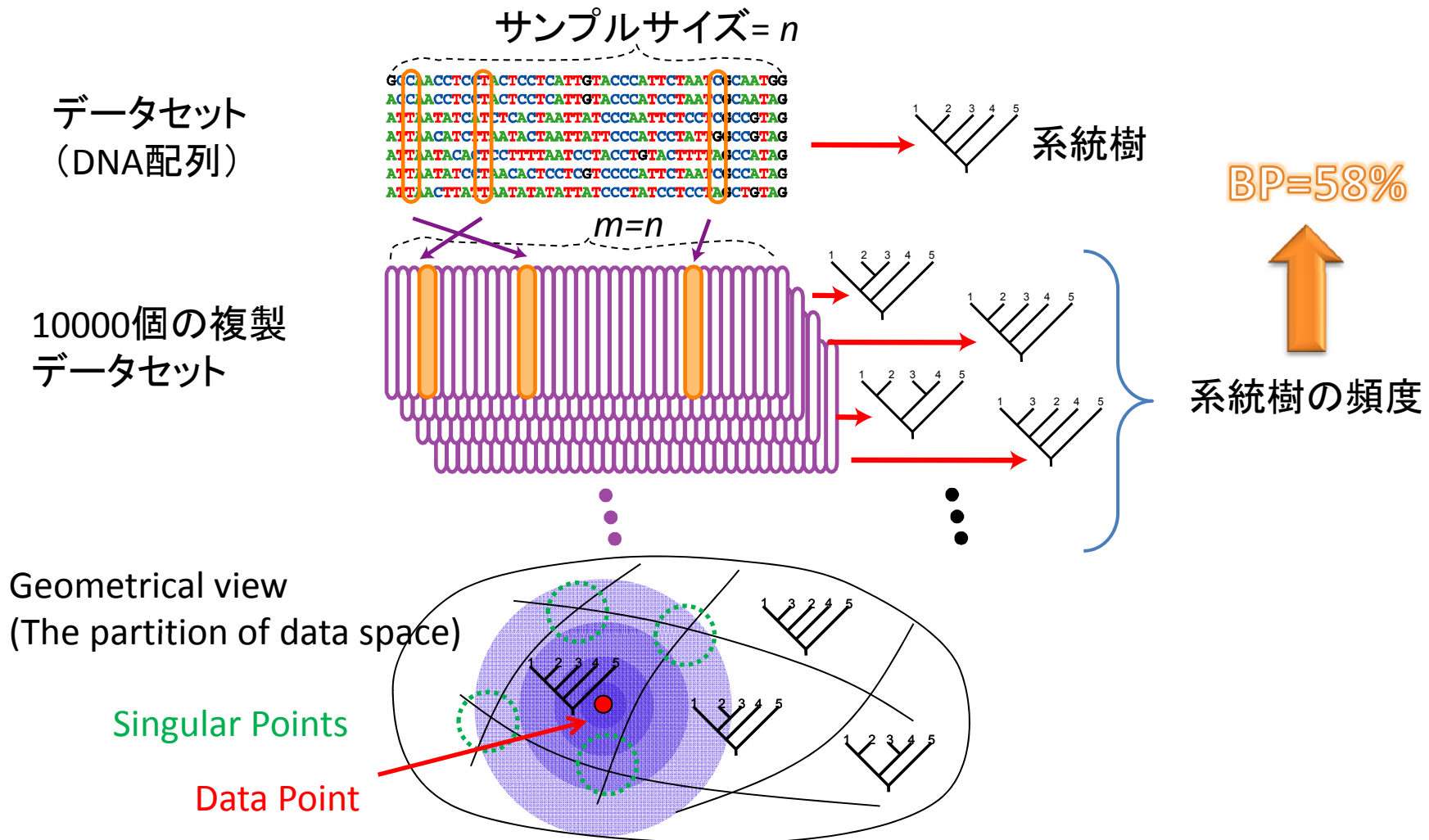
$\hat{\theta}^* = 0.418, 0.458, \dots, 0.420$
を大きくなる順にソート
小さい方から数えて2.5%と97.5%の値
 $\hat{\theta}_{0.025} = 0.394$ $\hat{\theta}_{0.975} = 0.482$

仮説検定

$\#\{\hat{\theta}^* > 0.5\} = 32$
帰無仮説 $\theta = 0.5$ の p -値は 0.0064



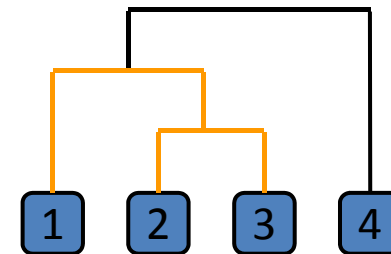
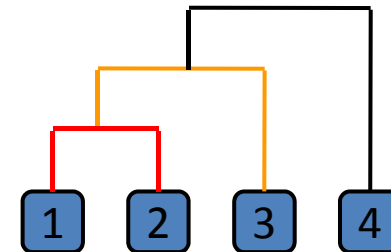
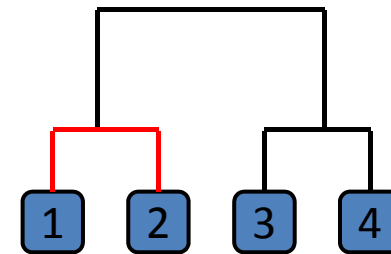
ブートストラップ確率 (Bootstrap Probability を略してBP)



Bootstrap probability = a Bayesian posterior probability (using flat prior dist.)

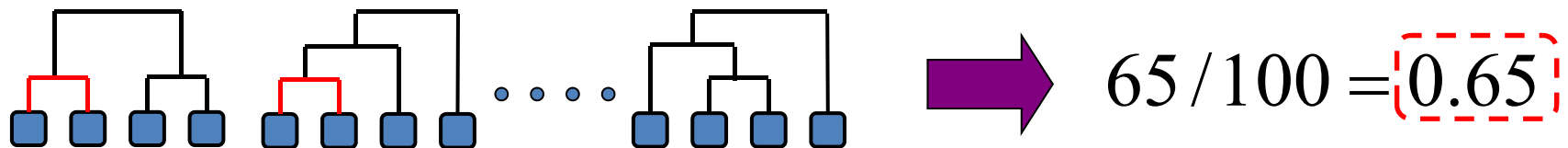
クラスター分析のバラツキ評価

- 樹形図に含まれる群に注目
- 右図の例
 - 2回出現した群
 - $\{1,2\}, \{1,2,3\}$
 - 1回出現した群
 - $\{3,4\}, \{2,3\}$
- 群の「出現頻度」を数える



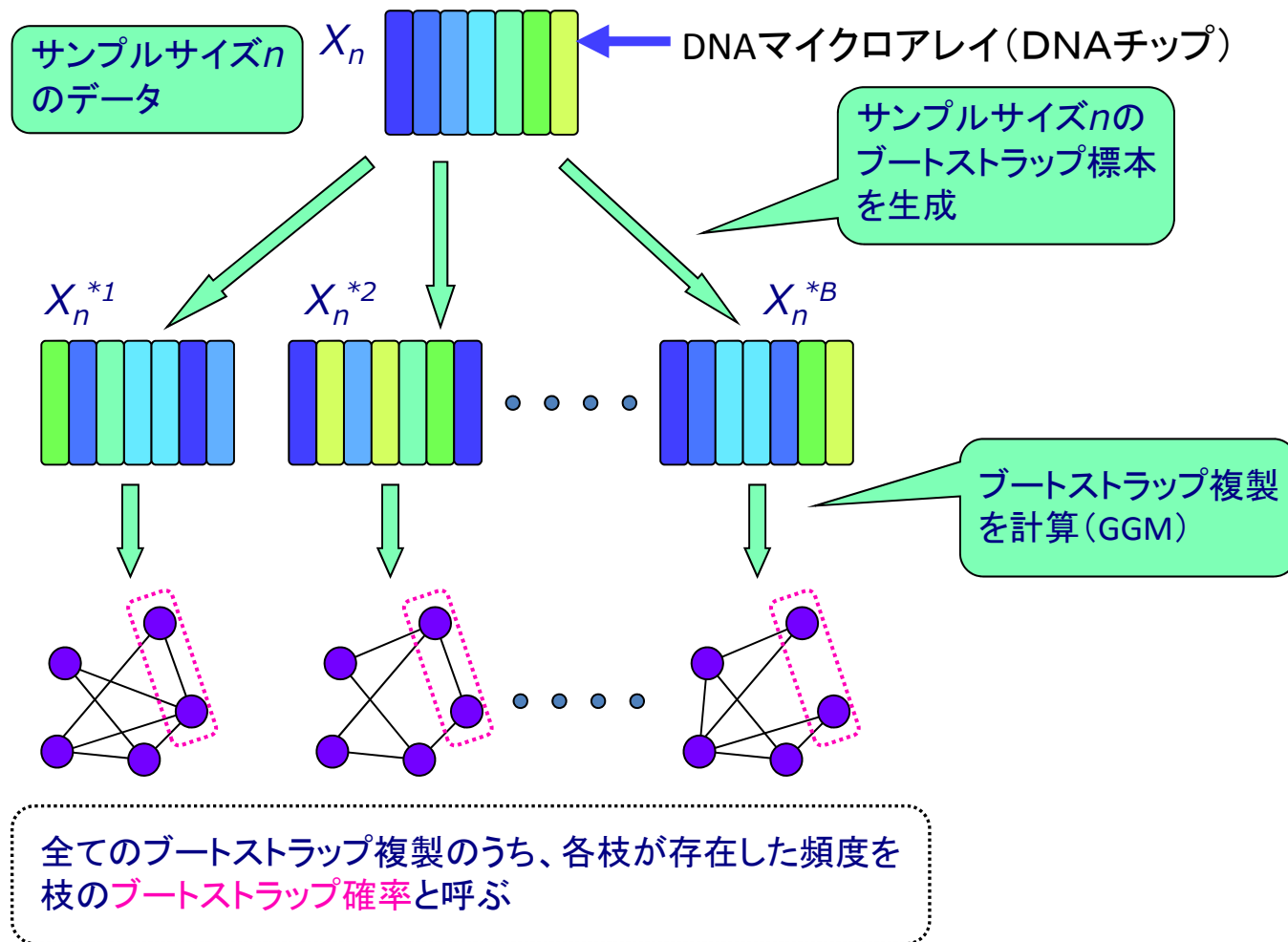
ブートストラップ確率

- ブートストラップ確率
 - モデルが仮説を満たした頻度
 - 例: 群が存在した頻度



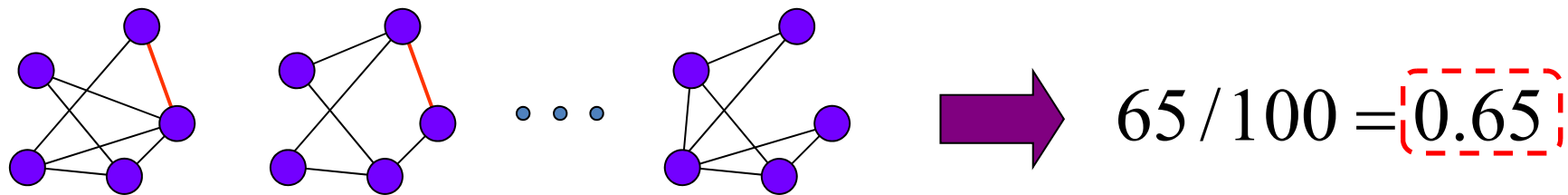
100回のブートストラップ中, ある群が65回存在

ネットワーク推定のランダムネス



ブートストラップ確率

- ブートストラップ確率
 - モデルが仮説を満たした頻度
 - 例: 枝が存在した頻度

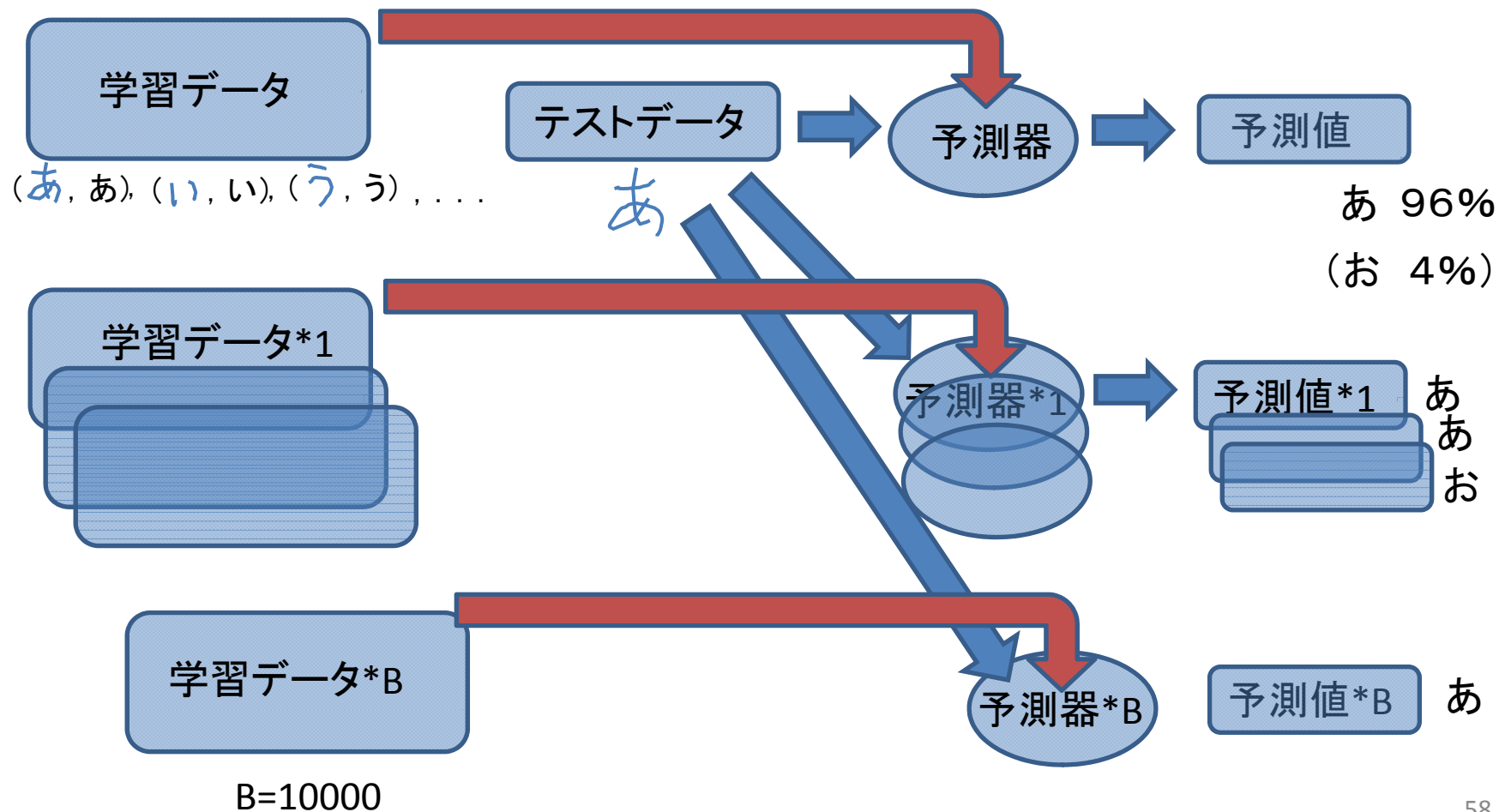


100回のブートストラップ中、ある枝が65回存在

バギング

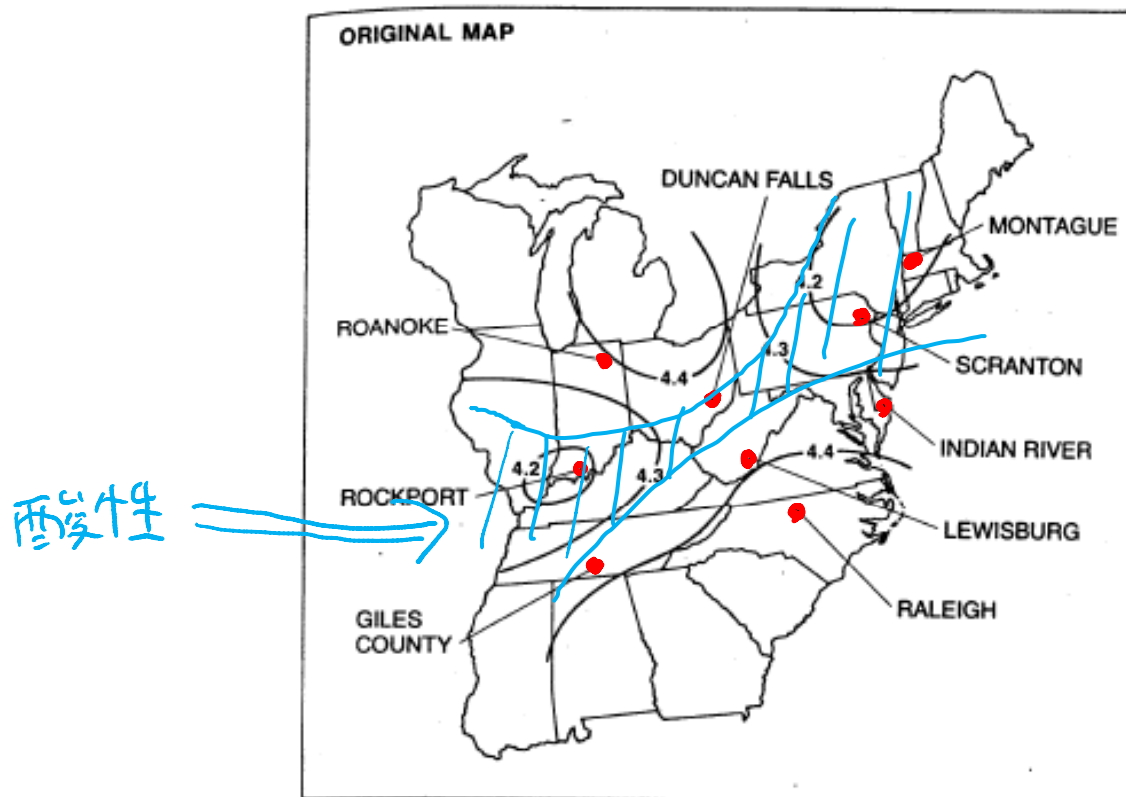
Bootstrap aggregating (bagging)

機械学習の判別問題 (例: 手書き文字認識)



酸性雨(pH)の分布

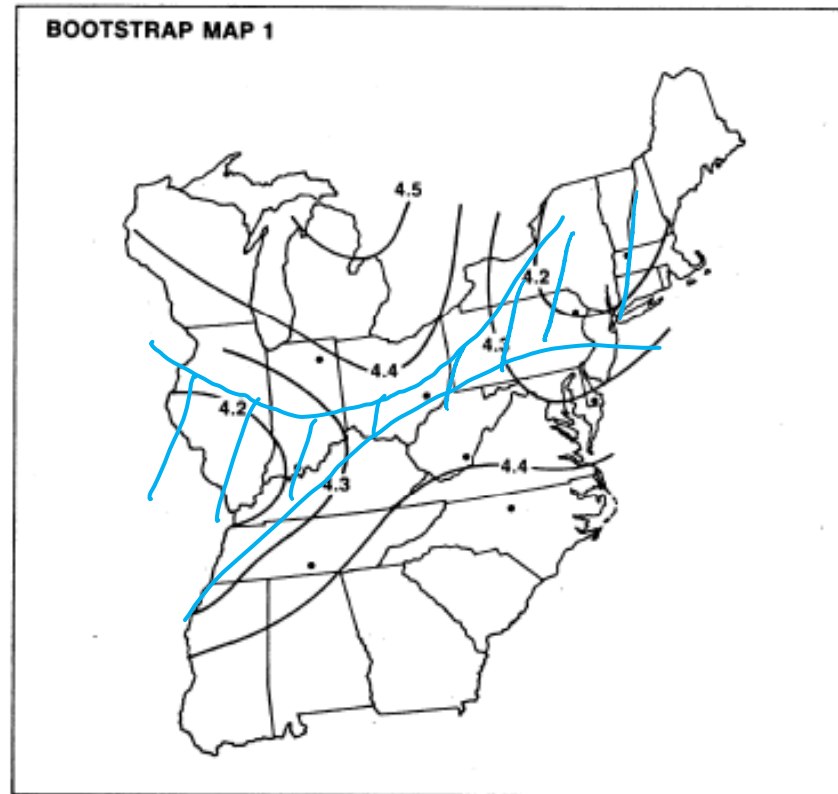
2年間(2000個)の測定値 @ 9カ所の測候所



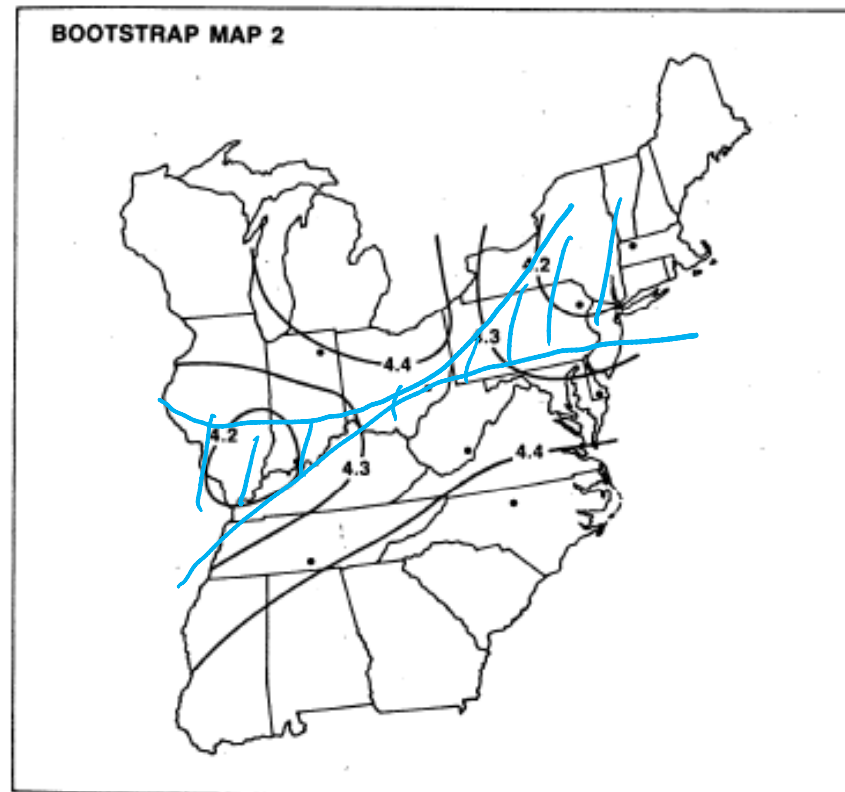
ダイアコニス, エフロン (1983)
「コンピュータがひらく新しい統計学」サイエンス (松原望 訳)

ブートストラップでデータセットを作って再分析する

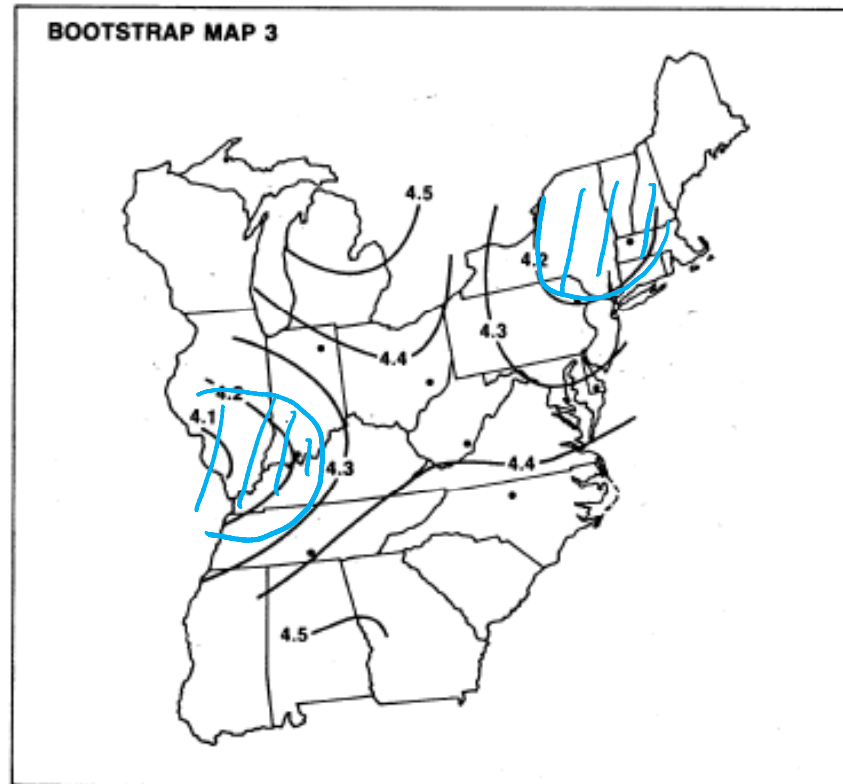
ブートストラップ標本(1)



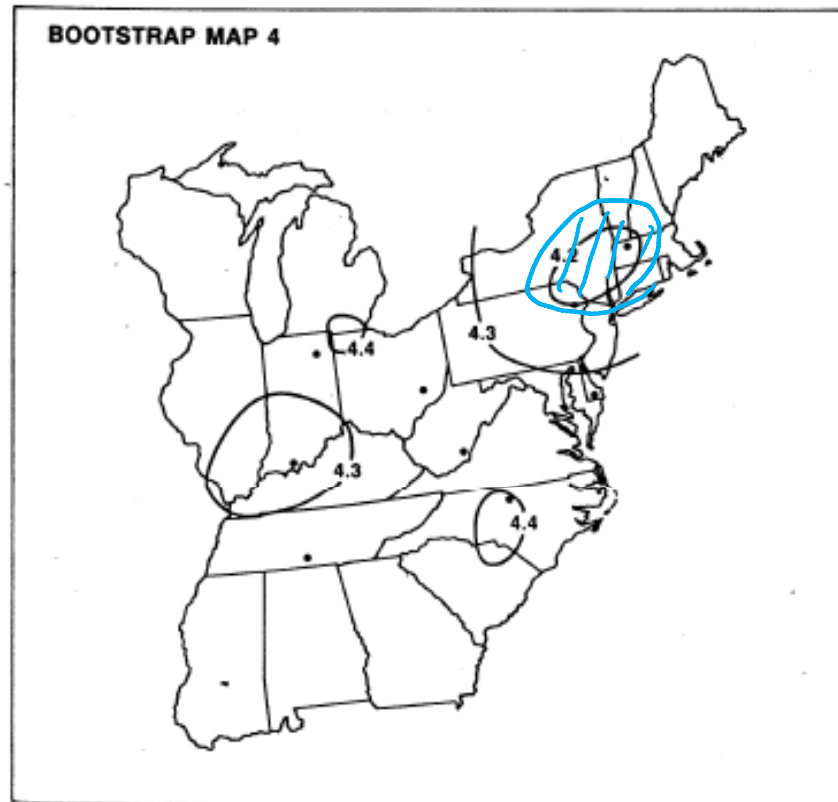
ブートストラップ標本(2)



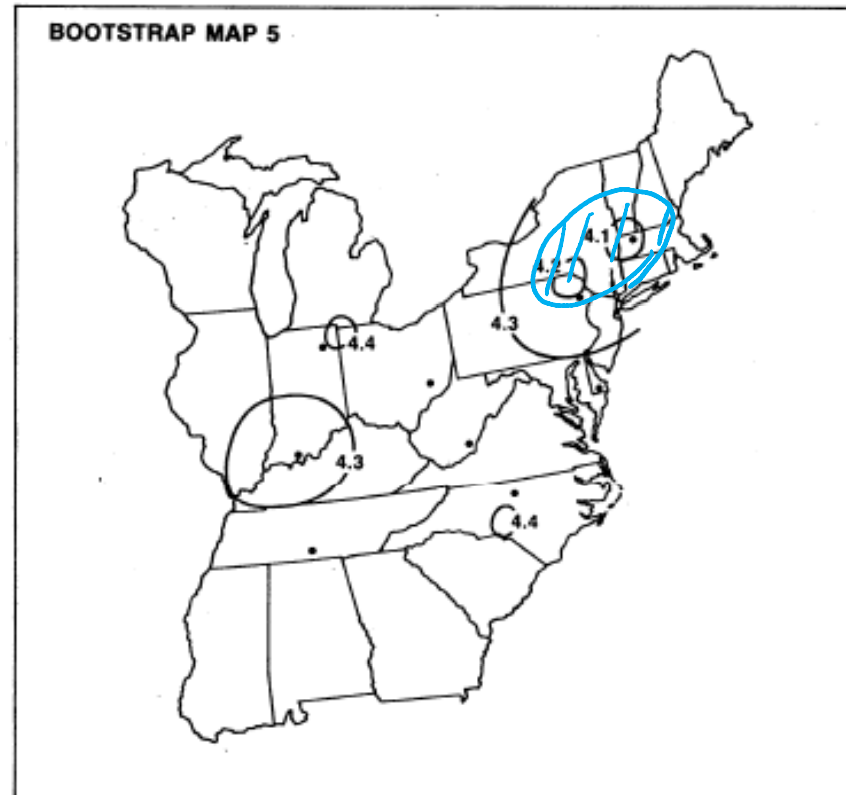
ブートストラップ標本(3)



ブートストラップ標本(4)



ブートストラップ標本(5)



系統樹のランダムネスを測る 3つの方法

- ブートストラップ確率
- Shimodaira-Hasegawa 検定
(多重比較法の一つ)
- マルチスケール・ブートストラップ
(サンプルサイズを $m=n$ にする)

新しい「発見」

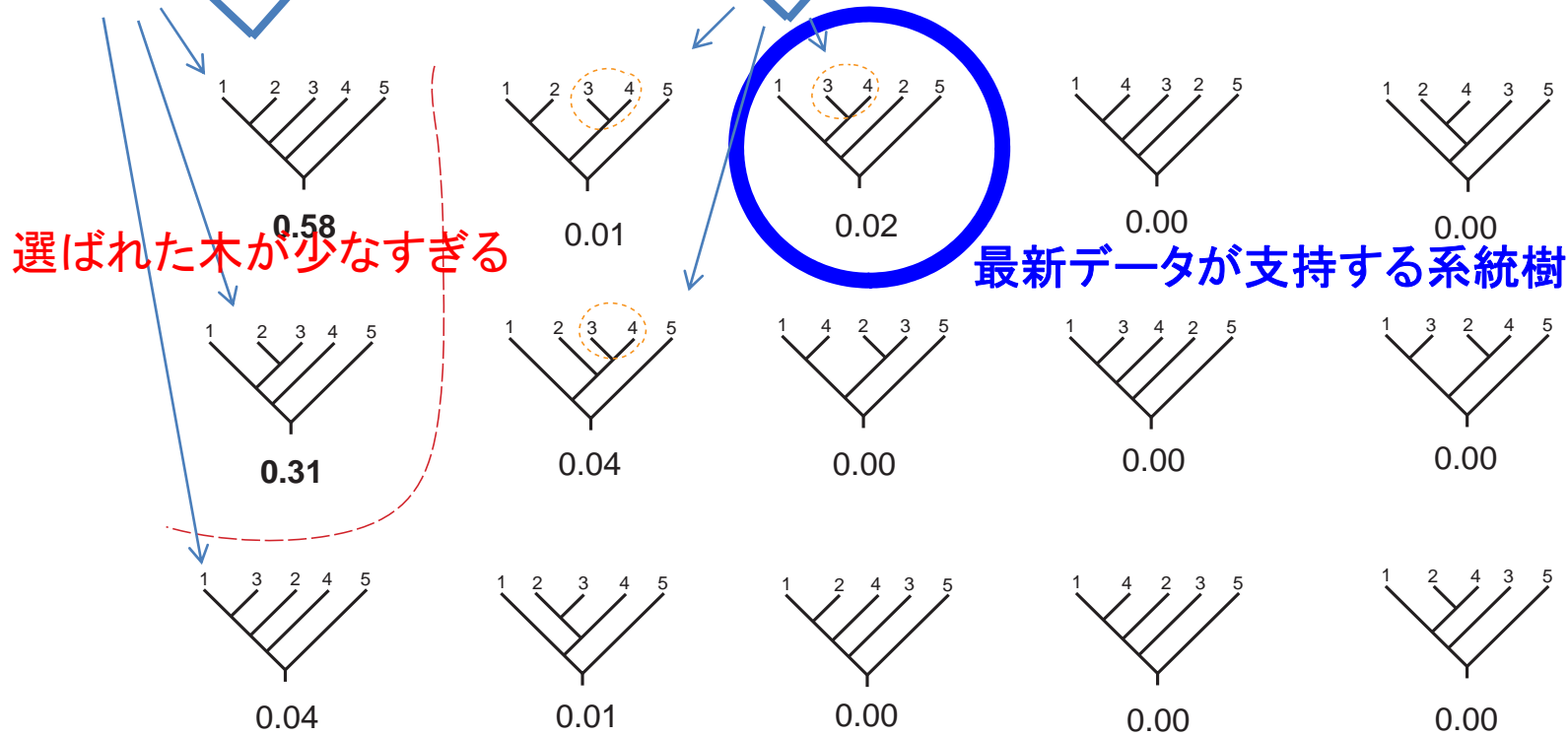


伝統的な仮説



ブートストラップ確率 (系統樹)

Efron (1979) *Ann. Stat.* Felsenstein (1985) *Evolution*



ブートストラップ確率はランダムネスを過小評価 ➡ 偽物の「発見」につながる (false positives)

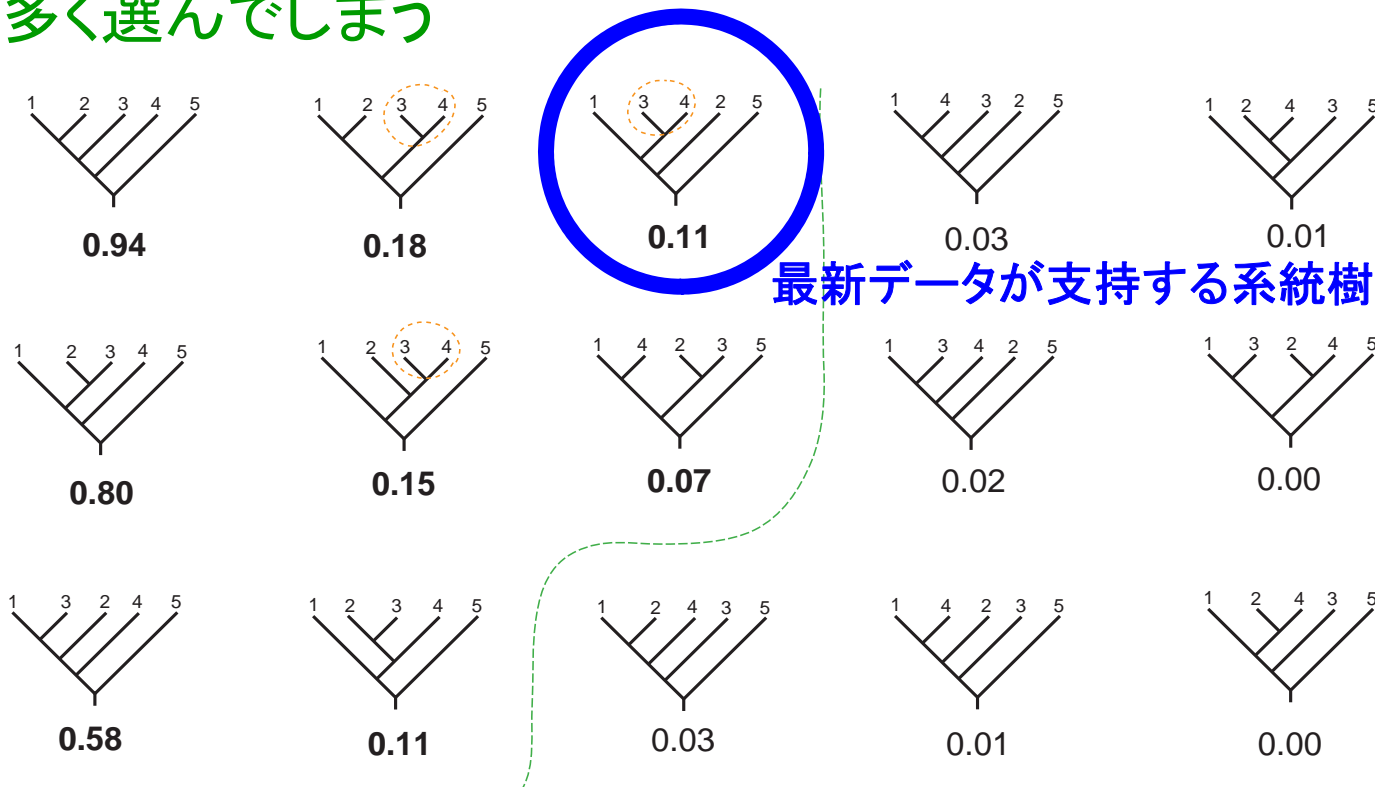
Mathematical reason: Hypothesis region is convex and the curvature is positive

Shimodaira-Hasegawa 検定 (多重比較法の一つ)

木を多く選んでしまう

Shimodaira (1998) *AISM*

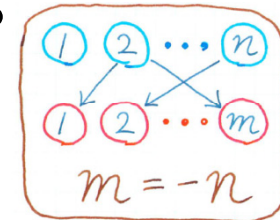
Shimodaira and Hasegawa (1999) *Mol. Biol. Evol.*



SH検定はランダムネスを過大評価 ➡ 本当の発見を見逃してしまう

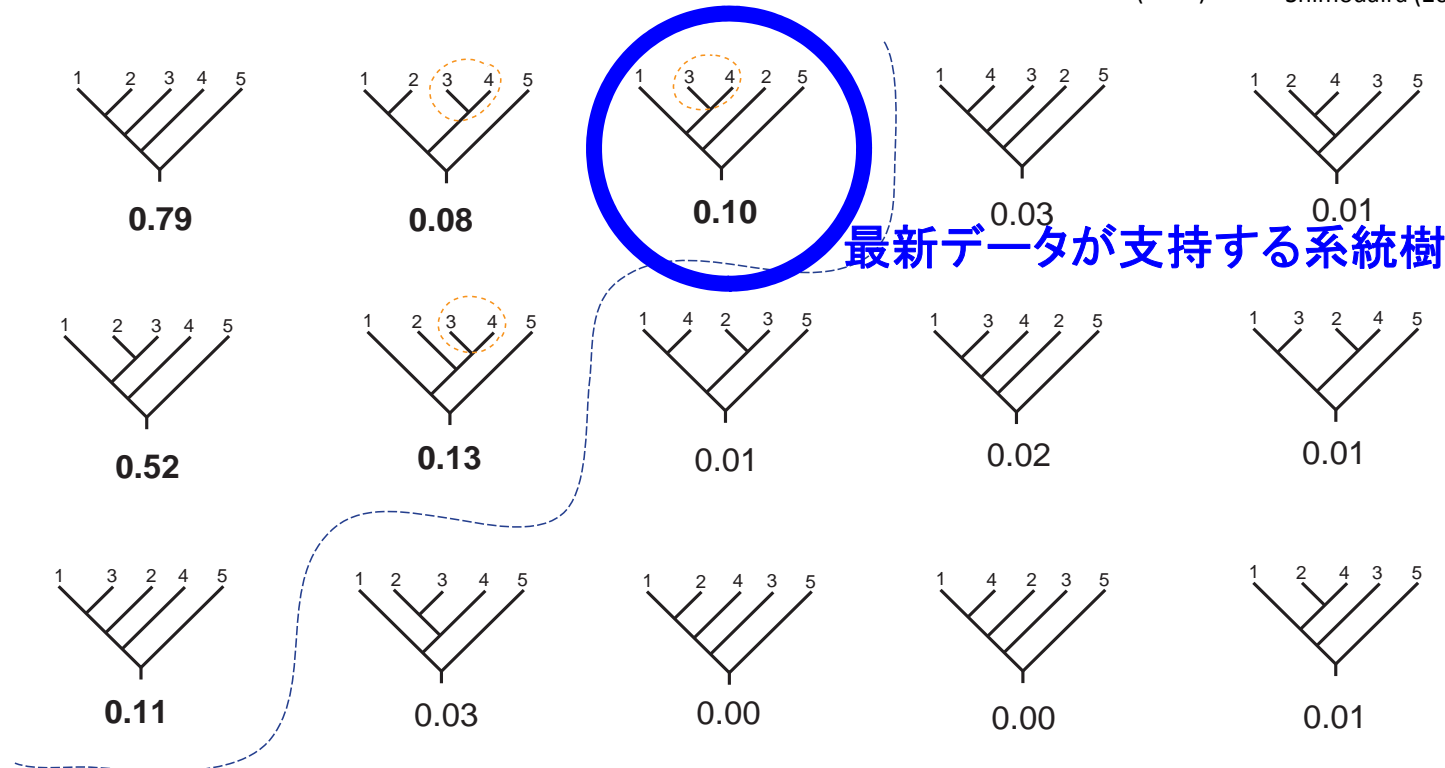
Mathematical reason: Multiple comparisons method evaluates the worst-case scenario and is conservative

マルチスケール・ブートストラップ (サンプルサイズを $m=n$ にする)



ちょうどいい個数の木が選ばれる

Shimodaira and Hasegawa (2001) *Bioinformatics*
 Shimodaira (2002) *Syst. Biol.* Shimodaira (2004) *Ann. Stat.*
 Shimodaira (2008) *JSPI* Shimodaira (2010) *AIMS*



マルチスケール・ブートストラップ法はランダムネスを適切に評価(バイアスゼロ)

➡ ランダムネスは消えないが, 無知の程度を正しく知ることができる.

32種の哺乳類の場合

AU: マルチスケール・ブートストラップ法

Table 2. p -values for the fifteen constrained candidate tree topologies.

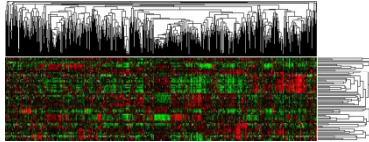
model	$\Delta\ell$	PP1	PP2	BP	AU	KH	SH	WSH	tree topology
T_1	0.0	0.28	0.61	0.23	0.69	0.55	0.97	0.95	((G1,G2),(G3,G4),G5)
T_2	1.5	0.49	0.14	0.28	0.60	0.46	0.83	0.86	((G1,(G2,G3)),G4,G5)
T_3	1.7	0.15	0.12	0.16	0.47	0.41	0.84	0.84	((G1,G2,G3),G4,G5)
T_4	1.9	0.06	0.09	0.13	0.45	0.33	0.84	0.81	(G1,(G2,(G3,G4)),G5)
T_5	2.6	0.01	0.04	0.09	0.37	0.27	0.80	0.73	((G1,(G3,G4)),G2,G5)
T_6	6.2	0.00	0.00	0.02	0.16	0.15	0.64	0.54	((G1,G2,G4),G3,G5)
T_7	6.8	0.00	0.00	0.03	0.25	0.28	0.58	0.61	((G1,G4),(G2,G3),G5)
T_8	8.3	0.00	0.00	0.01	0.08	0.23	0.51	0.40	(G1,((G2,G3),G4),G5)
T_9	8.7	0.00	0.00	0.04	0.25	0.21	0.50	0.66	((G1,G4),G2),G3,G5)
T_{10}	9.9	0.00	0.00	0.02	0.14	0.18	0.43	0.59	((G1,G4),G3),G2,G5)
T_{11}	12.7	0.00	0.00	0.00	0.00	0.10	0.29	0.20	((G1,G3),G2),G4,G5)
T_{12}	15.9	0.00	0.00	0.00	0.01	0.05	0.17	0.27	((G1,(G2,G4)),G3,G5)
T_{13}	18.6	0.00	0.00	0.00	0.00	0.03	0.09	0.13	(G1,((G2,G4),G3),G5)
T_{14}	18.8	0.00	0.00	0.00	0.00	0.02	0.09	0.09	((G1,G3),G4),G2,G5)
T_{15}	21.5	0.00	0.00	0.00	0.00	0.01	0.04	0.10	((G1,G3),(G2,G4),G5)

Recent Results

← ③ } Nishihara et al (2008)
 ← ② }
 ← ① } Genome-scale analysis of 1 Mbp

Note: Only the fifteen candidate tree topologies are considered; the subtree topologies for G1, ..., G5 are specified in Fig. 1. $\Delta\ell$ denotes the log-likelihood difference from the ML topology. The trees are numbered by increasing order of $\Delta\ell$. PP1 denotes the PP calculated by the MCMCMC using MrBayes with clade constraints, and PP2 denotes the PP calculated by the BIC approximation. p -values ≥ 0.05 are in boldface.

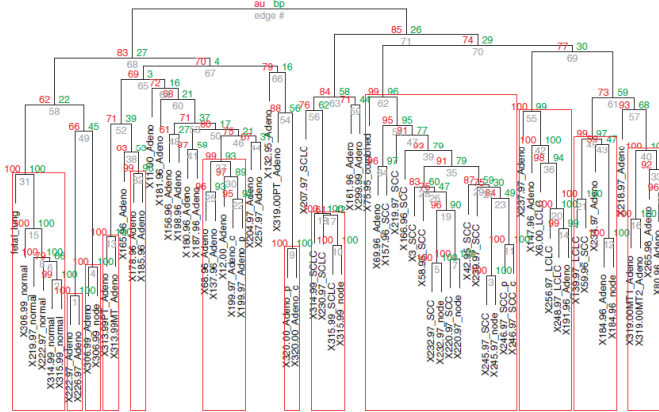
マルチスケール・ブートストラップの応用



遺伝子発現データ
(DNAチップ)

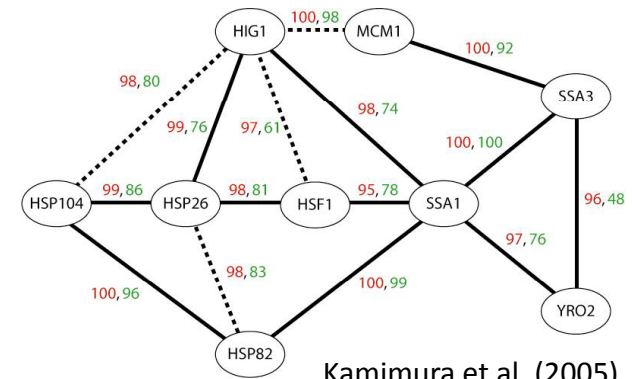


癌のタイプ分類



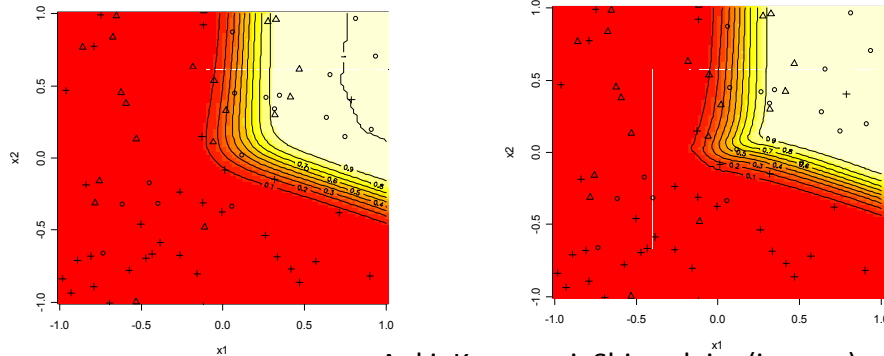
Suzuki and Shimodaira (2006)

遺伝子制御ネットワーク



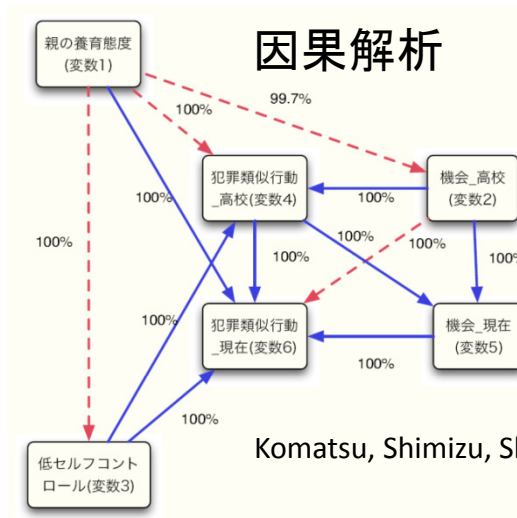
Kamimura et al. (2005)

機械学習



Aoki, Kanamori, Shimodaira (in prep)

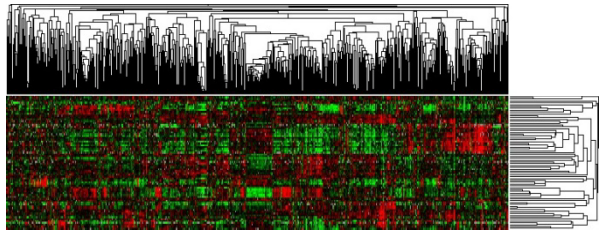
因果解析



Komatsu, Shimizu, Shimodaira (2010)

並列計算でスピードアップ (下平研究室 2008/04/10)

問題: DNAチップから癌の診断



東工大GSICのスパコン
TSUBAME



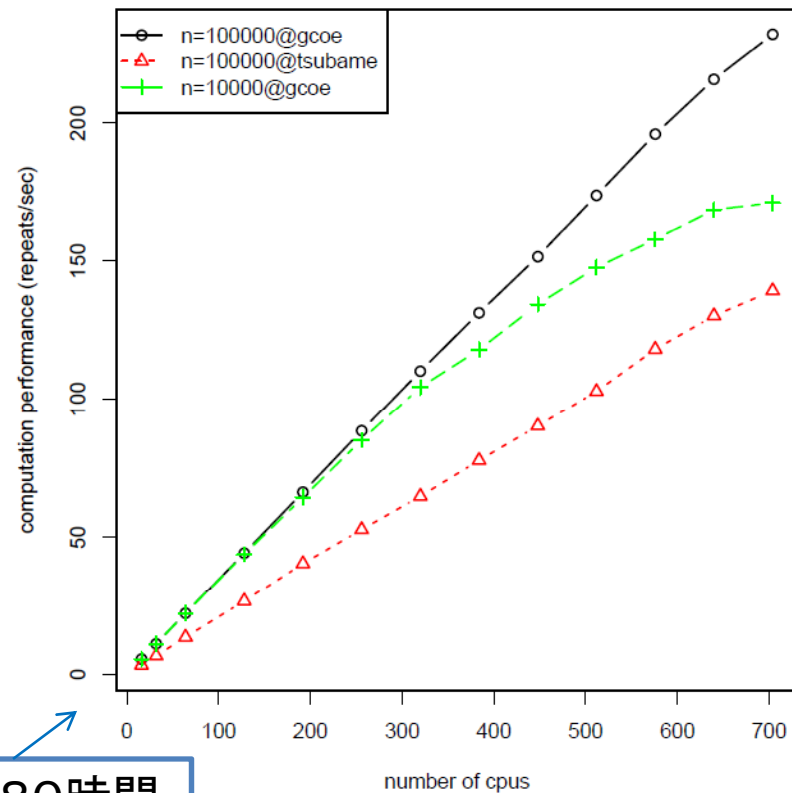
10000コア

毎秒あたりの計算量

1個 = 80時間

700個 = 7分

bootstrap computation (pvclust)



計算プロセッサ(コア)の個数

Data and text mining

Pvclust: an R package for assessing the uncertainty in hierarchical clustering

Ryota Suzuki^{1,2,*} and Hidetoshi Shimodaira¹

¹Department of Mathematical and Computing Sciences, Tokyo Institute of Technology, 2-12-1 Ookayama, Meguro-ku, Tokyo 152-8552, Japan and ²E-prime, Inc., 2-17-5 Nihonbashi-Kayabacho, Chuo-ku, Tokyo 103-0025, Japan

Received on August 30, 2005; revised on March 6, 2006; accepted on March 25, 2006
 Advance Access published April 4, 2006
 Associate Editor: Satru Vajnsu

Software developed by a master course student at Shimodaira-lab

“pvclust” already in use

下平研の学生が作成したソフトウェアが、すぐに他の研究で使われた
 ES細胞の研究(話題の iPS細胞に関係する)

PNAS January 24, 2006

ES cells derived from cloned...

ES cells derived from cloned and fertilized blastocysts are transcriptionally and functionally indistinguishable
 Tobias Brambrink¹, Konrad Hochdinger¹, George Ball², and Rudolf Jaenich^{1*}

Abstract
 Reproductive cloning is universally rejected as a valid technology in human because of the severely abnormal phenotypes seen in cloned animals. Gene expression abnormalities observed in tissues of cloned animals have also raised concerns regarding the therapeutic application of “somatic” embryonic stem (ES) cells derived by nuclear transplantation (NT) from a somatic somatic cell. Although previous experiments in mice have demonstrated that the developmental potential of ES cells derived from cloned blastocysts are as competent as those of ES cells derived from fertilized blastocysts, a systematic molecular characterization of NT ES cell lines is lacking. To investigate whether transcriptional alterations similar to those observed in tissues of cloned mice also occur in ES cells, we have compared transcriptional profiles of mouse NT- and fertilization-derived ES cell lines. We report here that ES cells from cloned and fertilized blastocysts are indistinguishable based on their transcriptional profiles, consistent with their normal developmental potential. Our results indicate that, in contrast to embryonic and fetal development of mice, the process of NT ES cell derivation rigorously selects for those somatic cells that have erased the “epigenetic memory” of the donor nucleus and, thus, become functionally equivalent. Our findings support the notion that ES cells derived from cloned or fertilized blastocysts have an identical therapeutic potential.

Introduction
 Nuclear transfer allows for the derivation of genetically identical ES cell lines from somatic cells of cloned mice, including a host of cell types (1–4). The feasibility of this approach, sometimes referred to as somatic cell nuclear transfer, has been demonstrated in animal models (5) and the clinical application of human nuclear transplantation (NT) ES cells represents an attractive prospect for the treatment of various medical conditions (6).

In animals, however, alterations in gene expression patterns, like the failure to activate genes essential for early embryonic development or to silence genes that are specific for the somatic donor cell type, affect the capacity of postimplantation NT embryos, resulting in a high frequency of embryonic and fetal lethality and widespread gene dysregulation in clones at both (4, 7–10). Furthermore, cloning often results in abnormal phenotypic and transcriptional abnormalities (11–14). It has been shown that the differentiation state of the donor nucleus can influence the abnormal gene expression patterns in newborn clones, suggesting retention of the “epigenetic memory” of the somatic donor genome after NT throughout embryonic and fetal development (5, 11–14).

The persistence of gene expression abnormalities in somatic tissues of cloned animals has raised the question whether ES cells derived by NT, in contrast to fertilization-derived ES cells, are prone to carry epigenetic alterations during transcriptional changes in response to tissue-specific gene-level promoter constructs. However,

Gene Expression Profiles (DNA microarray)

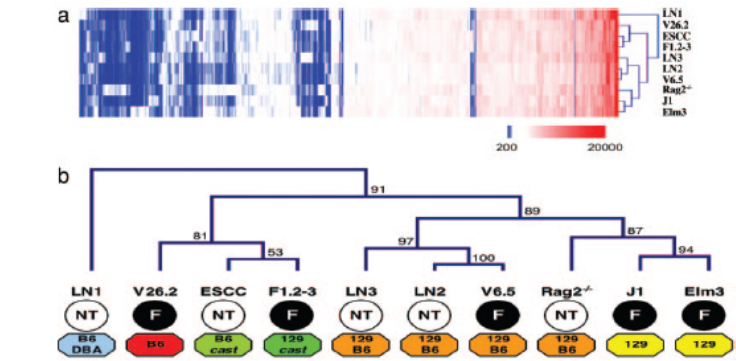


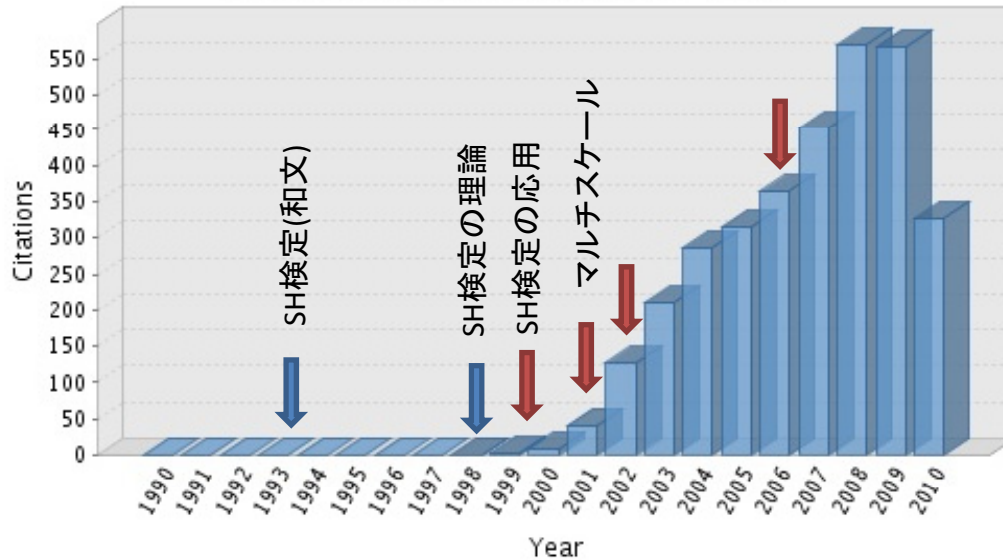
Fig. 2. Hierarchical clustering of individual ESC line expression profiles. (a) Heat map of clustering results (blue, no or very low expression; white, low expression; red, high expression). (b) Sample tree obtained from hierarchical clustering. ES cell line expression profiles cluster by genetic background (colored octagons) rather than by type of donor blastocyst (NT, cloned; F, fertilized; numbers next to nodes display multiscale bootstrap resampling probability based on 10,000 replications).

Brambrink et al.

論文被引用回数

<http://www.researcherid.com/citation/B-9127-2008>

Citation Distribution by year



- Title: [Multiple comparisons of log-likelihoods with applications to phylogenetic inference](#)
Author(s): SHIMODAIRA, H; HASEGAWA, M
Source: MOLECULAR BIOLOGY AND EVOLUTION Volume: 16 Issue: 8 Pages: 1114-1116 Published: AUG 1999
Times Cited: 1735
- Title: [An approximately unbiased test of phylogenetic tree selection](#)
Author(s): SHIMODAIRA, H
Source: SYSTEMATIC BIOLOGY Volume: 51 Issue: 3 Pages: 492-508 Published: MAY-JUN 2002
Times Cited: 505
DOI: 10.1080/10635150290069913
- Title: [CONSEL: for assessing the confidence of phylogenetic tree selection](#)
Author(s): SHIMODAIRA, H; HASEGAWA, M
Source: BIOINFORMATICS Volume: 17 Issue: 12 Pages: 1246-1247 Published: DEC 2001
Times Cited: 500
- Title: [Mitochondrial genome variation in Eastern Asia and the peopling of Japan](#)
Author(s): TANAKA, M; CABRERA, VM; GONZALEZ, AM; et al.
Source: GENOME RESEARCH Volume: 14 Issue: 10A Pages: 1832-1850 Published: OCT 2004
Times Cited: 156
DOI: 10.1101/gr.2286304
- Title: [Pyclust: an R package for assessing the uncertainty in hierarchical clustering](#)
Author(s): SUZUKI, R; SHIMODAIRA, H
Source: BIOINFORMATICS Volume: 22 Issue: 12 Pages: 1540-1542 Published: JUN 15 2006
Times Cited: 72
DOI: 10.1093/bioinformatics/btl117

Total Articles in Publication List: 19
Articles With Citation Data: 19
Sum of the Times Cited: 3267
Average Citations per Article: 171.95
h-index: 12
Last Updated: 06/21/2010 12:54 GMT

合計回数

Letter to the Editor
Multiple Comparisons of Log-Likelihoods with Applications to Phylogenetic Inference
SHIMODAIRA, H; HASEGAWA, M
Source: MOLECULAR BIOLOGY AND EVOLUTION
Volume: 16 Issue: 8 Pages: 1114-1116 Published: AUG 1999
Times Cited: 1735

An Approximately Unbiased Test of Phylogenetic Tree Selection
SHIMODAIRA, H
Source: SYSTEMATIC BIOLOGY
Volume: 51 Issue: 3 Pages: 492-508 Published: MAY-JUN 2002
Times Cited: 505
DOI: 10.1080/10635150290069913

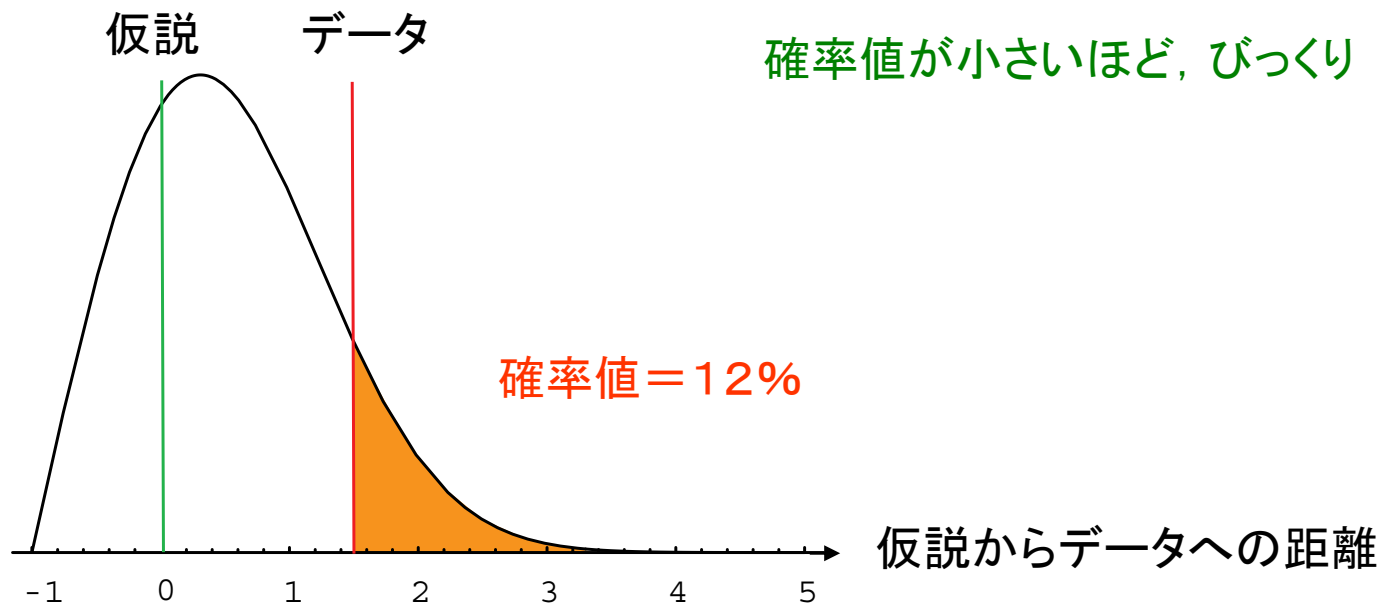
BIOINFORMATICS APPLICATIONS NOTE
CONSEL: for assessing the confidence of phylogenetic tree selection
SHIMODAIRA, H; HASEGAWA, M
Source: BIOINFORMATICS
Volume: 17 Issue: 12 Pages: 1246-1247 Published: DEC 2001
Times Cited: 500
DOI: 10.1093/bioinformatics/btl117

マルチスケール・ブートストラップ法の理論

- 仮説検定の確率値
- 確率値とブートストラップ確率の関係
- マルチスケール・ブートストラップ法の原理
- ブートストラップ確率のスケーリング則

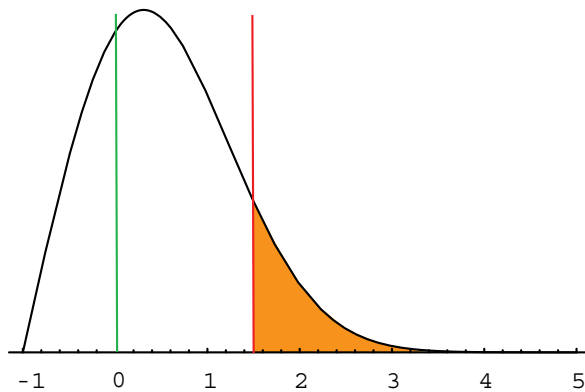
仮説検定

- データのバラツキによる見せかけか否かを判定する統計的手続き
- 「驚き」の程度 \Rightarrow 確率値 (p-value, p-値)



確率値とブートストラップ確率の関係

仮説 データ



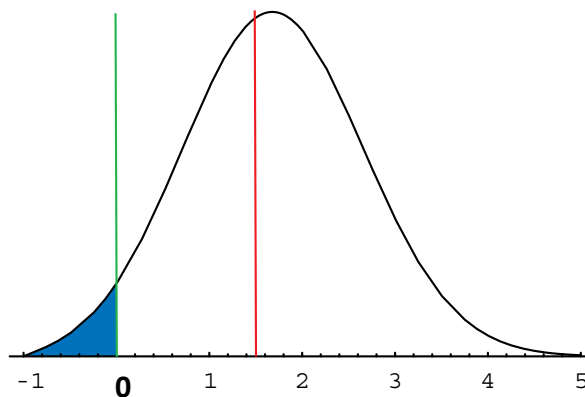
検定の確率値: 12%



もし分布が左右対称ならば両者は一致する



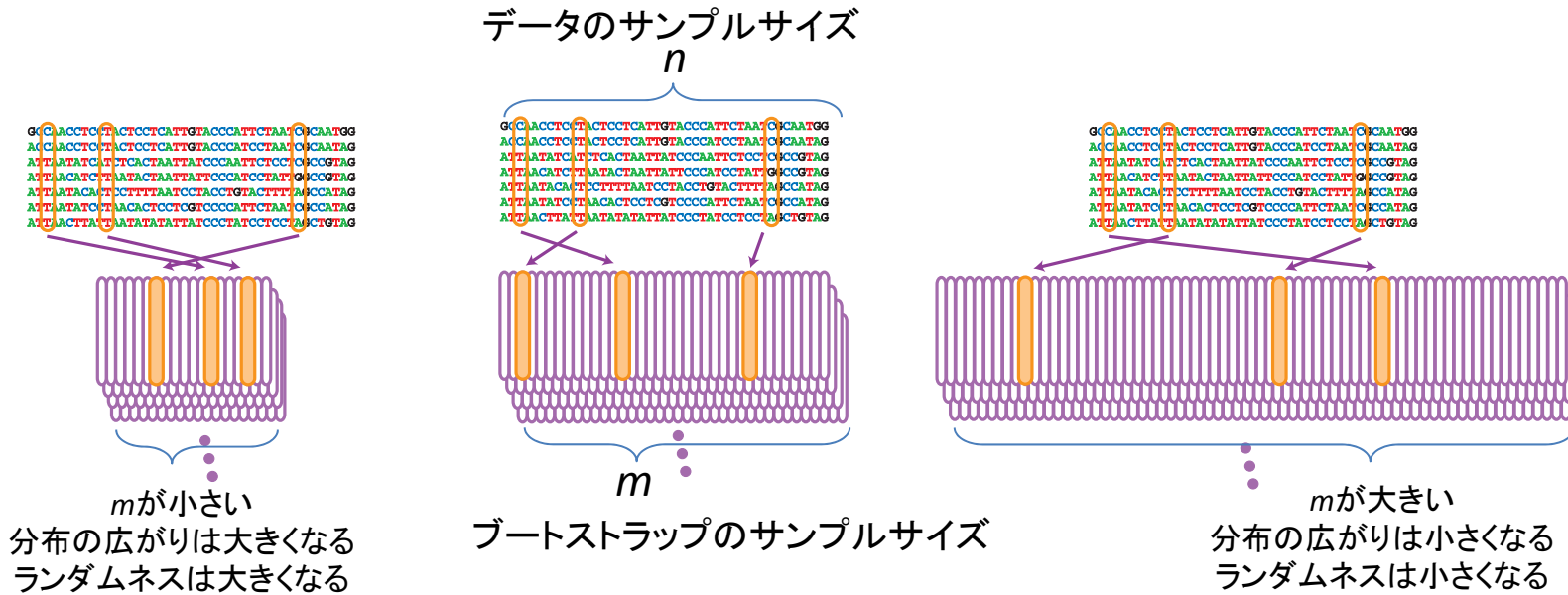
ブートストラップ確率: 3%



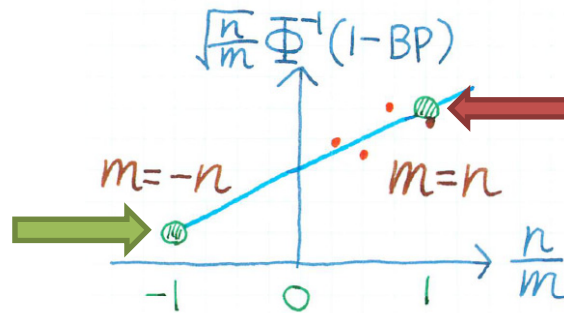
もし下図の分布を左右反転させてからブートストラップ確率を計算すれば上図に近づく

マルチスケール・ブートストラップの原理

$m = -n$ への外挿 (“ m out of n bootstrap”)

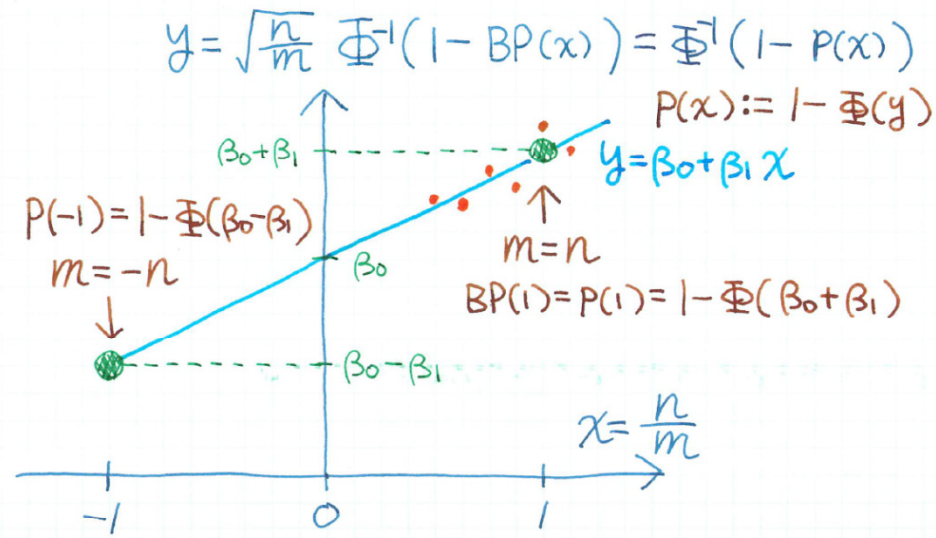


ココが $m = -n$ への外挿
分布が左右反転する？
ランダムネスがひっくり返る？

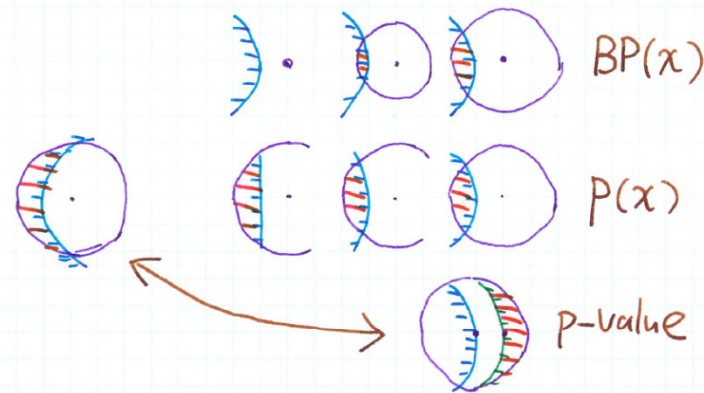


ココが $m = n$ の通常の
ブートストラップ

マルチスケール・ブートストラップ法の幾何学



$m = -n$



仮説境界を表す曲面の曲率とデータまでの距離が関係している

Example: exact $p=0.05$

Normalized bootstrap z-value

$$\Psi(\sigma^2) = -\sigma \Phi^{-1}(BP(\sigma^2))$$

$$\Phi(-\Psi(\sigma^2))$$

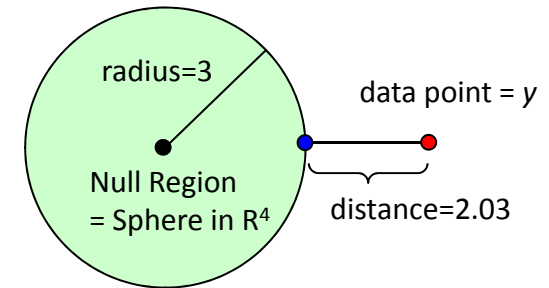
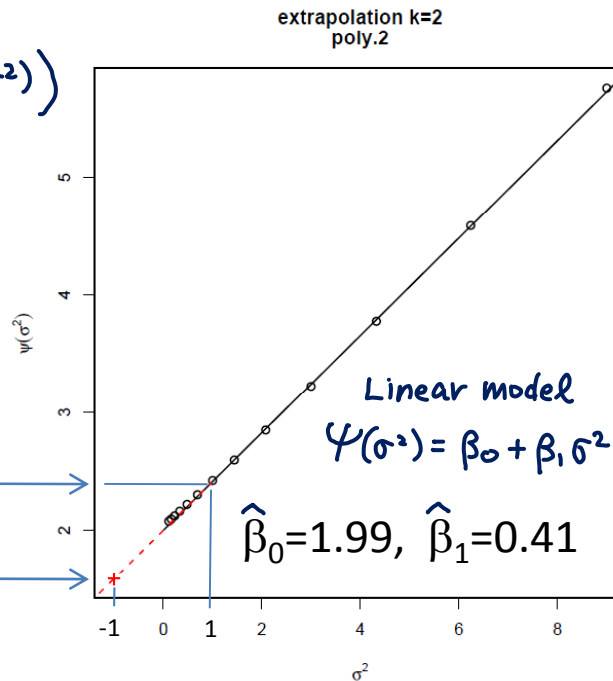
$$BP=0.0078 \quad (k=1)$$

$$AU=0.057 \quad (k=2)$$

$$\Psi(\sigma^2)$$

$$2.40$$

$$1.58$$



Quadratic model

$$\Psi(\sigma^2) = \beta_0 + \beta_1 \sigma^2 + \beta_2 \sigma^4$$

$$AU=0.052 \quad (k=3)$$

$$\hat{\beta}_0=2.02, \hat{\beta}_1=0.39, \hat{\beta}_2=0.024$$

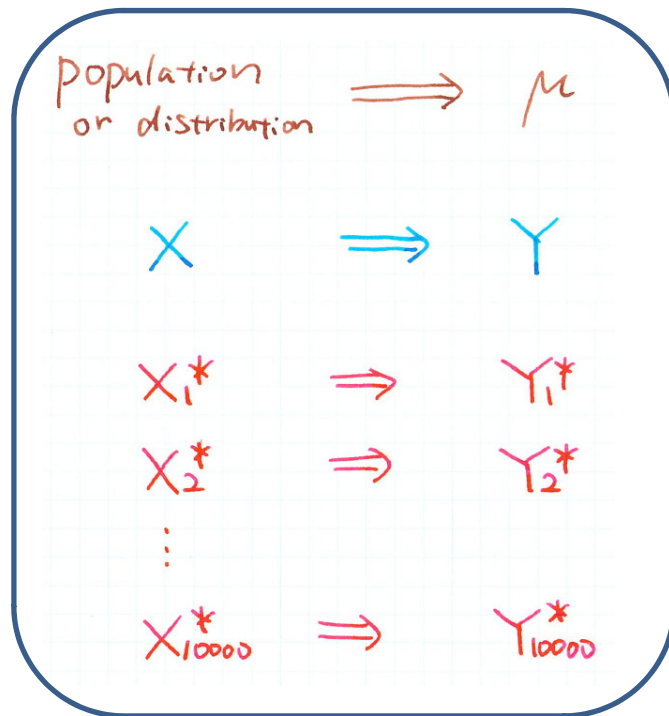
Geometrical Interpretations of the Coefficients

β_0 : signed distance

β_1 : mean curvature

β_2 : mean 4th derivatives

A Simplified Model



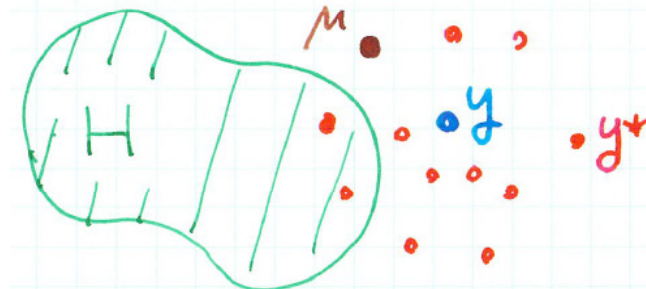
Multivariate Normal Model

$$Y \sim N(\mu, I)$$

$$Y^* | y \sim N(y, \sigma^2 I)$$

$$\sigma^2 = \frac{n}{m}$$

$$BP = P_{\sigma^2}(Y^* \in H | y)$$



Why $m=-n$?

(only for illustration purpose)

Intuitive Explanation of $m=-n$

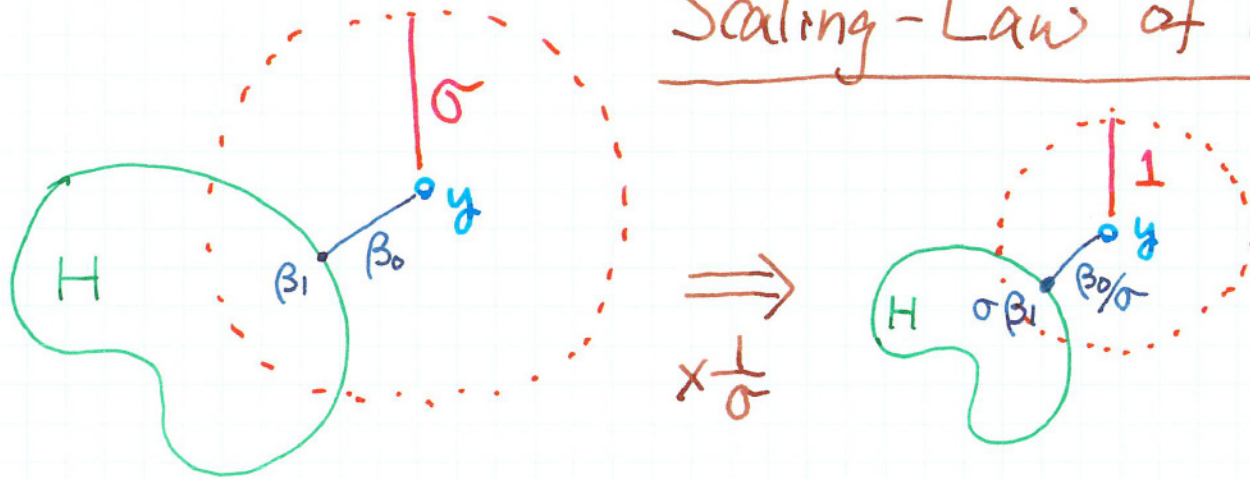
$$\left. \begin{array}{l} Y \sim N(\mu, I) \\ Y^* - Y \sim N(0, \sigma^2 I) \end{array} \right\} \Rightarrow Y^* \sim N(\mu, \underline{(1+\sigma^2)I})$$

$$\therefore Y^* \equiv \mu \text{ for } \sigma^2 = -1 \text{ or } m = -n.$$

Rescaling the Randomness: bridging Bayesian to Frequentist

ブートストラップ確率のスケールリング則

Scaling-Law of BP



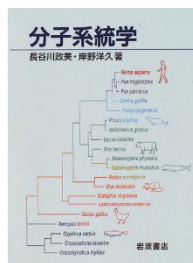
$$BP = 1 - \Phi[\text{distance} + \text{curvature} + \dots] \quad \text{for } \sigma = 1$$

$$= 1 - \Phi\left[\frac{\beta_0}{\sigma} + \sigma\beta_1 + \dots\right] \quad \text{for } \forall \sigma > 0$$

$$= 1 - \Phi\left[\frac{1}{\sigma} \{ \beta_0 + \sigma^2\beta_1 + \dots \} \right]$$

$$\therefore \sigma \Phi^{-1}(1 - BP) = \beta_0 + \sigma^2\beta_1 + \dots$$

参考文献



- 分子系統学 長谷川政美・岸野洋久
岩波書店 1996
- そのほか多数あります

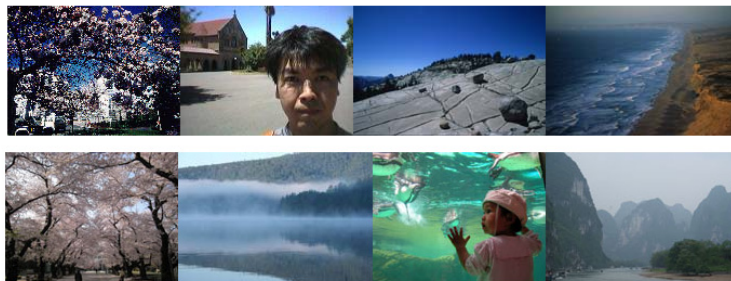
- 分子進化と分子系統学 根井 正利, S. クマー, Sudhir Kumar, 大田 竜也 (2006/7)
 新品: ¥ 7,350
 12時間以内に「お急ぎ便」でご注文いただくと、2010/6/27 日曜日までにお届けします。
 ★★★★★ (1) Amazonプライム
 和書: 全2冊品を見る
- 分子進化—解析の技法とその応用 西田 隆 (1998/7)
 新品: ¥ 3,675
 1冊品 ¥ 3,675より 2冊品 ¥ 2,390より
 12時間以内に「お急ぎ便」でご注文いただくと、2010/6/27 日曜日までにお届けします。
 ★★★★★ (3) Amazonプライム
 2ページの引用: "この時期には、分子進化学における知識が飛躍的に拡大し、他分野への応用も盛んに行われた。研究の進んでいくと、現存 ... たどる"分子系統進化学,の流れと、生物やその構成要素である遺伝子やタンパク 質の進化機構の解明を目指す"分子 ...
 和書: 全2冊品を見る
- 分子系統学 長谷川 政美 岸野 洋久 (1996/3)
 2冊品 ¥ 5,000より
 和書: 全2冊品を見る
- 分子系統学への統計的アプローチ—計算分子進化学 Ziheng Yang, 藤 博幸, 加藤 和典, 大安 裕美 (2009/3/20)
 新品: ¥ 6,090
 1冊品 ¥ 6,090より 1冊品 ¥ 12,180より
 12時間以内に「お急ぎ便」でご注文いただくと、2010/6/27 日曜日までにお届けします。
 Amazonプライム
 著作権の引用: "共著(2006年)分子系統学への統計的アプローチ 計算分子進化学 Computational Molecular Evolution 2009 ...
 和書: 全2冊品を見る
- ゲノム進化学入門 CD-ROM付 斎藤 成也 (2007/12/20)
 新品: ¥ 3,570
 1冊品 ¥ 3,570より 2冊品 ¥ 3,890より
 12時間以内に「お急ぎ便」でご注文いただくと、2010/6/27 日曜日までにお届けします。
 ★★★★★ (1) Amazonプライム
 索引の引用: "... 分枝過程(branching process)81 ,83 ,92 ,106 分子系統学"
 和書: 全2冊品を見る
- 生物系統学 (Natural History) 三中 徳宏 (1997/12)
 新品: ¥ 5,880
 1冊品 ¥ 5,880より 2冊品 ¥ 5,100より

amazon.co.jpで「分子系統学」を検索

下平 ホームページ

<http://www.is.titech.ac.jp/~shimo/index-j.html>

下平英寿 (しもだいらひでとし)



統計科学 & バイオインフォマティクス

略歴 (Linkedin), 文献リスト (被引用回数ランキング)

下平英寿 (しもだいらひでとし)

博士 (工学), 准教授

〒152-8552 東京都目黒区大岡山 2-12-1-W8-46
[東京工業大学](http://www.is.titech.ac.jp) [情報理工学研究科](http://www.is.titech.ac.jp) [数理・計算科学専攻](http://www.is.titech.ac.jp)

shimo (あつとまーく) [is.titech.ac.jp](http://www.is.titech.ac.jp)
<http://www.is.titech.ac.jp/~shimo/>