

先端生命科学実験 I

バイオインフォマティクス

遺伝的アルゴリズムによるタンパク質立体構造の重ね合わせ

藤 博幸

この実験の目的

- (1) プログラミング言語Rの使い方を学ぶ。
- (2) 進化の分子機構と**遺伝的アルゴリズム**を理解する。
- (3) 遺伝的アルゴリズムを具体的な生物の問題であるタンパク質立体構造の重ね合わせに応用する。(世代数、突然変異率の影響など)

OUTLINE

1. 遺伝的アルゴリズム
2. タンパク質の立体構造
3. Rの復習
4. 相同性の説明
5. 遺伝的アルゴリズムによる相同タンパク質の立体構造の重ね合わせ
6. レポートの構成

1. 遺伝的アルゴリズム

生物システムの情報科学への応用

- ニューラルネットワーク
- 遺伝的アルゴリズム
- 進化プログラミング
- 人工免疫システム
- 粒子群最適化
- 蟻コロニー最適化
- 人工生命

遺伝的アルゴリズム

菊川怜の卒論テーマ

「遺伝的アルゴリズムを適用したコンクリートの要求性能型の調合設計に関する研究」



遺伝的アルゴリズム

Genetic Algorithm (GA)

1975年 John H. Hollandにより提案

生物の遺伝の仕組みを模倣した最適化の手法

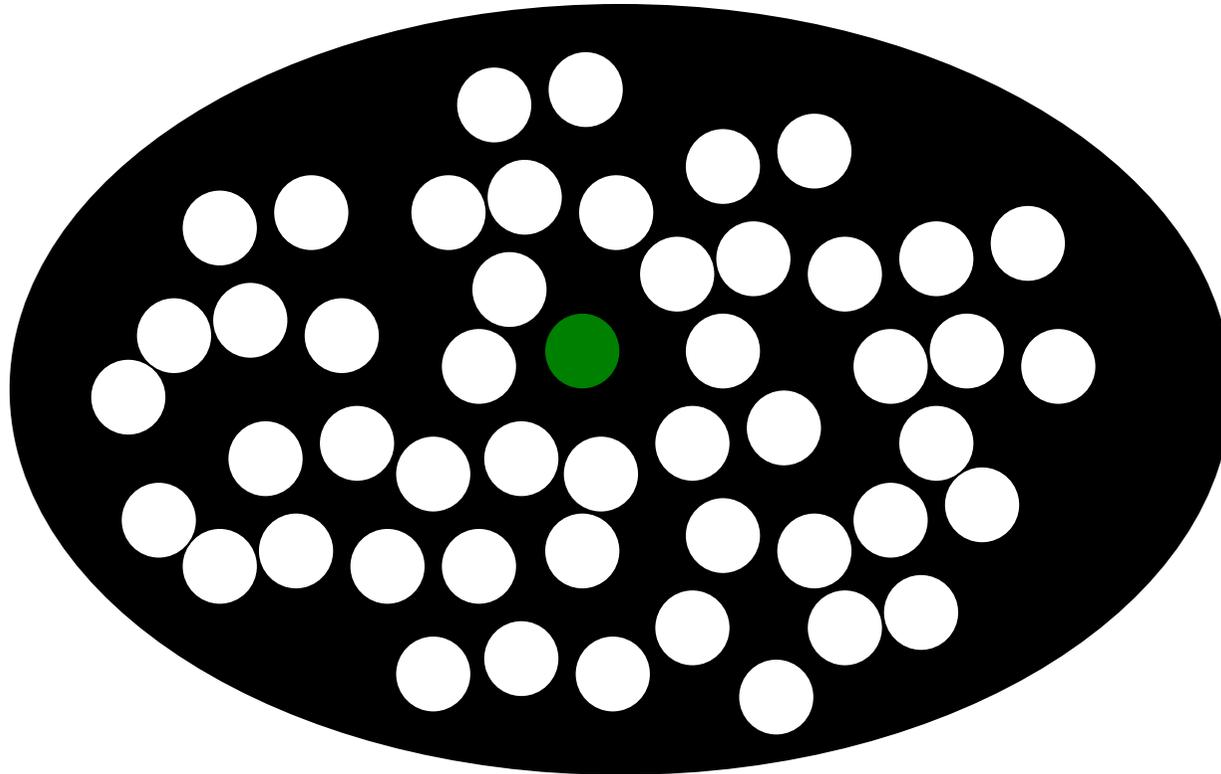
たとえば

ナップザック問題、ハミルトン閉路問題、

巡回セールスマン問題、施設配置問題

などの最適化問題に適用できる

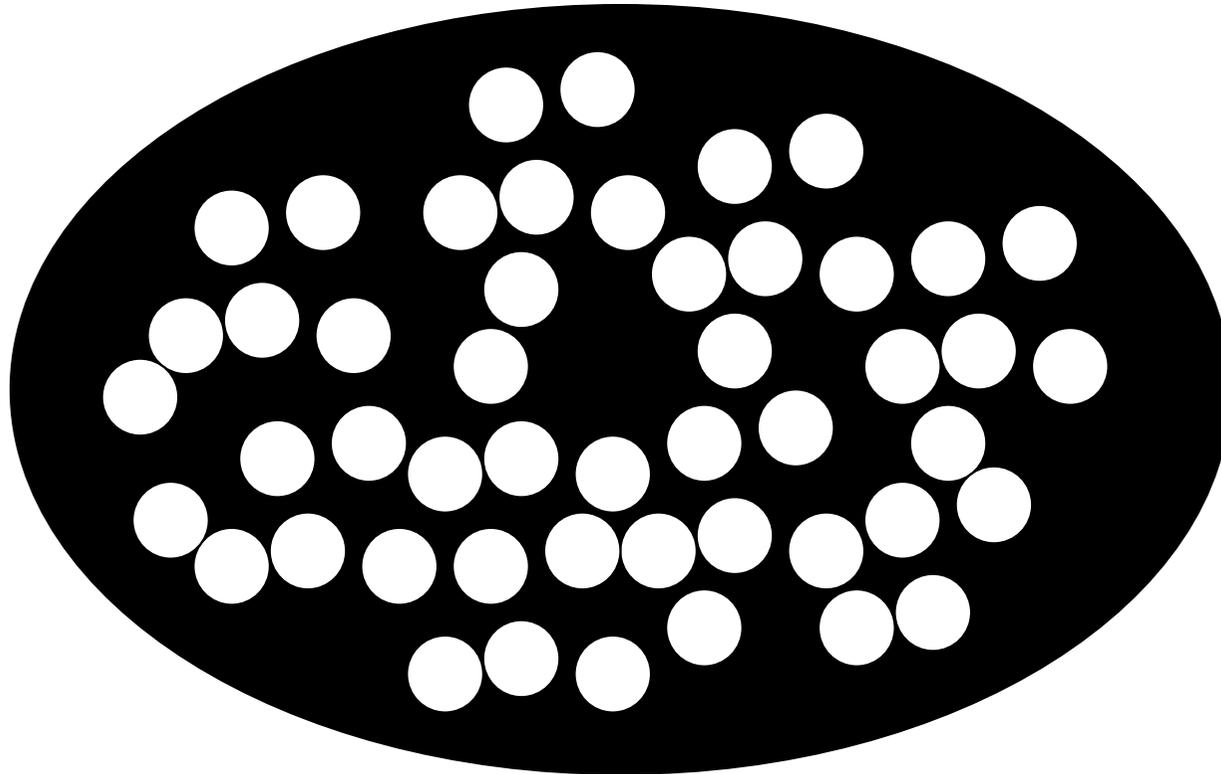
突然変異(mutation)と置換(substitution)



突然変異は集団中の個体に生じる

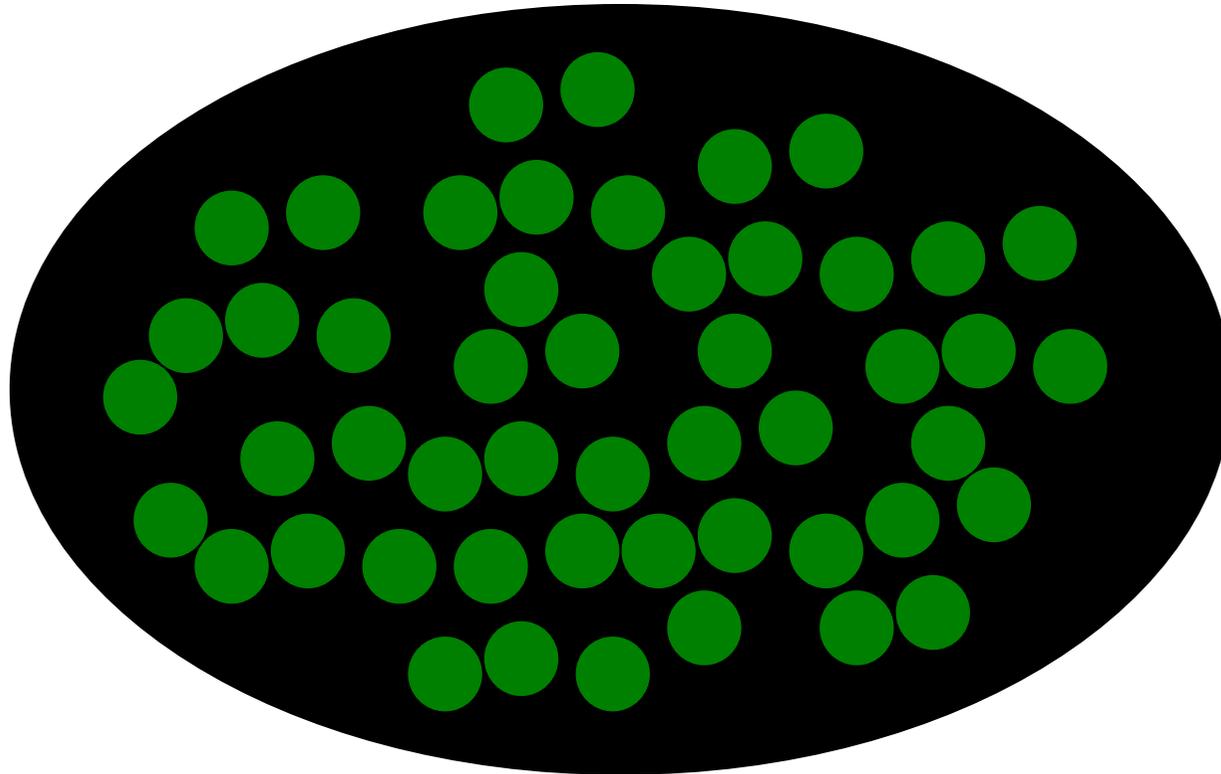
※ 進化に寄与するのは体細胞ではなく生殖系列
の細胞に生じる突然変異

突然変異(mutation)と置換(substitution)



有害な突然変異は、集団から除去される
(**負の選択** or **純化淘汰**)

突然変異(mutation)と置換(substitution)



有利な突然変異は急速に集団中に広まり
集団全体がその突然変異遺伝子で置き換
えられる（**正の選択**）

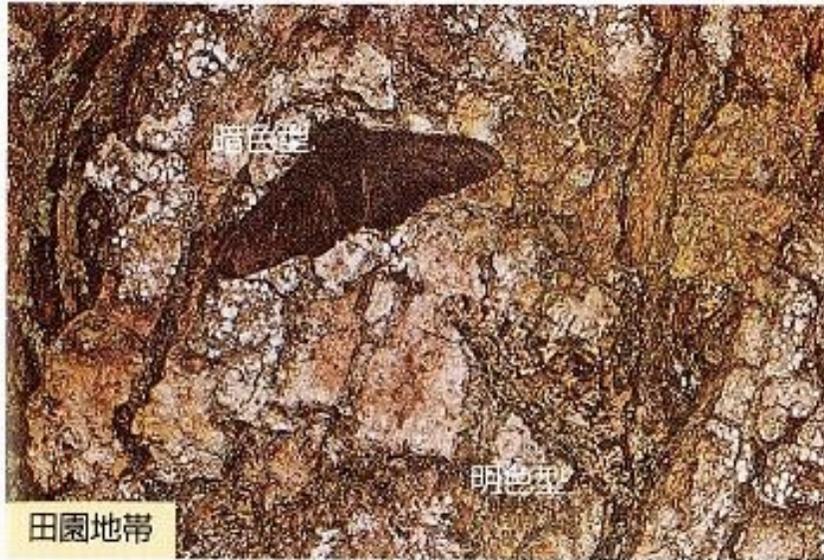
工業暗化 (industrial melanism)



もともと淡色型と暗色型の体色を持つオオモリエダシャクが存在

19世紀後半から、ヨーロッパの工業都市が発展するにつれて、その付近に生息するガ(蛾)に暗色の個体が増加した

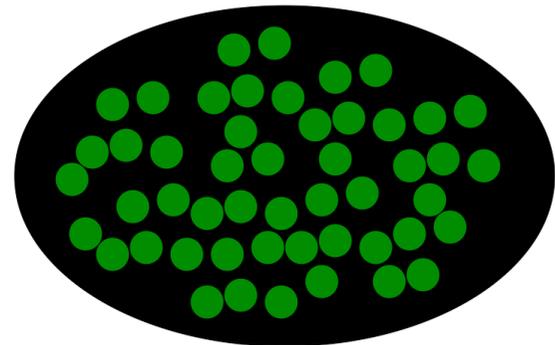
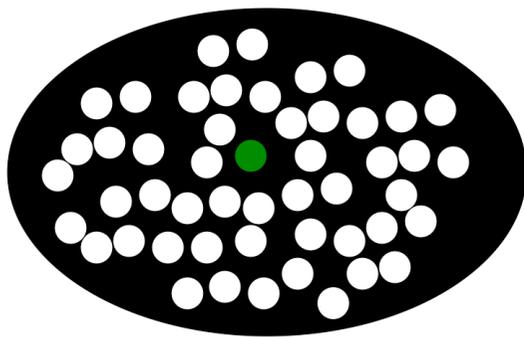
田園都市では白っぽい地衣類が木にはえていて淡色型の方が保護色となり鳥の補色を免れて生き残るのに対して、工業化に伴い地衣類が枯れ、煤煙で木が黒くすすけてしまったため暗色型のほうが保護色となるため



樹皮は白っぽい地衣植物でおおわれており、
明るい色をしている。



大気汚染のために地衣植物がなくなり、
樹皮は暗い色をしている。



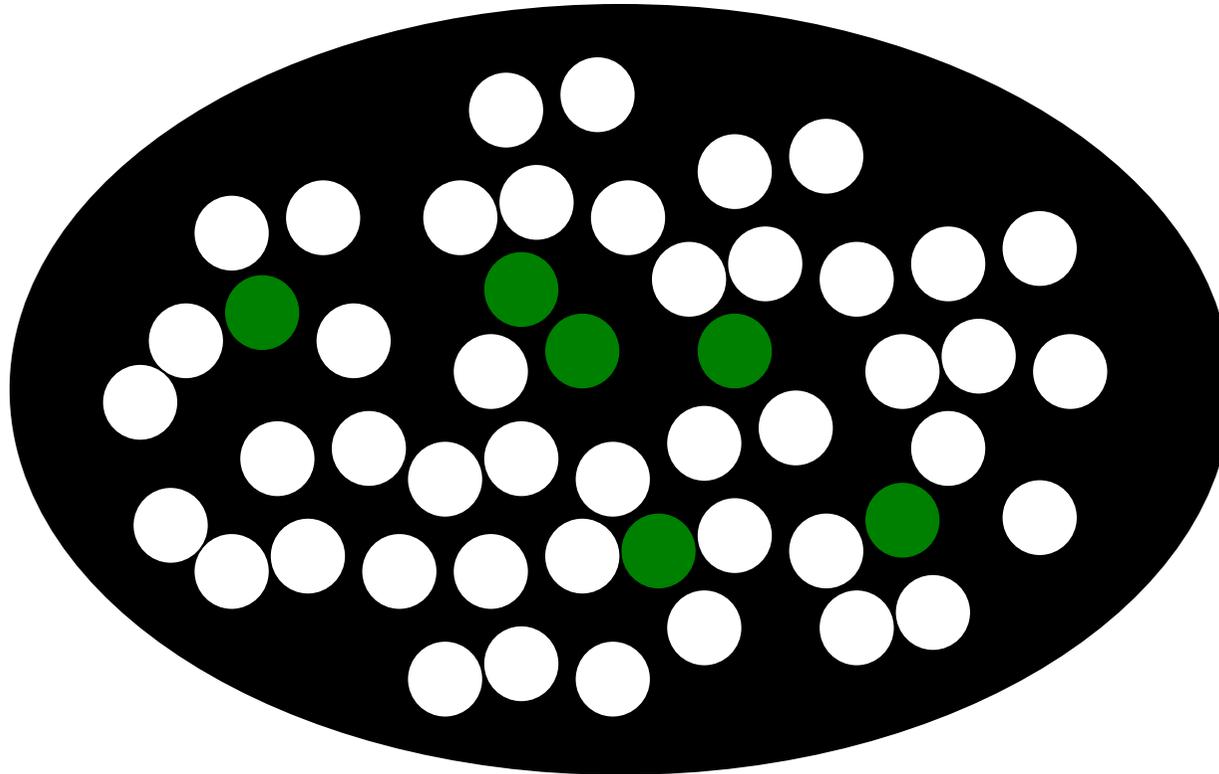
補足

大進化と小進化

工業暗化では、種内での形質の変化は起きているが
“種分化”はおきていない。このようなレベルの進化を小進化とよぶ

これに対し、新しい種あるいは種より高次の分類群の形成、また絶滅
などのレベルの進化を大進化とよぶ

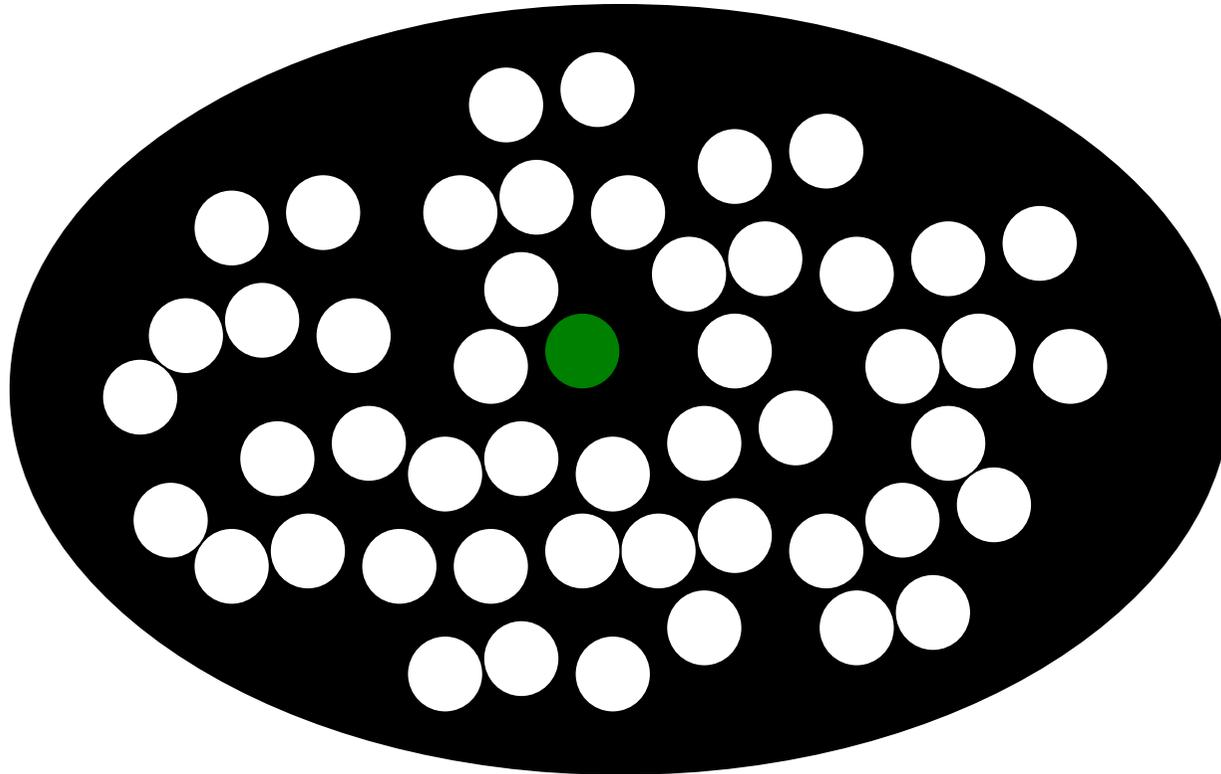
突然変異(mutation)と置換(substitution)



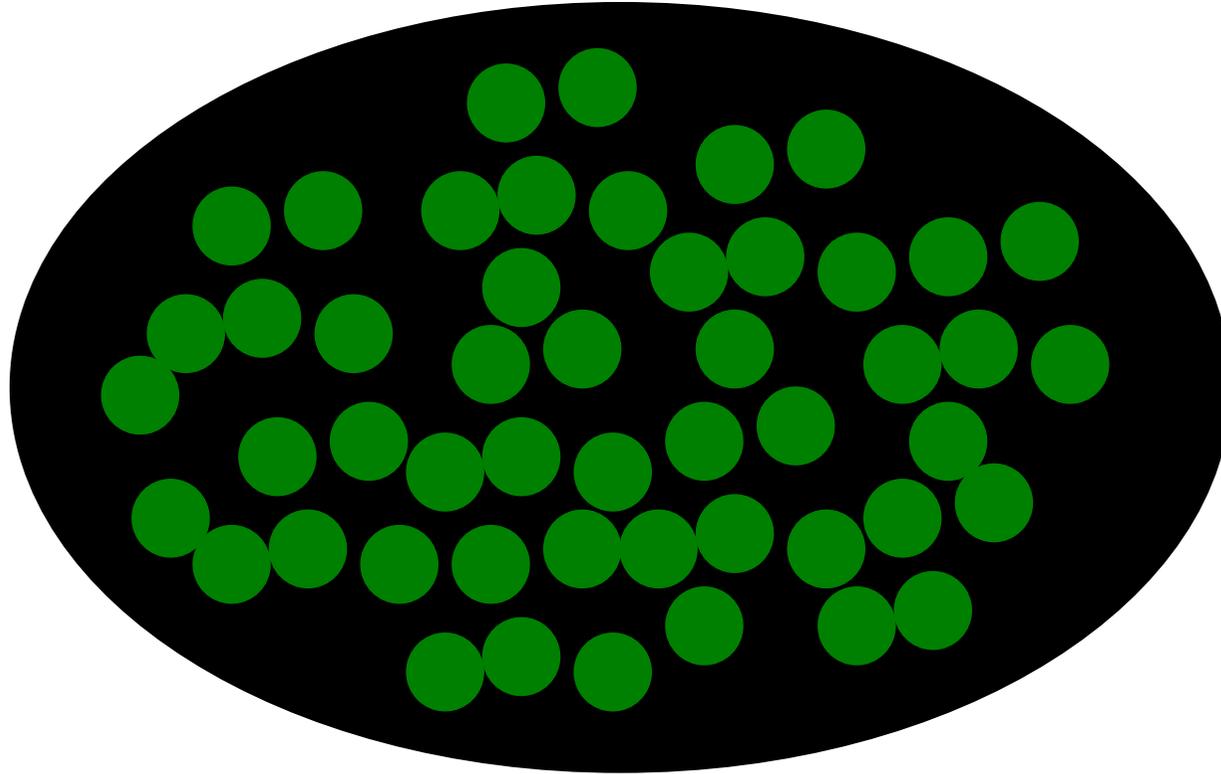
中立な突然変異の頻度はランダムに浮動し
確率的に集団中に固定（選択的に中立）

固定にいたる途中の過程 多型(polymorphism)

突然変異(mutation)と置換(substitution)



突然変異(mutation)と置換(substitution)



個体に生じた突然変異が集団全体に広まること: **固定**

突然変異が固定されること: **置換**

The industrial melanism mutation in British peppered moths is a transposable element

Arjen E. van't Hof^{1*}, Pascal Campagne^{1*}, Daniel J. Rigden¹, Carl J. Yung¹, Jessica Lingley¹, Michael A. Quail², Neil Hall¹, Alistair C. Darby¹ & Ilik J. Saccheri¹

Discovering the mutational events that fuel adaptation to environmental change remains an important challenge for evolutionary biology. The classroom example of a visible evolutionary response is industrial melanism in the peppered moth (*Biston betularia*): the replacement, during the Industrial Revolution, of the common pale *typica* form by a previously unknown black (*carbonaria*) form, driven by the interaction between bird predation and coal pollution¹. The *carbonaria* locus has been coarsely localized to a 200-kilobase region, but the specific identity and nature of the sequence difference controlling the *carbonaria*-*typica* polymorphism, and the gene it influences, are unknown². Here we show that the mutation event giving rise to industrial melanism in Britain was the insertion of a large,

involved in wing pattern development or melanization. By extending the association mapping approach to a larger population sample and more closely spaced genetic markers (see Methods), we narrowed the *carbonaria* candidate region to about 100 kb (Fig. 1a). The candidate region resides entirely within the span of one gene — the orthologue of *Drosophila cortex* (*cort*), the only known function of which is as a cell-cycle regulator during meiosis¹¹. In *B. betularia*, *cortex* consists of eight non-first exons, multiple alternative first exons (of which only two, 1A and 1B, are strongly expressed in developing wing discs), and a very large first intron (Fig. 1b).

The rapid spread of *carbonaria* gave rise to strong linkage disequilibrium², such that many sequence variants are associated with the *carbonaria* phenotype. This poses a challenge for isolating the specific

オオシモフリエダシャクの17番染色体 cortex遺伝子
ショウジョウバエでは減数分裂時の細胞周期の制御に関与
黒い蛾のcortex 遺伝子第一イントロンにトランスポゾン

羽化時に発現が上昇するが、トランスポゾン挿入によるさらに発現が上昇

まだ、この遺伝子が工業暗化をもたらすことの機構までは解明されていない

適応度 (fitness)

ある個体の子供の内の、繁殖年齢まで達したものの数

遺伝的アルゴリズム

Genetic Algorithm (GA)

1975年 John H. Hollandにより提案

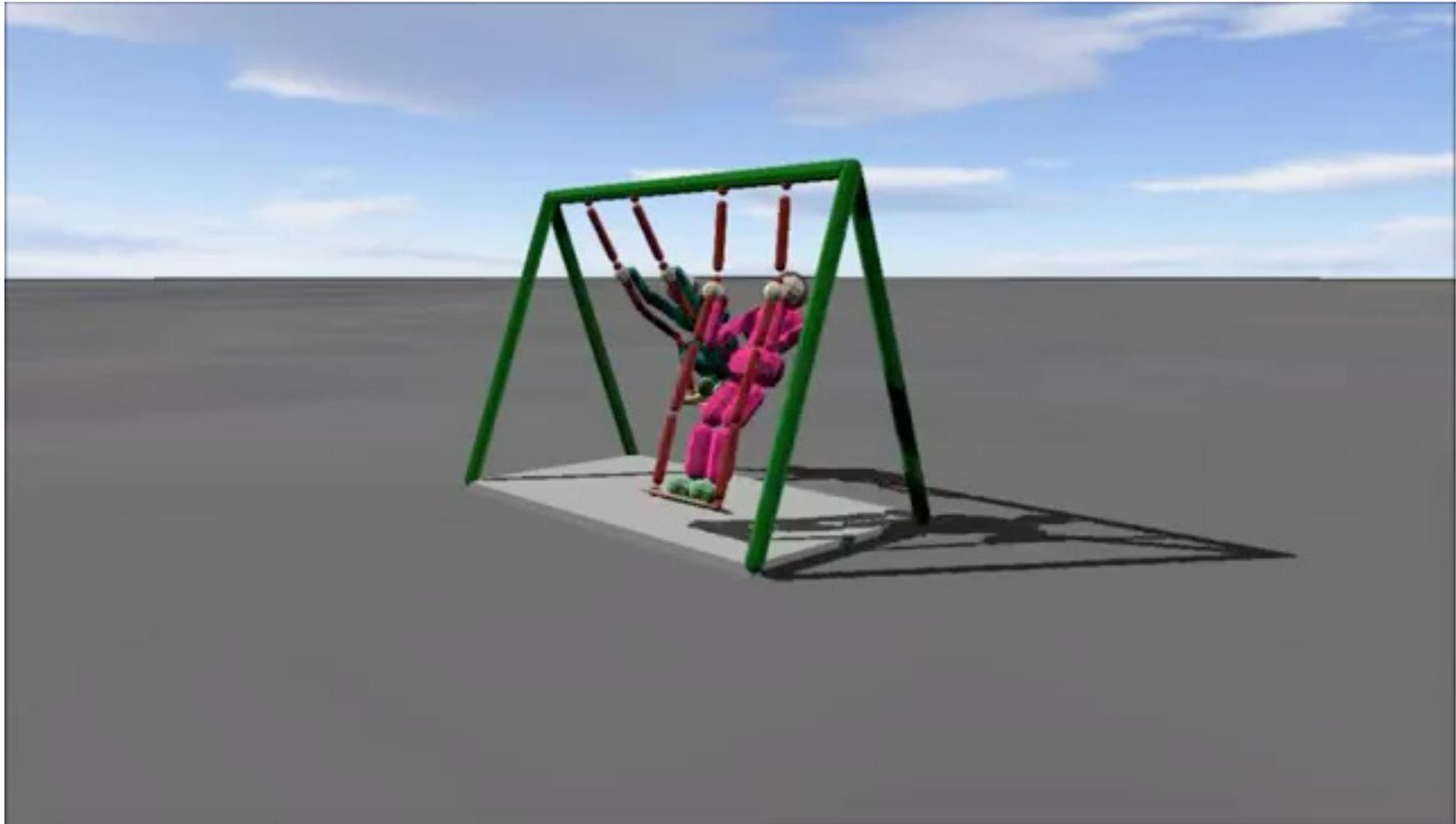
生物の遺伝の仕組みを模倣した最適化の手法
たとえば

ナップザック問題、ハミルトン閉路問題、
巡回セールスマン問題、施設配置問題
などの最適化問題に適用できる

巡回セールスマン問題を例として遺伝的アルゴリズムを学ぶ

遺伝的アルゴリズムでブランコの漕ぎ方を学習させた。

<http://www.youtube.com/watch?v=8vzTCC-jbwM#t=14>



物理エンジン

物理エンジン: コンピュータ上の仮想的な3次元空間において重力、摩擦力などを自動計算して物体の動きをシミュレーションするソフトウェア

Chipmunk Physics

<http://chipmunk-physics.net/>

PhysX

http://developer.nvidia.com/object/physx_downloads.html

BULLET

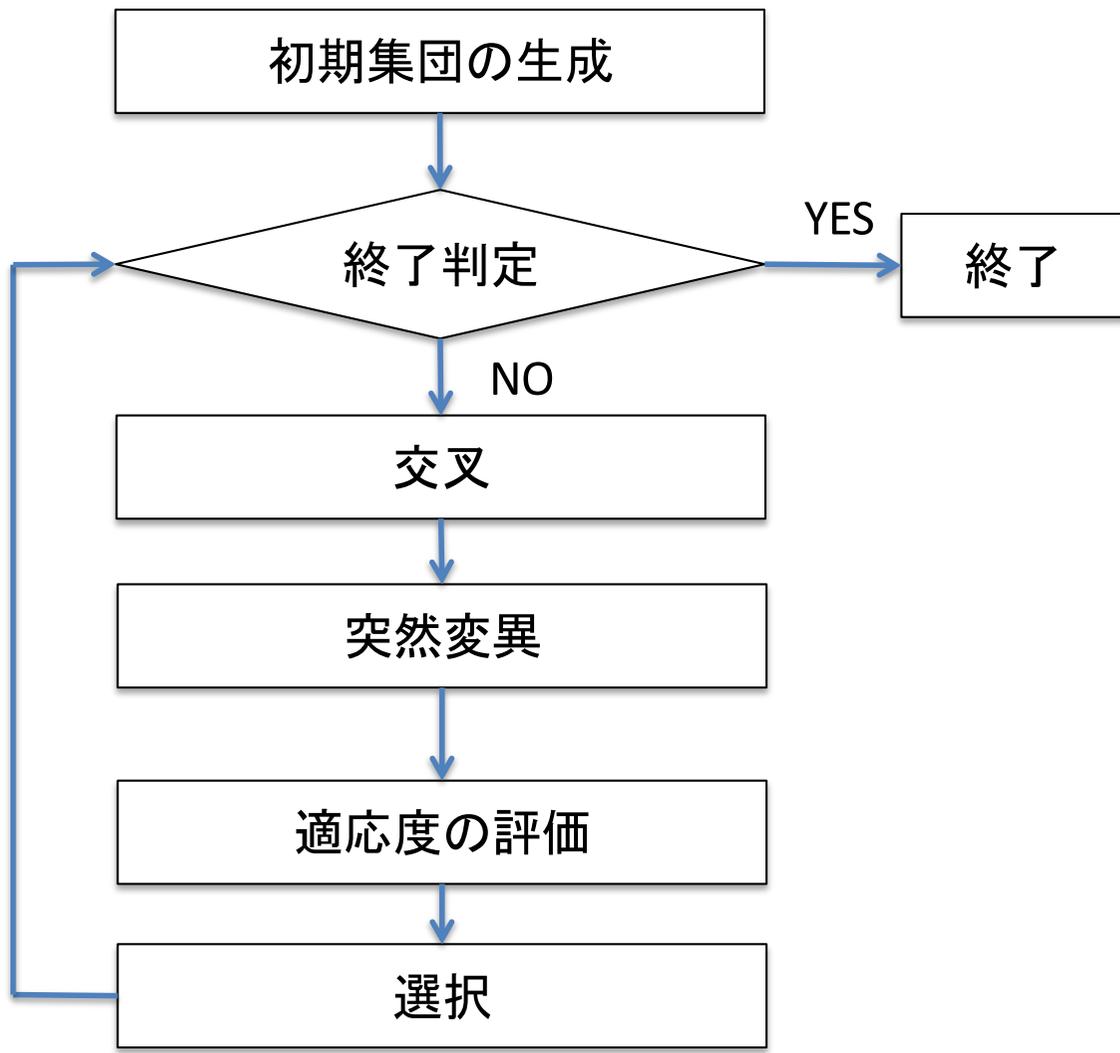
<http://code.google.com/p/bullet/downloads/list>

Dynamo

<http://home.iae.nl/users/starcat/dynamo/dynamo.zip>

Springhead

<http://springhead.info/wiki/index.php?plugin=attach&refer=Springhea...>



初期集団の生成

終了判定

YES

終了

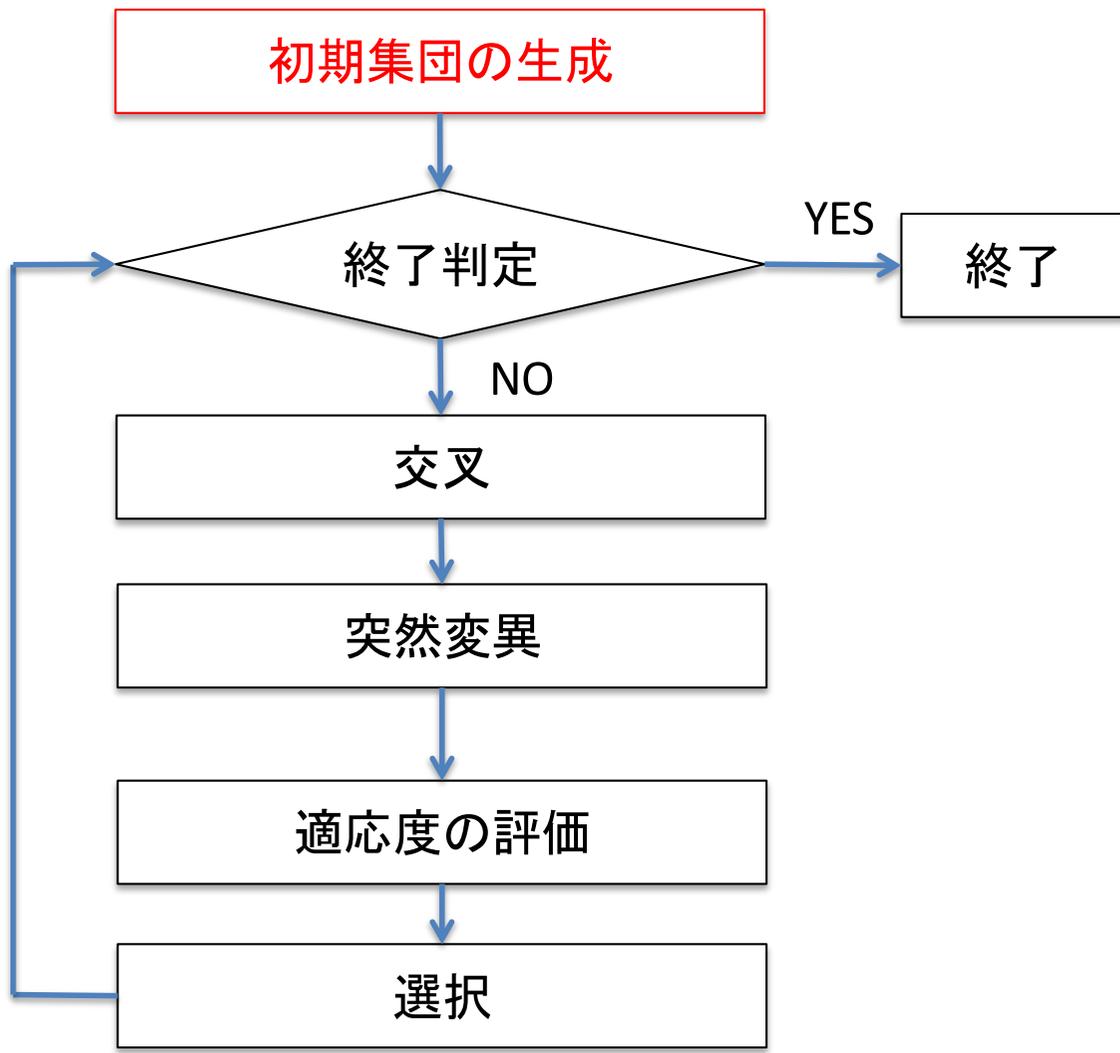
NO

交叉

突然変異

適応度の評価

選択



初期集団の生成

終了判定

YES

終了

NO

交叉

突然変異

適応度の評価

選択

初期集団の生成

集団サイズを3とする

3個体を生成し、ブランクの立った座ったを10でランダムに表現

個体1: 1 1 1 0 0 1 ... 0 1 0 1 1

個体2: 0 1 0 0 1 1 ... 1 0 0 1 0

個体3: 0 0 1 1 1 0 ... 1 1 0 0 0

用語の説明

個体 (Individual) : 染色体(後で説明)によって特徴づけられた自律的な個

集団 (Population) : 個体の集まり

集団サイズ (Population Size) : 集団内の個体数

用語の説明

- **遺伝子 (Gene)** : 個体の形質を規定する基本構成要素
個体2の場合 0 1 0 0 1 1 ... 1 0 0 1 0の1, 0のそれぞれが遺伝子
- **染色体 (Chromosome)** : 複数の遺伝子の集まり
個体2の場合、0 1 0 0 1 1 ... 1 0 0 1 0の全体が染色体
- **遺伝子座 (Locus)** : 染色体上の遺伝子の位置
- **対立遺伝子 (Allele)** : 遺伝子がとりうる値
遺伝子座2は、0, 1の2つの値をとる。
この例では、個体1と2では1、個体3で0

個体1:	1	1	1	0	0	1	...	0	1	0	1	1
個体2:	0	1	0	0	1	1	...	1	0	0	1	0
個体3:	0	0	1	1	1	0	...	1	1	0	0	0

用語の説明

- **遺伝子型** (Genotype) : 染色体の内部表現
- **表現型** (Phenotype) : 染色体によって規定される形質の
外部的表現. 巡回セールスマン問題の場合、巡路

表現型	遺伝子型
個体1: ブランコの漕ぎ方1	11100 1 ... 01011
個体2: ブランコの漕ぎ方2	01001 1 ... 10010
個体3: ブランコの漕ぎ方3	00111 0 ... 11000

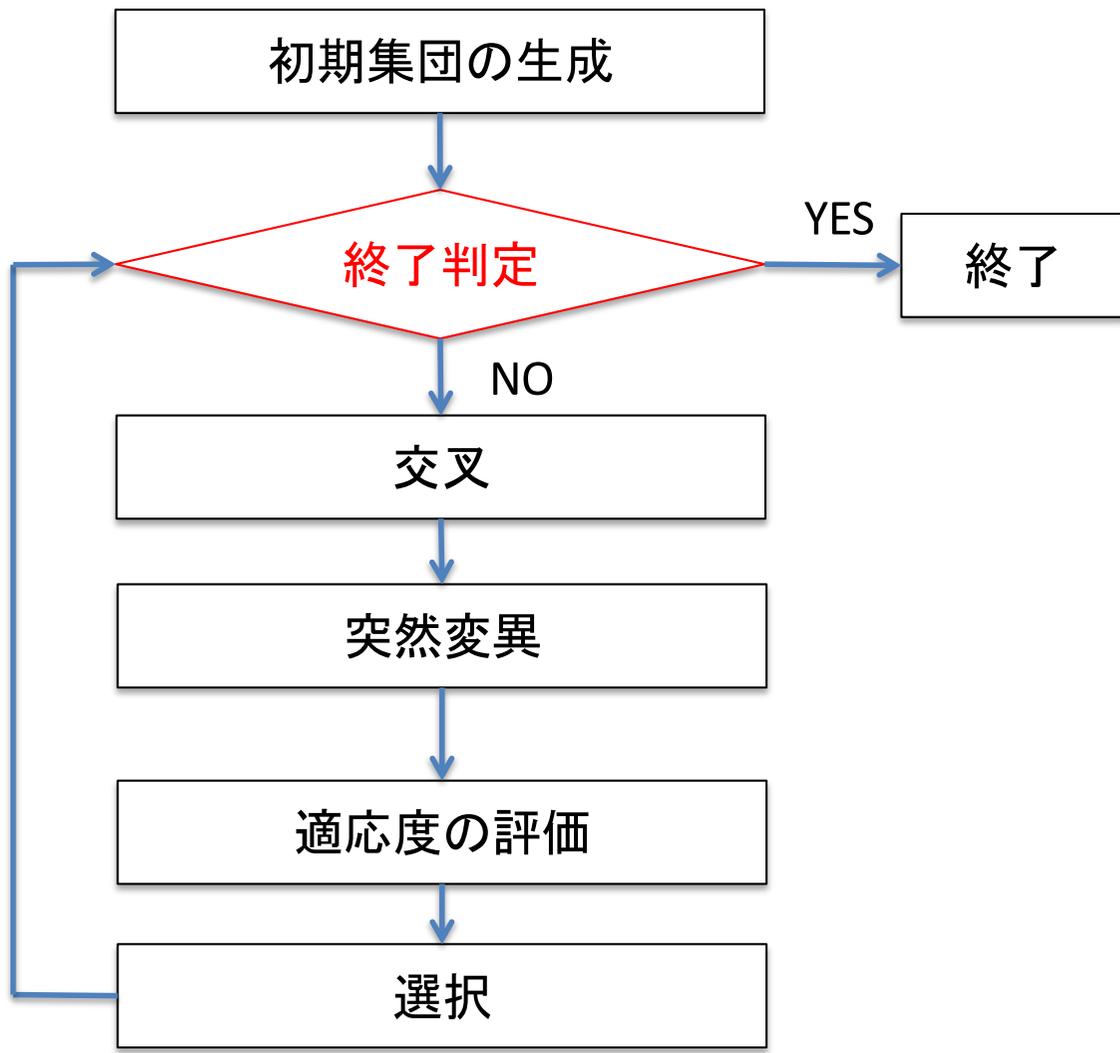
- **コード化** (Coding) : 表現型から遺伝子型へ変換すること
- **デコード化** (Decoding) : 遺伝子型から表現型変換すること

初期集団完成

個体1:1 1 1 0 0 1 ... 0 1 0 1 1

個体2:0 1 0 0 1 1 ... 1 0 0 1 0

個体3:0 0 1 1 1 0 ... 1 1 0 0 0



初期集団の生成

終了判定

YES

終了

NO

交叉

突然変異

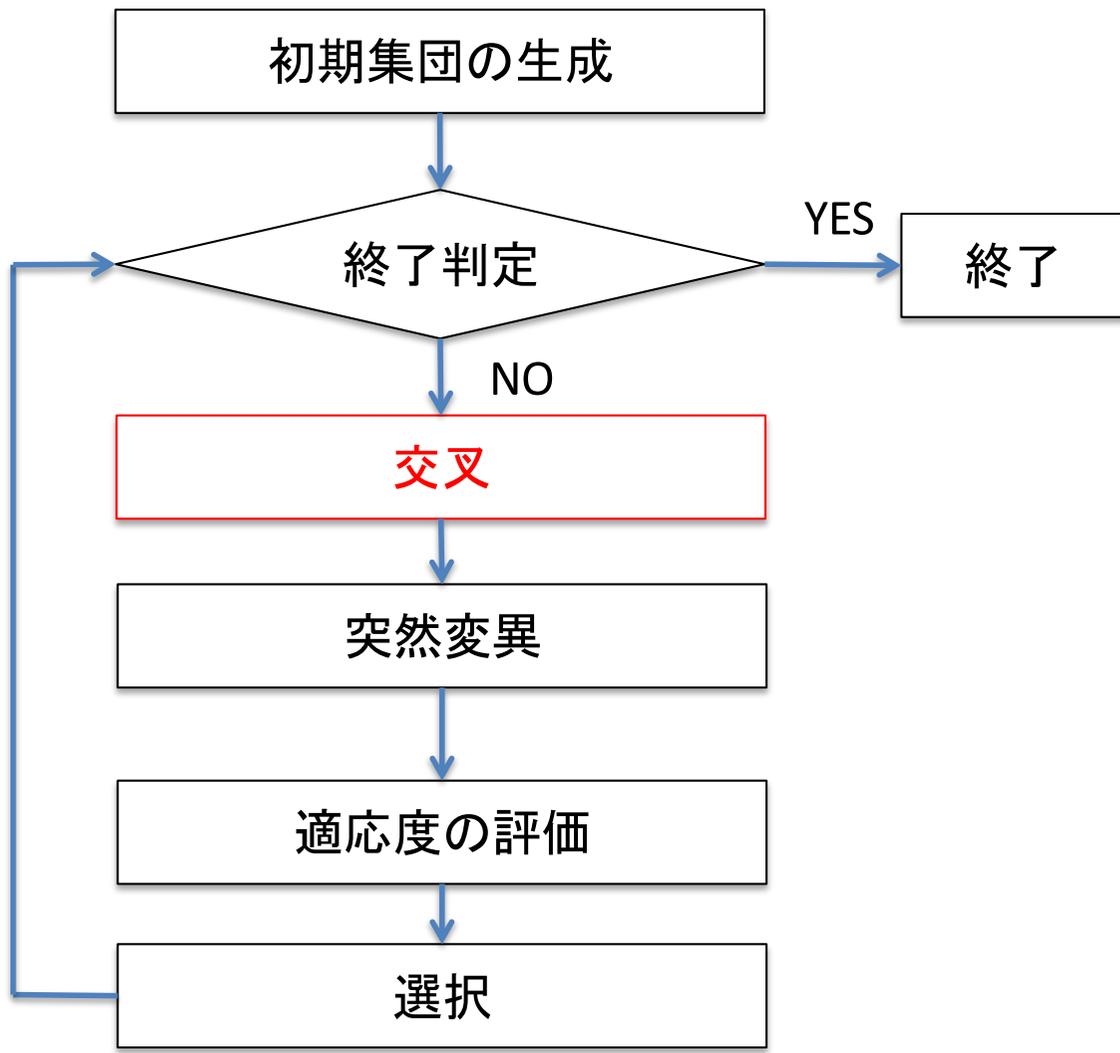
適応度の評価

選択

終了の判定

判定方法としていくつかのものが考えられる

- 集団中の最大の適応度が、ある閾値より大きくなった場合
- 集団全体の平均適応度が、ある閾値より大きくなった場合
- 集団の適応度の増加率がある閾値以下になる世代が一定期間続いた場合
- 世代交代の回数が規定の回数を超えた場合



初期集団の生成

終了判定

YES

終了

NO

交叉

突然変異

適応度の評価

選択

交叉

- N 個の親集団からランダムに $2M$ 個の個体を選択
- ペアごとに組替え処理
- $N + 2M$ 個の個体が生成される

今、 $M = 1$ とし、乱数を利用して次の2個体を選択されたとする、

個体2: 0 1 0 0 1 1 ... 1 0 0 1 0

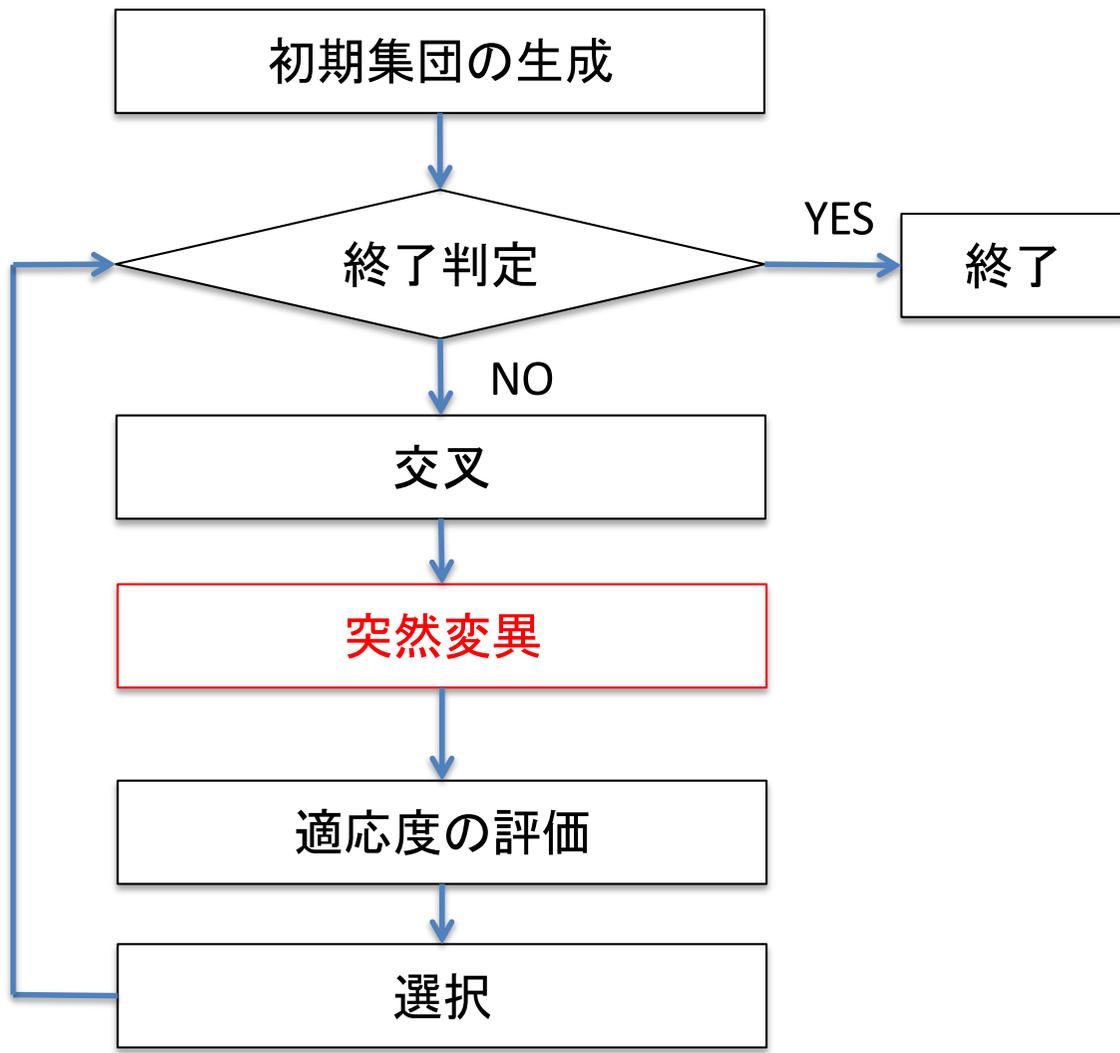
個体3: 0 0 1 1 1 0 ... 1 1 0 0 0

一様交叉を選択された個体2と個体3に適用
マスクが011011...01101とおく

011011...01101
個体2:001010...11010
個体3:010111...10000

が得られる

子集団に交叉で得られた集団(サイズ2)を加える。



初期集団の生成

終了判定

YES

終了

NO

交叉

突然変異

適応度の評価

選択

突然変異

ある与えられた確率 p_m (突然変異率) で、突然変異を起こさせる。
突然変異率は通常小さな値 (たとえば5%) に設定。

一般的な方法として、各遺伝子をランダムに対立遺伝子に置き換える。

今回の例の場合、親集団 (サイズ3) の各遺伝子座 (5個) について、0~1の一様乱数を発生させ、それが0.05より小さい場合、その遺伝子座の可能な対立遺伝子のいずれかに置換する。

- 個体1: 1 1 1 0 0 1 ... 0 1 0 1 1
- 個体2: 0 1 0 0 1 1 ... 1 0 0 1 0
- 個体3: 0 0 1 1 1 0 ... 1 1 0 0 0

ここで乱数が0.05より小さかったとする

個体3: 0 0 1 1 1 0 ... 1 1 0 0 0

遺伝子座2を、0から1に変更

0 1 1 1 1 0 ... 1 1 0 0 0

という個体を生成し、子集団に加える。

0 1を入れ替える突然変異 = bit flip mutation

子集団

交叉によって生じた

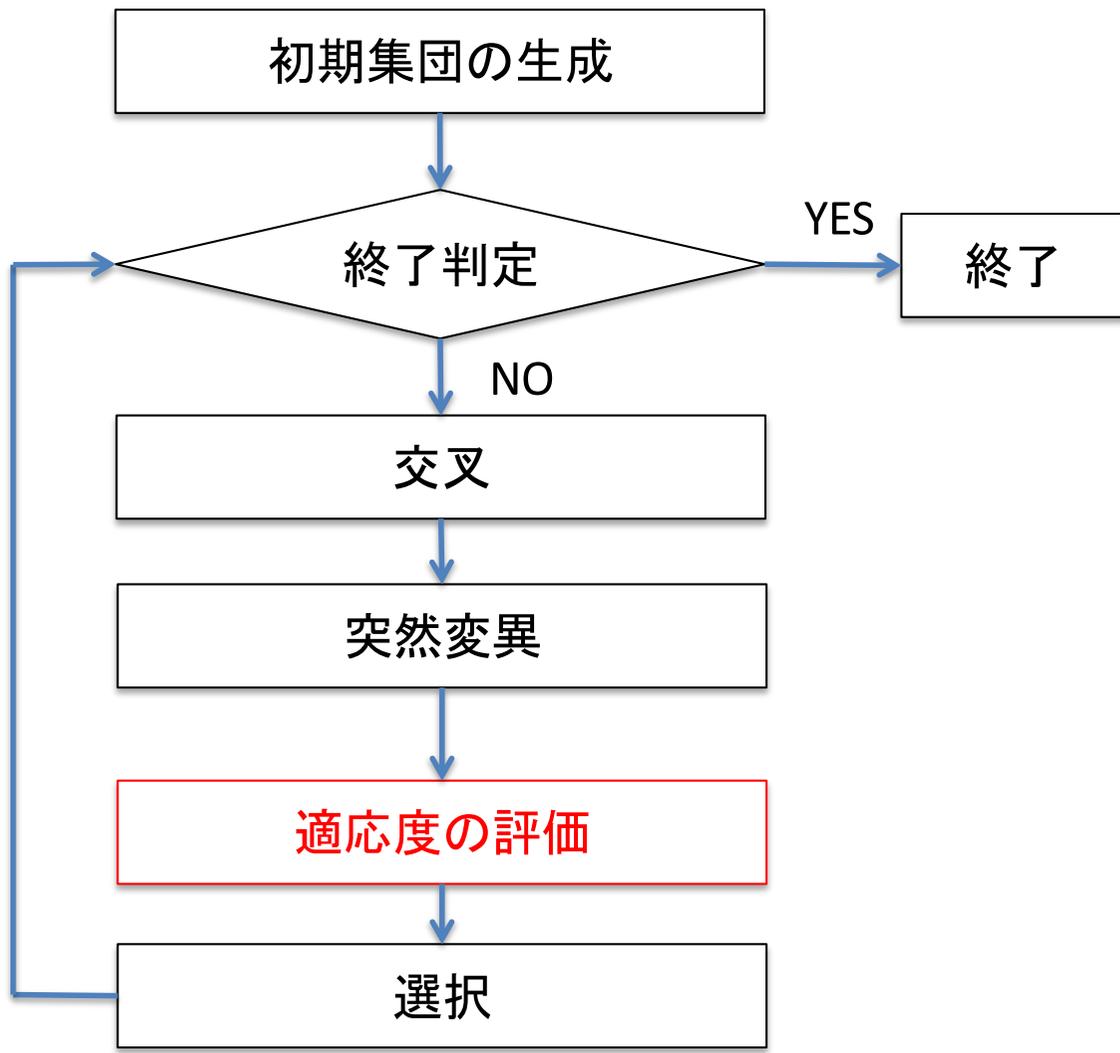
0 0 1 0 1 0 ... 1 1 0 1 0

0 1 0 1 1 1 ... 1 0 0 0 0

突然変異によって生じた

0 1 1 1 1 0 ... 1 1 0 0 0

の3個体が含まれている。



適応度 (fitness)

ある個体の子供の内の、繁殖年齢まで達したものの数

各個体の評価

親集団

11100 1... 01011

01001 1... 10010

00111 0... 11000

+

子集団

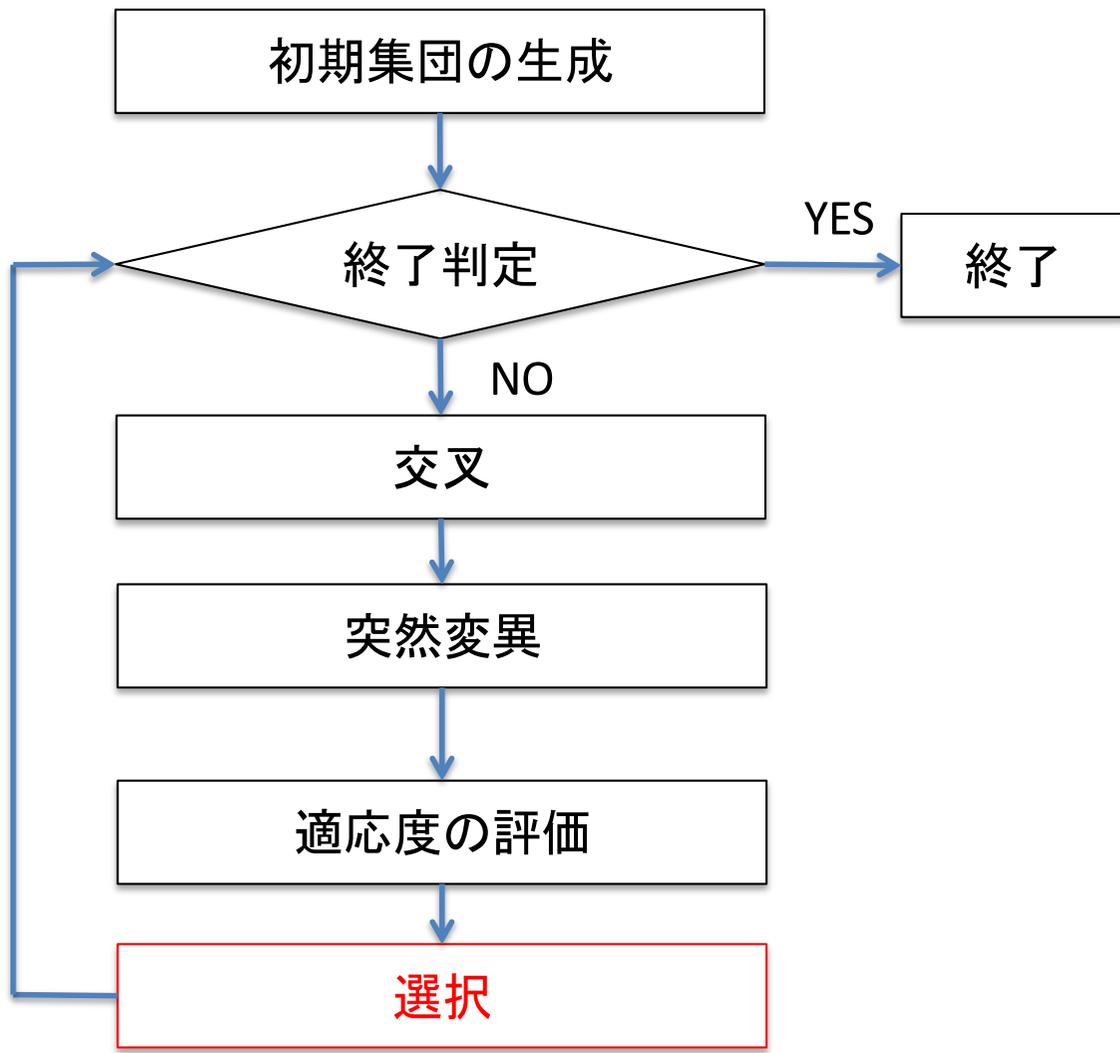
001010... 11010

010111... 10000

011110... 11000



ブランクの振れ幅に基づき各個体の適応度を計算



初期集団の生成

終了判定

YES

終了

NO

交叉

突然変異

適応度の評価

選択

選択

親集団(サイズ N) + 子集団(サイズ M)から次世代の集団を構成する N 個体を選ぶ

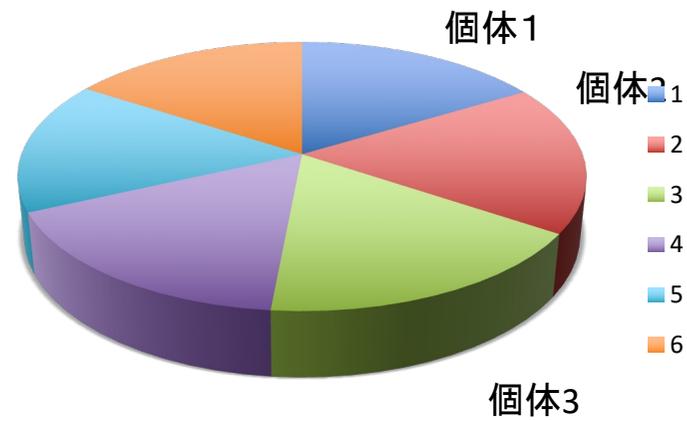
選択の方法は様々なものが考案されている

- (1)ルーレット方式
- (2)ランキング方式
- (3)定常状態選択
- (4)ボルツマン選択
- (5)トーナメント方式

...

		適応度
個体1	1 1 1 0 0 1 ... 0 1 0 1 1	20
個体2	0 1 0 0 1 1 ... 1 0 0 1 0	23
個体3	0 0 1 1 1 0 ... 1 1 0 0 0	22
個体a	0 0 1 0 1 0 ... 1 1 0 1 0	18
個体b	0 1 0 1 1 1 ... 1 0 0 0 0	17
個体c	0 1 1 1 1 0 ... 1 1 0 0 0	21

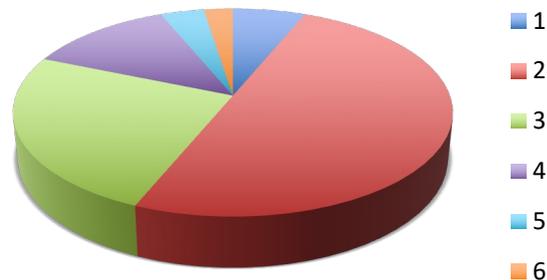
各個体をルーレット方式で選択する



ランキング方式

各個体を適応度によってランク付けして、「1位なら確率 p_1 , 2位なら確率 p_2 , 3位なら...」というふうにランクごとにあらかじめ確率を決めておく方式

		適応度	順位	確率
個体1	1 1 1 0 0 1 ... 0 1 0 1 1	20	4	0.06
個体2	0 1 0 0 1 1 ... 1 0 0 1 0	23	1	0.5
個体3	0 0 1 1 1 0 ... 1 1 0 0 0	22	2	0.25
個体a	0 0 1 0 1 0 ... 1 1 0 1 0	18	5	0.04
個体b	0 1 0 1 1 1 ... 1 0 0 0 0	17	6	0.02
個体c	0 1 1 1 1 0 ... 1 1 0 0 0	21	3	0.13



ルーレット方式の場合

- (1) 先の例の場合のように適応度にあまり差がつかない個体間に選択に関する差を導入
- (2) 最も適合度の高い個体がルーレットの90%を占める場合、他の個体は選ばれなくなる

エリート主義 (elitism)

最も高い適応度を示す個体あるいは複数の上位の適応度を示す個体を次世代に残し、残りをルーレット方式やランキング方式などで選択。

見つかった中で最も良いものを失わずにすむので、GAのパフォーマンスを向上させることができる

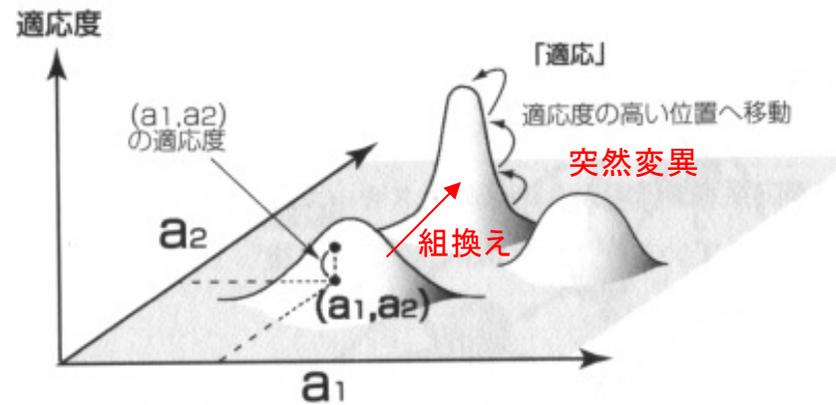
ここではランキング方式で、第二世代として以下が選択されたとする

- 010011... 10010 **23** **1** **0.5**
- 001110... 11000 **22** **2** **0.25**
- 010011... 10010 **23** **1** **0.5**

ここから終了判定移行し、終了と判定されなければ、第二世代について、交叉、突然変異、選択を実施し、第三世代を作成する。

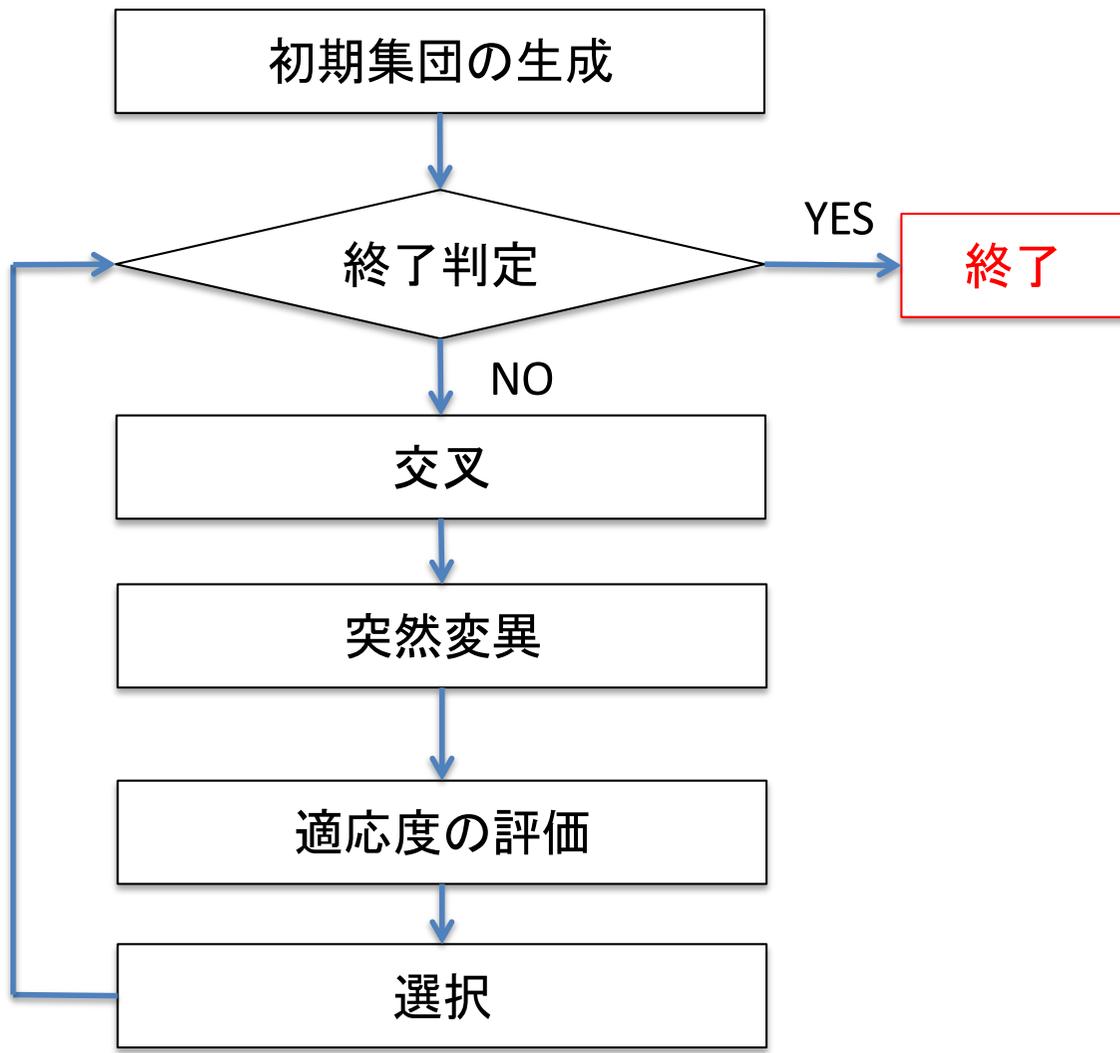
終了するまでこの操作を繰り返す。

適応度地形(フィットネス・ランドスケープ) で捉える「適応」(Adaptation)



「モデリングシミュレーション入門」
第13回 2005/01/14 井庭 崇 より

交叉によって、局所的な最適解に陥るのをふせぎ、大域的な最適解を探索できる。



初期集団の生成

終了判定

YES

終了

NO

交叉

突然変異

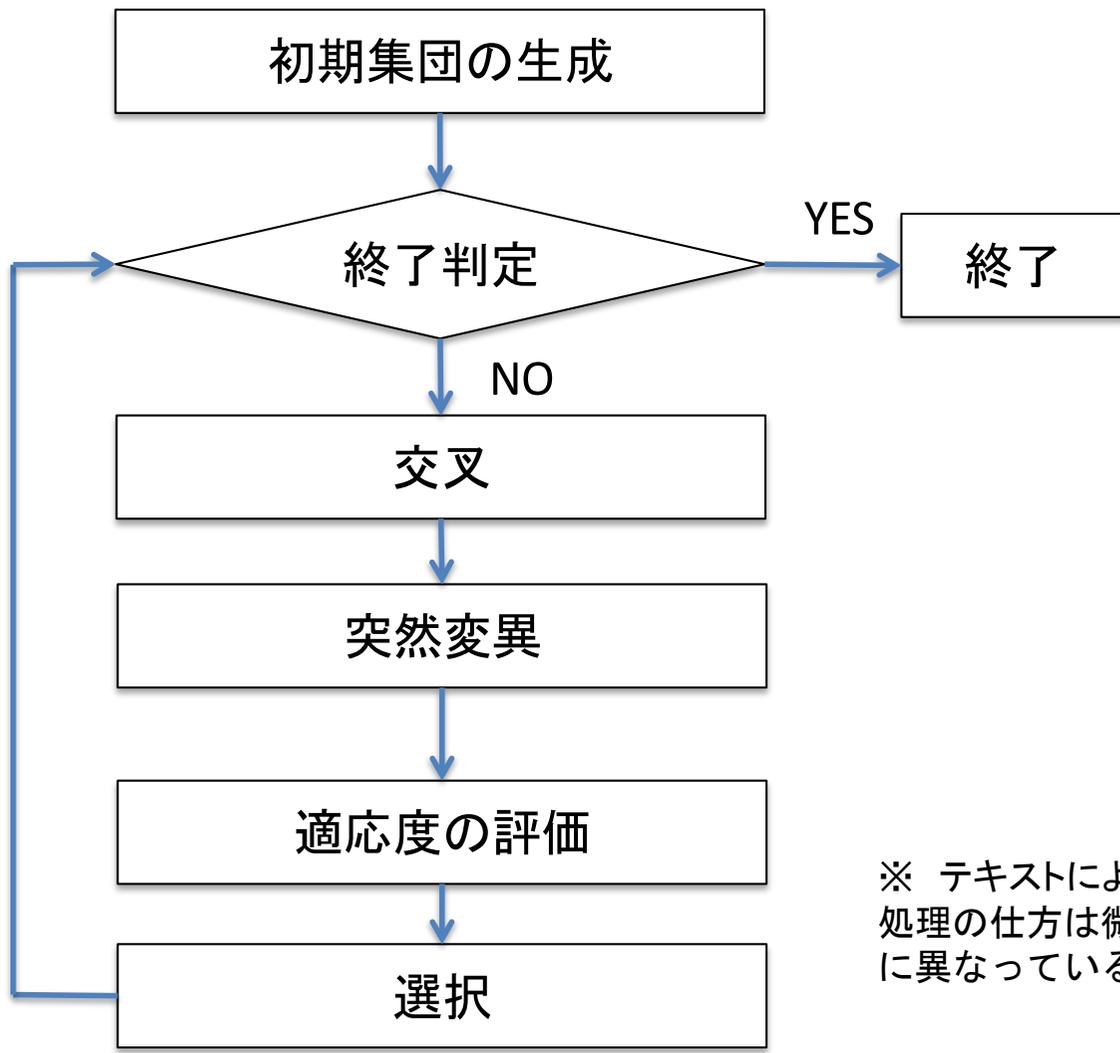
適応度の評価

選択

終了時処理

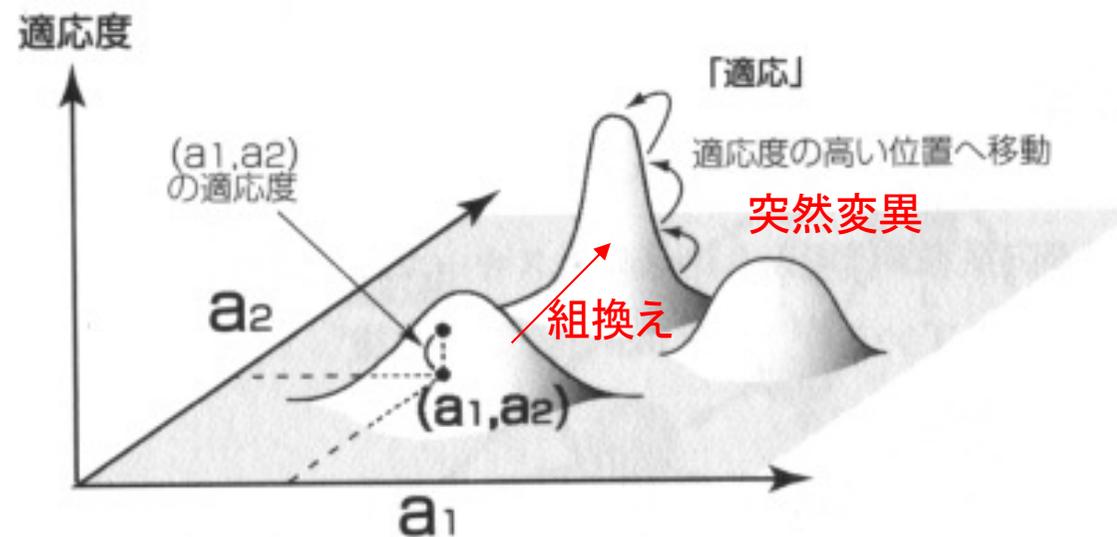
最終的に「現世代」の中で最も適応度の高い個体を「解」として出力する。

それまでの世代で生成された最も適応度の高い個体を「解」として出力する方法もある。



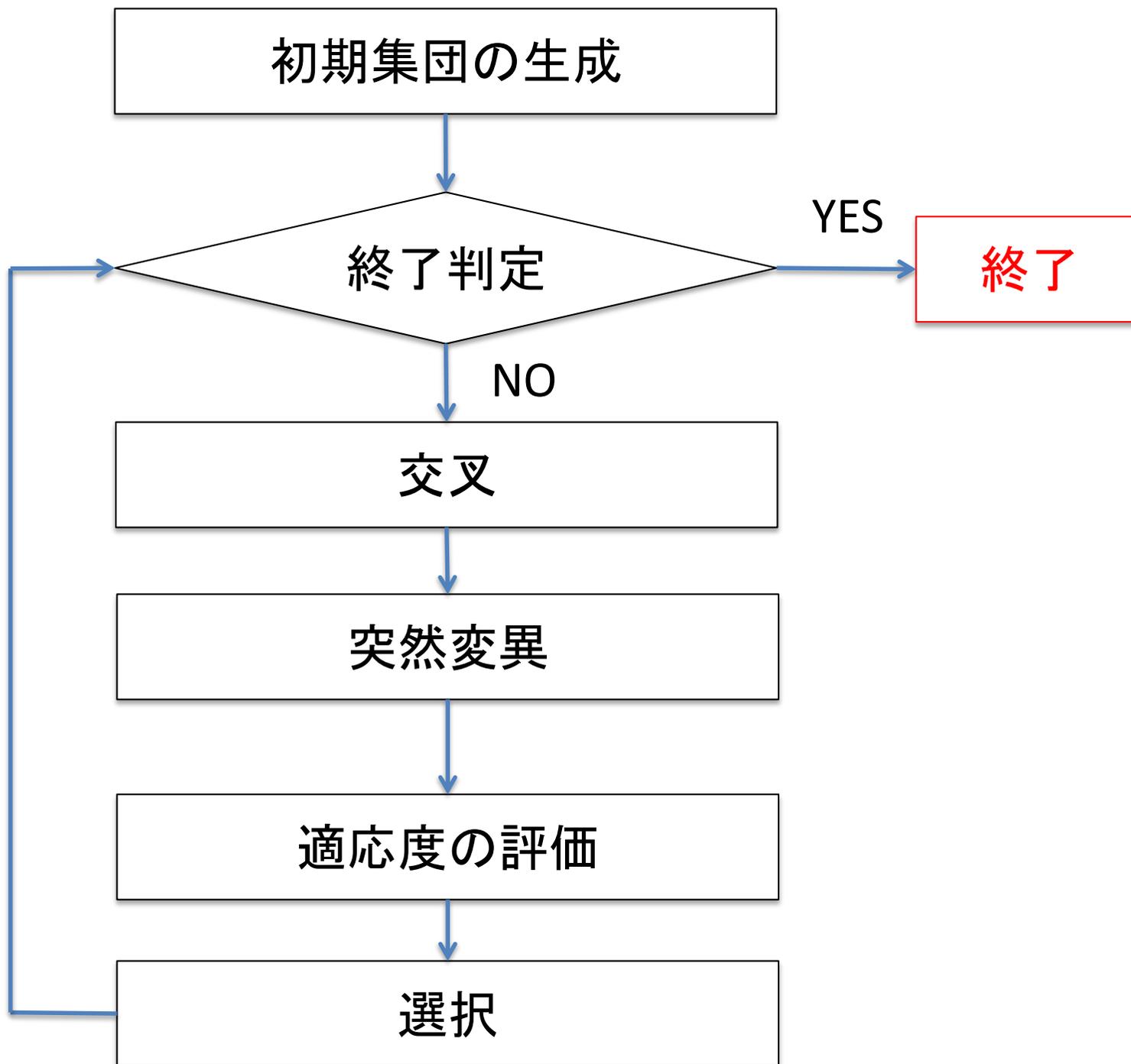
※ テキストによって
処理の仕方は微妙
に異なっている

適応度地形(フィットネス・ランドスケープ) で捉える「適応」(Adaptation)



「モデリングシミュレーション入門」
第13回 2005/01/14 井庭 崇 より

交叉によって、局所的な最適解に陥るのをふせぎ、大域的な最適解を探索できる。



初期集団の生成

終了判定

YES

終了

NO

交叉

突然変異

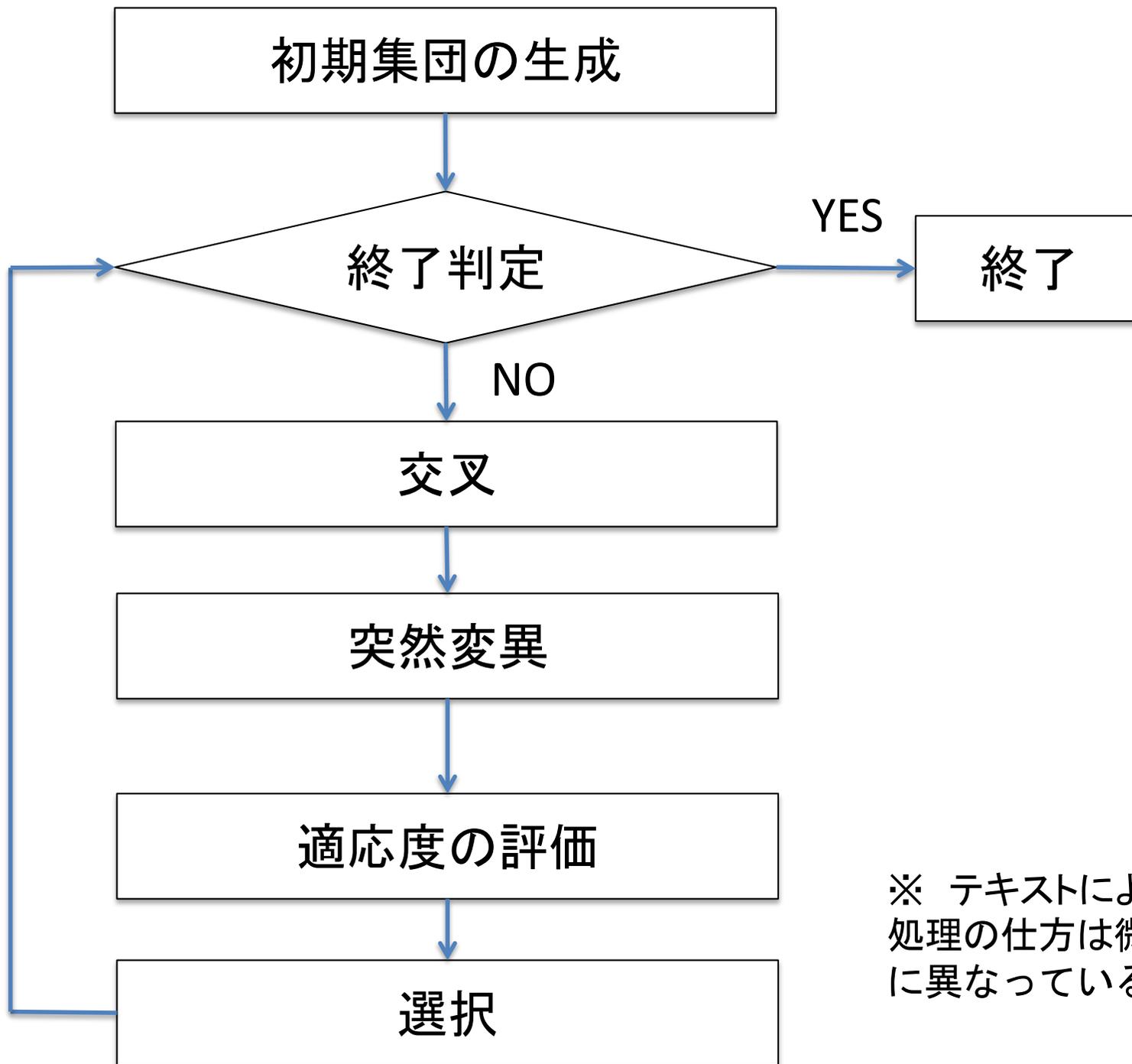
適応度の評価

選択

終了時処理

最終的に「現世代」の中で最も適応度の高い個体を「解」として出力する。

それまでの世代で生成された最も適応度の高い個体を「解」として出力する方法もある。



※ テキストによって
処理の仕方は微妙
に異なっている

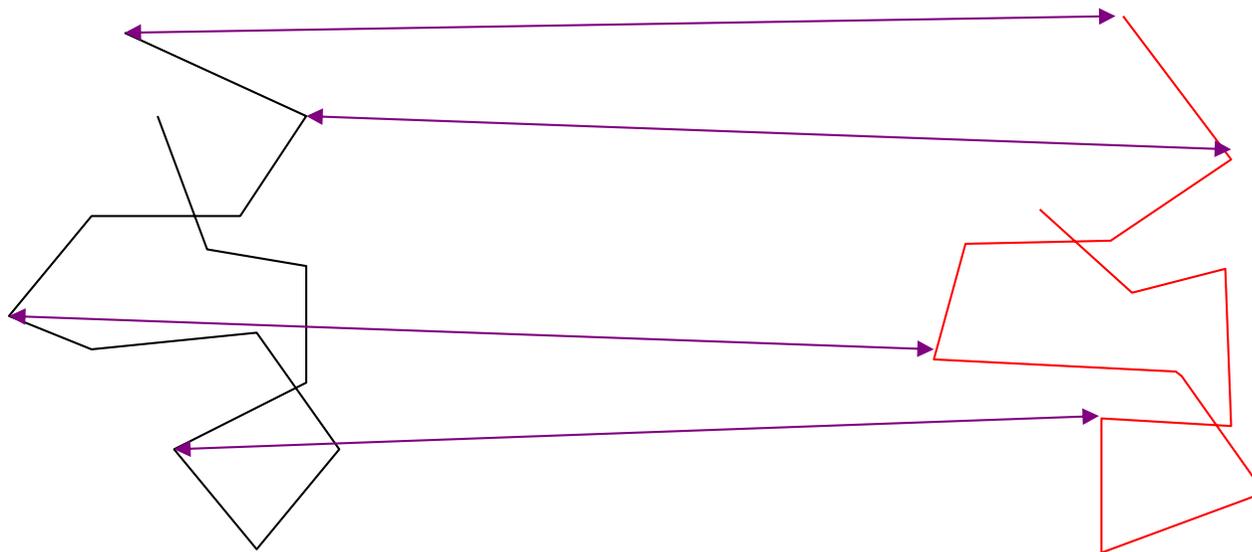
2. タンパク質の立体構造

2つの**相同**な**タンパク質**の立体構造の
重ね合わせを遺伝的アルゴリズムで行う

「相同」と「重ね合わせ」については後で説明

まず「タンパク質」について復習しよう

構造比較の原点 - 重ね合わせ (superposition) -

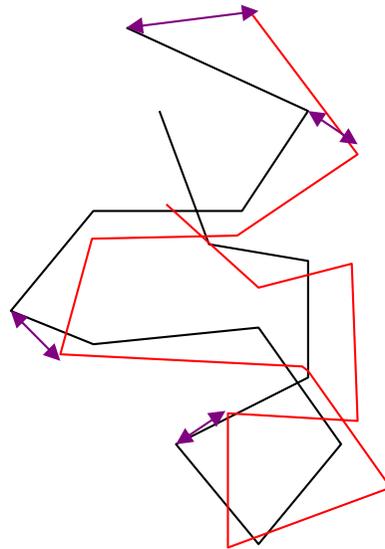


対応するC α 原子間距離が最小になるように
二つの鎖を重ね合わせる (平行移動と回転)

McLachlan, A.D. (1972) *Nature New Biol.* 240, 83-85.

RMSD

$$\text{rmsd (root mean square distance)} = \sqrt{\frac{1}{n} \sum (\text{dist}(A(i), B(i)))^2}$$



残基間対応が最初に
与えられていると
計算は容易

3.1.2 立体構造の計測法とそのデータ形式

生体高分子の立体構造データは、主に、X線結晶解析、核磁気共鳴法（NMR）、電子顕微鏡の三つの方法で決められている（参考文献 1, 2）。各手法の特徴を表 3.1 にまとめた。生成された立体構造データの登録受付と配布は、米国の RCSB-PDB、欧州の PDBe、日本の PDBj（大阪大学の蛋白質研究所）の三つの拠点による共同組織 wwPDB（worldwide Protein Data Bank）が行っている。

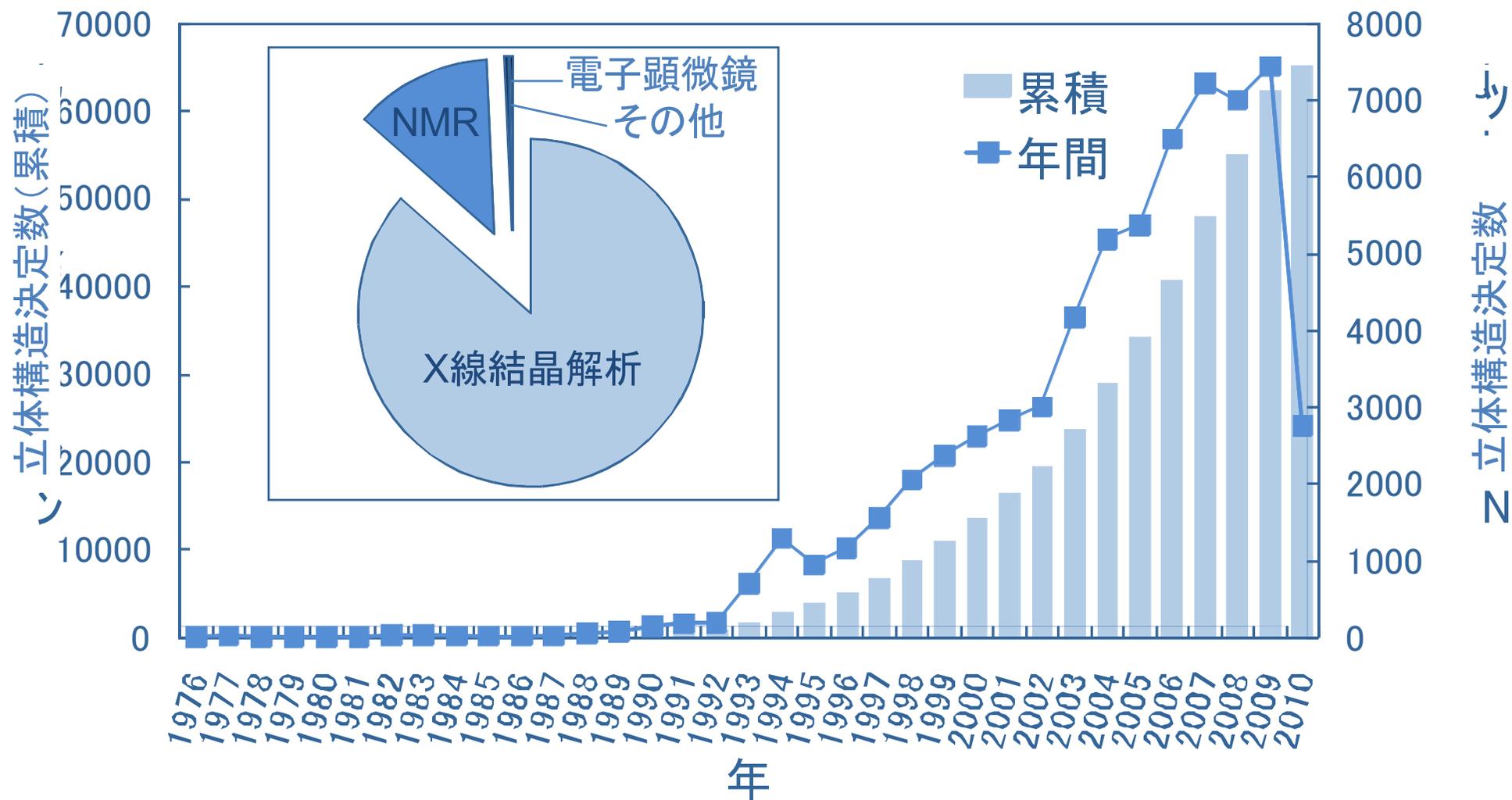
表 3.1 生体高分子立体構造を決定する実験手法

実験手法	X線結晶解析	核磁気共鳴（NMR）	電子顕微鏡
PDB データベースに占める割合*	89.5%	8.7%	1.6%
試料	結晶状態	同位体標識された高濃度の溶液状態	急速冷却によりグリッドに氷包埋された状態
データ測定	放射光などによる X 線回折実験	NMR 装置によるシグナル測定	電子顕微鏡による 2D 画像群の撮影
計算過程	フーリエ逆変換・位相解決・原子モデルの構築	シグナル帰属、原子間距離の推定、原子モデルの構築	単粒子解析やトモグラフィによる 3D マップ再構築、原子モデルの構築
特徴	解像度は比較的高く、大きな分子量でも解析可能。解像度が特に高い場合を除き、水素原子は観測できない。	分子量の小さな分子の解析に向く。帰属に用いるため、水素原子の座標も記載されている。複数の候補モデルが書かれていることが多い。	比較的大きな分子の解析に向く。現状では中程度の解像度（3-5 Å）のデータが多い。

* 2018 年 7 月 4 日の PDB の 141842 データの統計に基づく。

構造データの推移

- 立体構造情報の増加



生体高分子の立体構造

遺伝子

タンパク質

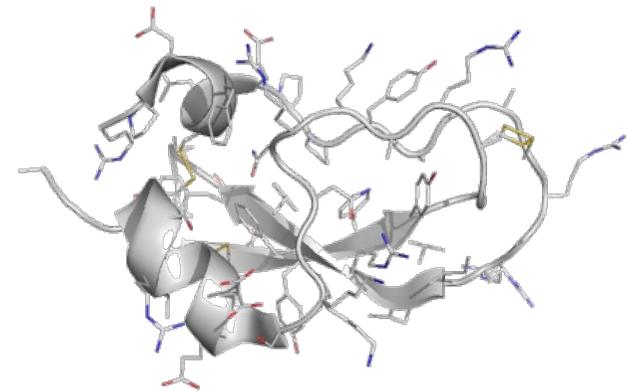
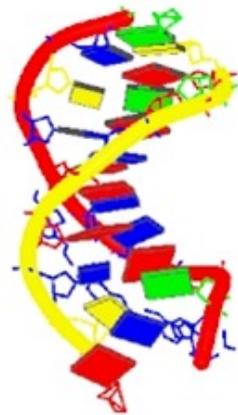
配列

>1AUL
gtctattagt
actaatagac

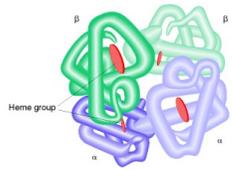
>1QLQ:A
RPDFCLEPPYAGACRARIIRYFYNA
KAGLCQTFVYGGCRAKRNNFKSAED
CLRTC GGA

折りたたみ
(エネルギー的に安定)

立体構造

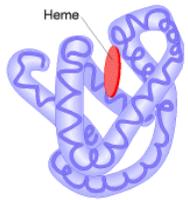


タンパク質立体構造の階層性



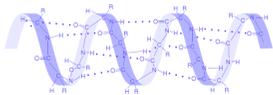
四次構造 →

立体構造をとる二本以上のペプチド鎖から形成される構造
(具体例:ヘモグロビンは4つのペプチド鎖が集まり機能 (α サブユニット $\times 2 + \beta$ サブユニット $\times 2$))

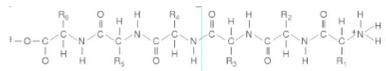


三次構造 →

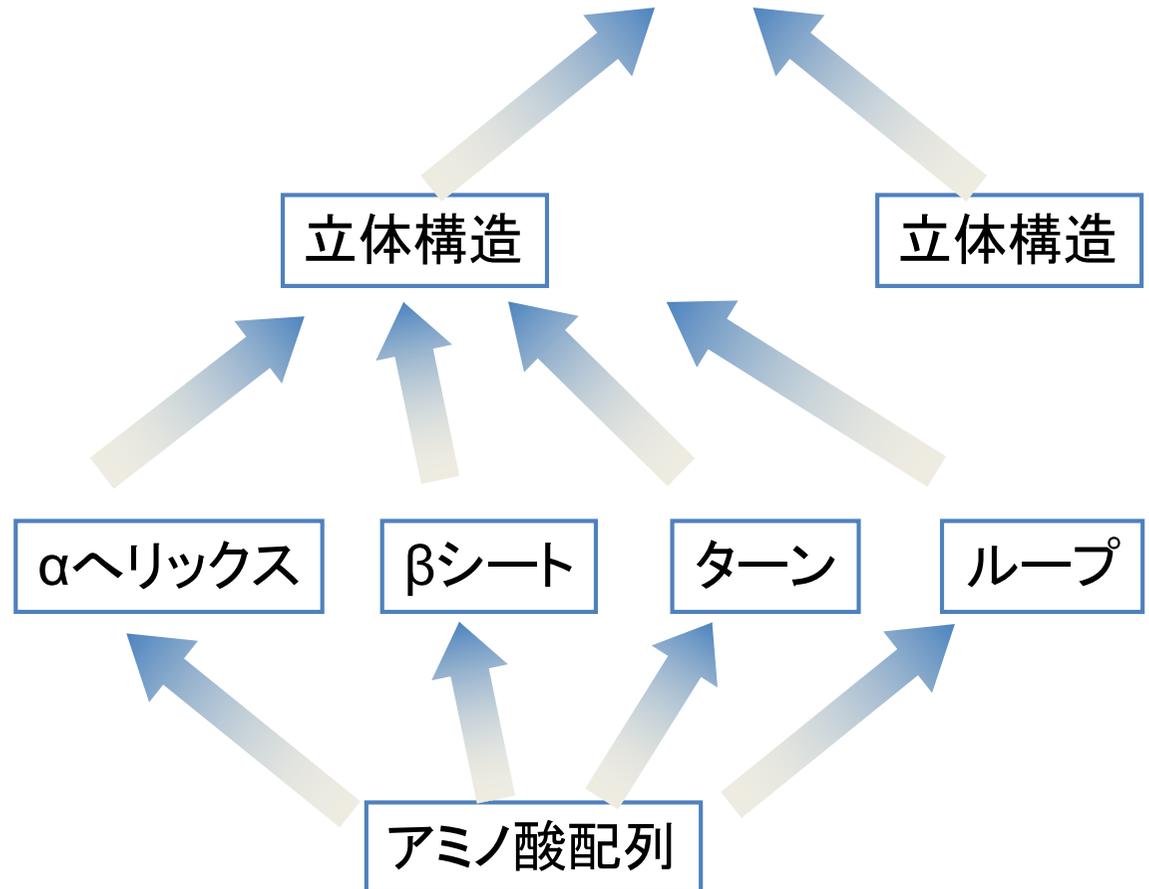
ドメイン構造 →



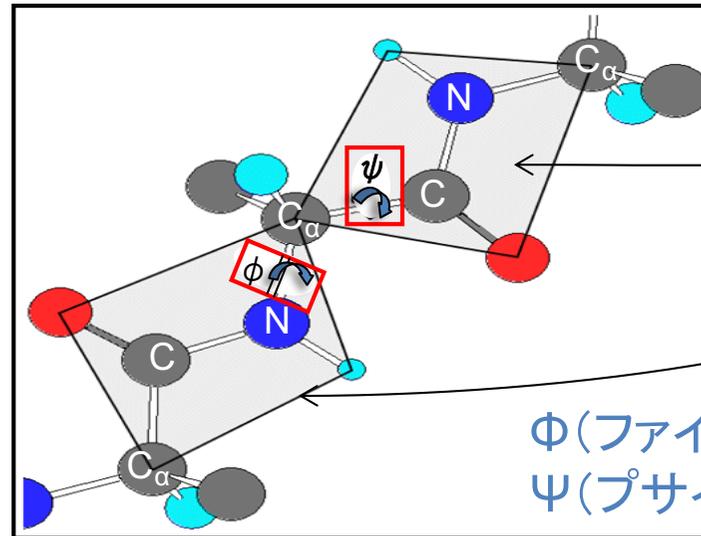
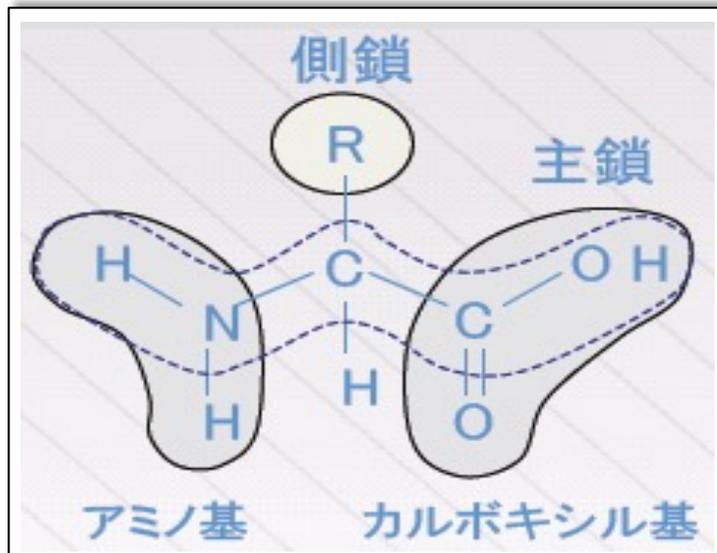
二次構造 →



一次構造 →



二面角（ねじれ角）

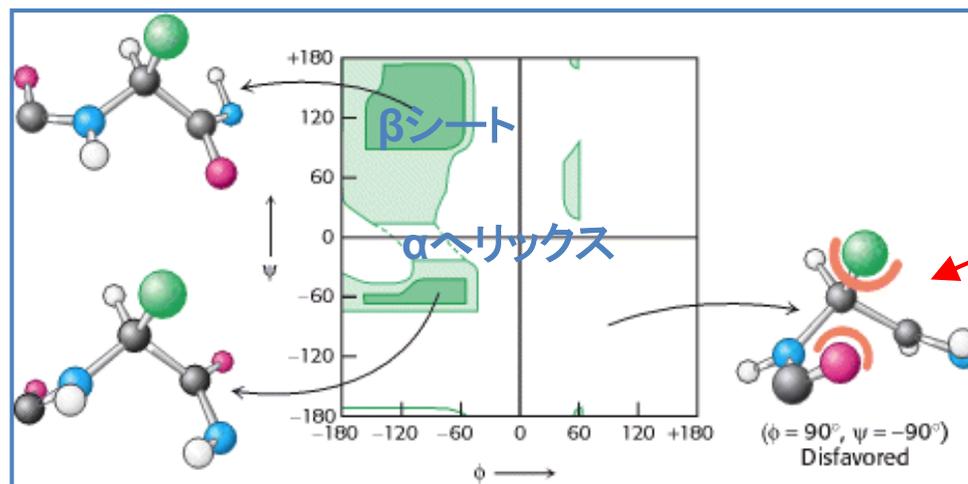


ペプチド結合は平面構造をとる

二面角の定義 = 2つの面が形成する角度

<http://www.ncbi.nlm.nih.gov/books/bv.fcgi?rid=bioinfo.figgrp.151>

ラマチャンドラプロット



立体障害がおきる

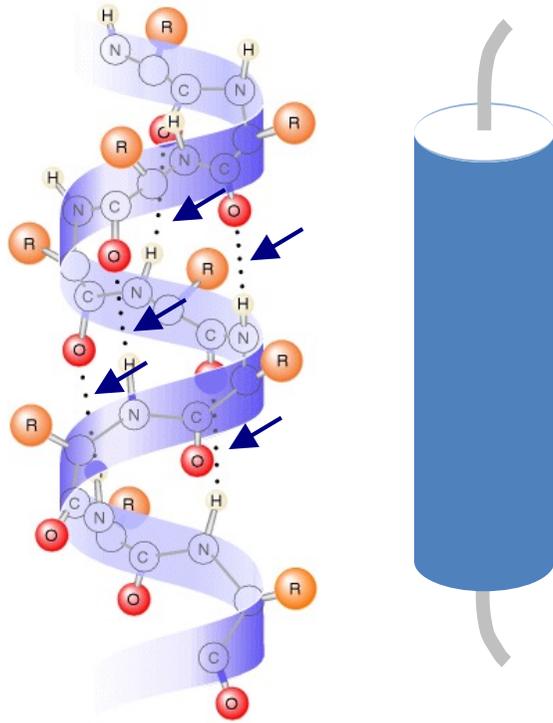
二面角がとる値には制限がある

<http://www.ncbi.nlm.nih.gov/books/bv.fcgi?rid=stryer.figgrp.319>

二次構造

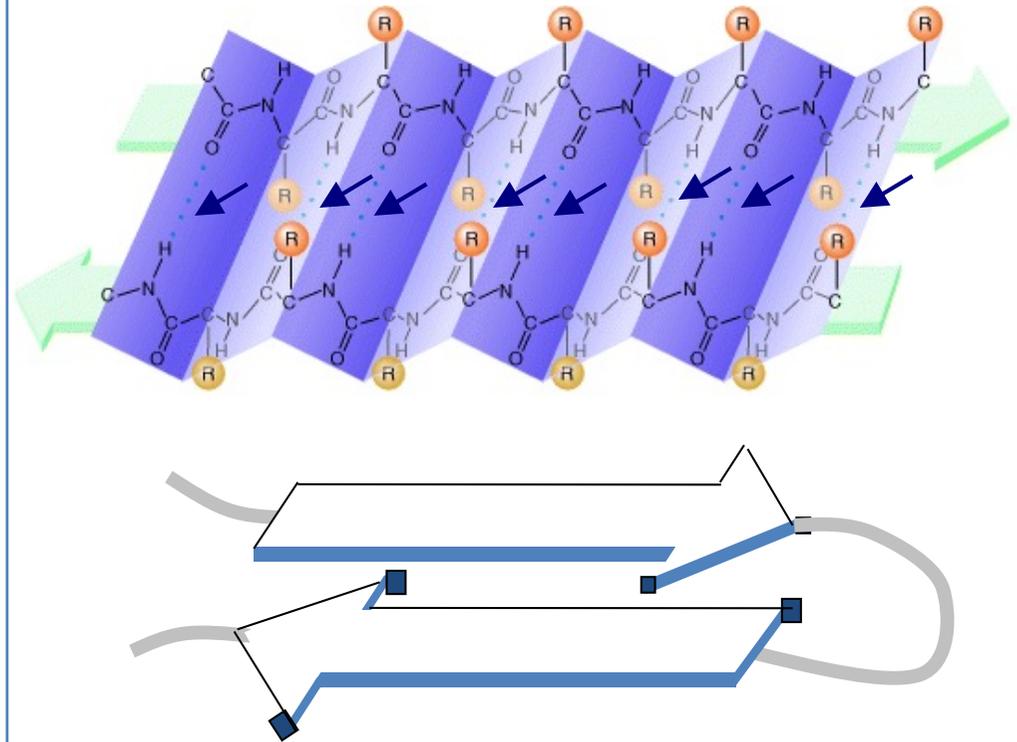
規則的な水素結合のパターンを持つ

αヘリックス



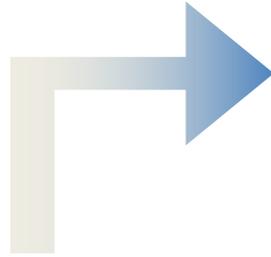
4残基先のアミノ酸と水素結合を作る

βシート



配列上離れた位置にアミノ酸と水素結合を作る

三次構造



内側：

疎水性残基が集まる

外側：

親水性残基が多い



独自の安定な構造をとる

ドメイン：

- ・比較的安定な構造をとる単位
- ・三次構造より
(たいていの場合)小さな構造
- ・ドメインは、様々な三次構造中に見られる

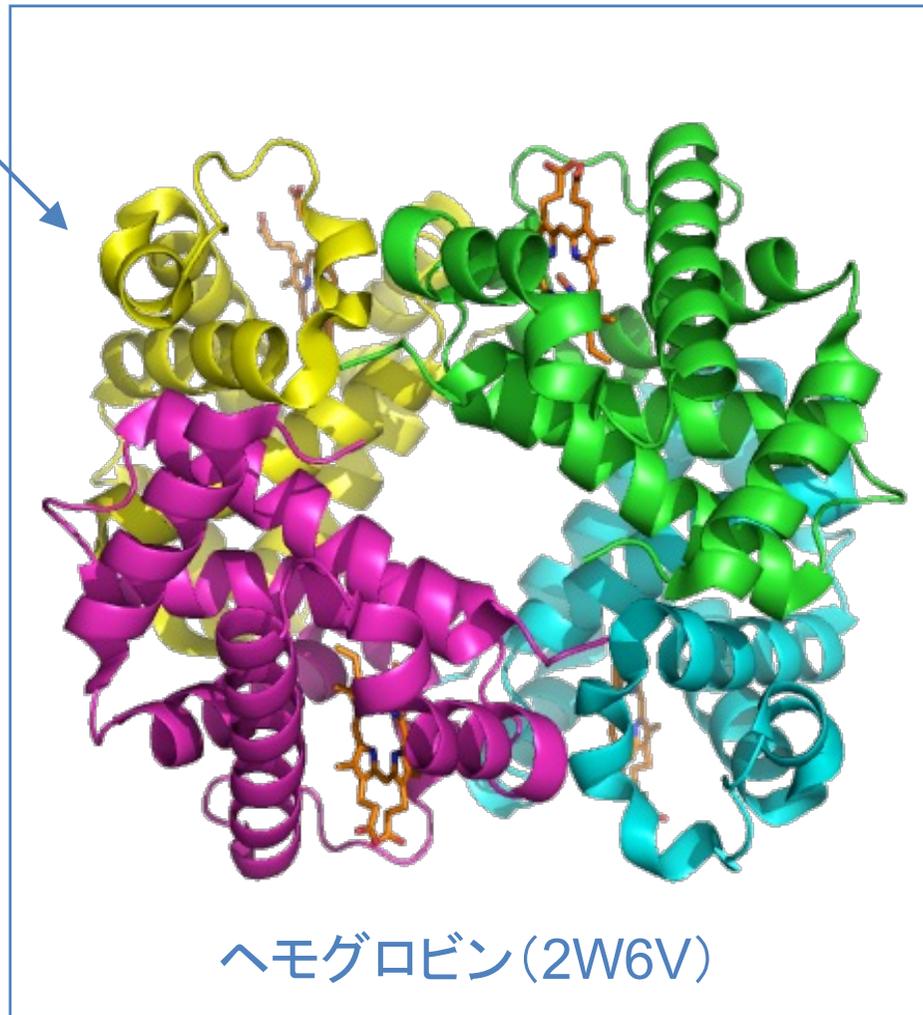
二次構造の組み合わせで、ある形に折りたたまれたもの。
一本のポリペプチドからなる。

立体構造保持に働く力

共有結合	ジスルフィド結合 (S-S結合)	アミノ酸配列中のシステイン残基が2個以上あり、それらが空間的に近い場合、 イオウ原子同士が形成 する結合。
	イオン結合	反対の電荷を持つ残基間の電気的な結合。陽イオンと陰イオンが クーロン力 で引き合う。
非共有結合	水素結合	タンパク質分子内のアミノ基の水素が他のタンパク質分子内の酸素と 電氣的に強く引きあって できる結合。
	ファンデルワールス力	あらゆる分子同士で働く、 互いに引き合う力 。
	疎水結合	疎水性のアミノ酸残基の側鎖が、水を嫌って 水と接触しないように集って生じる 結合。

四次構造

サブユニット



多量体の呼び方

- 1本 : モノマー
- 2本 : ダイマー
- 3本 : トリマー
- 4本 : テトラマー
- 5本 : ペンタマー
- 6本 : ヘキサマー

生体内で機能しているタンパク質の多くは、四次構造を形成している

複数のポリペプチドによって形成される立体構造

立体構造情報（座標データとして表現）

```
ATOM      3  C   MET A   1      26.913  26.639  3.531  1.00  9.62      C
ATOM      4  O   MET A   1      27.886  26.463  4.263  1.00  9.62      O
ATOM      5  CB  MET A   1      25.112  24.880  3.649  1.00 13.77      C
ATOM      6  CG  MET A   1      25.353  24.860  5.134  1.00 16.29      C
ATOM      7  SD  MET A   1      23.930  23.959  5.904  1.00 17.17      S
ATOM      8  CE  MET A   1      24.447  23.984  7.620  1.00 16.11      C
ATOM      9  N   GLN A   2      26.335  27.770  3.258  1.00  9.27      N
ATOM     10  CA  GLN A   2      26.850  29.021  3.898  1.00  9.07      C
ATOM     11  C   GLN A   2      26.100  29.253  5.202  1.00  8.72      C
ATOM     12  O   GLN A   2      24.865  29.024  5.330  1.00  8.22      O
ATOM     13  CB  GLN A   2      26.733  30.148  2.905  1.00 14.46      C
ATOM     14  CG  GLN A   2      26.882  31.546  3.409  1.00 17.01      C
ATOM     15  CD  GLN A   2      26.786  32.562  2.270  1.00 20.10      C
ATOM     16  OE1  GLN A   2      27.783  33.160  1.870  1.00 21.89      O
ATOM     17  NE2  GLN A   2      25.562  32.733  1.806  1.00 19.49      N
ATOM     18  N   ILE A   3      26.849  29.656  6.217  1.00  5.87      N
ATOM     19  CA  ILE A   3      26.235  30.058  7.497  1.00  5.07      C
ATOM     20  C   ILE A   3      26.882  31.428  7.862  1.00  4.01      C
ATOM     21  O   ILE A   3      27.906  31.711  7.264  1.00  4.61      O
ATOM     22  CB  ILE A   3      26.344  29.050  8.645  1.00  6.55      C
ATOM     23  CG1  ILE A   3      27.810  28.748  8.999  1.00  4.72      C
ATOM     24  CG2  ILE A   3      25.491  27.771  8.287  1.00  5.58      C
ATOM     25  CD1  ILE A   3      27.967  28.087 10.417  1.00 10.83      C
```

座標情報

>2W36:F

cgatctgtagc (塩基配列)

>2W36:B

MDYRQLHRWDLPPAEAIVQNELRKKIKLTPYEPEYVAGVALSFPGKEEGLAVIVVLEYPFKILEVVSERGEITFPYIP
GLLAFREGPLFLKAWEKLRTPDVVVFDGQGLAHPKLG IASHMGLFIEIPTIGVAKSRLYGTFKMPEDKRCSSWSYLYDGEE
IIGCVIRTKEGSAPIFVSPGHLMDVSSKRLIKAFTLPGRRIPEPTRLAHIYTQRLKKGLF (アミノ酸配列)

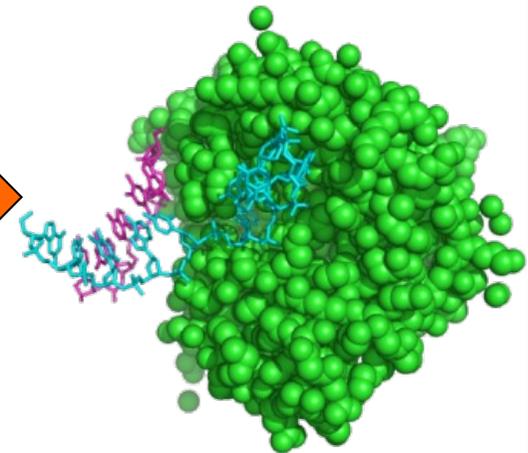
文字列情報

立体構造情報（座標データとして表現）

```
ATOM 3 C MET A 1 26.913 26.639 3.531 1.00 9.62 C
ATOM 4 O MET A 1 27.886 26.463 4.263 1.00 9.62 O
ATOM 5 CB MET A 1 25.112 24.880 3.649 1.00 13.77 C
ATOM 6 CG MET A 1 25.353 24.860 5.134 1.00 16.29 C
ATOM 7 SD MET A 1 23.930 23.959 5.904 1.00 17.17 S
ATOM 8 CE MET A 1 24.447 23.984 7.620 1.00 16.11 C
ATOM 9 N GLN A 2 26.335 27.770 3.258 1.00 9.27 N
ATOM 10 CA GLN A 2 26.850 29.021 3.898 1.00 9.07 C
ATOM 11 C GLN A 2 26.100 29.253 5.202 1.00 8.72 C
ATOM 12 O GLN A 2 24.865 29.024 5.330 1.00 8.22 O
ATOM 13 CB GLN A 2 26.733 30.148 2.905 1.00 14.46 C
ATOM 14 CG GLN A 2 26.882 31.546 3.409 1.00 17.01 C
ATOM 15 CD GLN A 2 26.786 32.562 2.270 1.00 20.10 C
ATOM 16 OE1 GLN A 2 27.783 33.160 1.870 1.00 21.89 O
ATOM 17 NE2 GLN A 2 25.562 32.733 1.806 1.00 19.49 N
ATOM 18 N ILE A 3 26.849 29.656 6.217 1.00 5.87 N
ATOM 19 CA ILE A 3 26.235 30.058 7.497 1.00 5.07 C
ATOM 20 C ILE A 3 26.882 31.428 7.862 1.00 4.01 C
ATOM 21 O ILE A 3 27.906 31.711 7.264 1.00 4.61 O
ATOM 22 CB ILE A 3 26.344 29.050 8.645 1.00 6.55 C
ATOM 23 CG1 ILE A 3 27.810 28.748 8.999 1.00 4.72 C
ATOM 24 CG2 ILE A 3 25.491 27.771 8.287 1.00 5.58 C
ATOM 25 CD1 ILE A 3 27.967 28.087 10.417 1.00 10.83 C
```

座標情報

可視化ソフト



立体構造
2w36(pdb)

>2W36:F

cgatctgtagc (塩基配列)

>2W36:B

MDYRQLHRWDLPPAEAIVQNELRKKIKLTPYEPEYVAGVALSFPGKEEGLAVIVVLEYPFKILEVVSERGEITFPYIP
GLLAFREGPLFLKAWEKLRTPDVVVFDGQGLAHPKLG IASHMGLFIEIPTIGVAKSRLYGTFKMPEDKRCSSWSYLYDGEE
IIGCVIRTKEGSAPIFVSPGHLMDVSSKRLIKAFTLPGRRIPEPTRLAHIYTQRLKKGLF (アミノ酸配列)

文字列情報

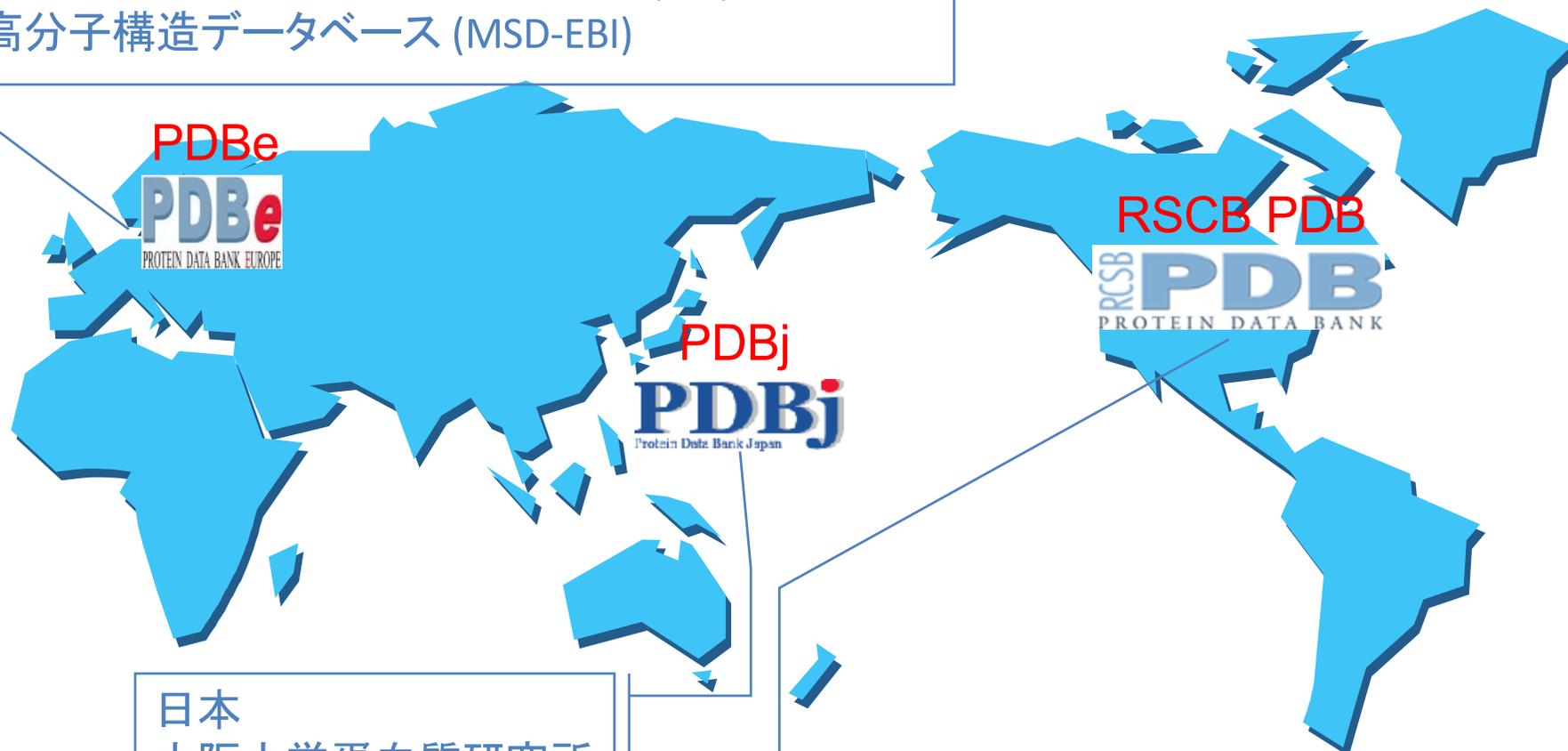
生体高分子の立体構造を扱ったデータベース

- 網羅的に立体構造情報を収集しているもの

ヨーロッパ

欧州バイオインフォマティクス研究所 (EBI)

高分子構造データベース (MSD-EBI)



日本

大阪大学蛋白質研究所

アメリカ

構造バイオインフォマティクス
研究共同体 (RCSB)

生体高分子の立体構造を扱ったデータベース

▶ 網羅的に立体構造情報を収集しているもの

The image displays three screenshots of the Protein Data Bank (PDB) website, illustrating its comprehensive collection of biological macromolecular structures.

Top Left Screenshot (PDBe): Shows the EBI Protein Data Bank in Europe (PDBe) homepage. It features a navigation menu with links to Databases, Tools, EBI Groups, Training, Industry, About Us, and Help. The main content area includes a welcome message, a link to the wwPDB Statement on Retraction of UAB PDB Entries, and a section for the Nobel Prize for Chemistry 2009. A search bar is visible at the top.

Top Right Screenshot (RCSB PDB): Shows the RCSB PDB homepage. It features a search bar and a navigation menu. The main content area includes a featured article titled "A Resource for Studying Biological Macromolecules" and a section for "Molecule of the Month: Concanavalin A and Circular Permutation".

Bottom Screenshot (PDB Japan): Shows the PDB Japan homepage. It features a search bar and a navigation menu. The main content area includes a search bar and a section for "64932 entries available on 28 Apr., 2010".

アメリカ
構造バイオインフォマティクス
研究共同体 (RCSB)

生体高分子の立体構造を扱ったデータベース

• PDBファイルの中身

```
HEADER      TRANSFERASE                               30-AUG-04   1XBB
TITLE      CRYSTAL STRUCTURE OF THE SYK TYROSINE KINASE DOMAIN WITH
TITLE      2 GLEEVEC
COMPND     MOL_ID: 1;
COMPND     2 MOLECULE: TYROSINE-PROTEIN KINASE SYK;
COMPND     3 CHAIN: A;
COMPND     4 SYNONYM: SPLEEN TYROSINE KINASE;
COMPND     5 EC: 2.7.1.112;
COMPND     6 ENGINEERED: YES
SOURCE     MOL_ID: 1;
SOURCE     2 ORGANISM_SCIENTIFIC: HOMO SAPIENS;
SOURCE     3 ORGANISM_COMMON: HUMAN;
SOURCE     4 ORGANISM_TAXID: 9606;
SOURCE     5 GENE: SYK;
SOURCE     6 EXPRESSION_SYSTEM: SPODOPTERA FRUGIPERDA;
SOURCE     7 EXPRESSION_SYSTEM_COMMON: FALL ARMYWORM;
SOURCE     8 EXPRESSION_SYSTEM_TAXID: 7108;
SOURCE     9 EXPRESSION_SYSTEM_CELL_LINE: SF9;
SOURCE    10 EXPRESSION_SYSTEM_VECTOR_TYPE: BACULOVIRUS
KEYWDS     GLEEVEC, STI-571, IMATINIB, SYK, SPLEEN TYROSINE KINASE,
KEYWDS     2 ACTIVE CONFORMATION, STRUCTURAL GENOMICS, STRUCTURAL
KEYWDS     3 GENOMIX, TRANSFERASE
EXPDTA     X-RAY DIFFRACTION
AUTHOR     V. L. NIENABER, S. ATWELL, J. M. ADAMS, J. BADGER, M. D. BUCHANAN,
AUTHOR     2 I. K. FEIL, K. J. FRONING, X. GAO, J. HENDLE, K. KEEGAN, B. C. LEON,
AUTHOR     3 H. J. MULLER-DEICKMANN, B. W. NOLAND, K. POST, K. R. RAJASHANKAR,
AUTHOR     4 A. RAMOS, M. RUSSELL, S. K. BURLEY, S. G. BUCHANAN
```

(省略)

生体高分子の立体構造を扱ったデータベース

• PDBファイルの中身

		ヘッダー	ボディ
HEADER	TRANSFERASE	HEADER	物質の分類名、日付、ファイル名
TITLE	CRYSTAL STRUCTURE	COMPND	タンパク質名
TITLE	2 GLEEVEC	SOURCE	物質が由来する生物名
COMPND	MOL_ID: 1;	AUTHOR	座標を作成した著者名
COMPND	2 MOLECULE: TYROSINASE	REVDAT	データが登録された日
COMPND	3 CHAIN: A;	JRNL	文献情報
COMPND	4 SYNONYM: SPL	REMARK	実験条件などに関する情報
COMPND	5 EC: 2.7.1.112	SEQRES	アミノ酸配列
COMPND	6 ENGINEERED: YES	HET	アミノ酸以外の原子
SOURCE	MOL_ID: 1;	HELIX	α -ヘリックスに関する情報
SOURCE	2 ORGANISM_SCIENTIFIC	SHEET	β -シートに関する情報
SOURCE	3 ORGANISM_COMMON	TURN	β -ターンに関する情報
SOURCE	4 ORGANISM_TAXID	SSBOND	ジスルフィド結合に関する情報
SOURCE	5 GENE: SYK;	ATOM	個々のアミノ酸の原子の座標
SOURCE	6 EXPRESSION_SYSTEM	HETATM	アミノ酸以外の原子の座標
SOURCE	7 EXPRESSION_SYSTEM	CONTACT	原子の結合に関する情報
SOURCE	8 EXPRESSION_SYSTEM	END	ファイルの終了
SOURCE	9 EXPRESSION_SYSTEM		
SOURCE	10 EXPRESSION_SYSTEM		
KEYWDS	GLEEVEC, STI-571		
KEYWDS	2 ACTIVE CONFORMATION		
KEYWDS	3 GENOMIX, TRANSDOMAIN		
EXPDTA	X-RAY DIFFRACTION		
AUTHOR	V. L. NIENABER, S. J. SWEENEY		
AUTHOR	2 I. K. FEIL, K. J. HAN		
AUTHOR	3 H. J. MULLER-DEBIL		
AUTHOR	4 A. RAMOS, M. RUSSELL		

(省略)

生体高分子の立体構造を扱ったデータベース

- 座標データ

1原子

ATOM	1	N	VAL	A	363	22.741	-1.397	11.729	1.00	33.32	N
ATOM	2	CA	VAL	A	363	21.557	-0.831	11.024	1.00	32.13	C
ATOM	3	C	VAL	A	363	20.954	-1.757	9.943	1.00	31.73	C
ATOM	4	O	VAL	A	363	19.737	-1.906	9.845	1.00	30.94	O
ATOM	5	CB	VAL	A	363	21.883	0.552	10.391	1.00	33.45	C
ATOM	6	N	TYR	A	364	21.798	-2.389	9.135	1.00	29.77	N
ATOM	7	CA	TYR	A	364	21.310	-3.035	7.928	1.00	27.96	C
ATOM	8	C	TYR	A	364	20.929	-4.485	8.157	1.00	26.43	C
ATOM	9	O	TYR	A	364	21.735	-5.274	8.650	1.00	29.31	O
ATOM	10	CB	TYR	A	364	22.349	-2.901	6.808	1.00	28.13	C
ATOM	11	CG	TYR	A	364	22.619	-1.461	6.442	1.00	28.30	C
ATOM	12	CD1	TYR	A	364	21.725	-0.751	5.638	1.00	30.33	C
ATOM	13	CD2	TYR	A	364	23.772	-0.815	6.893	1.00	30.80	C
ATOM	14	CE1	TYR	A	364	21.959	0.588	5.285	1.00	31.59	C
ATOM	15	CE2	TYR	A	364	24.019	0.523	6.541	1.00	32.22	C
ATOM	16	CZ	TYR	A	364	23.101	1.208	5.733	1.00	32.48	C
ATOM	17	OH	TYR	A	364	23.330	2.524	5.389	1.00	36.29	O

1残基

1残基

原子の名前

アミノ酸の名前

チェーンの名前

アミノ酸残基番号

y座標

x座標

z座標

温度因子

占有率

原子

次に、立体構造データがどういう形式でデータベースに格納されるのかを説明する。まず、各立体構造データには、“1mbd”“4hhb”といった「1文字の数字 + 3文字の英数字」の **PDB ID コード** がついている。各データは、立体構造が計測された一つのまとまりが単位であり、一つのタンパク質鎖が必ずしも一つの PDB ID に対応するわけではない。タンパク質の一部だけが入っている場合、複数のタンパク質が含まれている場合もある。 DNA や RNA などの高分子や、結合している低分子化合物の立体構造も含まれている。図 3.1a にヒト・ヘモグロビンの立体構造 (PDB ID: 1bzl) をリボンモデルで示した。このデータには4つのタンパク質鎖と4つのヘム分子が入っている。図 3.1b には最初の三つのアミノ酸の原子構造を示した。この立体構造に対応するデータを、**旧 PDB フォーマット** (図 3.2) と **mmCIF フォーマット** (図 3.3) の二つの形式で示した。どちらも、立体構造は、1原子が1行に記載され、原子の中心点の XYZ 座標の値が書かれている。座標値の単位は Å (オングストローム) である ($1 \text{ \AA} = 0.1 \text{ nm} = 10^{-10} \text{ m}$)。また、各原子には原子番号、残基番号、3文字以下の残基名 (例えば、MET, ALA) と4文字以下の原子名 (CA, N, C, O, CB など) が割り当てられている。3文字の残基名は、アミノ酸以外の分子にも割り当てられている。例えば、ヘム分子は“HEM”, ATP 分子は“ATP”, カルシウムイオンは“CA”, イレッサは“IRE”などと決められている。一つのデータに複数の分子が入っている場合には、それぞれ、A, B, C といった鎖識別子 (chain identifier) がついている。

旧 PDB フォーマットと 2014 年から標準フォーマットとなった mmCIF フォーマットの違いについて簡単に説明しよう。1970 年代から使われていた旧 PDB フォーマットは、わかりやすい反面、1 行 80 文字の固定幅であるため、原子番号、残基番号、鎖識別子の文字数に制限が設定され、記載できる原子数、鎖数に明らかな上限がある。例えば、鎖識別子は 1 文字であるため、小文字や数字を使っても、タンパク質鎖数が 100 を超える巨大分子は記載できないことになる。そこで導入されたのが、mmCIF フォーマットである。このフォーマットは、スペースで区切られた形式を基本としているため、原子番号や鎖識別子の文字数に制限はない。また、登録者が入力する残基番号などの情報とデータベース管理者が入力する情報の両方を記載することで利便性と統一性を実現している。現在 PDB では、大部分のデータでは旧 PDB フォーマットでもダウンロードできるが、一部の巨大分子の立体構造（例えば、リボソーム：4v42、光化学系 II：4v62 など）は旧 PDB フォーマットでは記載できないため、mmCIF フォーマットだけで公開されている。

data_1B21

```

#
#_entry.id 1B21
#
#_pdbx_database_status.recvd_initial_deposition_date 1998-11-05
#
#_struct.title 'HEMOGLOBIN (ALPHA + MET) VARIANT'
#_struct.pdbx_descriptor 'PROTEIN (HEMOGLOBIN) VARIANT (CHAIN A, C, ADDITIONAL NH2-TERMINAL MET)'
```

キー・バリュー形式

```

loop_
  _atom_site.group_PDB
  _atom_site.id
  _atom_site.type_symbol
  _atom_site.label_atom_id
  _atom_site.label_alt_id
  _atom_site.label_comp_id
  _atom_site.label_asym_id
  _atom_site.label_entity_id
  _atom_site.label_seq_id
  _atom_site.pdbx_PDB_ins_code
  _atom_site.Cartn_x
  _atom_site.Cartn_y
  _atom_site.Cartn_z
  _atom_site.occupancy
  _atom_site.B_iso_or_equiv
  _atom_site.pdbx_formal_charge
  _atom_site.auth_seq_id
  _atom_site.auth_comp_id
  _atom_site.auth_asym_id
  _atom_site.auth_atom_id
  _atom_site.pdbx_PDB_model_num
```

表形式

データベース管理者が入力する

原子座標

登録者が入力する

ATOM	1	N	N	.	MET	A	1	1	?	15.774	28.408	41.946	1.00	87.06	?	1	MET	A	N	1
ATOM	2	C	CA	.	MET	A	1	1	?	17.105	28.442	42.578	1.00	83.35	?	1	MET	A	CA	1
ATOM	3	C	C	.	MET	A	1	1	?	18.021	29.477	41.921	1.00	75.77	?	1	MET	A	C	1
ATOM	4	O	O	.	MET	A	1	1	?	17.695	30.170	40.943	1.00	76.91	?	1	MET	A	O	1
ATOM	5	C	CB	.	MET	A	1	1	?	17.757	27.078	42.696	1.00	100.06	?	1	MET	A	CB	1
ATOM	6	C	CG	.	MET	A	1	1	?	18.477	26.602	41.467	1.00	108.30	?	1	MET	A	CG	1
ATOM	7	S	SD	.	MET	A	1	1	?	19.741	25.379	41.972	1.00	109.12	?	1	MET	A	SD	1
ATOM	8	C	CE	.	MET	A	1	1	?	18.713	24.184	42.861	1.00	108.71	?	1	MET	A	CE	1
ATOM	9	N	N	.	VAL	A	1	2	?	19.204	29.547	42.527	1.00	63.43	?	2	VAL	A	N	1
ATOM	10	C	CA	.	VAL	A	1	2	?	20.227	30.471	42.079	1.00	50.48	?	2	VAL	A	CA	1
:																				
HETATM	4401	C	CHA	.	HEM	E	3	.	?	18.729	18.645	20.255	1.00	12.22	?	143	HEM	A	CHA	1
HETATM	4402	C	CHB	.	HEM	E	3	.	?	21.007	20.618	24.060	1.00	16.53	?	143	HEM	A	CHB	1
HETATM	4403	C	CHC	.	HEM	E	3	.	?	18.570	17.704	27.086	1.00	11.28	?	143	HEM	A	CHC	1

```

# 原子番号 原子名 残基名 残基番号 X座標 Y座標 Z座標 占有率 温度因子 残基番号 asym_id(鎖識別子)
# 元素名 asym_id 残基名 原子名
```

図 3.3 mmCIF フォーマットによる立体構造データの例

正式には PDBx/mmCIF フォーマットという。PDB ID: 1bz1 の、いくつかの項目と、原子座標の一部を示す。このフォーマットは、項目名と値のペアを記述するシンプルな「キー・バリュー形式」(オレンジの点線枠内)と、項目名を最初に示しデータが行として多数並ぶ「表形式」(青色の点線枠内)のどちらかで記述されている。値の区切り文字はスペースで、値の長さは可変である。項目名はカテゴリ名とアイテム名がピリオド(.)で組み合わせられている。例えば、_atom_site.Cartn_x では、atom_site がカテゴリ名、Cartn_x がアイテム名である。表形式の場合、カテゴリ名は表の名前、アイテム名は列の名前に相当する。atom_site の表が、旧 PDB フォーマットの ATOM 行と HETATM 行に相当する。残基番号 (seq_id)、残基名 (comp_id)、鎖識別子 (asym_id)、原子名 (atom_id) は、データベース管理者が入力する label_ で始まる項目(緑色)と登録者が入力する auth_ で始まる項目(赤色)の両方が入っている。これらはまったく同一である場合も多いが、異なることも少なくない。この例では、HETATM 行の HEM の asym_id は、データベース入力の label_asym_id は E だが、登録者入力の auth_asym_id は A となっている。データベース管理者は、このヘム分子は 5 つ目の分子であるから E としたが、登録者はタンパク質の A 鎖に結合している分子だから A としたのだと考えられる。

演習

- ヒトの“チロシンキナーゼ (SYK: spleen tyrosine kinase)”の座標データを入手する

立体構造情報を手に入れる

RCSB PDBの トップページ

The image shows a screenshot of the RCSB PDB homepage with three red callout boxes providing instructions for searching. The browser address bar shows the URL <https://www.rcsb.org/>. The search bar contains the text "spleen tyrosine kinase" and a "Go" button. The main content area features a "Welcome" message, a "Deposit" button, a "Search" button, and a "Video: How Enzymes Work" section. A "November Molecule" section highlights "Phospholipase A2". The bottom of the page includes "Latest Entries" (featuring "6HNN"), "Features & Highlights" (with a "doi" logo), and "News" (with a "New EM map validation in OneDep" article).

① <https://www.rcsb.org/>と入力する

② キーワードに"spleen tyrosine kinase"と入力

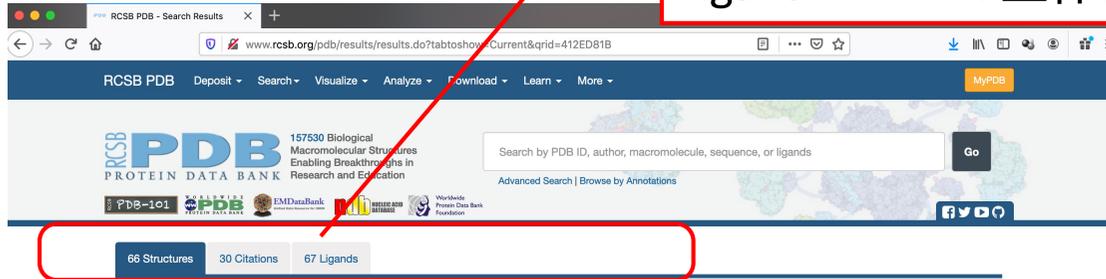
③ "Go"ボタンをクリックする

立体構造情報を手に入れる

検索結果画面(トップ)

- Structures : 立体構造をリスト
- Citations : 立体構造に関する文献をリスト
- Ligand : 立体構造中に含まれるリガンドをリスト

66構造ヒット



Structure タブがデフォルトで示される

Search Parameter:
Text Search for: syk tyrosine kinase domain

- Refinements
- ORGANISM
 - Homo sapiens (65)
 - Mus musculus (2)
 - Bos taurus (1)
 - UNIPROT MOLECULE NAME
 - Tyrosine-protein kinase SYK (57)
 - Tyrosine-protein kinase Z ... (8)
 - E3 ubiquitin-protein liga ... (3)
 - Ubiquitin-conjugating enz ... (2)
 - Proto-oncogene vav (2)
 - T-cell surface glycoprote ... (2)
 - T-cell surface glycoprote ... (1)
 - TAXONOMY
 - Eukaryota only (66)
 - EXPERIMENTAL METHOD
 - X-ray (61)
 - Solution NMR (5)
 - X-RAY RESOLUTION
 - less than 1.5 Å (6)
 - 1.5 - 2.0 Å (21)
 - 2.0 - 2.5 Å (22)
 - 2.5 - 3.0 Å (10)
 - 3.0 and more Å (2)
 - RELEASE DATE
 - before 2000 (3)
 - 2000 - 2005 (6)
 - 2005 - 2010 (7)
 - 2010 - 2015 (22)
 - 2015 - today (29)

1エントリー

1XBA
Crystal structure of apo syk tyrosine kinase domain
Atwell, S., Adams, J.M., Badger, J., Buchanan, M.D., Felli, I.K., Froning, K.J., Gao, X., Hendle, J., Keegan, K., Leon, B.C., Muller-Deickmann, H.J., Nienaber, V.L., Noland, B.W., Post, K.W., Rajashankar, K.R., Ramos, A., Russell, M., Burley, S.K., Buchanan, S.G.
(2004) J Biol Chem 279 55827-55832
Released: 11/2/2004
Method: X-ray Diffraction
Resolution: 2.0 Å
Residue Count: 291
Macromolecule: Tyrosine-protein kinase SYK (protein)
Unique Ligands: --
Search term match score: 722.59

Matched fields in 1XBA.cif:
_citation.title: A novel mode of Gleevec binding is revealed by the structure of spleen tyrosine kinase.
_entity.pdbx_description: Tyrosine-protein kinase SYK
_entity.name.com.name: Spleen tyrosine kinase
_struct.title: Crystal structure of apo syk tyrosine kinase domain
_struct.keywords.text: spleen tyrosine kinase , active conformation, structural genomics, Structural GenomIX, Transferase

1XBB
Crystal structure of the syk tyrosine kinase domain with Gleevec
Atwell, S., Adams, J.M., Badger, J., Buchanan, M.D., Felli, I.K., Froning, K.J., Gao, X., Hendle, J., Keegan, K., Leon, B.C., Muller-Deickmann, H.J., Nienaber, V.L., Noland, B.W., Post, K.W., Rajashankar, K.R., Ramos, A., Russell, M., Burley, S.K., Buchanan, S.G.
(2004) J Biol Chem 279 55827-55832
Released: 11/2/2004
Method: X-ray Diffraction
Resolution: 1.57 Å
Residue Count: 291
Macromolecule: Tyrosine-protein kinase SYK (protein)
Unique Ligands: STI
Search term match score: 722.59

立体構造情報を手に入れる

The screenshot shows the RCSB PDB website interface. The search bar contains the text "Text Search for: syk tyrosine kinase domain". The "Citation" tab is selected, showing a list of 67 citations. The first citation is highlighted, showing the 3D structure, related structure (5TIU), and citation information.

3D Structure	Related Structure	Citation Information
	5TIU	Carboxamide Spleen Tyrosine Kinase (Syk) Inhibitors: Leveraging Ground State Interactions To Accelerate Optimization. Ellis J.M., Altman M.D., Cash B., Haidle A.M., Kubiak R.L., Maddess M.L., Yan Y., Northrup A.B. (2016) ACS Med Chem Lett 7: 1151-1155 Pubmed Article: 27994755
	5T68	Synthesis and optimization of furano[3,2-d]pyrimidines as selective spleen tyrosine kinase (Syk) inhibitors. Hoemann M., Wilson N., Argiadi M., Barsch D., Burchat A., Calderwood D., Clapham B., Cox P., Duignan D.B., Konopacki D., Somal G., Vasudevan A. (2016) Bioorg Med Chem Lett 26: 5562-5567 Pubmed Article: 27789138
	5LMA	Optimisation of a novel series of potent and orally bioavailable azanaphthyridine SYK inhibitors. Garton N.S., Barker M.D., Davis R.P., Douault C., Hooper-Greenhill E., Jones E., Lewis H.D., Liddle J., Lugo D., McCleary S., Preston A.G.S., Ramirez-Molina C., Neu M., Shipley T.J., Somers D.O., Watson R.J., Wilson D.M. (2016) Bioorg Med Chem Lett 26: 4606-4612 Pubmed Article: 27578246
	4XG2	Crystal structures of spleen tyrosine kinase in complex with novel inhibitors: structural insights for design of anticancer drugs. Lee S.J., Choi J.S., Han B.G., Kim H.S., Song H.J., Lee J., Nam S., Goh S.H., Kim J.H., Koh J.S., Lee B.I. (2016) FEBS J 283: 3613-3625 Pubmed Article: 27504936
	4WNM	A Novel Triazolopyridine-Based Spleen Tyrosine Kinase Inhibitor That Arrests Joint Inflammation. Ferguson G.D., Delgado M., Plantevin-Krenitsky V., Jensen-Pargakes K., Bates R.J., Torres S., Celeridad M., Brown H., Burnett K., Nadozny L., Tehrani L., Packard G., Pagarigan B., Haelewyn J., Jensen-Pargakes K., Bates R.J., Torres S., Celeridad M., Brown H., Burnett K., Nadozny L., Chamberlain P., LaBrain L., Xia W., Bennett B., Blease K. (2016) PLoS One 11: e0145705 Pubmed Article: 26756335

Citation タブ

The screenshot shows the RCSB PDB website interface. The search bar contains the text "Text Search for: syk tyrosine kinase domain". The "Ligand" tab is selected, showing a list of 67 ligands. The first ligand is highlighted, showing the ligand structure, ID/Formula/Name, and structures with specific ligands.

Ligand Structure	ID / Formula / Name	Structures with Specific Ligands
	057 C21 H22 N4 O2 N-(2-hydroxy-1,1-dimethylethyl)-1-methyl-3-(1H-pyrrolo[2,3-b]pyridin-2-yl)-1H-indole-5-carboxamide	1 PDB Structures contains 057 (3FQH ...)
	0JE C17 H21 F N6 O2 3-[5-(5-ethoxy-6-fluoro-1H-benzimidazol-2-yl)-1H-pyrazol-4-yl]-1,1-diethylurea	1 PDB Structures contains 0JE (3VF8 ...)
	0VE C22 H23 N5 O2 3-[8-(4-{[ethyl(2-hydroxyethyl)amino]imino}imidazo[1,2-a]pyrazin-5-yl]phenol	1 PDB Structures contains 0VE (4FYN ...)
	0VF C25 H24 N6 O5 S 4-[[[3S]-1-(7-[[3,4-dimethoxyphenyl]amino][1,3]thiazolo[5,4-d]pyrimidin-5-yl)pyrrolidin-3-yl]carbamoyl]benzoic acid	1 PDB Structures contains 0VF (4FYO ...)
	0VG C22 H22 N6 N-[6-[(2S)-2-methylpyrrolidin-1-yl]pyridin-2-yl]-6-phenylimidazo[1,2-b]pyridazin-8-amine	1 PDB Structures contains 0VG (4FZ8 ...)
	0VH C18 H25 N7 O 6-[[[1R,2S)-2-aminocyclohexyl]amino][4-[[6-(pyridin-2-yl)amino]pyridazine-3-carboxamide	1 PDB Structures contains 0VH (4FZ7 ...)
	0XF C24 H24 N6 O2	1 PDB Structures contains 0XF (4GFG ...)

Ligand タブ

Structure タブに戻す

立体構造情報を手に入れる

RCSB PDB - Search Results

www.rcsb.org/pdb/results/results.do?tabto=show=Current&qid=412ED81B

RCSB PDB Deposit Search Visualize Analyze Download Learn More

157530 Biological Macromolecular Structures Enabling Breakthroughs in Research and Education

Search by PDB ID, author, macromolecule, sequence, or ligands

66 Structures 30 Citations 67 Ligands

Search Parameter: Text Search for: syk tyrosine kinase domain

Refinements

ORGANISM: Homo sapiens (65), Mus musculus (2), Bos taurus (1)

UNIPROT MOLECULE NAME: Tyrosine-protein kinase SYK (57), Tyrosine-protein kinase Z ... (6), E3 ubiquitin-protein liga ... (3), Ubiquitin-conjugating enz ... (2), Proto-oncogene vav (2), T-cell surface glycoprote ... (2), T-cell surface glycoprote ... (1)

TAXONOMY: Eukaryota only (66)

EXPERIMENTAL METHOD: X-ray (61), Solution NMR (5)

X-RAY RESOLUTION: less than 1.5 Å (6), 1.5 - 2.0 Å (21), 2.0 - 2.5 Å (22), 2.5 - 3.0 Å (10), 3.0 and more Å (2)

RELEASE DATE: before 2000 (3), 2000 - 2005 (5), 2005 - 2010 (7), 2010 - 2015 (22), 2015 - today (29)

Currently showing 1 - 25 of 66 Page: 1 of 3

View: Detailed Reports: Select a Report Sort: Match score: Higher to Lower

1XBA Crystal structure of apo syk tyrosine kinase domain

Atwell, S., Adams, J.M., Badger, J., Buchanan, M.D., Fell, I.K., Froning, K.J., Gao, X., Hendle, J., Keegan, K., Leon, B.C., Muller-Deickmann, H.J., Nienaber, V.L., Noland, B.W., Post, K.W., Rajashankar, K.R., Ramos, A., Russell, M., Burley, S.K., Buchanan, S.G.

(2004) J Biol Chem 279 55827-55832

Released: 11/2/2004 Method: X-ray Diffraction Resolution: 2.0 Å Residue Count: 291

Macromolecule: Tyrosine-protein kinase SYK (protein) Unique Ligands: -- Search term match score: 722.59

Matched fields in 1XBA.cif:

- _citation.title: A novel mode of Gleevec binding is revealed by the structure of spleen tyrosine kinase.
- _entity.pdbx_description: Tyrosine-protein kinase SYK
- _entity.name.com.name: Spleen tyrosine kinase
- _struct.title: Crystal structure of apo syk tyrosine kinase domain
- _struct.keywords.text: syk tyrosine kinase , active conformation, structural genomics, Structural GenomIX, Transferrin

1XBB Crystal structure of the syk tyrosine kinase domain with Gleevec

Atwell, S., Adams, J.M., Badger, J., Buchanan, M.D., Fell, I.K., Froning, K.J., Gao, X., Hendle, J., Keegan, K., Leon, B.C., Muller-Deickmann, H.J., Nienaber, V.L., Noland, B.W., Post, K.W., Rajashankar, K.R., Ramos, A., Russell, M., Burley, S.K., Buchanan, S.G.

(2004) J Biol Chem 279 55827-55832

Released: 11/2/2004 Method: X-ray Diffraction Resolution: 1.57 Å Residue Count: 291

Macromolecule: Tyrosine-protein kinase SYK (protein) Unique Ligands: STI Search term match score: 722.59

1XBBをクリック

立体構造情報を手に入れる

エントリー: 1XBBのデータ画面(トップ)

The screenshot shows the RCSB PDB website interface for entry 1XBB. The browser address bar shows 'www.rcsb.org/structure/1XBB'. The page header includes the RCSB PDB logo and navigation menus. The main content area features a 3D ribbon diagram of the protein structure on the left. To the right, the entry title '1XBB' is displayed, followed by a description: 'Crystal structure of the syk tyrosine kinase domain with Gleevec'. Below this, there is a 'wwPDB Validation' section with a bar chart showing various metrics like Rfree, Clashscore, Ramachandran outliers, Sidechain outliers, and RSRZ outliers. The 'Structure Summary' tab is highlighted with a red arrow.

- Structure Summary : エントリーの概要
- 3D View : 構造グラフィクス
- Annotation : ドメイン、ファミリーの注釈付
- Sequence : 配列の特徴
- Seq. Similarity : 類似した配列の情報
- Str. Similarity : 類似した立体構造の情報
- Experiment : 構造決定の実験条件

デフォルトはStructure Summary

立体構造情報を手に入れる

エントリー: 1XBBのデータ画面(ページ上部のタブをクリックすると。。。)

”Sequence”タブをクリックする

1XBB
Crystal structure of the syk tyrosine kinase domain with Gleevec

Sequence Display for the Entities in PDB 1XBB

The graphical representation below shows this entry's sequences as reported in UniProtKB, in the sample (SEQRES), or as observed in the experiment (ATOM). Different 3rd party annotations can be graphically mapped on the sequence and displayed in the Jmol viewer. Read more about the sequence display on our help pages.

The structure 1XBB has in total 1 chains.

Display Options

Currently viewing all chains

Show unique chains only

Sequence & Structure Relationships. Enable Jmol to view annotations in 3D.

Display Jmol

Redundancy Reduction and Sequence Clustering. View the clustering results for 1XBB.

Chain A: Tyrosine-protein kinase SYK

Chain Downloadable Files

Download FASTA File

View Sequence & DSSP Image

Download Sequence Chain Image

Chain Info

Polymer: 1
Length: 291 residues
Chain Type: polypeptide(L)
Reference: UniProtKB (P43405)
Up-to-date UniProt IDs are provided by the SIFTS project

Display Parameters

Currently displayed
SEQRES sequence.
Display external (UniProtKB) sequence

Mouse over an annotation to see more details. Click on any annotation to enable Jmol.

Add an Annotation

Select

Annotations

Annotations	Details
Domain	Tyrosine-protein kinase SYK: 288 residues
Assignment: SCOP	[view] [reference]
Secondary Structure: DSSP	38% helical (11 helices; 111 residues) 21% beta sheet (14 strands; 63 residues)
Structural Feature: Site Record	BINDING SITE FOR RESIDUE STA A 1 (Software)

Sequence Chain View

Site Record Legend

BINDING SITE FOR RESIDUE STA A 1 (Software)

DSSP Legend

- no secondary structure assigned
- B: beta bridge
- S: strand
- T: turn
- E: beta strand
- G: 3/10-helix
- H: alpha helix

二次構造情報

立体構造情報を手に入れる

エントリー: 1XBBのデータ画面(ページ上部のタブをクリックすると。。。)

"Seq. Similarity"タブをクリックする

1XBB
Crystal structure of the syk tyrosine kinase domain with Gleevec

Sequence Similarity Clusters for the Entities in PDB 1XBB

Sequence Similarity Cutoff	Rank	Chains in Cluster	Cluster ID / Name	Structural variation in cluster
100 %	3	30	1303	
95 %	9	65	794	Flexibility: Low Max RMSD: 3.0, Avg RMSD: 0.9
90 %	9	65	824	
70 %	9	65	861	
50 %	9	66	890	
40 %	61	707	13	
30 %	111	1146	8	

配列類似度

類似する配列数

Instructions
In the table for each entity, view a list of similar sequences by selecting the link associated with the percentage cutoff.
View Table Legend
View more detailed documentation on the redundancy reduction and sequence clustering procedure used by RCSB PDB.
You can also use the structure comparison tool to compare any 2 given structures

Contact Us

RCSB PDB (citation) is managed by two members of the Research Collaboratory for Structural Bioinformatics: Rutgers and UCSD/SDSC

RUTGERS | UC San Diego | SDSC

RCSB PDB is a member of the PDB | EMDatabank

The RCSB PDB is funded by a grant (DBI-1338415) from the National Science Foundation, the National Institutes of Health, and the US Department of Energy.

構造が決定されているものについて、1xbbに対して配列の類似するものの個数が、類似度のレベルごとに示されている

各行のいずれかの数字をクリック

立体構造情報を手に入れる

Sequence Similarity - 1XBB: CRY X

www.rcsb.org/pdb/explore/sequence

RCSB PDB 157530 Biological Macromolecular Structures Enabling Breakthroughs in Research and Education

Search by PDB ID, author, macromolecule, sequence, or ligands

Go

Advanced Search | Browse by Annotations

Structure Summary 3D View Annotations Sequence **Sequence Similarity** Structure Similarity Experiment

1XBB

Crystal structure of the syk tyrosine kinase domain with Gleevec

Sequence Similarity Clusters for the Entities in PDB 1XBB

Entity #1 | Chains: A
Tyrosine-protein kinase SYK protein, length: 291 (BLAST)

Legend

Sequence Similarity Cutoff	Rank	Chains in Cluster	Cluster ID / Name	Structural variation in cluster
100 %	3	30	1303	
95 %	9	65	784	Flexibility: Low Max RMSD: 3.0, Avg RMSD: 0.9 PDBFlex
90 %	9	65	824	
70 %	9	65	861	
50 %	9	66	890	
40 %	61	707	13	
30 %	111	1146	8	

Chains in cluster: 30 | Cluster ID:1303

Instructions

In the table for each entity, view a list of similar sequences by selecting the link associated with the percentage cutoff.

View Table Legend

View more detailed documentation on the redundancy reduction and sequence clustering procedure used by RCSB PDB.

You can also use the structure comparison tool to compare any 2 given structures

Contact Us

ACTION - (A) Select for download / view details OR (B) Select two chains for comparison

Download PDB Files View Structure Details --- Select Comparison Method --- Submit

<input type="checkbox"/>	Rank	PDB ID	Entity ID	Chains	Description	Details	Taxonomy	EC Number
<input type="checkbox"/>	1	4FY0	1	A	Tyrosine-protein kinase SYK	RESIDUES 356-635, PROTEIN KINASE DOMAIN	9606	2.7.10.2 Details
<input type="checkbox"/>	2	5TIU	1	A	Tyrosine-protein kinase SYK	UNP residues 356-635	9606	2.7.10.2 Details
<input type="checkbox"/>	3	1XBB	1	A	Tyrosine-protein kinase SYK		9606	2.7.10.2 Details
<input type="checkbox"/>	4	4FZ7	1	A	Tyrosine-protein kinase SYK	RESIDUES 365-635, PROTEIN KINASE DOMAIN	9606	2.7.10.2 Details
<input type="checkbox"/>	5	4FL1	1	A	Tyrosine-protein kinase SYK		9606	2.7.10.2 Details
<input type="checkbox"/>	6	5Y5T	1	A	Tyrosine-protein kinase SYK	UNP residues 356-635	9606	2.7.10.2 Details
<input type="checkbox"/>	7	4I0T	1	A	Tyrosine-protein kinase SYK	protein kinase domain (UNP residues 356-635)	9606	2.7.10.2 Details
<input type="checkbox"/>	8	4XG7	1	A	Tyrosine-protein kinase SYK	protein kinase domain, residues 356-635	9606	2.7.10.2 Details
<input type="checkbox"/>	9	4RSS	1	A	Tyrosine-protein kinase SYK	UNP RESIDUES 356-635	9606	2.7.10.2 Details
<input type="checkbox"/>	10	4FZ6	1	A	Tyrosine-protein kinase SYK	RESIDUES 356-635, PROTEIN KINASE DOMAIN	9606	2.7.10.2 Details
<input type="checkbox"/>	11	4I0S	1	A	Tyrosine-protein kinase SYK	protein kinase domain (UNP residues 356-635)	9606	2.7.10.2 Details
<input type="checkbox"/>	12	1XBA	1	A	Tyrosine-protein kinase SYK		9606	2.7.10.2 Details

クリックした行に対応した立体構造の情報が表示される

立体構造情報を手に入れる

エントリー: 1XBBのデータ画面(トップ)からデータのダウンロード

“Download Files” をクリックし、
“PDBx/mmCIF File” を選択する

RCSB PDB Deposit Search Visualize Analyze Download Learn More MyPDB

157530 Biological Macromolecular Structures Enabling Breakthroughs in Research and Education

Search by PDB ID Advanced Search

Structure Summary 3D View Annotations Sequence **Sequence Similarity** Structure Similarity Experiment

1XBB

Crystal structure of the syk tyrosine kinase domain with Gleevec

Sequence Similarity Clusters for the Entities in PDB 1XBB

Entity #1 | Chains: A
Tyrosine-protein kinase SYK protein, length: 291 ([BLAST](#))

Sequence Similarity Cutoff	Rank	Chains in Cluster	Cluster ID / Name	Structural variation in cluster
100 %	3	30	1303	
95 %	9	65	784	Flexibility: Low Max RMSD: 3.0, Avg RMSD: 0.9 PDBFlex
90 %	9	65	824	
70 %	9	65	861	
50 %	9	66	890	
40 %	61	707	13	
30 %	111	1146	8	

Display Files **Download Files**

- FASTA Sequence
- PDB File (Text)
- PDB File (gz)
- PDBx/mmCIF File**
- PDBx/mmCIF File (gz)
- PDBML/XML File
- PDBML/XML File (gz)
- Structure Factor (Text)
- Structure Factor (gz)
- Biological Assembly (gz) (A)

Instructions:
In the table for list of similar structures, selecting the line the percentage
[View Table Legend](#)
[View more details](#) on the redundancy reduction and sequence clustering procedure used by RCSB PDB.
You can also use the [structure comparison tool](#) to compare any 2 given structures

Chains in cluster: 30 | Cluster ID:1303

立体構造情報を手に入れる

エントリー: 1XBBのデータ画面(トップ)からデータのダウンロード

“ファイルを保存”をクリック
ダウンロードフォルダに1xbb.cifが
保存される

1XBB

Crystal structure of the syk tyrosine kinase domain with Gleevec

Sequence Similarity Clusters for the Entities in PDB 1XBB

Entity #1 | Chains: A
Tyrosine-protein kinase SYK protein, length: 291 ([BLAST](#))

Sequence Similarity Cutoff	Rank	Chains in Cluster	Cluster ID / Name	Structure
100 %	3	30	1303	
95 %	9	65	784	Flexib Max P PDBFl
90 %	9	65	824	
70 %	9	65	861	
50 %	9	66	890	
40 %	61	707	13	
30 %	111	1146	8	

1xbb.cif を開く

次のファイルを開こうとしています:

- 1xbb.cif
ファイルの種類: cif File (72.0 KB)
ファイルの場所: <https://files.rcsb.org>

このファイルを保存しますか?

キャンセル **ファイルを保存**

on the redundancy reduction and sequence clustering procedure used by RCSB PDB.

You can also use the [structure comparison tool](#) to compare any 2 given structures

Contact Us

座標データを見ている何も情報が得られない。



立体構造の可視化ツールが必要



分子グラフィクス

PyMol, jViewer, MolMol etc.

今回は、MolMilで分子構造を見る

molmilの使い方 1 - 1個の構造の表示 -

molmil: PDBjで開発されたタンパク質立体構造表示ビューア

ブラウザを立ち上げて、molmil PDBjで検索

pdbj.org › help › molmil ▼

[Molmil 分子ビューア - Help - 日本蛋白質構造データバンク - PDBj](#)

2013/12/19 - Molmilはインターネットに接続されたウェブ環境で利用するためにPDBjで開発した新しい分子閲覧ソフトで、できるだけ多くの環境で利用できるように設計されています。

OpenGLによるGPUハードウェアアクセラレーションを用い、美しい ...

このページに 4 回アクセスしています。前回のアクセス: 20/03/12

クリック

pdbj.org › molmil ▼ [このページを訳す](#)

[Molmil viewer <<< Pymol-like command interface bound. <](#)

Molmil viewer. <<<. Pymol-like command interface bound. <

このページに 5 回アクセスしています。前回のアクセス: 19/12/08

pdbj.org › help › molmil-manual ▼

[Molmilユーザマニュアル - Help - 日本蛋白質構造データ ... - PDBj](#)

Molmil 分子ビューア

このページの他言語版もあります: [English](#)

[\[はじめに\]](#) [\[利用方法\]](#) [\[利用環境\]](#) [\[トラブルシューティング\]](#)

はじめに

Molmilはインターネットに接続されたウェブ環境で利用するためにPDBJで開発した新しい分子閲覧ソフトで、できるだけ多くの環境で利用できるように設計されています。OpenGLによるGPUハードウェアアクセラレーションを用い、美しいグラフィックスと最適なパフォーマンスを実現します。Molmilビューアは既存のサービスに組み込むだけでなく、jVと同じように [スタンドアロンビューア](#) としても用いることができます。

スタンドアロンビューアを使って、自身で作成したPDB、PDBML、mmCIF、mmJSONの各フォーマットファイルを読み込むことができるだけでなく、PDB IDやChem Comp IDを指定して既存の構造を読み込むこともできます。更に、jV用 [ポリゴンXMLファイル](#) も読み込むことができます。将来更に機能を追加していく予定です。また、[機能追加のリクエスト](#) もお待ちしております。また、Molmilを自分で作ったウェブサイトに埋め込むこともできます。その方法については [こちら](#) をご覧ください。

利用方法

Molmilの使い方については下記ページを参照下さい。

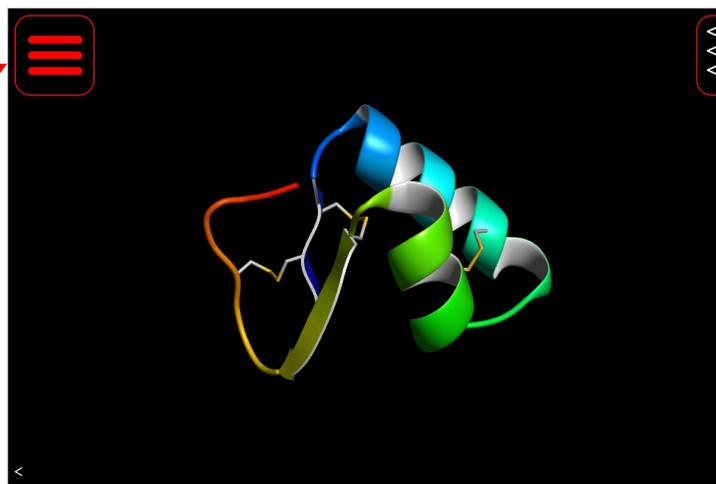
- [Molmilユーザマニュアル](#)
- [基本的な使い方](#) (GitHub)
- [ウェブページへの埋め込み方](#) (GitHub)
- [詳細な使い方 \(Molmil APIの解説\)](#) (GitHub)
- [Molmil FAQ](#) (GitHub)

利用環境

Molmilを利用する際、ご利用のブラウザを最新版に更新しておくことをお勧めします。最新ではないSafariやOperaの中には、WebGLを有効化する必要があるものもあります。また、ドライバが古すぎない（2010年12月31日以前ではない）ことも確認して下さい。

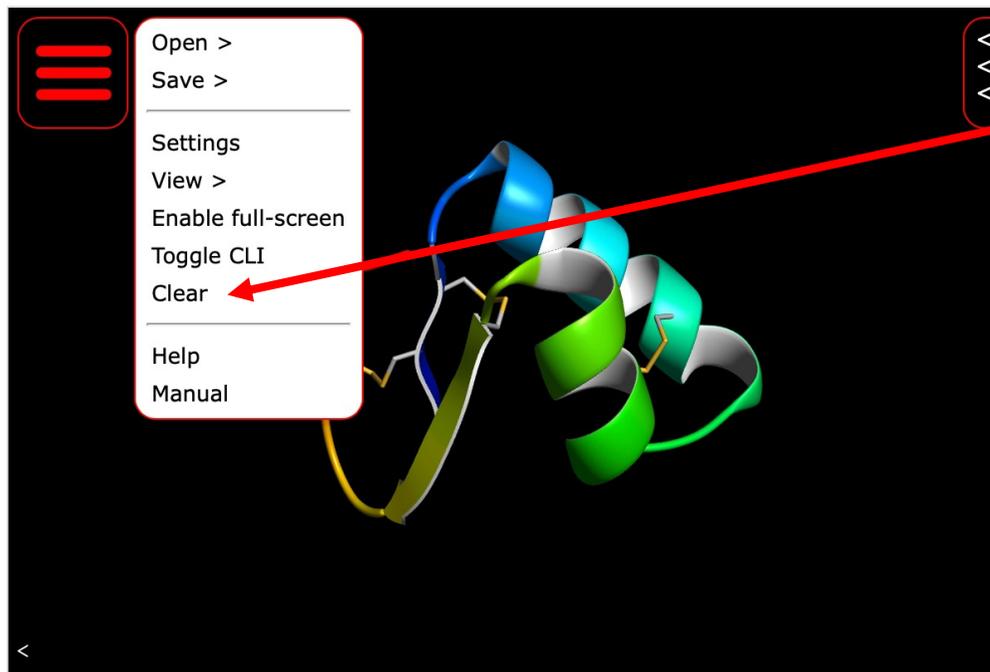
モバイル端末（スマートフォン、タブレット）のうちAndroid、Blackberry場合、Chrome、Firefox、Opera、Blackberryブラウザのいずれかで最新のものをお使い下さい。iOS（iPhone、iPad）については、現在iOS 8.0は対応しています。

ブラウザがWebGLをサポートしている場合、この下に新しいMolmilビューアを表示することができるでしょう。



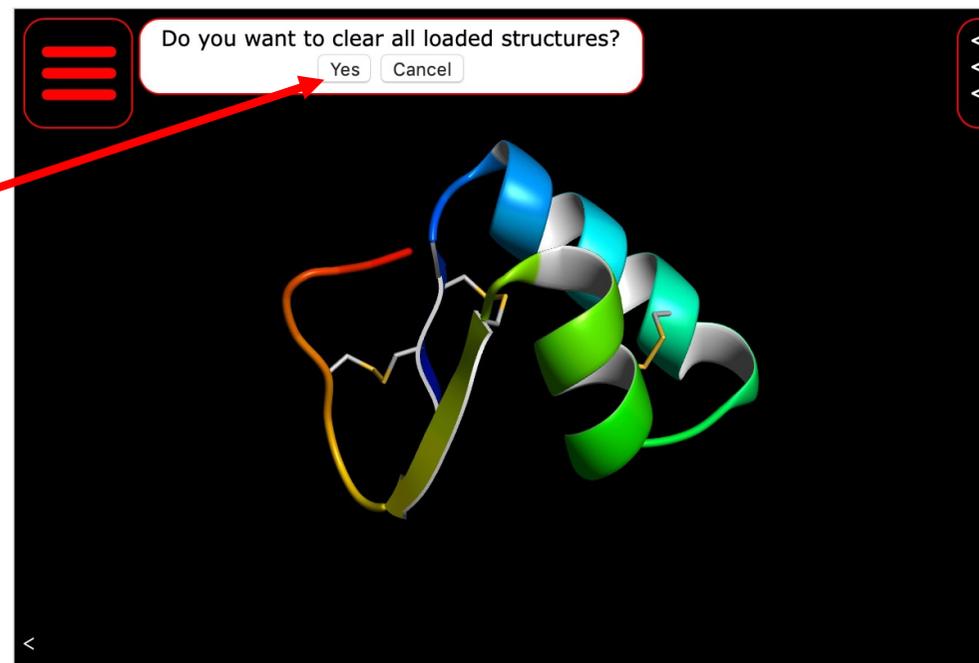
クリック

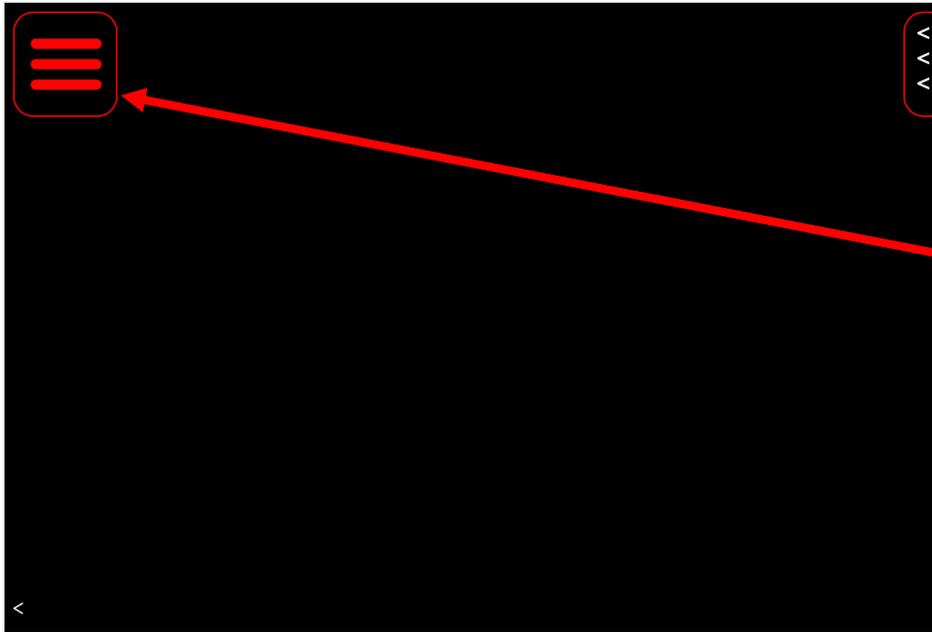
PDBID [1cm](#): 原子レベルの解像度で解いた疎水性たんぱく質の水構造。克蘭ピン (crambin) の結晶内における水分子の5員環。



Clearをクリック

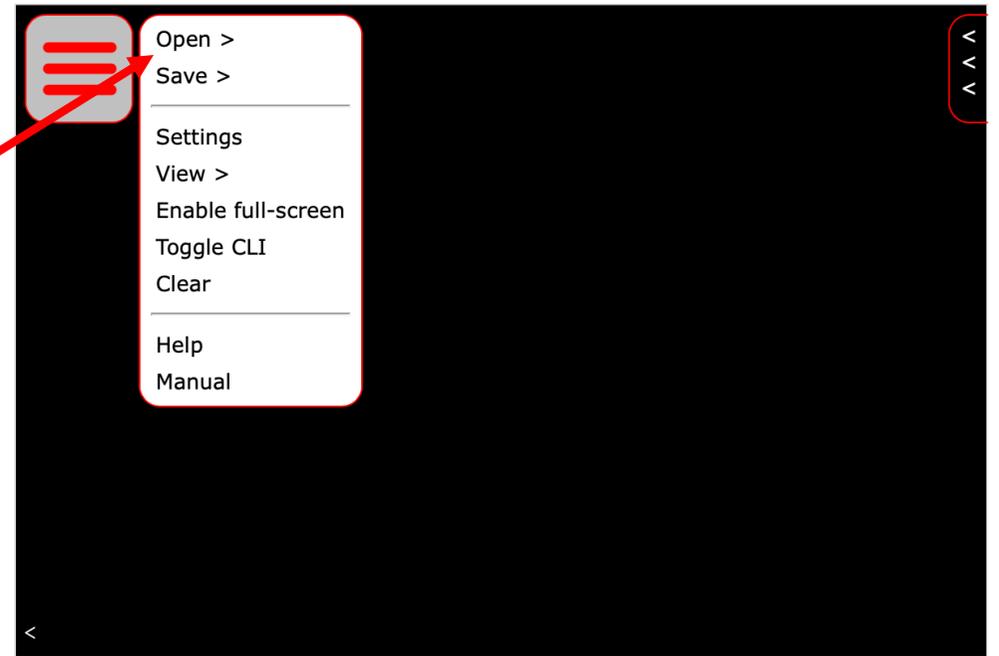
Yesをクリック
デフォルトで表示されている
構造をクリア

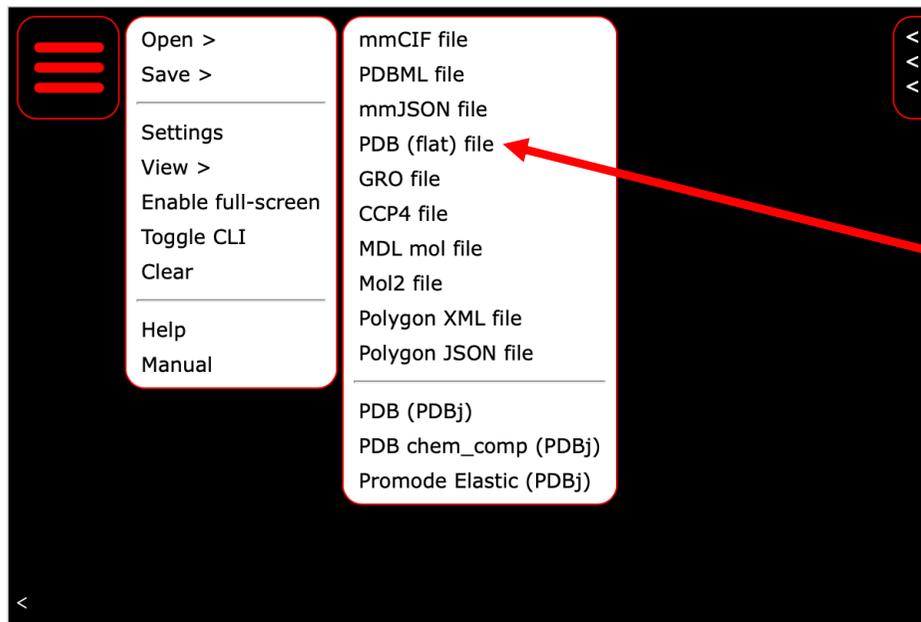




クリック

Openをクリック

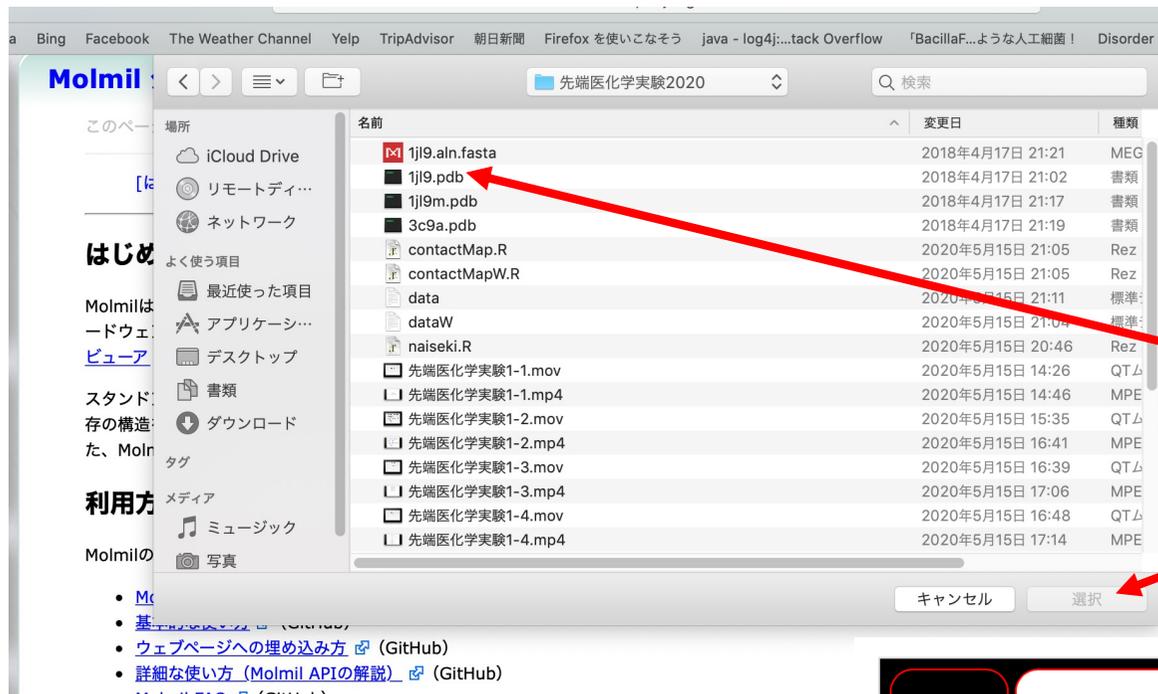




PDB (flat) file クリック

ファイルを選択をクリック

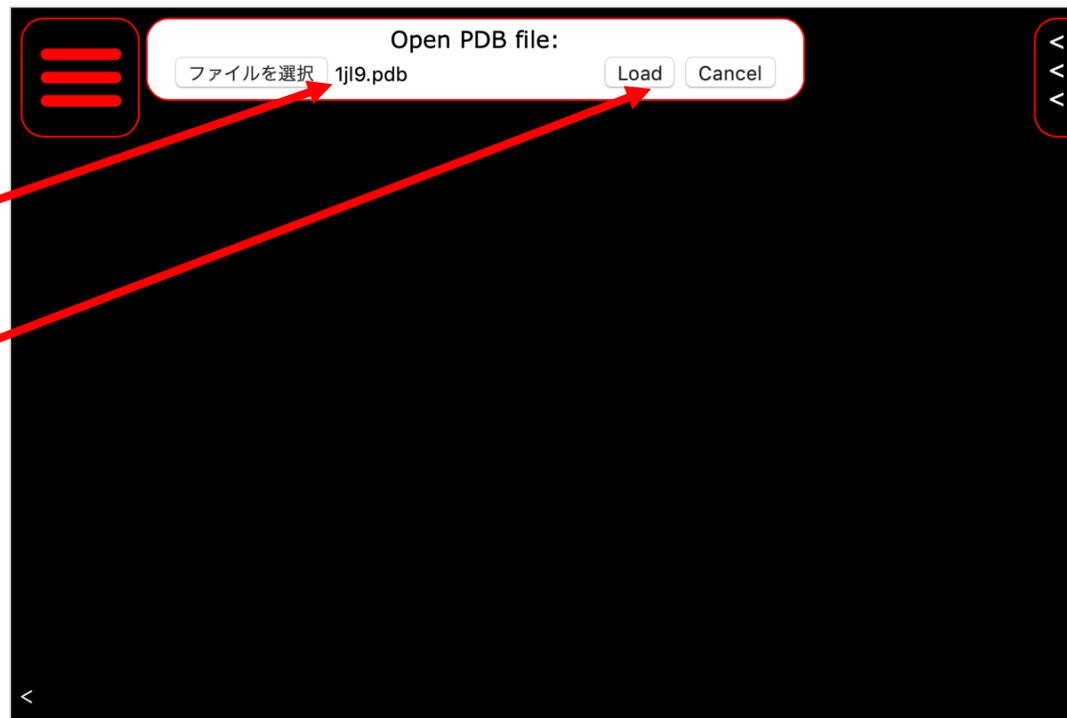


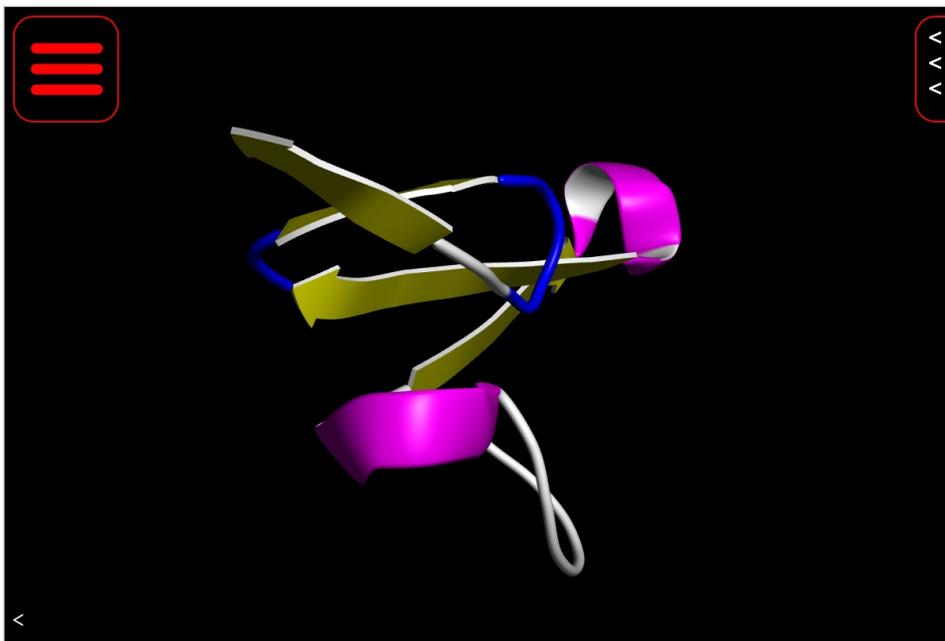


ファイルチューザが開くので
ディレクトリを移動して
1j19.pdb を選択し、“選択”を
クリック

ファイルが選択されている
ことを確認

Loadをクリック





1jl9の構造が表示される。

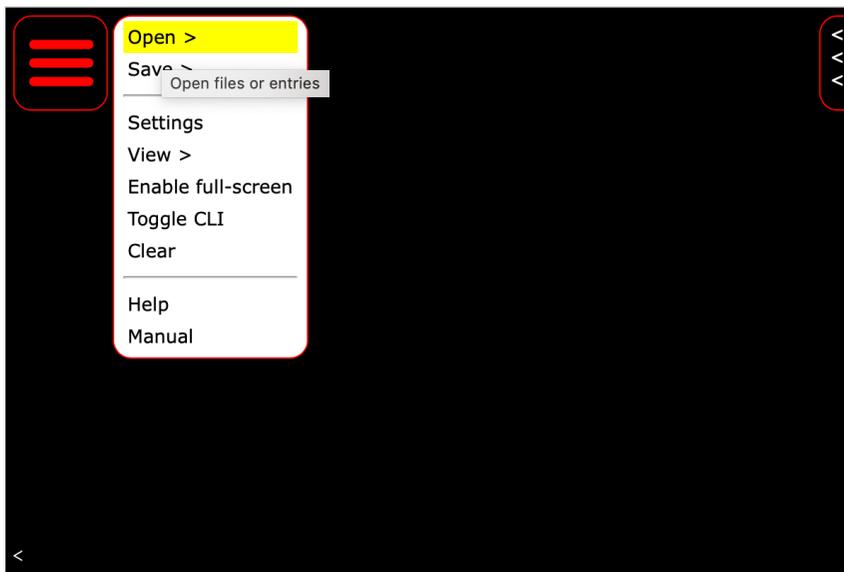
画面を適当にクリックし
クリックしたままドラッグすると
回転することを確認

Shiftキーを押したままで、
クリックしたままドラッグすると
平行移動することを確認

PDBID [1jcn](#): 原子レベルの解像度で解いた疎水性たんぱく質の水構造。克蘭ピン (crambin) の結晶内における水分子の5員環。

Clearで構造を消す

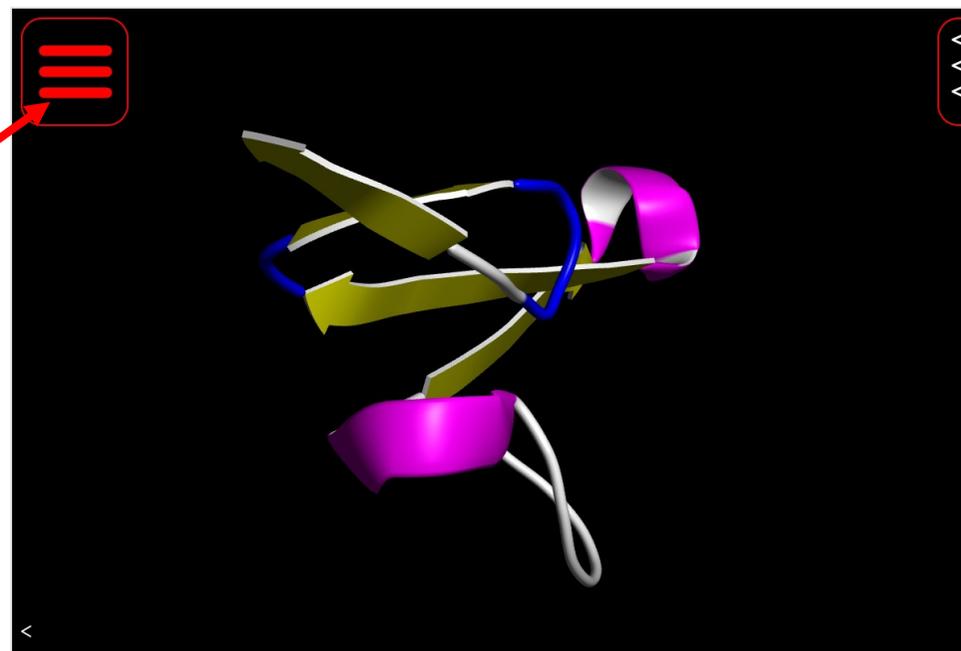
molmilの使い方 2 - 2個の構造の表示 -

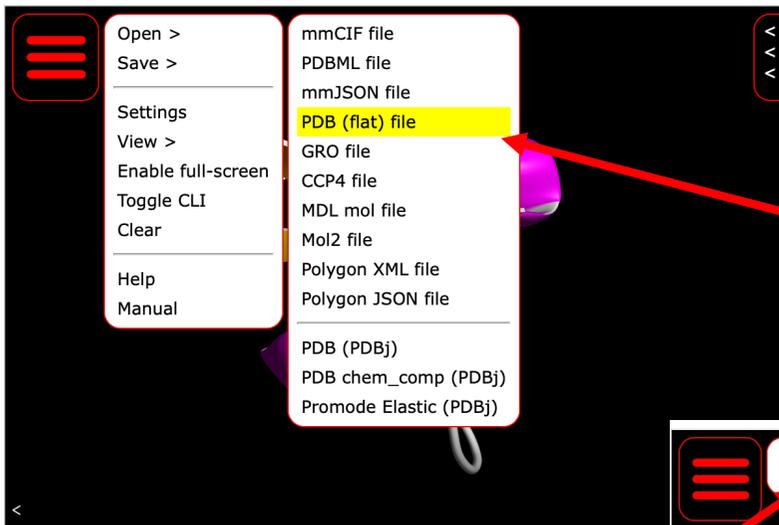


まず、先ほど説明したやり方で
1個の構造(1jl9.pdb)を読み込む

操作は同じなので、読み込み部分は省略

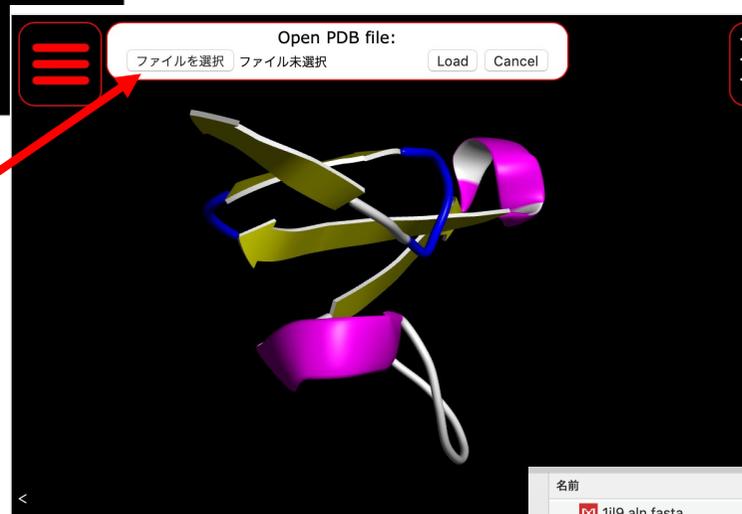
読み込めたらここをクリックして
別の構造おの読み込み操作を行う





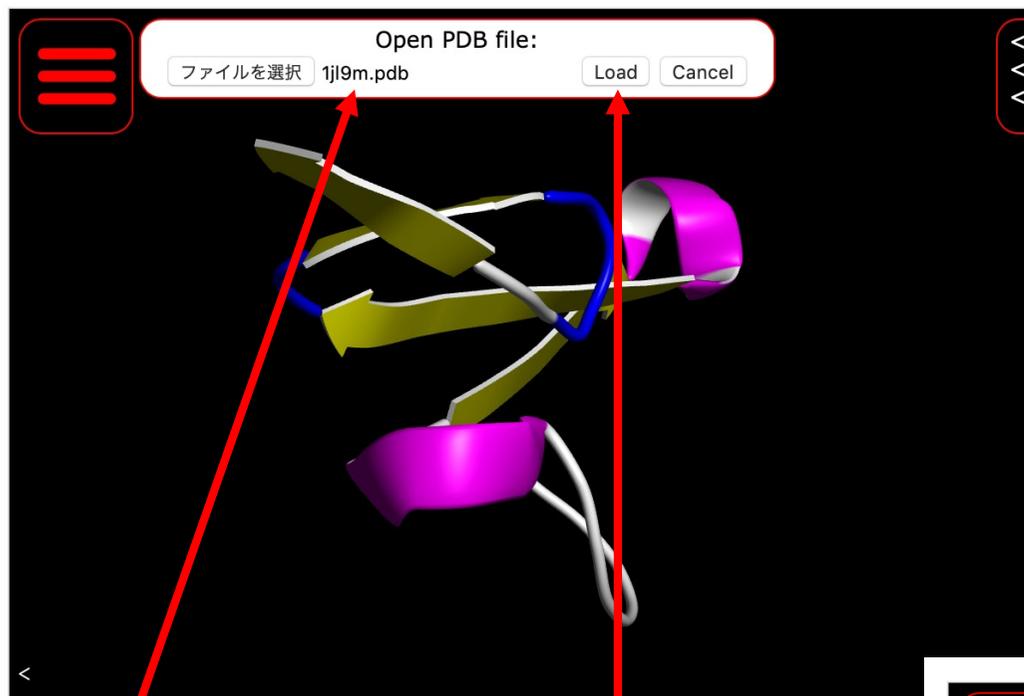
Open → PDB (flat) file を選択

ファイルを選択を
クリック



ファイルチューザから
1jl9m.pdbを選択し、
選択をクリック

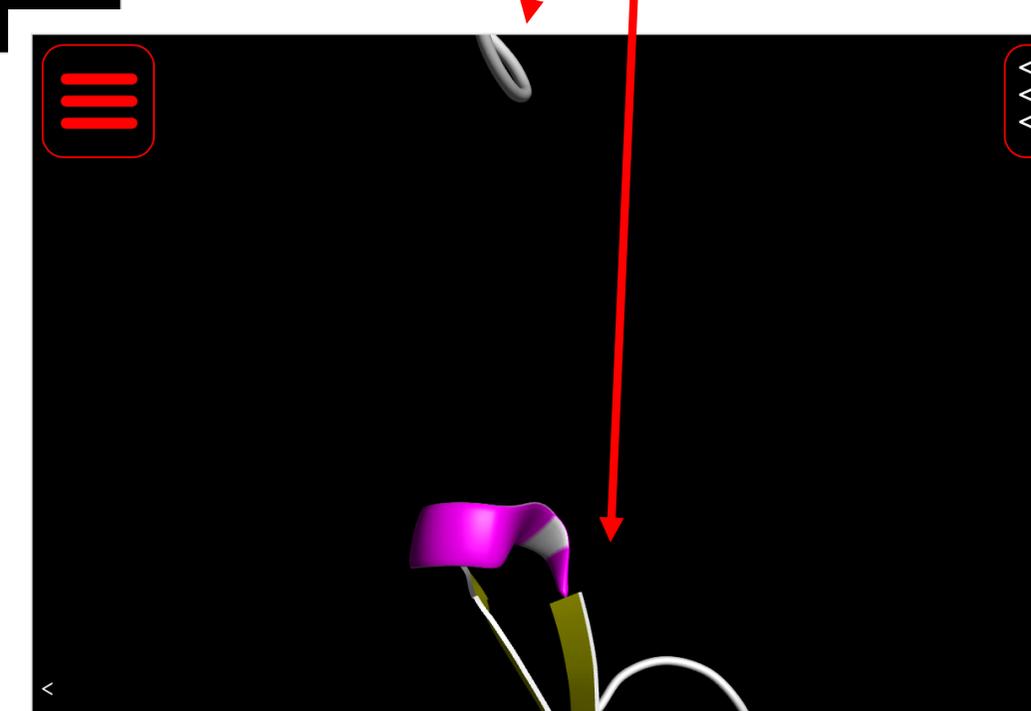


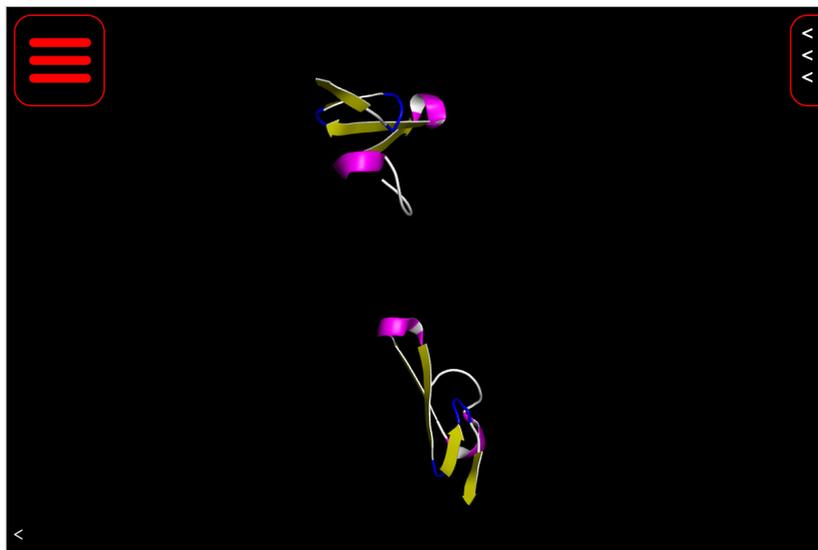


1jl9m.pdbが選択されていることを確認

Loadをクリック

2つの構造が表示されているがスケールがウィンドウサイズにあわないため全体像が表示されていない





マウスのホイールを回転させて
拡大、縮小
あるいはマウスの右ボタンをクリックした
ままドラッグして拡大縮小

平行移動、回転も同様にできる。

1j19m.pdbは、1j19.pdbの座標を適当に回転させた後に、平行移動したものの

今回、遺伝的アルゴリズムでこの二つの構造の重ね合わせを行う。
もともと同一構造なのでピッタリ重なるはず。

1j19.pdbと遺伝的アルゴリズムで回転させた座標データをmolmillに表示し、本当に
重ね合わせできてるかを確認する

3. Rの復習

今回の実習に使うプログラムはRで書かれている。

プログラムの動作を理解し、またプログラムを書き直すためにRについて復習する

Rとは

統計解析のためにデザインされたプログラミング言語

オープンソースのフリーソフトウェア（無料）

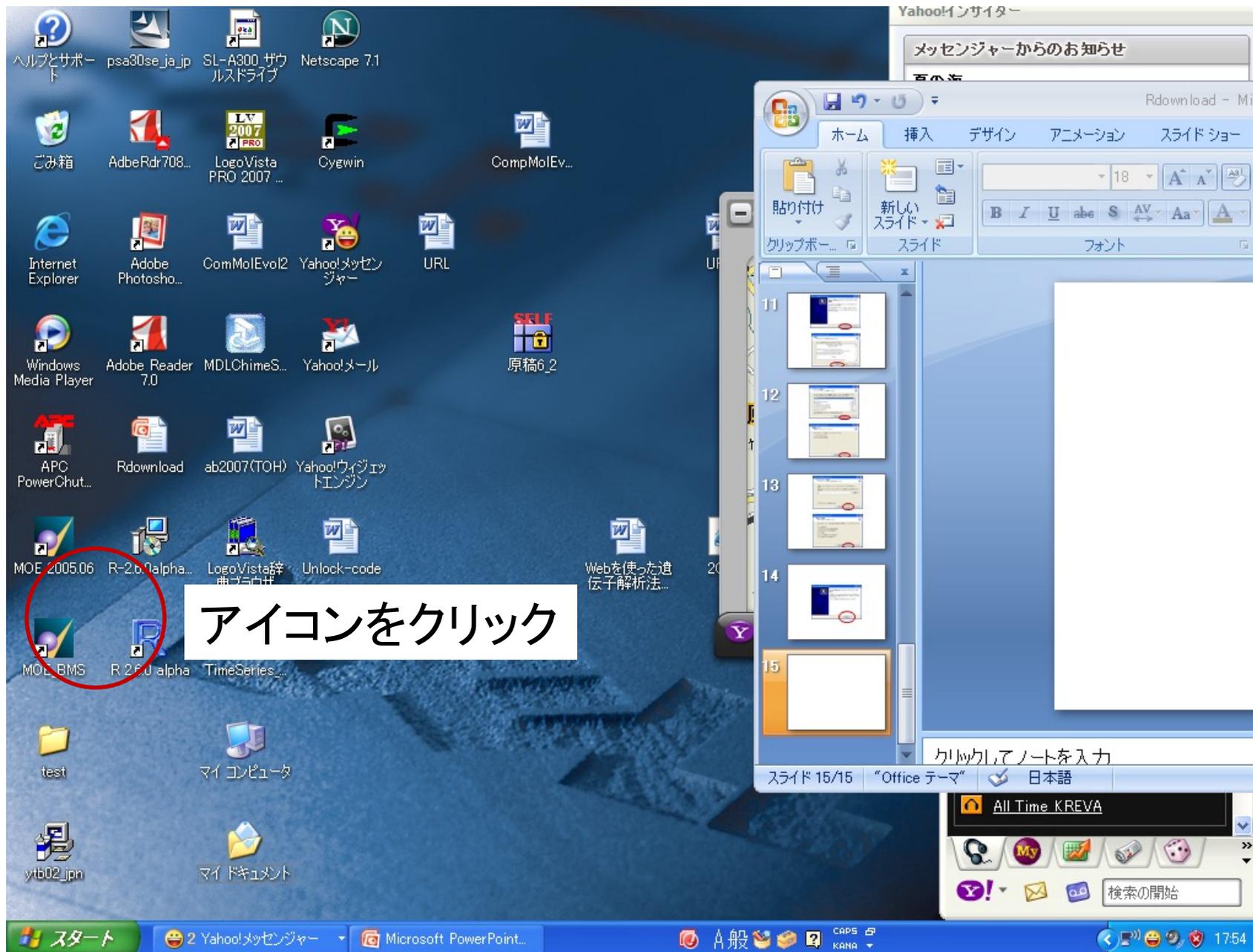
マルチプラットフォーム

(Windows, Mac, UNIXなどのOSで同一の作業ができる)

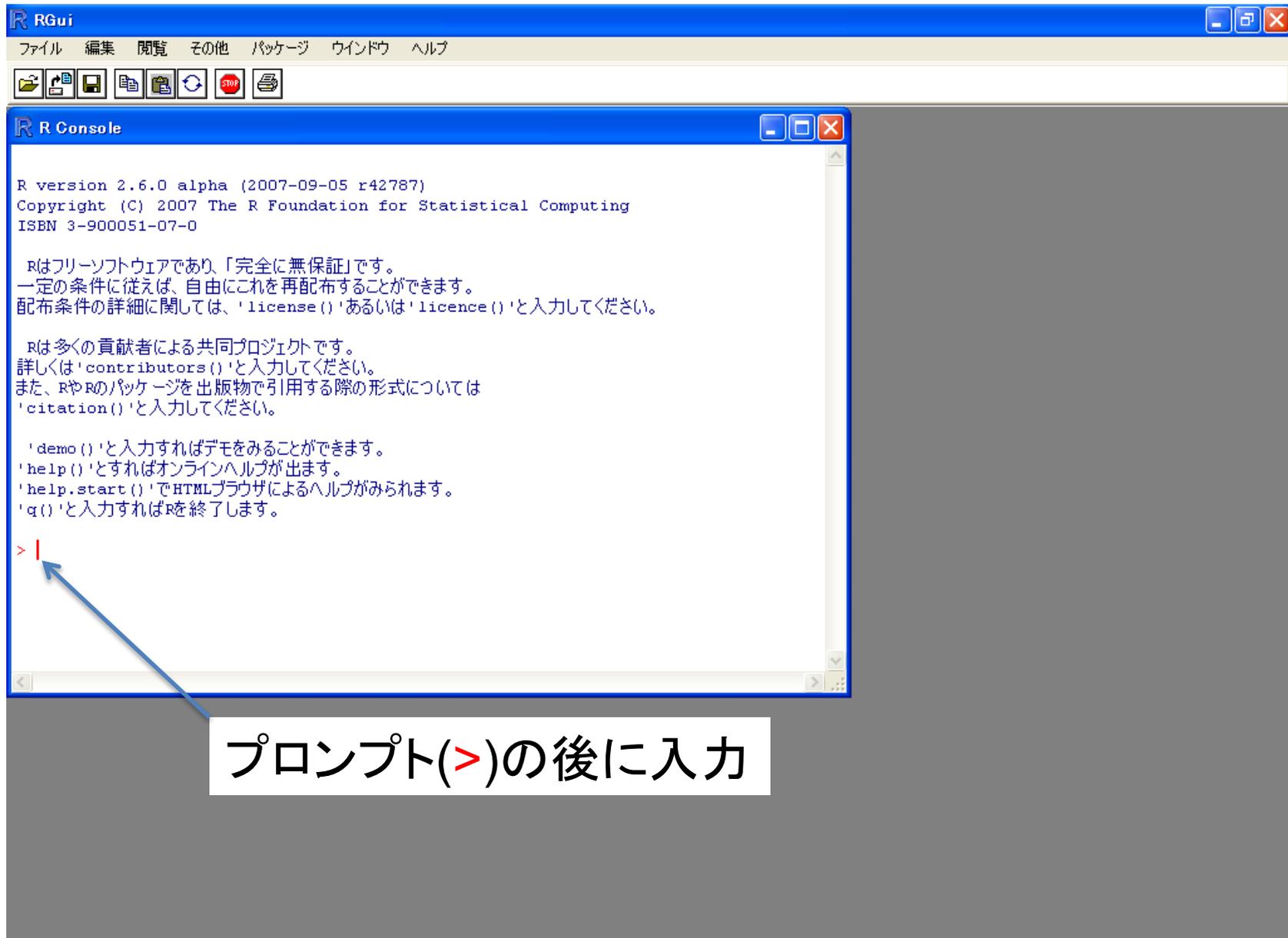
世界各地のRユーザが開発したプログラム（**パッケージ**）が
CRAN (The Comprehensive R Archive Network)を介して
配信されている。（最新の手法が無料で利用可能）

生物系、特にゲノムワイドなデータの解析のパッケージも
充実している（Bioconductor）

Rの起動



コンソール画面が立ち上がる



プロンプトの後ろにコマンドを入力してEnterを押すと実行

```
> 3 + 4
```

加

```
[1] 7
```

```
> 2.3-2.1
```

減

```
[1] 0.2
```

```
> 3.14*2.0
```

乗

```
[1] 6.28
```

```
> 3/4
```

除

```
[1] 0.75
```

```
> 1.5^3
```

冪

```
[1] 3.375
```

```
> 3.2%%3
```

剰余

```
[1] 0.2
```

```
> 9%/4
```

切り捨て

```
[1] 2
```

変数への値の代入と演算

```
> a <- 2
> a
[1] 2
> print(a)
[1] 2
> b <- 9.7
> b^a
[1] 94.09
> c <- b^a + b/a
> c
[1] 98.94
> b/a
[1] 4.85
```

変数名は、Rの関数名などとして予約されているものでなければ、半角英数字を使って自由に指定できる
ただし、変数名を数字ではじめてはいけない
大文字、小文字は区別される

```
> a <- 2
> A <- 3
> print(a)
[1] 2
> print(A)
[1] 3
> ips <- 5
> a.Boss <- 3
> ips <- ips + a.Boss
> ips
[1] 8
```

```
> sqrt(2)
[1] 1.414214
> exp(2)
[1] 7.389056
> exp(1)
[1] 2.718282
> 2.718282^2
[1] 7.389057
> sqrt(exp(2))
[1] 2.718282
```

sqrt

平方根

exp

指数関数

関数の引数に関数をいれることも
できる

```
> sin(pi/2)
[1] 1
> cos(pi/2)
[1] 6.123234e-17
> sin(pi)
[1] 1.224647e-16
> cos(pi)
[1] -1
> tan(0)
[1] 0
> tan(pi/2)
[1] 1.633124e+16
> sin(pi/2)/cos(pi/2)
[1] 1.633124e+16
```

三角関数

sin, cos, tan

数值的に計算されるので
正確には0であるが、
非常に小さな値として表示

$$6.123234e-17 = 6.123234 \times 10^{-17}$$

本来、無限大になるはず
だが、数值的に計算され
ているので大きな値として
表示

四捨五入

```
> pi
[1] 3.141593
> round(pi)
[1] 3
> round(pi, 4)
[1] 3.1416
> round(pi, 3)
[1] 3.142
> round(pi, 2)
[1] 3.14
```

切り捨て trunc
最大整数 floor
切り上げ ceiling

```
> trunc(-5.678)
```

```
[1] -5
```

```
> trunc(5.678)
```

```
[1] 5
```

```
> floor(-5.678)
```

```
[1] -6
```

```
> floor(5.678)
```

```
[1] 5
```

```
> ceiling(-5.678)
```

```
[1] -5
```

```
> ceiling(5.678)
```

```
[1] 6
```

Rではベクトルはc()で作成

```
> y <- c(9.02, 20.124, -0.998, 4.787)
```

```
> y
```

```
[1] 9.020 20.124 -0.998 4.787
```

最大、最小 max, min

```
> min(y)
```

```
[1] -0.998
```

```
> max(y)
```

```
[1] 20.124
```

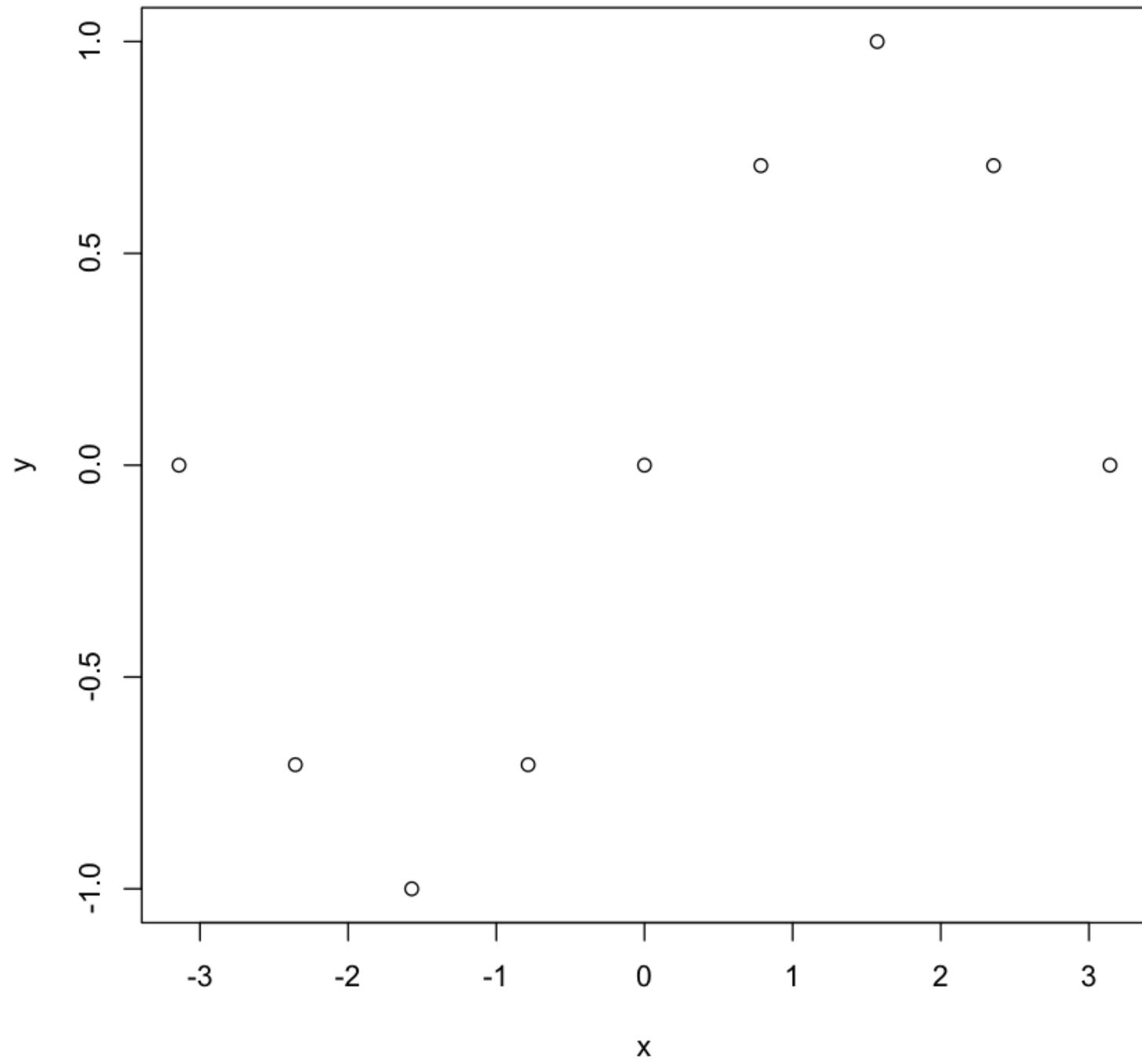
```
> x <- c(-pi, -pi*0.75, -pi*0.5, -pi*0.25, 0, pi*0.25, pi*0.5,
pi*0.75, pi)
> y <- sin(x)
> y
[1] -1.224647e-16 -7.071068e-01 -1.000000e+00 -7.071068e-01
0.000000e+00 7.071068e-01 1.000000e+00 7.071068e-01
[9] 1.224647e-16
> plot(x, y)
```

次ページ図

もっと、間隔を細かくとって、より滑らかな曲線を描きたい。
しかし、手作業で x 座標を書くのはめんどくさい。

次の方法で解決

seq ある区間を、等しい間隔で区切ったベクトルを発生させる



```
> seq(1,5,0.5)
[1] 1.0 1.5 2.0 2.5 3.0 3.5 4.0 4.5 5.0
> seq(3,5,1)
[1] 3 4 5
```

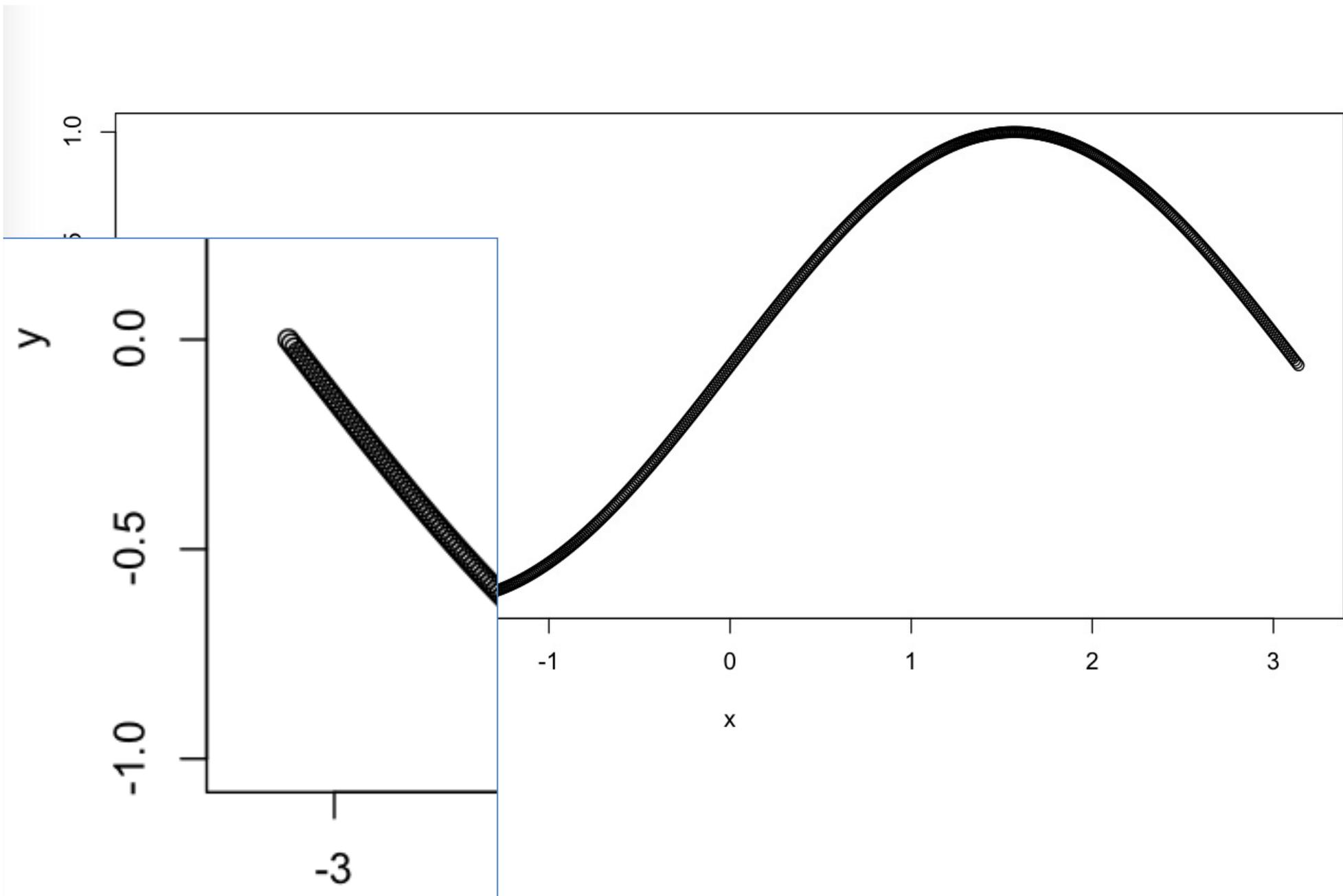
`seq(start, end, 間隔)`

```
> x <- seq(-pi,pi,0.01)
> x
 [1] -3.141592654 -3.131592654 -3.121592654 -3.111592654
...
[628] 3.128407346 3.138407346
```

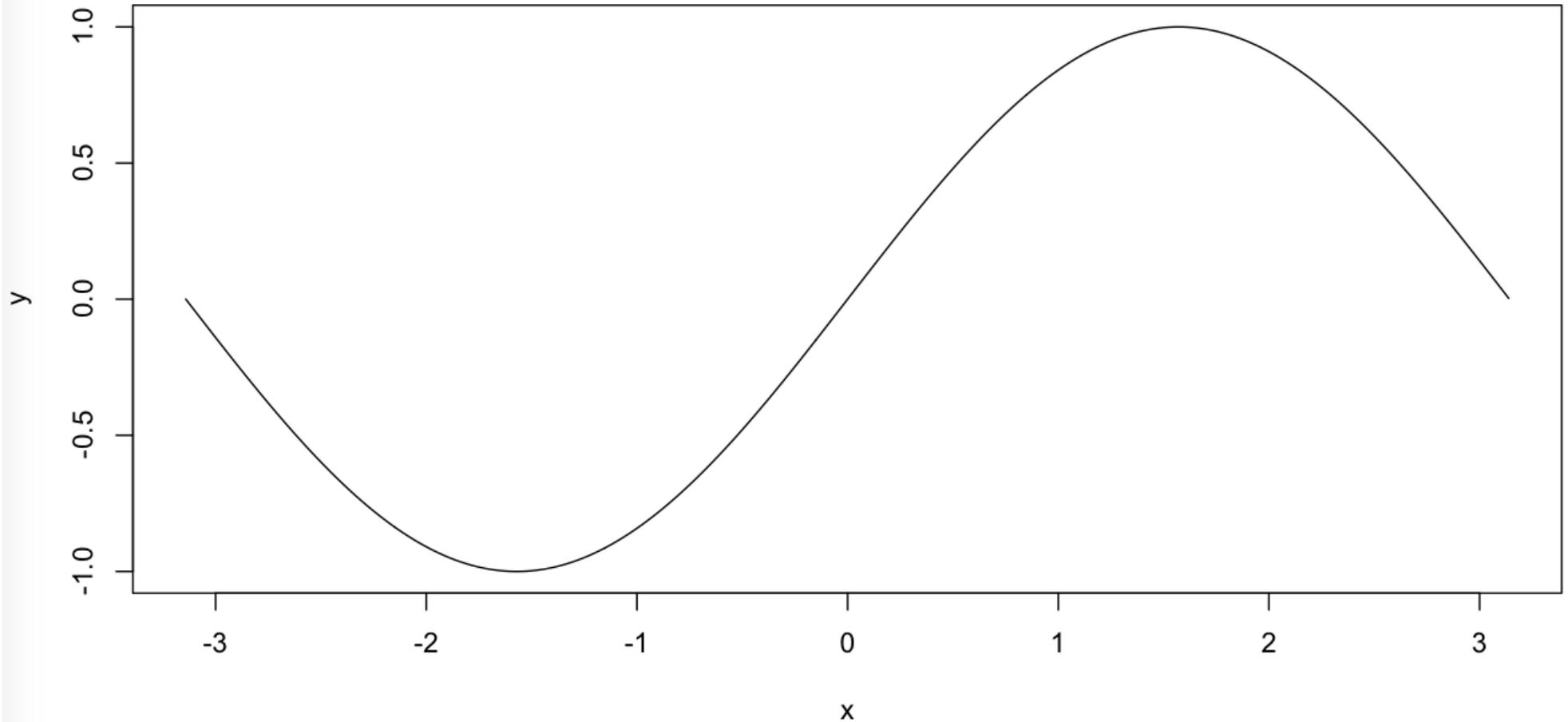
```
> y <- sin(x)
> plot(x,y)
```

図は次ページ

なめらかになったが、ドットで表示するのは見栄えがよくない。
点を打って線でつなぎたい。

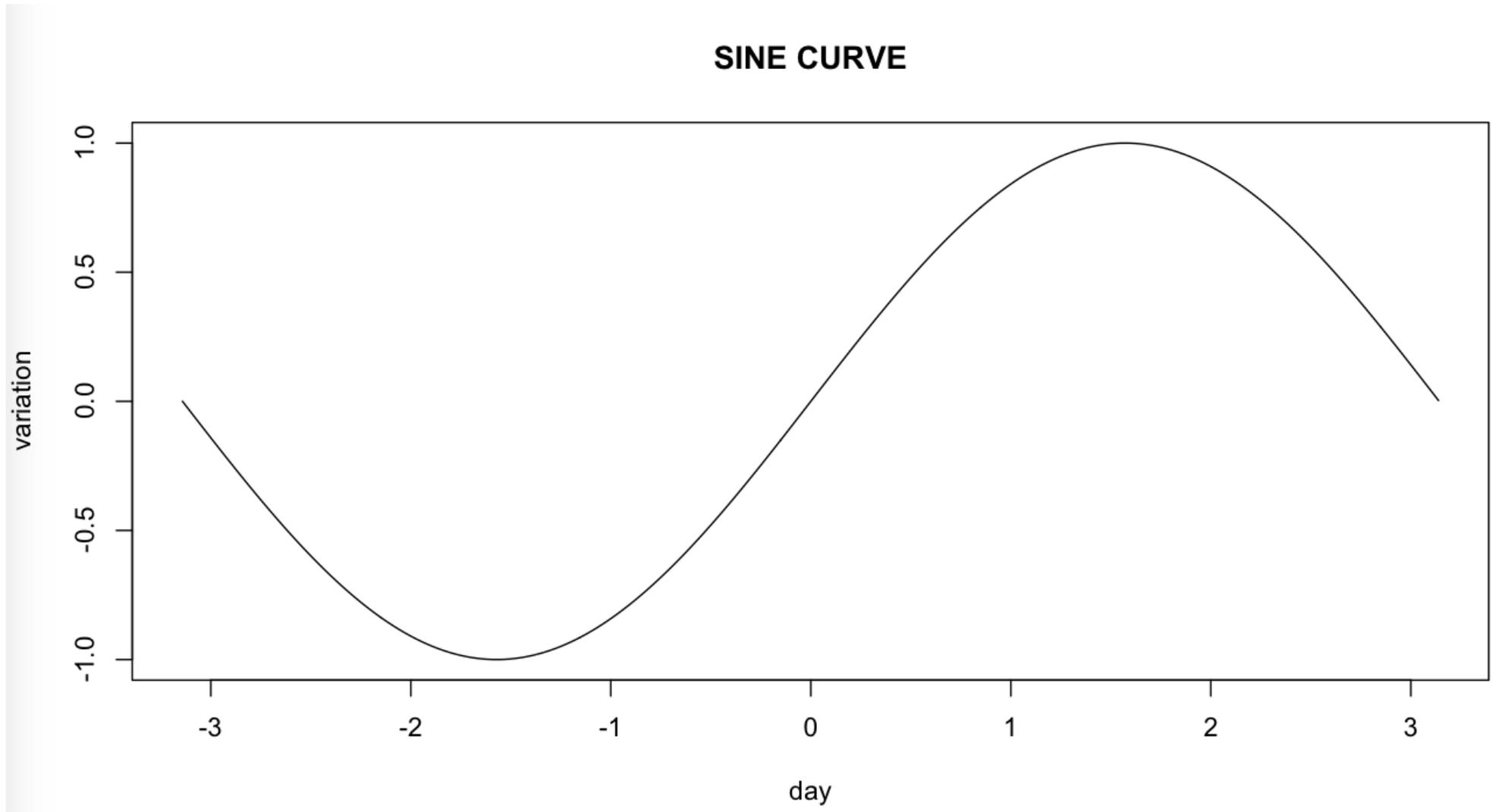


```
> plot(x, y, ty='l')
```



x軸のラベルをday、y軸のラベルをvariation, 全体のタイトルをSINE CURVEとする。

```
> plot(x, y, ty='l', xlab="day", ylab="variation", main="SINE CURVE")
```



生成されたグラフィクスの保存

画像ファイルの形式を、pngあるいはepsで保存したい場合
グラフィクスコマンドを実行する前に
png (ファイル名)、あるいはpostscript (ファイル名)を実行

グラフィクスコマンド(plotやcurve)を実行した後に

dev.off()を実行

Rを実行している作業ディレクトリ(後述)内に、画像ファイル
が作成される。

```
> png("test.png")
> plot(x, y, pch=23)
> dev.off()
quartz
  2
```

スクリーン上にはグラフィクスは表示されない

test.pngをクリックすると図が表示される。

```
getwd()
```

でディレクトリを確認

関数quit()で終了



```
> a
[1] 2
> print(a)
[1] 2
> b <- 9.7
> b^a
[1] 94.09
> c <- b^a + b/a
> c
[1] 98.94
> b/a
[1] 4.85
> A <-1
> a
[1] 2
> a.b <-3
> a.b
[1] 3
> 3 -> b
> b
[1] 3
> 4 = c
以下にエラー 4 = c : 代入の左辺が不正 (do_set) です
> c = 4
> a <- 2
> A <- 3
> print(a)
[1] 2
> print(A)
[1] 3
> ips <- 5
> a.Boss <- 3
> ips <- ips + a.Boss
> ips
[1] 8
> a <- "atgcttgaccgtaat"
> aaseq <- "MLKPAQWELLTGRS"
> a
[1] "atgcttgaccgtaat"
> aaseq
[1] "MLKPAQWELLTGRS"
> aaseq[1]
[1] "MLKPAQWELLTGRS"
> substr(aaseq,2:3)
以下にエラー substr(aaseq, 2:3) :
  引数 "stop" がありませんし、省略時既定値もありません
> greeting <- "Hello World!"
> print(greeting)
[1] "Hello World!"
> quit()
```

Rセッションを終了
ワークスペースのイメージファイルを保存しますか？

保存しない キャンセル 保存

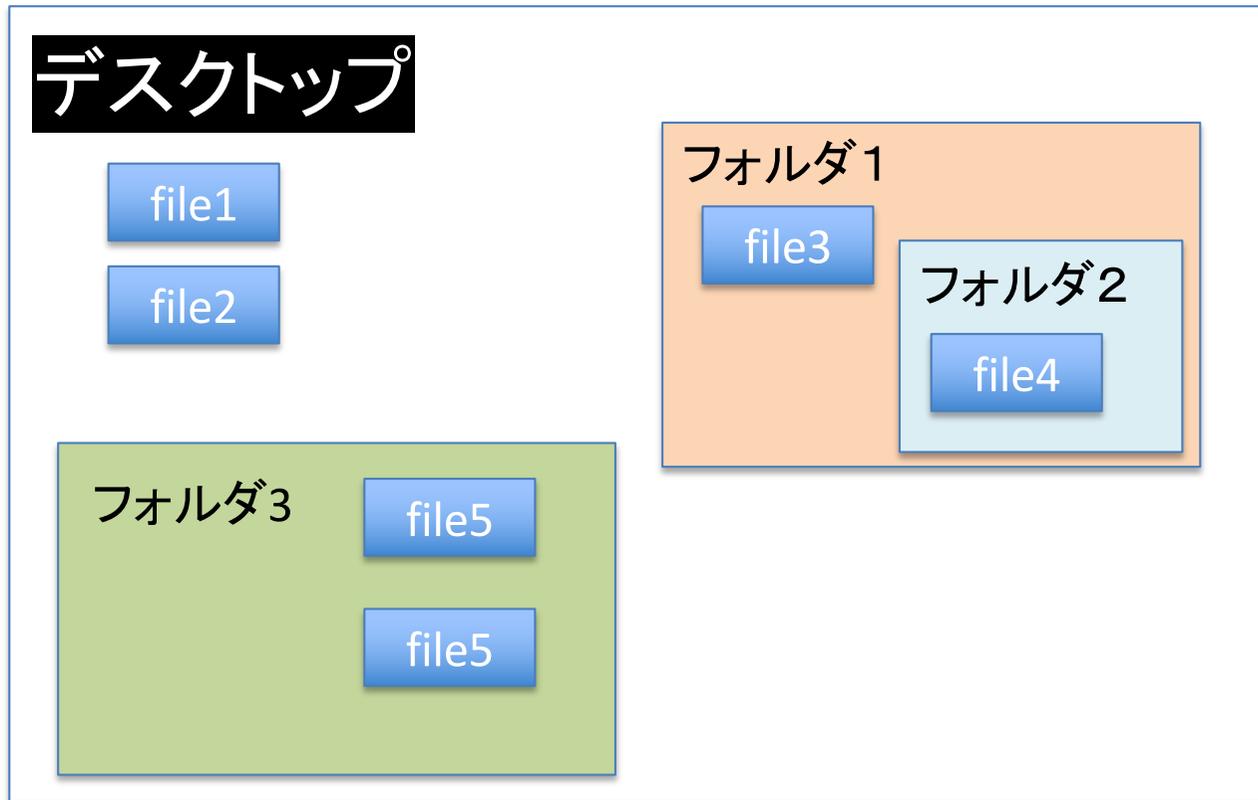
“保存しない”を押して
終了してから、再度
Rを起動

作業ディレクトリの確認と設定

フォルダ = ディレクトリ

フォルダの中には、ファイルだけでなくフォルダも作れる

※ デスクトップもフォルダの一つ



Rでは、どのフォルダで作業するかを選択できる

`getwd()` : 現在のディレクトリ名を出力

`list.files()` : 現在のディレクトリ内の
ファイルの一覧を出力

`dir()` : `list.files()`と同じ

`setwd("フォルダ名")` : ワークスペースを二重引用符
で囲まれたフォルダに変更
←-- 作業フォルダを移動するのに使う

`ls()` : 現在Rが持っているオブジェクト(ベクトルや変数)
の一覧

setwdの使用例

Macの場合

```
setwd("/Users/toh/Desktop")
```

Windowsの場合

```
setwd("c:/mydocuments")
```

今回、作成したRのプログラムやファイルをダウンロードして使用する。
そのプログラムやファイルが置かれているフォルダを作業ディレクトリとすること

ベクトル (vector)

`c()`, `seq` 以外の方法

```
> 1:50
 [1]  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24
[25] 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48
[49] 49 50
>
> rep(3, 12)                # 3を12回繰り返す
 [1] 3 3 3 3 3 3 3 3 3 3 3 3
>
> rep(c(5,7), 2)           # 5と7を2回繰り返す
 [1] 5 7 5 7
>
> rep(3:5, 1:3)           # 3を1回、4を2回、5を3回繰り返す
 [1] 3 4 4 5 5 5
>
> x <- 1:50
> x
 [1]  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24
25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41
[42] 42 43 44 45 46 47 48 49 50
> x <- seq(1,50)
> x
 [1]  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24
25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41
[42] 42 43 44 45 46 47 48 49 50
```

```
> x <- c(2, 4.5, 1.9, 1.1, 3.3)
> x
[1] 2.0 4.5 1.9 1.1 3.3
> y <- c(4.1, 5.5, 1, 2.1, 2.2)
> y
[1] 4.1 5.5 1.0 2.1 2.2
> length(x) ベクトルのサイズ
[1] 5
> length(y)
[1] 5
> z <- 5*x ベクトルの実数倍
> z
[1] 10.0 22.5 9.5 5.5 16.5
> x %*% y ベクトルの内積
      [,1]
[1,] 44.42
```

```
> x[1]
```

```
[1] 2
```

ベクトル x の要素1を参照

```
> y[2]
```

```
[1] 5.5
```

ベクトル y の要素1を参照

```
> x[1]*y[1]+x[2]*y[2]+x[3]*y[3]+x[4]*y[4]+x[5]*y[5]
```

```
[1] 44.42
```

```
> sum(x)
```

```
[1] 12.8
```

```
> sum(y)
```

```
[1] 14.9
```

```
> x <- matrix(c(1,9,2,8,3,7,4,6), ncol=4)
```

```
> x
```

```
      [,1] [,2] [,3] [,4]  
[1,]    1    2    3    4  
[2,]    9    8    7    6
```

```
> x <- matrix(c(1,9,2,8,3,7,4,6), ncol=4, byrow=TRUE)
```

```
> x
```

```
      [,1] [,2] [,3] [,4]  
[1,]    1    9    2    8  
[2,]    3    7    4    6
```

```
> Z <- matrix(1:6, ncol=3) → 1 3 5
>
> dim(Z) # 行列のサイズ 2 4 6
[1] 2 3
>
> ncol(Z) # 列の数
[1] 3
>
> nrow(Z) # 行の数
[1] 2
>
> Z[1,2] # 要素の取り出し
[1] 3
>
> Z[1,] # 行の取り出し
[1] 1 3 5
>
> Z[,2] # 列の取り出し
[1] 3 4
>
> length(Z) # 全要素の数
[1] 6
```

繰り返し処理の書き方 (1) for ループ

```
for ( i in ベクトル ) {  
    繰り返し処理の内容  
}
```

制御変数 i はベクトルの要素の値を毎回とる
(変数名は i である必要はない)

ベクトルの次元数だけ処理が繰り返される

通常 ベクトルの与え方は $1:10$
などのような i が1つつ増えるように書くことが多い。

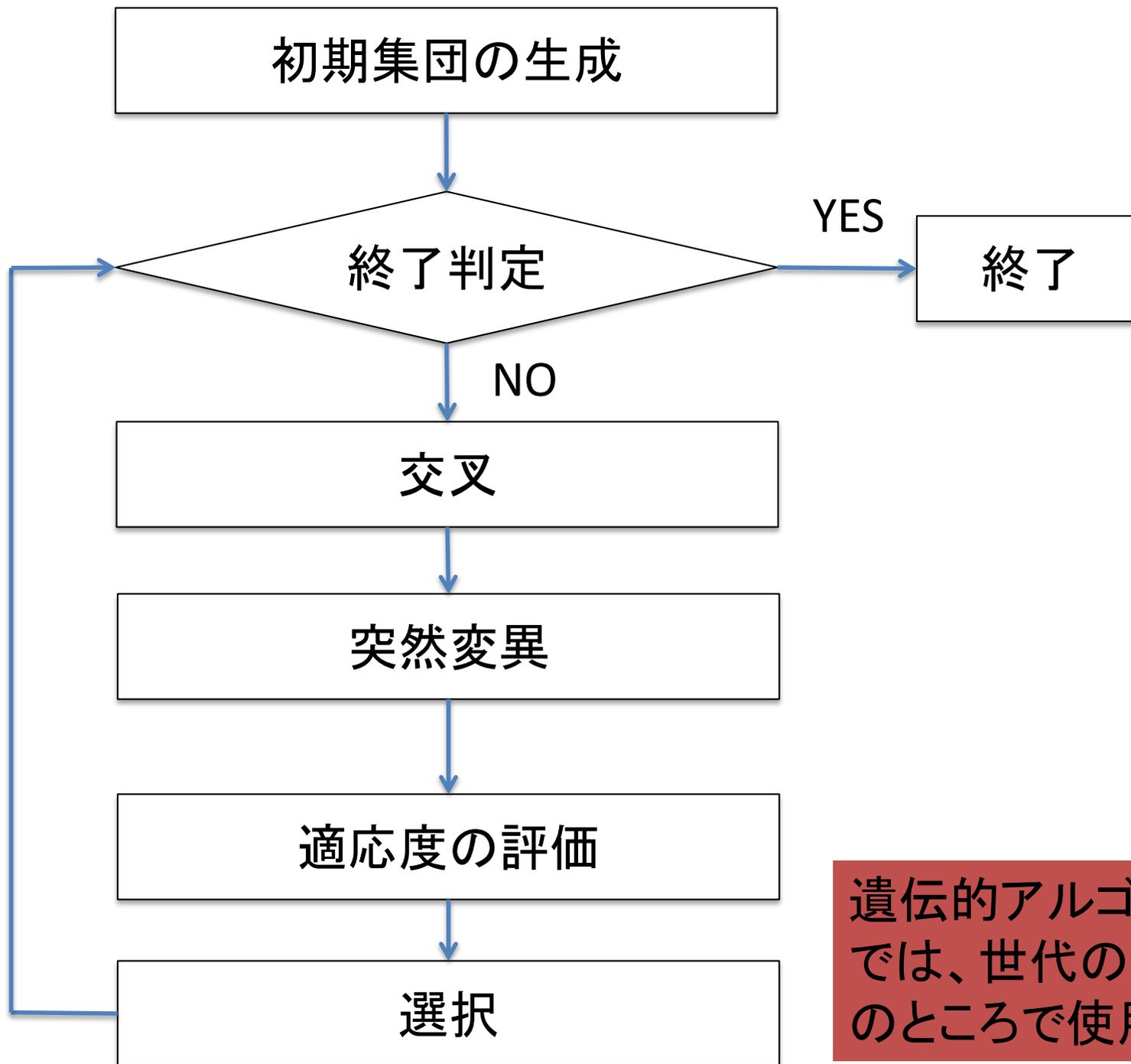
```
j <- 0
for (i in 1:5) {
  j <- j + i
}
print(j)
```

1から5までの数字を 足す処理

for 文で i は1から5まで
変化し、各 i について
{ } 内の処理が行われる。

i	左辺のj		右辺のj		i
i = 1	1	<-	0	+	1
i = 2	3	<-	1	+	2
i = 3	6	<-	3	+	3
i = 4	10	<-	6	+	4
i = 5	15	<-	10	+	5
	更新されたj		前のステップのj		

R上で 1+2+3+4+5 を実行して同じになるか確認しよう



遺伝的アルゴリズム
では、世代の繰り返しの
ところで使用

条件分岐

```
if (条件) {  
    条件が真 (TRUE) の時の処理  
}  
else {  
    条件が偽 (FALSE) の時の処理  
}
```

条件の書き方

$x == y$

一致

$x < y$

x が y より小

$x <= y$

x が y 以下

$x > y$

x が y より大

$x >= y$

x が y 以上

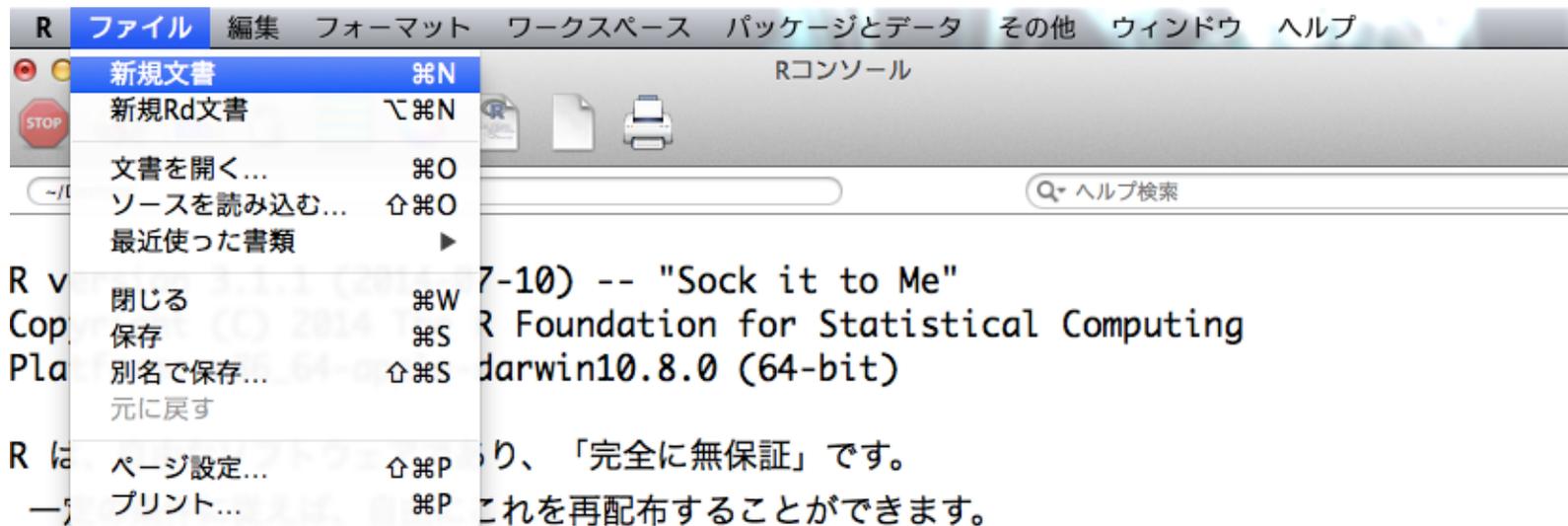
$x != y$

x は y と異なる

== と != は数値だけでなく、文字列についても使える

forループや今回学ぶ条件分岐など、複数行にわたる
スクリプトを書く場合、メモ帳やこれから説明する
Rのエディタに一旦書いて、それをRのコンソールに
コピーすると良い。

間違っているても修正が簡単



配布条件の詳細に関しては、`'license()'` あるいは `'licence()'` と入力してください。

R は多くの貢献者による共同プロジェクトです。

詳しくは `'contributors()'` と入力してください。

また、R や R のパッケージを出版物で引用する際の形式については `'citation()'` と入力してください。

`'demo()'` と入力すればデモをみることができます。

`'help()'` とすればオンラインヘルプが出ます。

`'help.start()'` で HTML ブラウザによるヘルプがみられます。

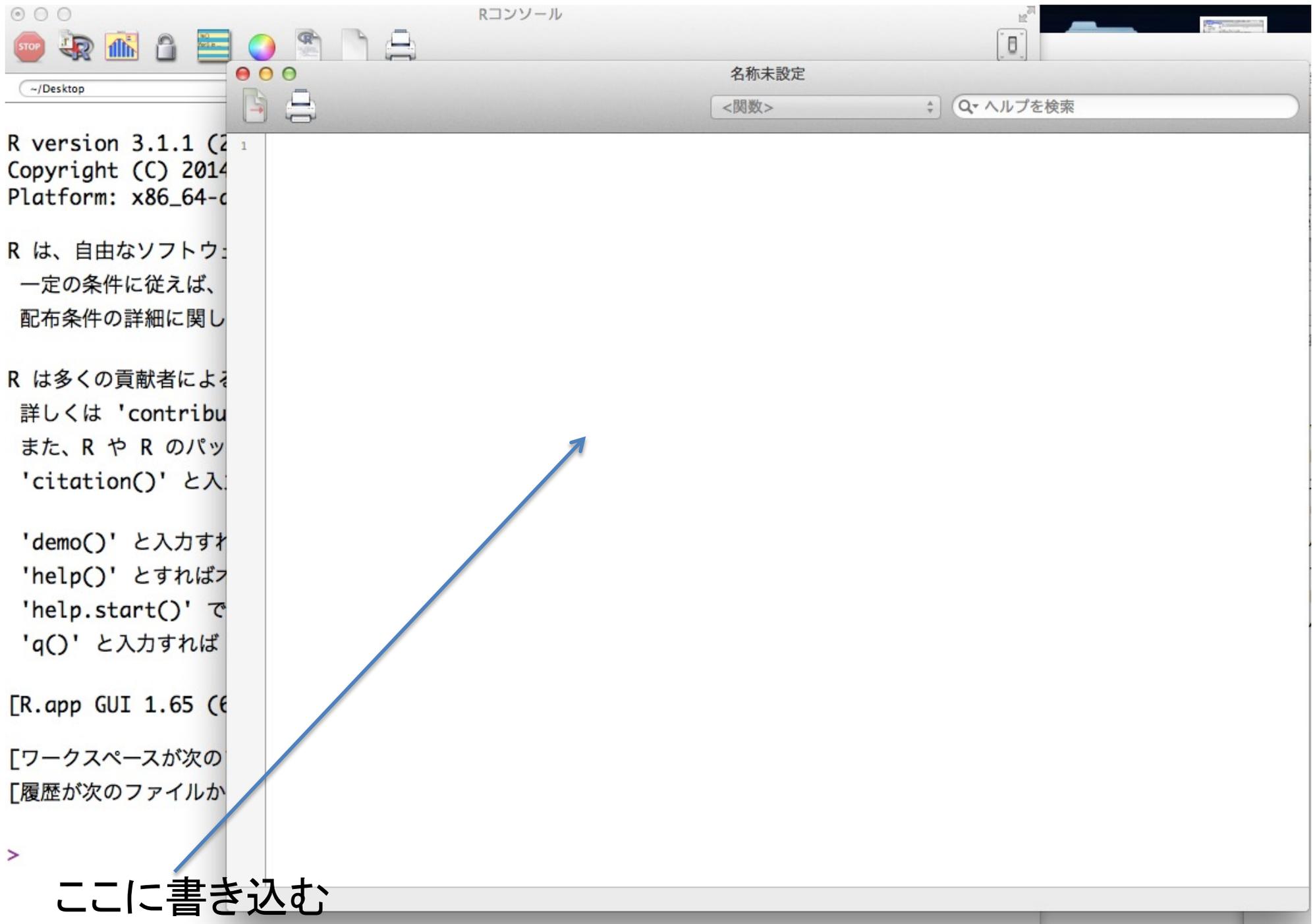
`'q()'` と入力すれば R を終了します。

```
[R.app GUI 1.65 (6784) x86_64-apple-darwin10.8.0]
```

```
[ワークスペースが次のファイルから読み込まれました /Users/toh/Desktop/.RData]
```

```
[履歴が次のファイルから読み込まれました /Users/toh/Desktop/.Rapp.history]
```

```
> |
```



R version 3.1.1 (2014-08-16)
Copyright (C) 2014 The R Foundation for Statistical Computing
Platform: x86_64-darwin13.0.0

R は、自由なソフトウェアです。
一定の条件に従えば、
配布条件の詳細に関し

R は多くの貢献者による
詳しくは 'contributor()' と入力すれば
また、R や R のパッケージについて
'citation()' と入力すれば

'demo()' と入力すれば
'help()' とすれば
'help.start()' で
'q()' と入力すれば

[R.app GUI 1.65 (64-bit)]

[ワークスペースが次のファイルに保存されています]
[履歴が次のファイルに保存されています]

>

ここに書き込む

The image shows a screenshot of an R environment. On the left is the R console, and on the right is the R script editor. The script editor contains the following code:

```
1 x <- c(120,10,34,909,21,3,67,221,87,112)
2 len <- length(x)
3 for ( i in 1:len) {
4   if (x[i] > 100) {
5     print(i)
6   }
7   else {
8     print("****")
9   }
10 }
```

The console shows the output of the script execution:

```
[1] 1
[1] "****"
[1] "****"
[1] 4
[1] "****"
[1] "****"
[1] "****"
[1] 8
[1] "****"
[1] 10
> x <- c(120,
> len <- leng
> for ( i in
+ if (x[i] >
+   print(i)
+ }
+ else {
+   print(x)
+ }
+ }
[1] 1
[1] 10
[1] 34
[1] 4
[1] 21
[1] 3
[1] 67
[1] 8
[1] 87
[1] 10
>
```

Two blue arrows point from the Japanese text to the script editor. One arrow points to the code block, and the other points to the console prompt.

コピーして

ペーストしてエンターをおす

```
x <- c(120,10,34,909,21,3,67,221,87,112)
len <- length(x)
for ( i in 1:len) {
  if (x[i] > 100) {
    print(i)
  }
  else {
    print("****")
  }
}
```

ベクトル `x`の要素が100より大なら
`i` を出力
さもなければ
`****` を出力

```
[1] 1
[1] "****"
[1] "****"
[1] 4
[1] "****"
[1] "****"
[1] "****"
[1] 8
[1] "****"
[1] 10
>
```

複数の条件を組み合わせる

条件1 && 条件2

条件1も条件2の真であれば真
その他は、偽

条件1 || 条件2

条件1か条件2のいずれかが真であれば真
両方が偽の時に偽

2条件だけでなく

(条件1 && 条件2) || 条件3

などのように複数の条件を組み合わせて利用できる

問題1

ベクトルの要素が偶数である以下、10以上70以下の時だけ、そのxの要素を書き出すコード

```
x <- c(120, 10, 34, 909, 21, 3, 67, 221, 87, 112)
len <- length(x)

for (i in 1:len) {
  if (  ) {
    print(x[i])
  }
  else {
    print(i)
  }
}
```

x[i] は偶数か、x[i]が10以上70以下
という条件はどうなるか？

偶数であること => 2で割った余りが0

Rには様々な関数が容易されているが自分のデータの処理には足りない場合がある。



自分のオリジナル関数を作成しよう

Rでは、自分の関数は次のように作成する。

```
関数名 <- function(引数1, 引数2, ...) {
```

括弧内の引数を使った処理を記述

```
  return (戻り値)
```

例えば x と y の積を求める関数は次のように書ける。

```
xy <- function(x, y) {  
  return (x*y)  
}
```

※ 引数を x , y にする必然性はなく、他の変数名でも構わない
処理の記述部分には引数として与えられた変数名で記述

プログラムのファイルは拡張子をRとしておく。
ドキュメントフォルダ (Zドライブ) におく。

`xy.R`を`source`で、実行中のコンソール画面から読み込むと
オリジナル関数`xy`を利用できるようになる。

```
> source("Z:xy.R")
> xy(4, 5)
[1] 20
> xy(1.41421356, 1.41421356)
[1] 2
```

問題2

forループを使って、与えられた2つのベクトルの内積をもとめるプログラムを作成してみる。

naiseki.Rを穴埋めする形で作成

```
naiseki <- function(a, b)  {  
  len <-   
  j <- 0  
  for (i in 1:len)  {  
    j <- j +   
  }  
  return (j)  
}
```

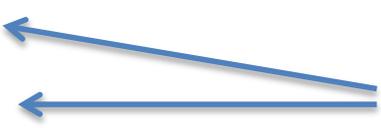
データフレーム

行列は1種類の型のみで構成される表形式のデータ

```
> x <- matrix(c(1,5,9,2),ncol=2)
> typeof(x)
[1] "double"
> x <- matrix(c("a","b","c","d","ef","gde"), ncol=3)
> typeof(x)
[1] "character"
> x
      [,1] [,2] [,3]
[1,] "a"  "c"  "ef"
[2,] "b"  "d"  "gde"
```

行列で数字と文字を混在させると、自動的に数字は文字型にされる

```
> x <- matrix(c("a", "b", "c", "d", 1, 2), ncol=3)
> x
      [,1] [,2] [,3]
[1,] "a"  "c"  "1"
[2,] "b"  "d"  "2"
> typeof(x)
[1] "character"
```



行列では文字型として扱われる

行列は1種類のデータ型のみで構成される!!

-----→ 複数のデータ型を混在させることのできる表形式のデータをRでは**データフレーム**とよぶ

データフレーム作成 1

(1) ベクトルとして作成したデータをつなげる

```
Ht <- c(165.6, 163.3, 172.4, 175.3)
Sx <- c("male", "female", "female", "male")
a <- data.frame(height=Ht, Sex=Sx)
```

a

	height	Sex
1	165.6	male
2	163.3	female
3	172.4	female
4	175.3	male

データフレーム作成 2

(2) ファイルからの読み込み

```
x <- read.table("data.txt")
```

```
x
```

	V1	V2	V3
1	1.95	3.3	RED
2	1.55	2.1	YELLOW
3	1.45	2.2	YELLOW
4	2.01	3.4	RED
5	2.13	3.0	RED
6	1.54	1.9	YELLOW
7	1.98	2.9	RED
8	1.47	1.8	YELLOW
9	1.34	1.7	YELLOW
10	1.36	1.8	YELLOW
11	2.23	3.1	RED

データフレーム作成 3

(2) データフレームの各列に名前をつける

```
names(x) <- c("activity", "size", "color")
```

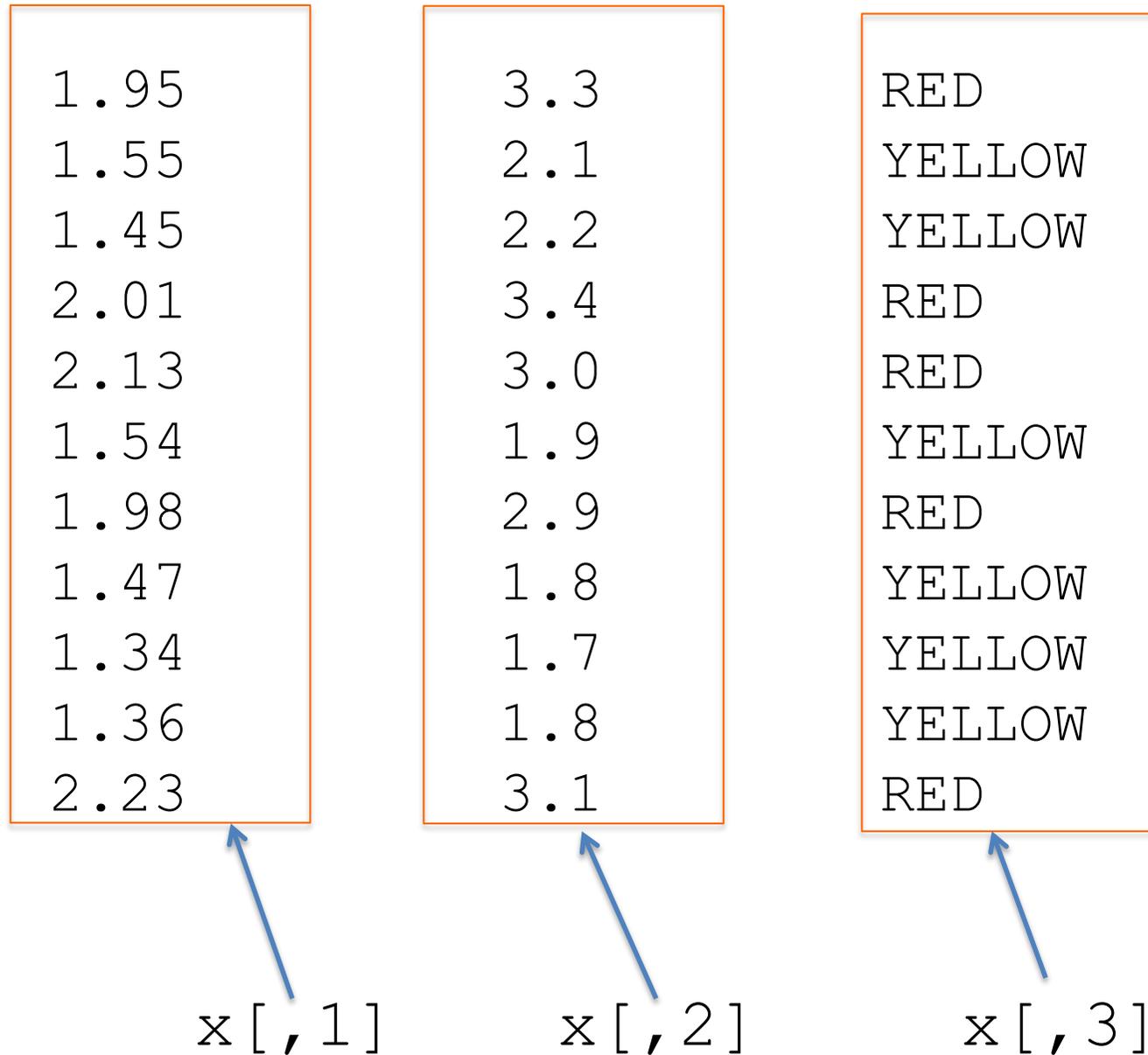
```
x
```

	activity	size	color
1	1.95	3.3	RED
2	1.55	2.1	YELLOW
3	1.45	2.2	YELLOW
4	2.01	3.4	RED
5	2.13	3.0	RED
6	1.54	1.9	YELLOW
7	1.98	2.9	RED
8	1.47	1.8	YELLOW
9	1.34	1.7	YELLOW
10	1.36	1.8	YELLOW
11	2.23	3.1	RED

```
names(x)
```

```
[1] "activity" "size"      "color"
```

データフレーム x



それぞれ、xの1列、2列、3列をベクトルとして抽出

データフレーム x

1.95	3.3	RED	$x[1,]$
1.55	2.1	YELLOW	
1.45	2.2	YELLOW	
2.01	3.4	RED	
2.13	3.0	RED	$x[5,]$
1.54	1.9	YELLOW	
1.98	2.9	RED	
1.47	1.8	YELLOW	
1.34	1.7	YELLOW	
1.36	1.8	YELLOW	
2.23	3.1	RED	$x[11,]$

x の1行、5行、11行をベクトルとして抽出

データフレーム x

x[1, 1]

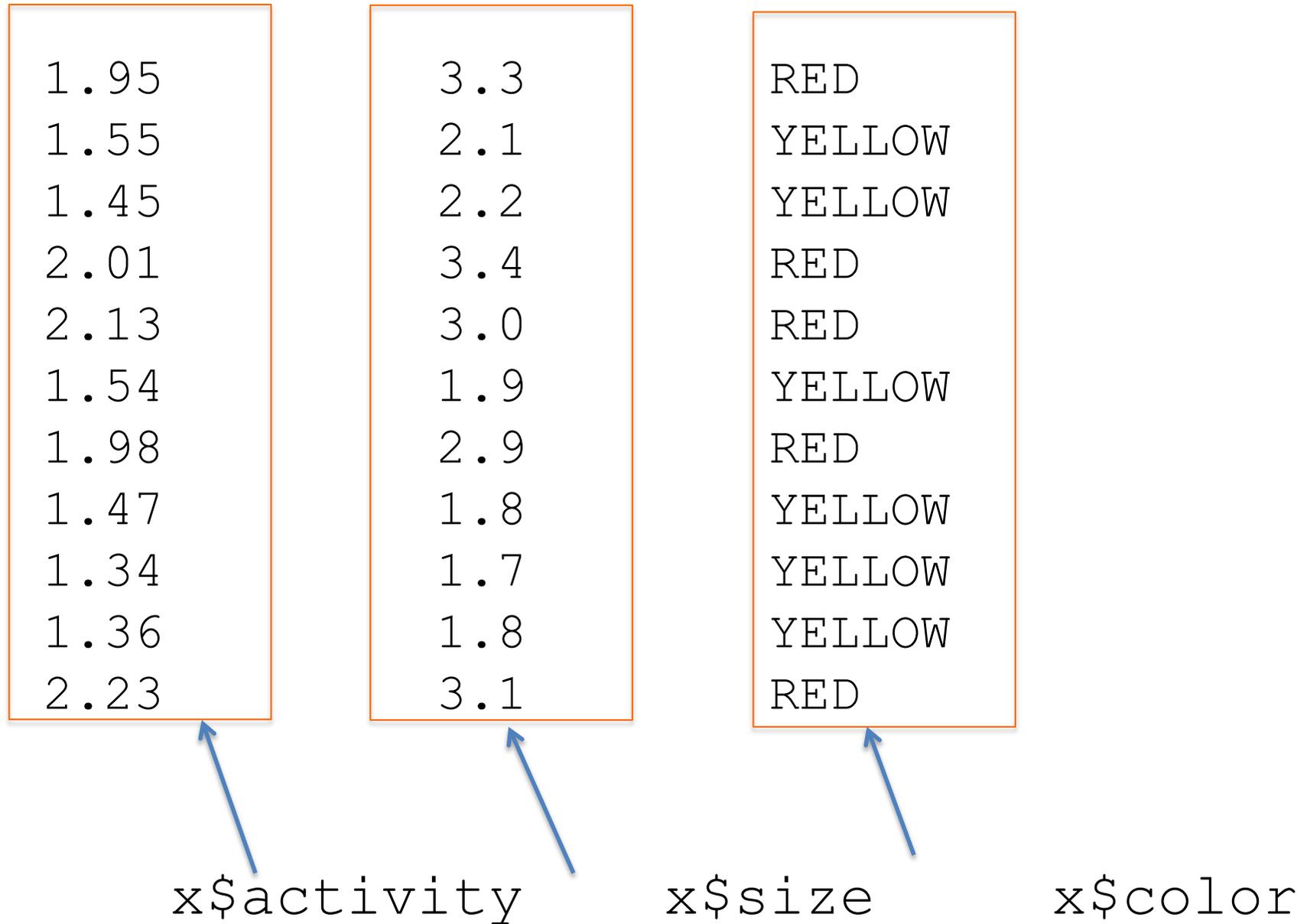
1.95	3.3	RED
1.55	2.1	YELLOW
1.45	2.2	YELLOW
2.01	3.4	RED
2.13	3.0	RED
1.54	1.9	YELLOW
1.98	2.9	RED
1.47	1.8	YELLOW
1.34	1.7	YELLOW
1.36	1.8	YELLOW
2.23	3.1	RED

x[5, 3]

x[11, 2]

添字でxの要素を添字を指定して抽出

データフレーム x



それぞれ、xの1列、2列、3列をベクトルとして抽出

よくある間違い

エラーが出た時にTAに聞く前に
まず、自分で確認してください！

- 大文字と小文字の書き間違い

例. 正source(“...”) 誤Source(“...”)

- スペルミス

例. 正source(“...”) 誤source(“...”)

- 0 (ゼロ) と O (オー)、1 (エル) と I (アイ) と 1 (数字)

例. 正“1j19.pdb” (1ジェイエル9) 誤“1j19.pdb” (エルジェイ19)

- 変更を保存せずに実行 → 変更が反映されません

→ プログラムに変更を加える度に保存する(CTRL + S)癖をつけてください

- ファイル名、ファイルの中では、なるべく**日本語/全角入力しない**こと

→ Rが正しくファイルを読み込んでくれません

例. 正source(“...”) 誤source(“...”)

↑
全角入力になっている

データフレームや行列では、行部分の条件式を書くとその条件を満たす行だけでできたデータフレームあるいは行列を作成できる

次のようにデータフレームdfを作成する

```
x <- c(1, 2, 2, 3)
y <- c('A', 'A', 'C', 'B')
z <- c(100, 101, 102, 103)
df <- data.frame(X=x, Y=y, Z=z)
```

ここで

```
df[df$Y=='A', ]
df[df$X==2, ]
```

と入力して、上の記述を確認しよう。また、この出力自体もデータフレームなので、

```
df[df$Y=='A', ]$X
df[df$X==2, ]$Y
```

のようにして、得られた行の中の必要な要素だけを取り出すことができる。

4. 相同性

今回、**相同**なタンパク質の立体構造の
重ね合わせを行う

相同なタンパク質とは何か？

立体構造の重ね合わせには相同配列の
アラインメントが必要

相同配列の形成

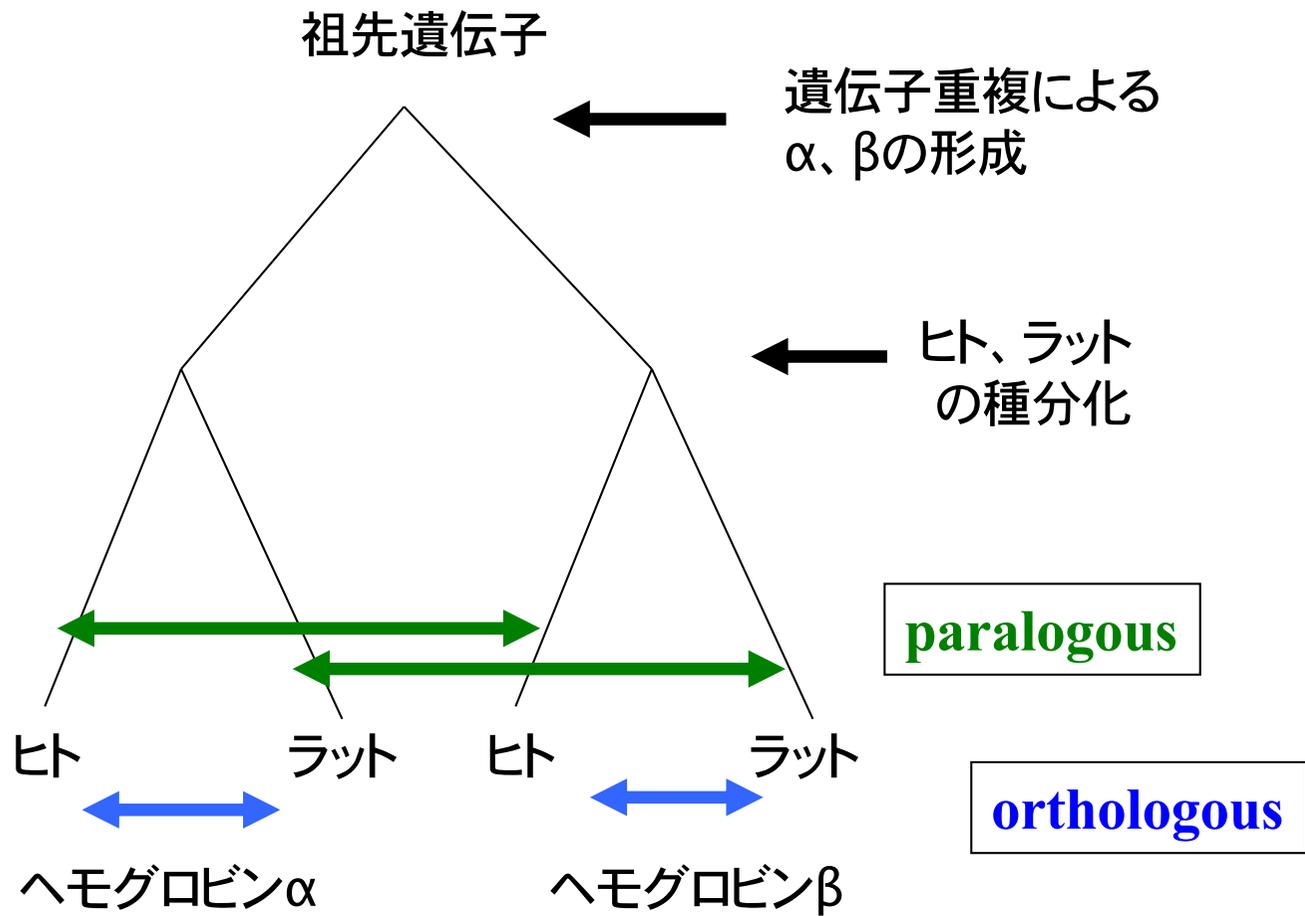
- ・種分化 ortholog
 - ・遺伝子重複 paralog
- 機能の多様化に特に重要

分子進化 (Molecular Evolution)

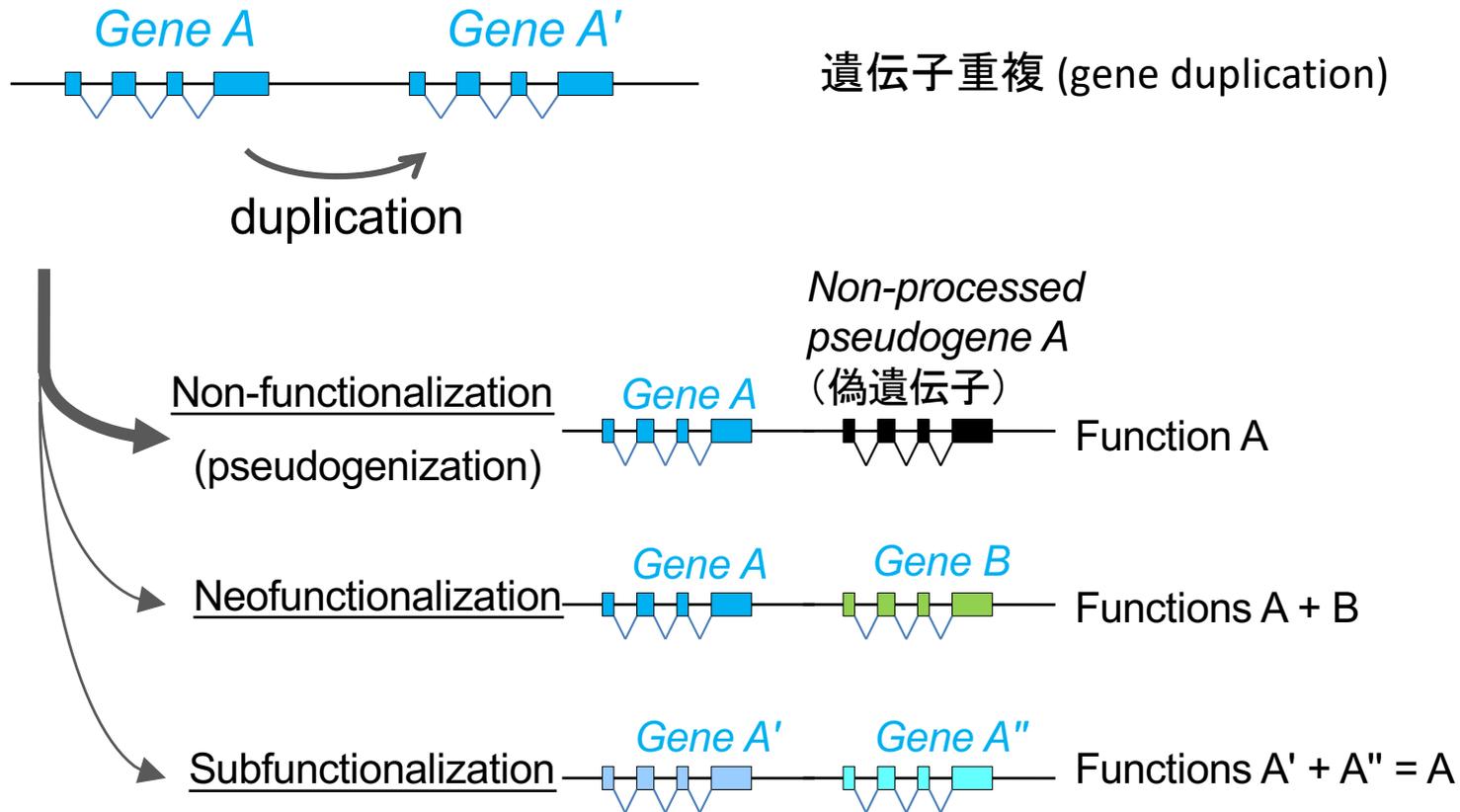
上の機構で分岐した遺伝子に突然変異が生じることで、配列が変化していくこと

基本ステップは

- 塩基(アミノ酸)置換
- 挿入/欠失

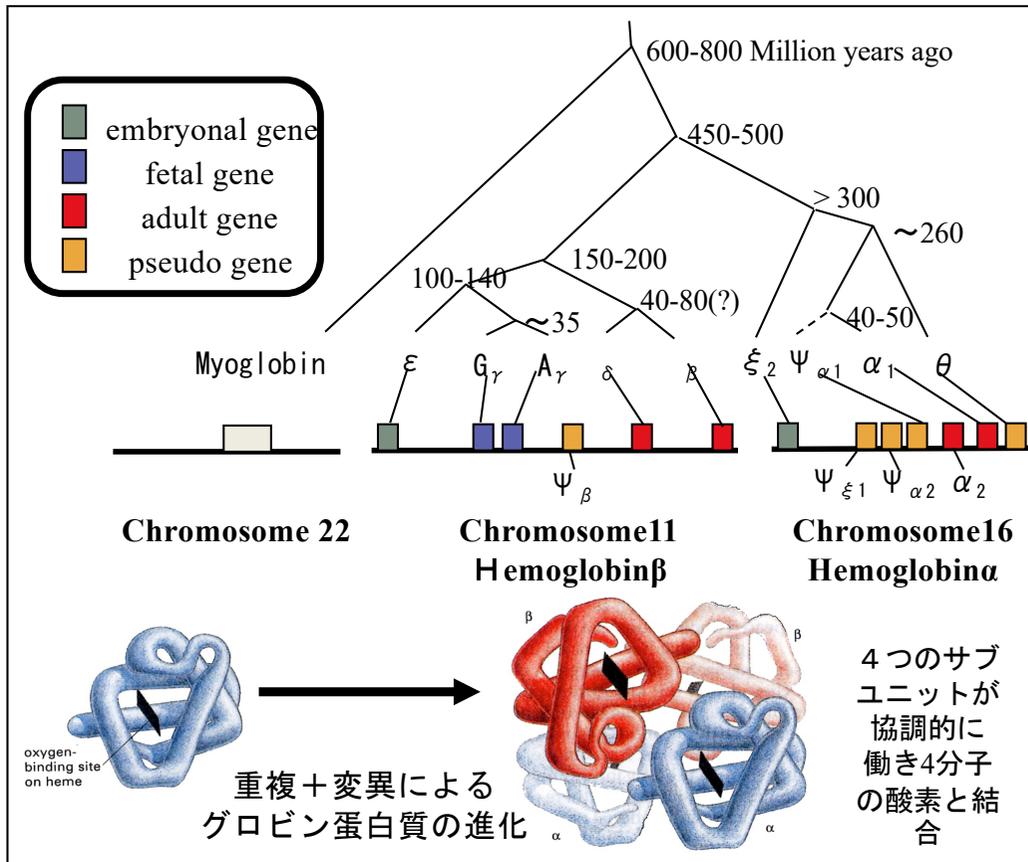


Evolutionary fate and functional consequence



●機能遺伝子の生成

・重複遺伝子が異なる機能を獲得した例 (globin superfamily)



* Myoglobin :

・モノマーで機能

* Hemoglobin :

・ヘテロテトラマーで機能

・各発生段階で発現するsubunitの組み合わせを変える

↓
酸素分子との親和性が変化し、段階に応じた機能が微調整される

ex.)

$\alpha_2\gamma_2$ (胎児期) >

$\alpha_2\beta_2, \alpha_2\delta_2$ (成体)

図1.5 点変異・挿入変異・欠失変異

(a) 点変異



(b) 挿入変異



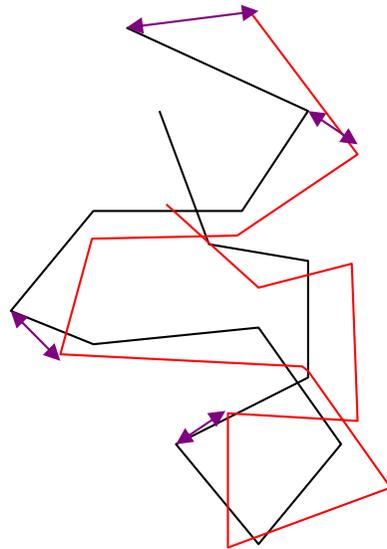
(c) 欠失変異



RMSD

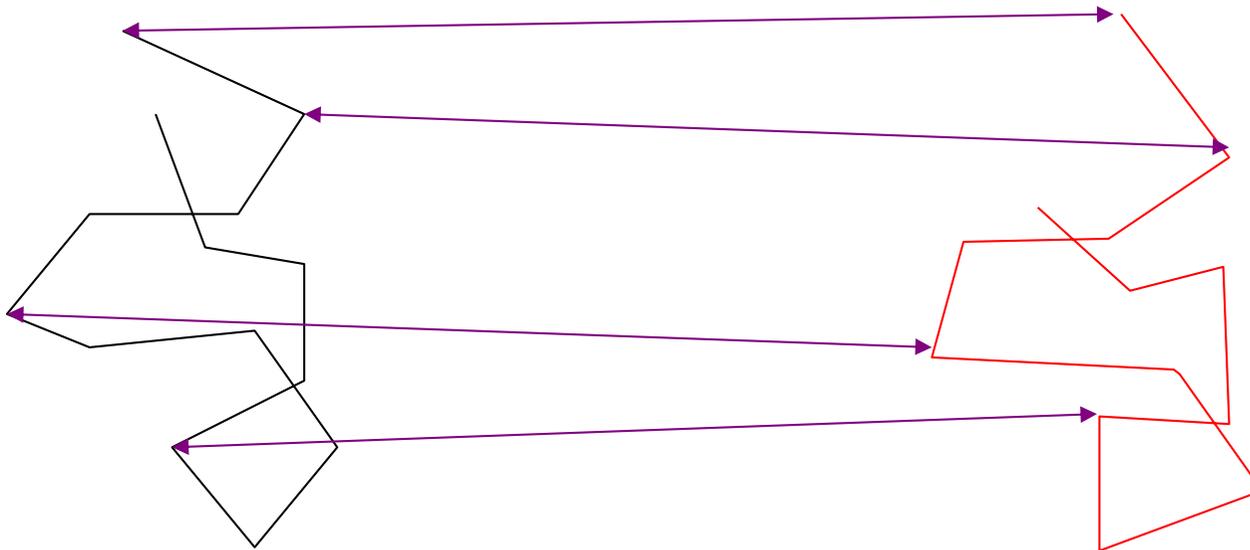
rmsd (root mean square distance) =

$$\sqrt{\frac{1}{n} \sum (dist(A(i), B(i)))^2}$$



残基間対応が最初に
与えられていると
計算は容易

構造比較の原点 - 重ね合わせ (superposition) -



対応するC α 原子間距離が最小になるように
二つの鎖を重ね合わせる (平行移動と回転)

McLachlan, A.D. (1972) *Nature New Biol.* 240, 83-85.

2. 配列アラインメント (sequence alignment)

進化の過程でのINDELを考慮しながら、相同な配列の間の対応する塩基(あるいはアミノ酸)を同じ位置に並べる操作あるいは、その操作によってできたもの。

INDELに対応してギャップ(gap)とよばれる空記号を挿入し位置をずらして、塩基やアミノ酸を対応づける。

通常、動的計画法(dynamic programming algorithm)や、そのバリエーションを用いて、配列間の類似度が高くとなるようにアラインメントが構築される。

アラインメントの原理は次回

マルチプルアラインメント (multiple alignment)

分子系統樹を構築するには、まず複数本の相同配列についてのマルチプルアラインメントを作成する。

リゾチームのアミノ酸配列

>LYC_HUMAN **ヒト** Lysozyme C
MKALIVLGLVLLSVTVQGVFERCELARTLKRLGMDGYRGISLANWMCLAKWESGYNTRATNYNAGDRST
DYGIFQINSRYWCNDGKTPGAVNACHLSCSALLQDNIADAVACAKRVVRDPQGI RAWVAWRNRCQNRDVR
QYVQCGV

>LYC1_BOVIN **ウシ** Lysozyme C 1
MKALIILGFLFLSVAVQGVFERCELARTLKKLGLDGYKGVSLANWLCLTKWESSYNTKATNYNPGSEST
DYGIFQINSKWWCNDGKTPNAVDGCHVSCSELMENDIAKAVACAKQIVSEQGITAWVAWKSHCRDHDVSS
YVEGCTL

>LYC_CHICK **ニワトリ** Lysozyme C
MRSLLILVLCFLPLAALGKVFGRCELAAMKRHGLDNYRGSYSLGNWVCAAKFESNFNTQATNRNTDGSTD
YGILQINSRWWCNDGRTPGSRNLCNIPCSALLSSDITASVNCAKKIVSDNGMNAWVAWRNRCKGTDVQA
WIRGCLR

>LYC2_ONCMY **マス** Lysozyme C II
MRAVVVLLLVAVASAKVYDRCELARALKASGMDGYAGNSLPNWVCLSKWESSYNTQATNRNTDGSTDYGI
FQINSRYWCDDGRTPGAKNVCGIRCSQLLTADLTVAIRC AKRVVLDPNGIGAWVAWRLHCQNQDLRSYVA
GCGV

>LYC_BOMMO **カイク** Lysozyme
MQKLIIFALVVLVCGSEAKTFTRCGLVHELKHKHGFEEENLMRNWVCLVEHESSRDTSKTNTNRNGSKDYGL
FQINDRYWCSKGASPGKDCNVKCSDLLTDDITKAAKCAKKIYKRHRFD AWYGWKNHCQGS L PDISSC

>LYSP_DROME **ハエ** Lysozyme P
MKAFLVICALTLTAVATQARTMDRCSLAREMSKLGVPRDQLAKWTCIAQHESSEFRGTGVGPANSNGSNDY
GIFQINNKYWCKPADGRFSYNECGLSCNALLTDDITNSVKCARKIQRQQGWTAWSTWKYCSGSLPSINSC
F

リゾチームのアミノ酸配列 マルチプル・アラインメント

二次構造情報を重ねて表示  α helix  β strand

ヒト
ウシ
マス
ニワトリ
カウ
ハエ



アラインメントを作成してみよう

MAFFTによるマルチプルアライメント

mafftは宮田研究室で開発され、加藤和貴によって継続的に開発されているマルチプルアライメントのフリーソフトウェア

海外の多くの研究機関で利用されている。

Web上でのアライメントサービスに加え、ダウンロードして自身のPC上で利用できる。Mac, Windows, Linuxなど様々なOSに対応している

ここでは、web serviceとして公開されているmafftを利用する。

<https://mafft.cbrc.jp/alignment/software/>

MAFFT version 7

Multiple alignment program for amino acid or nucleotide sequences



Download version

- [Mac OS X](#)
- [Windows](#)

Contact email address, kazutaka.katoh@aist.go.jp, is temporarily unavailable from 2018/Feb/7. If you sent an email to this address but have received no response, then please re-send the email to katoh@ifrec.osaka-u.ac.jp.

- [Linux](#)
- [Source](#)
- [Online version](#)
- [Alignment](#)
- [mafft --add](#)
- [Merge](#)
- [Phylogeny](#)
- [Rough tree](#)
- [Merits / limitations](#)
- [Algorithms](#)
- [Tips](#)
- [Benchmarks](#)
- [Feedback](#)
- [Follow](#)

ABOUT

MAFFT is a multiple sequence alignment program for unix-like operating systems. It offers a range of multiple alignment methods, L-INS-i (accurate; for alignment of <-200 sequences), FFT-NS-2 (fast; for alignment of <-30,000 sequences), etc.

Download and Installation

- [Mac OS X](#)
- [Linux](#)
- [Windows](#)
- [Source](#)
- [Changelog](#)

The latest version is 7.471, 2020/Jul.

A bug in parsing input filename has been fixed. Please use 7.470 or higher. (2020/Jun)

Input Format

Fasta format. [example1 \(LSU rRNA\)](#), [example2 \(protein\)](#)

The type of input sequences (amino acid or nucleotide) is automatically recognized.

Usage

```
% mafft [arguments] input > output
```

An alias for an accurate option (L-INS-i) for an alignment of up to ~200 sequences x ~2,000 sites:

```
% mafft-linsi input > output
```

A fast option (FFT-NS-2) for a larger sequence alignment:

```
% mafft input > output
```

If not sure which option to use,

```
% mafft --auto input > output
```

- [Manual \(v6.240\)](#)
- [Tips \(not yet included in the manual\)](#) for large alignment, ncRNA alignment, profile alignment, etc.

Related Resources

- [MAFFT server](#) at EBI
- [MAFFT server](#) at the MPI Bioinformatics Toolkit
- [ClustalW / MAFFT / RDPH at GenomeNet](#)

Online Versionをクリック

- [Linux](#)
- [Source](#)
- [Online version](#)
- [Alignment](#)
- [mafft --add](#)
- [Merge](#)
- [Phylogeny](#)
- [Rough tree](#)
- [Merits / limitations](#)
- [Algorithms](#)
- [Tips](#)
- [Benchmarks](#)
- [Feedback](#)
- [Follow](#)

ABOUT

MAFFT i
<-30,00

Downlo

- [M](#)
- [L](#)
- [W](#)
- [S](#)
- [C](#)

The

Download version

- [Mac OS X](#)
- [Windows](#)
- [Linux](#)
- [Source](#)

Online version

- Alignment**
- [mafft --add](#)
- [Merge](#)
- [Phylogeny](#)
- [Rough tree](#)
- [Merits / limitations](#)
- [Algorithms](#)
- [Tips](#)
- [Benchmarks](#)
- [Feedback](#)



To avoid overload, try [a light-weight option](#), for MSA of full-length SARS-CoV-2 genomes (2020/Apr).

For a large number of short sequences, try [an experimental service](#).

[Experimental service for aligning raw reads \(2019/Aug\)](#)

Multiple sequence alignment and NJ / UPGMA phylogeny

Input:

Paste protein or DNA sequences in fasta format. [Example](#)

ファイルを選択をクリック

or upload a **plain text** file:

- Use [DASH](#) to add homologous structures (protein only) **New! 2018/Dec/23**
 - Output original plus DASH sequences
 - Output original sequences only
- Give structural alignment(s) externally prepared
- Allow unusual symbols (Selenocysteine "U", Inosine "i", non-alphabetical characters, etc.) [Help](#)

UPPERCASE / lowercase:

- Same as input
- Amino acid → UPPERCASE / Nucleotide → lowercase

Direction of nucleotide sequences: [Help](#)

- Same as input
- Adjust direction according to the first sequence (accurate enough for most cases)
- Adjust direction according to the first sequence (only for highly divergent data; **extremely slow**)

Output order:

- Same as input
- Aligned

Job name (optional; used as output file name and subject of emails):

(basic Latin alphabet, number and space only)

Notify when finished (optional; recommended when submitting large data):

Email address:

ファイルチューザからRetroProtease.fastaを選択

MAFFT version 7

Multiple alignment program for amino acid or nucleotide s

Download version

[Mac OS X](#)

[Windows](#)

[Linux](#)

[Source](#)

Online version

Alignment

[mafft --add](#)

[Merge](#)

[Phylogeny](#)

[Rough tree](#)

[Merits / limitations](#)

[Algorithms](#)

[Tips](#)

[Benchmarks](#)

[Feedback](#)

[Follow](#)



To avoid overlo

For a large number of short

[Experimental service for alig](#)

Multiple sequence alig

Input:

Paste protein or DNA s

or upload a **plain text** file: ファイル未選択

名前	変更日	種類	サイズ
7HVP.pdb	2015年9月22日 16:08	書類	105 KB
Bioinformatics2	2019年10月8日 14:46	Adobe...cument	7 MB
Bioinformatics2pptx	2019年10月2日 0:58	PowerP...(pptx)	5.2 MB
FoxP2.aln	2015年9月22日 16:38	書類	12 KB
FoxP2.fasta	2015年9月22日 16:33	MEGAX書類	8 KB
FOX2.PDF	2015年1月12日 4:46	Adobe...cument	2.4 MB
retrop.pse	2015年9月22日 16:22	書類	215 KB
RetroProtease.aln	2015年9月22日 16:03	書類	1 KB
RetroProtease.fasta	2015年9月22日 16:02	MEGAX書類	879 バイト

MAFFT version 7

Multiple alignment program for amino acid or nucleotide sequences

[Download version](#)

[Mac OS X](#)

[Windows](#)

[Linux](#)

[Source](#)

Online version

Alignment

[mafft --add](#)

[Merge](#)

[Phylogeny](#)

[Rough tree](#)

[Merits / limitations](#)

[Algorithms](#)

[Tips](#)

[Benchmarks](#)

[Feedback](#)

[Follow](#)



To avoid overload, try [a light-weight option](#), for MSA of full-length SARS-CoV-2 genomes (2020/Apr).

For a large number of short sequences, try [an experimental service](#).

[Experimental service for aligning raw reads \(2019/Aug\)](#)

Multiple sequence alignment and NJ / UPGMA phylogeny

Input:

Paste protein or DNA sequences in fasta format. [Example](#)

or upload a **plain text** file: RetroProtease.fasta

Use **DASH** to add homologous structures (protein only) **New! 2018/Dec/23**

Output original plus DASH sequences Output original sequences only

Give structural alignment(s) externally prepared

Allow unusual symbols (Selenocysteine "U", Inosine "i", non-alphabetical characters, etc.) [Help](#)

UPPERCASE / lowercase:

Same as input

Amino acid → UPPERCASE / Nucleotide → lowercase

Direction of nucleotide sequences: [Help](#)

Same as input

Adjust direction according to the first sequence (accurate enough for most cases)

Adjust direction according to the first sequence (only for highly divergent data; **extremely slow**)

Output order:

Same as input

Aligned

Job name (optional; used as output file name and subject of emails):

(basic Latin alphabet, number and space only)

Notify when finished (optional; recommended when submitting large data):

Email address:

①ファイルが選択
されていることを
確認

②パラメータは
デフォルトのまま
Submitをクリック

CLUSTAL 形式 でアラインメント が表示される

[Clustal format](#) | [Fasta format](#) | [MAFFT result](#) | [View](#) | [Tree](#) | [Refine dataset](#) | [Return to home](#)

View

Reformat to GCG, PHYLIP, MSF, NEXUS, uppercase/lowercase, etc. with Readseq

GUIDANCE2 computes the residue-wise confidence scores and extracts well-aligned residues.

Refine dataset

Phylogenetic tree

MAFFT-L-INS-i Result

```
CLUSTAL format alignment by MAFFT (v7.471)

gi|443546|pdb|7 PQITLW-----QRPLVTIRIGGQ-----LKEALLDTGADDTVLEEMNLPG
simian          ---SLW-----NRPTTVVEIEGQ-----KVEALLDTGADDTVIKDLDLKG
HIV2           -----VTAYIEDQ-----PVEVLDTGADDSIVAGIELGD
gi|4389337|pdb| LAMTMEHK-----DRPLVRVILTNTGSHPVKQRSVYITALLDTGADDTVISEEDWPT
gi|224443|prf|| ---TLDDQGGQGEPPPEPRITLKVGGQ-----PVTFLVDTGAQHSVLTQNPGL
                : .                               *:***:..::

gi|443546|pdb|7 KW-----KPKMIGGIGGFIKVRQ---YDQIPVEIXGHKAIGTVL---VGPTPVNIIGR
simian          NW-----KPQIIGGIGGS INVKQ---FFNCKVTIAGKTTHASVL---VGPTPVNIVGR
HIV2           NY-----TPKIVGGIGGFINTKE---YKNVEIKVLNKRVRATIM---TGDTPINIFGR
gi|4389337|pdb| DWPVMEANPQ-IHGIGGGIPVRKSRDMELGVINRDGSLERPLLFPVAMTPVNIIGR
gi|224443|prf|| SD-----KSAWVQGATGGKRYRW---TDRKVHLATGKVTHSFLH---VPDCPYPLGR
                .   .. : * * : : : : .. . * :.*

gi|443546|pdb|7 NLLTQIGXTLN-----F
simian          NVLKKLGCTLN-----
HIV2           NILT-----
gi|4389337|pdb| DCLQGLGLRLT-----NL
gi|224443|prf|| DLLTKLKAQIHFEQSGAQVMGPMGQPLQVL
                : *
```

Method

L-INS-i (Probably most accurate, very slow)

```
% mafft --reorder --auto input
```

LUSTAL format alignment by MAFFT L-INS-i (v7.130b)

```
gi|443546|pdb|7 PQITLW-----QRPLVTIRIGGQL-----KEALLDTGADDTVLEEMNLP  
HIV2              -----VTAYIEDQP-----VEVLLDTGADDSIVAGIELGD  
simian           ---SLW-----NRPTTVVEIEGQK-----VEALLDTGADDTVIKDLDLK  
gi|4389337|pdb| LAMTMEHK-----DRPLVRVILTNTGSHPVKQRSVYITALLDTGADDTVISEEDWPT  
gi|224443|prf|| ---TLDDQGGQGQEPPEPRITLKVGGQP-----VTFVLVDTGAQHSVLTQNPGPL  
                  : .                               *:****:..::  
  
gi|443546|pdb|7 KW-----KPKMIGGIGGFIKVRQ---YDQIPVEIXGHKAIGTVL----VGPTPVNIIGR  
HIV2              NY-----TPKIVGGIGGFINTKE---YKNVEIKVLNKRVRATIM----TGDTPINIFGR  
simian           NW-----KPQIIGGIGGSINVKQ---FFNCKVTIAGKTTHASVL----VGPTPVNIVGR  
gi|4389337|pdb| DWPVMEANPQ-IHGIGGGIPVRKSRDIELGVINRDGSLERPLLLFPLVAMTPVNIIGR  
gi|224443|prf|| SD-----KSAWVQGATGGKRYRW---TTDRKVHLATGKVTHSFLH---VPDCPYPLIGR  
                  .      .. : * * : : : .. : . * : **  
  
gi|443546|pdb|7 NLLTQIGXTLN-----F  
HIV2              NILT-----  
simian           NVLKKLGCTLN-----  
gi|4389337|pdb| DCLQGLGLRLT-----NL  
gi|224443|prf|| DLLTKLKAQIHFEGSGAQVMGPMGQPLQVL  
                  : *
```

Clustal形式のアラインメント

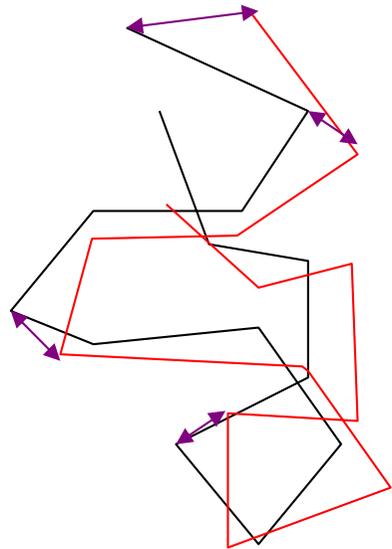
強く保存しているセグメント(モチーフ)が2ヶ所見いだされる

5. 遺伝的アルゴリズムによる 相同タンパク質の立体構造の重ね合わせ

RMSD

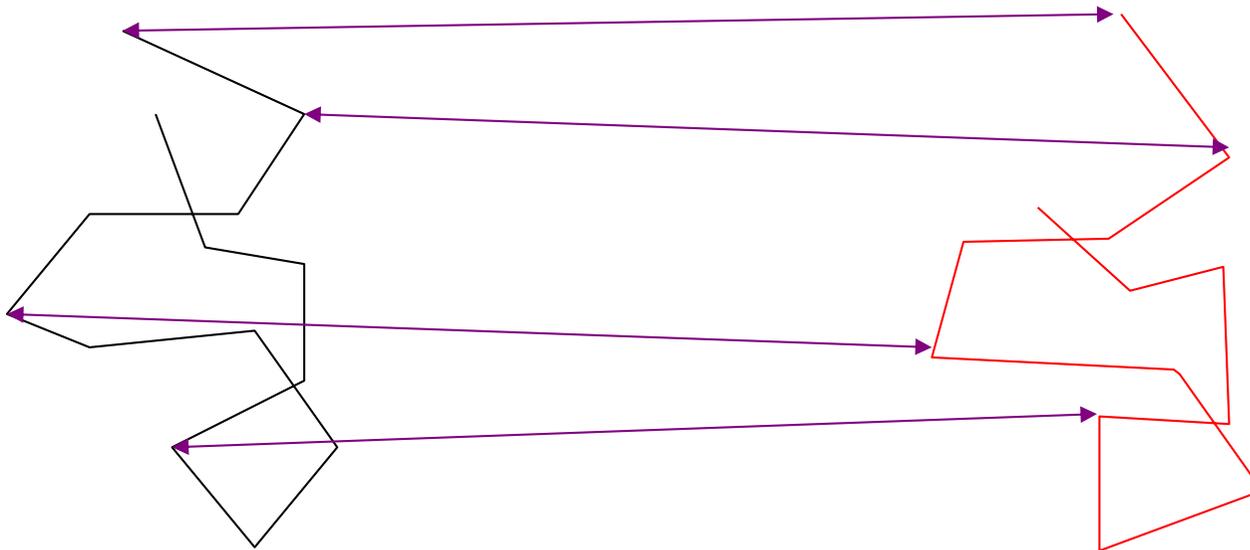
rmsd (root mean square distance) =

$$\sqrt{\frac{1}{n} \sum (dist(A(i), B(i)))^2}$$



残基間対応が最初に
与えられていると
計算は容易

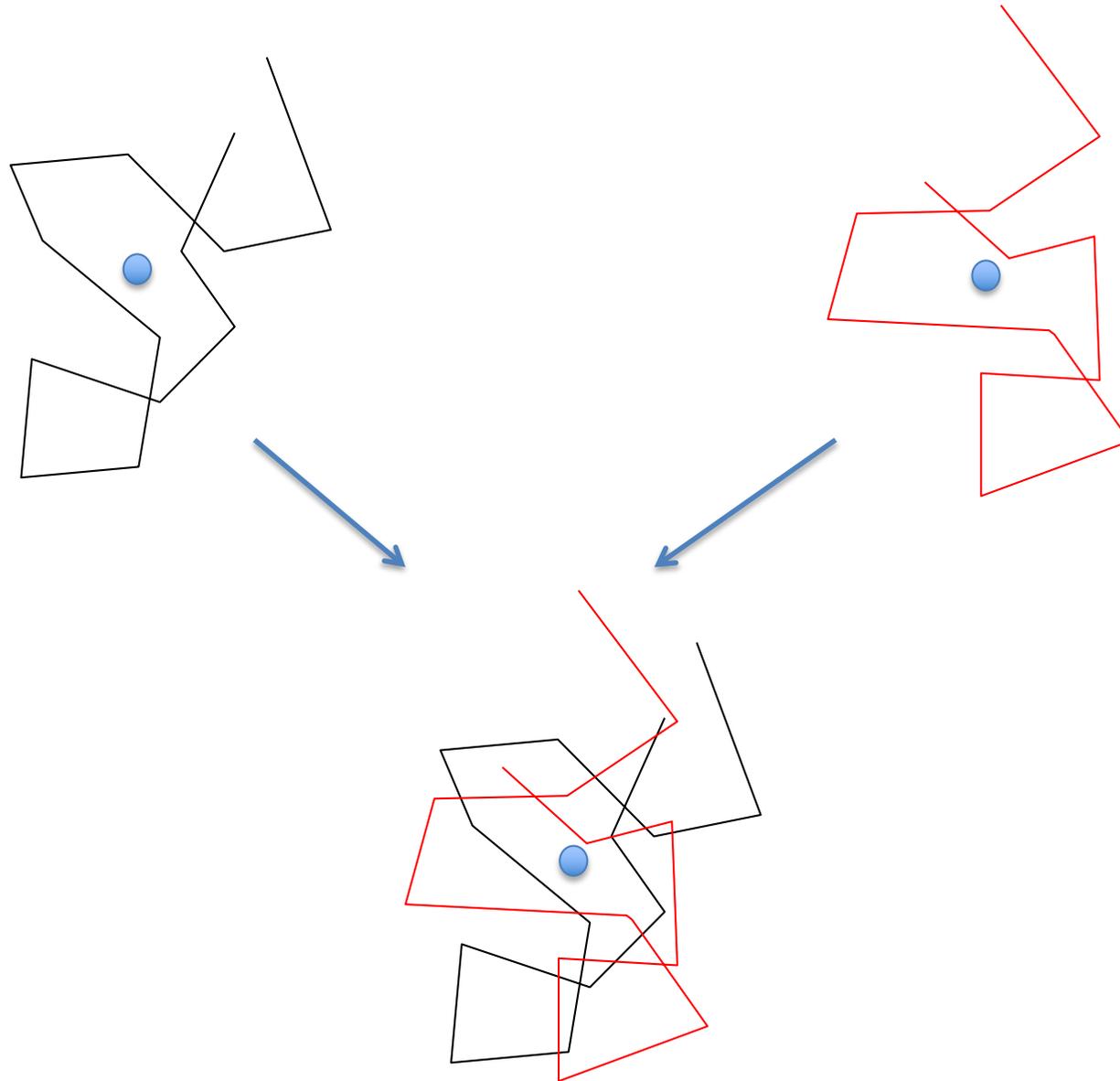
構造比較の原点 - 重ね合わせ (superposition) -

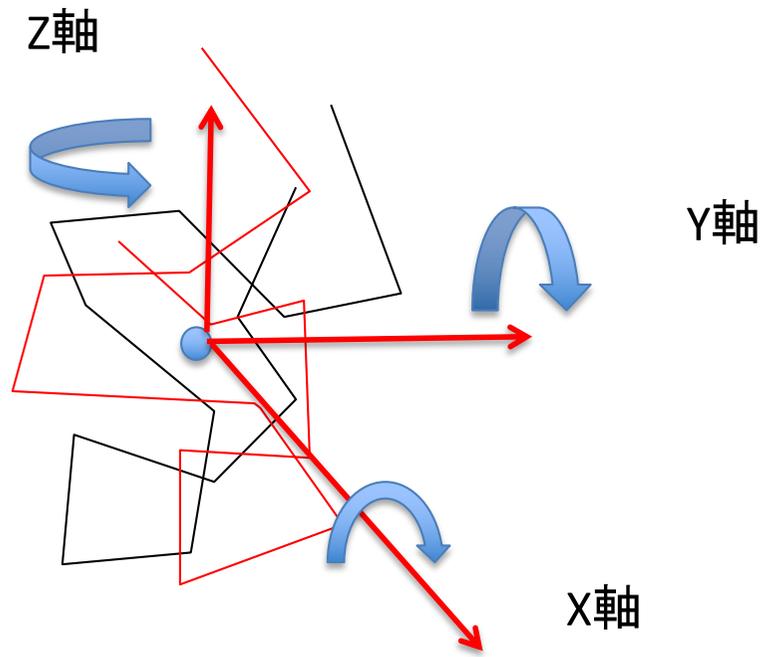


対応するC α 原子間距離が最小になるように
二つの鎖を重ね合わせる (平行移動と回転)

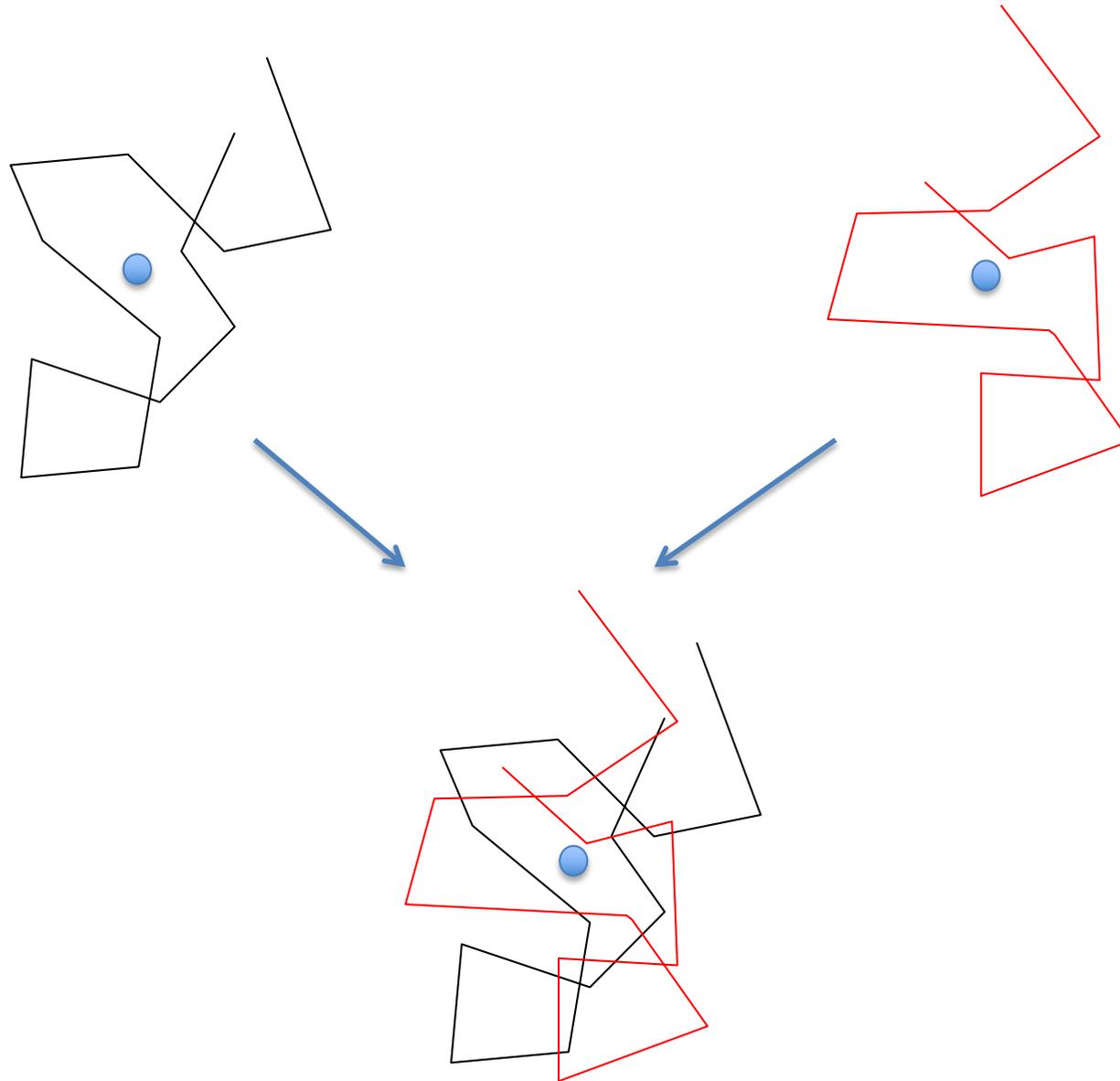
McLachlan, A.D. (1972) *Nature New Biol.* 240, 83-85.

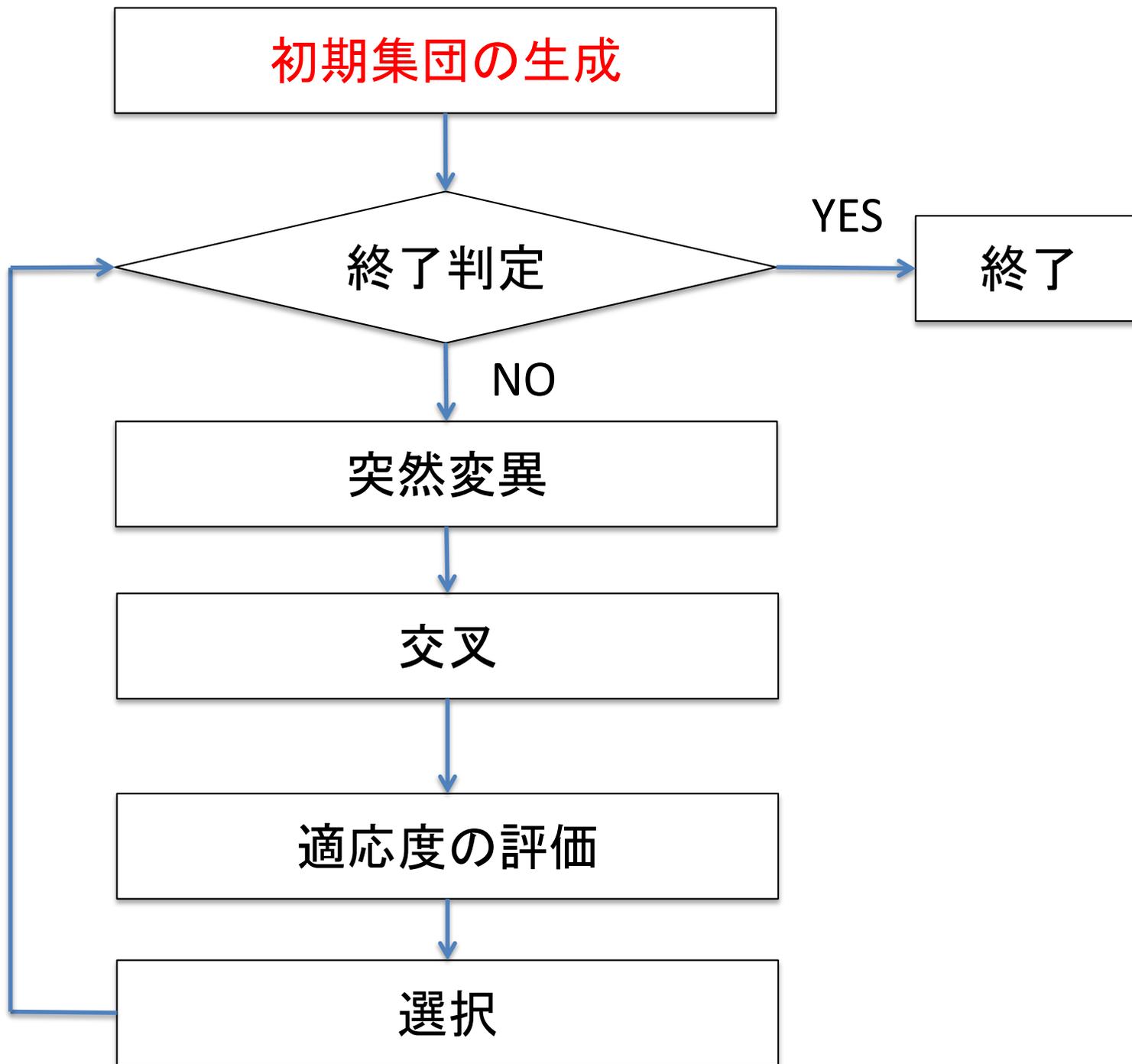
$C\alpha$ の幾何重心を求め、それぞれ幾何重心の位置を原点とする





$C\alpha$ の幾何重心を求め、それぞれ幾何重心の位置を原点とする





遺伝子の設計

実数コード

X軸周りの 回転角度
Y軸周りの 回転角度
Z軸周りの 回転角度

遺伝子1

遺伝子2

遺伝子3

染色体

集団の生成

```
population <- matrix(runif(popsiz*3)*2*pi,ncol=3)
```

個体1

X1	Y1	Z1
----	----	----

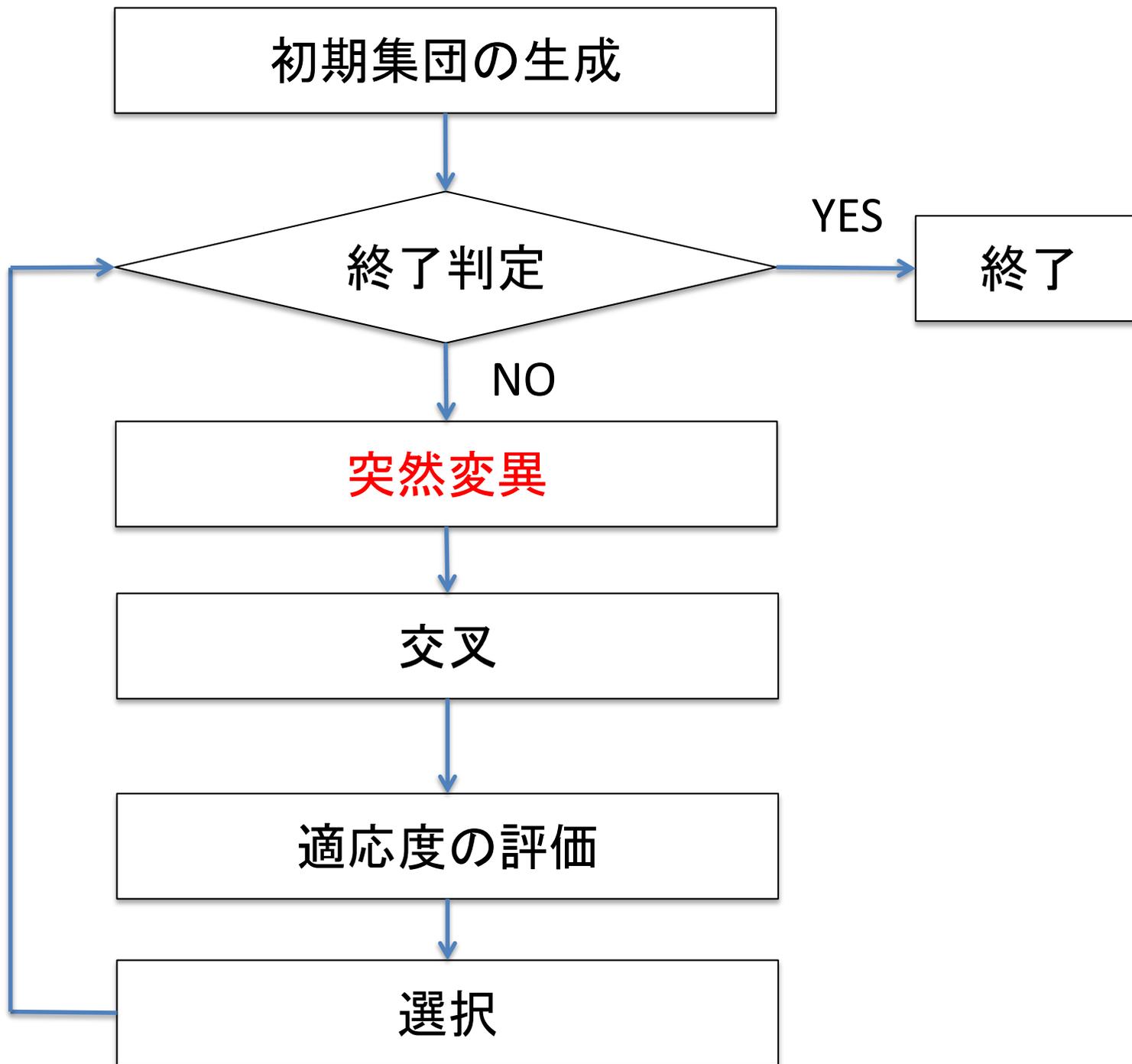
個体2

X2	Y2	Z2
----	----	----

...

個体popsiz

X_{popsiz}	Y_{popsiz}	Z_{popsiz}
---------------------	---------------------	---------------------



初期集団の生成

終了判定

YES

終了

NO

突然変異

交叉

適応度の評価

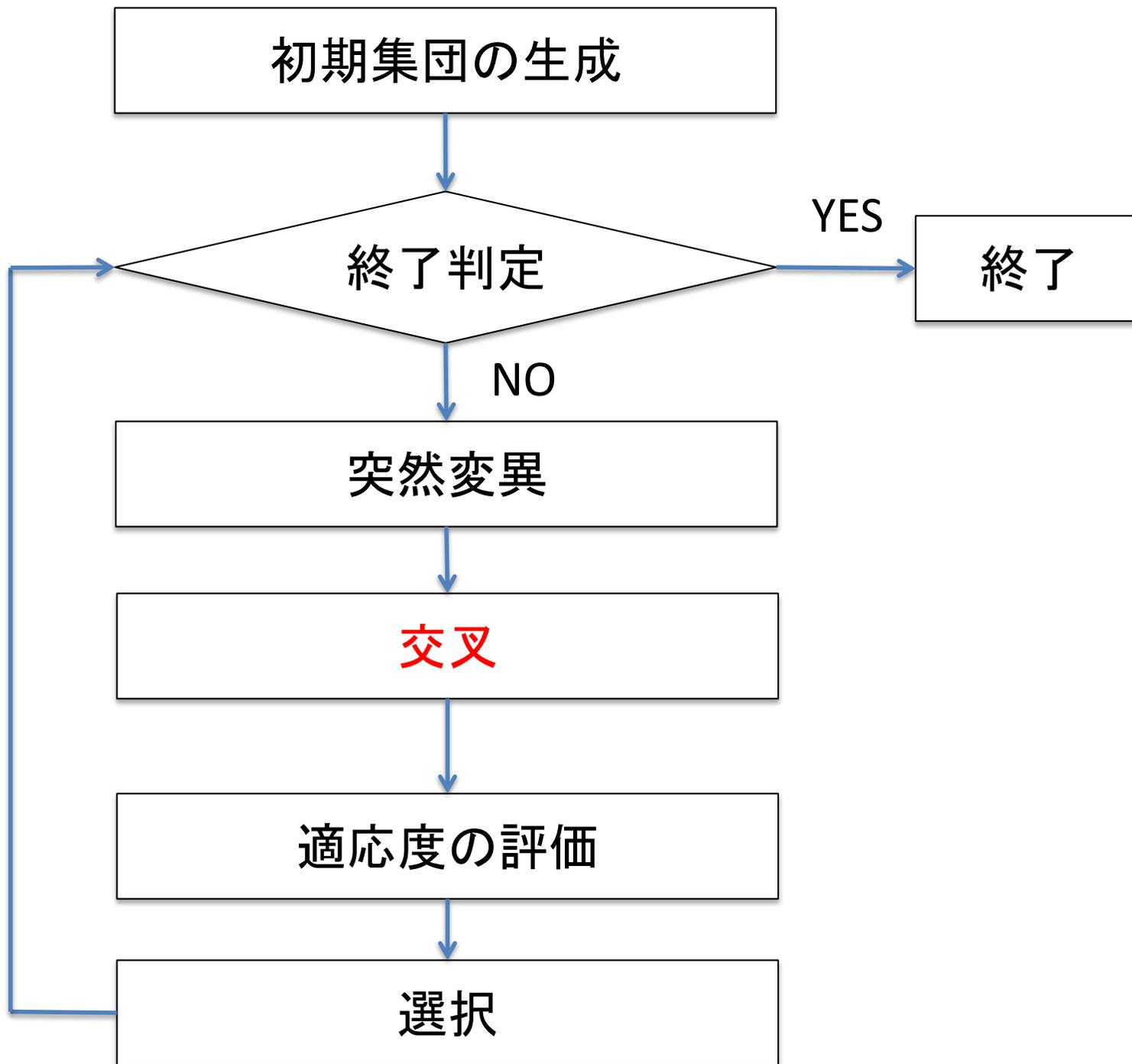
選択

個体 i

X_i	Y_i	Z_i
-------	-------	-------

突然変異個体 j

$X_i + \text{delta}$	Y_i	$Z_i - \text{delta}$
----------------------	-------	----------------------



個体 i

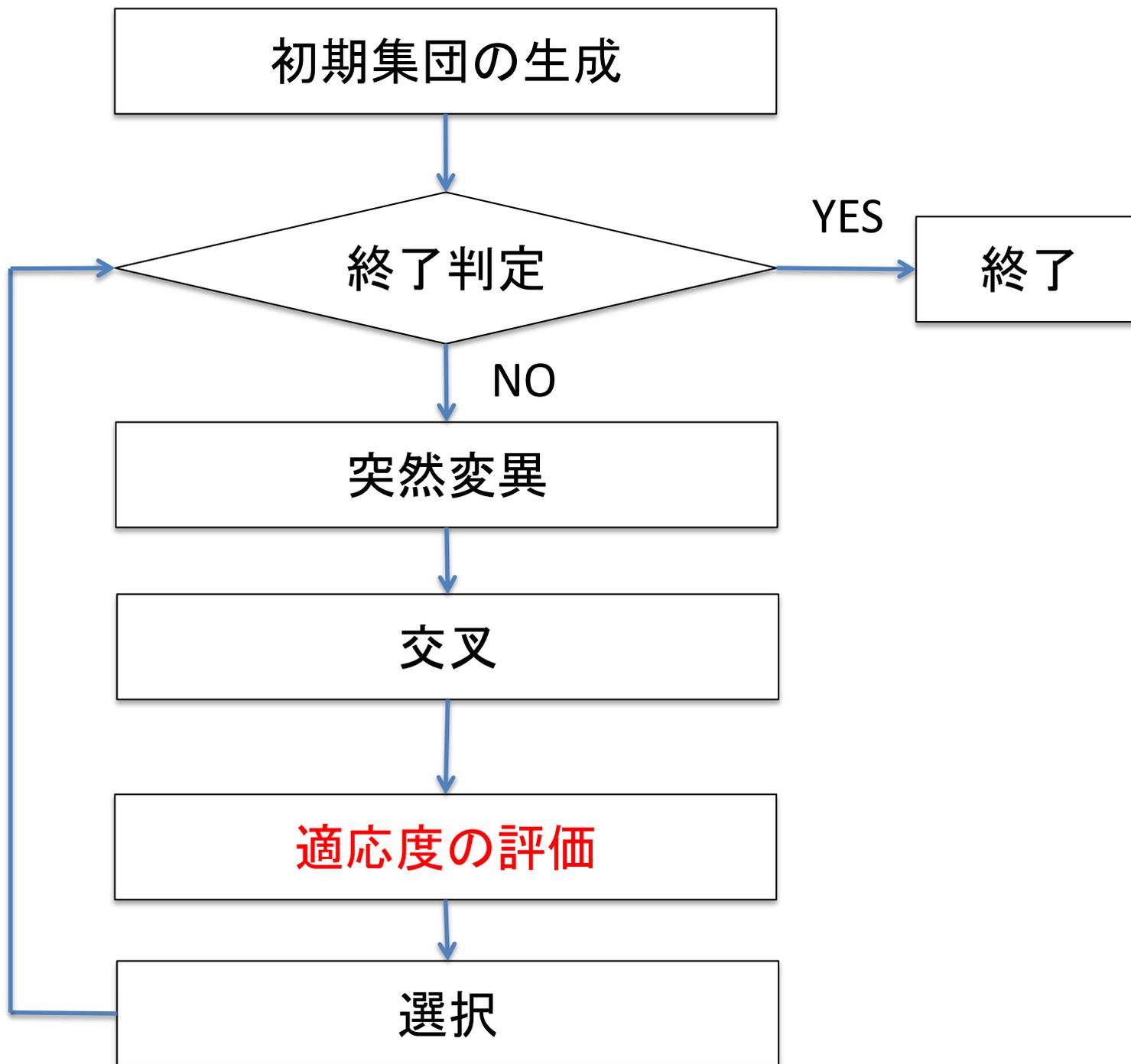
x_i	y_i	z_i
-------	-------	-------

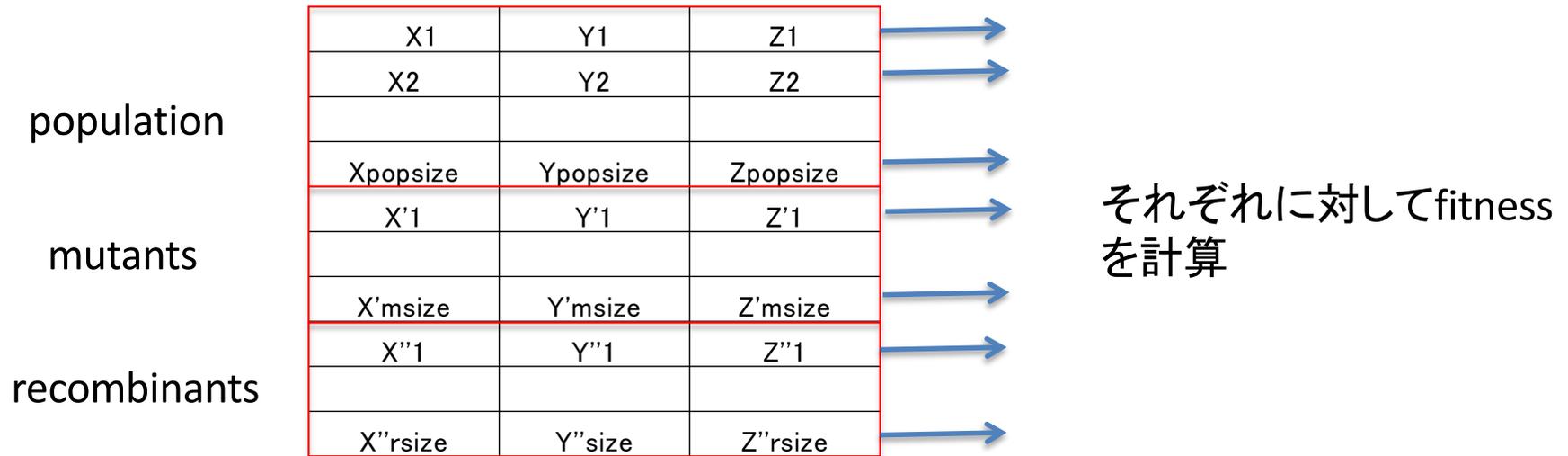
個体 j

x_j	y_j	z_j
-------	-------	-------

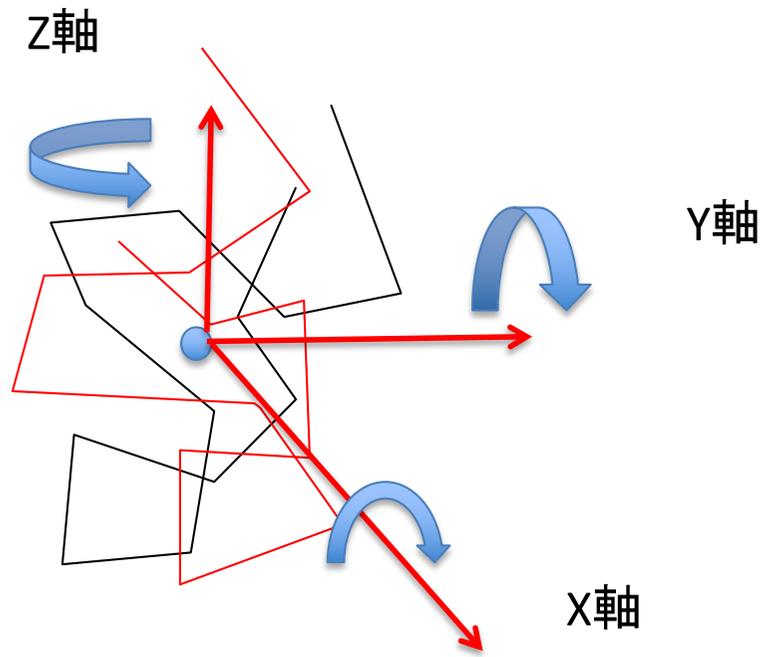
組換え体 k

x_i	y_i	z_j
-------	-------	-------





$$\text{Fitness} = \frac{1}{\text{各個体の回転操作による構造間のRMSD} + 0.1}$$



3次元の回転行列 [\[編集\]](#)

各軸周りの回転 [\[編集\]](#)

3次元空間でのx軸、y軸、z軸周りの回転を表す回転行列は、それぞれ次の通りである：

$$R_x(\theta) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \theta & -\sin \theta \\ 0 & \sin \theta & \cos \theta \end{bmatrix}$$

$$R_y(\theta) = \begin{bmatrix} \cos \theta & 0 & \sin \theta \\ 0 & 1 & 0 \\ -\sin \theta & 0 & \cos \theta \end{bmatrix}$$

$$R_z(\theta) = \begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

ここで回転の方向は、 R_x はy軸をz軸に向ける方向、 R_y はz軸をx軸に向ける方向、 R_z はx軸をy軸に向ける方向である。

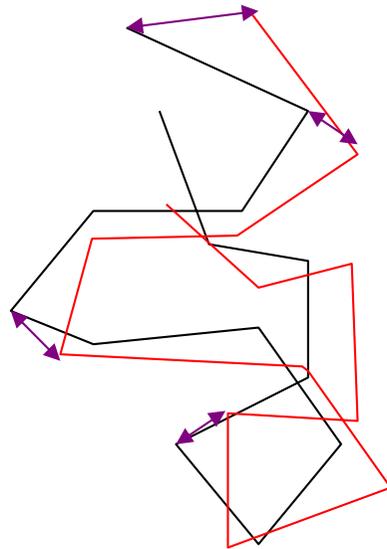
`rx`, `ry`, `rz` は、それぞれx軸、y軸、z軸の周りで、引数として与えられた角度だけ回転させる回転行列を作成し、それを返り値とする関数

`FunctionLibrary.R` 中に記述されており、`output`と`calcFitness`の二つの関数で使用されている。

RMSD

rmsd (root mean square distance) =

$$\sqrt{\frac{1}{n} \sum (dist(A(i), B(i)))^2}$$



残基間対応が最初に
与えられていると
計算は容易

注意事項

今回はGAを使用した発見的な方法で重ねあわせを行っているが、重ねあわせは厳密に求める方法があるので、実際の使用にあたっては、そちらを利用する方が良い

今回はGAを理解する手段として重ねあわせを用いてるだけである。

課題の説明

[1] 1j19.pdbと1j19m.pdbの重ね合わせ

[1-1] GAに及ぼす集団サイズの効果

[1-2] GAに及ぼす世代数の効果

[1-3] GAに及ぼす突然変異率の効果

[1-4] GAに及ぼす組換え率の効果

[2] 1j19.pdbと3c9a.pdbの重ね合わせ

[2-1] GAに及ぼす集団サイズの効果

[2-2] GAに及ぼす世代数の効果

[2-3] GAに及ぼす突然変異率の効果

[2-4] GAに及ぼす組換え率の効果

[3] プログラムの修正による1j19.pdbと3c9a.pdbの重ね合わせの改良

[3-1] 修正前に実行

[3-2] 修正後に実行

重ね合わせデータ

今回、

(1) 1j19.pdbと1j19m.pdbの重ね合わせ
と

(2) 1j19.pdbと3c9a.pdbの重ね合わせ
を行う。

1j19.pdb human epidermal growth factor (EGF)

42アミノ酸

3c9a.pdb *Drosophila melanogaster* Spitzタンパク質の

EGFドメイン

48アミノ酸

重ね合わせには、一方の構造のどのアミノ酸と他方の構造のどのアミノ酸が対応するかを決めておかないと、RMSDが計算できない。

対応関係はアラインメントファイルで与える。

アラインメント

(1) 1j19.pdbと1j19m.pdbの重ね合わせ用

1j19.aln.fasta

```
>1JL9A:EPIDERMAL GROWTH FACTOR
CPLSHDGYCLHDGVCMYIEALDKYACNCVVG YIGERCQYRDL
>1JL9B:EPIDERMAL GROWTH FACTOR
CPLSHDGYCLHDGVCMYIEALDKYACNCVVG YIGERCQYRDL
```

(2) 1j19.pdbと3c9a.pdbの重ね合わせ用

seq.aln.fasta

```
>3C9A:C|PDBID|CHAIN|SEQUENCE
PTYKCPETFD AWYCLNDAHCF AVKIADLPVYSCECAIGFMGQRCEYKE-
>1JL9B:EPIDERMAL GROWTH FACTOR
----CPLSHDG-YCLHDGVCMYIEALD--KYACNCVVG YIGERCQYRDL
```

進化の過程での挿入/欠失を考慮して、配列間の類似度が最大になるように並置

上記のファイル形式をFASTA形式と呼ぶ(“>”で始まる行は注釈行、それ以外は配列データ)

‘-’はギャップと呼ばれる空記号。挿入/欠失に対応させて、アミノ酸の位置をずらすために使われる。

赤:一致、青:置換されているが物理化学的に類似

使用するデータファイル

(1) 立体構造

1j19.pdb, 1j19m.pdb, 3c9a.pdb

の3つ

(2) アライメント

1j19.aln.fasta, seq.aln.fasta

の2つ

テキスト形式で、改行コードはWindowsにしてある

もし、MacやLinuxユーザでファイルの改行がおかしかったり、Rでうまく読み込めない場合は、藤 (tohhir@kwansei.ac.jp)

まで連絡すること

プログラム

ExecuteGA.R

FunctionLibraryX.R

FunctionLibraryY.R

ExecuteGA.R

第一パート

GAのパラメータや入力ファイルや出力ファイルを設定する。

第二パート

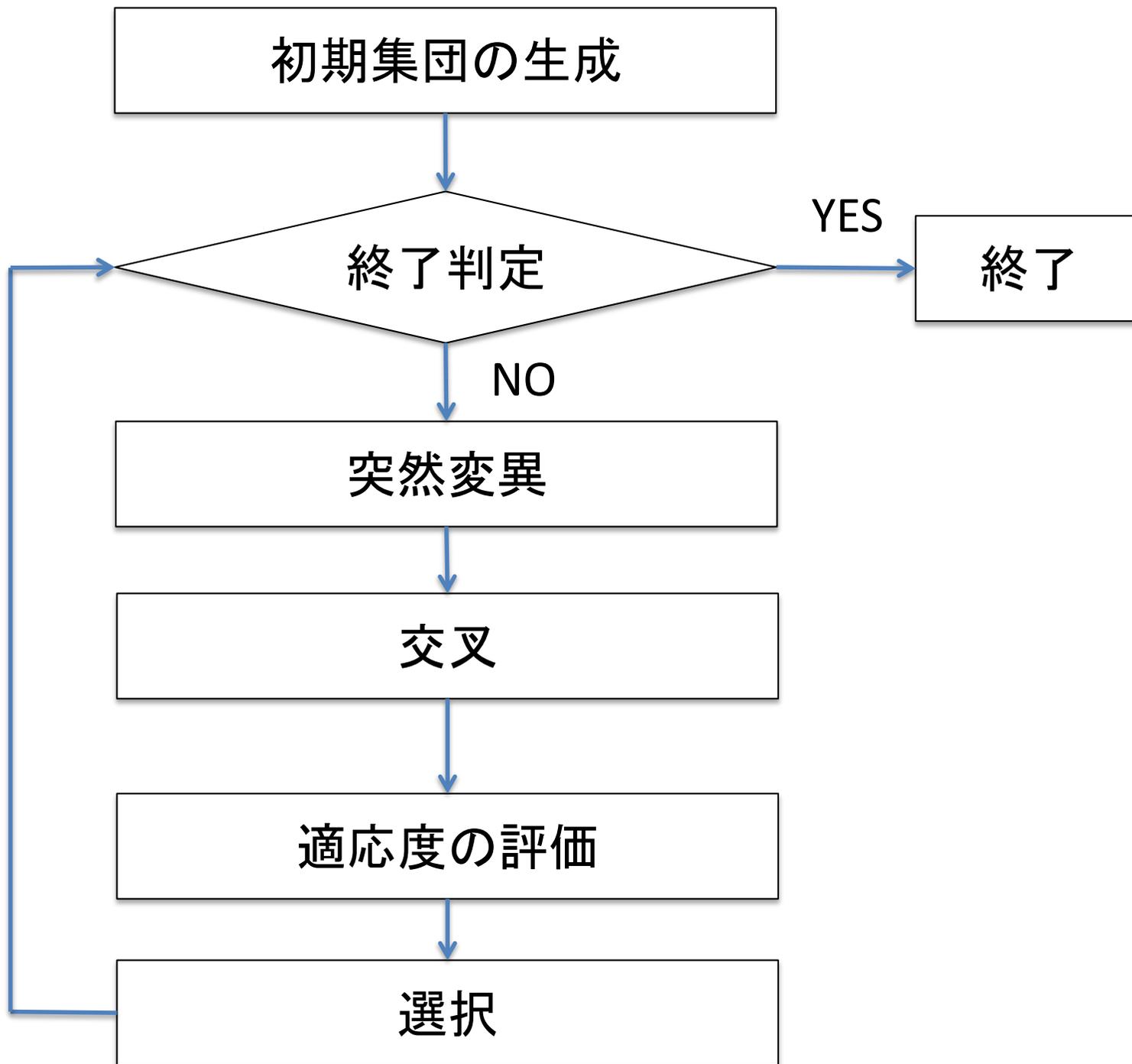
GAの実施部分

第一パートの構造データの指定の順番と
アラインメントデータの順番が一致するように気をつけること

```
seq.aln.fastaでは、アラインメントが
>3C9A:C|PDBID|CHAIN|SEQUENCE
PTYKCPETFDAWYCLNDAHCFVAVKIADLPVYSCECAIGFMGQRCEYKE-
>1JL9B:EPIDERMAL GROWTH FACTOR
----CPLSHDG-YCLHDGVCMYIEALD-KYACNCVVG YIGERCQYRDL
```

と与えられている。この時、ExecuteGA.Rでは、

```
#入力立体構造1
InputFileName1 <- '3c9a.pdb'
#入力立体構造2
InputFileName2 <- '1j19m.pdb'
#アラインメントファイル
AlignmentFile <- 'seq.aln.fasta'
#変異率
MutationRate <- 0.9
#組み換え率
RecombinationRate <- 0.9
#redex: 適応度の更新が見られなかった時に、どれくらい回転角度の範囲を狭めるのかを示すパラメータ
redx <- 1.0
#出力ファイル名1(回転後の立体構造1)
OutputFileName1 <- '3c9a_500_100_0.9_0.9_1.pdb'
#出力ファイル名2(回転後の立体構造2)
OutputFileName2 <- '1j19_500_100_0.9_0.9_1.pdb'
#出力ファイル名3(経過プロット)
OutputPlotName <- 'plot_3c9a_1j19_500_100_0.9_0.9_1.png'
```



初期集団の生成

終了判定

YES

終了

NO

突然変異

交叉

適応度の評価

選択

第一パート

#集団サイズ このパラメータを書き換えます

```
PopulationSize <- 10
```

#世代数 このパラメータを書き換えます

```
GenerationNumber <- 100
```

#入力立体構造1 この入力ファイル名を書き換えます

```
InputFileName1 <- '~/Desktop/先端実習2020/1j19.pdb'
```

#入力立体構造2 この入力ファイル名を書き換えます

```
InputFileName2 <- '~/Desktop/先端実習2020/1j19m.pdb'
```

#アラインメントファイル この入力ファイル名を書き換えます

```
AlignmentFile <- '~/Desktop/先端実習2020/1j19.aln.fasta'
```

#変異率 このパラメータを書き換えます

```
MutationRate <- 1.0
```

#組み換え率 このパラメータを書き換えます。

```
RecombinationRate <- 1.0
```

#redex: 適応度の更新が見られなかった時に、どれくらい回転角度の範囲を狭めるのかを示すパラメー

```
redx <- 1.0        # 固定にします、今回は使いません。
```

第一パート (続き)

```
#出力ファイル名1 (回転後の立体構造1) この出力ファイル名を書き換える(設定したパラメータを反映する名前にする  
OutputFileName1 <- '1j19_10_100_1.0_1.0.pdb'
```

```
#出力ファイル名2 (回転後の立体構造2) この出力ファイル名を書き換える(設定したパラメータを反映する名前にする  
OutputFileName2 <- '1j19m_10_100_1.0_1.0.pdb'
```

```
#出力ファイル名3 (経過プロット) この出力ファイル名を書き換える (設定したパラメータを反映する名前にする  
OutputPlotName <- 'plot_1j19_1j19m_10_100_1.0_1.0.png'
```

注意: 同じパラメータで実行する際には、出力ファイル名を変更すること。
変更しないと読み込み時に停止する (関数prepareGAの中の処理)

ExecuteGA.Rを実行する際の注意

- 今回の実習の後半で扱う遺伝的アルゴリズムのプログラムは、出力ファイルが上書きされないように、すでにフォルダにあるファイルと同じファイル名を出力しようとする、下のような**エラーが出る**ようにしています

コンソール画面:

```
> source('~/Desktop/2017 sentan/Rprograms/ExecuteGA.R', chdir = TRUE)
Error in prepareGA(InputFileName1, InputFileName2, AlignmentFile, PopulationSize) :
  ERROR::Same OutputFileName modified1j19.pdb was detected in this directory...
This Program have been stopped.  Rename Your OutputfileName1.
```

→ ExecuteGA.Rを実行する度に、ExecuteGA.Rのプログラムの出力する**ファイル名を変更**してください

第二パート

関数ライブラリを読み込みます。

```
source("FunctionLibraryY.R")
```

遺伝的アルゴリズム構築#####

以下の preparationGA, mutation, recombination, calculationFitness, selection,

output は**FunctionLibraryY.R**の中で定義されている関数

初期集団を作成します。

```
Population <- prepareGA(InputFileName1, InputFileName2, AlignmentFile,  
                        PopulationSize)
```

繰り返し操作を実行します。

```
for(i in 1:GenerationNumber){
```

突然変異を行います。

```
Mutated <- mutation(Population, MutationRate)
```

Mutatedには、親集団と突然変異で作成した子集団を一緒にした集団が格納

組み換え(交叉)を行います。

```
Recombined <- recombination(Mutated, RecombinationRate)
```

Recombinedには、MutatedとMutatedから組換えで作成した子集団を一枝にした集団が格納

適応度の評価を行います。

```
Fitness <- calculationFitness(Recombined)
```

集団recombinedの各個体に対する適応度(fitness)を計算

適応度にしたがって個体の選択を行います。

```
Population <- selection(Recombined, Fitness)
```

Fitnessに従いRecombinedから次世代のPopulationを作成する

```
}
```

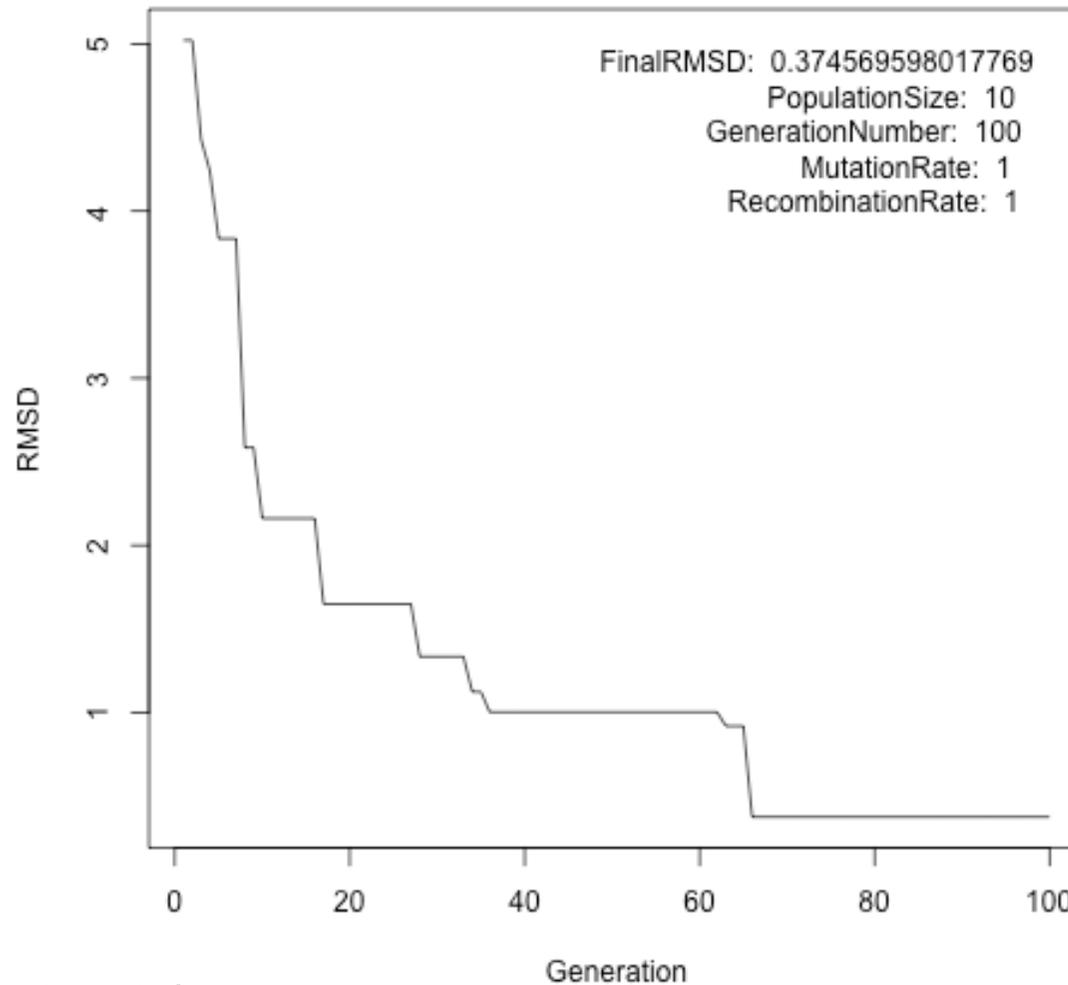
終了結果をファイルに出力します。

```
output(Population, Fitness, OutputFileName1, OutputFileName2, OutputPlotName)
```

出力例

出力ファイル名3で指定されたpngファイル

Filename: EXAMPLEPLOT



横軸:世代数

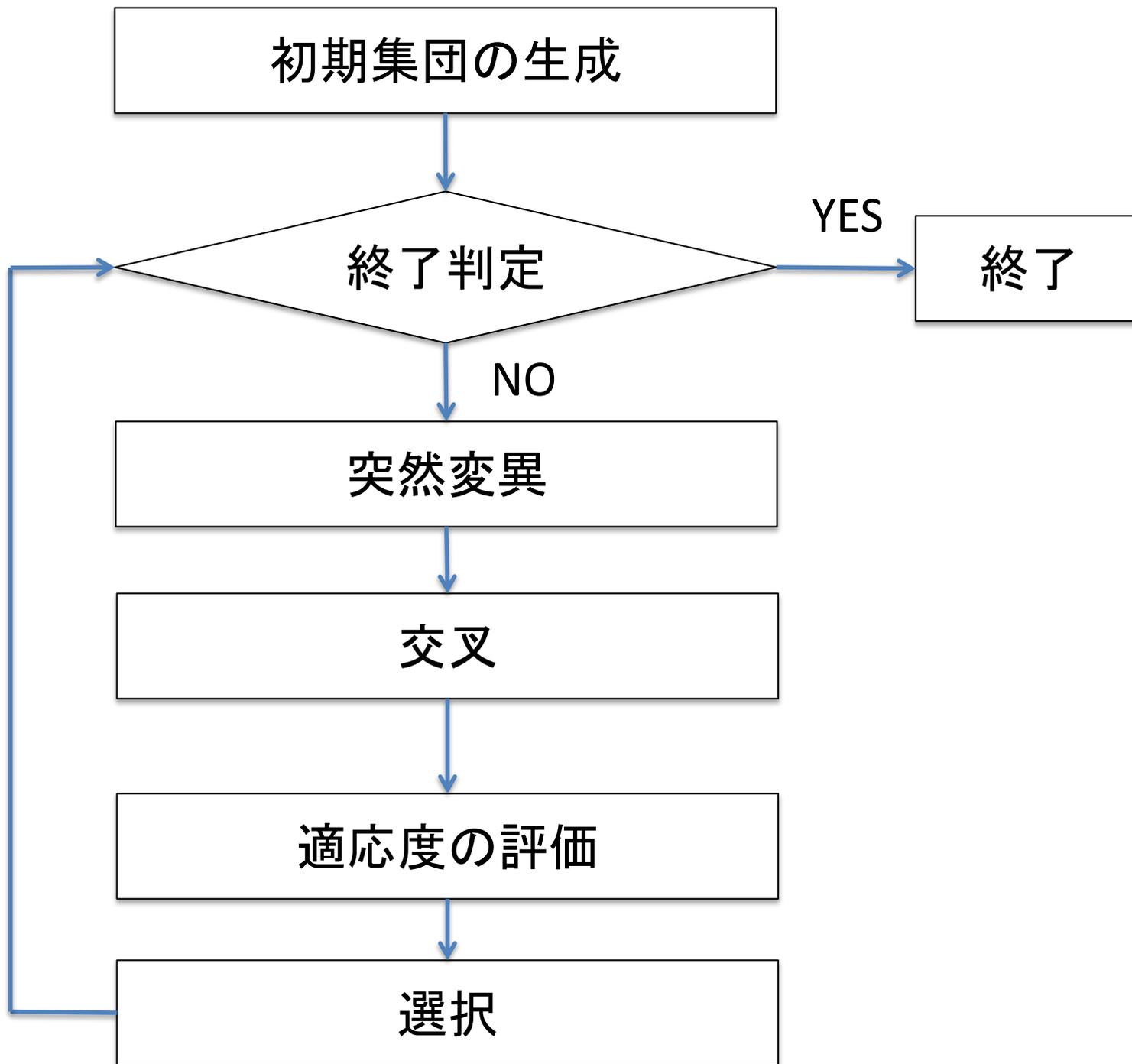
縦軸:各世代の重ね合わせのRMSD

この図で確認すべきポイント

世代数、集団サイズ、突然変異率、組換え率を変更した時、どれだけ早くRMSDが小さくなるか

RMSDが小さい = 構造の重ね合わせが良い

RMSDを減少させるのに、どのパラメータが大きく寄与するかを調べる。



初期集団の生成

終了判定

YES

終了

NO

突然変異

交叉

適応度の評価

選択

[1] 1j19.pdbと1j19m.pdbを重ね合わせる

[1-1] 集団サイズ3と集団サイズ100の比較

世代数は100、突然変異率0.7, 組換え率0.7に固定

executeGA.Rのパラメータを変更して実行

同じパラメータで5回実行

~~set.seedの括弧内を回数値に設定する~~

最終のステップのRMSDの平均

--- 重ね合わせがどれだけうまくいっているかの指標

RMSDがはじめて1.0以下になった時のステップ数の平均

--- どれだけ速く重ね合わせできているかの指標

RMSDの世代数に対するプロット(集団サイズ3と100)、それぞれ
5回実施したものの一つ

重ね合わせ前の1j19.pdbと1j19m.pdbの構造のmolmilによる表示

重ね合わせ後の構造のmolmilによる表示(集団サイズ3と100)

それぞれ5回実施した時のどれか一つについてスクリーンショットを撮る

集団サイズが重ねあわせに及ぼす影響を、上記のデータに基づいて議論

※ 1j19.pdbと1j19m.pdbは、同じ構造なので、重ね合わせのRMSDは
本来0になるはず

スクリーンショットの撮り方の調べ方

Googleなどで検索

キーワードは “スクリーンショット”と”Windows”

あるいは

“スクリーンショット”と”Mac”

自分が使っているPCによってWindowsかMacを変えて検索

[1] 1j19.pdbと1j19m.pdbを重ね合わせる

[1-2] 世代数5と世代数100の比較

集団サイズは100、突然変異率0.7, 組換え率0.7に固定
executeGA.Rのパラメータを変更して実行

同じパラメータで5回実行

~~set.seedの括弧内を回数値に設定する~~

最終のステップのRMSDの平均

--- 重ね合わせがどれだけうまくいっているかの指標

RMSDの世代数に対するプロット(世代数5と100)、それぞれ
5回実施したものの一つ

重ね合わせ前の1j19.pdbと1j19m.pdbの構造のmolmilによる表示
重ね合わせ後の構造のmolmilによる表示(世代数5と世代数100)
それぞれ5回実施した時のどれか一つについてスクリーンショットを撮る

世代数が重ねあわせに及ぼす影響を、上記のデータに基づいて議論

※ 1j19.pdbと1j19m.pdbは、同じ構造なので、重ね合わせのRMSDは
本来0になるはず

[1] 1j19.pdbと1j19m.pdbを重ね合わせる

[1-3] 突然変異率0と突然変異率0.7の比較

集団サイズは100、世代数100, 組換え率0.7に固定

executeGA.Rのパラメータを変更して実行

同じパラメータで5回実行

~~set.seedの括弧内を回数値に設定する~~

最終のステップのRMSDの平均

--- 重ね合わせがどれだけうまくいっているかの指標

RMSDがはじめて1.0以下になった時のステップ数の平均

--- どれだけ速く重ね合わせできているかの指標

RMSDの世代数に対するプロット(~~世代数5と100, 突然変異率0と0.7~~)、
それぞれ5回実施したものの一つ

重ね合わせ前の1j19.pdbと1j19m.pdbの構造のmolmilによる表示

重ね合わせ後の構造のmolmilによる表示(突然変異率0と0.7)

それぞれ5回実施した時のどれか一つについてスクリーンショットを撮る

突然変異率が重ねあわせに及ぼす影響を、上記のデータに基づいて議論

※ 1j19.pdbと1j19m.pdbは、同じ構造なので、重ね合わせのRMSDは本来0になるはず

[1] 1j19.pdbと1j19m.pdbを重ね合わせる

[1-4] 組換え率0と組換え率0.7の比較

集団サイズは100、世代数100、突然変異率0.7に固定
executeGA.Rのパラメータを変更して実行
同じパラメータで5回実行

~~set.seedの括弧内を回数値に設定する~~

最終のステップのRMSDの平均

- 重ね合わせがどれだけうまくいっているかの指標
RMSDがはじめて1.0以下になった時のステップ数の平均
- どれだけ速く重ね合わせできているかの指標
1.0に到達しなければステップ数は100とする

RMSDの世代数に対するプロット(~~世代数5と100~~組換え率0と0.7)、
それぞれ5回実施したものの一つ

重ね合わせ前の1j19.pdbと1j19m.pdbの構造のmolmilによる表示
重ね合わせ後の構造のmolmilによる表示(組換え率0と組換え率0.7)
それぞれ5回実施した時のどれか一つについてスクリーンショットを撮る

組換え率が重ねあわせに及ぼす影響を、上記のデータに基づいて議論

※ 1j19.pdbと1j19m.pdbは、同じ構造なので、重ね合わせのRMSDは
本来0になるはず

[2] 3c9a.pdbと1j19.pdbを重ね合わせる

[2-1] 集団サイズ3と集団サイズ100の比較

世代数は100、突然変異率0.7, 組換え率0.7に固定

executeGA.Rのパラメータを変更して実行

同じパラメータで5回実行

~~set.seedの括弧内を回数値に設定する~~

最終のステップのRMSDの平均

--- 重ね合わせがどれだけうまくいっているかの指標

RMSDがはじめて3.0以下になった時のステップ数の平均

--- どれだけ速く重ね合わせできているかの指標

RMSDの世代数に対するプロット(集団サイズ3と100)、それぞれ
5回実施したものの一つ

重ね合わせ前の1j19.pdbと3c9a.pdbの構造のmolmilによる表示

重ね合わせ後の構造のmolmilによる表示(集団サイズ3と100)

それぞれ5回実施した時のどれか一つについてスクリーンショットを撮る

集団サイズが重ねあわせに及ぼす影響を、上記のデータに基づいて議論

[2] 1j19.pdbと3c9a.pdbを重ね合わせる

[2-2] 世代数5と世代数100の比較

集団サイズは100、突然変異率0.7, 組換え率0.7に固定
executeGA.Rのパラメータを変更して実行

同じパラメータで5回実行

~~set.seed~~の括弧内を回数値に設定する

最終のステップのRMSDの平均

--- 重ね合わせがどれだけうまくいっているかの指標

RMSDの世代数に対するプロット(世代数5と100)、それぞれ
5回実施したものの一つ

重ね合わせ前の1j19.pdbと3c9a.pdbの構造のmolmilによる表示
重ね合わせ後の構造のmolmilによる表示(世代数5と世代数100)
それぞれ5回実施した時のどれか一つについてスクリーンショットを撮る

世代数が重ねあわせに及ぼす影響を、上記のデータに基づいて議論

[2] 1j19.pdbと3c9a.pdbを重ね合わせる

[2-3] 突然変異率0と突然変異率0.7の比較

集団サイズは100、世代数100, 組換え率0.7に固定

executeGA.Rのパラメータを変更して実行

同じパラメータで5回実行

~~set.seedの括弧内を回数値に設定する~~

最終のステップのRMSDの平均

--- 重ね合わせがどれだけうまくいっているかの指標

RMSDがはじめて3.0以下になった時のステップ数の平均

--- どれだけ速く重ね合わせできているかの指標

RMSDの世代数に対するプロット(~~世代数5と100~~ 突然変異率0と0.7)、それぞれ5回実施したものの一つ

重ね合わせ前の1j19.pdbと3c9a.pdbの構造のmolmilによる表示

重ね合わせ後の構造のmolmilによる表示(突然変異率0と0.7)

それぞれ5回実施した時のどれか一つについてスクリーンショットを撮る

突然変異率が重ねあわせに及ぼす影響を、上記のデータに基づいて議論

[2] 1j19.pdbと3c9a.pdbを重ね合わせる

[2-4] 組換え率0と組換え率0.7の比較

集団サイズは100、世代数100, 突然変異率0.7に固定
executeGA.Rのパラメータを変更して実行
同じパラメータで5回実行

~~set.seed~~の括弧内を回数~~の値~~に設定する

最終のステップのRMSDの平均

--- 重ね合わせがどれだけうまくいっているかの指標

RMSDがはじめて3.0以下になった時のステップ数の平均

--- どれだけ速く重ね合わせできているかの指標

3.0に到達しなければステップ数は100とする

RMSDの世代数に対するプロット(~~世代数5と100~~ 組換え率0と0.7)、それぞれ
5回実施したものの一つ

重ね合わせ前の1j19.pdbと3c9a.pdbの構造のmolmilによる表示

重ね合わせ後の構造のmolmilによる表示(組換え率0と組換え率0.7)

それぞれ5回実施した時のどれか一つについてスクリーンショットを撮る

組換え率が重ねあわせに及ぼす影響を、上記のデータに基づいて議論

[3] 3c9a.pdbと1j19.pdbを重ね合わせる

[3-1] 集団サイズは500、世代数100, 突然変異率0.7, 組換え率0.7
で重ね合わせを実施

5回set.seedの値を変えながら実行

同じパラメータで5回実行

最終のステップのRMSDの平均

---- 重ね合わせがどれだけうまくいっているかの指標

RMSDがはじめて3.0以下になった時のステップ数の平均

---- どれだけ速く重ね合わせできているかの指標

3.0に到達しなければステップ数は100とする

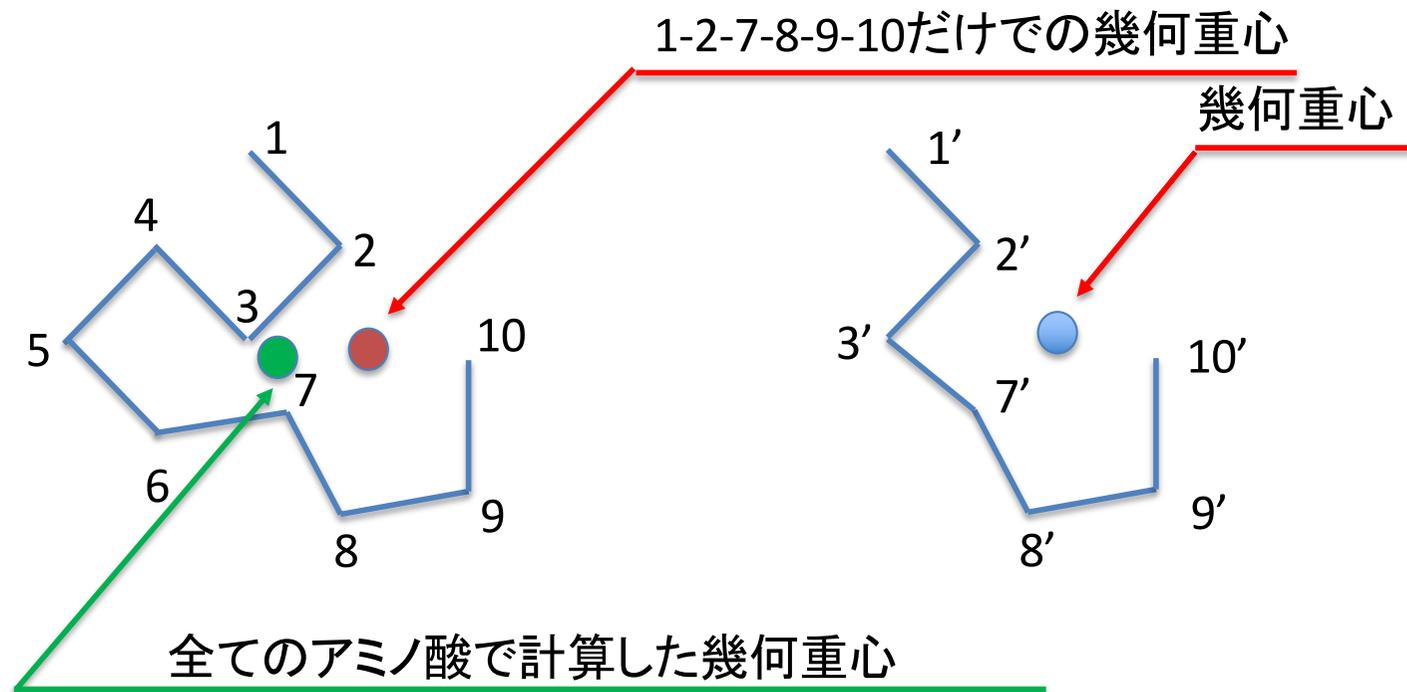
重ね合わせ前の3c9a.pdbと1j9.pdbの構造のmolmilによる表示

重ね合わせ後の構造のmolmilによる表示

それぞれ5回実施した時のどれか一つについてスクリーンショットを撮る

初期配置に比べると重なっているが、ズレが大きいことを確認

----→ プログラムを書き換えて重ね合わせを改善する。



アラインメント	1	2	3	4	5	6	7	8	9	10
	1'	2'	3'	-	-	-	7'	8'	9'	10'

重心あわせの後、アラインメント上対応する残基で回転させてRMSDを最小化
 現在のプログラムでは全てのアミノ酸のCaで幾何重心を計算
 アミノ酸4-5-6 (挿入)のせいで対応するアミノ酸の重心の位置がずれる

重心計算の際、アラインメントで対応するアミノ酸のCaだけを使えば
 より重ね合わせが良くなる？

[3] 3c9a.pdbと1j19.pdbを重ね合わせる

[3-2] FunctionLibraryY.Rの該当箇所を変更して、[3-1]と同じパラメータ値でexecuteGA.Rを実行

同じパラメータで5回実行

最終のステップのRMSDの平均

--- 重ね合わせがどれだけうまくいっているかの指標

RMSDがはじめて**3.0**以下になった時のステップ数の平均

--- どれだけ速く重ね合わせできているかの指標

3.0に到達しなければステップ数は100とする

重ね合わせ前の3c9a9.pdbと1j9.pdbの構造のmolmilによる表示

重ね合わせ後の構造のmolmilによる表示

それぞれ5回実施した時のどれか一つについてスクリーンショットを撮る

[3-1]の結果に比べて重ね合わせが改善されていることを確認

FunctionLibraryY.Rの

#1つ目のPDBのファイルの幾何重心を計算 (77行 - 79行) をコメントアウト

```
ca1gx <- mean(ca1$x)
```

```
ca1gy <- mean(ca1$y)
```

```
ca1gz <- mean(ca1$z)
```

#2つ目のPDBファイルの幾何重心を計算 (120行 - 122行) をコメントアウト

```
ca2gx <- mean(ca2$x)
```

```
ca2gy <- mean(ca2$y)
```

```
ca2gz <- mean(ca2$z)
```

ca1gx, ca1gy, ca1gz, ca2gx, ca2gy, ca2gzの計算において、変数alignから対応する残基のみを取り出し、その残基だけで重心を計算する

※ コメントアウトとは、コマンド行をコメントとして実行から外すこと上の例では、行頭に#を入れると良い。例えば77行の場合

```
#ca1gx <- mean(ca1$x)
```

ヒント: 36行でアラインメントのFASTAファイルを読み込んでいる変数alignを使う

align\$ali

```
          [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10] [,11] [,12] [,13] [,14] [,15] [,16] [,17] [,18]
3C9A:C|PDBID|CHAIN|SEQUENCE "P" "T" "Y" "K" "C" "P" "E" "T" "F" "D" "A" "W" "Y" "C" "L" "N" "D" "A"
1JL9B:EPIDERMAL            "-" "-" "-" "-" "C" "P" "L" "S" "H" "D" "G" "-" "Y" "C" "L" "H" "D" "G"
          [,19] [,20] [,21] [,22] [,23] [,24] [,25] [,26] [,27] [,28] [,29] [,30] [,31] [,32] [,33] [,34]
3C9A:C|PDBID|CHAIN|SEQUENCE "H" "C" "F" "A" "V" "K" "I" "A" "D" "L" "P" "V" "Y" "S" "C" "E"
1JL9B:EPIDERMAL            "V" "C" "M" "Y" "I" "E" "A" "L" "D" "-" "-" "K" "Y" "A" "C" "N"
          [,35] [,36] [,37] [,38] [,39] [,40] [,41] [,42] [,43] [,44] [,45] [,46] [,47] [,48] [,49]
3C9A:C|PDBID|CHAIN|SEQUENCE "C" "A" "I" "G" "F" "M" "G" "Q" "R" "C" "E" "Y" "K" "E" "-"
1JL9B:EPIDERMAL            "C" "V" "V" "G" "Y" "I" "G" "E" "R" "C" "Q" "Y" "R" "D" "L"
```

align\$ali[1,]に3c9aのアラインメントが得られる

align\$ali[2,]に1j19のアラインメントが得られる

align\$ali[1,4]はK、align\$ali[2,4]は'-'を

3c9aの重心を求める際、赤字の残基は、1j19に相手がいないので (gapなので)、重心の計算に使わない

1j19の重心を求める際、青字の残基は、3c9aに相手がいないので (gapなので)、重心の計算に使わない

アラインメントの長さは length(align\$ali[1,]) あるいは length(align\$ali[2,]) で得られる

length(align\$ali[1,]) も length(align\$ali[2,]) も、同じ値

ヒントの続き

立体構造のCaの座標データを含む行は56行、57行で取り出されている

```
ca1 <- p1$atom[p1$atom[,3]=='CA',]
```

```
ca2 <- p2$atom[p2$atom[,3]=='CA',]
```

立体構造の1番目のアミノ酸のx、y、z座標は `ca1[1,]$x`, `ca1[1,]$y`, `ca1[1,]$z` で取り出せる。」

例えば、3c9aの重心を求める時 (81行から97行の先頭の#を削除してコメント行から復帰させる)

アラインメントのサイト番号は、立体構造の残基番号と対応していないことに注意

変数 `siteNo` を用意しておく。 `siteNo <- 0`

空ベクトル `ax1`, `ay1`, `az1`を用意しておく。例: `ax1 <- c()`

for 文で アラインメントの一番目のサイトから最後のサイトまで順番に見ていく

(1) if文を使う。 3c9aが'- 'でない時、`siteNo <- siteNo + 1`

(2) 相手側(1j19)のアラインメントサイトが'- 'でない時だけ、x, y, zの座標値を重心の計算に利用
← if文を使う

```
siteNo <- 0
ax1 <- c()
ay1 <- c()
az1 <- c()
for (i in ????) {
  if (3c9aのi番目のアラインメントサイトが'- 'でない場合) {
    siteNo <- siteNo + 1
    if (1j19のi番目のアラインメントサイトが'- 'でない場合) {
      ax1 <- c(ax1, ca1のsiteNo番目のx座標)
      ay1 <- c(ay1, ca1のsiteNo番目のy座標)
      az1 <- c(az1, ca1のsiteNo番目のz座標)
    }
  }
}
```

ヒントの続き

得られたax1, ay1, az1を使って重心を求める

```
ca1gx <- mean(ax1)
```

```
ca1gy <- mean(ay1)
```

```
ca1gz <- mean(az1)
```

1j19についても同様の操作を行う。

コメントになっている106行から122行を復帰させて実施

```
siteNo <- 0
ax2 <- c()
ay2 <- c()
az2 <- c()
for (i in ?????) {
  if (1j19のi番目のアラインメントサイトが'-'でない場合) {
    siteNo <- siteNo + 1
    if (3c9aのi番目のアラインメントサイトが'-'でない場合) {
      ax2 <- c(ax2, ca2のsiteNo番目のx座標)
      ay2 <- c(ay2, ca2のsiteNo番目のy座標)
      az2 <- c(az2, ca2のsiteNo番目のz座標)
    }
  }
}
ca2gx <- mean(ax2)
ca2gy <- mean(ay2)
ca2gz <- mean(az2)
```

6. レポートの構成

レポートの構成

1. 背景
2. 方法
3. 結果
4. 考察

レポートの作成

1. Introduction (背景)

(1) 立体構造の重ね合わせの説明

(2) 遺伝的アルゴリズムの一般論の説明

(3) Rの説明

(4) 目的

(4-1) 立体構造の重ね合わせを題材として遺伝的アルゴリズムの各種パラメータの及ぼす影響を調べる

(4-2) 重心計算の部分を改良して、重ね合わせを改善する

2. 方法

(1) 立体構造の重ねあわせのためのGAの設定の説明

(1-1) 染色体や遺伝子の構成

(1-2) 初期集団はどのように生成したか

(1-3) 突然変異や組換えはどのように行われるか

(1-4) 適応度はどのように定義されているか

(1-5) 選択はどのように行われているか

(1-6) 終了条件

(2) 準備されたExecuteGA.Rと~~FunctionLibraryY.RとFunctionLibraryX.R~~について説明

FunctionLibraryY.RとFunctionLibraryX.Rの中の各関数と(1)の処理との対応

(3) 各種パラメータを変更してGAを実施し、パラメータのGAに及ぼす影響の説明

(4) 現在のスクリプトの重心計算の問題点と改善の方法の説明

3. 結果

(1) 1j19.pdbと1j19m.pdbの重ね合わせ

[1-1] ~ [1-4] (スライドのp197 - 201)の結果を
図を含めて説明

(2) 3c9a.pdbと1j19.pdb の重ね合わせ

[2-1] ~ [2-4] (スライドのp.202 ~ p.205)の結果を図を含
めて説明

(3) プログラムを修正して3c9a.pdbと1j19.pdb
を重ね合わせ

[3-1] ~ [3-2] (スライドのp.215 ~ p.221)の
結果を図を含めて説明

4. 考察

(1) 3の結果の(1)(2)に基づき、各種パラメータのGAのパフォーマンスに及ぼす影響を議論
何故そうなるのか？

(2) 3の結果の(3)に基づき、重心計算の変更が重ね合わせに及ぼす影響や理由を議論