

距離行列 類似度行列 固有値分解 グラフ

法数学勉強会

2019/01/26

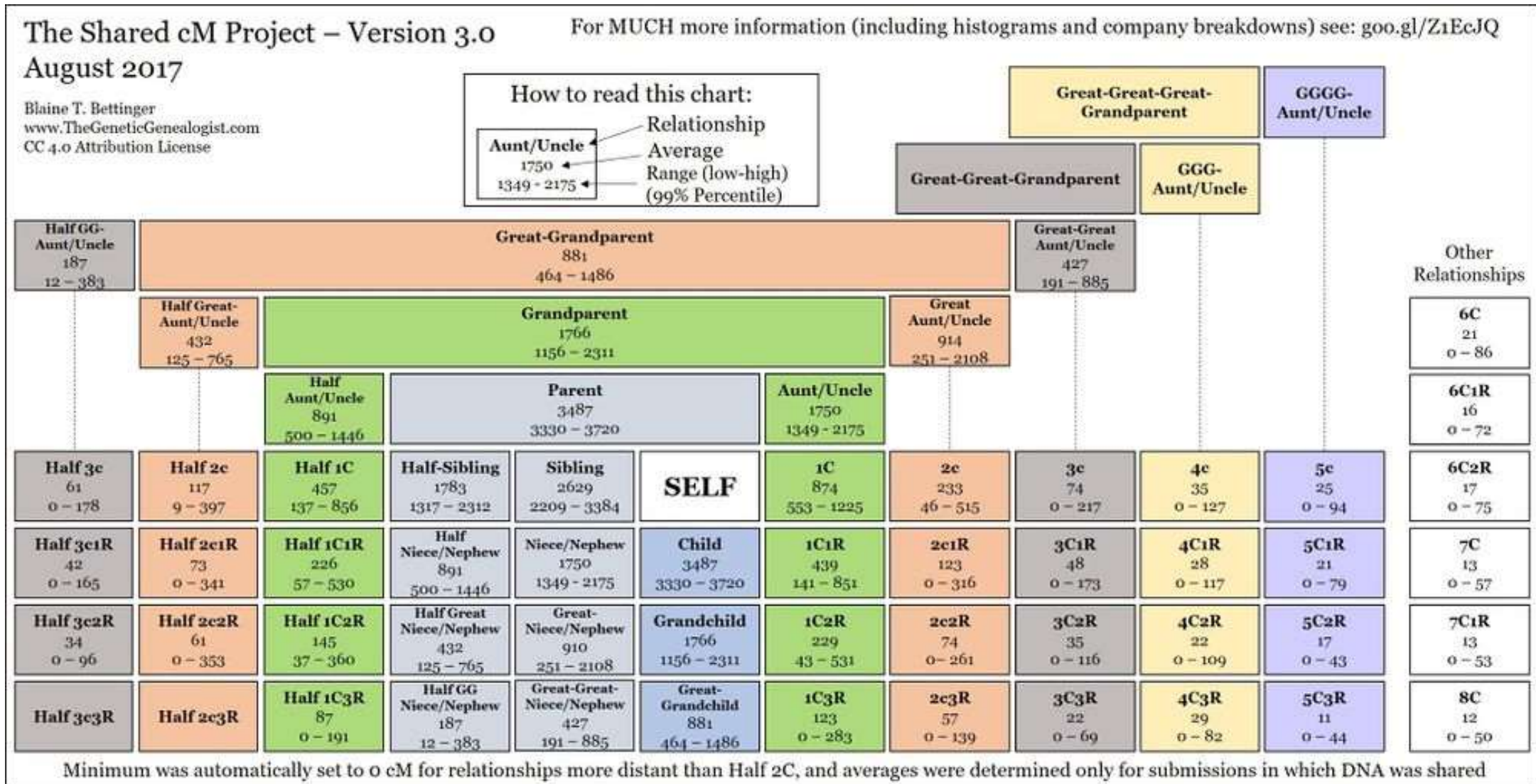
京都大学(医)統計遺伝学分野

山田 亮

前回の話

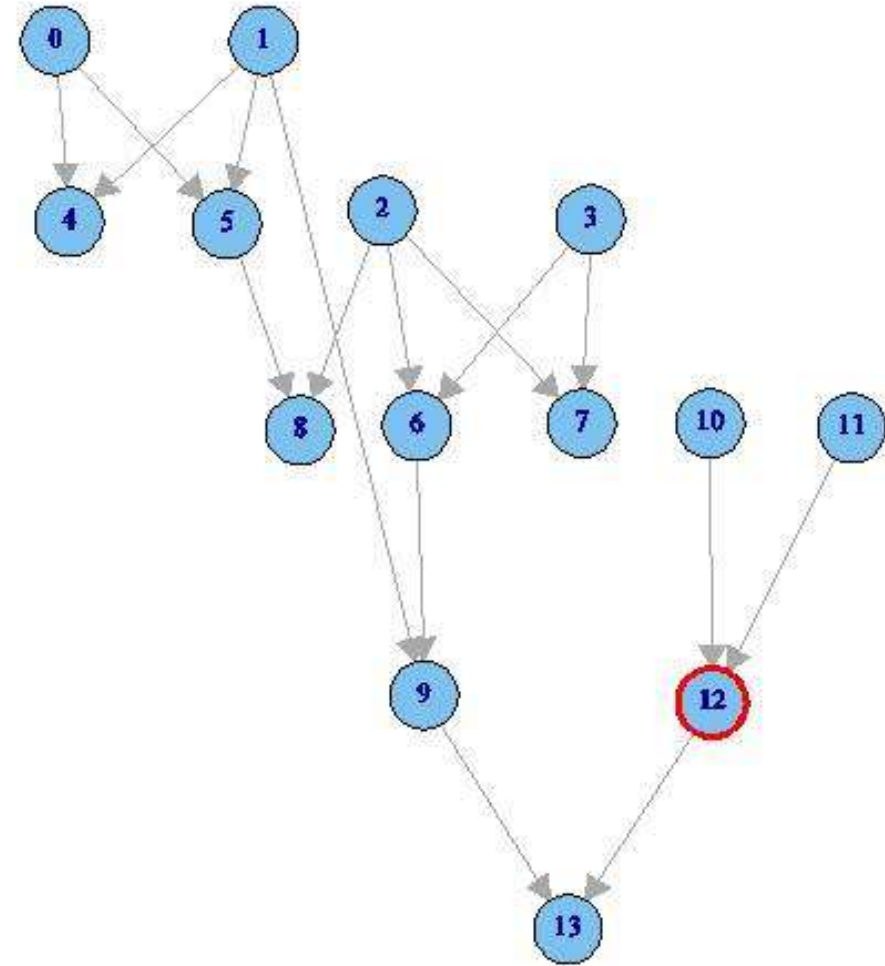
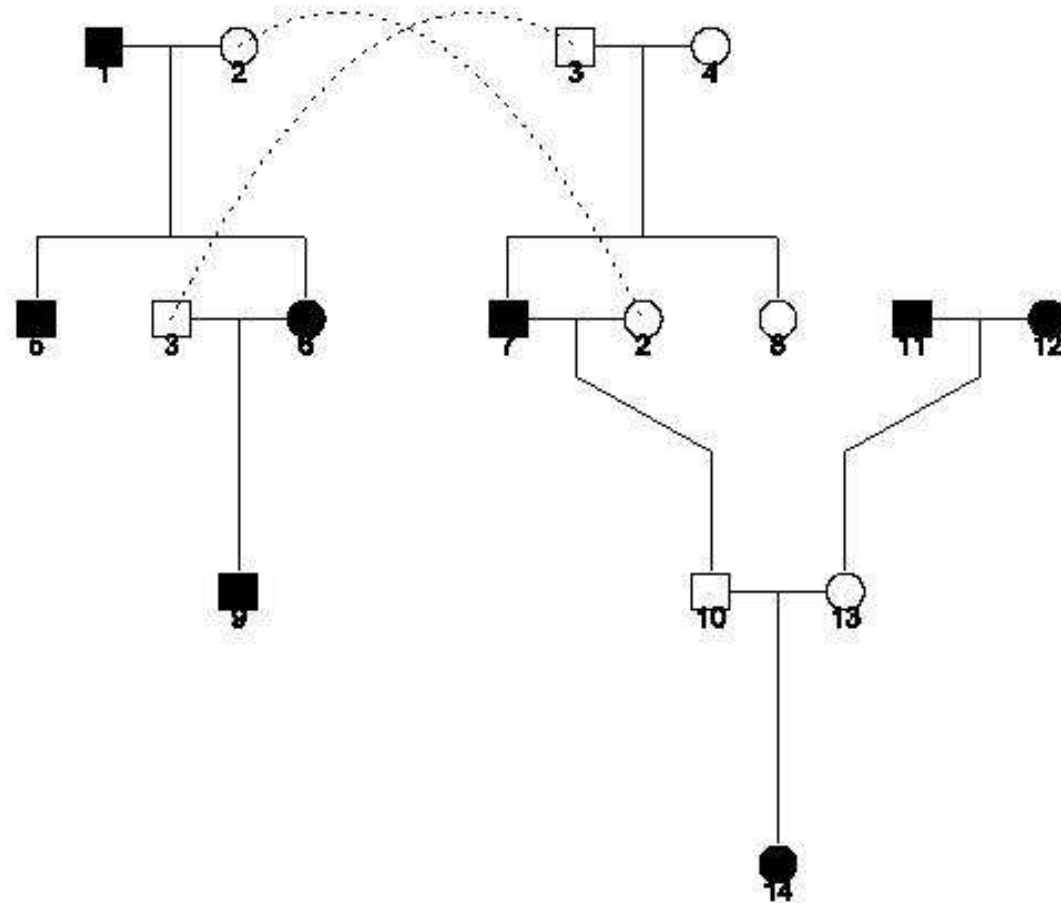
- ゴールデン・ステート・キラー
- DNA多型情報を用いて
- 2人間の血縁関係を推定し
- (犯人を特定した)
- 「家系図を復元したか?」、「単に登録して2名の関係の推定の列挙をしただけか?」
 - → 多分、後者…

Autosomal DNA match



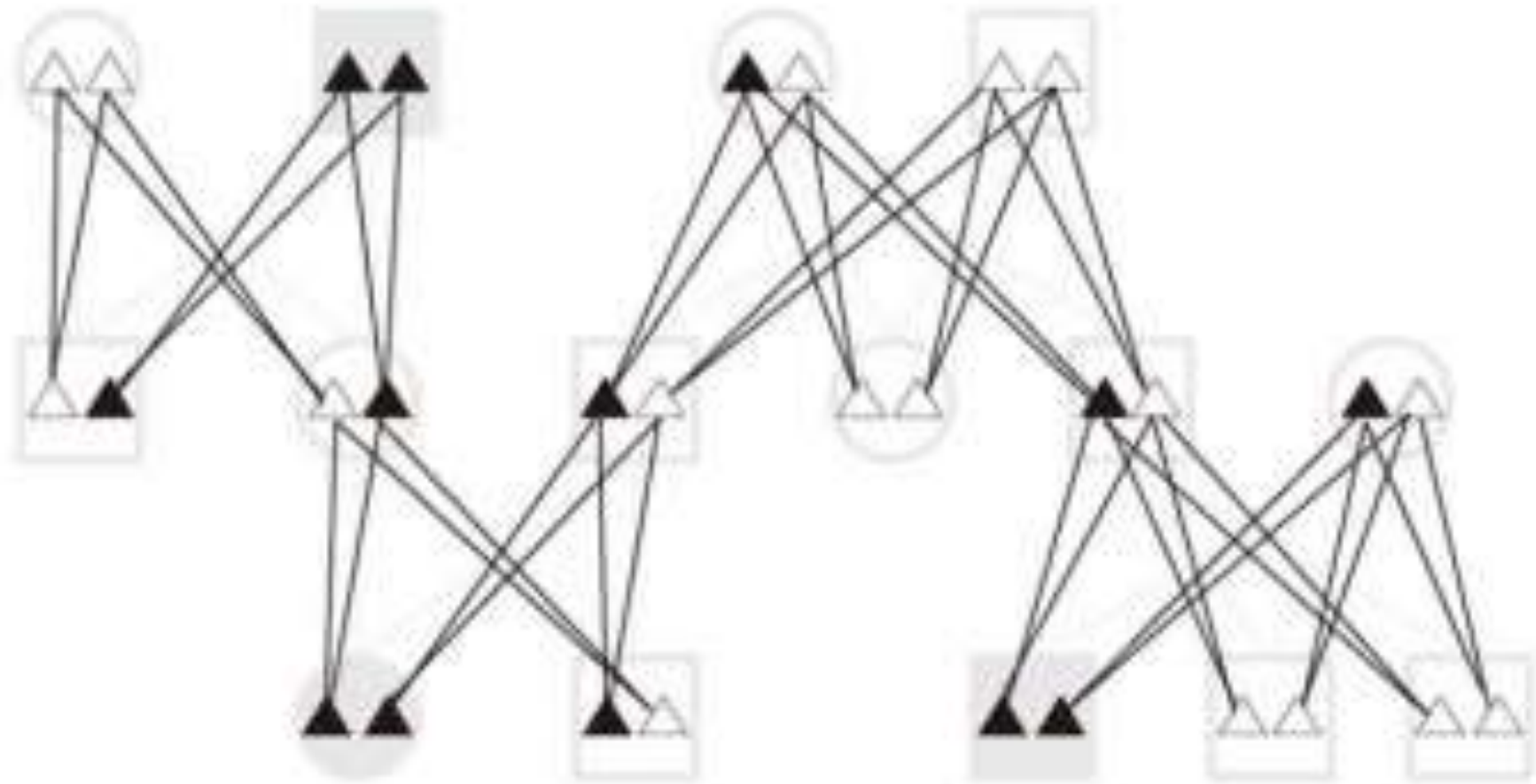
N人の関係を、2人ペアの関係の集まりと
して見る話し

家系図



親子は1/2の関係

子によって「かすがい」された「夫婦」は、1/4の関係なのか？



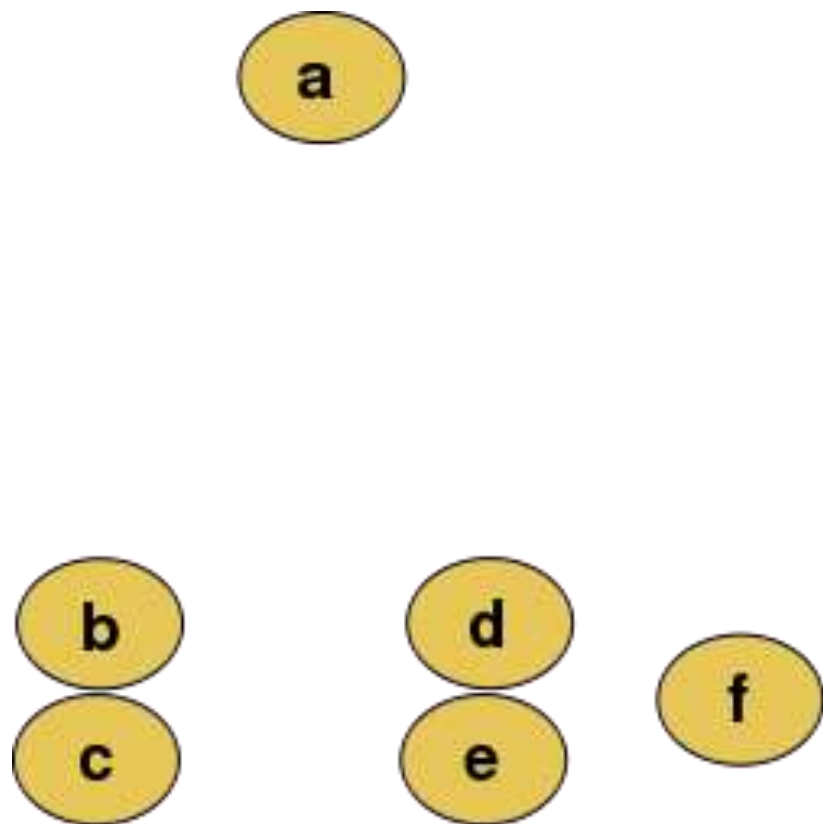
つながっているのは「染色体」
夫婦(の染色体)は子を介してつながっていない
孫(の染色体)を介してつながっている

たくさんの対象が作り上げる世界

- 二項定理
- 対象がN個あるとき
 - N個、それぞれについての情報： $N = {}_N C_1$
 - ペアに関する情報： ${}_N C_2$
 - トリオに関する情報： ${}_N C_3$
 - 4つ組に関する情報： ${}_N C_4$
 - ...
 - N個すべての組に特有の情報： ${}_N C_N$
 - ${}_N C_1 + {}_N C_2 + \dots + {}_N C_N = 2^N - 1$
 - $(1+x)^N = {}_N C_0 x^0 + {}_N C_1 x^1 + \dots + {}_N C_N x^N$
 - $x = 1$ のとき...

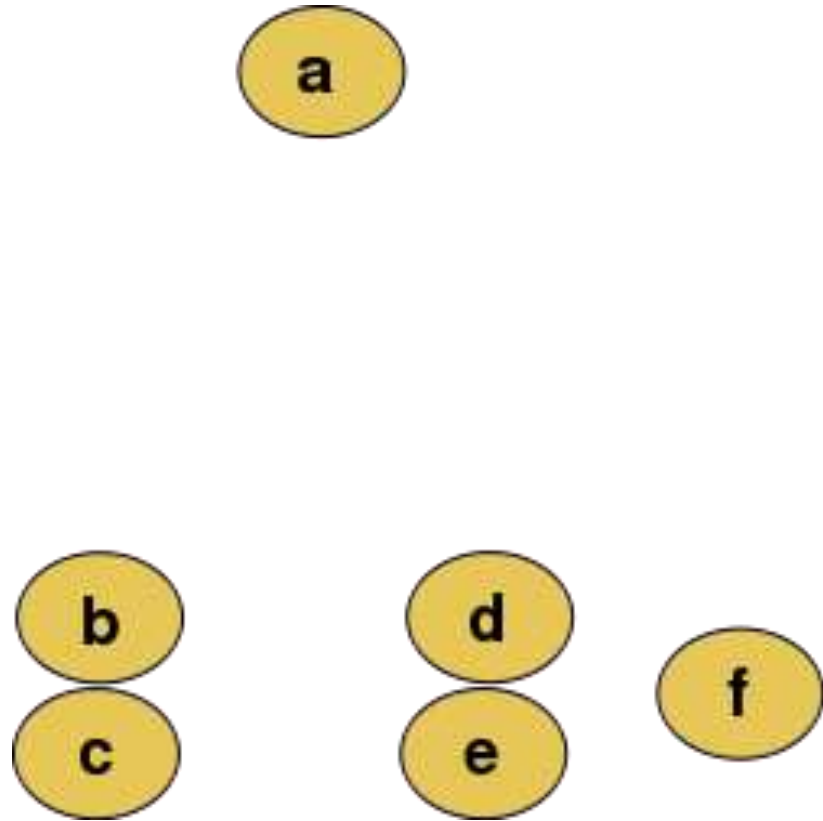
2^N は大変、 ${}_N C_2 = N(N-1)/2$ で済ませる

距離行列



	a	b	c	d	e	f
a	0	184	222	177	216	231
b	184	0	45	123	128	200
c	222	45	0	129	121	203
d	177	123	129	0	46	83
e	216	128	121	46	0	83
f	231	200	203	83	83	0

対称行列、ペア数は三角部分



	a	b	c	d	e	f
a	0	184	222	177	216	231
b	184	0	45	123	128	200
c	222	45	0	129	121	203
d	177	123	129	0	46	83
e	216	128	121	46	0	83
f	231	200	203	83	83	0

距離行列から、「全体像」を復元する

- 多次元尺度構成法 Multi-Dimensional Scaling (MDS)
 - 多数の対象が空間に置かれている
 - 個々の対象の座標は調べられないが
 - ペアワイズの距離は調べ上げられる
 - 距離行列から、オリジナルの座標を復元する

MDS

- 第1点を原点に取る
- 第2点を、第1点との距離に応じて、 X_1 軸上にとる
- 第3点を、第1点と第2点からの距離を満足するように、 (X_1, X_2) 平面上にとる
- ...
- 第 k 点を、第 $1, 2, \dots, k-1$ 点からの距離を満足するように (X_1, \dots, X_{k-1}) 次元空間にとる
- $N-1$ 次元空間にきちんと配置できる

	a	b	c	d	e	f
a	0	184	222	177	216	231
b	184	0	45	123	128	200
c	222	45	0	129	121	203
d	177	123	129	0	46	83
e	216	128	121	46	0	83
f	231	200	203	83	83	0

MDS

- N個の対象はN-1次元空間座標に、うまく納められる
- 場合によってはもっと低い次元に納められる
- 低次元に納められれば
 - 次元縮約に成功した、ことになる

	a	b	c	d	e	f
a	0	184	222	177	216	231
b	184	0	45	123	128	200
c	222	45	0	129	121	203
d	177	123	129	0	46	83
e	216	128	121	46	0	83
f	231	200	203	83	83	0

MDS

- N個の対象はN-1次元空間座標に、うまく納められる
- 場合によってはもっと低い次元に納められる
- 低次元に納められれば
 - 次元縮約に成功した、ことになる
- 次元縮約を目的にした場合には、「固有値分解」して、できるだけ説明力の強い軸を取り出す工夫をしながら、座標配置を実現する

	a	b	c	d	e	f
a	0	184	222	177	216	231
b	184	0	45	123	128	200
c	222	45	0	129	121	203
d	177	123	129	0	46	83
e	216	128	121	46	0	83
f	231	200	203	83	83	0

MDS

- N個の対象はN-1次元空間座標に、うまく納められる
- 場合によってはもっと低い次元に納められる
- 低次元に納められれば
 - 次元縮約に成功した、ことになる
- 次元縮約を目的にした場合には、「固有値分解」して、できるだけ説明力の強い軸を取り出す工夫をしながら、座標配置を実現する

	a	b	c	d	e	f
a	0	184	222	177	216	231
b	184	0	45	123	128	200
c	222	45	0	129	121	203
d	177	123	129	0	46	83
e	216	128	121	46	0	83
f	231	200	203	83	83	0

『もしも、元の空間がユークリッド空間だったなら』

MDS

- オリジナルの対象は、ユークリッド空間に置かれていたかどうかの保証はない
- ペアワイズの距離も、「ユークリッド距離」で測ったものとは限らない
- MDSをすると、固有値がすべて非負とは限らない
- 非負の固有値が出るということは、「空間をユークリッド的な平らな空間」に納めきれないということ
- 距離の測り具合が曲がっている（非ユークリッド的になっている）空間にあるとみなすこともできる

例

- 染色体の配列の違いを距離にする
 - 0010100
 - 0101100
 - 1101110 ...
 - 多次元(0,1)立方体の頂点に配置されているとみなせる
 - わざわざユークリッド空間に配置しなおしたい？
 - →線形代数的な解析の対象に乗せることができる
- 配列の違いの距離の測り方に、「マンハッタン距離」と言われるものがある
- 「ユークリッド距離」ではなく「マンハッタン距離」なので、ユークリッド空間に納めきれない
- MDSの固有値に負のものが現れる

非ユークリッド空間(球面)に納める

- 染色体の配列の違いを『内積』的にとらえる
 - 似ているほど内積は大きい
- 類似度行列を作る
 - 以下の配列をベクトルとみて、対応する値が同じならカウントし異なればカウントしないことにする
 - まったく同一の配列の内積が1となるように補正する
 - 0010100
 - 0101100
 - 1101110 ...
- 自配列との内積はすべて1 →すべての配列は球面上の点となる
- 内積は $\cos(\theta)$ になるので、角度が配列の違いに対応する
 - 0000000 と 1111111とは、直角をなす

内積

$$V = \begin{pmatrix} v1 \\ v2 \\ \dots \\ vn \end{pmatrix}$$

$$V^T = (v1, v2, \dots, vn)$$

- $V^T V = (v1, v2, \dots, vn) \begin{pmatrix} v1 \\ v2 \\ \dots \\ vn \end{pmatrix} = v1^2 + v2^2 + \dots + vn^2$

內積行列

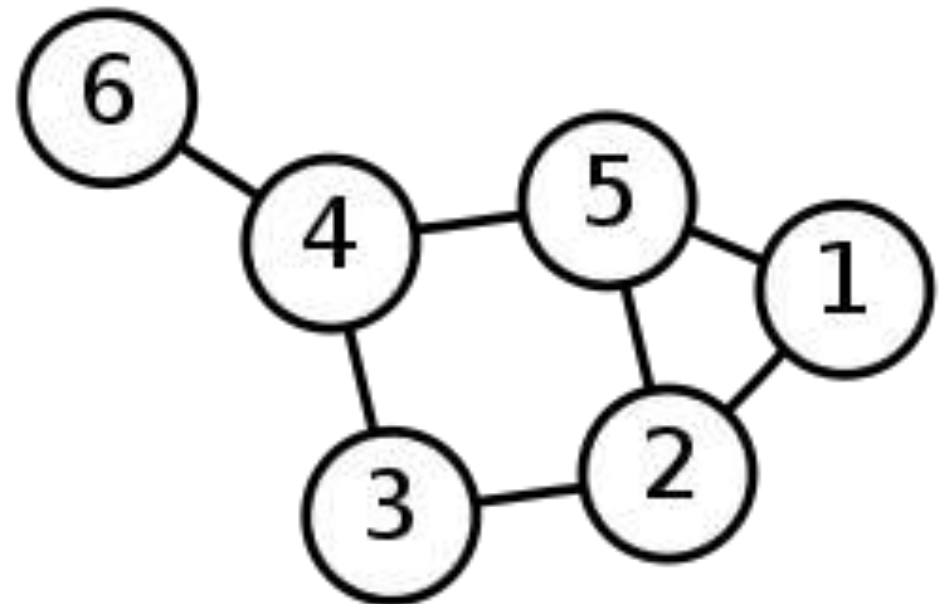
$$\begin{pmatrix} V_1 \\ V_2 \\ \dots \\ V_n \end{pmatrix} (V_1, V_2, \dots, V_n) = \begin{pmatrix} V_1^T V_1, V_1^T V_2, \dots, V_1^T V_n \\ V_2^T V_1, V_2^T V_2, \dots, V_2^T V_n \\ \dots, \dots, \dots, \dots, \\ V_n^T V_1, V_n^T V_2, \dots, V_n^T V_n \end{pmatrix}$$

内積行列から座標を復元する

- 固有値分解
- $H = V^{-1} \Sigma V$
- Σ は対角行列。対角成分は固有値
- 座標は $S V$
 - ただし、 $\Sigma = S S$ 、 S は固有値の平方根を対角成分とする対角行列
- 固有値がすべて非負なら、「ちゃんとした座標」に復元できる
- すべての対象は n 次元球面上の点となる
- 配列の例の場合は、うまくこのようにできる
- 場合によっては負の固有値が出る → 球面よりさらに曲がった空間を考えることになる

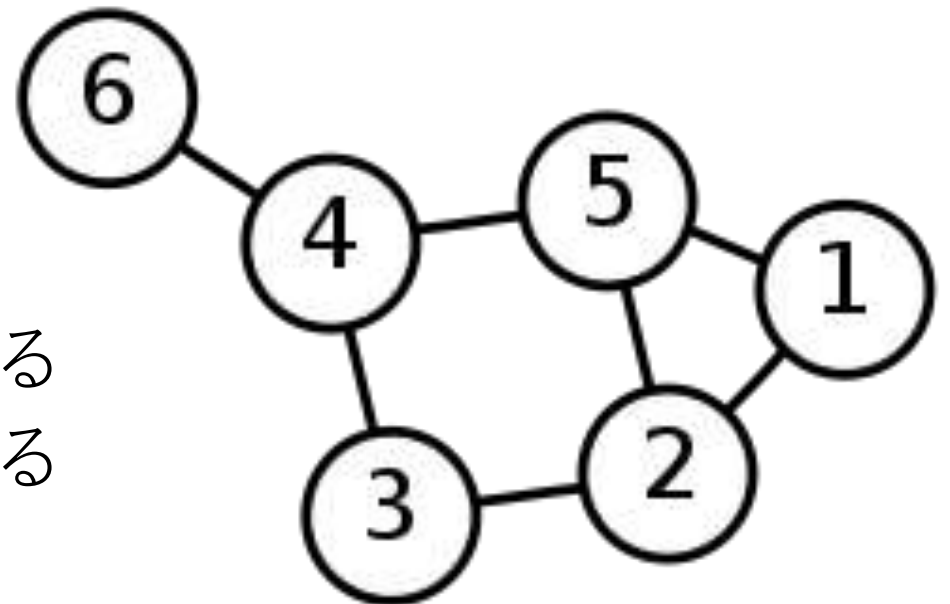
グラフ

- グラフは、ノードの集合とエッジの集合 からできている
- エッジはノードのペアに関して定まる
- グラフは「複数の対象 = ノード」と、その「ペア = エッジ」を現したものの
- ${}_N C_2 = N(N-1)/2$ で済ませる仕組み



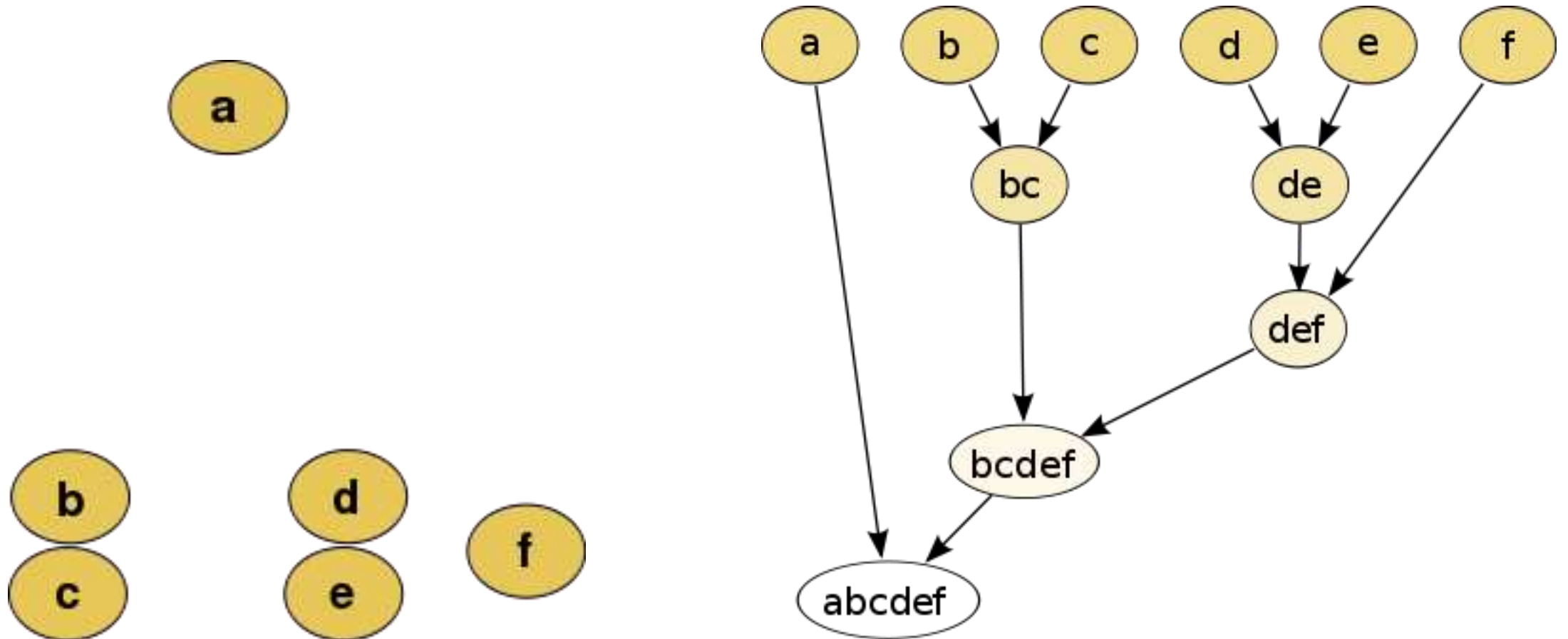
グラフ

- グラフは、ノードの集合とエッジの集合 からできている
- エッジはノードのペアに関して定まる
- グラフは「複数の対象 = ノード」と、その「ペア = エッジ」を現したものの
- ${}_N C_2 = N(N-1)/2$ で済ませる仕組み
- エッジには長さを持たせることができる
- グラフ上の頂点間には距離が定義できる



距離行列からグラフを復元する

- 階層的クラスタリング

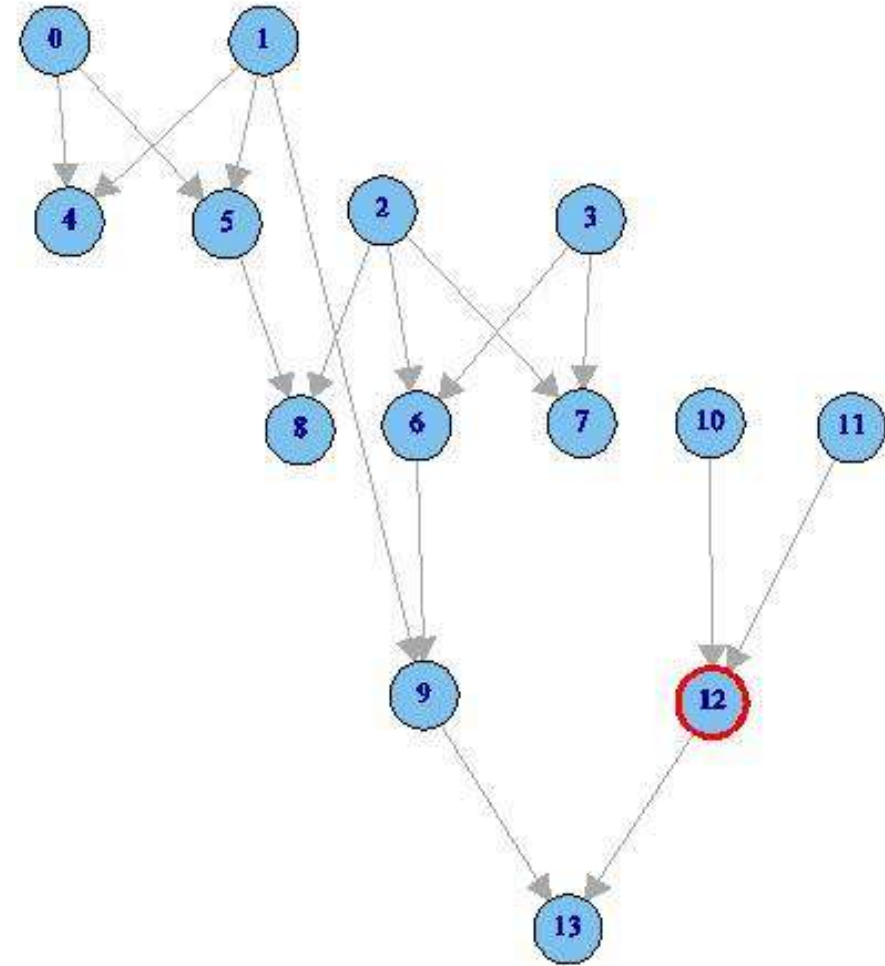
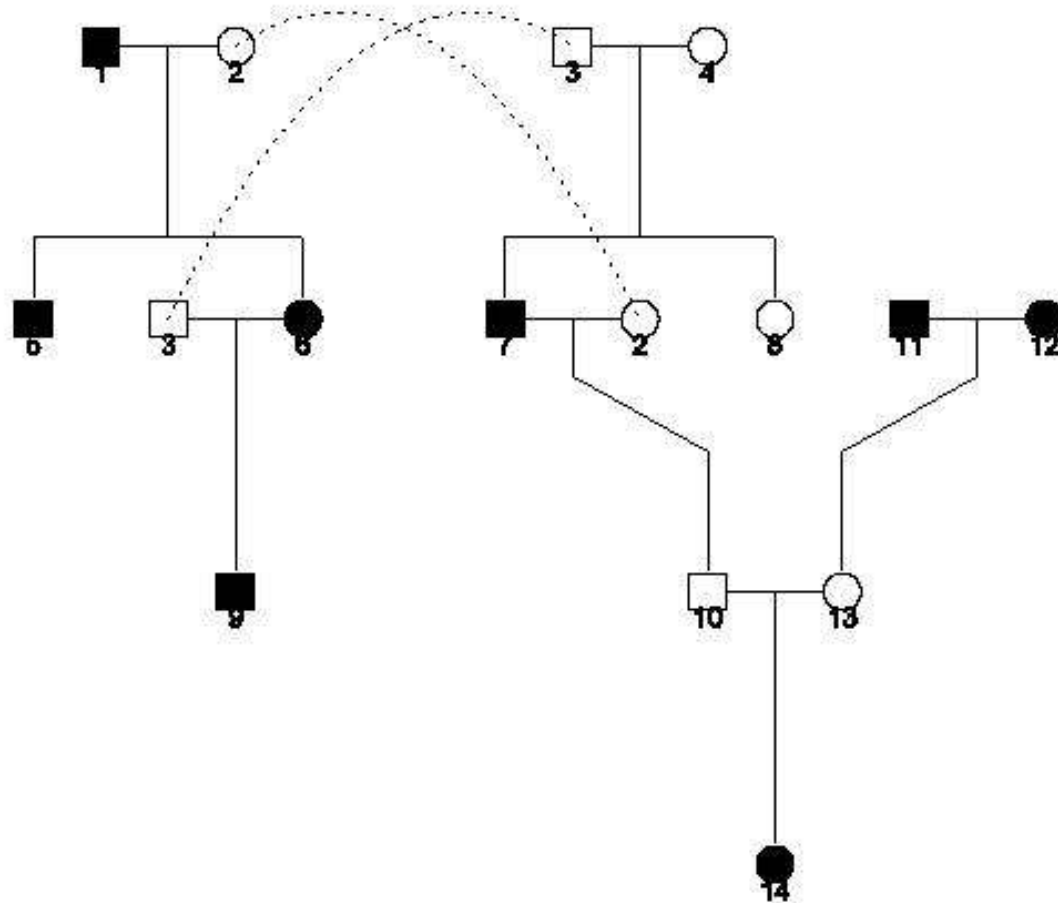


オリジナルが木構造で、その木(グラフ)上の
グラフ距離が距離行列になっていれば、
完璧な木グラフが復元できる

オリジナルがユークリッド空間に配置されて
いればMDSで完璧なユークリッド座標が復元
できた

- グラフ距離が、木グラフ上の距離を満足していない場合には、
ずれをどこかで吸収する必要がある
- 階層的クラスタリングのアルゴリズムがその不一致を吸収する

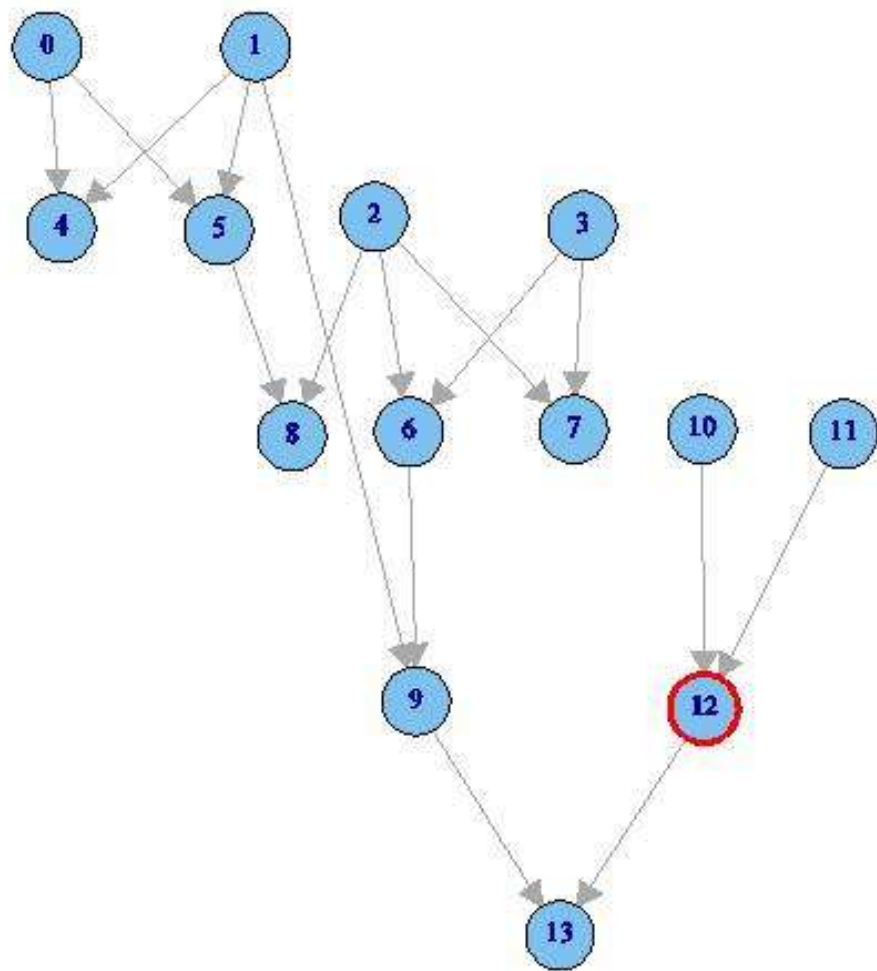
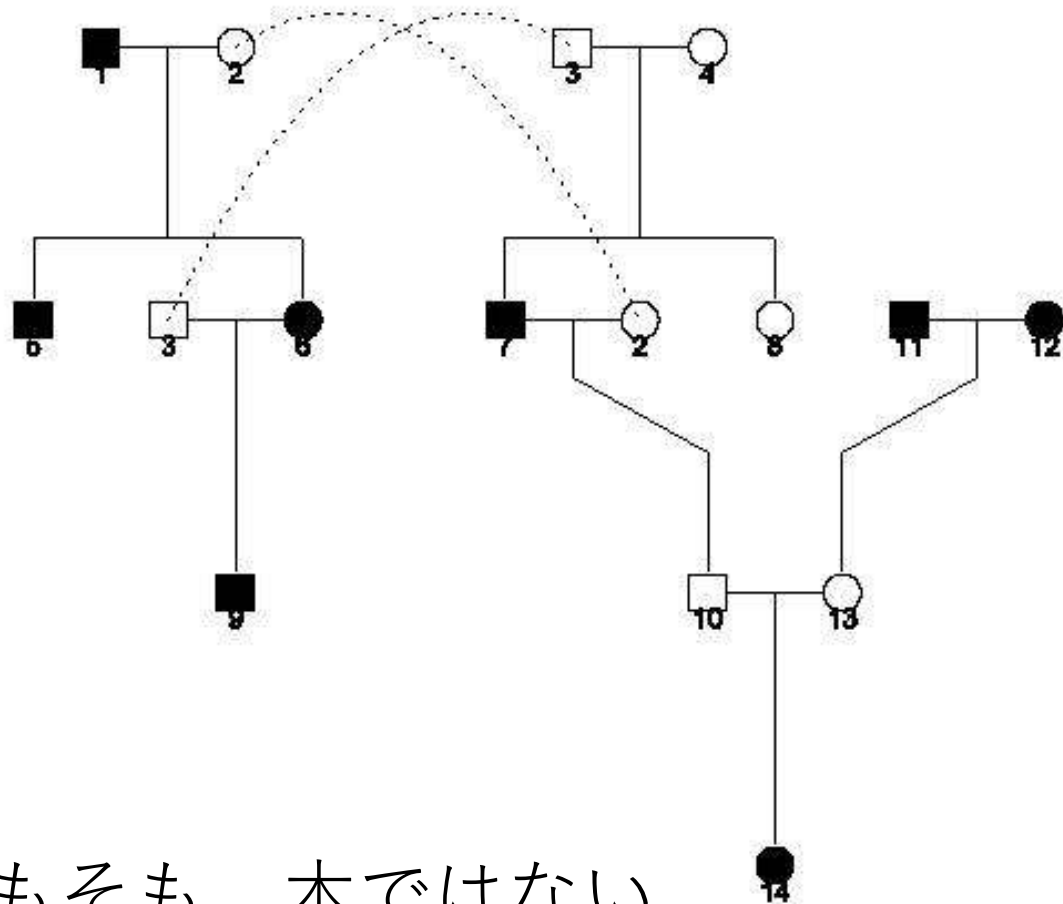
問題、再登場



親子は1/2の関係

子によって「かすがい」された「夫婦」は、1/4の関係なのか？

家系図



そもそも、木ではない
木は

$$\text{エッジの本数} = \text{ノードの数} - 1$$

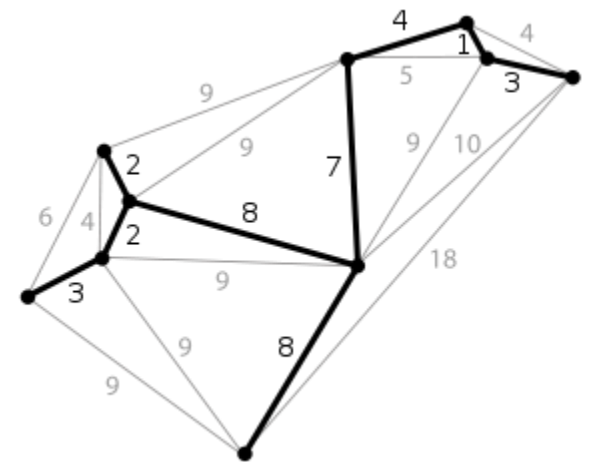
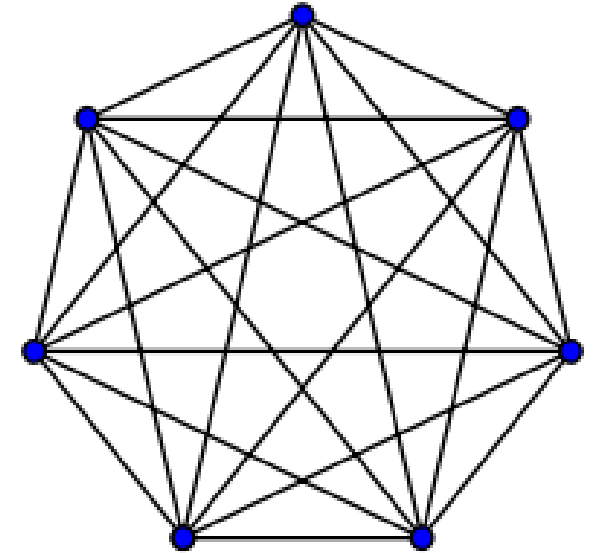
染色体・配列なら木にできる

- 親の2本の染色体配列と、子の染色体配列とは、1/2一致の関係
- 祖父母の染色体配列と、孫の染色体配列とは、1/4一致の関係
- 配列一致程度、そのものを使うと、「グラフ上の距離」との関係がうまくいかない
- 対数を取ると…
 - 親子の染色体の違いは $\log(1/2) = -1$
 - 祖父母-孫の違いは、 $\log(1/4) = -2$

配列の共有度の対数が作る距離行列(非類似度行列)から「木」を作る？

- 最小全域木

- すべての配列ペアに「距離」がある
- 配列をノードにして、すべての配列ペアにエッジを持たせる
 - すべてのノードペアにエッジがあるグラフを完全グラフと呼ぶ
- すべてのノードを含む木であって、エッジの長さの総和が最小になるような木を最小全域木という
- 最小全域木を見つけることは(計算機にとって)それほど大変ではない



配列の非類似度から最小全域木

- 祖父母・親子3代の配列があるとき、世代間を結ぶことは最小全域木の連結の実現になる(ことが多い)
- 親世代が抜けているときには、「世代を飛び越して」つないでやることもあり
- 非血縁者が結ばれることは、(ほぼ)ない

- すくなくとも、大まかな、配列伝達関係の骨格を、かなり単純に作り出せそう
- 血族婚により、染色体・配列の木構造が破綻しサイクルが生じる
 - 最小全域木を作った後で、木に採用されなかった、「強い類似度」のエッジを復活させることで血族婚に由来するサイクルの復活も可能か…

まとめ

- 多数の対象をペア間の類似度・非類似度の観察・定量することは、様々な場面で採用される方法
- オリジナルたちが、収まっていた構造(ユークリッド空間・木構造)を忠実に反映した距離行列が得られていれば、完璧な復元ができる
- そうでない場合は、歪みが生じる(負の固有値、曲がった空間、クラスタリングアルゴリズムによる残差処理)
- 血縁関係の場合は、2倍体による特殊性に注意
- 配列の場合は、対数を導入することで距離化できる
- 血族婚はサイクルを作る