

確率モデルによる配列解析

津本 周作

島根大学医学部医療情報学講座

平成16年7月21日

Outline

分子生物学の基本事項

Bioinformaticsに必要な確率統計の基礎

配列アラインメントとスコア関数

隠れマルコフモデル

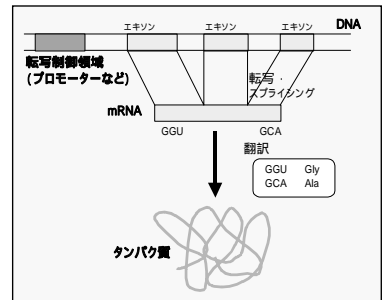
EMアルゴリズム

ゲノム研究と情報科学

- ゲノム解析計画の急速な進展
 - 既に数十種の微生物の全塩基配列が決定
 - 線虫(1000細胞, C.elegans)他、ショウジョウバエも決定
 - ヒトも概要配列が公開済み
- 情報解析の必要性
 - DNA配列 プログラムのオブジェクトコード
 - 意味の解析が必要
- DNA配列: A,C,G,Tの4文字からなる文字列
 - 様々な情報技術が利用可能

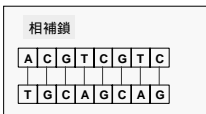
遺伝子と蛋白質

- 遺伝情報の流れ
 - DNA RNA タンパク
- 遺伝子
 - DNA配列中で直接的に機能する部分
- ゲノム
 - 染色体全体(半数体)
 - 遺伝情報の総体
- タンパク質
 - アミノ酸(20種類)の鎖



DNAとアミノ酸

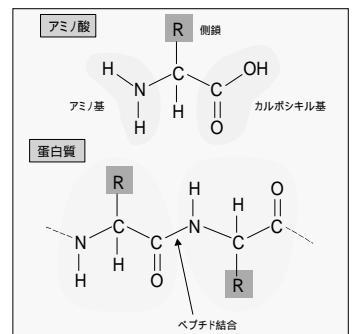
- DNAはA,C,G,Tの4文字の並び
- DNAは二重らせん構造 相補鎖
- 塩基: DNA1文字、残基: アミノ酸1文字
- DNA3文字がアミノ酸1文字に対応 (アミノ酸は20種類)



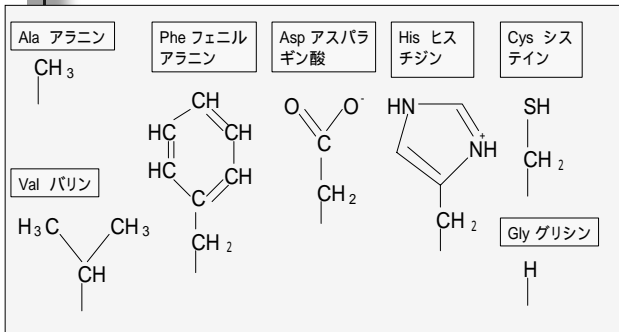
コード表		2文字目			
		T	C	A	G
1文字目	T	TTT F TTC TTA L TTG	TCT S TCC TCA TCG	TAT Y TAC TAA stop TAG	TGT C TGC TGA stop TGG W
	C	CTT L CTC CTA CTG	CCT P CCG CCA CCG	CAT H CAC CAA CAG	CGT R CGC CGA CGG
	A	ATT I ATC ATA ATG	ACT T ACC ACA ACG	AAT N AAC AAA AAG	AGT S AGC AGA AGG
G	GTT V GTC GTA GTG	GCT A GCC GCA GCG	GAT D GAC GAA GAG	GGT G GGC GGA GGG	

アミノ酸と蛋白質

- アミノ酸: 20種類
- 蛋白質: アミノ酸の鎖(短いものはペプチドと呼ばれる)
- 側鎖によって性質が異なる.



側鎖の例



アミノ酸コード表

Ala	A	アラニン	Leu	L	ロイシン
Arg	R	アルギニン	Lys	K	リシン
Asn	N	アスパラギン	Met	M	メチオニン
Asp	D	アスパラギン酸	Phe	F	フェニルアラニン
Cys	C	システイン	Pro	P	プロリン
Gln	Q	グルタミン	Ser	S	セリン
Glu	E	グルタミン酸	Thr	T	トレオニン
Gly	G	グリシン	Trp	W	トリプトファン
His	H	ヒスチジン	Tyr	Y	チロシン
Ile	I	イソロイシン	Val	V	バリン

突然変異とアミノ酸置換

- TTT: F(フェニルアラニン)
- TTC: F
- 三番目の文字がC,T間で置換しても、アミノ酸としては変わらない。
- TTA, TTG: L(ロイシン)
- 側鎖の大きさは異なるが疎水性という性質は変わらない。
- 以上から、TT*というコードは、3文字目が置換してもタンパク質の性質を変えるものではない。
- ATT (I: イソロイシン), AAT(N: アスパラギン)では、疎水性と親水性でかなり異なるため、2文字目が置換すると、タンパク質の性質を変える可能性がある。

		2文字目			
		T	C	A	G
1文字目	T	TTT F TTC F TTA L TTG L	TCT S TCC S TCA S TCG S	TAT Y TAC Y TAA stop TAG stop	TGT C TGC C TGA stop TGG W
	C	CTT L CTC L CTA L CTG L	CCT P CCC P CCA P CCG P	CAT H CAC H CAA Q CAG Q	CGT R CGC R CGA R CGG R
	A	ATT I ATC I ATA I ATG M	ACT T ACC T ACA T ACG T	AAT N AAC N AAA K AAG K	AGT S AGC S AGA R AGG R
G	GTT V GTC V GTA V GTG V	GCT A GCC A GCA A GCG A	GAT D GAC D GAA E GAG E	GGT G GGC G GGA G GGG G	

Bioinformaticsにおける確率統計

- 重要なのはデータからのモデル(もしくはパラメータ)の推定
 - 最尤法
 - ベイズ推定
 - 最大事後確率推定 (MAP)
- 通常、データが持っている情報よりも多い量の情報(モデル)を推測することになる
 - このデータを満たしうるモデルが数多く存在しうる。
 - パラメータの推定: モデルの選択 EM Algorithm

ベイズの定理

- 重要なのはデータからのモデル(もしくはパラメータ)の推定

$$P(\theta, D) = P(\theta | D)P(D) = P(D | \theta)P(\theta)$$

θ : パラメータ(モデル), D : データ

$P(\theta, D)$ は、データ D とモデルのパラメータ θ が同時に起こりうる確率

ベイズの定理 (2)

- 重要なのはデータからのモデル(もしくはパラメータ)の推定

$$P(\theta, D) = P(\theta | D)P(D) = P(D | \theta)P(\theta)$$

$P(D | \theta)$: 尤度: モデルが成立した時のデータが起こる条件付き確率

$P(\theta | D)$: 事後確率: データが観察された時にモデルが成立する条件付き確率

最尤推定

- $P(D|\theta)$ (尤度)
 - モデルパラメータ θ のもとでのデータ D の出現確率
- 最尤法
 - $P(D|\theta)$ を最大化する θ を選ぶ
- 例
 - コインを5回投げて、表が3回出た後、裏が2回出た
 - $p(\text{表})=a, p(\text{裏})=1-a$ とすると $P(D|\theta)=_5C_3 a^3 (1-a)^2$
 - $a=3/5$ の時、 $P(D|\theta)$ は最大
 - 一般に表が出る頻度を f とすると $a=f$ で尤度は最大

ベイズ推定とMAP推定

- ベイズ推定: 尤度とモデル(パラメータ)の事前確率から、ベイズの定理により、事後確率を推定

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

ただし、 $P(D) = \int_0^1 P(D|\theta')P(\theta')$ (θ が連続値の時)

- 最大事後確率 (MAP) 推定
 - $P(D|\theta)P(\theta)$ を最大化する θ を計算
 - $P(\theta)$ が一様分布なら最尤推定と同じ

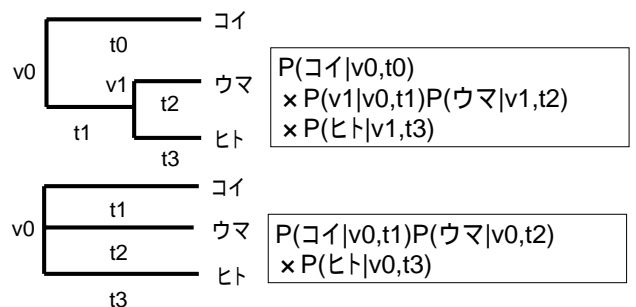
不正サイコロのベイズ推定

- 公正サイコロと不正サイコロ
 - 公正: $P(i|\text{公正})=1/6$
 - 不正: $P(6|\text{不正})=1/2, P(i|\text{不正})=1/10$ for $i \neq 6$
 - $P(\text{公正})=0.99, P(\text{不正})=0.01$
- 6が3回続けて出た場合の事後確率

$$P(\text{不正} | 666) = \frac{P(666 | \text{不正})P(\text{不正})}{P(666)}$$

$$= \frac{(0.5)^3(0.01)}{(0.5)^3(0.01) + (\frac{1}{6})^3(0.99)} = 0.21$$

モデルから尤度の計算 (2)



推定が必要なパラメータは上が6個, 下が4個

確率モデルによる配列解析

- 配列データのレコード数は少ない
 - 例えば, 前の例だと3レコード
- レコードが含んでいる項目数は多い
 - ヘモグロビンだと, 184個のアミノ酸
- 変数の数は多いが, 方程式の数が少ない
 - 不定要素が多い. (一意的に決定できない)
- 確率モデルでのパラメータ推定
 - 不定要素が多いため, 当てはまるモデルが多数存在する.
 - 不定要素が多い中でのパラメータ推定

配列解析

異種の生物の配列解析: 遺伝子の対応づけの必要性 (アライメント)

対応づけが済んだあとの配列の解析

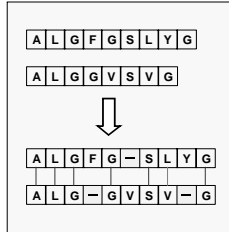
系統樹推定

隠れマルコフモデル

Neural Network 他

アラインメント

- 配列解析の最重要技術
- 2個もしくは3個以上の配列の類似性の判定に利用
- 文字間の最適な対応関係を求める (最適化問題)
- 2個の配列長を同じにするように、ギャップ記号 (挿入、欠失に対応) を挿入



アラインメント

- 遺伝子の欠失, 挿入は進化の過程で起こりうる
しかし, 単に欠失, 挿入だけでは, タンパク質の構造は変わらない.
- 遺伝子の配列から逆に欠失, 挿入の可能性を推定しよう.
- ABCDEF ABCDEF -BCDEF- -BCDEF-
- CDF -- CDE- AB-EEFG ABEEFG
- 文字間の最適な対応関係を求める (最適化問題)
- 2個の配列長を同じにするように、ギャップ記号 (挿入、欠失に対応) を挿入

アラインメント (例:ヘモグロビン)

```
> human
MVHLTPEEKSAVTALWGKVNVEVGGEEALGRLLVYYPWTQRFESFGDLSTPDAVMGNPK
VKAHGKKVLGAFSDGLAHLNLTGKFTALSELHCDKLVDPENFRLLGNVLVCLAHHFG
KEFTPPVQAAYQKVVAGVANALAHKYH
> gorilla
MVHLTPEEKSAVTALWGKVNVEVGGEEALGRLLVYYPWTQRFESFGDLSTPDAVMGNPK
VKAHGKKVLGAFSDGLAHLNLTGKFTALSELHCDKLVDPENFRLLGNVLVCLAHHFG
KEFTPPVQAAYQKVVAGVANALAHKYH
> horse
VQLSGECKAVALALWDKVNVEEVEGGEEALGRLLVYYPWTQRFESFGDLSTPDAVMGNPKV
KAHGKVLHSGEGVHLDNLKGTFAALSELHCDKLVDPENFRLLGNVLVWLRHFGK
DFTPELQASQKVVAGVANALAHKYH
> chicken
MVHWTAEKQLITGLWGVNVAECGAELARLLVYYPWTQRFESFGDLSTPDAVMGNPK
VRAHGKVLTSFGDAVKNDLNKFTSLSLHCDKLVDPENFRLLGDLIIVLAAHFS
KDFTEPCQAQWQKLVVVAHALARKYH
> eastern gray kangaroo
VHLTAEKNAITSLWGKVAIEQTGGEEALGRLLVYYPWTSRFFDFHGDLSNAKAVMANPKV
LAHGAKVLVAFGDAIKNDLNKGTFAALSELHCDKLVDPENFRLLGNVIICLAEHFGK
EFTIDTQVAWQKLVAGVANALAHKYH
```

アラインメント (例:ヘモグロビン)

```
human
MVHLTPEEKSAVTALWGKVNVEVGGEEALGRLLVYYPWTQRFESFGDLSTPDAVMGNPK
gorilla
MVHLTPEEKSAVTALWGKVNVEVGGEEALGRLLVYYPWTQRFESFGDLSTPDAVMGNPK
horse
-VQLSGECKAVALALWDKVNVEEVEGGEEALGRLLVYYPWTQRFESFGDLSTPDAVMGNPK
eastern
-VHLTAEKNAITSLWGKVAIEQTGGEEALGRLLVYYPWTSRFFDFHGDLSNAKAVMANPK
chicken
MVHWTAEKQLITGLWGVNVAECGAELARLLVYYPWTQRFESFGDLSTPDAVMGNPK
*: : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : *
```

スコア行列 (置換行列)

- 残基間 (アミノ酸文字間) の類似性を表す行列
 - PAM250, BLOSUM45 など

	A	R	N	D	C	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	5	-2	-1	-2	-1	-1	0	-2	-1	-2	-1	-3	-1	0	-3	-2	0		
R	-2	7	-1	-2	-4	1	0	-3	0	-4	-3	3	-2	-3	-3	-1	-1	-3	-1
N	-1	-1	7	2	-2	0	0	0	1	-3	-4	0	-2	-4	-2	1	0	-4	-2
D	-2	-2	2	8	-4	0	2	-1	-1	-4	-4	-1	-4	-5	-1	0	-1	-5	-3
C	-1	-4	-2	-4	13	-3	-3	-3	-2	-2	-3	-2	-4	-1	-1	-5	-3	-1	
Q	-1	1	0	0	-3	7	2	-2	-1	-3	2	0	-4	-1	0	-1	-1	-1	-3
E																			
G																			
H																			
I																			
L																			
K																			
M																			
F																			
P																			
S																			
T																			
W																			
Y																			
V																			

BLOSUM50 スコア行列 (置換行列) の一部分

ギャップ無しアラインメントのスコア

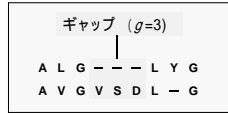
$$S = \log \frac{P(x,y|M)}{P(x,y|R)} = \log \frac{\prod_i p_{x_i y_i}}{\prod_i q_{x_i} \cdot \prod_i q_{y_i}}$$

$$= \log \prod_i \frac{p_{x_i y_i}}{q_{x_i} q_{y_i}} = \sum_i s(x_i, y_i) \quad \text{但し, } s(a,b) = \log \left(\frac{p_{ab}}{q_a q_b} \right)$$

- スコア: $P(x,y|M)$ と $P(x,y|R)$ の対数オッズ比
- $P(x,y|R)$: ランダムモデル R において、配列 x,y が独立に生じる確率
- q_a : 文字 a の R における出現確率
- $P(x,y|M)$: 一致モデル M において、配列 x,y が生じる確率
- p_{ab} : 文字ペア a,b の M における出現確率

ギャップペナルティ

- 線形スコア
 - $(g) = -gd$
(g はギャップの長さ、 d は定数)
- アファインギャップスコア
 - $(g) = -d - e(g-1)$ (d :ギャップ開始ペナルティ, e :ギャップ伸長ペナルティ)
- ギャップの確率的解釈
 - $P(gap) = f(g) \quad q_{xi} \rightarrow (g) = \log(f(g))$
 - ギャップ長の確率の対数



ペアワイズ・アラインメント

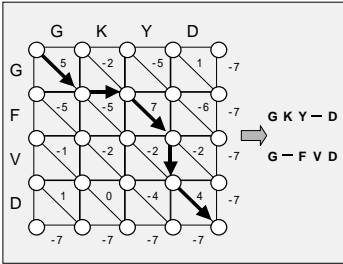
- 配列が2個の場合でも可能なアラインメントの個数は指数オーダー
- しかし、スコア最大となるアラインメント(最適アラインメント)は動的計画法により、 $O(mn)$ 時間で計算可能 (m, n :入力配列の長さ)

入力配列	AGCT, ACGCT	最適アラ インメント	
アラインメント	AGCT - AG - CT ACGCT ACGCT	A - GCT ACGCT	- AGC - - T AC - - GCT
スコア	-3	1	3 -5

(同じ文字の時1, 違う文字の時-1, ギャップ1文字-1)

動的計画法による最適アラインメント

最長経路計算による最適アラインメント



DP (動的計画法)による最長経路の計算

$$F(0, j) = -jd, \quad F(i, 0) = -id$$

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) - d \\ F(i, j-1) - d \end{cases}$$

行列表からの経路の復元は、 $F(m, n)$ からmaxでととなっている $F(i, j)$ を逆にたどることを行う (トレースバック)

局所アラインメント

- 配列の一部のみに共通部分があることが多い
- 共通部分のみのアラインメント

$$F(i, j) = \max \begin{cases} 0 \\ F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) - d \\ F(i, j-1) - d \end{cases}$$

	H	E	A	W	G	E	H
G	0	0	0	0	0	0	0
F	0	0	0	0	0	1	0
A	0	0	0	1	0	0	0
W	0	0	0	0	2	-1	0
E	0	0	1	0	1	0	2
D	0	0	0	0	0	1	1

AWGE
AW - E

(スコア: ギャップ-1, 置換-1, 一致1)

アラインメント (例:ヘモグロビン)

```

human
MVHLTPPEEKSAVTALWGKVVNDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPK
gorilla
MVHLTPPEEKSAVTALWGKVVNDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPK
horse
-VQLSGEELAAVLAALWVKVNEEEVGGGEALGRLLVVYPWTQRFFESFGDLSPGAVMGNPK
eastern
-VHLTAEKNAI TSLWIKVA I EQTGGEALGRLLI VYPWTQRFFDFHFGDLSSNAKAVMANPK
chicken
MVHWITAEKQL I TGLWIKVNVAECEGAEALARLL I VYPWTQRFFASFGNLSSTPA I LGNPM
* : : * * : : * * * : : * * * * : * * * * * * * * * * * : : * *
    
```

スコアのベイズ論的解釈

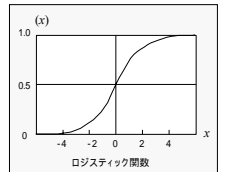
- $P(M|x, y)$ をベイズの定理に基づき導出
- $P(M)$: 2個の配列に関連性がある事前確率
- $P(R) = 1 - P(M)$: ランダムモデルの事前確率
- 以下の式より尤度比にオッズ比を掛け合わせたものを0と比較すべきであることがわかる

$$P(M|x, y) = \frac{P(x, y|M)P(M)}{P(x, y)}$$

$$= \frac{P(x, y|M)P(M)}{P(x, y|M)P(M) + P(x, y|R)P(R)}$$

$$= \frac{P(x, y|M)P(M) / P(x, y|R)P(R)}{1 + P(x, y|M)P(M) / P(x, y|R)P(R)} = \sigma(S')$$

但し、 $S' = \log \left(\frac{P(x, y|M)}{P(x, y|R)} \right) + \log \left(\frac{P(M)}{P(R)} \right)$, $\sigma(x) = \frac{e^x}{1 + e^x}$



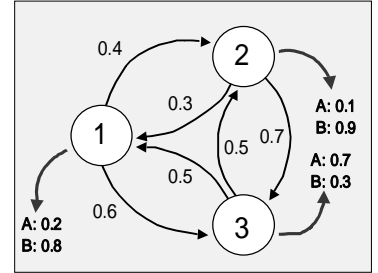
スコア行列の導出

- 基本的には頻度の比の対数をスコアとする
- BLOSUM行列
 - 既存のスコア行列を用いて多くの配列のアラインメントを求め、ギャップ無しの領域(ブロック)を集める
 - 残基がL%以上一致しているものを同一クラスタに集める
 - 同じクラスタ内で残基aが残基bにアラインされる頻度 A_{ab} を計算
 - $q_a = \sum_b A_{ab} / \sum_{cd} A_{cd}$, $p_{ab} = A_{ab} / \sum_{cd} A_{cd}$ を求め、 $s(a,b) = \log(p_{ab}/q_a q_b)$ としたのち、スケーリングし近傍の整数値に丸める

隠れマルコフモデル(HMM)

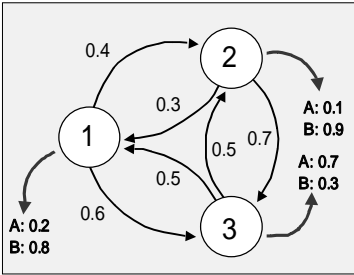
- HMM=有限オートマトン+確率
- 定義

- 出力記号集合
- 状態集合 $S = \{1, 2, \dots, n\}$
- 遷移確率($k \rightarrow l$) a_{kl}
- 出力確率 $e_k(b)$
- (開始状態= 終了状態= 0)

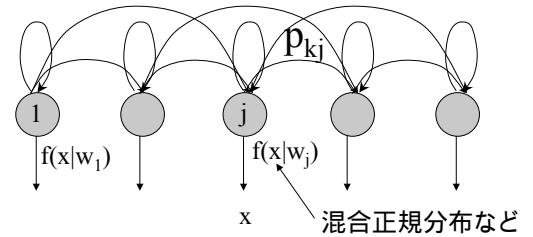


隠れマルコフモデル(HMM) (2)

- 音声認識について革命的な変化をもたらした。
- それまでは認識率は50%も出ればよい方。
- HMMによって、80%以上の正答率が実現した。

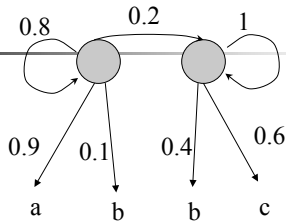


隠れマルコフモデルとは



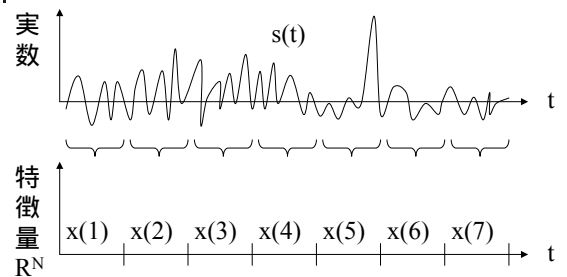
外部からは状態については不明
 $x_1 x_2 x_3 x_4 \dots x_t \dots$ (時系列)

例



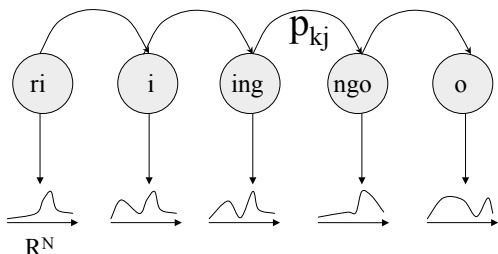
$$\begin{aligned}
 &P(\text{aaabccc}) \\
 &= 0.9 * 0.8 * 0.9 * 0.8 * 0.9 \\
 &\quad * (0.8 * 0.1 * 0.2 + 0.2 * 0.4) * 0.6 * 0.6 * 0.6
 \end{aligned}$$

使い方(音声認識)



使い方 (音声認識)

「りんご」 = 「ri i ing ngo o」を発生するモデル



使い方 (音声認識)

「りんご」 $p(x_1, x_2, \dots, x_T | p_{kj}, w_j)$

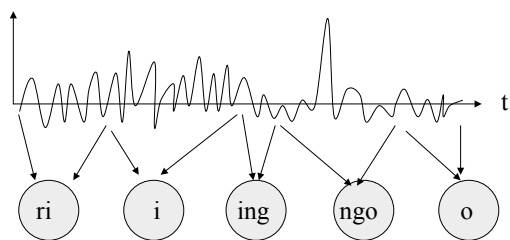
「るんご」 $q(x_1, x_2, \dots, x_T | p_{kj}, w_j)$

「りんぐ」 $r(x_1, x_2, \dots, x_T | p_{kj}, w_j)$

特徴量を入力して発生確率の高い候補を認識結果として選ぶ。

利点は？

音声は、時間方向に伸縮がおおきい：
変化に対応できる りーんご、りんご、りーんごー、



隠れマルコフモデルの意義

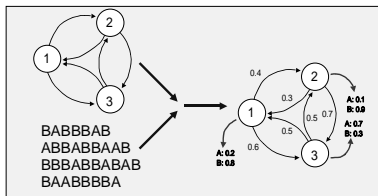
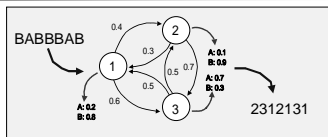
形式言語

有限オートマトン
確率オートマトン
確率的な情報発生装置

実世界情報

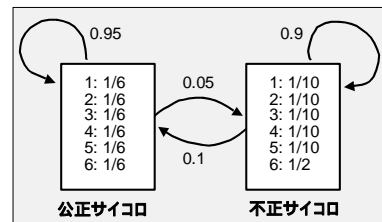
HMMにおける基本アルゴリズム

- Viterbiアルゴリズム
 - 出力記号列から状態列を推定
 - Parsing (構文解析)
- Baum-Welchアルゴリズム (EMアルゴリズム)
 - 出力記号列からパラメータを推定
 - Learning (学習)



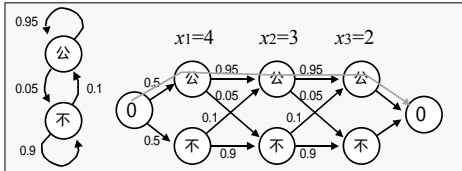
時々いかさまをするカジノ

- サイコロの出目だけが観測可能、どちらのサイコロを振っているかは観測不可能
- サイコロの出目から、どちらのサイコロを振っているかを推定
- 6,2,6,6,3,6,6,6, 4,6,5,3,6,6,1,2 不正サイコロ
- 6,1,5,3,2,4,6,3, 2,2,5,4,1,6,3,4 公正サイコロ
- 6,6,3,6,5,6,6,1, 5,4,2,3,6,1,5,2 途中で公正サイコロに交換



Viterbi アルゴリズム(1)

- 観測列 (出力配列データ) $x=x_1 \dots x_L$ と状態列 $= 1 \dots L$ が与えられた時、その同時確率は $P(x, \pi) = a_0 \prod_{i=1}^L e_i(x_i) a_{i+1}$ 但し、 $L+1=0$
- x が与えられた時の、最も尤もらしい状態列は $\pi^* = \text{argmax}_{\pi} P(x, \pi)$
- 例: どちらのサイコロがいつ使われたかを推定



$$\max_{\pi} P(x_1, x_2, x_3, \pi) = P(x_1, x_2, x_3, \text{公公公}) = 0.5 \cdot \frac{1}{6} \cdot 0.95 \cdot \frac{1}{6} \cdot 0.95 \cdot \frac{1}{6}$$

Viterbi アルゴリズム(2)

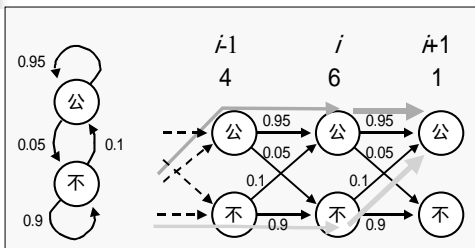
- x から、 $\pi^* = \text{argmax}_{\pi} P(x, \pi)$ を計算
- そのためには $x_1 \dots x_i$ を出力し状態 k に至る確率最大の状態列の確率 $v_k(i)$ を計算

$$v_k(i) = \max_{\pi} a_{0\pi_j} \prod_{j=1}^{i-1} e_{\pi_j}(x_j) a_{\pi_j \pi_{j+1}}$$

- $v_k(i)$ は以下の式に基づき動的計画法で計算

$$v_k(i+1) = e_i(x_{i+1}) \max_k (v_k(i) a_{kj})$$

Viterbi アルゴリズム(3)



$$v_{\text{公}}(i+1) = \max \{ e_{\text{公}}(1) \cdot 0.95 \cdot v_{\text{公}}(i), e_{\text{公}}(1) \cdot 0.1 \cdot v_{\text{不}}(i) \}$$

EM (Expectation Maximization) アルゴリズム

- 「欠けているデータ」のある場合の最尤推定のための一般的アルゴリズム

x : 観測データ、 y : 欠けているデータ、
 θ : パラメータ集合
 目標: $\log P(x|\theta) = \log \sum_y P(x, y|\theta)$ の最大化

- 最大化は困難であるので、反復により尤度を単調増加させる (t より $t+1$ を計算)
- HMMの場合、「欠けているデータ」は状態列

EM アルゴリズムの導出

$\log P(x|\theta) = \log P(x, y|\theta) - \log P(y|x, \theta)$
 両辺に $P(y|x, \theta')$ をかけて y についての和をとり、
 $\log P(x|\theta) = \sum_y P(y|x, \theta') \log P(x, y|\theta) - \sum_y P(y|x, \theta') \log P(y|x, \theta)$
 右辺第1項を $Q(\theta|\theta')$ とおくと、
 $\log P(x|\theta) - \log P(x|\theta') =$
 $Q(\theta|\theta') - Q(\theta'|\theta') + \sum_y P(y|x, \theta') \log \frac{P(y|x, \theta')}{P(y|x, \theta)}$
 最後の項は相対エントロピーで常に正なので、
 $\log P(x|\theta) - \log P(x|\theta') \geq Q(\theta|\theta') - Q(\theta'|\theta')$
 よって、 $\text{argmax}_{\theta} Q(\theta|\theta')$ をみたとす θ^{t+1} をとれば尤度は増大

EM アルゴリズムの一般形

- 初期パラメータ θ^0 を決定。 $t=0$ とする。
- $Q(\theta|\theta^t) = \sum_y P(y|x, \theta^t) \log P(x, y|\theta)$ を計算。 (Expectation)
- $Q(\theta|\theta^t)$ を最大化する θ^* を計算し、
 $\theta^{t+1} = \theta^*$ とする。 $t=t+1$ とする。
 (Maximization)
- Q が増大しなくなるまで、2, 3 を繰り返す。

前向きアルゴリズム

- 配列 x の生成確率 $P(x) = P(x_1, \dots, x_L)$ を計算

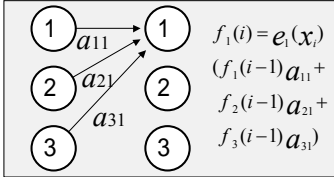
$$f_0(0) = 1, f_k(0) = 0$$

$$f_i(i) = e_i(x_i) \sum_k f_k(i-1) a_{ki}$$

$$P(x) = \sum_k f_k(L) a_{k0}$$

- Viterbi アルゴリズムと類似

- $f_k(i) = P(x_1, \dots, x_i, i=k)$ を DP により計算



後向きアルゴリズム

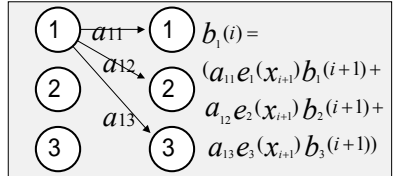
- $b_k(i) = P(x_{i+1}, \dots, x_L | i=k)$ を DP により計算

$$b_k(L) = a_{k0}$$

$$b_k(i) = \sum_l a_{kl} e_l(x_{i+1}) b_l(i+1)$$

$$P(x) = \sum_k a_{k0} e_k(x_1) b_k(1)$$

- $P(i=k|x) = f_k(i) b_k(i) / P(x)$



HMM に対する EM アルゴリズム (Baum-Welch アルゴリズム)

A_{kl} : a_{kl} が使われる回数の期待値 x^j : j 番目の配列
 $E_k(b)$: 文字 b が状態 k から現れる回数の期待値

$$A_{kl} = \sum_j \frac{1}{P(x^j)} \sum_i f_k^j(i) a_{kl} e_l(x_{i+1}^j) b_l^j(i+1)$$

$$E_k(b) = \sum_j \frac{1}{P(x^j)} \sum_{\{i|x_i^j=b\}} f_k^j(i) b_k^j(i)$$

パラメータの更新式

$$\hat{a}_{kl} = \frac{A_{kl}}{\sum_{l'} A_{kl'}} \quad \hat{e}_k(b) = \frac{E_k(b)}{\sum_{b'} E_k(b')}$$

Baum-Welch の EM による解釈

$$P(x, \pi | \theta) = \prod_{k=1}^M \prod_b [e_k(b)]^{E_k(b, \pi)} \prod_{k=0}^M \prod_{l=1}^M a_{kl}^{A_{kl}(\pi)} \quad \text{および}$$

$$Q(\theta | \theta') = \sum_{\pi} P(\pi | x, \theta') \log P(x, \pi | \theta) \quad \text{より、}$$

$$Q(\theta | \theta') = \sum_{\pi} P(\pi | x, \theta') \times \left[\sum_{k=1}^M \sum_b E_k(b, \pi) \log e_k(b) + \sum_{k=0}^M \sum_{l=1}^M A_{kl}(\pi) \log a_{kl} \right]$$

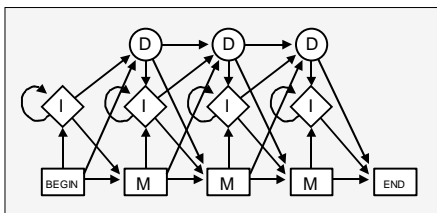
$$= \sum_{k=1}^M \sum_b E_k(b) \log e_k(b) + \sum_{k=0}^M \sum_{l=1}^M A_{kl} \log a_{kl}$$

ここで $\sum_l p_l \log q_l$ は $q_l = p_l$ の時、最大より、

$$e_k(b) = \frac{E_k(b)}{\sum_{b'} E_k(b')}, \quad a_{kl} = \frac{A_{kl}}{\sum_{l'} A_{kl}}$$

プロフィール HMM (1)

- 配列をアラインメントするための HMM
- 一致状態 (M)、欠失状態 (D)、挿入状態 (I) を持つ

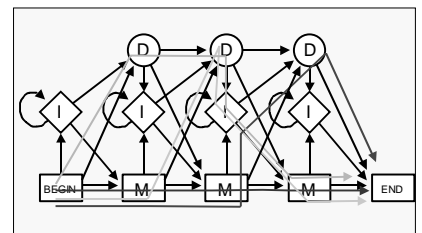


プロフィール HMM (2)

マルチプルアラインメント

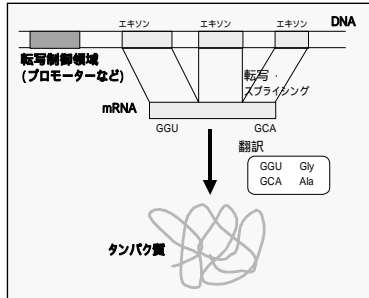
	M	M	...	M		
こうもり	A	G	-	-	C	
ラット	A	-	A	G	-	C
ネコ	A	G	-	A	A	-
ハエ	-	-	A	A	A	C
ヤギ	A	G	-	-	-	C

プロフィール HMM



隠れマルコフモデルの適用

- マルチプルアライメント
- エキソン, イントロン, プロモーターの検出
- スプライス部位の検出
- モチーフの検出
(タンパク質の機能単位:
例えば, ヘモグロビンの鉄結合部位)



まとめ

- DNA, アミノ酸配列における突然変異
- 突然変異の蓄積を確率過程としてとらえる
- 適切な確率モデルを構築して:
 - 配列の類似性
 - 配列の規則性 を検出する .
 - そこでは, ベイズ推定, 確率過程が重要と思われる .
- 配列アライメントとスコア関数
- 隠れマルコフモデル
- EMアルゴリズム