

\*

DNA

.%90

2010/ 12/7 :

---

/ / \*

2010/ 10/ 12:

...

[2]

---

## Human deoxyribonucleic acid (DNA) sequence

### Abstract

The problem of finding genetic mutations is an issue of special importance for its relationship with genetics evolution of living organisms. In this paper we present three statistical methods proposed for the creation of genetic mutations, the first depends mainly on the Markovian model, the second on the autoregressive models, and the third depends on the algorithm proposed based on the hidden Markov model. The three methods are applied to DNA for sequence with shortage oxygen of humans, and shows the practical application show that of these methods give good results in the detection of genetic mutations, reached in some cases to 90%.

:

mutation

( )

.(Watson et al. ,2003)

DNA

:

: substitution mutation

)

mRNA

( messenger

m

DNA

.(Drake ,1970)

:frame shift mutations

ب.

DNA

RNA

RNA

)

(2010

:

... \_\_\_\_\_ [4]  
: -1

.  
: -2

.DNA

Alberts and )

.(other ,2002

:  
G A : : .  
C T : : .  
DNA T U mRNA U

DNA

.A,T,C,G

:

:( )

.p DNA

:

.1

p

$\hat{x}_{T+1}$   $T + 1$  .2

$x_T, x_{T-1}, \dots, x_{T-p+1}$   $T - p + 1 \dots T - 1$

$i_T, i_{T-1}, \dots, i_{T-p+1}$

$i_T, i_{T-1}, \dots, i_{T-p+1}$

$S = \{1, 2, \dots, N\}$

$\underline{P}$

j

$\hat{x}_{T+1} = j$

$i_T, i_{T-1}, \dots, i_{T-p+1}$

$$\hat{x}_{T+1} = \max_{j=1,2,\dots,N} p_{i_T, i_{T-1}, \dots, i_{T-p+1}} \quad (12)$$

$S = \{A, T, C, G\}$

64

( ) 3

:( )

autoregressive model of  $p$

$p$  difference equation (order  $p$ )

$y_1, y_2, \dots, y_T$  .AR( $p$ )

$p$

:(2010 )

$$y_n = a_1 y_{n-1} + a_2 y_{n-2} + \dots + a_p y_{n-p} + e_n; \quad n = p, p+1, \dots \quad (13)$$

$e_n$  ( )  $a_1, a_2, \dots, a_p$

.Noise

:(MATLAB , 2009 )

$$A(q) y_n = e_n, \quad (14)$$

:  $p$  Polynomial  $A(q)$

$$A(q) = 1 - a_1 q^{-1} - a_2 q^{-2} - \dots - a_p q^{-p} \quad (15)$$

backward shift operator  $q^{-1}$

$$q^{-k} y_n = y_{n-k} \quad q^{-2} y_n = y_{n-2} \quad q^{-1} y_n = y_{n-1} :$$

$\{X_n; n = 0, 1, 2, \dots\}$   $x_1, x_2, \dots, x_T$

$$S = \{s_1, s_2, \dots, s_N\}$$

$$y_1, y_2, \dots, y_T \quad 1, 2, \dots, N$$

$$p \quad \text{ar}(y, p) \quad \text{ar}$$

: DNA

.	.	.	$x(n-2)$	$x(n-1)$	$x(n)$	$x(n+1)$	$x(n+2)$	.	.	.
---	---	---	----------	----------	--------	----------	----------	---	---	---

$$: \quad x(n)$$

$$.( \quad )$$

$$. \hat{x}(n)$$

$$14071 \quad :$$

$$2500$$

AR

$$) \quad \text{BIC} \quad \text{AIC} \quad \text{FPE}$$

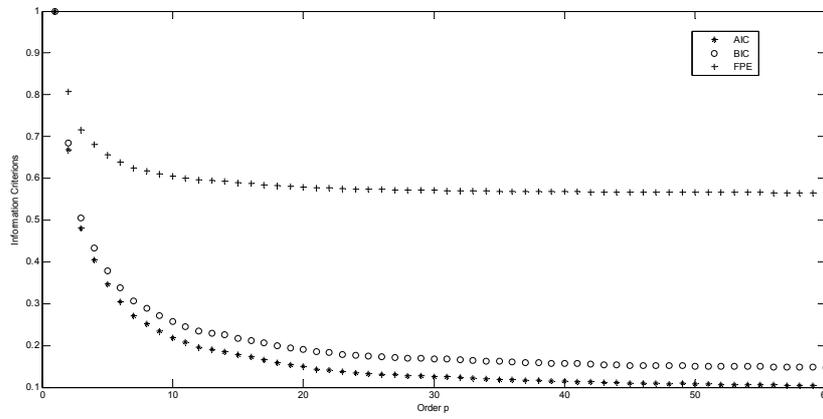
.(1

≡ Akaike Information Criterion

AIC

FPE ≡ Final Prediction Error Approach

BIC ≡ Bayesian Information Criterion



.BIC AIC FPE : (1)

AR(3)

:

$$y(n) = 0.3195 y(n-1) + 0.2805 y(n-2) + 0.3383 y(n-3) \quad (16)$$

$$.FPE = 1.35225 \quad \text{Loss function} = 1.35167$$

:

DNA

:(1)

		y(n)	
A	G	G	8
T	T	T	57
T	T	T	129
A	A	A	235
C	C	C	308
T	T	T	484
T	C	G	528
C	C	C	820
T	C	T	987
A	C	C	1001
T	A	A	1063
C	C	C	1184
T	A	C	1374
C	C	C	1490
C	G	C	1751
C	G	C	1991
A	A	A	2076
T	T	T	2121
T	A	A	2427
T	T	T	2467

:

... [10]

.%75

.%70

:( )

(states)

(transitions)

( )

.(Aazami ,2008)

(Protein

Secondary Structure rediction)

.(Thomsen ,2001 )

(MMs)

(HMMs) .(Fonzo et al 2007)

(O) (N, M)

$$\lambda = (A, B, \pi)$$

DNA

Viterbi

16571

:

( ) : (1)

X : (1)

. DNA

: (2)

$$. G \equiv 4 \quad C \equiv 3 \quad T \equiv 2 \quad A \equiv 1$$

:

: (3)

... \_\_\_\_\_ [12]

M) M A,T,C,G  
.(Mutation

M A :

M T :

M C :

M G :

$\pi \quad \lambda=(A, B, \pi) \quad :(4)$

A .N=4 1\*N

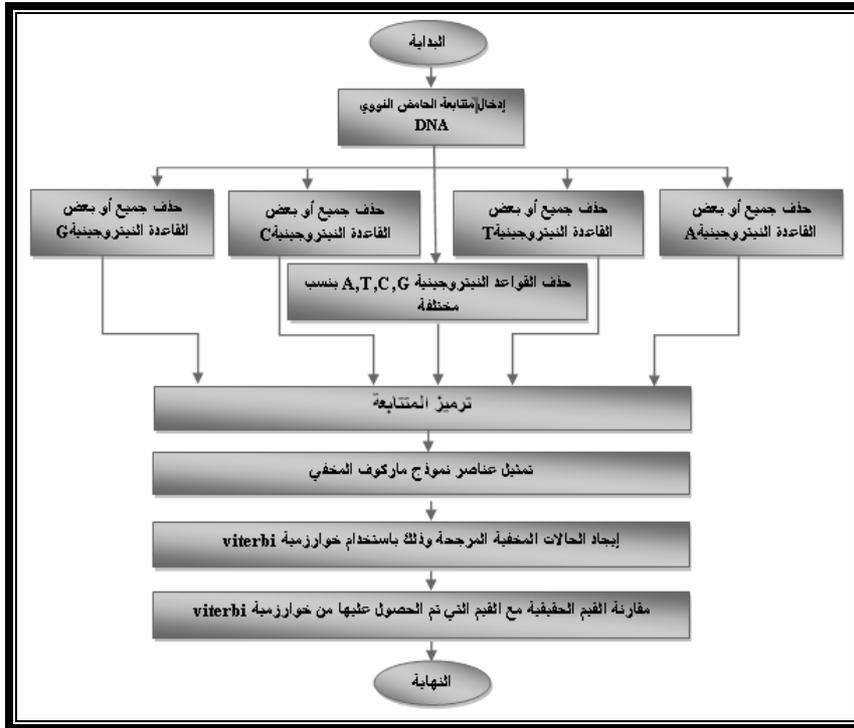
B .(N\*N)

.M=5 (N\*M)

viterbi :(5)

(5) :(6)

.(5)



: (2)

:

(1)

-1

16571

.G=2180 C=5192 T=4086 A=5113 :

M

( )

...

[14]

: Viterbi

**Transition Matrix=**

$$\begin{pmatrix} 0.3118 & 0.2392 & 0.2924 & 0.1567 \\ 0.3348 & 0.2457 & 0.2942 & 0.1253 \\ 0.2958 & 0.2770 & 0.3426 & 0.0846 \\ 0.2822 & 0.1932 & 0.3286 & 0.1960 \end{pmatrix}$$

**Emission Matrix =**

$$\begin{pmatrix} 0.9679 & 0 & 0 & 0 & 0.0321 \\ 0 & 0 & 0 & 0 & 1.0000 \\ 0 & 0 & 0.9682 & 0 & 0.0318 \\ 0 & 0 & 0 & 0.9606 & 0.0394 \end{pmatrix}$$

**mse\_seq = 0.0406**

**sum\_mutation = 4501**

**error = 415**

**Successful = 90.7798**

A, T, C, G

-2

.M

Viterbi

:(a-2)

(%)			
96.9	2	65	A
97.3	2	73	T
96.5	2	57	C
98.7	1	77	G

:(b-2)

(%)			
100%	0	5113	A
100%	0	4086	T
100%	0	5192	C
100%	0	2180	G

...

\_\_\_\_\_

[16]

:

"

"

.%90

1. الخياط ، باسل يونس ذنون (2010). "النمذجة الماركوفية مع تطبيقات عملية"، دار الكتب للطباعة والنشر، الموصل.
2. قاري، سمير بن حسن محمد و جبر، جميل فوزي جميل (2010). " مدخل إلى الوراثة البشرية" دار الفكر للطباعة والنشر , مكة المكرمة .
3. Aazami, Farshideh Einsele ,2008, " **Recognition of ultra low resolution, Anti-aliased text with small font sizes**", Unpublished Ph.d. thesis, Scientarium informaticarum, Faculty of science, University of Fribourg, Switzerland.
4. Alberts, B., Johnson, A., Lewis, J., Raff, M., Roberts, K. and Walter ,P. (2002). " **Molecular Biology of the Cell**", 4th Ed , Garland Science, New York, USA
5. Drake, J.W. (1970).” **Molecular Basis of Mutation**”, Holdeh – Day.San Frameisco.
6. Fonzo, V. , Aluffi-Pentini, F., and Parisi, V,2007, " **Hidden Markov Models in Bioinformatics**", Vol. 2, No. 1, Euro. Bio. Park, University di Roma, Roma, Italy.
7. Thomsen, R. (2001). “ **Evolving the Topology of Hidden Markov Models Using Evolutionary Algorithms** ” < thomsen @ daimi.au.dk>
8. Watson, James D and Berry, Andrew (2003). ” **DNA : the Secret of Life** “, Alfred A. Knopf, New York.

...

	<b>A</b>	<b>T</b>	<b>C</b>	<b>G</b>
<b>AAA</b>	0.2989	0.2107	0.3448	0.1456
<b>AAT</b>	0.3333	0.2312	0.3118	0.1237
<b>AAC</b>	0.3422	0.2281	0.3727	0.0570
<b>AAG</b>	0.2679	0.1914	0.3589	0.1818
<b>ATA</b>	0.3151	0.2247	0.2438	0.2164
<b>ATT</b>	0.3505	0.2810	0.2598	0.1088
<b>ATC</b>	0.3896	0.2125	0.2725	0.1253
<b>ATG</b>	0.1698	0.4151	0.2201	0.1950
<b>ACA</b>	0.1592	0.2377	0.2758	0.3274
<b>ACT</b>	0.4015	0.2044	0.2968	0.0973
<b>ACC</b>	0.2799	0.2529	0.3861	0.0811
<b>ACG</b>	0.3417	0.1750	0.3167	0.1667
<b>AGA</b>	0.3708	0.1292	0.3034	0.1966
<b>AGT</b>	0.3988	0.2209	0.2454	0.1350
<b>AGC</b>	0.3158	0.2246	0.4105	0.0491
<b>AGG</b>	0.3371	0.1657	0.3657	0.1314
<b>TAA</b>	0.3130	0.2934	0.3007	0.0929
<b>TAT</b>	0.2539	0.3313	0.2724	0.1424
<b>TAC</b>	0.2334	0.3660	0.3024	0.0981
<b>TAG</b>	0.2046	0.1815	0.3707	0.2432
<b>TTA</b>	0.3018	0.2622	0.2988	0.1372

<b>TTT</b>	0.2817	0.2500	0.3333	0.1349
<b>TTC</b>	0.2968	0.3355	0.2452	0.1226
<b>TTG</b>	0.2982	0.2544	0.2018	0.2456
<b>TCA</b>	0.3086	0.2967	0.3014	0.0933
<b>TCT</b>	0.2754	0.2656	0.3082	0.1508
<b>TCC</b>	0.2318	0.3883	0.2989	0.0810
<b>TCG</b>	0.2397	0.1570	0.3884	0.2149
<b>TGA</b>	0.2857	0.2328	0.2910	0.1905
<b>TGT</b>	0.4000	0.2400	0.2400	0.1200
<b>TGC</b>	0.3033	0.2623	0.3770	0.0574
<b>TGG</b>	0.2000	0.2100	0.4100	0.1800
<b>CAA</b>	0.3147	0.2263	0.3341	0.1250
<b>CAT</b>	0.3029	0.2476	0.3341	0.1154
<b>CAC</b>	0.3194	0.2467	0.3612	0.0727
<b>CAG</b>	0.2079	0.2475	0.3416	0.2030
<b>CTA</b>	0.2893	0.2950	0.2280	0.1877
<b>CTT</b>	0.3249	0.2114	0.3817	0.0820
<b>CTC</b>	0.3612	0.2129	0.3565	0.0694
<b>CTG</b>	0.3867	0.1934	0.2597	0.1602
<b>CCA</b>	0.2959	0.2808	0.3348	0.0886
<b>CCT</b>	0.3856	0.1845	0.3007	0.1292
<b>CCC</b>	0.2579	0.2975	0.3703	0.0744
<b>CCG</b>	0.2606	0.1972	0.3803	0.1620
<b>CGA</b>	0.3790	0.1855	0.2419	0.1935
<b>CGT</b>	0.3846	0.2692	0.2436	0.1026
<b>CGC</b>	0.3185	0.2739	0.3694	0.0382
<b>CGG</b>	0.3250	0.1625	0.3125	0.2000
<b>GAA</b>	0.1859	0.3417	0.2864	0.1859
<b>GAT</b>	0.2946	0.3125	0.2143	0.1786
<b>GAC</b>	0.2601	0.2832	0.3295	0.1272
<b>GAG</b>	0.2061	0.1985	0.3435	0.2519
<b>GTA</b>	0.2876	0.2353	0.2353	0.2418
<b>GTT</b>	0.2404	0.2788	0.3077	0.1731

...

[20]

---

<b>GTC</b>	0.2991	0.3178	0.3084	0.0748
<b>GTG</b>	0.2982	0.1579	0.3333	0.2105
<b>GCA</b>	0.2488	0.2679	0.2392	0.2440
<b>GCT</b>	0.3556	0.2889	0.2167	0.1389
<b>GCC</b>	0.2657	0.3100	0.3358	0.0886
<b>GCG</b>	0.3036	0.1786	0.3214	0.1964
<b>GGA</b>	0.2602	0.1707	0.2927	0.2764
<b>GGT</b>	0.2250	0.2875	0.3000	0.1875
<b>GGC</b>	0.2105	0.2697	0.3289	0.1908
<b>GGG</b>	0.2500	0.2361	0.3056	0.2083