

استخدام نماذج ماركوف المخفية في المعلوماتية الحيوية

اسماء حسين سمير **

* أ.د. عبدالرحيم خلف راهي

المُسْنَدُ لِـ :

انصب اهتمام هذا البحث على توظيف نماذج ماركوف المخفية (Hidden Markov models(HMMs) في المعلوماتية الحيوية Bioinformatics ، لتحديد متسلسلة الحامض النووي DNA للإنسان والمكون من سلسلة من القواعد النيتروجينية التي يطلق عليها النوكليوتيدات Nucleotides وهي الأدينين Adenine ويرمز له بالحرف a ، السيتوسين Cytosine ويرمز له بالحرف c ، الثايمين Thymine ويرمز له بالحرف t ، الكوانين Guanine ويرمز له بالحرف g ، تم التعامل مع مشكلة تحليل متسلسلة الحامض النووي DNA كمشكلة لغوية تتكون من اربعة حرف هي (a , t , c , g) والتباين بمتسسلة الحامض النووي DNA باستخدام خوارزمية فيتربي (Viterbi) وتم التوصل الى ان متسلسلة المتباين كانت مطابقة لمتسسلة الحامض النووي DNA للجين CCR5 delta-24 allele .

الكلمات الأساسية : نماذج ماركوف المخفية ، خوارزمية فيتربي ، المعلوماتية الحيوية

Using Hidden Markov Models in Bioinformatics

Abstract :

This research focused on the recruitment of Hidden Markov models (HMMs) in Bioinformatics, to identify DNA sequence of humans and consisting of a sequence of nitrogenous bases, called nucleotides, Adenine are signified by the letter a, Cytosine and signified by the letter c, Thymine and signified by the letter t, Guanine signified by the letter g, is dealing with the problem of analyzing DNA sequence DNA as a linguistic problem consists of four letters (a , c , t , g) and Ltinaya a series of DNA using an algorithm Viterbi .

Keywords: Hidden Markov models (HMMs), Viterbi Algorithm, Bioinformatics .

المبحث الأول

المعلوماتية الحيوية Bioinformatics

1-1 المقدمة [3][7] :

ان التطور الهائل في علم الاحياء ادى الى انتاج قواعد بيانات ضخمة تتعلق بتركيب الحامض النووي DNA ، تركيب الموروثات ، وظائف البروتينات . تم توظيف التقنيات الالكترونية والبرمجيات بما يتلاءم ونوع هذه البيانات للحفظ عليها والاستفادة منها ، ان وجود الانترنت ساعد ان تكون هذه البيانات متاحة للدارسين والباحثين .

* الجامعة المستنصرية / كلية الادارة والاقتصاد .

** باحثة .

تأريخ استلام البحث 2015/12/15

تأريخ قبول النشر 2016/1/19

مستل من اطروحة دكتوراة

هذا التطور ادى الى نشوء علم المعلوماتية الحيوية **Bioinformatics** وهو علم جديد يجمع بشكل اساس ما بين علم الحاسوبات وعلم الاحياء والذي يعرف على انه ادارة المعلوماتية الحيوية باستخدام التقانة الحاسوبية والمعلوماتية ، كذلك عرف المركز الوطني لمعلومات التكنولوجيا الحيوية (NCBI) National Center for Biotechnology Information علم المعلوماتية الحيوية "العلم الذي نتج عن اندماج علم الحاسوب الالي والاحياء وتكنولوجيا المعلومات لتكوين مجال علمي واحد " استخدمت المعلوماتية الحيوية على نطاق واسع في ابحاث الجينوم البشري لتحديد سلسلة الحامض النووي DNA الكاملة للإنسان ، كما ان احد مهام هذا العلم والتي تمثل التحدي الاكبر لحد هذه اللحظة هي مشكلة ايجاد الجين Gene finding في متسلسلة الحامض النووي DNA ، يعتمد علم المعلوماتية الحيوية Bioinformatics على علوم الرياضيات والحواسوب والاحصاء والطب والكيمياء للاستفادة من هذه العلوم في تحليل البيانات .

تهدف المعلومات الحيوية الى الآتي :-

1. استحداث اساليب تقنية وبناء خوارزميات وبرامج تساهمن في الحصول على المعلومات من مجموعة ضخمة من البيانات .
2. تحليل وتفسير الانواع المختلفة من البيانات التي تتضمن سلسلة الاحماض الامينية والقطع والبني البروتينية .
3. تحليل وتفسير بيانات المورثات البشرية .
4. تطوير وتنفيذ ادوات تساعد على ادارة فعالة لالاماط المختلفة في المعلومات .
5. الحصول على معلومات جديدة تخص ترميز البروتينات ووظيفة البروتينات والعديد من الوظائف الحيوية الأخرى من مجموعة البيانات الخام ، هذه المعلومات الجديدة مطلوبة لتصميم الادوية والتشخيص الطبي والعلاج الطبي فضلاً عن الميادين البحثية الأخرى .

ان المجالات البحثية للمعلوماتية الحيوية كثيرة منها على سبيل المثال :

1- تقسيي الجين Gene finding

وهو معرفة الجينات على سلسلة الحامض النووي DNA وكذلك التنبؤ بتركيب الجين اي تحديد الاكسونات Exons (وهي المناطق المشفرة Coding Sites) ومناطق اللصق (splice sites) ، والانترنونات Entrons (وهي مناطق غير مشفرة non-Coding Sites) .

2- محاذات السلسل Sequence Alignment

مقارنة السلسلة التي نحصل عليها مع سلسلة او سلسل اخرى معرفة سابقا فإذا كان هناك تشابه فهذا يعني وجود وظيفة مشتركة او مشابهة بمعنى اخر ان السلسلة الجديدة تؤدي نفس الوظيفة .

3- تركيب البروتينات Protein foldin

معرفة التركيب الثنائي والثلاثي والرباعي (ان وجد) للبروتينات لمعرفة وظيفة البروتين وذلك لوجود ارتباط وثيق بين شكل البروتين ووظيفته .

4- تحديد موقع ارتباط عوامل النسخ Transcription Factor Binding Site Identification وهي موقع صغيرة جدا على سلسلة الحامض النووي DNA ، ان اغلب البيانات الخام المستخدمة في المعلوماتية الحيوية Bioinformatics هو عبارة عن متسلسلة من النوكليوتيدات Nucleotides تمثل الهيكل الاساسي للحامض النووي DNA ، ومتسلسلة من الاحماض الامينية المناظرة للهيكل الاساسي للبروتينات .

1-2 الهدف :

توظيف نماذج ماركوف المخفية Hidden Markov models(HMMs) في علم المعلوماتية الحيوية Bioinformatics ، لتحديد متسلسلة الحامض النووي DNA .

المبحث الثاني

نماذج ماركوف المخفية (HMMs)

2-1 المقدمة :

تعد نماذج ماركوف المخفية Hidden Markov Models (HMMs) تعينا لنماذج ماركوف الاعتيادية وتطويرا لعمليات ماركوف ، ففي نماذج ماركوف الاعتيادية Markov Models تكون الحالات مرئية بشكل مباشر لذا تكون احتمالات انتقال من حالة الى حالة اخرى معلوم بينما في نماذج ماركوف المخفية (HMMs) تكون الحالات غير مرئية .

ان نماذج ماركوف المخفية (HMMs) هي عبارة عن مجموعة منتهية من الحالات كل حالة تقترن بتوزيع احتمالي ، بشكل عام تولد الحالة الناتجة طبقاً للاحتمالات المقترنة بالحالة حيث توجد احتمالات ناتجة فقط ولا توجد حالة ظاهرة يمكن مشاهدتها ، لذا تكون الحالات مخفية وهذا هو مفهوم نماذج ماركوف المخفية (HMMs) بشكل عام .

تعد نماذج ماركوف المخفية (HMMs) اداة احصائية قوية للتنبؤ بسلسلة الحالة من خلال سلسلة المشاهدات ، تم الاستفادة منها وتطبيقاتها في علم المعلوماتية الحيوية Bioinformatics التي تهتم بقواعد البيانات الحيوية والوراثية وادارتها وتطويرها .

2-2 معلمات نماذج ماركوف المخفية (HMMs)

[5] [4] وان معلمات نماذج ماركوف المخفية (HMMs) هي كالتالي :-

1- مصفوفة الاحتمالات الانتقالية ويرمز لها $D(nxn)$ يرمز للعنصر في الصف i والعمود j بالرمز d_{ij} وهو يمثل احتمال الانتقال من الحالة i الى الحالة j ويجب ان يحقق الشرط الاتي :-

$$\sum_{i=0}^n d_{ij} = 1 \quad (1)$$

Where $i=1 \dots n, j=1 \dots n$

بالنسبة لمتسلسلة الجين التي سيتم دراستها تتكون من اربع حالات هي (a , t , c , g) اي ان $n=4$

2- مصفوفة العلاقات (الاحتمالات المقترنة بالحالة) ويرمز لها $Z(nxm)$ يرمز للعنصر في الصف i والعمود k بالرمز Z_{ik} وهو احتمال رمز المشاهدة للمؤشر k المنبعث من الحالة i ويجب ان يحقق الشرط الاتي

$$\sum_{k=0}^m Z_{ik} = 1 \quad (2)$$

Where $i=1 \dots n, k=1 \dots m$

في نموذج البحث هناك اربعة احتمالات مقترنة في كل حالة اي ان $m=4$

3- متوجه التوزيع الابتدائي ويرمز له بالرمز π ويجب ان يحقق الشرط الاتي :-

$$\sum_{i=0}^n \pi_{ij} = 1 \quad (3)$$

$\lambda = (D, Z, \pi)$ ويعبر عن الانموذج

2-3 المسائل الاساسية لنماذج ماركوف المخفية

1- مسألة التقويم (Evolution Problem) وهي عملية ايجاد $P(O|\lambda)$ اي ايجاد الامكان الاعظم لمتسلسلة المشاهدات $O=(o_1, o_2, \dots, o_T)$ عندما يكون الانموذج $\lambda = (D, Z, \pi)$

2- مسألة حل الشفرة (Decoding problem) وهي عملية ايجاد متابعة الحالة المثلثى للانموذج $\lambda = (D, Z, \pi)$ ومتسلسلة المشاهدات $O=(o_1, o_2, \dots, o_T)$

3- مسألة التدريب (Training Problem) وهي عملية تعديل معلمات الانموذج $P(O|\lambda) = \lambda$ لتعظيم ويتم حل هذه المسائل من خلال الخوارزميات الآتية

- الخوارزمية الامامية الخلفية (Forward-Backward Algorithm)
- خوارزمية فيتربي (Viterbi Algorithm)
- خوارزمية بوم ولتش (Beaum-Welch)

عند توظيف نماذج ماركوف المخفية في المعلوماتية الحيوية فان الزمن Time للمتسلسلات الحيوية يعني موقع النوكليوتيدات على طول متسلسلة الحامض النووي DNA .

1-3-2 الخوارزمية الإمامية الخلفية (Forward-Backward Algorithm)

اولاً : الخوارزمية الإمامية Forward-Algorithm

تعرف الاحتمالية الإمامية $\alpha_t(i)$ او ما يدعى α -Pass على انها الاحتمالية لمتسلسلة المشاهدات الجزئية والحالة $(q_t=S_i)$ عندما يكون المعطى هو النموذج (λ) والتي عبر عنها بالصيغة الآتية :-

$$\alpha_t(i) = P(O_1, O_2, \dots, O_t, q_t = S_i | \lambda)$$

ان هذه الدالة الاحتمالية يمكن ان تحل لـ (N) من الحالات و (T) من المشاهدات بشكل تكراري كالتالي :-

1- البداية (Initialization)

$$\alpha_1(i) = \prod_{j=1}^N b_j(O_j) \quad i=1, 2, \dots, N$$

2- التعاقب (Induction)

$$\alpha_t(j) = \left[\sum_{i=1}^N \alpha_{t-1}(i) a_{ij} \right] z_{ij}(O_t) \quad j=1, 2, \dots, N \quad t=2, 3, \dots, T$$

3- النهاية (Termination)

$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i)$$

ثانياً : الخوارزمية الخلفية Backward Algorithm

يعرف المتغير الخلفي (β_t) او ما يدعى β -Pass على ان الاحتمالية لمتسلسلة المشاهدات الجزئية $O_{T+1}, O_{T+2}, \dots, O_T$ عندما يكون المعطى هو الحالة (S_i) عند الزمن (t) والنماذج (λ) والذي تم التعبير عنه بالصيغة الآتية :-

$$\beta_t(i) = P(O_{T+1}, O_{T+2}, \dots, O_T | q_t = S_i)$$

ان هذا الاجراء مشابه للإجراء الإمامي لكن يكون سريان الحالة بالاتجاه الخلفي فيأخذ المشاهدة الاخيرة في اللحظة الزمنية (T) لحين الوصول الى اللحظة الاولى وكالتالي :-

1 ← ----- ← T-2 ← T-1 ← T

ان هذه الدالة الاحتمالية يمكن ان تحل لـ (N) من الحالات و (T) من المشاهدات بشكل تكراري كالتالي :-

1- البداية (Initialization)

$$\beta_T(i) = 1 \quad i=1, 2, \dots, N$$

2- التعاقب (Induction)

$$\beta_t(i) = \sum_{j=1}^N d_{ij} z_j(O_{t+1}) \beta_{t+1}(j) \quad j=1, 2, \dots, N \quad t=T-1, T-2, \dots, 1$$

3- النهاية (Termination)

$$P(O|\lambda) = \sum_{i=1}^N \alpha_t(i) \beta_t(i)$$

2-3-2 خوارزمية فيتربi Vfeterb Algorithm

هي طريقة البرمجة الديناميكية Dynamic Programming وتعمل على ايجاد مسار الحالة الانتقالية الاكثر احتمالاً عندما يكون المعطى النموذج (λ) وعدد الحالات هو N ومتتابعة المشاهدات $O = (O_1, O_2, \dots, O_T)$ ذات الطول T وان المتغير الاساسي لهذه الخوارزمية هو المتغير (i) حيث يمثل اعلى احتمالية على طول المسار الوحد في الحالة (i) عند الزمن (t) والذي يساوي متسلسلة الحالة الجزئية الاكثر احتمالاً بالنسبة لمتسلسلة المشاهدات المنتهية في الحالة (i) ، وهذه الخوارزمية مشابهة للخوارزمية الإمامية الا انها تأخذ اعلى Max احتمالية للمسارات السابقة بينما الخوارزمية الإمامية تأخذ مجموع Sum احتمالية المسارات السابقة ، كذلك خوارزمية فيتربi لها مؤشرات ترجعية (Back-Pointers) لاتملکها الخوارزمية الإمامية ، ويتم ايجاد التسلسل الامثل للحالة وذلك بالاحتفاظ بمسار الحالات المخفية التي تقود لكل حالة ، ثم تتبع افضل مسار (Back-Trace) الى البداية .
اما خطوات خوارزمية فيتربi فهي كالتالي :

1- البداية Initialization

$$\delta_1(j) = \pi_i z_j(O_1), \quad j = 1, 2, \dots, N \quad \Psi_1(j) = 0$$

حيث ان $\Psi_t(j)$ هو متغير حفظ تتبع الآخر (keep track).

2- التعاقب Induction

$$\delta_t(j) = \max_{i=1}^N [\delta_{t-1}(i) d_{ij}] z_j(O_t) \quad j = 1, 2, \dots, N \quad t = 2, 3, \dots, T$$

$$\Psi_t(j) = \arg \max_{i=1}^N [\delta_{t-1}(i) d_{ij}] \quad j = 1, 2, \dots, N \quad t = 2, 3, \dots, T$$

3- النهاية Termination

$$P^* = \max_{i=1}^N [\delta_T(i)]$$

$$q^* = \arg \max_{i=1}^N [\delta_t(i)]$$

4- التعاقب المعاكس Back Tracking

$$q^*_{-t} = \Psi_{t+1}(q^*_{t+1}) \quad t = T-1, T-2, \dots, 2, 1$$

ثم تقرأ المتسلسلة الأفضل للحالات من متجهات (Ψ_t) ولهذا فان خوارزمية فيتربي تؤدي الى نتيجتين مفيدتين هما :-

1- يكون الاختيار (من بين كل المسارات) هو المسار الأفضل $Q^* = [q_1^*, q_2^*, \dots, q_T^*]$ والذي يطابق متسلسلة الحالة التي تعطي الامكان الاعظم لمتسلسلة المشاهدات .

2- الامكان على طول المسار الأفضل يكون $P(O, Q^*/\lambda) = P^*(O/\lambda)$.

بالمقارنة بالخوارزمية الامامية فان خوارزمية فيتربي تقوم بحساب الامكان على طول المسار الأفضل .

3-3-2 خوارزمية بوم ولتش (Beaum-Welch)

تستخدم هذه الخوارزمية لعادة تقدير معلمات النموذج (D, Z, Γ) من اجل تعظيم الاحتمالية لمتسلسلة المشاهدات وذلك بالاعتماد على المتغيرين الآتيين :-

المتغير الاول :-

هو (i) γ_t وهو يمثل كون العملية في الحالة (S_i) عند الزمن (t) عندما يكون المعطى هو المشاهدات والنماذج ، ويعبر عنه كالتالي

$$\gamma_t(i) = \frac{P(q_t = s_i, O/\lambda)}{P(O/\lambda)}$$

$$\gamma_t(i) = \frac{\alpha_t(i) \beta_t(i)}{\sum_{i=1}^N \alpha_t(i) \beta_t(i)}$$

المتغير الثاني :-

هو ((i, j)) ξ وهو يمثل كون العملية في الحالة (S_i) عند الزمن (t) والحلة (S_j) عند الزمن ($t+1$) ، ويعبر عنه كالتالي :-

$$\xi_t(i, j) = \frac{P(q_t = s_i, q_{t+1} = s_j, O/\lambda)}{P(O/\lambda)}$$

$$\xi_t(i, j) = \frac{\alpha_t(i) d_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \alpha_t(i) d_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}$$

4-2 استخدام (HMMs) في المعلوماتية الحيوية [2][3][7]

منذ أن تم اعتماد نظرية ماركوف المخفية بصورة رسمية في أواخر السنتينيات القرن الماضي ، ان الكثير من العلماء طبق النظرية في المشاكل اللغوية والكلام لكونها مناسبة للمشكلة ، حيث وجود كميات كبيرة من البيانات والمعرفة المستمدبة من الملاحظة وهذا هو الحال ايضاً في الجانب الحيوي .

استخدمت المعلوماتية الحيوية على نطاق واسع في ابحاث الجينوم البشري لتحديد سلسلة الحامض النووي DNA للانسان والذي هو عبارة عن سلسلة من القواعد النيتروجينية التي يطلق عليها التكليوتيدات وهي الادنین Nucleotides ويرمز له بالحرف a ، السيتوسين Cytosine ويرمز له بالحرف t ، الكوانين Guanine ويرمز له بالحرف c ، الثايمين Thymine ويرمز له بالحرف t ، الكوانين Guanine ويرمز له بالحرف g . اثنان

من اهم مشاكل تحليل التسلسل الاكثر دراسة هي :-

- 1- تمييز الكلام Speech Recognition
- 2- معالجة اللغة Language Processing

ان سلسلة الحامض النووي DNA هي مماثلة لمشكلة تمييز الكلام Speech Recognition ومعالجة اللغة Language Processing وهي تمثل أنموذج الآيسير الایمن ، لذلك يمكن التعامل مع مشكلة تحليل متسلسلة الحامض النووي DNA كمشكلة لغوية تتكون من اربعة احرف هي (a , t , c . g) ، فمثلاً مشكلة ايجاد الجين Gene Finding في الحامض النووي DNA كتحليل اللغة الى كلمات وجمل ذات معنى وعليه يمكن ان نتعامل بنفس التقنيات المستخدمة للتعرف على الكلم ومعالجة اللغة . مثلاً استخدمت نظرية HMM كنماذج رياضية للغة يمكن استخدامها كذلك كنماذج رياضية لتحليل متسلسلة الحامض النووي .

ان النهج القائم على اساس (HMMs) هو نهج مناسب للمشاكل الحيوية الآتية :-

- 1- محاذات السلسل Sequence Alignment وهي مشكلة متكررة في المعلوماتية الحيوية .

2- رسم الخرائط الوراثية Genetic Mapping وهو من اول تطبيقات (HMMs) في المعلوماتية الحيوية وهو تقدير بعض انواع المسافات الفاصلة بين الواقع الوراثي على طول الكروموسوم .

3- التنبؤ بالتركيب الثنائي للبروتين Secondary Structure Protein Prediction تعمل (HMMs) على التنبؤ بالهيكل الثنائي للبروتين وهي خطوة مهمة للتنبؤ بهيكل ثلاثي الابعاد.

4- ايجاد الجين Gene finding وهو كما بين سابقاً معرفة الجينات على سلسلة الحامض النووي DNA وكذلك التنبؤ بتركيب الجين اي تحديد الاكسونات Exons والانترونات .

المبحث الثالث

الجانب النظري والنجزي

3-1 المقدمة [8]

ان المشاهدات التي استخدمت في البحث تم الحصول عليها من موقع المركز الوطني لمعلومات التكنولوجيا الحيوية National Center for Biotechnology Information (NCBI) وهي عبارة عن متسلسلة من القواعد النيتروجينية في متسلسلة الحامض النووي DNA تمثل الجين CCR5 delta-24 allele الذي يؤثر في مقاومة الجسم لفيروس الايدز HIV وكما مبين في الشكل رقم (1) والشكل رقم (2)

The screenshot shows the NCBI Nucleotide search results for the **Homo sapiens CCR5 (CCR5) gene, CCR5 delta-24 allele, partial cds**. The GenBank ID is GQ121035.1, and the sequence ID is 289466094. The page displays detailed information about the gene, including its definition, accession number, version, keywords, source, organism, reference, authors, title, journal, and submission details. On the right side, there are sections for 'Analyze this sequence' (with options like Run BLAST, Pick Primers, Highlight Sequence Features, and Find in this Sequence), 'Articles about the CCR5 gene' (listing several academic papers), and 'Reference sequence'.

(1) الشكل

This screenshot shows the same search results as above, but it highlights the 'ORIGIN' section, which displays the DNA sequence of the CCR5 gene. The sequence starts with: 1 cggagccctg caaaaaaatc aatgtgaagg aaatcgacg ccgcctctg cttccgcctt... and continues for several lines. To the right of the sequence, there is additional information about the gene's function, homologs, external resources, and related information.

(2) الشكل

ان هذه المتسلسلة هي عبارة عن متسلسلة من الاحرف (a , t , c . g) والتي يبلغ عددها (500) قاعدة نيتروجينية ، هذه القواعد النتروجينية تم تحويل الرموز الحرفية لها الى ارقام كالتالي :-

$$A \equiv 1 , \quad t \equiv 2 , \quad c \equiv 3 , \quad g \equiv 4$$

ان القواعد النتروجينية في متسلسلة الحامض النووي DNA تعتمد على مواقعها في المتسلسلة المطلوب التنبؤ بمتسلسلة الحالات من خلال متسلسلة الرموز الرقمية ومقارنتها بمتسلسلة الحالة الحقيقة. عند استخدام نماذج ماركوف المخفية (HMMs) في بيانات المعلومات الحيوية يجب ان نحدد الأنماذج المستخدم وعند اختيار او تحديد الأنماذج يجب الاخذ بنظر الاعتبار الآتي :-

1- عدد الحالات والانتقالات

ان النموذج المستخدم يتكون من عدد محدد من الحالات ينبع من كل حالة اربعة رموز بتوزيع احتمالي معين لكل رمز ، للأنماذج مجموعة من الانتقالات بين الحالات والتي تسمح بتغير الحالة بعد انباع الرموز .

$$N = (a , t , c , g) \quad \text{الحالات (States) هي :-}$$

الاحتمالات الانتقالية (Transition Probability) هي :-

$$A = \begin{bmatrix} 0.320388 & 0.271845 & 0.23301 & 0.174757 \\ 0.092105 & 0.256579 & 0.328947 & 0.322368 \\ 0.270073 & 0.423358 & 0.240876 & 0.065693 \\ 0.185185 & 0.250000 & 0.277778 & 0.287037 \end{bmatrix}$$

الاحتمالات المنبعثة (Emission Probability) هي :-

$$B = \begin{bmatrix} 0.97 & 0.01 & 0.01 & 0.01 \\ 0.01 & 0.97 & 0.01 & 0.01 \\ 0.01 & 0.01 & 0.97 & 0.01 \\ 0.01 & 0.01 & 0.01 & 0.97 \end{bmatrix}$$

متسلسلة المشاهدات هي :-

$$O = (3441433324 \ 3311111123 \ 1124241143 \ 1112343143 \ 3343323324 \ 3323343232 \\ 1323132442 \ 4223123222 \ 4422224244 \ 4311312432 \ 4423123323 \ 1233241211 \\ 1324311114 \ 4324114143 \ 1241324131 \ 2321332432 \ 3113324433 \ 1232324133 \\ 2422222332 \ 2322132423 \ 3332232444 \ 3231321243 \ 2433433314 \ 2444132224 \\ 4111213112 \ 4242311323 \ 2241314443 \ 2321222212 \ 1443223223 \ 2324411232 \\ 2322312312 \ 3323324131 \ 1234121442 \ 1332443242 \ 3423312432 \ 4242224322 \\ 2111143314 \ 4134423133 \ 2224444244 \ 2413114242 \ 4123132244 \ 4244244324 \\ 2422243423 \ 2323331441 \ 1231232221 \ 3314123231 \ 1111411442 \ 3223122131 \\ 3324314323 \ 2312222331)$$

2- الزمن Time

الزمن في متسلسلة الحامض النووي DNA يعني موقع المشاهدات على طول المتسلسلة. (2) وبتطبيق خوارزمية فيتربي Viterbi على المعلومات (متسلسلة المشاهدات، الاحتمالات الانتقالية، الاحتمالات المنبعثة) باستخدام برنامج MATLAB تم التنبؤ بمتسلسلة الحالة الآتية :-

```
cggagccctg ccaaaaaatc aatgtgaagc aaatcgacgc ccgcctcctg cctccgcct
actcaactggt gttcatctt gggttgtgg gcaacatgc ggtcatcctc atcctgataa
actgaaaag gctgaagagc atgactgaca tctacctgct caacctggcc atctctgacc
tgttttcct tcttactgtc cccttctggg ctcactatgc tgccgccccag tgggactttg
gaaatacaat gtgtcaactc ttgacagggc tctattttat aggcttcttc tctggaaatct
tcttcatcat cctccgtaca atcgataggt acctggctgt cgccatgct gtgttgctt
taaaagccag gacggtcacc ttgggggtgg tgacaagtgt gatcactgg gtgggtgg
tgtttgcgtc tctcccagga atcatctta ccagatctca aaaagaaggt cttcattaca
cctgcagctc tcattttcca
```

ان المتسلسلة المتنبئ بها هي مطابقة لمتسلسلة الحامض النووي للجين CCR5 delta-24 allele

3-2 الاستنتاجات

- 1- ان نماذج ماركوف المخفية (HMMs) نماذج احصائية قوية للتتبؤ بمتسلسلة الحامض النووي في المعلوماتية الحيوية DNA . Bioinformatics
- 2- من خلال البحث تم التوصل الى ان المتسلسلة المتنبأ بها هي مطابقة لمتسلسلة الحامض النووي للجين CCR5 delta-24 allele.
- 3- كلما كانت معلمات النموذج مثلى او قريبة من الامثلية تكون متسلسلة الحالة المتنبأ بها صحيحة.

3-3 التوصيات

- 1- نوصي باختيار النموذج المناسب للبيانات وذلك من خلال تحديد الغرض من استخدام (HMMS) ، حيث هناك عدة نماذج لـ (HMMs) في المعلوماتية الحيوية Bioinformatics .
- 2- نوصي باهتمام بحثي اكبر بالتعاون مع الاختصاصات ذات العلاقة في علم المعلوماتية الحيوية Bioinformatics .

المصادر

- 1- Birney,E.(2001)."Hidden Markov Models in Biological Sequence Analysis" ,IBM J.RES & DEV , Vol. 45 , No.3/4 ,P.449.
 - 2- Valeria De Fonzo, Filippo Aluffi-Pentini and Valerio Parisi (2007) "Hidden Markov Models in Bioinformatics"Current Bioinformatics,Vol. 2,P. 49-61.
 - 3- Smith , Kaleigh(2002) " Hidden Markov Models in Bioinformatics with Application to Gene Finding in Human DNA " January 17, 2002.
 - 4- Rabiner.L.R and Juang.B.H. (1986)." An Introduction to Hidden Markov Models " IEEE ASSP MAGAZINE JANUARY 1986, 0740-7467/86/0100-0004\$01.00@198I6E EE,P.4-16.
 - 5- Stamp , M. (2012),"A Revealing Introduction to Hidden Markov Models", Associate Professor, Department of Computer Science, San Jose State University, Email: stamp@cs.sjsu.edu
 - 6- Sean R Eddy(2004)" What is a hidden Markov model?" Nature Biotechnology Volume 22 Number 10 October 2004 .
 - 7- Xiong , Jin (2006), " Essential Bioinformatics " Cambridge university press.
 - 8- www.ncbi.nlm.nih.gov
-
.....
.....