

## Using Genetic Algorithm to Estimate (RNA) Estimator

\*Ban Ahmed Mitras

\*\*Farah Saad Nashat

[\\*dr.banah.mitras@gmail.com](mailto:dr.banah.mitras@gmail.com)

College of Computer Sciences and Mathematics  
University of Mosul

\*\*College of Education. University of Dohuk

Received on: 16/08/2010

Accepted on: 10/11/2010

### ABSTRACT

In this paper the genetic algorithm has been used to estimate the parameter  $\theta$  which exist in Boltzmann Distribution which controls the structure of the Ribo Nucleic Acid (RNA). Two algorithms have been suggested. The first found the value of the estimator which maximizes the likelihood function of Boltzmann Distribution. The second minimized the generation constraint of Boltzmann Distribution by using the genetic algorithm. Matlab (7.0) has been used in writing the programs of algorithms and achieved the following results: The maximum value for the likelihood estimator for Boltzmann Distribution appear at the value -4.1614 where the value of  $\theta$  is 0.1457, and the minimum value for the Constraint Generation for Boltzmann Distribution appear at the value  $0.951039101 \times 10^7$  where the value of  $\theta$  is -4.4066.

**Keyword:** Genetic Algorithm, Boltzmann Distribution, RNA Estimator.

إستخدام الخوارزمية الجينية لتقدير معلمة الحامض النووي الرايبوي (RNA)

فرح سعد نشاط

كلية التربية

جامعة دهوك

بان احمد حسن متراس

كلية علوم الحاسوب والرياضيات

جامعة الموصل

تاريخ قبول البحث: 16/08/2010

تاريخ استلام البحث: 10/11/2010

### الملخص

تم في هذا البحث استخدام الخوارزمية الجينية لتقدير المعلمة  $\theta$  الموجودة في توزيع بولتزمان الذي يخضع له تركيب الحامض النووي الرايبوي (RNA) إذ تم اقتراح خوارزمتان أولهما تقوم بإيجاد قيمة المقدر الذي يعظم دالة الترجيح لتوزيع بولتزمان، والثانية تقوم بتقليل دالة قيد الجيل لتوزيع بولتزمان باستخدام الخوارزمية الجينية، ولقد استخدم نظام Matlab (7.0) في كتابة برامج الخوارزميات والتي حصلنا من خلالها على النتائج التالية: أعظم قيمة لمقدر الترجيح لتوزيع بولتزمان ظهرت عند القيمة -94.1614 حيث كانت قيمة  $\theta$  هي 0.1457، وأقل قيمة لدالة قيد الجيل لتوزيع بولتزمان ظهرت عند القيمة  $0.951039101 \times 10^7$  حيث كانت قيمة  $\theta$  هي -4.4066.

**الكلمات المفتاحية:** الخوارزمية الجينية، توزيع بولتزمان، مقدر RNA

## 1. مقدمة: Introduction

شهد القرن العشرين تقدماً هائلاً في الأساليب العلمية المستخدمة في البحث العلمي في حقول المعرفة كافة ويشمل هذا التطور علم الإحصاء الذي أصبح علماً مستقلاً له أهميته بوصفه وسيلة للبحث العلمي وأداته، ونظراً لأهميته واستخدامه في العديد من مجالات الحياة، فقد اهتم العديد من الباحثين به وقاموا بدراسته وادخاله في بعض المجالات الحياتية ومنها دراسة بعض الأحماض النووية التي تدخل في التركيب الداخلي لجسم الإنسان ومن أهمها الـ RNA والـ DNA، إذ أن سلاسل هذه الأحماض تعتبر متغيرات عشوائية تخضع لتوزيعات احتمالية [6] [7].

والمقصود بتوزيع المتغير العشوائي (أو التوزيع الإحتمالي لمتغير عشوائي) بأنه تعريف السلوك الرياضي لإحتمالات قيم المتغير. أي أن التوزيع الإحتمالي للمتغير العشوائي هو القانون الإحتمالي لذلك المتغير [2]. أن لنظرية التقدير أهمية كبيرة في تطبيقات النظرية الإحصائية في الجوانب العلمية، إذ أنها توفر قواعد يتم بموجبها تقدير معلمات مجهولة لظاهرة معينة في الواقع العملي ولمختلف المجالات. ويمكن تجزئتها إلى جزئين متكاملين، الأول يهتم بالبحث عن أفضل تقدير لمعلمة مجهولة في المجتمع ويسمى بالتقدير النقطي (Point Estimation)، ويهتم الجزء الثاني بالبحث عن أفضل فترة يمكن حصر قيمة المعلمة المجهولة خلالها وهذا ما يسمى بتقدير الفترة (Interval Estimation) [5]. لقد كان هناك اهتمام كبير في السنوات الأخيرة على طريقة الترجيح الأعظم لكونها أكفأ وأفضل طريقة لتقدير المعلمات، وفي عام (1996) قدم الباحثان Lehr و Lii الخوارزمية الأساسية للترجيح الأعظم (Maximum Likelihood Estimation) التي تعطي تقديرات ثابتة وكفاءة للمعلمات، إن تقدير العينات الكبيرة يمكن أن تكون مكلفة أو حتى مستحيلة ولهذا طور كل من Xu و Vogl في عام (2000) طريقة الترجيح الأعظم وتمكنا من تقدير معلمات هذه العينات.

وهناك عدة بحوث تختص بتقدير المعلمات باستخدام توزيع بولتزمان الذي يمثل تركيب سلسلة الـ RNA وفيما يلي بعض من هذه الأعمال و البحوث:

وصف Pond وآخرون في عام (2006) طريقة أساسها دالة الإمكان لإختيار النموذج التطوري باستخدام الخوارزمية الجينية (GA) ليستكشف بسرعة مجموعة كبيرة متحدة لكل نماذج ماركوف القابلة للانعكاس مع عدد ثابت لإحلال النسب. وحققوا أيضاً فائدة لعدة مقاييس مسافة للمقارنة والتغاير للنماذج التطورية المستتجة [10].

وفي العام نفسه وصف Ding وآخرون إجراء بديل لتمثيل تركيب الحامض الرايبوي المرسل mRNA، الذي اختبرت فيه التراكيب من مجموعة بولتزمان (Boltzmann) المرجحة لتراكيب سلسلة الـ RNA الثانوية المتجمعة، بالاعتماد على عينة عشوائية من طول طبيعي للإنسان وأخذ الـ mRNA منه، واطهروا خصائص المستوى العنقودي بأنها قابلة للإنتاج بشكل إحصائي في مقارنة بين هيكل mRNA و RNA [8].

كما قدم Andronescu وآخرون في العام نفسه طريقة قيد الجيل (CG)، وهي الطريقة الحسابية الأولى لتقدير معلمة الطاقة الحرة لسلسلة الـ RNA التي يجب أن تدرب بشكل كفاءة على هيكلية المجموعات الكبيرة فضلاً عن بيانات الديناميكية الحرارية. وبينوا أن طريقة قيد الجيل (Constraint Generation) توظف مخططاً تكرارياً مبتكراً، إذ تم حساب الطاقة أولاً كحل لمسألة الأمثلية، ثم استخدموا معلمات الطاقة المحسوبة لتجديد القيود لدالة الأمثلية، لكي يحسنوا أمثلية معلمات الطاقة في التكرار القادم [3].

وقد استعمل الباحثان Chan و Ding في عام (2008) أسلوب العينات وجمعوا عناقيد لمجموعة بولتزمان (Boltzmann) لتكوين RNA الثانوي، للتحقق في هيئة المجموعة المعروضة للسلاسل الحيوية التي ميزت من تعديلاتهم العشوائية، وقاموا بتحليل أساس تركيب الـ MicroRNA المتمثل بـ (miRNA) لـ 9 مجموعات مميزة [4].

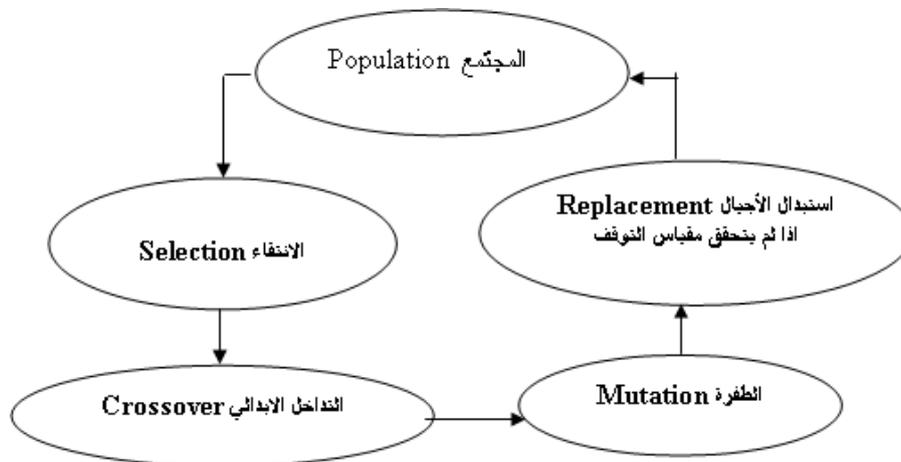
في هذا البحث تم تقدير معلمة توزيع بولتزمان المتمثل بمعلمة RNA باستخدام الخوارزمية الجينية وذلك بـ

:

1. تعظيم دالة الترجيح لتوزيع بولتزمان باستخدام الخوارزمية الجينية.
2. تقليل قيد الحيل لتوزيع بولتزمان باستخدام الخوارزمية الجينية.

## 2. الخوارزمية الجينية: Genetic Algorithm

تتضمن الخوارزمية الجينية عدداً من الخطوات الأساسية، هذه الخطوات ثابتة لمختلف المسائل ولكل التطبيقات ويكون الاختلاف في صياغة كل خطوة من الخطوات وتطبيقها حسب المسألة أو مجال تطبيقها، إن خطوات هذه الخوارزمية تكون مترابطة بعضها مع البعض الآخر، ولا يمكن تطبيق هذه الخوارزمية على أية مسألة ما لم تطبق جميع هذه الخطوات وإلا تفقد الخوارزمية قيمتها وفائدتها في إيجاد أو تحسين الحل. [13] [1] [12] والشكل (1) يوضح المخطط العام للخوارزمية الجينية:



الشكل (1) يمثل خطوات الخوارزمية الجينية

### خطوات عمل الخوارزمية الجينية: Working Steps in Genetic Algorithm

- 1- البداية Start: توليد مجتمع عشوائي من الكروموسومات، أي بمعنى إيجاد حلول مناسبة للمسألة.
- 2- اللياقة Fitness: هي تحويل دالة الهدف (Objective Function) إلى دالة مناسبة للحل في الخوارزمية الجينية.
- 3- مجتمع جديد New Population: توليد جيل جديد بتكرار الخطوات الآتية إلى أن يكتمل الجيل، وتتضمن:

- الاختيار Selection: يتم اختيار اثنين من الكروموسومات والدين (Parents Chromosomes) من المجتمع الابتدائي استناداً إلى دالة اللياقة (أفضل القيم التي لها فرص أكبر للاختيار).
  - التداخل الابدالي Crossover: إجراء إحدى عمليات التعابر للحصول على الذرية (Offspring) ويكون بين كروموسومين.
  - الطفرة Mutation: مع احتمالية وجود الطفرة يتم عمل الطفرة للسلف الجديد بموقع معين في الكروموسوم، وتجري بين الجينات في الكروموسوم الواحد.
  - الاستبدال (Replacement): عملية وضع السلف الجديد المتكون في الجيل الجديد للحلول محل المجتمع الابتدائي.
  - الاختبار Test: عند تحقق شرط التوقف، فإن الخوارزمية الجينية تتوقف وتعيد الحل الجيد من آخر جيل متكون.
  - الدورة Loop: يتم الرجوع إلى الخطوة 2.
- سإن كل تكرار لهذه العملية يسمى بالجيل (Generation)، وبعد نهاية التنفيذ يقوم الباحث بتقديم تقرير عن الحقائق التي تم التوصل إليها.

### 3. الأحماض النووية في الخلية: Nucleic of the Cell

توجد الأحماض النووية في الخلية، إذ أن الخلية هي وحدة البناء والوظيفة في الكائن الحي وهي على نوعين: الحامض النووي الرايبوسيدي (Ribo Nucleic Acid) ويرمز له اختصاراً RNA والحامض النووي الرايبوسيدي منقوص الأوكسجين (Deoxyribo Nucleic Acid) ويرمز له DNA [11].

إن موضوع دراستنا يتناول الحامض النووي الرايبوسومي (RNA) وهو يوجد في كل من النواة والساييتوبلازم، حيث يوجد في النوية وفي الرايبوسومات وفي تراكيب أخرى.

ويمكن تعريف (الساييتوبلازم) بأنه جزء المادة الحية للخلية الذي يقع خارج النواة. وتعرف (النواة) على أنها أكبر عضوية متميزة داخل الخلية وللنواة أهمية كبيرة في نقل الصفات الوراثية وفي النشاط الأيضي للخلية [9].

تتكون جزيئة الحامض النووي الرايبوسيدي (RNA) من شريط واحد، إذ أن النيوكليوتيدات الداخلة في تركيب الـ RNA تتألف من ثلاث جزيئات بسيطة مرتبطة بعضها ببعض مباشرة وهي:

1- قاعدة نيتروجينية: وهي مركب حلقي يحتوي على النيتروجين فضلاً عن الكربون والهيدروجين والأوكسجين عدا الأدينين الذي لا يحتوي بدوره على الأوكسجين، ويوجد نوعان منها:

أ- البريميدينات (Pyrimidines): وتتكون من حلقة واحدة وتشمل القواعد الآتية:

1- الساييتوسين (Cytosine) ويرمز لها (C).

2- اليوراسيل (Uracil) ويرمز لها (U).

ب- البيورينات (Purines): وتتكون من حلقتين وتشمل:

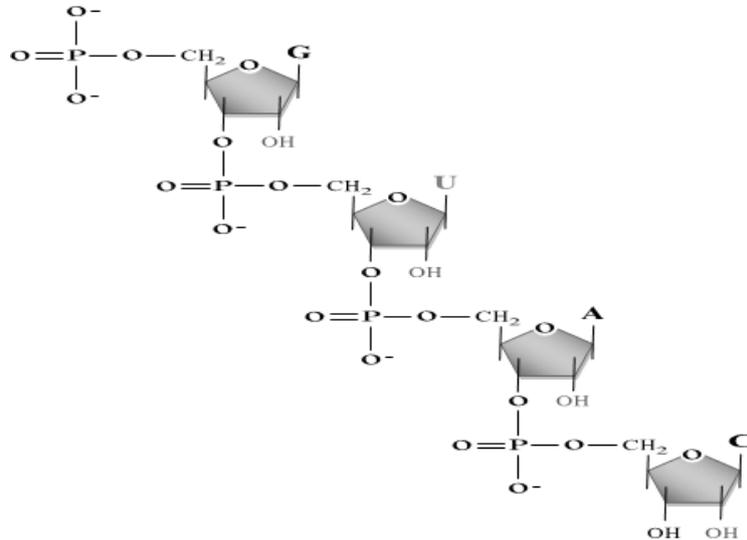
1- الأدينين (Adenine) ويرمز لها (A).

2- الكوانين (Guanine) ويرمز لها (G).

فالأدينين يرتبط دائماً مع اليوراسيل بأصرتين هيدروجينيتين  $A=U$  والساييتوسين يرتبط دائماً مع الكوانين بثلاث أوأصر هيدروجينية  $C \equiv G$ .

- 2- سكر خماسي الكربون وهو الرايبوز وصيغته الجزيئية هي  $(C_5H_{10}O_5)$ .  
3-حامض الفسفوريك.

والشكل الآتي يبين العرض الهيكل لسلسلة الـ RNA



الشكل ( 2 ) يمثل العرض الهيكل لسلسلة الـ RNA

وفي حالات نادرة قد يكون الحامض النووي الرايبوسومي RNA هو المادة الوراثية كما في بعض الرواشح (الفيروسات)، ولكن أهميته أعظم وأعم في الأحياء بسبب الدور الذي يقوم به في عملية بناء البروتين بما في ذلك الأنزيمات.

ويمكن تمييز ثلاثة أنواع من الـ RNA وكلها تصنع في النواة وتنتقل إلى الساييتوبلازم لتشارك في صنع البروتين بشكل خاص وهي:

- 1- الحامض الرايبوزي المرسل (mRNA) Messenger RNA: ويقوم بنقل المعلومات الوراثية من DNA النواة إلى مناطق نشاطه في الساييتوبلازم.
- 2- الحامض الرايبوزي الريبوسومي (rRNA) Ribosomal RNA: ويدخل في تركيب الريبوسومات إذ يشترك مع البروتين في بنائها.
- 3- الحامض النووي الرايبوزي الناقل (tRNA) Transfer RNA: ويوجد في الساييتوبلازم ويقوم بنقل الأحماض الأمينية إلى الريبوسومات.

#### 4. توزيع بولتزمان: Boltzmann Distribution

تعرف دالة كثافة الاحتمال لتوزيع بولتزمان بالشكل الآتي: [3]

$$P(x; \theta) = \begin{cases} \exp \frac{-\Delta G(x; \theta)}{kT} & \text{if } \Delta G(x; \theta) > 0 \\ 1 & \text{if } \Delta G(x; \theta) \leq 0 \end{cases} \dots(1)$$

إذ أن:  $\Delta G$  : التغير في الطاقة، وتعرف بالصيغة الآتية:

$$\Delta G = \sum_x \frac{1}{y!} \sum_{i=0}^{y-1} (-1)^i \frac{y!}{i!(y-i)!} (y-i)^x \theta$$

k: ثابت بولتزمان و T : درجة الحرارة المطلقة.

## 5. تقدير معلمة توزيع بولتزمان باستخدام الخوارزمية الجينية:

### The Estimate Parameter of the Boltzmann Distribution Using Genetic Algorithm

سوف يتم في هذا البحث اقتراح خوارزمتين لتقدير معلمة توزيع بولتزمان باستخدام الخوارزمية الجينية

التي بدورها سوف تقوم بالمهام التالية:

(أ) تعظيم دالة الترجيح لتوزيع بولتزمان.

(ب) تقليل قيد الجبل لتوزيع بولتزمان.

#### 5.1 تعظيم دالة الترجيح لتوزيع بولتزمان باستخدام الخوارزمية الجينية:

##### Maximizing the Likelihood Function of Boltzmann Distribution by Using Genetic Algorithm

سوف يتم إيجاد تقدير المعلمة  $\theta$  لتوزيع بولتزمان الشرطي باستخدام طريقة الترجيح الأعظم، وهنا سنعرف سلسلة الـ (RNA) بـ (x)، ونعرف الاحتمالية لتكوين سلسلة الـ (RNA) بـ (y)، ونستخدم النموذج الخطي- اللوغاريتمي الشرطي لتوزيع بولتزمان (Boltzmann Distribution) ثم نقوم بتعظيم دالة الترجيح لتوزيع بولتزمان باستخدام الخوارزمية الجينية.

#### 5.2 اشتقاق الـ MLE لتوزيع بولتزمان

فيما يأتي الاشتقاق الكامل لتقدير الترجيح الأعظم (M.L.E.) لتوزيع بولتزمان الشرطي، إذ أن دالة

كثافة الاحتمال لتوزيع بولتزمان الشرطي (Conditional Boltzmann Distribution) تعرف كالآتي: [3]

$$p(y/x, \theta) = \frac{1}{z(x, \theta)} \exp\left(\frac{-1}{kT} \Delta G(x, y, \theta)\right) \quad \dots(2)$$

إذ أن  $\Delta G$  تعرف بالصيغة الآتية:

$$\Delta G = \sum_x \frac{1}{y!} \sum_{i=0}^{y-1} (-1)^i \frac{y!}{i!(y-i)!} (y-i)^x \theta \quad \dots(3)$$

وأن الدالة الجزئية  $z(x, \theta)$  تمثل بـ:

$$z(x, \theta) = \sum_y \exp\left(\frac{-1}{kT} \Delta G(x, y, \theta)\right) \quad \dots(4)$$

$$p(y/x, \theta) = \frac{\exp\left(\frac{-1}{kT} \sum_x \frac{1}{y!} \sum_{i=0}^{y-1} (-1)^i \frac{y!}{i!(y-i)!} (y-i)^x \theta\right)}{\sum_y \exp\left(\frac{-1}{kT} \sum_x \frac{1}{y!} \sum_{i=0}^{y-1} (-1)^i \frac{y!}{i!(y-i)!} (y-i)^x \theta\right)} \quad \dots(5)$$

ولإيجاد مقدر الترجيح الأعظم لتوزيع بولتزمان الشرطي نتبع الخطوات الآتية:

أولاً: نجد دالة الإمكان  $L(\theta)$  كالآتي:

$$L(\theta) = \prod_{r=1}^n p(y/x, \theta)$$

$$\prod_{r=1}^n = \sum_{m=1}^{16} \sum_{m=1}^{16} \quad \text{إذ أن:}$$

$$= \prod_{r=1}^n \frac{\exp\left(\frac{-1}{kT} \sum_x \frac{1}{y!} \sum_{i=0}^{y-1} (-1)^i \frac{y!}{i!(y-i)!} (y-i)^x \theta\right)}{\sum_y \exp\left(\frac{-1}{kT} \sum_x \frac{1}{y!} \sum_{i=0}^{y-1} (-1)^i \frac{y!}{i!(y-i)!} (y-i)^x \theta\right)} \dots (6)$$

$$= \left( \frac{\exp \sum_{r=1}^n \left( \frac{-1}{kT} \sum_x \frac{1}{y!} \sum_{i=0}^{y-1} (-1)^i \frac{y!}{i!(y-i)!} (y-i)^x \theta \right)}{\prod_{r=1}^n \sum_y \exp\left(\frac{-1}{kT} \sum_x \frac{1}{y!} \sum_{i=0}^{y-1} (-1)^i \frac{y!}{i!(y-i)!} (y-i)^x \theta\right)} \right) \dots (7)$$

ثانياً: نأخذ اللوغاريتم لدالة  $L(\theta)$  لتصبح المعادلة كالتالي:

$$\ln L(\theta) = \log \left( \frac{\exp \sum_{r=1}^n \left( \frac{-1}{kT} \sum_x \frac{1}{y!} \sum_{i=0}^{y-1} (-1)^i \frac{y!}{i!(y-i)!} (y-i)^x \theta \right)}{\prod_{r=1}^n \sum_y \exp\left(\frac{-1}{kT} \sum_x \frac{1}{y!} \sum_{i=0}^{y-1} (-1)^i \frac{y!}{i!(y-i)!} (y-i)^x \theta\right)} \right) \dots (8)$$

وبتبسيط المعادلة أعلاه نحصل على ما يلي:

$$\log L(\theta) = - \sum_{x=1}^4 \sum_{y=1}^4 \left( \ln \sum_{m=1}^{16} \sum_{m=1}^{16} \exp\left(\frac{-1}{kT} \sum_x \frac{1}{y(m)!} \sum_{i=0}^{y-1} (-1)^i \frac{y(m)!}{i!(y(m)-i)!} (y(m)-i)^{x(m)} \theta\right) - \dots (9) \right)$$

$$\sum_{m=1}^{16} \sum_{m=1}^{16} \left( \frac{-1}{kT} \sum_x \frac{1}{y(m)!} \sum_{i=0}^{y-1} (-1)^i \frac{y(m)!}{i!(y(m)-i)!} (y(m)-i)^{x(m)} \theta \right)$$

إن دالة كثافة الاحتمال لتوزيع بولتزمان الشرطي (Conditional Boltzmann Distribution) تعرف

كالتالي: [3].

$$p(y/x, \theta) = \frac{1}{z(x, \theta)} \exp\left(\frac{-1}{kT} \Delta G(x, y, \theta)\right)$$

### 5.3 الخوارزمية الجينية المقترحة (1):(1) Proposed Genetic Algorithm

إيجاد قيمة المقدر الذي يعظم دالة الترجيح لتوزيع بولتزمان

#### Finding the Value of the Estimator which Maximizes the Likelihood Function for Boltzmann Distribution

تم استخدام دالة الترجيح الأعظم لتوزيع بولتزمان في الخطوات المقترحة للخوارزمية الجينية، والخطوات

موضحة في أدناه :

الخطوة الأولى:- البيانات الأولية (Initial Data): وهي قراءة لمجموعة من المتغيرات التي استخدمت في

الخوارزمية:

- العداد m: وهو يمثل الطول لكل من سلسلتي x و y .

•  $x(m)$ : يمثل متجه لأحد النيوكليوتيدات الداخلة في تركيب الـ RNA ويمثل القواعد النتروجينية: (أدينين A، يوراسيل U، سايتوسين C، كوانين G) وكل عنصر مكرر أربع مرات، إذ تم تشفير هذه التراكيب ورمزنا للأدينين بالرمز (000)، ولليوراسيل بالرمز (001)، وللسايتوسين بالرمز (111)، والكوانين بالرمز (010).

•  $y(m)$ : يمثل متجه التراكيب الناتجة من تداخل هذه القواعد النتروجينية وللتكرارات أعلاه (فقط التراكيب المقبولة التي فيها يرتبط الأدينين A مع اليوراسيل U بأصرتين هيدروجينيتين  $A=U$  ويرتبط السايتوسين C مع الكوانين G بثلاث أوامر هيدروجينية  $C\equiv G$ )، وبهذا سوف يصبح تركيب الـ AU هو (000001)، و تركيب الـ CG هو (111010)، و تركيب الـ UA هو (001000)، و تركيب الـ GC هو (010111).

•  $k$ : يمثل ثابت الغاز.

•  $T$ : تمثل درجة الحرارة المطلقة.

**الخطوة الثانية:- إنشاء الجيل الابتدائي (Initial Generation):** الجيل الابتدائي في هذه الخوارزمية يتكون من كروموسومين احدهما يمثل القواعد النتروجينية والآخر يمثل تراكيب هذه القواعد النتروجينية ، وقد تم وضع القيم الحقيقية في جينات الكروموسومين أي أن التشفير كان تشفيراً حقيقياً للكروموسومات.

**الخطوة الثالثة:- قيمة الجودة (Fitness Value):** إن قيمة الجودة في هذه الخوارزمية تمثل مُقدر الترجيح الأعظم لتوزيع بولتزمان، إذ تم استخدام اللوغاريتم لأنه الأساس في استخراج المقدرات التي تجعل المقدر أعظم ما يمكن.

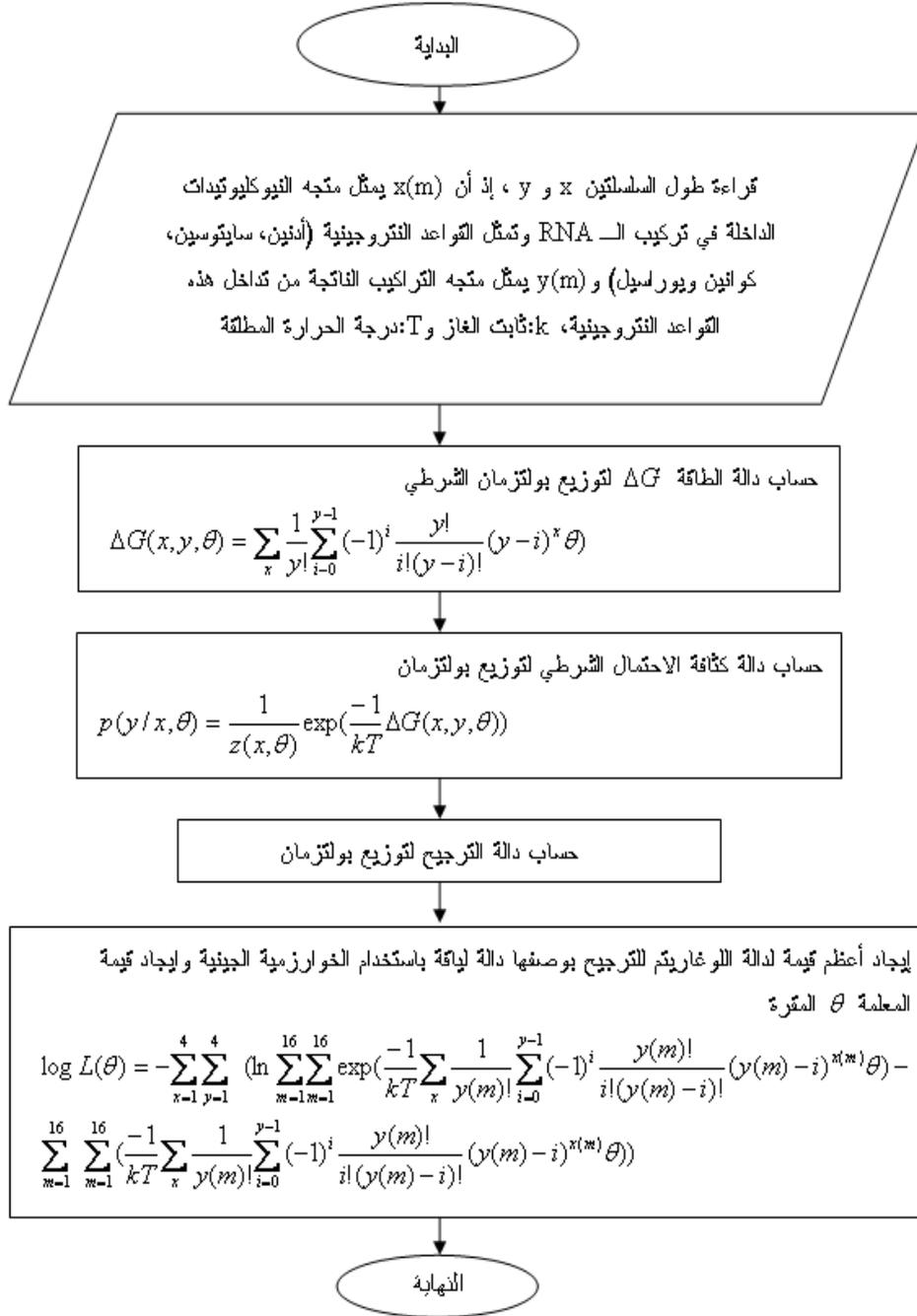
**الخطوة الرابعة:-** تم استخدام الانتقاء من النوعين الآتيين (Uniform، Roulette).

**الخطوة الخامسة:-** تم استخدام التعابر من الأنواع الآتية (Scattered، Intermediate، Single Point) Heuristic

**الخطوة السادسة:-** تم استخدام الطفرة من الأنواع الآتية (Gaussian، Uniform).

وتم التحديد المسبق لعدد الأجيال في توقف الخوارزمية الجينية.

والشكل الآتي يمثل المخطط الانسيابي للخوارزمية الجينية المقترحة:



الشكل (3) يمثل المخطط الانسيابي لتقدير قيمة المعلمة  $\theta$

الذي يعظم دالة الترجيح لتوزيع بولتزمان

تم تصميم برنامج بلغة (MATLAB) ثم تم ايجاد قيمة دالة الجودة (أعظم قيمة لدالة الترجيح) ومن خلاله تم تحديد المعلمة التي ستحقق القيمة العظمى للدالة وتم الاعتماد على التحديد المسبق لعدد الأجيال لبيان مدى التقرب من الحل الأمثل وقمنا باختيار عدة أنواع لكل من (الانتقاء ، التداخل الابدالي والطفرة) وكذلك عدد مرات توليد الأجيال وملاحظة النتائج بعد كل جيل إذ كانت أفضل النتائج لدى توليد الأجيال

ما بين 25 إلى 1500 جيل ولقد حققت هذه الخوارزمية أعظم قيمة لمقدر الترجيح لتوزيع بولتزمان عدد القيمة (-94.1614) وقيمة  $\theta$  هي (0.1457) والتي ظهرت عند الجيل (1387) إذ كان:

1. الانتقاء من نوع (Roulette).
2. التداخل الابدالي من نوع (Intermediate).
3. الطفرة من نوع (Uniform).
- كما موضح في الجدول (1) .

جدول (1) يمثل تقدير المعلمة  $\theta$  بأسلوب تعظيم دالة الترجيح لتوزيع بولتزمان

عدد الأجيال	نوع الانتقاء	نوع التداخل الابدالي (النحائر)	نوع الطفرة	قيمة $\theta$	أعظم قيمة لمعقد M.L.E	
1.	100	Roulette	Intermediate	Uniform	-21.0133	-110.018
2.	77	Roulette	Intermediate	Uniform	-17.1375	-102.0289
3.	80	Uniform	Intermediate	Gaussian	-8.0851	-95.5772
4.	29	Uniform	Scattered	Gaussian	-4.1539	-94.7458
5.	100	Uniform	Scattered	Gaussian	-12.9348	-97.6835
6.	25	Uniform	Single Point	Uniform	0.0229	-94.1751
7.	52	Uniform	Intermediate	Uniform	0.044	-94.1727
8.	100	Roulette	Scattered	Gaussian	-23.5084	-117.5638
9.	73	Roulette	Intermediate	Uniform	-18.9013	-105.1023
10.	100	Roulette	Single Point	Uniform	0.0084	-94.1767
11.	171	Uniform	Scattered	Gaussian	-7.7767	-95.4946
12.	100	Uniform	Single Point	Gaussian	-13.0772	-97.7796
13.	100	Uniform	Scattered	Gaussian	-15.6705	-100.1023
14.	100	Roulette	Heuristic	Uniform	-29.8607	-142.7702
15.	100	Stochastic Uniform	Scattered	Gaussian	-20.5965	-108.9377
16.	84	Roulette	Intermediate	Uniform	-18.2274	-103.8222
17.	59	Uniform	Heuristic	Gaussian	-17.3304	-102.3224
18.	99	Roulette	Intermediate	Uniform	-19.8354	-107.105
19.	35	Uniform	Scattered	Gaussian	-5.3125	-94.9497
20.	49	Uniform	Single Point	Uniform	0.1016	-94.1663
21.	136	Roulette	Scattered	Uniform	0.0815	-94.1685
22.	234	Uniform	Scattered	Gaussian	-1.8355	-94.4026
23.	567	Roulette	Intermediate	Uniform	0.0060	-94.1770
24.	784	Roulette	Single Point	Gaussian	0.0326	-94.1740
25.	923	Uniform	Scattered	Uniform	0.0435	-94.1728
26.	419	Roulette	Single Point	Gaussian	-45.8077	-216.4093
27.	1097	Roulette	Heuristic	Uniform	-2.1757	-94.4485
28.	1265	Uniform	Scattered	Gaussian	-5.7313	-95.0304
29.	1387	Roulette	Intermediate	Uniform	0.1457	-94.1614
30.	1500	Uniform	Single Point	Gaussian	-4.2841	-94.7614

#### 5.4 تقليل قيد الجيل لتوزيع بولتزمان باستخدام الخوارزمية الجينية: Minimizing Constraint Generation (CG) of Boltzmann Distribution Using Genetic Algorithm

الطريقة البديلة لطريقة الترجيح في تقدير المعلمة هي طريقة قيد الجيل والتي اقترحت من قبل Andronescu وآخرون في عام (2006) والطريقة تتلخص بتقدير قيمة  $\theta$  لتوزيع بولتزمان بحل نظام القيود الموضح في المعادلة الآتية: [10]

$$G(x, y_x, \theta) \leq \Delta G(x, y, \theta) \quad \dots(10)$$

$$\text{لكل } (x, y_x) \in S \text{ و } y_x \in Y_x$$

إذ أن:

S: تمثل مجموعة بيانات تتكون من مجموعة أزواج  $(x, y_x)$ .

$y_x$ : تمثل مجموعة كل التراكيب الثانوية للمتسلسلة x.

هذه القيود تضمن أن لكل سلسلة x جميع تراكيبها الثانوية y غير المثلى والتي يرمز لها بـ  $(y_x)$  لها أعلى طاقة حرة دنيا. (من خلال ذلك نفرض انه لا يوجد تركيب آخر له نفس الطاقة الحرة الدنيا (Minimum Free Energy) (MFE) كالتركيب المعروف لدينا).

قد يحدث أن هذا النظام من القيود يكون غير قابل للتطبيق، أي انه لا يوجد حل للمعلمة  $(\theta)$  الموجودة في توزيع بولتزمان والتي تعرف كل القيود بشكل آني. وللتعامل مع التطبيقات غير الممكنة، سنضيف المتغيرات البطيئة  $\delta_{x,y} \geq 0$  إلى القيود، والتي سوف يتم تقليل قيمتها، وهذا يؤدي إلى استرخاء (Relaxation) القيود (أي تحسين قيمتها) كما موضح في المعادلة الآتية:

$$\Delta G(x, y_x, \theta) \leq \Delta G(x, y, \theta) + \delta_{x,y} \quad \dots(11)$$

وعلى اعتبار أن دالة الطاقة  $\Delta G$  معرفة سابقاً في المعادلة (3)، ويمكن تحويل المتراجحة أعلاه إلى

معادلة معرفة كالآتي:

$$\Delta G(x, y_x, \theta) - \Delta G(x, y, \theta) = \delta_{x,y} \quad \dots(12)$$

وبالتالي يجب أن تتحقق دالة قيد الجيل والممثلة بالمعادلة الآتية:

$$\text{minimize } \|\delta\|^2 \quad \dots(13)$$

مع تحقق القيد:

$$\delta \geq 0 \quad \dots(14)$$

تم استخدام المعادلة (13) مع القيد (14) كدالة هدف للخوارزمية الجينية المقترحة التالية والتي يتم فيها

تقليل دالة الهدف مع تحقق قيدها.

#### 5.5 الخوارزمية الجينية المقترحة (2): Proposed Genetic Algorithm (2)

إيجاد قيمة المقدر الذي يقلل قيد الجيل لتوزيع بولتزمان

#### Finding the Value of the Estimator which Minimizes the Constraint Generation for Boltzmann Distribution

استخدمت دالة قيد الجيل (CG) لتوزيع بولتزمان في الخطوات المقترحة للخوارزمية الجينية، والخطوات

موضحة في أدناه:

الخطوة الأولى والخطوة الثانية كما في مثلتيهما في الخوارزمية السابقة.

الخطوة الثالثة:- قيمة الجودة (Fitness Value): إن قيمة الجودة في هذه الخوارزمية ممثلة بالمعادلة (3.14) إذ استُخدم الرمز  $z$  ليمثل ناتج طرح دوال الطاقة في الحالتين: عند تكرار عدد التراكيب  $\Delta G(x,y, \theta)$  ، وعند ثبوت عدد التراكيب التي رُمز لها بـ  $\Delta G(x,y, \theta)$  ، أي أن:

$$z = \Delta G(x,y, \theta) - \Delta G(x,y, \theta)$$

وفي هذه الحالة يجب تحقق الشرط التالي وهو أن  $z$  أكبر أو تساوي الصفر، عندها سوف يتم حساب معيار الطول (Norm) لـ  $z$  التي تمثل المعادلة (3.14) بعدها سوف يتم تقليل قيمتها.

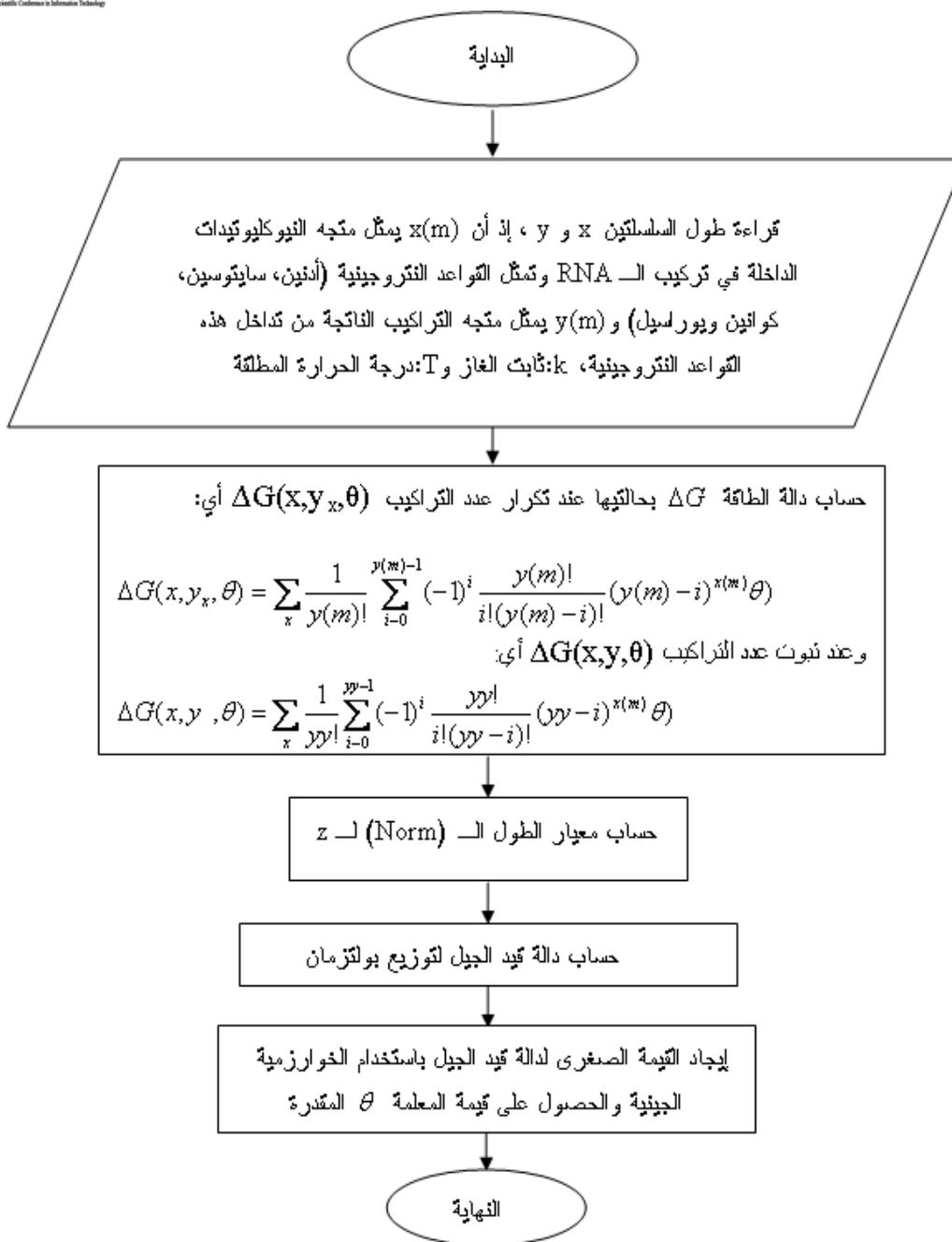
الخطوة الرابعة:- تم استخدام الانتقاء من النوعين الآتيين (Uniform، Roulette).

الخطوة الخامسة:- تم استخدام التعابر من الأنواع الآتية (Scattered، Intermediate، Single Point)، Heuristic، Two Point

الخطوة السادسة:- تم استخدام الطفرات الآتية (Gaussian، Uniform).

وتم التحديد المسبق لعدد الأجيال في توقف الخوارزمية الجينية.

والشكل التالي يمثل المخطط الانسيابي للخوارزمية الجينية المقترحة:



الشكل (4) يمثل المخطط الانسيابي لإيجاد قيمة المعلمة  $\theta$  الذي يقلل دالة قيد الجيل لتوزيع بولتزمان

بعد أن تم إيجاد قيمة الجودة (أقل قيمة لدالة قيد الجيل)، وذلك بتصميم برنامج بلغة MATLAB ومن خلاله تم تحديد المعلمة التي ستحقق القيمة الدنيا للدالة وتم اختيار عدة أنواع لكل من (الانتقاء، التداخل الابدالي والطفرة) وكذلك عدد مرات توليد الأجيال وملاحظة النتائج بعد كل جيل إذ كانت أفضل النتائج لدى توليد الأجيال ما بين 25 إلى 1500 جيل وتم الاعتماد على التحديد المسبق لعدد الأجيال لبيان مدى التقرب من الحل الأمثل ولقد حققت هذه الخوارزمية أقل قيمة لدالة قيد الجيل لتوزيع بولتزمان عند القيمة  $(95.1039101 * 10^7)$  وقيمة  $\theta$  هي (-4.4066) والتي ظهرت عند الجيل (80) إذ كان:

1. الانتقاء من نوع (Uniform).
  2. التداخل الابدالي من نوع (Intermediate).
  3. الطفرة من نوع (Gaussian).
- كما موضح في الجدول (2).

جدول (2) يمثل تقدير المعلمة  $\theta$  بأسلوب تقليل قيد الجيل CG لتوزيع بولترمان

رقم	عدد الأجيال	نوع الانتقاء	نوع التداخل الابدالي (التعابر)	نوع الطفرة	قيمة $\theta$	أقل قيمة لمقدر قيد الجيل (CG)
1.	64	Roulette	Two Point	Uniform	0.9267	$-0.18502 \cdot 10^1$
2.	100	Stochastic Uniform	Scattered	Gaussian	-21.6482	$-0.22952 \cdot 10^1$
3.	49	Uniform	Single Point	Uniform	0.9944	$-0.21303 \cdot 10^1$
4.	29	Uniform	Scattered	Gaussian	-4.6912	$-0.10778 \cdot 10^1$
5.	99	Roulette	Intermediate	Uniform	0.9895	$-0.21096 \cdot 10^1$
6.	100	Roulette	Intermediate	Uniform	0.9776	$-0.20589 \cdot 10^1$
7.	100	Uniform	Scattered	Gaussian	-10.9584	$-0.58814 \cdot 10^1$
8.	123	Uniform	Scattered	Gaussian	-3.4385	$-0.579081265 \cdot 10^7$
9.	80	Uniform	Intermediate	Gaussian	-4.4066	$-0.951039101 \cdot 10^7$
10.	100	Uniform	Single Point	Gaussian	-12.4559	$-0.75987 \cdot 10^1$
11.	52	Uniform	Intermediate	Uniform	0.9609	$-0.19892 \cdot 10^1$
12.	73	Roulette	Intermediate	Uniform	0.8569	$-0.15820 \cdot 10^1$
13.	25	Uniform	Single Point	Uniform	0.9916	$-0.21184 \cdot 10^1$
14.	100	Uniform	Scattered	Gaussian	-6.5909	$-0.21276 \cdot 10^1$
15.	84	Roulette	Intermediate	Uniform	0.8834	$-0.16815 \cdot 10^1$
16.	59	Uniform	Heuristic	Gaussian	-26.4713	$-0.34319 \cdot 10^1$
17.	350	Uniform	Scattered	Gaussian	-1.7094	$-0.143105685 \cdot 10^7$
18.	100	Roulette	Heuristic	Uniform	1.3858	$-0.41378 \cdot 10^1$
19.	77	Roulette	Intermediate	Uniform	0.9625	$-0.19961 \cdot 10^1$
20.	100	Roulette	Scattered	Gaussian	-21.5456	$-0.22736 \cdot 10^1$
21.	178	Uniform	Heuristic	Gaussian	-50.4029	$-0.12442 \cdot 10^1$
22.	232	Roulette	Intermediate	Gaussian	-7.8907	$-0.304943 \cdot 10^1$
23.	489	Uniform	Scattered	Uniform	0.9889	$-0.21066 \cdot 10^1$
24.	678	Uniform	Intermediate	Gaussian	-0.5562	$-0.15151389 \cdot 10^6$
25.	81	Uniform	Intermediate	Gaussian	-8.5679	$-0.35953 \cdot 10^1$
26.	1010	Uniform	Scattered	Gaussian	-1.8179	$-0.161848358 \cdot 10^7$
27.	1215	Roulette	Scattered	Uniform	-5.2872	$-0.13691 \cdot 10^1$
28.	1365	Uniform	Intermediate	Gaussian	-1.4349	$-0.100842706 \cdot 10^7$

## 6. الاستنتاجات:

إيجاد قيمة المقدر الذي يعظم دالة الترجيح لتوزيع بولتزمان باستخدام الخوارزمية الجينية المقترحة الأولى استنتج ما يأتي:

أعظم قيمة لمقدر الترجيح لتوزيع بولتزمان ظهرت:

1. عند القيمة (-94.1614) حيث كانت قيمة  $\theta$  هي (0.1457).

2. عند الجيل (1387).

3. الانتقاء من نوع (Roulette).

4. كذلك التداخل الابدالي من نوع (Intermediate).

5. الطفرة من نوع (Uniform).

إيجاد قيمة المقدر الذي يقلل قيد الجيل لتوزيع بولتزمان باستخدام الخوارزمية الجينية المقترحة الثانية استنتج ما يأتي:

أقل قيمة لقيد الجيل لتوزيع بولتزمان ظهرت:

1. عند القيمة ( $-0.95103910 * 10^7$ ) حيث كانت قيمة  $\theta$  هي (-4.4066).

2. عند الجيل (80).

3. الانتقاء من نوع (Uniform).

4. التداخل الابدالي من نوع (Intermediate).

5. الطفرة من نوع (Gaussian).

استنتجنا مما سبق بأن الخوارزمية الجينية كانت الأسلوب الأفضل لتحقيق أعظم قيمة لدالة الترجيح والحصول على أمثل قيمة للمعلمة  $\theta$  وكذلك حققت أدنى قيمة لدالة قيد الجيل والحصول أيضاً على أمثل قيمة للمعلمة  $\theta$ .

### المصادر

- [1]. ثابت، همسة معن. (2005)، "بعض تطبيقات الخوارزمية الجينية في حل مسائل الأمثلية"، رسالة ماجستير، جامعة الموصل، كلية علوم الحاسبات والرياضيات، قسم الإحصاء.
- [2]. ذنون، باسل يونس (1988)، "الإحتمالية والمتغيرات العشوائية"، مديرية دار الكتب والنشر، جامعة الموصل، وزارة التعليم العالي والبحث العلمي.
- [3]. Andronescu, M.; Condon, A.; Hoos, H. H.; Mathews, D. H. and Murphy, K., (2006), " **Efficient Parameter Estimation for RNA Secondary Structure Prediction**", © Oxford University Press.
- [4]. Chan, C. Y. and Ding, Y., (2008), " **Boltzmann Ensemble Features of RNA Secondary Structures: a Comparative Analysis of Biological RNA Sequences and Random Shuffles**", C. Y. Chan · Y. Ding Wadsworth Center, New York State Department of Health, Center for Medical Science, 150 New Scotland Avenue, Albany, NY 12208, USA.
- [5]. Devore, J. R., (2000), "**Probability and Statistics**", Brooks, Cole.
- [6]. Ding, Y. and Lawrence, C. E., (2003), "**A statistical Sampling Algorithm for RNA Secondary Structure Prediction. Nucleic Acids Research**", 31: pp. 7280–7301.
- [7]. Ding, Y., (2006), "**Statistical and Bayesian Approaches to RNA Secondary Structure Prediction**", Published by Cold Spring Harbor Laboratory Press, 12: pp. 323–331.
- [8]. Ding, Y.; Chan, C. Y. and Lawrence, C. E., (2006), " **Clustering of RNA Secondary Structures with Application to Messenger RNAs**", Elsevier Ltd, 0022-2836.
- [9]. Irwin, H. H., (1967), "**Basic Principle of Molecular Genetice**", Nelson Ltd. London.
- [10]. Pond, K. S. L.; Mannino, F. V. ; Gravenor, M. B.; Muse, S. V. and Frost, S. D. W., (2006), "**Evolutionary Model Selection with a Genetic Algorithm: A Case Study Using Stem RNA**", Oxford University Press on Behalf of the Society for Molecular Biology and sEvolution. Moscow.
- [11]. William, T. T. and Richared, J. W., (1968), "**General Biology**", 2<sup>nd</sup> Edition, Van Nostrand Reinhold Company.
- [12]. "**Genetic Algorithms**", the web site at  
<http://www.cs.felk.cvut.cz/~xobitko/ga/main.html> Obitko M. (2009)
- [13]. Marczyk, A., (2004), "**Genetic Algorithm and Evolutionary Computation**", April, Vol. 23.  
<http://www.talkor:gins.org/Faqs/genalg/genalg.html/#introduction>.