

بِسْمِ اللَّهِ الرَّحْمَنِ  
الرَّحِيمِ



دانشکده فنی و مهندسی

گروه مهندسی کامپیوتر  
گرایش: نرم افزار

عنوان:

پیش بینی انواع لوسمی مبتنی بر ترکیب مدل ها با  
استفاده از الگوریتم ژنتیک

استاد راهنما:

جناب آقای دکتر عباس عکاسی

محققین:

کمیل آقابابایی، خاطره اصغری، امیر جعفری زاده،  
نوید چمنی، مونا رجبی فرد، مژگان زرافشان، سهیلا  
سالاری، مرتضی سهرابی، محبوبه صادقی، سارا عباس  
زاده، سعید ناصری، مریم یوسفی زاده

آذر ماه 1397

با تشكر فراوان از جناب آقاى دكتر عباس عكاسى كه مارا يارى  
كردند.

تقديم به او كه جهان در انتظارش است.  
تقديم به او كه منجی عالم بشریت است.  
اللهم عجل لوليک فرج ...

## چکیده

لوسمی یکی از شایع ترین علل مرگ و میر در سراسر جهان است. پیش بینی دقیق نوع لوسمی برای استفاده از داروهای خاص مهم است. از این رو پژوهش در این زمینه ادامه دارد. تحلیل داده های ریز آرایه<sup>1</sup> یک روش کارآمد برای پیش بینی نوع لوسمی است، که بدون استفاده از الگوریتم های داده کاوی دشوار است. هدف این مقاله طبقه بندی<sup>2</sup> انواع لوسمی لنفوبلاستی حاد<sup>3</sup> و لوسمی میلوئید<sup>4</sup> حاد با استفاده از ترکیب 100 طبقه بند است. در این مطالعه، توصیفی از داده های بیان 7129 ژن مربوط به 72 بیمار مبتلا به لوسمی استفاده شد. سپس یک الگوریتم از بین الگوریتم های ماشین بردار پشتیبان<sup>5</sup>، الگوریتم بیزین<sup>6</sup> و الگوریتم درخت تصمیم<sup>7</sup> بر اساس رای گیری اکثریت<sup>8</sup>، به عنوان طبقه بند پایه انتخاب شد. با الگوریتم طبقه بند پایه 100 مدل ساخته شد. ترکیب 100 مدل با استفاده از الگوریتم ژنتیک<sup>9</sup> انجام شد. نتایج حاصل از مقایسه ی روش پیشنهادی با سایر روش های ترکیب، عملکرد خوب الگوریتم ژنتیک را نشان داد.

کلمه کلیدی: سرطان لوسمی، داده کاوی، الگوریتم ژنتیک، طبقه بندی، ارزیابی، یادگیری تجمعی

---

<sup>1</sup> Microarray

<sup>2</sup> Classification

<sup>3</sup> Acute lymphoblastic

<sup>4</sup> Acute myeloid

<sup>5</sup> Support vector machine

<sup>6</sup> Naïve byes algorithm

<sup>7</sup> Decision tree algorithm

<sup>8</sup> Majority voting

<sup>9</sup> Genetic algorithm

1..... فصل اول: کلیات تحقیق و اصول کلی.....	1-1
2..... مقدمه.....	1-1-1
3..... پیشنهادی.....	2-1
3-1-آماده..... سازی.....	2-1-2
4..... داده.....	3-1
4-1-انتخاب الگوریتم پایه و ساخت 100 طبقه بند با استفاده از آن 4	3-1-1
4-1-1-الگوریتم درخت تصمیم.....	3-1-1-1
4-1-2-الگوریتم ماشین بردار پشتیبان.....	3-1-1-2
4-1-3-الگوریتم بیزین.....	3-1-1-3
5-1-..... روش.....	3-1-1-4
5-1-..... کار.....	3-1-1-5
6-1-ارزیابی..... مدل.....	3-1-1-6
7-1-یادگیری..... ترکیبی.....	3-1-1-7
6-1-7-1-انتخاب طبقه بندها.....	3-1-1-7-1-1
6-1-7-2-ترکیب طبقه بندها.....	3-1-1-7-1-2
6-1-7-2-1-ترکیب همه ی طبقه بندها.....	3-1-1-7-1-2-1
7-1-7-2-2-انتخاب رو به جلو.....	3-1-1-7-1-2-2
7-1-7-2-3-حذف رو به عقب.....	3-1-1-7-1-2-3
8-1-7-2-4-ترکیب با استفاده از الگوریتم ژنتیک.....	3-1-1-7-1-2-4
9-1-7-2-5-ساختار الگوریتم ژنتیک.....	3-1-1-7-1-2-5
11-1-7-2-6-روند الگوریتم ژنتیک.....	3-1-1-7-1-2-6
12-1-7-2-7-روشهای انتخاب.....	3-1-1-7-1-2-7
12-1-7-2-8-شرط پایان الگوریتم.....	3-1-1-7-1-2-8
12-1-7-2-9-برخی از کاربرد الگوریتمهای ژنتیکی.....	3-1-1-7-1-2-9
12-1-7-2-10-شرط..... پایان.....	3-1-1-7-1-2-10
12-1-7-2-11-الگوریتم.....	3-1-1-7-1-2-11

**14..... فصل دوم : نتایج**

1-2-نتایج جداسازی مجموعه داده 15.....

2-2-نتایج ارزیابی مدل های ساخته شده با الگوریتم های ماشین بردار پشتیبان و بیز و درخت تصمیم روی development data : 15.....

2-3-ایجاد 100 مدل با الگوریتم پایه 15.....

2-4-مقایسه ی روش های ترکیب 16.....

**17..... منابع**

## فهرست شکل‌ها

عنوان  
صفحه

---

شکل 1-1- گردش کار روش پیشنهادی.....	3
شکل 1-2- آمیزش تک نقطه ایی [16].....	10
شکل 1-3- عمل جهش [19].....	11
شکل 1-4- گردش کار الگوریتم ژنتیک [21].....	11
شکل 1-5- بردار راه حل.....	13



## فهرست جداول

عنوان  
صفحه

---

جدول 1-1- تعداد نمونه ها و ویژگی های انواع لوسمی.....	4
جدول 1-2- ترکیب مدل ها بر مبنای رای گیری اکثریت.....	6
جدول 1-3- عملکرد روش انتخاب رو به جلو.....	7
جدول 1-4- عملکرد روش حذف رو به عقب.....	8
جدول 1-2- جدا سازی مجموعه داده .....15	15
جدول 2-2- مقایسه طبقه بند ماشین بردار پشتیبان ، درخت تصمیم و بیز	15
جدول 2-3- نتایج ارزیابی مدل های ایجاد شده با الگوریتم پایه.....	15
جدول 2-4- نتایج روش های ترکیب.....	16

# فصل اول کلیات تحقیق و اصول کلی

## 1-1- مقدمه

طبقه بندی انواع لوسمی طی 30 سال گذشته بهبود زیادی یافته است و تلاش های بسیار زیادی در این زمینه انجام شده است، اما روش های کمی برای اختصاص دادن تومور به کلاس های شناخته شده (پیش بینی کلاس) وجود دارد. برای غلبه بر این چالش و رسیدن به راه حلی برای طبقه بندی لوسمی، استفاده از طبقه بندی لوسمی صرفاً بر پایه بیان ژن می باشد. این روش یک استراتژی کلی برای کشف و پیش بینی لوسمی است و می تواند برای انواع دیگر سرطان، مستقل از دانش بیولوژیکی قبلی مورد استفاده قرار بگیرد. [1]

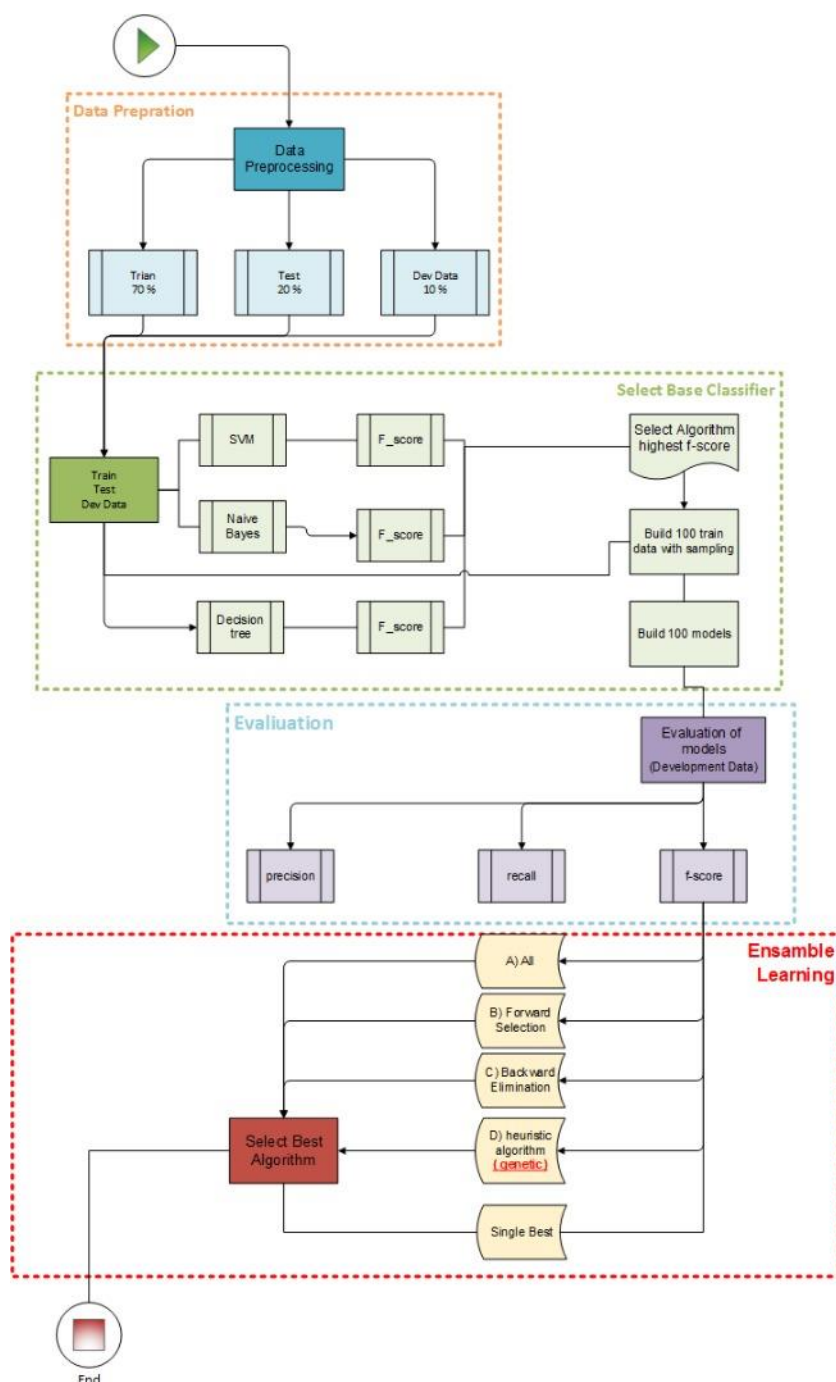
تکنولوژی ریز آرایه باعث می شود زیست شناسان قادر به نظارت بر بیان هزاران ژن در یک آزمایش واحد، در یک تراشه کوچک باشند. داده های مجموعه ای از بیان ژن ریز آرایه در اینترنت در دسترس عموم هستند. این مجموعه داده ها شامل تعداد زیادی از مقادیر بیان ژن هستند، از این رو نیاز به یک روش دقیق برای استخراج دانش و اطلاعات مفید از مجموعه داده های بیان ژن احساس می شود. [2]

نحوه طبقه بندی مهمترین چالش بعد از انتخاب نوع داده است. در این زمینه نیز تاکنون تلاش های بسیار زیادی انجام شده است که از جمله آنها می توان به دسته بندی لوسمی به وسیله الگوریتم SVM [3]، الگوریتم درخت تصمیم [4] و الگوریتم بیزین [5] اشاره کرد. مهمترین مشکل استفاده از این روش ها پائین بودن قابلیت اطمینان آنهاست.

در این مطالعه برای بالابردن قابلیت اطمینان بجای استفاده از یک داده آموزش و یک طبقه بند از 100 داده ی آموزش و 100 طبقه بند استفاده شد، برای انتخاب بهترین زیر مجموعه از بین طبقه بند ها از الگوریتم ژنتیک استفاده شد.

در ادامه به بررسی روش پیشنهادی و در فصل دوم نتایج و در آخر فهرست منابع آورده می شود.

## 2-1- روش پیشنهادی



شکل 1-1- گردش کار روش پیشنهادی

با توجه به شکل 1-1، روش پیشنهادی به چهار بخش تقسیم گردید، در هر بخش ابتدا مفاهیم مورد استفاده تعریف می شود، سپس روش کار بیان می شود.

### 3-1- آماده سازی داده

#### • مجموعه داده<sup>1</sup>

مجموعه داده leukemia از معتبرترین و استانداردترین مجموعه داده های معرفی شده در زمینه تشخیص نوع سرطان می باشد که در پژوهش های بسیاری مورد استفاده قرار گرفته است. این مجموعه داده یک رویکرد عمومی برای طبقه بندی سرطان خون بر اساس نظارت بر ژن توسط Microarrays DNA توصیف می کند. این مجموعه داده دارای 72 نمونه می باشد که از این تعداد، 47 نمونه مربوط به نوع لوسمی حاد لنفوبلاستیکی (ALL) و 25 نمونه مربوط به لوسمی حاد میلوئید (AML) می باشد. ساختار دیتاست متفاوت با ساختار دیتاست های مرسوم می باشد و هر نمونه به جای سطر به صورت ستونی ذخیره شده اند و برچسب تمام نمونه ها در سطر اول ذکر شده است. حداکثر تعداد ویژگی های هر نمونه برابر با 7128 می باشد که در نهایت یک ماتریس 7128\*72 را ایجاد می کند. [6]

جدول 1-1- تعداد نمونه ها و ویژگی های انواع لوسمی

	ALL	AML
تعداد نمونه	47	25
تعداد ویژگی	7128	7128

#### • پیش پردازش مجموعه داده

ابتدا مجموعه داده مسئله فراخوانی شد، سپس جداسازی مجموعه داده به بخش های داده های آموزش<sup>2</sup>، داده های آزمون<sup>3</sup> و داده های توسعه<sup>4</sup> به ترتیب با اندازه های 70%، 20% و 10% گرفت. این بخش ها با رعایت توزیع کلاس انتخاب شد، به این صورت که برای هر بخش به میزان مساوی از هر کلاس انتخاب شد.

### 4-1- انتخاب الگوریتم پایه و ساخت 100 طبقه بند با استفاده از آن

از الگوریتم های طبقه بندی درخت تصمیم، ماشین بردار پشتیبان و الگوریتم بیزین استفاده شد.

#### 1-4-1 الگوریتم درخت تصمیم

این الگوریتم با طرح این سوال که چه صفتی باید در ریشه درخت آزمایش شود آغاز می شود. برای پاسخ به این سوال، با استفاده از یکی از انواع آزمایش های آماری برای تعیین مناسب ترین صفت برای دسته بندی مثال های آموزشی، تصمیم براساس هر صفت نمونه را ارزیابی می کند. سپس بهترین صفت را انتخاب کرده و به عنوان تست در گره ریشه درخت استفاده می کند. برای هر مقدار ممکن صفت تست شده در ریشه، یک گره متناظر ایجاد شده و مثال های آموزشی براساس مقادیر صفت تست، بین این گره ها افراز می شوند. تمام فرایند ذکر شده،

<sup>1</sup> Data set

<sup>2</sup> Train data

<sup>3</sup> Test data

<sup>4</sup> Development data

با استفاده از مثال های آموزشی نسبت داده شده به هر گره، برای انتخاب بهترین صفت در گره های درخت تکرار می شود. این روش جستجوی حریصانه را برای یک درخت تصمیم قابل قبول ارائه می دهد که در این الگوریتم، هیچ گاه برای در نظر گرفتن دوباره انتخاب های قبلی، به عقب برگشت نمی شود.

درخت تصمیم که هدف اصلی آن، دسته بندی داده هاست، مدلی در داده کاوی است که ساختاری درخت مانند را جهت اخذ تصمیم و تعیین کلاس و دسته یک داده خاص به ما ارائه می کند. همان طور که از نام آن مشخص است، این درخت از تعدادی گره و شاخه تشکیل شده است به گونه ای که برگ ها کلاس ها یا دسته بندی ها را نشان می دهند و گره های میانی هم برای تصمیم گیری با توجه به یک یا چند صفت خاصه به کار می روند [7].

#### 1-4-2- الگوریتم ماشین بردار پشتیبان

این الگوریتم یک جداکننده است که با معیار قرار دادن بردارهای پشتیبان، بهترین دسته بندی و تفکیک بین داده ها را برای ما مشخص می کند. بردارهای پشتیبان به زبان ساده، مجموعه ای از نقاط در فضای  $n$  بعدی داده ها هستند که مرز دسته ها را مشخص می کنند و مرز بندی و دسته بندی داده ها براساس آنها انجام می شود و با جابجایی یکی از آنها، خروجی دسته بندی ممکن است تغییر کند. در فضای دو بعدی بردارهای پشتیبان، یک خط، در فضای سه بعدی یک صفحه و در فضای  $n$  بعدی یک ابر صفحه را شکل خواهند داد. در این الگوریتم فقط داده های قرار گرفته در بردارهای پشتیبان مبنای یادگیری ماشین و ساخت مدل قرار می گیرند و این الگوریتم به سایر نقاط داده حساس نیست و هدف آن هم یافتن بهترین مرز در بین داده هاست به گونه ای که بیشترین فاصله ممکن را از تمام دسته ها (بردارهای پشتیبان آنها) داشته باشد [8].

#### 1-4-3- الگوریتم بیزین

تئوری بیز یکی از روش های آماری برای رده بندی به شمار می آید. در این روش کلاس های مختلف، هر کدام به شکل یک فرضیه دارای احتمال در نظر گرفته می شوند. هر رکورد آموزشی جدید، احتمال درست بودن فرضیه های پیشین را افزایش و یا کاهش می دهد و در نهایت، فرضیاتی که دارای بالاترین احتمال شوند، به عنوان یک کلاس در نظر گرفته شده و برچسبی بر آن ها زده می شود. این تکنیک با ترکیب تئوری بیز و رابطه سببی بین داده ها، به طبقه بندی می پردازد [9].

#### 1-5- روش کار

این بخش به سه قسمت تقسیم شد:

- هر سه الگوریتم با داده های آموزشی، آموزش داده شد، سپس طبقه بند به دست آمده توسط داده های توسعه ارزیابی شد و برای هر الگوریتم  $F\_score$ ،  $precision$  و  $recall$  محاسبه شد، در نهایت الگوریتمی که دارای بهترین  $F\_score$  بود، به عنوان الگوریتم پایه انتخاب شد.
- 100 مجموعه داده آموزشی متفاوت به صورت رندوم با استفاده از نمونه برداری و با جایگذاری، ساخته شد.

- 100 مجموعه داده آموزشی به الگوریتم پایه داده شد و 100 مدل از آنها ساخته شد.

## 6-1- ارزیابی مدل ها

100 مدل ساخته شده در بخش قبلی روی داده های توسعه ارزیابی شد ، برای هر مدل F\_score , precision و recall محاسبه شد . مدل ها در یک آرایه بر مبنای F\_score به صورت صعودی مرتب شد و بهترین مدل<sup>1</sup> انتخاب شد.

## 7-1- یادگیری تجمعی<sup>2</sup>

یادگیری تجمعی ترکیبی از چندین مدل برای تولید یک مدل پیش بینی مطلوب است. [10]

این بخش شامل دو مرحله می شود :

### 1-7-1- انتخاب طبقه بندها

در یادگیری تجمعی فرض می کنیم استخری از طبقه بندها را داریم و می خواهیم زیر مجموعه ای از آن ها را انتخاب کنیم به گونه ای که ترکیب اعضای این زیرمجموعه بهترین نتیجه را بدهد. در این بخش ارزیابی روی داده های توسعه و داده های آزمایش مورد ارزیابی قرار گرفت.

### 2-7-1- ترکیب طبقه بندها

برای ترکیب طبقه بندها روش های متفاوتی وجود دارد ، در این مطالعه روش های ترکیب زیر مورد بررسی قرار داده شد.

### 1-2-7-1- ترکیب همه ی طبقه بندها

در این نوع ترکیب برای هر نمونه با توجه به صورت مسئله تعدادی طبقه بند وجود دارد که با توجه به خروجی هر طبقه بند، خروجی که بیشترین تکرار را داشته باشد به عنوان خروجی نهایی در نظر گرفته می شود.

به طور مثال 5 نمونه و 4 طبقه بند داریم که نتیجه را به صورت رای اکثریت نشان می دهد و در مواردی که تعداد هر دو کلاس با هم برابر باشد به صورت رندم یکی را به عنوان خروجی نهایی در نظر می گیرد.

جدول 1-2- ترکیب مدل ها بر مبنای رای گیری اکثریت

نمونه	E1	E2	E3	E4	Out
1	A	A	B	A	A
2	A	B	B	B	B
3	A	A	B	B	A
4	A	B	A	A	A
5	B	B	A	B	B

<sup>1</sup> Single best

<sup>2</sup> Ensemble learning

### 1-2-7-2- انتخاب رو به جلو<sup>1</sup>

انتخاب رو به جلو یک نوع رگرسیون گام به گام است که با یک مدل خالی شروع می شود و طبقه بندی را که به صورت نزولی بر مبنای یکی از معیارهای ارزیابی مرتب شده اند را یک به یک اضافه می کند و بهترین متغیر، با برخی از معیارهای از پیش تعیین شده تعیین می شود و به مدل اضافه می گردد. معیار استفاده شده برای تعیین اینکه کدام مدل اضافه شود، متفاوت است. در هر مرحله رو به جلو، شما یک طبقه بندی را اضافه می کنید که بهترین نتیجه را برای شما به ارمغان می آورد این یکی از دو روش معمول استفاده از رگرسیون گام به گام است. تغییر بیش از حد زمانی اتفاق می افتد که ما متغیرهای بیشتری را در اختیارمان قرار دهیم تا اینکه برای مدل مناسب باشد. به طور معمول یک اطلاعات بسیار دقیق از داده های مورد استفاده در رگرسیون را نشان می دهد، اما این مدل از نقاط داده های اضافی دور خواهد ماند و برای اینترپولسیون خوب نیست. [11]

• در این روش مدل ها بر مبنای  $F\_score$  به صورت نزولی مرتب شد، در ابتدا اولین مدل مورد ارزیابی قرار گرفت و سپس مدل دوم با مدل اول ترکیب شد و ترکیب این دو، مورد ارزیابی قرار گرفت اگر  $f\_score$  ترکیب این دو مدل بهتر از مدل اول به تنهایی باشد که ترکیب دو مدل را ذخیره می کنیم و مدل بعد را اضافه می کنیم در غیر این صورت مدل دوم حذف می شود و سومین مدل ترکیب می شود، به همین ترتیب تا آخرین مدل (100 مدل) انتخاب گردید.

با فرض داشتن سه مدل  $E1, E2, E3$  که به صورت نزولی بر مبنای  $F\_score$  مرتب شده اند جدول 2 و 3 را بدست آوردیم.

جدول 1-3- عملکرد روش انتخاب رو به جلو

Best classifier	Best $F\_score$
E2	$F(E2)$
$E2+E1$	$F(E2+E1)$

IF  $F(E2+E1) > F(E2)$  Select  $F(E2+E1)$ , else remove  $E1$

### 1-2-7-3- حذف رو به عقب<sup>2</sup>

در این روش ابتدا مدل ها را بر مبنای  $F\_score$  به صورت صعودی مرتب کردیم سپس همه ی مدل ها با هم ترکیب شد. ترکیب مدل ها ارزیابی شد و اولین مدل را حذف کردیم و نتیجه را مورد ارزیابی قرار دادیم به همین ترتیب تا آخرین مدل ادامه دادیم. [11]

<sup>1</sup> Forward selection

<sup>2</sup> Backward elimination



با فرض داشتن دو مدل، E1، E2 که به صورت صعودی بر مبنای F\_score مرتب شده اند جدول 4 را بدست آوردیم.

جدول 1-4- عملکرد روش حذف رو به عقب

Best classifier	Best F_score
E2+E1	F(E2+E1)
E1	F(E1)

remove E2 if  $F(E2+E1) < F(E1)$ , else select E2+E1

#### 4-2-7-1 ترکیب با استفاده از الگوریتم ژنتیک

الگوریتم ژنتیک، الهامی از علم ژنتیک و نظریه تکامل داروین است و بر اساس بقای برترین‌ها یا انتخاب طبیعی استوار است. یک کاربرد متداول الگوریتم ژنتیک، استفاده از آن بعنوان تابع بهینه‌کننده است. الگوریتم ژنتیک ابزار سودمندی در بازشناسی الگو، انتخاب ویژگی، درک تصویر و یادگیری ماشینی است [12]. در الگوریتم ژنتیک، نحوه تکامل ژنتیکی موجودات زنده شبیه‌سازی می‌شود.

این الگوریتم برای کاربردهای مهندسی و به صورت امروزی آن نخستین بار توسط جان هلند<sup>1</sup> متخصص علوم کامپیوتر دانشگاه میشیگان در سال 1975 پیشنهاد گردید. کار وی آغاز تمامی کوشش‌ها برای کاربرد الگوریتم ژنتیک در مهندسی است.

در یک الگوریتم ژنتیک یک جمعیت از افراد طبق مطلوبیت آنها در محیط بقا می‌یابند. افرادی با قابلیت‌های برتر، شانس ازدواج و تولید مثل بیشتری را خواهند یافت. بنابراین بعد از چند نسل فرزندی با کارایی بهتر بوجود می‌آیند. در الگوریتم ژنتیک هر فرد از جمعیت بصورت یک کروموزوم معرفی می‌شود. کروموزوم‌ها در طول چندین نسل کاملتر می‌شوند. در هر نسل کروموزوم‌ها ارزیابی می‌شوند و متناسب با ارزش خود امکان بقا و تکثیر می‌یابند. تولید نسل در بحث الگوریتم ژنتیک با عملگرهای تولید مثل<sup>2</sup> و جهش<sup>3</sup> صورت می‌گیرد. والدین برتر بر اساس یک تابع برازندگی انتخاب می‌شوند [13].

در هر مرحله از اجرای الگوریتم ژنتیکی، یک دسته از نقاط فضای جستجو مورد پردازش‌های تصادفی قرار می‌گیرند. به این صورت که به هر نقطه دنباله‌ای از کاراکترها نسبت داده می‌شود و بر روی این دنباله‌ها، عملگرهای ژنتیکی اعمال می‌شود. در آخر براساس این که تابع هدف در هر یک از نقاط چه مقدار باشد، احتمال شرکت نمودن آنها در مرحله بعد تعیین می‌گردد. [14]

این الگوریتم را می‌توان یک روش بهینه‌سازی تصادفی جهت‌دار دانست که به تدریج به سمت نقطه بهینه حرکت می‌کند. در مورد ویژگی‌های الگوریتم ژنتیک در مقایسه با دیگر روش‌های بهینه‌سازی می‌توان گفت که الگوریتمی است که بدون داشتن هیچ گونه اطلاعی از مسئله و هیچ گونه محدودیتی بر نوع متغیرهای آن برای هر گونه مسئله ای قابل اعمال است و دارای

<sup>1</sup> John Holland

<sup>2</sup> Cross over

<sup>3</sup> Mutation

کارآیی اثبات شده‌ای در یافتن بهینه کلی<sup>1</sup> می‌باشد. توانایی این روش در حل مسائل پیچیده بهینه‌سازی، است که روش‌های کلاسیک یا قابل اعمال نیستند و یا دریافتن بهینه کلی قابل اطمینان نیستند [14].

## 1-7-2-5- ساختار الگوریتم ژنتیک

به طور کلی، الگوریتم‌های ژنتیکی از اجزاء زیر تشکیل می‌شوند:

- **کروموزوم<sup>2</sup>**

در الگوریتم‌های ژنتیکی، هر کروموزوم نشان دهنده یک نقطه در فضای جستجو و یک راه حل ممکن برای مسئله مورد نظر است. خود کروموزوم‌ها (راه حل‌ها) از تعداد ثابتی ژن<sup>3</sup> (متغیر) تشکیل می‌شوند. برای نمایش کروموزوم‌ها، معمولاً از کدگذاری‌های دودویی (رشته‌های بیتی) استفاده می‌شود [14].

- **جمعیت<sup>4</sup>**

مجموعه‌ای از کروموزوم‌ها یک جمعیت را تشکیل می‌دهند. با تأثیر عملگرهای ژنتیکی بر روی هر جمعیت، جمعیت جدیدی با همان تعداد کروموزوم تشکیل می‌شود.

- **تابع برازندگی<sup>5</sup>**

به منظور حل هر مسئله با استفاده از الگوریتم‌های ژنتیکی، ابتدا باید یک تابع برازندگی برای آن مسئله ابداع شود. برای هر کروموزوم، این تابع عددی غیر منفی را برمی‌گرداند که نشان دهنده شایستگی یا توانایی فردی آن کروموزوم است.

- **عملگرهای ژنتیکی**

در الگوریتم‌های ژنتیکی، در طی مرحله تولید مثل<sup>6</sup> از عملگرهای ژنتیکی استفاده می‌شود. با تأثیر این عملگرها بر روی یک جمعیت، نسل<sup>7</sup> بعدی آن جمعیت تولید می‌شود. عملگرهای انتخاب<sup>8</sup>، آمیزش<sup>9</sup> و جهش<sup>10</sup> معمولاً بیشترین کاربرد را در الگوریتم‌های ژنتیکی دارند.

- **عملگر انتخاب**

این عملگر از بین کروموزوم‌های موجود در یک جمعیت، تعدادی کروموزوم را برای تولید مثل انتخاب می‌کند. کروموزوم‌های برآورده‌تر شانس بیشتری دارند تا برای تولید مثل انتخاب شوند.

- **عملگر آمیزش**

---

<sup>1</sup> Global Optimum

<sup>2</sup> Chromosome

<sup>3</sup> Gene

<sup>4</sup> Population

<sup>5</sup> Fitness Function

<sup>6</sup> Reproduction

<sup>7</sup> Generation

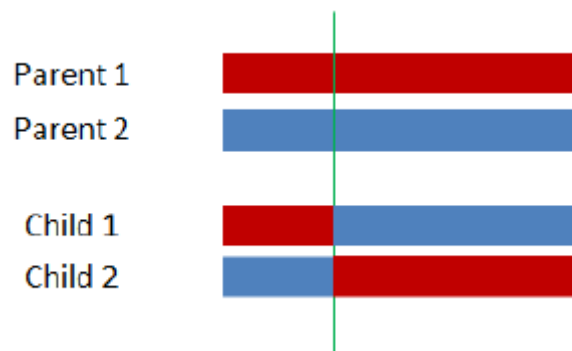
<sup>8</sup> Selection

<sup>9</sup> Crossover

<sup>10</sup> Mutation

عملگر آمیزش بر روی یک زوج کروموزوم از نسل مولد عمل کرده و یک زوج کروموزوم جدید تولید می‌کند. عملگرهای آمیزش متعددی از قبیل، آمیزش تک نقطه‌ای<sup>1</sup> و آمیزش دو نقطه‌ای<sup>2</sup> وجود دارد.

در آمیزش تک نقطه‌ای، یک موقعیت تصادفی بین دو ژن در نظر گرفته می‌شود. سپس تمامی ژن‌های طرف راست یا طرف چپ این موقعیت در کروموزوم‌های والد با یکدیگر جابجا می‌شوند تا کروموزوم‌های جدید بدست آیند. در شکل 2 آمیزش تک نقطه‌ای نشان داده شده است. در آمیزش دو نقطه‌ای، دو موقعیت به صورت تصادفی انتخاب می‌شود و تمامی ژن‌های بین این دو موقعیت در کروموزوم‌های والد با یکدیگر جابجا می‌شوند. لازم به ذکر است که آمیزش معمولاً بر روی همه زوج کروموزوم‌های انتخاب شده برای جفت‌گیری به کار برده نمی‌شود. معمولاً احتمال آمیزش برای هر زوج کروموزوم بین 0/6 تا 0/95 در نظر گرفته می‌شود که به این عدد نرخ آمیزش<sup>3</sup> یا احتمال آمیزش<sup>4</sup> گفته می‌شود و با  $P_c$  نمایش داده می‌شود. در صورتی که بر روی یک زوج کروموزوم عمل آمیزش صورت نگیرد، فرزندان با تکرار نمودن والدین تولید می‌شوند [15].



شکل 1-2- آمیزش تک نقطه‌ای [16]

#### • عملگر جهش

پس از اتمام عمل آمیزش، عملگر جهش بر روی کروموزوم‌ها اثر داده می‌شود. این عملگر یک ژن از یک کروموزوم را به طور تصادفی انتخاب نموده و سپس محتوای آن ژن را تغییر می‌دهد. اگر ژن از جنس اعداد دودویی باشد، آن را به وارونش تبدیل می‌کند و چنانچه متعلق به یک مجموعه باشد، مقدار یا عنصر دیگری از آن مجموعه را به جای آن ژن قرار می‌دهد. در شکل 3 چگونگی جهش یافتن پنجمین ژن یک کروموزوم نشان داده شده است [17].

احتمال انجام عمل جهش بر روی هر کروموزوم را نرخ جهش<sup>5</sup> یا احتمال جهش<sup>6</sup> می‌گویند و با  $P_m$  نمایش می‌دهند. معمولاً این عدد را بسیار کوچک (مثلاً 0/001) در نظر می‌گیرند. پس از اتمام عمل جهش، کروموزوم‌های تولید شده به عنوان نسل جدید شناخته شده و برای دور بعد اجرای الگوریتم ارسال می‌شوند [18].

<sup>1</sup> One-point Crossover

<sup>2</sup> Two-point Crossover

<sup>3</sup> Crossover Rate

<sup>4</sup> Crossover Probability

<sup>5</sup> Mutation Rate

<sup>6</sup> Mutation Probability

## Before Mutation

A5 

1	1	1	0	0	0
---	---	---	---	---	---

## After Mutation

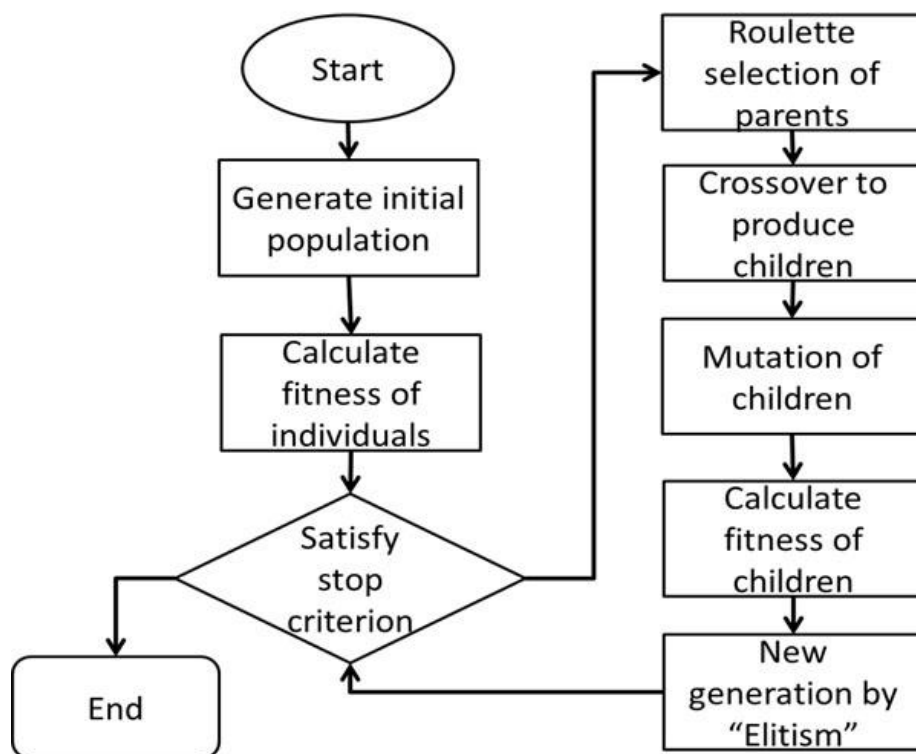
A5 

1	1	0	1	1	0
---	---	---	---	---	---

شکل 1-3- عمل جهش [19]

### 1-2-7-6- روند الگوریتم ژنتیک

در شکل 4 نمودار گردش الگوریتم‌های ژنتیکی نشان داده شده است. قبل از این که یک الگوریتم ژنتیکی بتواند اجرا شود، ابتدا باید نمایش مناسبی برای مسئله مورد نظر پیدا شود. همچنین یک تابع برازندگی نیز باید ابداع شود تا به هر راه حل کدگذاری شده ارزشی را نسبت دهد. در طی اجرا، والدین برای تولید مثل انتخاب می‌شوند و با استفاده از عملگرهای آمیزش و جهش با هم ترکیب می‌شوند تا فرزندان جدیدی تولید کنند. این فرآیند چندین بار تکرار می‌شود تا نسل بعدی جمعیت تولید شود. سپس این جمعیت بررسی می‌شود و در صورتی که ضوابط همگرایی برآورده شوند، فرآیند فوق خاتمه می‌یابد [20].



شکل 1-4- گردش کار الگوریتم ژنتیک [21]

#### 1-7-2-7- روش‌های انتخاب

روش‌های انتخاب متعددی برای استفاده در الگوریتم‌های ژنتیکی پیشنهاد شده‌اند که در ادامه این بخش، برخی از این روش معرفی می‌شوند [22].

##### • نمونه‌برداری به روش چرخ رولت

در این روش، به هر فرد قطعه‌ای<sup>1</sup> از یک چرخ رولت مدور اختصاص داده می‌شود. اندازه این قطعه متناسب با برازندگی آن فرد است. چرخ  $N$  بار چرخانده می‌شود که  $N$  تعداد افراد در جمعیت است. در هر چرخش، فرد زیر نشانگر چرخ انتخاب می‌شود و در مخزن والدین نسل بعد قرار می‌گیرد.

##### • انتخاب تورنمنت

در انتخاب تورنمنت<sup>2</sup> دو فرد از جمعیت به صورت تصادفی انتخاب می‌شوند. سپس، یک عدد تصادفی  $r$  بین  $0$  و  $1$  انتخاب می‌شود. اگر  $r < k$  (که  $k$  یک پارامتر است، برای مثال  $0/75$ ) باشد، فرد برآورده‌تر و در غیر این صورت فردی که برازندگی کمتری دارد، به عنوان والد انتخاب می‌شود. این دو سپس به جمعیت اولیه بازگردانده می‌شوند و دوباره در فرآیند انتخاب شرکت داده می‌شوند [23].

#### 1-8-2-7- شرط پایان الگوریتم

1- رسیدن به جواب بهینه

2- رسیدن به تعداد تکرار از پیش تعیین شده

#### 1-9-2-7- برخی از کاربرد الگوریتم‌های ژنتیکی

الگوریتم‌های ژنتیکی در حل بسیاری از مسائل علمی و مهندسی به کار گرفته شده‌اند. برخی از موارد کاربرد این الگوریتم‌ها عبارتند از:

بهینه‌سازی، برنامه‌سازی خودکار، یادگیری ماشین، تکامل تدریجی و یادگیری، اکولوژی و سیستم‌های اجتماعی.

#### 1-10-2-7- شرط پایان الگوریتم

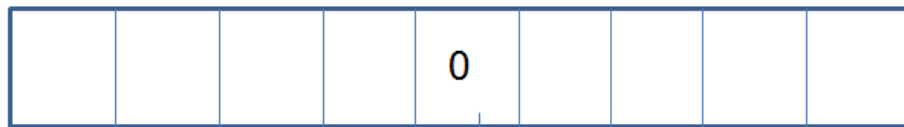
- رسیدن به تعداد تکرار مورد نظر
- بدست آوردن جواب بهینه

#### 1-11-2-7- روش کار

نگاشت مسئله به الگوریتم ژنتیک به این صورت است که یک کروموزم به طول تعداد مدل‌ها (100) ایجاد می‌کنیم، متغیرهای این بردار با مقادیر  $0$  و  $1$  پر می‌شوند، به این معنی که هر خانه‌ی این بردار متناظر با یک مدل است. مقدار  $0$  به معنی عدم حضور مدل مربوطه در ترکیب و مقدار  $1$  به معنی حضور مدل در ترکیب است. در نهایت همه‌ی مدل‌هایی که خانه‌ی متناظر با آنها با مقدار یک پر شود در ترکیب حضور دارند.

<sup>1</sup> Slice

<sup>2</sup> Tournament Selection



0 به معنی عدم حضور طبقه بند  
و 1 به معنی حضور طبقه بند می  
باشد.

شکل 1-5- بردار راه حل

جمعیت اولیه الگوریتم ژنتیک، 300 مقدار دهی شد. (سه برابر طول بردار راه حل) نرخ ضریب جهش 0.01 در نظر گرفته شد، نرخ ضریب آمیزش 0.45 مقدار دهی شد. برای انجام عمل cross over از روش انتخابی roulette wheel برای انتخاب والدین استفاده شد و ارث بری فرزندان از والدین به شیوه ی single point انجام گرفت. برای محاسبه ی تابع ارزیابی (fitness function) بعد از ترکیب مدل ها f\_score بدست آورده شد.

# فصل دوم

## نتایج

## 1-2- نتایج جداسازی مجموعه داده

جدول 1-2- جدا سازی مجموعه داده

	تعداد نمونه	تعداد ویژگی
X_train	48	7128
Y_train	48	-
X_validate	3	7128
Y_validate	3	-
X_test	18	7128
Y_test	18	-

مراحل پیش پردازش داده ها منجر به جابجایی سطر و ستون های مجموعه داده شد و در نهایت کلاس ها در ستون آخر قرار گرفت و سپس جدا سازی داده ها با توزیع یکسان به صورت تصادفی به 3 دسته ی (داده های آموزش، داده های آزمایش و داده های توسعه) تقسیم شد که از این بین 70% برای داده های آموزش، 20% داده های آزمایش و 10% داده های توسعه در نظر گرفته شد و در نهایت به برچسب AML عدد 0 و برچسب ALL عدد 1 اختصاص داده شد. تعداد نمونه ها برای داده های آموزش ، آزمایش و داده های توسعه به ترتیب 48 و 18 و 3 بدست آمد.

## 2-2- نتایج ارزیابی مدل های ساخته شده با الگوریتم های ماشین بردار پشتیبان و بیز و درخت تصمیم روی development data :

جدول 2-2- مقایسه طبقه بند ماشین بردار پشتیبان ، درخت تصمیم و بیز

	Precision	Recall	F_score
SVM	0.726	1.0	0.741
naïve byes	0.95	1.0	0.976
decision tree	0.95	1.0	0.962

با توجه به نتایج بدست آمده از ارزیابی طبقه بندهای SVM و Naïve byes و decision tree بر مبنای F\_score ، الگوریتم naïve byes به عنوان الگوریتم پایه برای ساخت 100 مدل انتخاب شد.

## 3-2- ایجاد 100 مدل با الگوریتم پایه

جدول 3-2- نتایج ارزیابی مدل های ایجاد شده با الگوریتم پایه



	precision	recall	F _ score
مدل 1 تا 92	1.0	1.0	1.0
مدل 93 تا 100	0.6666	0.6666	0.6666

100 مدل با استفاده از الگوریتم پایه ساخته شد ، سپس ارزیابی آن ها روی داده های توسعه انجام شد . که این نتایج در فایل allModels\_F-score.xls ضمیمه شده است . با توجه به تعداد کم اعضای داده ی توسعه که ارزیابی 100 مدل ایجاد شده روی آنها انجام شد ، نتایج ارزیابی از مدل 1 تا 92 کاملاً یکسان و از مدل 93 تا 100 نیز مشابه بدست آمد.

#### 4-2- مقایسه ی روش های ترکیب

جدول 4-2- نتایج روش های ترکیب

	precision	recall	F _ score
ترکیب همه مدل ها			
انتخاب رو به جلو			
حذف رو به عقب			
الگوریتم ژنتیک			
بهترین مدل			

- [1] <http://bioinformatics.upmc.edu/Help/UPITTGED.html>
- [2] <http://www-genome.wi.mit.edu/cgi-bin/cancer/datasets.cgi>
- [3] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik. "Gene selection for cancer classification using support vector machines", *Machine Learning*, 2000.
- [4] Yuh-Jye Lee and O. Mangasarian. "DT: A smooth Decision Tree for classification", *Computational Optimization and Applications*, 20, 5-22, 2001.
- [5] L. Kaufman. "Solving the quadratic programming problem arising in support vector classification", in *Nave-Bayes Classification*, B. Schölkopf, C. J. C. Burges, and A. J. Smola, eds., MIT Press, 47-167, 1999.
- [6] [http://web.stanford.edu/~hastie/CASI\\_files/DATA/leukemia.html](http://web.stanford.edu/~hastie/CASI_files/DATA/leukemia.html)
- [7] <https://medium.com/deep-math-machine-learning-ai/chapter-4decision-trees-algorithms-b93975f7a1f1>
- [8] <http://www.statsoft.com/textbook/support-vector-machines>
- [9] <http://www.statsoft.com/textbook/naive-bayes-classifier>
- [10] <https://towardsdatascience.com/ensemble-methods-in-machinelearning-what-are-they-and-why-use-them-68ec3f9fef5f>
- [11] Brusco Cradit, Steinly. An exact algorithm for hierarchically well-formulated subsets in second-order polynomial regression. *Technometrics*, 51(3):306-315, 2009
- [12] Elattar, E. E. (2015). A hybrid genetic algorithm and bacterial foraging approach for dynamic economic dispatch problem. *International Journal of Electrical Power & Energy Systems*, 69, 18–26.
- [13] Man, K.F., Tang, K.S., Kwong, S.: *Genetic Algorithms*, pp. 5–10. Springer, Concepts Des. (2000) .
- [14] Bagchi T (1999) *Multiobjective scheduling by genetic algorithms*. Kluwer, Boston
- [15] Goldberg DE (1989) *Genetic algorithms for search, optimization, and machine learning*. Addison-Wesley, Reading .
- [16] [https://www.researchgate.net/figure/Single-pointcrossover\\_fig21\\_265969600](https://www.researchgate.net/figure/Single-pointcrossover_fig21_265969600)
- [17] <http://web.cs.ucdavis.edu/~vemuri/classes/ecs271/Genetic%20Algorithms%20Short%20Tutorial.htm>
- [18] Coello CAC, Toscano G (2000) A micro-genetic algorithm for multi-objective optimization. Technical report Lania-RI-2000–06, Laboratoria Nacional de Informatica Avanzada, Xalapa, Veracruz.
- [19] Zahhad Abo, M., Ahmed Triumph No, S., Sasaki, S.: The new energy-efficient protocol for adaptive genetic algorithm to the collection month, and improve wireless sensor networks. *Int. J. Energy Inf. Commun.* 5(3), 47–72 (2014) .
- [20] Goldberg DE (1989) *Genetic algorithms for search, optimization, and machine learning*. Addison-Wesley, Reading
- [21] [https://www.researchgate.net/figure/Flowchart-of-the-Genetic-algorithm-workflow-The-first-population-was-generated-randomly\\_fig9\\_50264725](https://www.researchgate.net/figure/Flowchart-of-the-Genetic-algorithm-workflow-The-first-population-was-generated-randomly_fig9_50264725)
- [22] Deb K (1999) Solving goal programming problems using multi-objective genetic algorithms. In: *Proceedings of the CEC*, Washington, and pp 77–84 .
- [23] Deb K, Agrawal S, Pratap A, Meyarivan T (2002) a fast and elitist multi-objective genetic algorithm: NSGA-II. *IEEE Trans Evol Comput* 6:182–197