

PromoterME: یک روش شناسایی ترویج کننده ها با استفاده از شبکه عصبی ترکیبی و استخراج ویژگی از توالی DNA بر اساس شبکه عصبی پرسپترون چند لایه ای

نیلوفر آقایی ابیانه^۱ و مهدی رحیمی^۲

^۱ دپارتمان الگوریتم ها و محاسبات، دانشکده فنی، دانشگاه تهران niloofaraghaie@ut.ac.ir

^۲ دپارتمان الگوریتم ها و محاسبات، دانشکده فنی، دانشگاه تهران Mehdi.rahimi@ut.ac.ir

چکیده

شناسایی ترویج کننده^۱ ها برای تنظیمات بیان ژن^۲ بسیار مهم می باشند. اگرچه الگوریتم هایی برای شناسایی ترویج کننده ها ایجاد شده است؛ اما همچنان توسعه یک الگوریتم با دقت بالا چالشی بزرگ محسوب می شود. در این مقاله، یک الگوریتم برای شناسایی ترویج کننده ها ارائه می شود که به آن PromoterME گفته می شود. در این الگوریتم تنها از ویژگی DNA عددی شده^۳ برای شناسایی ترویج کننده ها استفاده می شود. و هر یک از چهار حرف A، C، G و T به عنوان یک ویژگی در نظر گرفته می شود. که این موضوع سبب می شود که فضای ویژگی خیلی بزرگ شود، لذا برای کاهش فضای ویژگی از شبکه عصبی پرسپترون چند لایه ای استفاده می شود. در آخر با

استفاده از روال یادگیری Adaboost عمل شناسایی ترویج کننده ها انجام می شود.

برای پیاده سازی این الگوریتم از پایگاه داده های EPD^۴ استفاده شده است.

۱. مقدمه

ترویج کننده ها قسمتی از توالی DNA می باشند که منجر به انجام عمل رونویسی^۵ از ژن می شود. هر ترویج کننده دارای مجموعه های مشخصه از توالی های کوچک حفاظت شده می باشند؛ این توالی ها همان ویژگی های ساختاری ترویج کننده ها می باشند. با این وجود، این ویژگی ها در همه ی ترویج کننده ها موجود نیست. از ویژگی های ساختاری ترویج کننده ها می توان به TATA-box که در نواحی غنی شده از AT قرار دارند، نواحی غنی شده از CpG و CAAT-box اشاره

^۱ promoter

^۲ Gene regulation

^۳ Digitized DNA

^۴ Eukaryotic Promoter Database

^۵ transcription

کرد. اگرچه بعضی از ترویج کننده ها فاقد TATA-box می باشند که به آنها ترویج کننده های فاقد TATA-box^۱ گفته می شود. از این ویژگی ها می توان برای شناسایی نواحی ترویج کننده ها و غیر ترویج کننده ها استفاده کرد [۱][۲][۳][۴].

الگوریتم های متعددی برای شناسایی نواحی ترویج کننده ها ارائه شده است مثلاً *Dragon Promoter Finder (DPF)*

PromoterExplorer

PromoterInspector

از یک دیدگاه می توان این الگوریتم ها را به دو دسته ی بر اساس ویژگی های مفهومی مانند الگوریتم *PromoterInspector* و بر اساس ویژگی های ساختار های مانند *PromoterExplorer* تقسیم بندی کرد [۵].

هر یک از این الگوریتم ها از شناسایی الگو های^۲ مختلفی برای این مسئله طبقه بندی و همچنین از ویژگی های مختلفی استفاده می کنند، مثلاً می توان به شبکه عصبی، تحلیل های خطی و مربعی^۳، مدل مارکوف^۴ و تحلیل مولفه مستقل^۵ اشاره کرد. [۶]. در DPF هر توالی از پنج نوکلئوتید پشت سر هم به عنوان فضای ویژگی در نظر گرفته می شود [۷]. در

PromoterInspector دو مجموعه از الیگونوکلئوتیدها به ترتیب مرتبط با ترویج کننده ها و غیر ترویج کننده ها در نظر گرفته می شود؛ همچنین جایگزین های متعددی در مکان های چندگانه معرفی می شوند. تعداد جایگزین ها و طول عناصر در الیگونوکلئوتیدها برای طبقه بندی بهینه می شوند و سپس برای تست کردن مورد استفاده قرار می گیرند [۸]. در *PromoterExplorer* توزیع محلی توالی های پنج تایی از نوکلئوتید ها، نواحی CpG و DNA عددی شده به عنوان فضای ویژگی در نظر گرفته می شوند. سپس با کمک الگوریتم طبقه بندی *AdaBoost* عمل طبقه بندی انجام می شود [۶][۹].

در اینجا کارایی الگوریتم ارائه شده در این مقاله روی پایگاه داده EPD^۶ محاسبه می شود. در روش ارائه شده در این مقاله تنها از ویژگی DNA عددی شده استفاده می شود و هر یک از باز های A, C, G و T و به عنوان یک ویژگی در نظر گرفته می شوند؛ که این به نوبه ی خود منجر به یک فضای ویژگی خیلی بزرگ می شود. برای حل این مشکل، با کمک شبکه عصبی چند لایه ای ابتدا فضای ویژگی کاهش داده می شود؛ سپس با استفاده از مدل ترکیبی طبقه بندی این مسئله دو کلاسه انجام داده می شود. در این مقاله از مدل های ترکیبی

^۱ TATA-box less promoter

^۲ Pattern recognition

^۳ Linear and quadratic discriminant analyses

^۴ Markov model

^۵ Independent component analyses

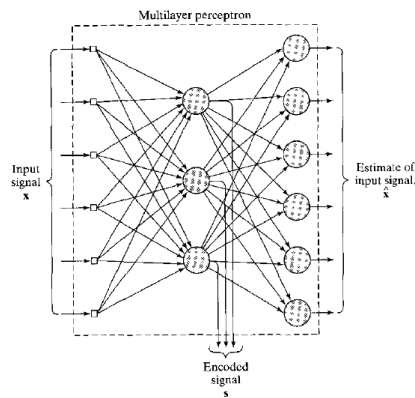
^۶ Eukaryotic Promoter Database (EPD)

۲. استخراج ویژگی از توالی DNA

همانطور که گفته شد در این مقاله از ویژگی DNA عددی شده استفاده شده است. در اینجا هر یک از باز های A, C, G و T به ترتیب معادل با ۰۰، ۰۱، ۱۰ و ۱۱ قرار داده شده، و نهایتاً به عنوان یک ویژگی در نظر گرفته می شوند [۱۰]. بنابراین اگر l طول توالی DNA باشد، فضای ویژگی 2^l می شود. برای کاهش فضای ویژگی از شبکه عصبی پرسپترون سه لایه ای استفاده می شود. برای این کار در این شبکه عصبی تعداد نرون های لایه ورودی و خروجی برابر با $2l$ قرار می گیرند و تعداد نرون های لایه میانی برابر با طول بردار ویژگی مورد انتظار در نظر گرفته می شوند. شکل ۱ این موضوع را نشان می دهد. به ازای هر نمونه آموزش بردار ورودی و خروجی یکسان قرار داده می شوند و وزن های این شبکه عصبی به صورت مد الگو^۱ طبق قاعده ی دلتا به روز می شوند.

$$w(n+1) = w(n) + \eta \cdot \varepsilon \cdot f' \cdot X$$

که $w(n)$ وزن ها در لحظه ی n ام، η نرخ یادگیری، ε خطا، f' مشتق تابع f و X بردار ورودی می باشند، فرآیند به روز رسانی وزن ها تا زمانی که خطای مجذور مربعات^۲ به مقداری مشخص برسد، ادامه می یابد. با این کار با داشتن وزن های لایه ورودی به میانی (با فرض اینکه شبکه عصبی ۳ لایه ای باشد) می توان برای هر نمونه با پیش پردازش روی نمونه، ابعاد آن را به مقدار مورد نظر کاهش داد [۱۱].



شکل ۱ - ساختار شبکه عصبی پرسپترون سه لایه برای استخراج ویژگی

در اینجا تعداد نرون های لایه میانی برابر با ۸ قرار میگیرد. در نتیجه پس از اینکه توالی DNA ورودی به صورت عددی شد فضای ویژگی برابر با 2^l می شود که با این عمل استخراج به ۸ کاهش می یابد.

^۲ Minimum Square Error (MSE)

^۱ Pattern mode

۳. آموزش طبقه بند ها برای عمل

طبقه بندی

Boosting یک الگوریتم یادگیری با نظارت^۱ است که برای بهبود صحت تصمیم گیری هر الگوریتم یادگیری به کار می رود. در این الگوریتم، یک طبقه بند که دقتش بر اساس مجموعه آموزش است بزرگتر از کارایی در یک حالت متوسط ایجاد می شود، و سپس طبقه بند های جدید به منظور ایجاد یک شکل ترکیبی اضافه می شوند؛ به طوریکه حاصل از این طبقه بندها دارای دقت بالایی باشد [۱۲][۱۳]. در این حالت گفته می شود که کارایی طبقه بندی ارتقا^۲ یافته است. به طور کلی آموزش های الگوریتم های طبقه بند پشت سر هم به همراه زیر مجموعه ای از کل داده آموزش که از همه با اهمیت تر است به مجموعه ای از طبقه بند های فعلی داده می شود.

Adaboost یک الگوریتم boosting است که یک یادگیرنده ی ضعیف^۳ را روی داده آموزش که هر بار کمی تغییر کرده است چندین بار اجرا می کند، و سپس فرضیات را به منظور به دست آوردن دقت بالاتری از فرضیاتی که یادگیرنده های ضعیف دارند با هم ترکیب می

کند. ایده اصلی *Adaboost* به این صورت است که هر نمونه ای از مجموعه آموزش برای تمایز در گام های آموزش متفاوت نقش متفاوتی را ایفا کنند. این نمونه ها که به آسانی شناسایی می شوند باید در آموزش به صورت زیر کمتر در نظر گرفته شوند؛ تا زمانی که نمونه ای اشتباه طبقه بندی می شود باید در مراحل بعدی به آن بیشتر توجه شود. برای این منظور یادگیرنده ضعیف باید روی نمونه های با اهمیت تر یا "سخت تر" تمرکز کند. اهمیت نمونه ها با وزن نشان داده می شود. در ابتدا همه وزن ها برابرند. سپس وزن ها هر بار باید بر اساس نتایج طبقه بندی تغییر کنند. تصمیم گیری نهایی ترکیبی از تصمیمات در هر مرحله است [۱۲] [۱۳].

اگر الگوریتم *AdaBoost* را به صورت یک مدل جمعی عام^۴ در نظر گرفته شود و سپس تابع هزینه ی^۵ رگرسیون لگاریتمی^۶ بر آن اعمال شود، الگوریتم *LogitBoost* حاصل می شود. به عبارت دیگر *LogitBoost* همان الگوریتم *AdaBoost* که تابع هزینه ی آن رگرسیون لگاریتمی باشد.

GentleBoost نیز نسخه ای تغییر یافته از *AdaBoost* است [۱۳].

^۴ generalized additive model

^۵ cost function

^۶ logistic regression

^۱ Supervised Learning Algorithm

^۲ boosted

^۳ Weak learner

۴. نتایج آزمایشگاهی

پس از فاز استخراج ویژگی یک مجموعه داده ۱۰ بعدی ایجاد می شود. در ادامه این مجموعه داده ۱۰ بعدی برای طبقه بند ترکیبی به کار می رود.

در این آزمایش ها به تعداد ۱۰۰۰ نمونه از توالی ترویج کننده ها از پایگاه داده EPD و ۱۰۰۰ نمونه از توالی غیر ترویج کننده ها انتخاب شده است، که ابتدا این ۲۰۰۰ نمونه در فاز استخراج ویژگی به یک مجموعه داده ای ۲۰۰۰ نمونه ای با ابعاد ۸ تبدیل می شود. پس از این فاز ویژگی ها به صورت عددی می باشند. از طرفی در هر آزمایش به میزان ۰/۸ از کل مجموعه داده ای برای آموزش و ۰/۲ باقیمانده برای تست و ارزیابی کارایی برنامه به کار می رود. اگرچه در هر آزمایش میزان داده های آموزش و تست ثابت در نظر گرفته می شود، اما این نمونه ها به صورت تصادفی انتخاب می شوند که بنابراین منجر به نتایج مختلف می شود. در این آزمایش ها *learners* نوع *Tree* و سه نوع *LogitBoost* *AdaBoostM1* *method* و *GentleBoost* مورد آزمایش قرار گرفته اند. همچنین هر آزمایش برای ۳ مقدار نرخ یادگیری ۰،۱، ۰،۵ و ۰،۹ روی داده های آموزش و تست ثابت بررسی می شود.

در هر آزمایش خطای جایگذاری^۱ بررسی می شود. تخمین این خطا براساس تفاوت بین مقادیر پیش بینی شده از یک مدل آزمایش و مقادیر مشاهده شده در آزمایش ها اندازه گیری می شود. همچنین معیار های حساسیت^۲ (S_n) و خاصیت (S_p) محاسبه می شوند. این معیار ها برابرند با:

$$S_n = \frac{TP}{TP + FN}$$

$$S_p = \frac{TP}{TP + Fp}$$

که TP ، شناسایی نواحی به عنوان ترویج کننده هایی که واقعاً ترویج کننده میباشند؛ FN ، شناسایی نواحی به عنوان غیر ترویج کننده هایی که واقعاً ترویج کننده میباشند؛ FP ، شناسایی نواحی به عنوان ترویج کننده هایی که واقعاً ترویج کننده نمی باشند.

در این آزمایش ها *learners* نوع *Tree* و سه نوع *LogitBoost* *AdaBoostM1* *method* و *GentleBoost* مورد آزمایش قرار گرفته اند.

برای شبیه سازی این طبقه بندی ترکیبی از تابع *fitensemble* در محیط متلب استفاده شده است. این تابع یک مدل ترکیبی را برای پیش بینی پاسخ به داده را می سازد. ساختار نحوی این تابع در محیط متلب به صورت زیر است.

$Ens = fitensemble(X, Y, method, nlearn, learners)$

^۱ Resubstitution Error

^۲ sensitivity

و پارامترهای تابع *fitensemble* به صورت زیر است. این ماشین بر اساس مدل های لیست شده در پارامتر *learners* ایجاد می شود.

توصیف پارامترهای اصلی

X: ماتریسی است که هر ستون آن مشخص کننده ی یک متغیر (ویژگی) است و هر سطر آن یک نمونه از مجموعه داده ای را نشان می دهد.

Y: در مسائل طبقه بندی مشخص کننده ی یک بردار ستونی است که یک کلاس را به ازای هر نمونه از مجموعه داده ای مشخص می کند.

method: یک رشته ی حساس به کوچک بزرگی حروف است که یکی از مقادیر زیر را بسته به تعداد کلاس ها می گیرد:

برای طبقه بندی دو کلاسه

- *AdaBoostM1*
- *LogitBoost*
- *GentleBoost*
- *RobustBoost*
- *Bag*
- *Subspace*

برای طبقه بندی سه کلاسه و بیشتر

- *AdaBoostM2*
- *Bag*
- *Subspace*

nlearn: یک عدد صحیح مثبت است که تعداد چرخه های یادگیری^۱ ترکیبی را نشان می دهد. در هر چرخه ی یادگیری تابع *fitensemble* بر روی همه ی الگوهای یادگیری لیست شده در پارامتر عبور می کند و یک یادگیر ضعیف را برای هر الگو آموزش می دهد. تعداد کل یادگیر های آموزش دیده برابر حاصل ضرب *nlearn* در تعداد *learners* است.

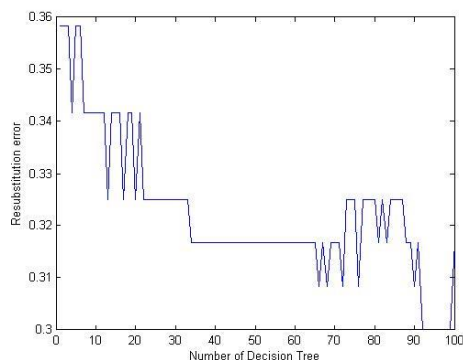
S: مشخص کننده ی نوع یادگیرهای ضعیف است. سه نوع از این یادگیرها که ما نیز در آزمایش ها از آن استفاده می شود. عبارت اند از

- *Discriminant*
- *KNN*
- *Tree*

در عین حال تعدادی پارامتر اختیاری ممکن برای همه و در برخی موارد برای انواع خاصی از *method* ها وجود دارد. یدر اینجا یکی از پارامترهایی که مورد آزمایش قرار گرفته *learning rate* است. در محیط متلب به صورت پیش فرض مقدار آن ۱ است [۱۵] [۱۴].

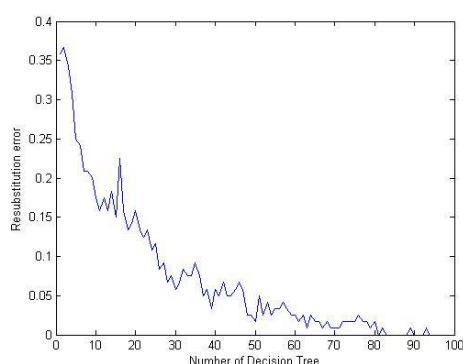
در این مقاله آزمایش های انجام شده به صورت زیر است.

^۱ Learning cycle



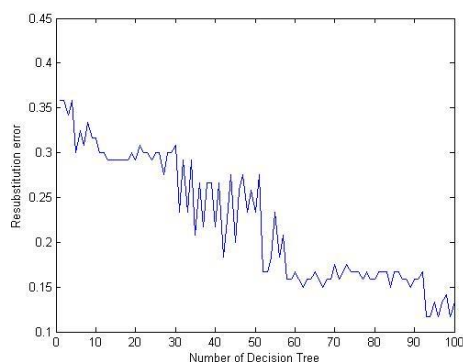
شکل ۳- نمودار خطای جایگزینی به ازای تعداد درخت ها در

آزمایش اول با $LR = 0.1$ و روش **LogitBoost**



شکل ۴- نمودار خطای جایگزینی به ازای تعداد درخت ها در

آزمایش اول با $LR = 0.1$ و روش **GentleBoost**



شکل ۵- نمودار خطای جایگزینی به ازای تعداد درخت ها در

آزمایش اول با $LR = 0.5$ و روش **AdaBoostM1**

	$LR = 0.1$	$LR = 0.5$	$LR = 0.9$
AdaBoostM1	۵۶/۶۶۶۷	۴۶/۶۶۶۷	۳۶/۶۶۶۷
LogitBoost	۴۶/۶۶۶۷	۴۳/۳۳۳۳	۴۳/۳۳۳۳
GentleBoost	۴۰	۴۰	۴۰

جدول ۱- مقادیر کارایی در آزمایش اول به ازای **LR** های ۰/۱،

۰/۵ و ۰/۹ با روش های **AdaBoostM1**، **LogitBoost**

GentleBoost

	$LR = 0.1$	$LR = 0.5$	$LR = 0.9$
AdaBoostM1	۰/۷۰۸۳	۰/۷۰۰۰	۰/۶۴۷۱
LogitBoost	۰/۷۰۰۰	۰/۶۸۴۲	۰/۶۸۴۲
GentleBoost	۰/۵۲۱۷	۰/۵۲۱۷	۰/۶۳۱۶

جدول ۲- مقادیر S_n در آزمایش اول به ازای **LR** های ۰/۱، ۰/۵

و ۰/۹ با روش های **AdaBoostM1**، **LogitBoost**

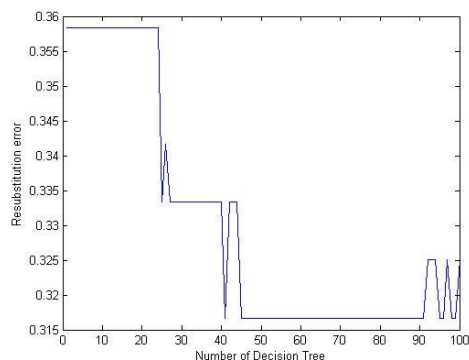
GentleBoost

	$LR = 0.1$	$LR = 0.5$	$LR = 0.9$
AdaBoostM1	۰/۷۳۹۱	۰/۵۸۳۳	۰/۴۵۸۳
LogitBoost	۰/۵۸۳۳	۰/۵۴۱۷	۰/۵۴۱۷
GentleBoost	۰/۶۳۱۶	۰/۶۳۱۶	۰/۵۲۱۷

جدول ۳- مقادیر S_p در آزمایش اول به ازای **LR** های ۰/۱، ۰/۵

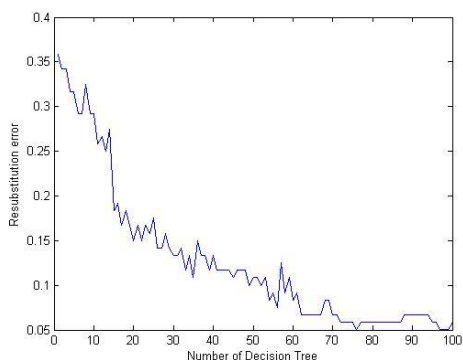
و ۰/۹ با روش های **AdaBoostM1**، **LogitBoost**

GentleBoost

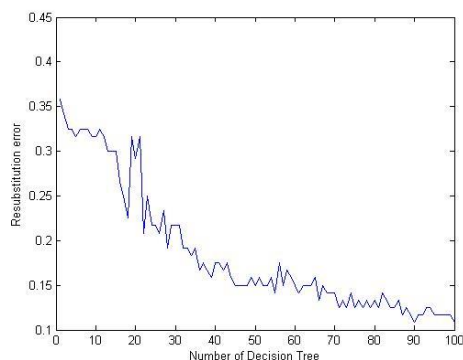


شکل ۲- نمودار خطای جایگزینی به ازای تعداد درخت ها در

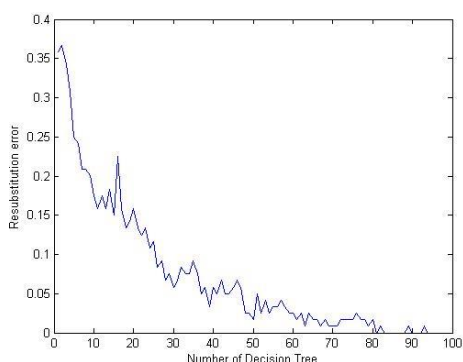
آزمایش اول با $LR = 0.1$ و روش **AdaBoostM1**



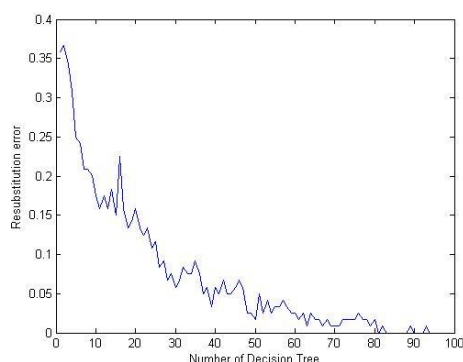
شکل ۹- نمودار خطای جایگزینی به ازای تعداد درخت ها در آزمایش اول با $LR = 0.9$ و روش LogitBoost



شکل ۶- نمودار خطای جایگزینی به ازای تعداد درخت ها در آزمایش اول با $LR = 0.5$ و روش LogitBoost



شکل ۱۰- نمودار خطای جایگزینی به ازای تعداد درخت ها در آزمایش اول با $LR = 0.1$ و روش LogitBoost



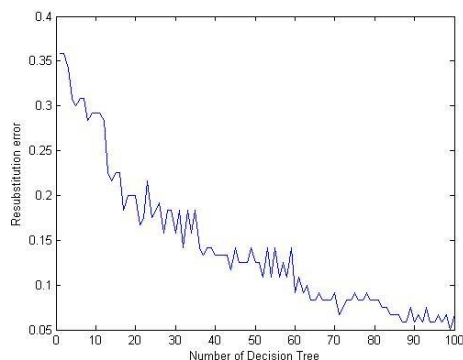
شکل ۷- نمودار خطای جایگزینی به ازای تعداد درخت ها در آزمایش اول با $LR = 0.05$ و روش LogitBoost

	$LR = 0.1$	$LR = 0.05$	$LR = 0.9$
AdaBoostM1	۵۰	۳۶/۶۶۶۷	۳۳/۳۳۳۳
LogitBoost	۵۰	۴۰	۳۳/۳۳۳۳
GentleBoost	۳۰	۳۰	۳۰

جدول ۴- مقادیر کارایی در آزمایش دوم به ازای LR های ۰/۱،

۰/۵ و ۰/۹ با روش های AdaBoostM1، LogitBoost،

GentleBoost



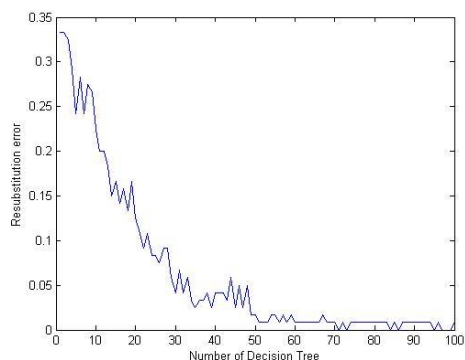
شکل ۸- نمودار خطای جایگزینی به ازای تعداد درخت ها در آزمایش اول با $LR = 0.9$ و روش AdaBoostM1

	$LR = 0.1$	$LR = 0.05$	$LR = 0.9$
AdaBoostM1	۰/۵۷۶۹	۰/۵۰۰۰	۰/۵۲۶۳
LogitBoost	۰/۵۷۶۹	۰/۵۴۵۵	۰/۵۰۰۰
GentleBoost	۰/۴۷۳۷	۰/۴۷۳۷	۰/۴۷۳۷

جدول ۵- مقادیر S_n در آزمایش دوم به ازای LR های ۰/۱، ۰/۵،

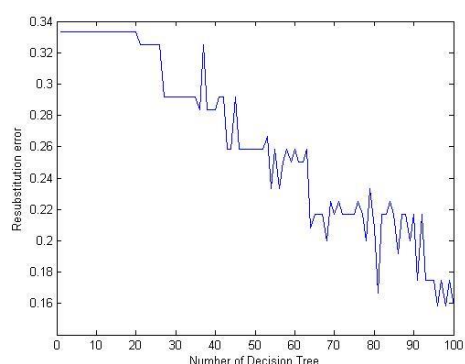
و ۰/۹ با روش های AdaBoostM1، LogitBoost،

GentleBoost



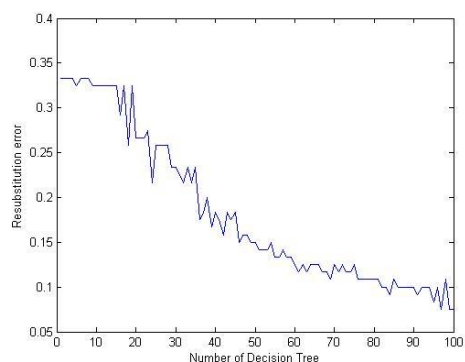
شکل ۱۳- نمودار خطای جایگزینی به ازای تعداد درخت ها در

آزمایش دوم با $LR = 0.1$ و روش **GentleBoost**



شکل ۱۴- نمودار خطای جایگزینی به ازای تعداد درخت ها در

آزمایش دوم با $LR = 0.5$ و روش **AdaBoostM1**



شکل ۱۵- نمودار خطای جایگزینی به ازای تعداد درخت ها در

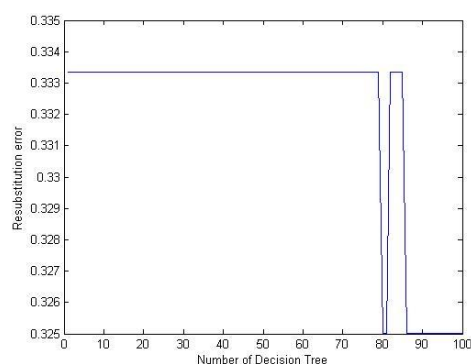
آزمایش دوم با $LR = 0.5$ و روش **LogitBoost1**

	$LR = 0.1$	$LR = 0.5$	$LR = 0.9$
AdaBoostM1	0.7391	0.5833	0.4583
LogitBoost	0.5833	0.5417	0.5417
GentleBoost	0.6316	0.6316	0.5217

جدول ۶- مقادیر S_p در آزمایش دوم به ازای LR های 0.1، 0.5،

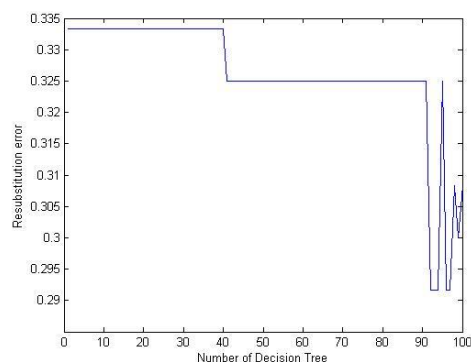
و 0.9 با روش های **LogitBoost**، **AdaBoostM1** و

GentleBoost



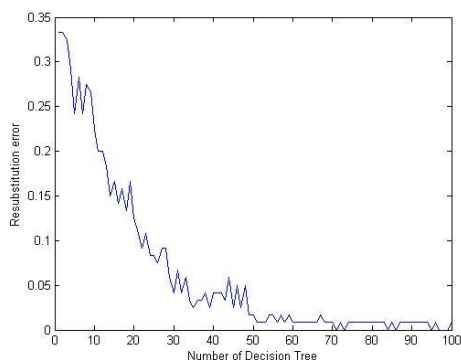
شکل ۱۱- نمودار خطای جایگزینی به ازای تعداد درخت ها در

آزمایش دوم با $LR = 0.1$ و روش **AdaBoostM1**



شکل ۱۲- نمودار خطای جایگزینی به ازای تعداد درخت ها در

آزمایش دوم با $LR = 0.1$ و روش **LogitBoost1**



شکل ۱۹ - نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش دوم با $LR = 0.9$ و روش GentleBoost

	$LR = 0.1$	$LR = 0.5$	$LR = 0.9$
AdaBoostM1	۵۰	۴۶/۶۶۶۷	۵۳/۳۳۳۳
LogitBoost	۵۰	۴۶/۶۶۶۷	۴۳/۳۳۳۳
GentleBoost	۵۰	۵۰	۵۰

جدول ۷ - مقادیر کارایی در آزمایش سوم به ازای LR های ۰/۱،

۰/۵ و ۰/۹ با روش های AdaBoostM1، LogitBoost،

GentleBoost

	$LR = 0.1$	$LR = 0.5$	$LR = 0.9$
AdaBoostM1	۰/۷۵۰۰	۰/۷۳۶۸	۰/۸۰۰۰
LogitBoost	۰/۷۵۰۰	۰/۷۷۷۸	۰/۷۶۴۷
GentleBoost	۰/۷۵۰۰	۰/۶۰۰۰	۰/۶۰۰۰

جدول ۸ - مقادیر S_n در آزمایش سوم به ازای LR های ۰/۱،

۰/۵ و ۰/۹ با روش های AdaBoostM1، LogitBoost،

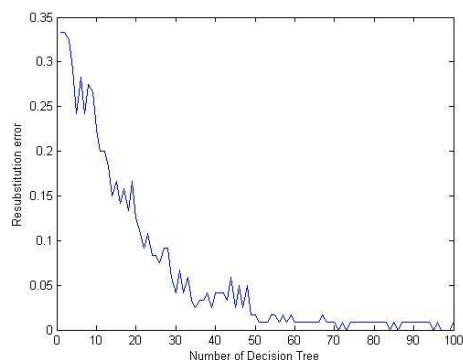
GentleBoost

	$LR = 0.1$	$LR = 0.5$	$LR = 0.9$
AdaBoostM1	۰/۶۰۰۰	۰/۵۶۰۰	۰/۶۱۵۴
LogitBoost	۰/۶۰۰۰	۰/۵۳۸۵	۰/۵۰۰۰
GentleBoost	۰/۶۰۰۰	۰/۷۵۰۰	۰/۷۵۰۰

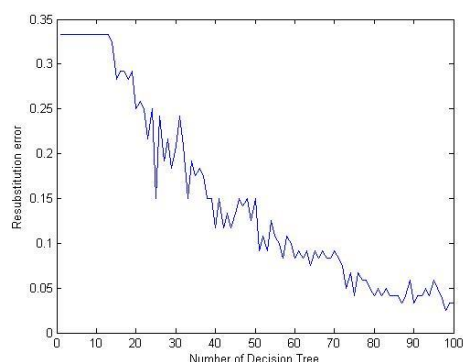
جدول ۹ - مقادیر S_p در آزمایش سوم به ازای LR های ۰/۱،

۰/۵ و ۰/۹ با روش های AdaBoostM1، LogitBoost،

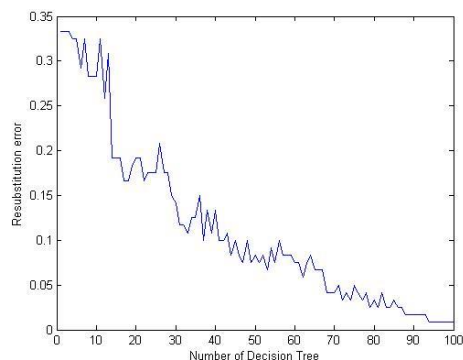
GentleBoost



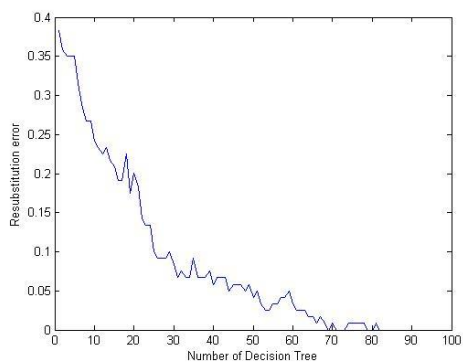
شکل ۱۶ - نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش دوم با $LR = 0.5$ و روش GentleBoost



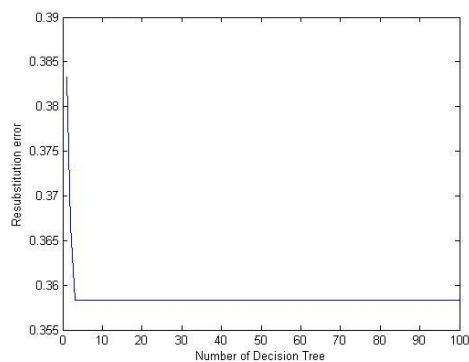
شکل ۱۷ - نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش دوم با $LR = 0.9$ و روش AdaBoostM1



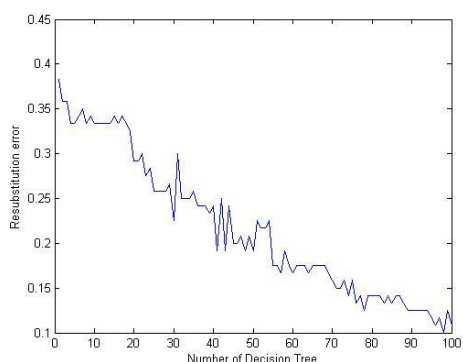
شکل ۱۸ - نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش دوم با $LR = 0.9$ و روش LogitBoost



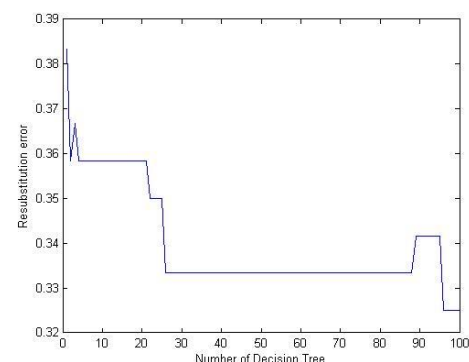
شکل ۲۳- نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش سوم با $LR = 0.5$ و روش AdaBoostM1



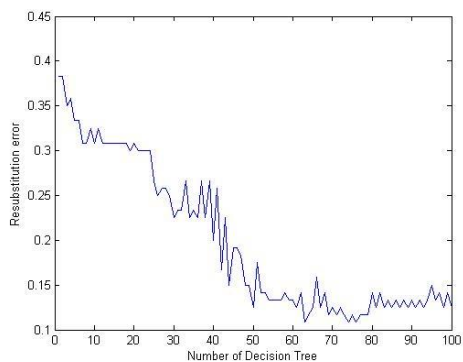
شکل ۲۰- نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش سوم با $LR = 0.1$ و روش AdaBoostM1



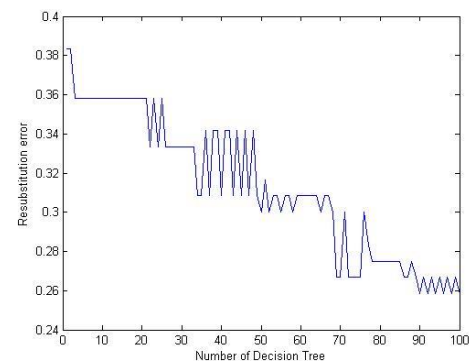
شکل ۲۴- نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش سوم با $LR = 0.5$ و روش LogitBoostM1



شکل ۲۱- نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش سوم با $LR = 0.1$ و روش LogitBoostM1



شکل ۲۵- نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش سوم با $LR = 0.5$ و روش GentleBoostM1



شکل ۲۲- نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش سوم با $LR = 0.1$ و روش GentleBoostM1

جدول ۱۰ - مقادیر کارایی در آزمایش چهارم به ازای LR های

۰/۱، ۰/۵ و ۰/۹ با روش های AdaBoostM۱، LogitBoost،

GentleBoost

	LR = ۰/۱	LR = ۰/۵	LR = ۰/۹
AdaBoostM۱	۰/۷۳۳۳	۰/۷۹۱۷	۰/۷۳۹۱
LogitBoost	۰/۷۸۵۷	۰/۸۳۳۳	۰/۷۲۷۳
GentleBoost	۰/۶۶۶۷	۰/۶۶۶۷	۰/۶۶۶۷

جدول ۱۱ - مقادیر S_H در آزمایش چهارم به ازای LR های ۰/۱،

۰/۵ و ۰/۹ با روش های AdaBoostM۱، LogitBoost،

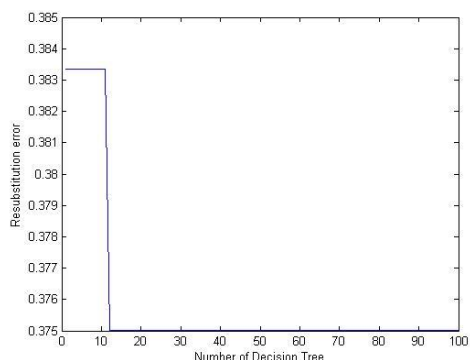
GentleBoost

	LR = ۰/۱	LR = ۰/۵	LR = ۰/۹
AdaBoostM۱	۱	۰/۷۶۰۰	۰/۷۰۸۳
LogitBoost	۰/۹۱۶۷	۰/۷۶۹۲	۰/۶۶۶۷
GentleBoost	۰/۶۰۸۷	۰/۶۰۸۷	۰/۶۰۸۷

جدول ۱۲ - مقادیر S_p در آزمایش چهارم به ازای LR های ۰/۱،

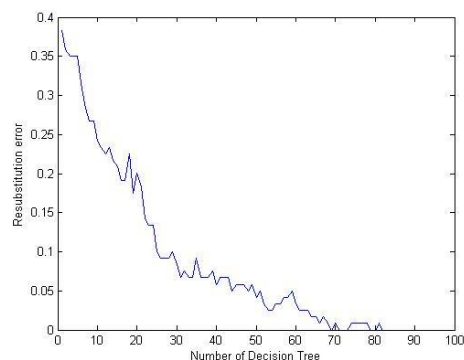
۰/۵ و ۰/۹ با روش های AdaBoostM۱، LogitBoost،

GentleBoost



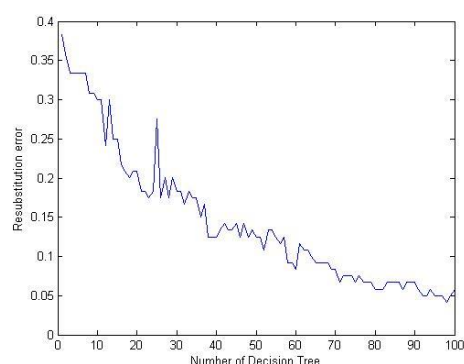
شکل ۲۹ - نمودار خطای جایگزینی به ازای تعداد درخت ها در

آزمایش سوم با LR = ۰/۱ و روش AdaBoostM۱



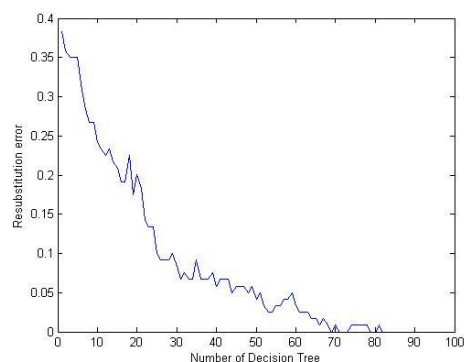
شکل ۲۶ - نمودار خطای جایگزینی به ازای تعداد درخت ها در

آزمایش سوم با LR = ۰/۹ و روش AdaBoostM۱



شکل ۲۷ - نمودار خطای جایگزینی به ازای تعداد درخت ها در

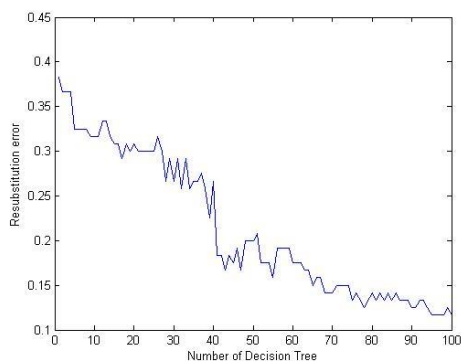
آزمایش سوم با LR = ۰/۹ و روش LogitBoost



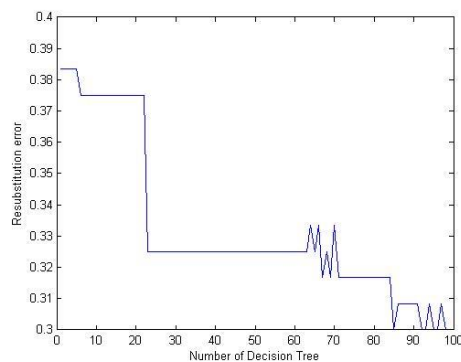
شکل ۲۸ - نمودار خطای جایگزینی به ازای تعداد درخت ها در

آزمایش سوم با LR = ۰/۹ و روش GentleBoost

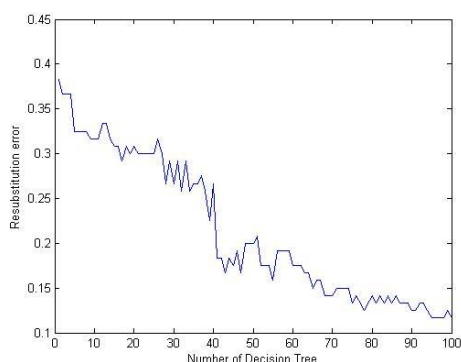
	LR = ۰/۱	LR = ۰/۵	LR = ۰/۹
AdaBoostM۱	۷۳/۳۳۳۳	۶۳/۳۳۳۳	۵۶/۶۶۶۷
LogitBoost	۷۳/۳۳۳۳	۶۶/۶۶۶۷	۵۳/۳۳۳۳
GentleBoost	۴۶/۶۶۶۷	۴۶/۶۶۶۷	۴۶/۶۶۶۷



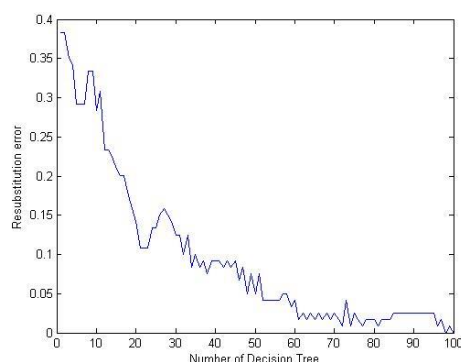
شکل ۳۳- نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش چهارم با $LR = 0.5$ و روش LogitBoost



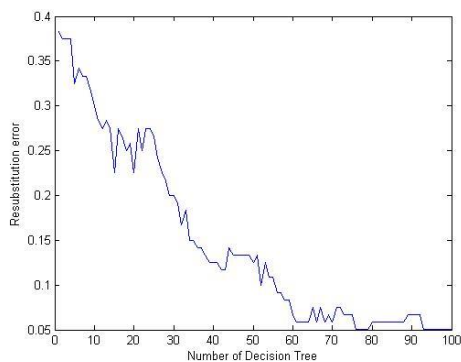
شکل ۳۰- نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش چهارم با $LR = 0.1$ و روش LogitBoost



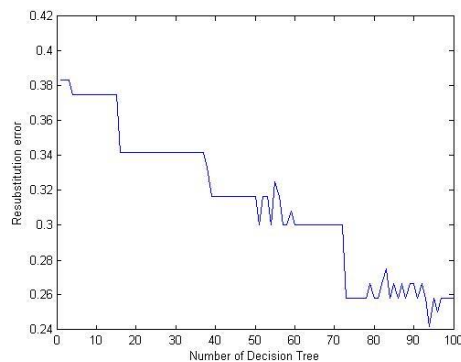
شکل ۳۴- نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش چهارم با $LR = 0.5$ و روش GentleBoost



شکل ۳۱- نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش چهارم با $LR = 0.1$ و روش GentleBoost



شکل ۳۵- نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش چهارم با $LR = 0.9$ و روش AdaBoostM1



شکل ۳۲- نمودار خطای جایگزینی به ازای تعداد درخت ها در
آزمایش چهارم با $LR = 0.5$ و روش AdaBoostM1

بیشتر بخصوص در قسمت استخراج ویژگی است. در واقع این مقاله تنها عرفی از الگوریتمی که به زودی نسخه ی نهایی آن منتشر خواهد شد می باشد. این الگوریتم نیاز به بهینه سازی در قسمت های مختلف از جمله در استخراج اطلاعات، تعیین پارامتر های بهینه برای شبکه های عصبی و... است. که به زودی انجام خواهد شد.

۶. منابع

[۱] Benjamin Lewin(۲۰۰۷), Genes, edition ۹, Jones and Bartlett Publishers.

[۲] James W. Fickett and Artemis G. Hatzigeorgiou(۱۹۹۷) Eukaryotic Promoter Recognition, Genome Res ۷: ۸۶۱-۸۷۸.

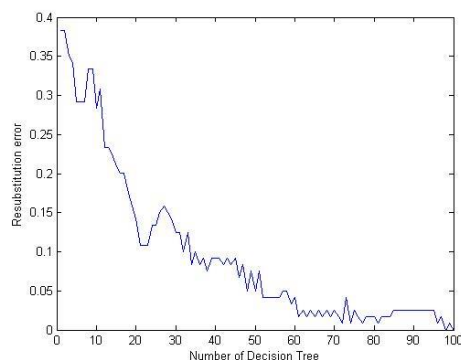
[۳] Anders Gorm Pedersen, Pierre Baldi, Yves Chauvin, and Søren Brunak(۱۹۹۹), The Biology of Eukaryotic Promoter Prediction—a Review.

[۴] Aditi Kanhere and Manju Bansal(۲۰۰۵), A novel method for prokaryotic promoter prediction based on DNA stability, BMC Bioinformatics ۲۰۰۵, ۶:۱
doi:۱۰.۱۱۸۶/۱۴۷۱-۲۱۰۵-۶-۱

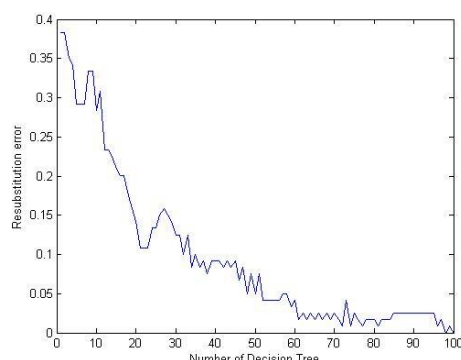
[۵] Jia Zeng, Shanfeng Zhu and Hong Yan(۲۰۰۹), Towards accurate human promoter recognition: a review of currently used sequence features and classification methods, BRIEFINGS IN BIOINFORMATICS. VOL ۱۰. NO ۵. ۴۹۸-۵۰۸

[۶] Xudong Xie, Shuanhu Wu, Kin-Man Lam and Hong Yan(۲۰۰۶), PromoterExplorer: an effective promoter identification method based on the AdaBoost algorithm. bioinformatic, Vol. ۲۲ no. ۲۲, ۲۷۲۲-۲۷۲۸

[۷] Vladimir B. Bajic, Seng Hong Seah, Allen Chong, Guaglan Zhang, Judice L.Y. Koh and Vladimir



شکل ۳۶- نمودار خطای جایگزینی به ازای تعداد درخت ها در آزمایش چهارم با $LR = ۰/۹$ و روش LogitBoost



شکل ۳۷- نمودار خطای جایگزینی به ازای تعداد درخت ها در آزمایش چهارم با $LR = ۰/۹$ و روش GentleBoost

۵. نتیجه گیری

با توجه به نتایج حاصل از چهار بار آزمایش که در جداول ۱ تا ۱۲ و شکل های ۲ تا ۳۷ آمده است. مشخص می شود که الگوریتم PromoterME نسبت به الگوریتم های Dragon Promoter Finder و PromoterInspector و PromoterExplorer از کارایی پایین خوردار است. که البته این نسخه نهایی آن نمی باشد و نیاز به بررسی و انجام آزمایش های

Brusic(2002), Dragon Promoter Finder:recognition of
vertebrate RNA polymerase II
Promoter,bioinformatics,vol 18,no. 1198-199

[8][http://www.genomatix.de/online_help/help_gems/
PromoterInspector_help.html](http://www.genomatix.de/online_help/help_gems/PromoterInspector_help.html)

[9] Xiaomeng Li, Jia Zeng, and Hong Yan(2008),
PCA-HPR: A principle component analysis model
for human promoter recognition, Bioinformation by
Biomedical Informatics Publishing Group.

[10] Pierre Baldi,Søren Brunak(2001), Bioinformatics
The Machine Learning Approach,second edition.
Massachusetts Institute of Technology.

[11] Ethem Alpaydin(2001),Introduction to machine
learning,second edition, The MIT Press Cambridge,
Massachusetts London, England.

[12] Robi Polikar(2006),Ensemble Based Systems in
Decision Making, 1531-636X/06/0201...©2006 IEEE

[13] Artur Ferreira(2007), Survey on Boosting
Algorithms for Supervised and Semi-supervised
Learning, Instituto de Telecomunicações

[14]<http://www.mathworks.com/help/stats/fitensemble.html>

[15]<http://www.mathworks.com/help/stats/classificationensemble.html>