



## تئوری بازی در تشخیص هرزنامه

سمانه قدس

استادیار، گروه علوم پایه، واحد سمنان، دانشگاه آزاد اسلامی، سمنان، ایران

s1ghods@gmail.com

### چکیده

امروزه در بسیاری از سیستم‌هایی که عمل طبقه‌بندی بر روی آن صورت می‌پذیرد، رقیب تغییراتی را بروی داده‌ها، به منظور کاهش دقت کلاسه‌بندی کننده انجام می‌دهند. مهمترین مثال از این نوع، تشخیص هرزنامه الکترونیکی می‌باشد. اسپم‌ها بطور معمول بروی هرزنامه‌ها تغییراتی را صورت می‌دهند تا فیلتر، آنها را بصورت نامه‌های درست تشخیص دهد. در این مقاله رفتار بین رقیب (اسپم‌ها) و فیلتر هرزنامه (کلاسه‌بندی کننده) بعنوان دو بازیکن، برای بازی پی در پی بررسی می‌گردد. در این روش فیلترها با یادگیری رفتار رقیب به وجود هرزنامه پی می‌برند، رقیب نیز با یادگیری پارامترهای فیلتر به فریب کلاسه‌بندی کننده می‌اندیشد. همچنین به کمک الگوریتم استراتژی تکاملی نقطه تعادل بازی را محاسبه می‌نماییم. نتایج آزمایشات موید این مطلب است که روش پیشنهادی در مقایسه با روشهای دیگر تشخیص هرزنامه به مراتب دقیق‌تر و کارآمدتر می‌باشد.

### اطلاعات مقاله

مقاله پژوهشی کامل

دریافت: ۱۴ مرداد ۱۳۹۶

پذیرش: ۲۴ شهریور ۱۳۹۶

ارائه در سایت: [۱۵ آذر ۱۳۹۶]

### کلیدواژگان

تشخیص هرزنامه

یادگیری رفتار رقیب

تئوری بازی

استراتژی تکاملی.

## Spam Detection by Game theory

Samaneh Ghods

Department of Engineering, Semnan Branch, Islamic Azad University, Semnan, Iran

s1ghods@gmail.com

### Article Information

Original Research Paper

Received 8 August 2017

Accepted 15 September 2017

Available Online 6 August 2017

### Keywords

Spam detection, adversarial learning, game theory, evolutionary strategy.

### ABSTRACT

There are number of datamining applications that are fighting with Adversaries, Spam filtering to intrusion detection is as an example. For reducing the classifier accuracy, Adversary intentionally manipulate data. Consequently, in all these applications initially successful classifiers will decline easily. In this paper, we model the interaction between the classifier and the adversary as a two players sequential game then we model the interaction as an optimization problem and solve it using evolutionary strategy. Finally, simulation results show the good performance of the proposed algorithm, and improves accuracy spam detection on several real world data sets.

Please cite this article using:

Samaneh Ghods, Spam Detection by Game theory, *Journal of Mechanical Engineering and Vibration*, Vol. 8, No. 3, pp. 49-56, 2017 (In Persian)

برای ارجاع به این مقاله از عبارت ذیل استفاده نمایید:

## ۱- مقدمه

ما اغلب هنگام استفاده از سرویس پست الکترونیکی تقریباً روزانه با پیام هایی در صندوق پستی خود مواجه می شویم که هیچگاه این پیام ها را نمی خواهیم و نمی دانیم که فرستنده ی آنها کیست. نوع آن پیام ها متنوعند. برخی شامل آگهی های تجاری اند، برخی دیگر شامل آگهی های فریبنده اند. ظاهراً زیر ساخت پست اینترنتی، به عنوان یک رسانه ی کارآمد برای توزیع اطلاعات، می تواند به راحتی برای سوءاستفاده به کار برده شود. به این پیام های انبوه ناخواسته که ۳۰ سال پیش برای اولین بار در شبکه آرپانت فرستاده شد، هرزنامه<sup>۱</sup> می گویند. امروزه هرزنامه ها به یک وسیله ی ساده برای اذیت کردن کاربران تبدیل شده اند و علت بسیاری از آسیب های اقتصادی هستند. برای مثال بر طبق آماری که ارائه شده است [۱]، در سال ۲۰۰۹، روزانه ۲۰۰ بیلیون هرزنامه فرستاده شده است و هزینه ای که هرزنامه به سازمان های معروف جهان تحمیل کرده است، ۵۰ میلیارد دلار بوده است. راهکارهای زیادی برای تشخیص هرزنامه ها بکاررفته است که از جمله می توان استفاده از قوانین دانش، فیلتر کردن دامنه ها و لیست های سفید و سیاه نام برد. که متأسفانه هنوز راهکار عملیاتی و موثری برای حذف هرزنامه ها وضع نشده است به همین دلیل امروزه تحقیقات فراوانی روی آن صورت می گیرد. با این وجود اجرایی ترین روشها در حال حاضر، فیلترینگ اتوماتیک است که با توجه به الگوریتم های یادگیری ماشین [۲]، تشخیص هرزنامه صورت می پذیرد.

مساله یادگیری رفتار رقیب اولین بار در [۴] بعنوان کلاس بندی رقیب (اسپمرها) مطرح شد. آن ها به بررسی یک بازی بین کلاسه بندی کننده و رقیب پرداختند و نتایج آزمایشات بر روی فیلترینگ هرزنامه انجام شد. با توجه به اینکه در این روش هر دو طرف بازی یعنی کلاس بندی کننده و رقیب باید دانش کافی در مورد یکدیگر داشته باشند و تغییرات اندکی در هر یک از آنها باعث غیر اجرایی بودن روش فوق می گردد لذا عملاً این مورد غیر واقعی است.

برای حل مشکل فوق در [۳ و ۵] فرستندگان هرزنامه ابتدا با یادگیری رقیب و مهندسی معکوس<sup>۲</sup> پارامترهای کلاسه بندی کننده (فیلترها) را شناسایی می نمایند، سپس حملات رقیب را بر این مینا پایه ریزی می کنند، در این روش هیچ نقطه تعادلی برای رقیب و کلاسه بندی کننده وضع نگردیده است.

اخیراً در [۶] مدل یادگیری رقیب را بر مبنای تئوری بازی طراحی نموده اند. ایده تئوری بازی، رسیدن به نقطه تعادل رقیب و کلاسه بندی کننده است و فرض بر آن است که رقیب و کلاسه بندی کننده، شناخت کاملی از فضای همدیگر دارند و برای حل این بازی در فضای نامحدود و رسیدن به نقطه تعادل از الگوریتم هیوریستیک، ژنتیک استفاده شده است.

در این مقاله از بازی جدیدی به نام استکلبرگ استفاده می شود و در آن رفتار بین رقیب و کلاسه بندی کننده، مدلسازی می شود. حال با توجه به اینکه فضای مدل یادگیری رقیب در بازی استکلبرگ با مقادیر حقیقی مشخص شده است لذا در این مقاله برای رسیدن به نقطه تعادل بازی از دیگر الگوریتم اکتشافی تحت نام استراتژی تکاملی (ES) استفاده شده است. آزمایشات و نتایج نیز نشان می دهند که عمل فیلترینگ هرزنامه در این مدل دقیق تر و کارآمد تر از روش [۶] می باشد.

در این مقاله به اجمال، به بررسی بخش های زیر می پردازیم.

در بخش ۲، بازی استکلبرگ معرفی می شود و در بخش ۳، مدل تئوری بازی را برای سامانه تشخیص هرزنامه بیان می نماییم. در بخش ۴ الگوریتم استراتژی تکاملی را مطرح نموده، سپس در بخش ۵ نتایج آزمایشات را بر روی چند مجموعه داده نشان می دهیم. سرانجام در بخش ۶ به نتیجه گیری مقاله می پردازیم.

## ۲- بازی استکلبرگ

در این بخش روابط بین رقیب و کلاسه بندی کننده را از نوع بازی استکلبرگ در نظر گرفته ایم، در جایی که یکی از طرفین نقش رهبر را ایفا می کند و می تواند استراتژی مدنظر خود را به

<sup>2</sup> Reverse Engineering

<sup>1</sup> Spam

پیرو نیز به صورت زیر عمل می نماید:

$$v^s = R_F(u^s) \quad (3)$$

که در آن  $(u^s, v^s)$  را تعادل استکلبرگ می نامیم.

نقطه تعادل استکلبرگ  $[v]$ ، نقطه تعادلی است که در آن هیچیک از دو بازیکن انگیزه‌ای برای تغییر استراتژی انتخابی خود ندارند، مادامی که بازیکن دیگر بر استراتژی انتخابی خود پایبند باشد.

### ۲-۳- تعادل استکلبرگ یک مساله رام نشدنی<sup>۱</sup>

رابطه (۲) در ریاضی به عنوان یک مساله مشکل و رام نشدنی مطرح می گردد.

$$\max_{u,v} J_L(u,v) \text{ به شرطی که}$$

$$g(u,v) \leq 0 \quad (4)$$

$$v \in \arg \max \{ J_F(u,v); h(u,v) \leq 0 \}$$

که در آن  $g, h$  محدودیت هایی هستند که در فضای  $U, V$  در نظر گرفته می شوند.

اگر تمام توابع مساله مذکور ، خطی باشند ، این مساله را رام نشدنی می گوئیم. همچنین اگر توابع مساله ما به صورت غیر خطی باشند ، محدودیت  $\arg \max$  بیان می گردد.

مساله‌ی فوق یک مساله برنامه‌ریزی دوسطحی<sup>۲</sup> می‌باشد که مروری خوب در این زمینه در [۷] ارائه گردیده است و باعث می شود که برای رسیدن به نقاط تعادل ، از مش اکتشافی استفاده شود که در بخش چهارم در مورد روش حل آن به طور کامل توضیح داده خواهد شد.

### ۲- تئوری بازی در تشخیص هرزنامه

طرف مقابل (پیرو) تحمیل نماید. ابتدا رهبر استراتژی خود را تعیین نموده، سپس پیرو عکس‌العمل خود را در قالب بهترین ۱۰ استراتژی با اطلاعات موجود انجام می‌دهد. در فرضیات مساله، فضای استراتژی برای این دو بازیکن ، را پیچیده و نامحدود بیان می نماییم. هدف نهایی بازی، رسیدن به نقطه تعادل است که با فرض اینکه بازیکن ها از توابع عایدی یکدیگر آگاهی دارند و با استفاده از توابع هدف بازیکن ها ، به این امر می پردازیم.

### ۲-۱- تعریف

اجزای یک بازی دو نفره استکلبرگ [۷] در زیر معرفی می شود:

الف- این بازی بین دو بازیکن رهبر ( $L$ ) و پیرو ( $F$ ) انجام می گیرد. در مدل ما رقیب، رهبر و کلاسه بندی کننده، پیرو می باشد. رهبر همیشه حرکت اول را انجام می دهد.

ب- هر بازیکن از مجموعه ای از استراتژی ها تشکیل شده است که این استراتژی ها را برای  $L$  و  $F$  ، به ترتیب  $U$  و  $V$  می نامیم.

ج- هر بازیکن دارای یک تابع عایدی مشتق پذیر  $J_i; (i = L, F)$  است که به صورت  $J_i(U, V) \rightarrow R$  تعریف می شود .

( $U$  و  $V$  کران دار و محدب هستند.)

حال ، هر بازیکن با توجه به حرکت بازیکن مقابل ، استراتژی خود را بصورت زیر انتخاب می کند:

$$R_L = \arg \max_v J_L(u, v) \quad (1)$$

$$R_F = \arg \max_u J_F(u, v)$$

### ۲-۲- نقطه تعادل در بازی

در بازی استکلبرگ ، رهبر اولین حرکت بازی را آغاز می کند و به دنبال آن پیرو ، استراتژی را انتخاب می کند که عایدی خود را با توجه به حرکت رهبر ، بیشینه نماید. در ریاضی ، یک مساله بهینه سازی به صورت زیر است :

$$u^s = \arg \max_{u \in U} J_L(u, R_F(u)) \quad (2)$$

<sup>1</sup> Np- Hard

<sup>2</sup> Bilevel Programming

$$D_{KL}(N_1 \setminus N_2) = \frac{1}{2} (\log_e (\frac{\det \Sigma_2}{\det \Sigma_1})) + tr(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \quad (5)$$

### ۳-۱- فرمول‌سازی عایدی رقیب

هدف رقیب که در اینجا رهبر نامیده می‌شود، ارسال پیام‌های هرنزنامه می‌باشد که از فیلتر (کلاسه بندی کننده) عبور نماید، بطوریکه تعداد این پیام‌ها را بیشینه نماید. لذا برای آن مدلی را در نظر می‌گیریم که پارامتر FNR (میزان نرخ منفی نادرست) در آن لحاظ شود و بر مبنای آن تابع هدف را محاسبه می‌نماییم.

**عایدی رقیب** = میزان پیام‌های هرنزنامه‌ای که از فیلتر بعنوان پیام درست عبور نماید - هزینه انتقال از یک فضای نرمال به فضای نرمال دیگر

که به صورت مدل ریاضی به عنوان تابعی از متغیرهای تصمیم رقیب به شکل زیر نشان داده می‌شود: (کلاسه بندی کننده خطی می‌باشد)

$$J_L(u, w) = FNR - \alpha KLD(\mu_1, \sigma_1, \mu_1 - u, \sigma_1) = F(w, \mu_1 - u, \sigma_1) - \alpha KLD(\mu_1, \sigma_1, \mu_1 - u, \sigma_1) \quad (6)$$

(که در آن پارامتر  $\alpha$ ، قدرت پناستی  $KLD$  را نشان می‌دهد) به ارزش عایدی رهبر، سود رقیب نیز می‌گویند.

با استفاده از تابع چگالی احتمال در فضای نرمال

$$N(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

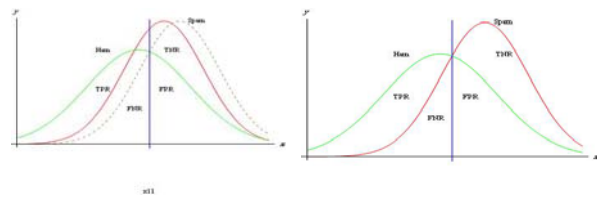
تابع چگالی تجمعی

$$F(t, \mu, \sigma) = \int_{-\infty}^t N(t, \mu, \sigma) dx$$

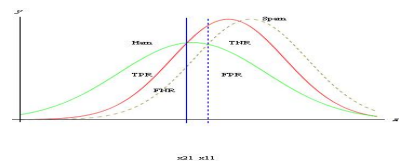
به دست می‌آید.

مدل بازی بین رقیب (اسپمر) و فیلتر هرنزنامه، با مساله طبقه بندی<sup>۱</sup> دو کلاسه مطرح می‌شود. برای سهولت، فضای ویژگی داده‌ها را به صورت یک بعدی در نظر می‌گیریم.

فرض کنیم فضای داده‌ای هرنزنامه‌ها فضای نرمال  $P \sim N(\mu_1, \sigma_1)$  و پیام‌های درست نیز در فضای نرمال  $Q \sim N(\mu_2, \sigma_2)$  قرار گرفته باشند. که در آن  $\mu_2 < \mu_1$ . رقیب بازی خود را با حرکت از  $\mu_1$  به  $\mu_1 - u$  آغاز می‌نماید (شکل ۱). سپس کلاسه بندی کننده، فضای تصمیم خود را از  $\frac{\mu_1 + \mu_2}{2}$  به سمت  $W$  (به طرف  $\mu_2$ ) سوق می‌دهد. فرض کنیم  $\mu_2 \leq W \leq \frac{\mu_1 + \mu_2}{2}$  و  $\mu_2 \leq u \leq \mu_1$ .



الف: حالت اولیه      ب: حرکت رقیب



ج: جابه‌جایی کلاسه بندی کننده با توجه به حرکت رقیب

شکل ۱ سه حالت سناریوی رفتار رقیب و کلاسه بندی کننده بر اساس نظریه بازی

برای تخمین تاثیر انتقال فضای هرنزنامه به اندازه  $u$  بر روی داده‌های اولیه از مفهومی به نام  $KLD$  (Kullback-Leibler Divergence) استفاده شده است که در [8] معرفی شده است.

تاثیر انتقال از  $N_1(\mu_1, \Sigma_1)$  به  $N_2(\mu_2, \Sigma_2)$  توسط  $KLD$  به صورت زیر محاسبه می‌گردد:

<sup>1</sup> Classification

## ۳-۲- فرمول سازی عایدی کلاسه بندی کننده

بعد از اینکه رهبر (رقیب) حرکت خود را مبنی بر فریب فیلتر آغاز نمود، فیلتر (پیرو) باید بر آن مبنا حرکت خود را انجام دهد یعنی درصد تشخیص درستی پیام های هرزنامه و پیام های درست را بیشینه نماید. لذا برای آن مدلی را در نظر می گیریم که پارامتر TPR (میزان نرخ مثبت درست) و TNR (میزان نرخ منفی درست) در آن لحاظ شود و بر مبنای آن تابع هدف را محاسبه می نماییم.

**عایدی کلاسه بندی کننده** = میزان پیام های هرزنامه ای که فیلتر به درستی تشخیص می دهد + میزان پیام های درستی که فیلتر صحیح تشخیص می دهد

که به صورت مدل ریاضی به عنوان تابعی از متغیرهای تصمیم کلاسه بندی کننده به شکل زیر نشان داده می شود: (طبقه بندی کننده خطی می باشد)

$$\begin{aligned} J_F(u, w) &= TPR + TNR \\ &= F(w, \mu_2, \sigma_2) - F(w, \mu - u_1, \sigma_1) + (1 - F(w, \mu_1 - u, \sigma_1)) - (1 - F(w, \mu_2, \sigma_2)) \\ &= 2F(w, \mu_2, \sigma_2) - 2F(w, \mu_1 - u, \sigma_1) \end{aligned} \quad (7)$$

$F$  تابع چگالی تجمعی به مانند مدل رقیب محاسبه می گردد...

معادلات (۶) و (۷) برای یک ویژگی از مجموعه داده ها بکار می رود حال برای محاسبه عایدی رقیب و کلاسه بندی کننده کل ویژگی های یک مجموعه داده، از روابط زیر استفاده می نماییم.

( فرض می نماییم تمام ویژگی های یک مجموعه داده مستقل از هم می باشند).

$$J_L(u, w) = \frac{1}{q} \sum_{i=1}^q J_{L_i}(u, w) \quad (8)$$

$$J_F(u, w) = \frac{1}{q} \sum_{i=1}^q J_{F_i}(u, w) \quad (9)$$

که  $q$  تعداد ویژگی های یک مجموعه داده می باشد.

حال با توجه به توابع عایدی ذکر شده و همچنین

معادلات (۲) و (۳) نقطه تعادل بازی که یک مساله بهینه سازی می باشد بیان می گردد. در بخش بعد توضیح می دهیم که چگونه استراتژی تکاملی (ES) این نقطه تعادل را بدست می آورد.

## ۴- استراتژی تکاملی

به منظور حل مدل رقیب-کلاسه بندی کننده استکلبرگ که یک مدل دو سطحی می باشد، بطوری که رقیب بعنوان سطح بالاتر با ارائه مقدار  $u$  به سطح پایین تر، کلاسه بندی کننده مقادیر بهینه  $w$  را محاسبه نموده و این مقادیر بهینه را به رقیب برای محاسبه تابع هدفش ارائه می دهد. حال برای بدست آوردن مقادیر بهینه  $u^*$  از الگوریتم متا هیوریستیک ES استفاده نموده ایم.

استراتژی های تکاملی عموماً برای مسائل بهینه سازی با مقادیر حقیقی به کار می روند و در آنها جهش<sup>۱</sup> نسبت به ترکیب<sup>۲</sup> (تقاطع) بیشتر مورد تاکید است. در حالت کلی یک کروموزوم می تواند به صورت رابطه زیر تعریف شود:

$$\langle X_1, \dots, X_n, \sigma_1, \dots, \sigma_n \rangle \quad \text{که در آن } X_i \text{ متغیر } i \text{ ام و } \sigma_i \text{ اندازه جهش می باشد.}$$

یکی از عملگرهایی که در سالهای اخیر در استراتژی تکاملی به کار گرفته شده است عملگر ترکیب است که پس از اعمال اندازه جهش والد، فرزندان را تولید می کند عملگر جهش بر اساس توزیع نرمال ارایه می شود و اندازه جهش برای هر متغیر به کمک روابط زیر محاسبه می شود.

$$\begin{aligned} \sigma' &= \sigma * \exp(\tau * N(0,1)) \\ x'_i &= x_i + \sigma' * N(0,1) \end{aligned} \quad (10)$$

که در آن  $N(0,1)$  مقداری تصادفی از توزیع نرمال با میانگین صفر و انحراف معیار ۱ می باشد.

در این روش، گزینش از اجتماع  $\mu$  والد و  $\lambda$  فرزند، تعداد  $\mu$  والد برای نسل بعد با توجه به برازندگی<sup>۳</sup> شان انتخاب می شوند.

<sup>1</sup> Mutation

<sup>2</sup> Recombination

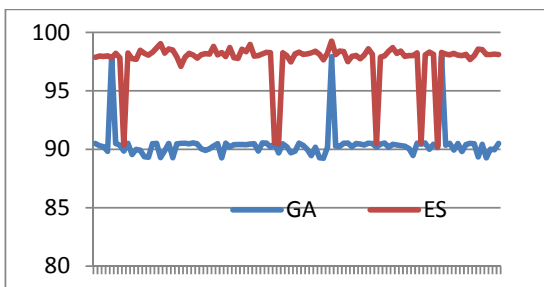
<sup>3</sup> Fitness

و تعداد نمونه ها نیز برای هر ماه متفاوت می باشد. در تمام آزمایشات از ارزیابی<sup>۱</sup> ده تایی<sup>۲</sup> استفاده می نماییم. مجموعه داده ها را بطور تصادفی به ۱۰ قسمت تقسیم کرده و آنرا ۱۰ بار تکرار می نماییم بطوریکه هر بار یک قسمت آنرا برای آزمایش و بقیه نه قسمت را برای آموزش استفاده می کنیم. هر الگوریتم تکاملی نیز برای ۱۰۰ نسل اجرا می شود.

#### ۵-۲- نتایج مرحله آموزش

همانطور که در بخش سوم عنوان شد، در مرحله آموزش، مدل بازی بین رقیب و فیلتر، برای تک تک ویژگیهای مجموعه داده ها اجرا می گردد و با الگوریتم ES نقطه تعادل بازی و در ادامه توابع عایدی و FNR و FPR محاسبه می شود.

شکل ۲ عایدی طبقه بندی کننده (فیلترینگ) و شکل ۳ عایدی رقیب را برای دو الگوریتم ES و GA مقایسه کرده است در این مثال از ویژگی اول مجموعه داده ژانویه استفاده شده است که فضای هرزنامه آن تابع نرمال  $P \sim N(0.15, 0.31)$  و فضای پیام درست نیز تابع نرمال  $Q \sim N(0.07, 0.29)$  می باشد. تابع طبقه بندی کننده خطی می باشد، متغیرهای مساله در ۱۰۰ نسل تکامل یافته اند و در هر اجرا میانگین نسل ها را محاسبه نموده ایم. مقدار پارامتر  $\alpha$  (قدرت پنداری  $KLD$ ) در معادله (۶)، ۰.۱ می باشد.



شکل ۲ عایدی طبقه بندی کننده (فیلترینگ) یک ویژگی براساس الگوریتم های ES و GA برحسب ۱۰۰ اجرا

حال برای حل معادله (۲) از استراتژی تکاملی استفاده می شود و بهترین انتقال فضای هرزنامه بعد از چند اجرا محاسبه می شود. جزئیات در الگوریتم ۱ بیان شده است.

#### الگوریتم ۱: استراتژی تکاملی برای حل تعادل بازی استکلبرگ

ورودی: تعداد  $K$  فرد  $u$  در یک نسل

خروجی: نقطه تعادل  $u^s$  برای بازی استکلبرگ

۱: انتخاب  $K$  بصورت تصادفی (تعداد جمعیت)

۲: موارد زیر را برای  $M$  نسل اجرا نماید

۱-۲: تعداد  $K$  بار، مقدار انتقال  $u$  از جمعیت انتخاب می کنیم و موارد ذیل را به ازای هر عضو جمعیت اجرا می نماییم.

۱-۲-۱:  $u_i$  انتقال فضای رقیب از جمعیت می باشد.

۲-۱-۲:  $R_F(u_i)$  را برای کلاس بندی کننده محاسبه می نماییم.

۲-۱-۳: عایدی رقیب به ازای  $u_i$  را از تابع هدف

$J_L(u_i, R_F^{u_i})$  محاسبه می نماییم.

۲-۲:  $u_i$  با کمترین عایدی رقیب را برای جمعیت مفروض انتخاب می نماییم.

۳-۲: عملگرهای ترکیب و جهش و تولید نسل جدید را بر روی جمعیت اجرا می نماییم.

۳-۳: انتقال  $u$  برای کمترین عایدی رقیب در  $M$  نسل،  $u^s$  را نشان می دهد.

#### ۵- نتایج

در این فصل نتایج محاسباتی الگوریتم های ژنتیک (GA) [۶] و استراتژی تکاملی (ES)، که برای مدل رقیب و کلاس بندی کننده در فاز آموزش و آزمایش ارائه گردید را از نظر کیفیت جوابها مورد بررسی قرار داده و سپس دقت فیلتر تشخیص هرزنامه را برای انواع مجموعه داده ها بیان می نماییم.

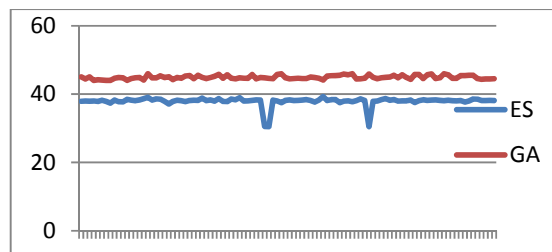
#### ۵-۱- مشخصات آزمایش

برای رسیدن به نتایج الگوریتم ها و مدل بیان شده از مجموعه داده هایی که از یک سرور پست الکترونیکی [۹] در شش ماه اول یک سال جمع آوری شده است، استفاده می نماییم. این مجموعه داده ها عددی می باشند و از نمونه های هرزنامه و پیام های درست تشکیل شده اند. مجموعه داده ها دارای ۲۰ ویژگی

<sup>1</sup> Cross Validation

<sup>2</sup> 10-Fold

Data Set	GA		ES	
	Spam	Ham	Spam	Ham
JAN	92.6%	94.6%	94%	97.5
FEB	66%	91.3%	67.3%	93%
MAR	64.5%	87.3%	64.9%	88%
APR	50.8%	95.6%	51.7%	96.3%
MAY	66.6%	94.8%	68.2%	95.6%
JUN	67.18%	95.7%	69.3%	96.2%



شکل ۳ عایدی رقیب یک ویژگی براساس الگوریتم ES و GA برحسب ۱۰۰ اجرا

همانطور که در بخش دوم آورده شد در مدل معرفی شده دو طرف بازی یعنی رقیب و فیلتر (کلاسه بندی کننده) دارای دو عایدی می باشند؛ برای مقایسه چند الگوریتم، هر چه عایدی فیلتر بیشتر (فیلتر هرزنامه بیشتری تشخیص دهد) و هر چه عایدی رقیب کمتر (رقیب کمتر می تواند فیلتر را فریب دهد) باشد، الگوریتم دقت بهتری دارد. با توجه به این نکته و شکل های ۲ و ۳، می توان نتیجه گرفت که ES از الگوریتم ژنتیک بهتر عمل می کند.

### ۵-۳- نتایج مرحله آزمایش

در این مرحله به کمک موقعیت نهایی کلاسه بندی کننده در نقطه تعادل که برای هر ویژگی با داده های مرحله آموزش تعیین گردید، دقت فیلتر را با مراحل ذیل محاسبه می نماییم:

- از داده های مرحله آزمایش هر مجموعه داده، یک نمونه را انتخاب و هر ویژگی از آن را به کمک کلاسه بندی کننده در مرحله آموزش، مشخص می نماییم که مقدار این ویژگی، نمونه را هرزنامه و یا پیام درست تشخیص می دهد.
  - برچسب کلاس را برای تمام ویژگیهای یک نمونه به روش قبل بدست می آوریم.
- حال با توجه به روش رای گیری، برچسب آن نمونه بستگی به بیشترین تعداد برچسبی دارد که ویژگیها مشخص نموده اند.

### ۶- نتیجه

در این مقاله با استفاده از نظریه تئوری بازی مسابقه بین رقیب و کلاسه بندی کننده مورد بحث و بررسی قرار گرفته است و تشخیص هرزنامه توسط کلاسه بندی کننده در این بازی بیان می شود. این دو بازیکن با توجه به عایدی خود بهترین استراتژی را انتخاب نموده و بعد از تکرار بازی به یک نقطه تعادل می رسند. نتایج آزمایشات نشان می دهد که روش استراتژی تکاملی توانایی حل نقطه تعادل بازی را داراست. و همچنین از روشهای دیگر الگوریتم های اکتشافی [۶] به مراتب دقیق تر می باشد. برای آینده کار در نظر داریم از توابع غیر خطی در کلاسه بندی کننده ها استفاده نموده تا بتوان عایدی و همچنین نقطه تعادل را با دقت بیشتری محاسبه نمود و منجر به عبور کمتری هرزنامه از کلاسه بندی کننده شویم.

### مراجع

- [1] Ferris Research, The Global Economic Impact of Spam, Report #409, 2009.
- [2] S. Heron, Technologies for spam detection, *Network Security*, Vol. 2009, No. 1, pp. 11-15, 2009.
- [3] D. Lowd, C., Meek, Good word attacks on statistical spam filters, In CEAS, Palo Alto, CA, 2005.
- [4] D. Nilesch, D. Pedr, Adversarial classification, In KDD, New York, NY, USA, pp. 99-108, 2004.

جدول 1 دقت تشخیص هرزنامه (Spam) و پیام های درست (Ham) بر روی مجموعه داده های مختلف براساس الگوریتم ES و GA

- [5] D. Lowd, C. Meek, Adversarial learning, In KDD, New York, NY, USA, pp. 641-647, 2005.
- [6] W. Liu, S. Chawla, A Game Theoretical Model for Adversarial Learning, IEEE on Data Mining, pp. 25-30, 2009.
- [7] K. Binmore, Playing for real: a text on game theory, Oxford University Press, USA, 2007.
- [8] S. Kullback, R.A. Leibler, On information and sufficiency, The Annals of Mathematical Statistics, pp. 79-86, 1951.