

## تحلیل نرخ دلار در شبکه اجتماعی توییتر با استفاده از الگوریتم های یادگیری

### ماشین

صبا فرهادی<sup>۱</sup>، محبوبه شمسی<sup>۲</sup>، عبدالرضا رسولی کناری<sup>۳</sup>

<sup>۱</sup>دانشجوی مهندسی کامپیوتر، دانشگاه صنعتی قم،

<sup>۲</sup>استادیار مهندسی کامپیوتر، دانشگاه صنعتی قم، [shamsi@Qut.ac.ir](mailto:shamsi@Qut.ac.ir)

<sup>۳</sup>استادیار مهندسی کامپیوتر، دانشگاه صنعتی قم، [Rasouli@qut.ac.ir](mailto:Rasouli@qut.ac.ir)

### چکیده

در سال های اخیر، تحقیقات تحلیل احساسات در توییتر، توییت ها را برای استخراج احساسات کاربر تحلیل می کنند که گسترش زیادی داشته است. بسیاری از محققان برای چنین تحلیل هایی تصمیم می گیرند از الگوریتم های یادگیری ماشین و یادگیری عمیق استفاده کنند. در این پژوهش با استفاده از الگوریتم های یادگیری ماشین و یادگیری عمیق، یک تحلیل احساسات جزئی از داده های توییتر فارسی مربوط به «نرخ دلار» در بازه زمانی سه ماهه تابستان ۱۳۹۷ انجام دادیم. این پژوهش شامل پیش پردازش توییت توسط یک رویکرد استخراج ویژگی است. سپس با طبقه بندی احساسی شامل ناامیدی، خوشحالی، ترس، ناراحتی و امید، ویژگی ها طبقه بندی می شوند. در چارچوب پیشنهادی، الگوریتم های رگرسیون لجستیک، ماشین بردار پشتیبان، درخت تصمیم و جنگل تصادفی و روش جدید ترکیبی رگرسیون لجستیک و ماشین بردار پشتیبان برای ارزیابی تحلیل احساسات مورد استفاده قرار می گیرد. همچنین از دو الگوریتم یادگیری عمیق شبکه عصبی پیچشی و شبکه عصبی بازگشتی استفاده شده است. مشاهدات تجربی نشان می دهد که احساسات مردم در مورد تغییر نرخ دلار در بازه زمانی سه ماهه تابستان بیشتر به سمت ناامیدی گرایش داشته است. همچنین روش ترکیبی ارائه شده برای ماه اول تابستان با مقدار ۹۹ درصد و دو ماه دیگر با مقدار ۹۷ درصد بیشترین دقت را نسبت به دیگر الگوریتم ها داشته است.

### کلمات کلیدی:

توییتر، یادگیری ماشین، یادگیری عمیق، احساسات، دلار، دقت

### ۱- مقدمه

به موجودیت هایی مثل عناوین، رویدادها، افراد خاص، مسائل، خدمات، محصولات، سازمان ها و موارد مرتبط با آنها است. تحلیل احساسات شاخه ای از یادگیری ماشین، داده کاوی<sup>۲</sup>، پردازش زبان طبیعی و زبان شناسی محاسباتی است، همچنین مفاهیمی را از علم جامعه شناسی

تحلیل احساسات<sup>۱</sup> گاهی اوقات از آن به نظر کاوی<sup>۲</sup> نیز یاد می کنند و یک حوزه تحقیقاتی است که هدف آن تحلیل احساسات یا نظرات افراد راجع

احساس ناامیدی، خوشحالی، ترس، ناراحتی و امید طبقه بندی می کنیم. سپس متون را به ماتریس تبدیل کرده و توسط الگوریتم های مختلف یادگیری ماشین و یادگیری عمیق و یک روش جدید ترکیبی برای طبقه بندی این توییت ها به کار برده شده است. سپس برای ارزیابی روش پیشنهادی پارامترهایی چون صحت<sup>۵</sup>، بازیابی<sup>۶</sup>، F-Measure و دقت<sup>۷</sup> را مورد توجه قرار داده ایم [۳].

## ۱-۱- کارهای گذشته

تحقیقات متعددی بر روی تحلیل داده های شبکه های اجتماعی تمرکز کرده اند، خصوصاً داده هایی که مختص به یک رویداد خاص هستند. این استفاده گسترده از شبکه های اجتماعی توجه زیادی را از سوی محققان به خود جلب کرده و بسیاری از تحقیقات برای دریافت اطلاعات مهم بر روی این رویدادها انجام شده است. در سال های اخیر تحقیقات گسترده ای در زمینه تحلیل احساسات در توییت انجام شده است.

Jain و Dandannavar چندین مرحله برای تحلیل احساسات بر روی داده های توییت با استفاده از الگوریتم های یادگیری ماشین را بررسی کردند. رویکرد آن ها در ابتدا داده ها را جمع آوری می کرد و بعد با استفاده از تکنیک های NLP پیش پردازش می کند. سپس، برای استخراج ویژگی های مربوط به احساس، عمل استخراج ویژگی بر روی آن ها انجام شد. در نهایت، یک مدل با استفاده از دسته بندی مانند بیز ساده<sup>۸</sup>، ماشین بردار پشتیبان و درخت تصمیم آموزش داده شد. فریمورک ارائه شده تحلیل احساسات را با استفاده از بیز ساده چندجمله ای<sup>۹</sup> و درخت تصمیم اجرا کرد. نتایج نشان داد که درخت تصمیم موثرتر عمل می کند و صحت، دقت، بازیابی و F1-score بهتری را ارائه می دهد [۴].

تحقیقات گسترده ای نیز توسط Go و همکارانش انجام شده است که یک راه حلی با استفاده از نظارت از راه دور<sup>۱۰</sup> برای تحلیل احساسات بر اساس توییت ها ارائه داده اند. در این روش، آن ها از داده های آموزشی شامل توییت های حاوی شکلک استفاده کردند که به عنوان برچسب های پرت<sup>۱۱</sup> عمل می کرد. آن ها مدلی بر اساس دسته بندی بیز ساده، حداکثر آنتروپی، و ماشین بردار پشتیبان ایجاد کردند. ویژگی های آن ها از یونیگرام<sup>۱۲</sup>، بایگرام<sup>۱۳</sup> و POS<sup>۱۴</sup> تشکیل می شد. آن ها نتیجه گرفتند که ماشین بردار پشتیبان نسبت به سایر مدل ها بهتر عمل می کرد و این که یونیگرام ها ویژگی موثرتری نسبت به بقیه بودند [۵].

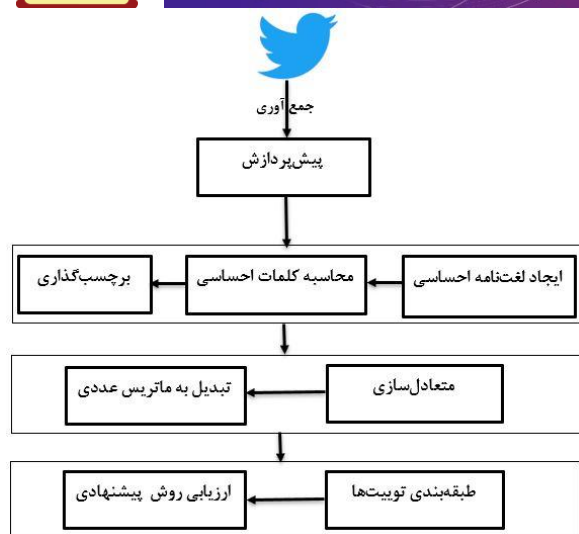
علاوه بر این، چندین پژوهش در زمینه SemEval، وظیفه دسته بندی توییت ها با چند صد شرکت کننده را مورد بررسی قرار دادند [۶]. در طول

و روان شناسی قرض می گیرد. از چند سال گذشته، رونق شبکه های اجتماعی، توسعه تحلیل احساسات را سرعت بخشیده است. [۱]  
تشخیص ساختار متن به سه شکل مختلف سند، جمله و کلمه امکان پذیر است. در سطح سند، اسناد به نظرات مثبت، منفی و خنثی طبقه بندی می شوند. در سطح جمله تعیین می کند که آیا جمله نظر مثبت، منفی و یا خنثی را بیان می کند و در سطح ویژگی، برای تعیین نظر مثبت، منفی و یا خنثی به جزئی ترین شکل ممکن (کلمه) انجام می گیرد و نظرات را در مورد ویژگی های خاص موضوع بررسی می کنیم [۲]. در این مطالعه، طبق داده های جمع آوری شده توییت<sup>۴</sup> با هشتگ دلار، تحلیل احساسات در سطح جمله مورد توجه قرار گرفته است.

در گذشته برای به دست آوردن نظرات کاربران از فرم های نظرسنجی استفاده می شد، ولی امروزه با گسترش وب و توسعه اینترنت این کار به راحتی از طریق پست گذاشتن در شبکه های اجتماعی و یا ارائه دیدگاه ذیل هر پست قابل انجام است. رشد اطلاعات به صورت آنلاین در شبکه های اجتماعی، تحلیل احساسات را ضروری تر کرده است. از لحاظ اقتصادی، تحلیل احساسات می تواند توصیه ها و پیشنهادات آنلاین برای مشتریان و فروشندگان داشته باشد. از سوی دیگر، این ترجیحات کاربری که داده ها نمایان می کنند می تواند به بسترهای فروش آنلاین کمک کند تا محصولات و خدمات شان را تحلیل کنند. از سوی دیگر، به علت ذات مجازی خرید آنلاین، بررسی جزئیات دقیق یک کالا و اینکه آیا مصرف کننده مایل است نظرات یا دیدگاه های سایر مصرف کننده ها را بداند آسان نیست.

از لحاظ سیاسی، تقاضاهای گسترده برای اطلاعات سیاسی می تواند به یک عامل مهم دیگر در نظر گرفته شود. کاربردهای تجاری، تنها انگیزه مردم برای بررسی و بیان دیدگاه ها به صورت آنلاین نیست. به عنوان مثال در تحلیل بحث ها در توییت قبل از انتخابات پارلمان اروپا، محققان بیش از ۱.۲ میلیون توییت به سه زبان (انگلیسی، فرانسوی و آلمانی) در طول دوره دو هفته ای جمع آوری کردند و میزان مثبت و منفی بودن نظرات مردم نسبت به انتخابات پارلمان را تحلیل کردند. [۱] تحلیل احساسات توییت در مقایسه با سایر منابع بخاطر وجود لغات عامیانه، اشتباهات تایپی و لغات حاوی طعنه و کنایه، دشوارتر است. ماکزیمم تعداد کاراکتر مجاز در توییت ۱۴۰ است. روش مبتنی بر پایگاه دانش و یادگیری ماشین، دو استراتژی به کار رفته در تحلیل احساسات متون هستند.

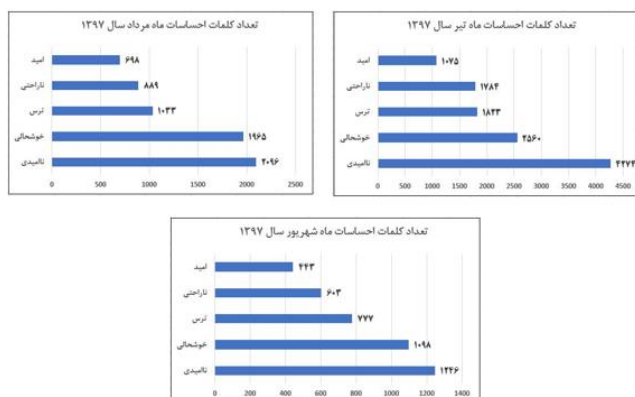
در این پژوهش، پست های توییت در رابطه با نرخ دلار را از اطلاعات بی معنا می زداییم و سپس بر اساس لغت نامه احساسی حاوی پنج



شکل (۱): نمایی از روش پیشنهادی

### ۳- ارزیابی روش پیشنهادی

در این پژوهش ما یک فرهنگنامه احساسی با عناوین ناامیدی، خوشحالی، ترس، ناراحتی و امید که هر کدام شامل کلمات مرتبط به آن‌ها است ایجاد کرده ایم و طبق این فرهنگنامه با شمارش کلمات احساسی هر توییت برجسب احساسی آن را تعیین و مشخص کرده ایم. در نهایت تعداد هر کلمه احساسی برای هر سه ماه تابستان در شکل (۲) مشاهده می کنید. در نتیجه ناامیدی در ماه تیر با مقدار ۴۲۷۴ و در ماه مرداد با مقدار ۲۰۹۶ و در ماه شهریور با مقدار ۱۲۴۶ جهت نرخ دلار بیشترین تاثیر را روی افراد داشته است و امیدواری در ماه تیر با مقدار ۱۰۷۵ و در ماه مرداد با مقدار ۶۹۸ و در ماه شهریور با مقدار ۴۴۳ کمترین تاثیر را داشته است.



شکل (۲): تعداد کلمات احساسی در سه ماه تابستان سال ۱۳۹۷

همچنین در این پژوهش برای ارزیابی روش پیشنهادی از پنج الگوریتم یادگیری ماشین از جمله رگرسیون لجستیک (LR)، ماشین بردار پشتیبان (SVM)، درخت تصمیم (DT)، جنگل تصادفی (RF) و

دهه گذشته توجه گسترده‌ای به تحلیل احساسات بر اساس داده‌های توییت و همچنین رگرسیون ترتیبی وجود داشته است. مسئله رگرسیون، یکی از اصلی‌ترین حوزه‌های پژوهش در زمینه یادگیری ماشین و داده‌کاوی است.

ساختار مقاله به شرح زیر می‌باشد: بخش ۲، پیاده‌سازی روش پیشنهادی تحلیل احساسات توضیح داده شده است. بخش ۳، ارزیابی روش پیشنهادی انجام شده است. بخش ۴، مقاله را به پایان می‌رساند و دامنه کارهای آینده را ارائه می‌دهد.

### ۲- پیاده‌سازی روش پیشنهادی

در این پژوهش، مجموعه توییت‌های فارسی در مورد نرخ دلار در بازه زمانی سه ماهه تابستان سال ۱۳۹۷ شامل ماه اول ۱۳۹۷ و توییت و ماه دوم ۱۱۶۵۰ توییت و ماه سوم ۷۶۷۰ توییت برای تحلیل احساسات در نظر گرفته شده است.

روش پیشنهادی برای بهبودی تحلیل احساسات یک فرهنگنامه از کلمات احساسی (ناامیدی، خوشحالی، ترس، ناراحتی و امید) توییت‌های فارسی ایجاد کرده ایم. ابتدا توییت‌ها را از کلمات بی معنا پاک‌سازی کرده و براساس فرهنگنامه احساسی موردنظر کلمات احساسی هر توییت را محاسبه می‌کنیم تا برجسب توییت‌ها را به درستی تشخیص دهیم و تحلیل احساسات در سطح جملات فارسی را بهبود می‌بخشیم. سپس برای متعادل سازی طبقه‌ها از روش نمونه‌گیری بیش از حد استفاده می‌کنیم و متون را به ماتریس عددی تبدیل کرده و با استفاده از الگوریتم‌های رگرسیون لجستیک ( $LR^{15}$ )، ماشین بردار پشتیبان ( $SVM^{16}$ )، درخت تصمیم ( $DT^{17}$ )، جنگل تصادفی ( $RF^{18}$ ) و روش جدید ترکیبی رگرسیون لجستیک و ماشین بردار پشتیبان ( $LR+SVM$ ) و شبکه عصبی پیچشی ( $CNN^{19}$ ) و شبکه عصبی بازگشتی ( $LSTM^{20}$ ) روش پیشنهادی مورد ارزیابی قرار می‌گیرد. مراحل روش پیشنهادی را در شکل (۱) مشاهده می‌کنید.

چارچوب پیشنهادی، الگوریتم های رگرسیون لجستیک، ماشین بردار پشتیبان، درخت تصمیم و جنگل تصادفی و روش جدید ترکیبی رگرسیون لجستیک و ماشین بردار پشتیبان و دو الگوریتم یادگیری عمیق از جمله الگوریتم های شبکه عصبی پیچشی و شبکه عصبی بازگشتی برای ارزیابی تحلیل احساسات مورد استفاده قرار می گیرد.

در نهایت به این نتیجه رسیدیم که احساسات مردم در مورد تغییر نرخ دلار در بازه زمانی سه ماهه تابستان سال ۱۳۹۷ بیشتر به سمت ناامیدی گرایش داشته است. همچنین روش ترکیبی ارائه شده برای ماه اول تابستان با مقدار ۹۹ درصد و دو ماه دیگر تابستان با مقدار ۹۷ درصد بیشترین دقت را نسبت به دیگر الگوریتم ها بهترین عملکرد را داشته است. همچنین الگوریتم های ماشین بردار پشتیبان (SVM) و درخت تصمیم (DT) و جنگل تصادفی (RF) و شبکه عصبی پیچشی (CNN) با مقدار ۹۸ درصد و رگرسیون لجستیک (LR) و شبکه عصبی بازگشتی (LSTM) با مقدار ۹۶ درصد بهترین دقت را به خود اختصاص داده اند. در زبان فارسی مشکلاتی چون عامیانه بودن کلمات و رعایت نکردن فاصله بین کلمات و رعایت نکردن قاعده ساختار جمله و استفاده از شکلک های تصویری می باشد که می توان این مشکلات را در کارهای آینده برای بهبودی تحلیل احساسات رفع کرد. همچنین می توان لغت- نامه احساسی را گسترش داد.

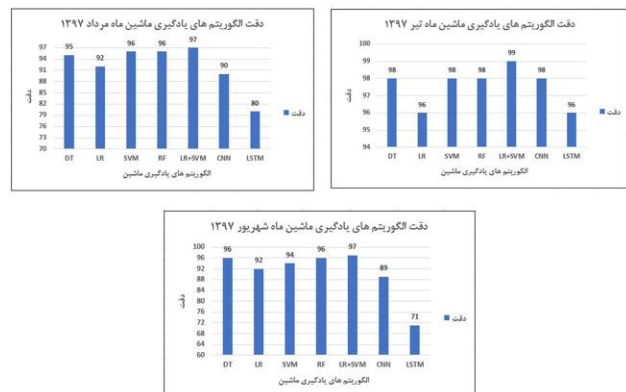
## مراجع

- [1] Yue, L., et al., *A survey of sentiment analysis in social media*. Knowledge and Information Systems, 2019.
- [2] Feldman, R., *Techniques and applications for sentiment analysis*. Communications of the ACM, 2013.
- [3] Zachman, John A., *A Framework for Information Systems Architecture*, IBM Systems Journal, Vol. 26, No. 3, 1987.
- [4] Jain, A.P., Dandannavar, P., *Application of machine learning techniques to sentiment analysis*. in 2016 2nd International Conference on Applied and Theoretical Computing and Communication Technology (iCATccT). IEEE, 2016.
- [5] Go, A., Bhayani, R., Huang, L., *Twitter sentiment classification using distant supervision*. CS224N project report, Stanford, 2009. 1(12): p. 2009.
- [6] Nakov, P., et al., SemEval-2016 task 4: *Sentiment analysis in Twitter*. arXiv preprint arXiv:1912.01973, 2019.

## زیر نویس ها

روش جدید ترکیبی رگرسیون لجستیک و ماشین بردار پشتیبان (LR+SVM) و از دو الگوریتم یادگیری عمیق از جمله شبکه عصبی پیچشی (CNN) و شبکه عصبی بازگشتی (LSTM) دقت آن را در سه ماه تابستان در شکل (۳) بیان شده است.

همان طور که در شکل (۳) مشاهده می کنید در ماه تیر روش ترکیبی LR+SVM با مقدار ۹۹ درصد بیشترین دقت و LR و LSTM با مقدار ۹۶ درصد کمترین دقت و سایر الگوریتم ها از جمله SVM و DT و RF و CNN با مقدار ۹۸ درصد بهترین دقت را دارند. در ماه مرداد روش ترکیبی LR+SVM با مقدار ۹۷ درصد بیشترین دقت و LSTM با مقدار ۸۰ درصد کمترین دقت و سایر الگوریتم ها از جمله SVM و RF با مقدار ۹۶ درصد و DT با مقدار ۹۵ درصد و LR با مقدار ۹۲ درصد و CNN با مقدار ۹۰ درصد می باشند. در ماه شهریور روش ترکیبی LR+SVM با مقدار ۹۷ درصد بیشترین دقت و LSTM با مقدار ۷۱ درصد کمترین دقت و سایر الگوریتم ها از جمله DT و RF با مقدار ۹۶ درصد و SVM با مقدار ۹۴ درصد و LR با مقدار ۹۲ درصد و CNN با مقدار ۸۹ درصد می باشند.



شکل (۳): نتایج دقت الگوریتم های یادگیری ماشین سه ماه تابستان سال ۱۳۹۷

## ۴- نتیجه گیری

در این پژوهش سعی کردیم با روش ارائه شده و لغت نامه موردنظر، احساس هر یک از توییت های استخراجی سه ماه تابستان سال ۱۳۹۷ مربوط به نرخ دلار را پیدا کنیم. علاوه بر این یک روش ترکیبی از رگرسیون لجستیک و ماشین بردار پشتیبان ارائه دادیم. توییت ها ابتدا از کلمات اضافی پاک سازی می شوند. سپس تحت یک طبقه بندی احساسی شامل احساسات ناامیدی، خوشحالی، ترس، ناراحتی و امید لغت نامه ای ایجاد می کنیم و طبق لغت نامه احساس هر توییت را می یابیم. در



<sup>3</sup> Data Mining

<sup>4</sup> Twitter

<sup>5</sup> Precision

<sup>6</sup> Recall

<sup>7</sup> Accuracy

<sup>8</sup> Naïve Bayes

<sup>9</sup> Multinomial Naïve Bayes

<sup>10</sup> Distant supervision

<sup>11</sup> Noisy labels

<sup>12</sup> Unigram

<sup>13</sup> Bigrams

<sup>14</sup> Part-Of-Speech

<sup>15</sup> Logistic Regression

<sup>16</sup> Support Vector Machine

<sup>17</sup> Diction Tree

<sup>18</sup> Random Forest

<sup>19</sup> Convolutional Neural Networks

<sup>20</sup> Long short-term memory