



یک چارچوب نیمه نظارتی مبتنی بر لغت نامه وفقی خودساخت جهت تحلیل نظرات فارسی

محسن نجف زاده^۱، سعید راحتى قوچانى^{۲*} و رضا قائمى^۳

^۱ گروه مهندسی کامپیوتر، واحد مشهد، دانشگاه آزاد اسلامی، مشهد، ایران

^۲ گروه مهندسی برق، واحد مشهد، دانشگاه آزاد اسلامی، مشهد، ایران

^۳ گروه مهندسی کامپیوتر، واحد قوچان، دانشگاه آزاد اسلامی، قوچان، ایران

چکیده

با معرفی وب ۲.۰ و ۳.۰ تعاملات کاربران در فضای مجازی، منجر به ایجاد انبوهی از نظرات ارزشمند شده است. با توجه به دشواری یا عدم امکان تحلیل و بررسی دستی این نظرات، تحلیل احساس متن و یا نظرکاوی به عنوان یکی از زیرمجموعه های پردازش زبان طبیعی مطرح شد. تلاش های محدودی در نظرکاوی فارسی نسبت به سایر زبان ها صورت گرفته است. در این مقاله برای نخستین بار، یک چارچوب نیمه نظارتی برای نظرکاوی فارسی ارائه شده است. در ضمن، از آنجاکه یکی از آخرین پیشرفت های علمی در نظرکاوی زبان فارسی الگوریتمی بر اساس استخراج الگوهای حسی وفقی (حساس به مجموعه داده) مبتنی بر خبره انسانی است، در این پژوهش ضمن ارتقای الگوریتم یادشده، تعیین برچسب های حاوی احساس به کمک یک لغت نامه خودساخت (بدون نیاز به خبره انسانی) وفقی انجام می گیرد؛ همچنین کاربرد دسته بند مدل مخفی مارکوف خودناظر بر روی خصیصه های یادشده در کنار قوانین مبتنی بر معیار شباهت برای فرآیند نظرکاوی بررسی شده است. در راستای خودآموزسازی هوشمند، روشی برای ارزیابی قابلیت اطمینان بالای خروجی، ارائه شده است که خودآموزی به شرط وجود آن انجام می پذیرد. روش پیشنهادی با اجرا بر روی دادگان مبنا نرخ صحت نود درصد (با وجود عدم نیاز به خبره انسانی) را که در مقایسه با روش های نظارتی و نیمه نظارتی مستقل از خبره موجود برتری قابل ملاحظه ای دارد، خروجی می دهد؛ همچنین این الگوریتم نیمه نظارتی هنگام استفاده از مجموعه آموزش کوچک با نسبت مجموعه دادگان آموزش/آزمون ده به نود نیز بررسی و با نرخ صحت ۸۰٪ قابلیت اطمینان آن به اثبات رسید.

واژگان کلیدی: نظرکاوی، یادگیری خودناظر، لغت نامه خودساخت، مدل مخفی مارکوف، لغت نامه وفقی

A Semi-supervised Framework Based on Self-constructed Adaptive Lexicon for Persian Sentiment Analysis

Mohsen Najafzadeh¹, Saeed Rahati Quchani^{2*} & Reza Ghaemi³

¹Department of Computer Engineering, Mashhad Branch, Islamic Azad University, Mashhad, Iran

²Department of Electronic Engineering, Mashhad Branch, Islamic Azad University, Mashhad, Iran

³Department of Computer Engineering, Quchan Branch, Islamic Azad University, Quchan, Iran

Abstract

With the appearance of Web 2.0 and 3.0, users' contribution to WWW has created a huge amount of valuable expressed opinions. Considering the difficulty or impossibility of manually analyzing such big data, sentiment analysis, as a branch of natural language processing, has been highly considered. Despite the other (popular) languages, a limited number of research studies have been conducted in Persian sentiment analysis. In this

* Corresponding author

* نویسنده عهده دار مکاتبات

فصلنامه



سال ۱۳۹۷ شماره ۲ پیاپی ۳۶

study, for the first time, a semi-supervised framework is proposed for Persian sentiment analysis. Moreover, considering that one of the most recent studies in Persian, is an algorithm based on extracting adaptive (dataset-sensitive) expert-based emotional patterns. In this research, extraction of the same state-of-the-art emotional patterns is proposed to be performed automatically. Moreover, application of the HMM classifier, by utilizing the mentioned features (as its states) is analyzed; and additionally, HMM-based sentiment analysis is upgraded by being combined with a rule-based classifier for the opinion assignment process. In addition, toward intelligent self-training, a criterion for evaluating, the high reliability of output is presented by which (assuming satisfaction of the criterion) the self-training process is performed in "lexicon-extraction" and "classifier," as learning systems. The proposed method, by being applied on the basis dataset, provides 90% of accuracy (despite its expert-independent lexicon generation nature), which in comparison with the supervised and semi-supervised methods in the state-of-the-art has a considerable superiority. Moreover, this semi-supervised method is evaluated by a 10/90 ratio of train/ test and its reliability is demonstrated by providing 80% of accuracy.

Keywords: Opinion Mining, Self-training, Self-constructed Lexicon, Hidden Markov Model, Adaptive Dictionary

پرهزینه و در بسیاری از موارد غیرممکن است. خوشبختانه هوش مصنوعی و به‌طور خاص پردازش زبان طبیعی^۳ و متن‌کاوی^۴ امکان کاوش خودکار نظرات را فراهم می‌آورد. به این فرآیند، به‌طور خاص، نظرکاوی^۵ یا تحلیل احساسات متن^۶ اطلاق می‌شود که با نام‌های اندیشه‌کاوی یا عقیده‌کاوی [2] نیز شناخته شده‌است و تاکنون توجه بسیاری از پژوهش‌گران را به خود جلب کرده است.

اغلب پژوهش‌های انجام‌شده در این حوزه و همچنین مجموعه داده‌های گردآوری شده بر روی زبان‌های انگلیسی، هلندی، چینی، اسپانیایی و ترکی است [15]. در دهه گذشته بیش از دو هزار پژوهش در حوزه نظرکاوی انجام پذیرفته است [13] و همچنین تعداد اسناد نمایه‌شده در این حوزه، توسط پایگاه دانش اسکوپوس تا پایان سال ۲۰۱۶ میلادی به حدود شش‌هزار پژوهش می‌رسد. در زبان فارسی، تلاش‌های پژوهشی به‌نسبه اندکی در زمینه نظرکاوی و همچنین ایجاد و گردآوری مجموعه داده‌های استاندارد در این زمینه صورت گرفته است.

از نگاهی دیگر، پژوهش‌های صورت‌گرفته برای حل مسئله نظرکاوی به سه دسته کلی نظارتی، نیمه‌نظارتی و بدون نظارت [15] تقسیم‌بندی می‌شوند. اغلب پژوهش‌های نظرکاوی به روش نظارتی از الگوریتم‌های یادگیری ماشین بردار پشتیبان [4]، بیز ساده [17] و مدل مخفی مارکوف [24] بهره جسته‌اند. روش‌های نظارتی، مناسب برای دادگان‌هایی هستند که در آن به میزان کافی نظرات برچسب‌گذاری شده احساسی از کاربران برای آموزش یک سامانه یادگیری ماشین موجود باشد؛ روش‌های نیمه‌نظارتی که دقت کمتری از آن‌ها نسبت به دسته قبلی انتظار می‌رود، مناسب مسائلی هستند که داده‌های برچسب‌گذاری شده اندکی که فقط برای راه‌اندازی

۱- مقدمه

در سال‌های اخیر با معرفی وب ۲،۰ و ۳،۰ که بر پایه مشارکت‌ها و تعاملات استوار است، شاهد گسترش فناوری‌های وب، رونق رسانه‌های اجتماعی و افزایش تعاملات کاربران در اغلب وب‌سایت‌ها هستیم که حجم زیادی از نظرات غنی و همچنین ارزان را ایجاد کرده‌اند. درصد قابل‌توجهی از این داده‌ها به‌صورت متن و صورت‌های دیگر رسانه‌ای نظیر صدا و تصویر نگهداری می‌شوند؛ اما به‌دلیل نبود یک استاندارد همه‌جانبه و دقیق در تنظیم متون و ثبت آن‌ها، این داده‌ها طبیعتی غیرساخت‌یافته و یا نیمه‌ساخت‌یافته دارند [10]. به‌طور کلی اطلاعات موجود در اسناد متنی را به دو دسته عینی^۱ و ذهنی^۲ می‌توان تقسیم‌بندی کرد. دسته عینی شامل واقعیت‌ها، دستورهای واقعی و قابل‌مشاهده درباره موجودیت‌های مستقل و اتفاقاتی است که در جهان می‌افتد؛ اما دسته ذهنی بازتاب عواطف انسانی و یا مشاهداتی است که مردم نسبت به دنیای خارج و اتفاقات آن دارند [12].

در این میان، تحلیل و بررسی دستی اطلاعات ذهنی بسیار موردعلاقه و توجه بنگاه‌های خرد و کلان اقتصادی، سیاست‌مداران و ... است. تحلیل نظرات مشتریان اهمیت زیادی برای صاحبان تجارت جهت گرفتن بازخورد از مشتریان و مصرف‌کنندگان کالاها یا خدمات ارائه‌شده‌شان دارد [9, 14]. همچنین بررسی افکار عمومی در بین منابع موجود در وب (همچون انجمن‌ها، شبکه‌های اجتماعی و ...) در مورد یک موضوع خاص مانند انتخابات [20] و یا بورس سهام [7] برای سیاست‌مداران ارزشمند است. به‌عنوان مثالی دیگر، تحلیل نقدهای یک فیلم [16] برای متصدیان این حوزه بسیار حائز اهمیت است.

اما تحلیل و بررسی این نظرات یا داده‌های ذهنی، بسیار

³ Natural Language Processing (NLP)

⁴ Text mining

⁵ Opinion mining

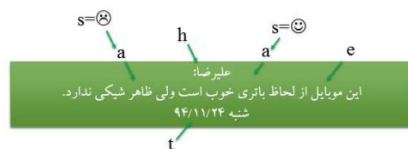
⁶ Sentiment analysis

¹ Objective

² Subjective

نظر کاوی ارائه می شود، نظر کاوی یا تحلیل احساسات، یک نظر به صورت زیر معرفی شده است:

نظر: یک نظر^۲ شامل پنج تایی (e, a, s, h, t) است که در آن e نام موجودیت، a جنبه های موجودیت، s لحن نظر بر روی جنبه a از موجودیت e، h نظر دهنده و t زمانی است که نظر دهنده h نظر خود را در مورد موجودیت e بیان می کند که در شکل (۱) قابل مشاهده است.



(شکل-۱): نمونه ای از یک نظر
(Figure-1): An example from a comment

۲-۲- چالش ها

چالش های فراوانی در حوزه نظر کاوی وجود دارد. از جمله بیان غیر صریح مانند: «برای خریدش بیشتر فکر کن»، وجود طعنه، کنایه یا شوخی مانند: «گوشی ش به درد گردو شکستن می خوره!»، زبان گفتاری (عامیانه یا غیر رسمی) مانند: «این گوشی رو به موقع از دست ندیدش» و وجود نظرات جعلی. در زبان فارسی نیز به سبب ویژگی های خاص آن و در عین حال نهادینه شدن سبک نگارش استاندارد، با چالش هایی روبه رو هستیم. نمونه هایی از این موارد عبارتند از: وجود یا عدم وجود تشدید «معین» و تنوین «واقعاً»، تنوع شیوه دگر نویسی (گردآوری / گردآوری) و پیوسته نویسی (سرهم یا با نیم فاصله) یا جدانویسی (کتاب شناسی / کتاب شناسی) و ... [3].

از دیگر چالش های این حوزه، به هزینه بر بودن تهیه داده های برچسب دار در روش های نظارتی می توان اشاره کرد. این چالش خود یکی از دلایل پدیدار شدن روش های نیمه نظارتی؛ است که هم از مزایای وجود تعداد کمی داده برچسب دار برای راه اندازی و آموزش اولیه بهره می جوید.

۲-۳- رویکردهای حل مسئله

به طور کلی برای حل مسئله نظر کاوی سه رویکرد استفاده از لغت نامه، استفاده از الگوریتم های یادگیری ماشین و رویکرد ترکیبی ارائه شده است [۱۵] که به بررسی آن ها خواهیم پرداخت.

اولیه یک سامانه به کار می آیند، وجود دارد و پس از آن، نظر کاو باید خودش را بهبود دهد. روش های بدون نظارت که دقت کمتر و یا سرعت کندتر نسبت به دو روش گذشته از آن ها انتظار می رود، مناسب برای مسائلی هستند که تنها نظرات بدون برچسب در آن ها وجود داشته باشد و چنین سامانه هایی به طور طبیعی داده ای برای یادگیری از بیرون ندارند.

تمرکز این مقاله بر روی دسته دوم یعنی روش های نیمه نظارتی است که پژوهش های زیادی در این حوزه انجام شده است؛ از جمله یکی از جدیدترین آن ها پژوهش داسیلوا و همکاران (۲۰۱۶) [22] در زبان انگلیسی است. اگرچه، در زبان فارسی، بنابر آخرین اطلاعات ما، پژوهشی در این حوزه صورت نپذیرفته است. در این پژوهش، الگویی نیمه نظارتی برای تحلیل احساسات متن ارائه می کنیم تا علاوه بر نظر کاوی با دقت بالا در زبان فارسی، دادگان^۱ با داده های برچسب خورده کم را نیز بتواند پشتیبانی کند.

به صورت خلاصه جنبه های نوآوری این مقاله به شرح زیر است:

- نخستین پژوهش در حوزه نظر کاوی نیمه نظارتی در زبان فارسی (به طور خاص در حوزه خود آموزی).
- خود کار سازی استخراج خصیصه های حسی به جای استفاده از خبره انسانی
- استفاده از چندنگاشت های عامیانه، به عنوان یک رویکرد فراکتشافی، برای پوشش دادن چالش های بررسی نشده زبان فارسی

- استفاده از ویژگی های زوج مرتبی (برچسب صرفی، برچسب احساسی) به عنوان وضعیت های طبقه بند مدل مخفی مارکوف برای افزایش دقت نظر کاوی در زبان فارسی.

ابتدا در بخش دو مروری بر پژوهش های پیشین خواهیم داشت؛ بخش سوم به معرفی روش پیشنهادی اختصاص یافته است؛ در بخش چهارم نتایج تجربی گزارش و ارزیابی شده و در نهایت در بخش پنجم به نتیجه گیری و ارائه کارهای آینده خواهیم پرداخت.

۲- مروری بر پژوهش های پیشین

۲-۱- تعاریف پایه

نظر کاوی به استخراج نظرات کاربران و تشخیص قطبیت^۲ آن ها درون متون ذهنی و نیز عینی می پردازد. اندیشه کاوی، تحلیل مستندات جهت استخراج نظرات، احساسات و خواسته های کاربران در یک حوزه خاص است. در تعریفی که در [12] برای

¹ Dataset

² Polarity

³ Opinion

۱-۳-۲- رویکرد استفاده از لغت‌نامه

در رویکرد استفاده از لغت‌نامه، از منابع کمکی برای طبقه‌بندی استفاده می‌شود این منابع می‌توانند یک فرهنگ لغت، اطلاعات فراهم‌شده توسط یک موتور جستجوگر یا مجموعه‌ای از صفات و قیود باشد. یکی از مدل‌های متنی متداول در این دسته از روش‌ها، استفاده از الگوهای ثابت مستخرج از برچسب جملات مانند جدول (۱) [11] است.

(جدول ۱): الگوهای ثابت مستخرج جهت نظرکاوی [11]
(Tabel-1): Fixed patterns, extracted for Opinion Mining [11]

واژه اول	واژه دوم	واژه سوم
JJ	NN or NNS	سایر
RB,RBR, or RBS	JJ	not NN n or NNS
JJ	JJ	not NN n or NNS
NN or NNS		not NN n or NNS
RB,RBR, or RBS	VB,VBD,VBN, or VBG	سایر
JJ - Adjective NN - Noun, singular or mass NNS - Noun, plural NNP - Proper noun, singular VB - verb, base form VBD - Verb, past tense VBG - Verb, gerund or present participle RB - Adverb RBR - Adverb, comparative RBS - Adverb, superlative		

به نظر می‌رسد نخستین کار حوزه نظرکاوی فارسی را شمس، شاکری و فیلی [18] ارائه دادند. در این مقاله از لغت‌نامه‌ای شامل ۸۰۲۷ کلمه استفاده شده که در اصل به زبان انگلیسی بوده و نویسندگان از ترجمه فارسی آن بهره جسته‌اند. در آن پژوهش دقت حدود هشتاد درصد برای طبقه‌بندی نظرات گزارش شده‌است. بصیری، نقشنیچی و آقایی [6] نیز چارچوبی برای نظرکاوی در زبان فارسی با استفاده از یادگیری غیرنظارتی ارائه کرده و با استفاده از آن ده درصد بهبود دقت را در مقایسه با روش مطرح‌شده در [۱۷] (که در بخش بعد به آن اشاره خواهیم نمود) بر مجموعه داده‌ای با نام BG Data گزارش کردند.

۲-۳-۲- رویکرد استفاده از الگوریتم‌های یادگیری ماشین

دسته دوم روش‌های حل مسئله نظرکاوی از الگوریتم‌های یادگیری ماشین و دسته‌بندی‌های مختلف مبتنی بر یادگیری ماشین استفاده می‌کنند. فرآیند یادگیری در طبقه‌بندی‌های استفاده‌شده در این دسته، با استفاده از ویژگی‌های متون و

نظرات (مانند $TF-IDF^1$) و در صورتی که یادگیری تحت نظارت^۲ باشد با استفاده از برچسب، نظر مربوط به هر متن انجام می‌شود.

سارایی و باقری [17] با تمرکز بر انتخاب ویژگی در نظرکاوی فارسی، روش نوینی برای نظرکاوی فارسی ارائه دادند که در آن (بعد از ریشه‌یابی کلمات) از ویژگی‌هایی مثل تکرار کلمات، واریانس فرکانس کلمه^۳ و ویژگی ارائه‌شده توسط اطلاعات متقابل^۴ (MI) استفاده شده و بعد از پالایش ویژگی‌ها به وسیله الگوریتم یادگیری ماشین بیز ساده (NB)^۵ مرحله آموزش ویژگی ایجاد مدل آغاز می‌شود؛ نویسندگان با پیاده‌سازی این روش بر روی دادگانی که با ۸۲۹ نظر فارسی در خصوص محصولات مختلف تلفن همراه گردآوری کردند، دقت ۸۵ درصد را ارائه دادند.

روش نیمه‌نظارتی در حوزه نظرکاوی نخستین بار در سال ۲۰۰۹ توسط رو و رابیندران [26] معرفی شد و بهترین نتایج یافت‌شده توسط بیکر، اهرات و همکارانشان [27] در سال ۲۰۱۳ با الگوریتم یادگیری Self-training و معیار دقت ۰/۶۴۱ حاصل شده‌است. تعدادی از پژوهش‌های نیمه‌نظارتی انجام‌شده در زبان انگلیسی در جدول (۲) آورده شده است.

۳-۳-۲- رویکرد ترکیبی

دسته سوم روش‌های حل مسئله نظرکاوی، روش‌های هر دو دسته یادشده قبلی را هم‌زمان و به صورت ترکیبی استفاده می‌کند. ایده اصلی روش‌های ترکیبی به کارگیری هم‌زمان مزایای رویکرد استفاده از لغت‌نامه مانند سرعت و قدرت پردازش و نیز مزایای رویکرد استفاده از یادگیری ماشین از جمله عدم وابستگی به ساختار دستور زبانی جمله است.

علیمردانی و آقایی [4] روشی ترکیبی برای نظرکاوی فارسی ارائه کرده‌اند. در این روش ابتدا معادل انگلیسی کلمات به وسیله لغت‌نامه SentiWordNet جستجو شده و سپس بر اساس یافتن معادل واژه اصلی، مثبت، منفی و یا خنثی بودن قطبیت آن گزارش می‌شود و سپس این لغات به عنوان ویژگی توسط الگوریتم ماشین بردار پشتیبان^۶ برای طبقه‌بندی نظرات مورد استفاده قرار می‌گیرد؛ بهترین نتیجه گزارش شده برای طبقه‌بندی نظراتی که توسط ایشان از داده‌های مربوط به یک تارنمای هتل گردآوری شده بود، ۸۳/۵٪ بوده است.

یکی از جدیدترین پژوهش‌های انجام‌شده در این حوزه،

¹ Term Frequency-Inverse Document Frequency (TF-IDF)

² Supervised Learning

³ Term Frequency Variance (TFV)

⁴ Mutual information (MI)

⁵ Naive Bayes (NB)

⁶ Support vector machine (SVM)

<p>در این مرحله در ابتدا بدون در نظر گرفتن پیش پردازش، برچسب گذاری اجزای کلام صورت می گیرد. کلمات با همان برچسب و چندنگاشت ها با برچسب فرضی $\langle Phrase \rangle$، بر مبنای معیاری از تعداد تکرار و دارا بودن شرایط خاص و تشخیص مثبت و یا منفی بودن، برچسب آن ها با پیشوند P_- و یا N_- ویرایش می شود.</p>	<p>آزمون</p>
<p>ممکن است، واژگانی در مرحله قبل برچسب گذاری نشده باشند؛ لذا کل جمله پیش پردازش شده و لم یابی و ریشه یابی نیز صورت می گیرد و سپس برچسب گذاری اجزای کلام انجام می شود</p>	
<p>در مرحله بعد معادل عددی هر برچسب به عنوان برداری عددی ایجاد می شود و به کمک الگوریتم خودآموز HMM تشخیص دسته بندی انجام می گیرد. چنانچه نرخ شباهت به دست آمده از مقدار آستانه ما بیشتر بود، نتیجه حاصل از الگوریتم یادگیری مان مورد تأیید قرار می گیرد؛ و در غیر این صورت دسته بر اساس قوانین سه گانه که در نظر گرفته شده است، تشخیص داده می شود.</p>	
<p>در این گام با بهره گیری از روش نیمه نظارتی خودیادگیر ($Self-training$) که در پایان هر مرحله آزمون انجام می گیرد، در صورت یکسان بودن نتایج حاصل از الگوریتم و قوانین، چندنگاشت های جمله جاری به لغت نامه وفقی پویا افزوده و بردار عددی حاصل نیز به داده های الگوریتم یادگیری ماشین افزوده می شود؛ و بازآموزی انجام و مدل جدید بازسازی خواهد شد.</p>	<p>پایانی</p>

به طور دقیق تر باید گفت در این مقاله امکان گسترش مجموعه آموزش با استفاده از نمونه های طبقه بندی شده با اطمینان بالا (خوبش آموزی) وجود دارد و نیز کلیه الگوهای حسی به صورت خودکار و هوشمند استخراج شده و نیز به جای تکنگاشت های استفاده شده در [1]، برای افزایش دقت، از چندنگاشت های حسی استفاده شده است و نیز به صورت پیش فرض از الگوریتم خودآموز HMM استفاده می کنیم (که از پتانسیل این داده های برچسب دار برای آموزش اولیه بهره می جوید) و تنها در صورتی که خروجی HMM ضریب اطمینان کافی را نداشته باشد از الگوریتم مبتنی بر قانون که ظرفیت کمتری برای دقت بالا دارد، استفاده می شود.

۳-۱- مراحل روش پیشنهادی

سامانه ارائه شده، نظر کاوی را مطابق شکل (۲)، طی سه مرحله انجام می دهد که شرح آن ها در ادامه ذکر می شود.

روشی است که آقایان اصغری، کاهانی و عسگریان [1] ارائه کرده اند. روش ایشان با استفاده از نقش های صرفی کلمات و نیز تعریف برچسب های خاص احساسی و ترکیب آن ها با یکدیگر، الگوهای کامل تری را از بیان حس در جملات نظری معرفی می کند؛ اما از آنجاکه این الگوها با توجه به نظرات کاربران در بخش آموزش و به صورت خبره آورد (مبتنی بر خبره) استخراج می شود، این روش تنها بعد از ایجاد لغت نامه احساسی خبره آورد بومی برای یک مجموعه داده، تنها بر روی همان مجموعه داده قابلیت استفاده خواهد داشت. در عوض، روش ارائه شده با آزمایش بر روی مجموعه داده آزمون ارائه شده توسط آزمایشگاه فناوری وب^۱، نرخ بالای صحت ۹۴٪ را در پی داشت.

(جدول-۲): نگاهی بر چند الگوریتم متداول در روش های

نیمه نظارتی نظر کاوی توییتهای انگلیسی [25]

(Table-2): Overview of literature on SSL for tweet sentiment analysis [25]

نام روش نیمه نظارتی	مجموعه داده	F-score
Self-training [26]	SemEval 2013	0.641
Self-training [27]	SemEval 2013	0.637
Co-training [28]	TREC 2011	در دسترس نیست
Co-training [29]	Taco Bell	نمونه های مختلف برآورد شده است

۳- روش پیشنهادی

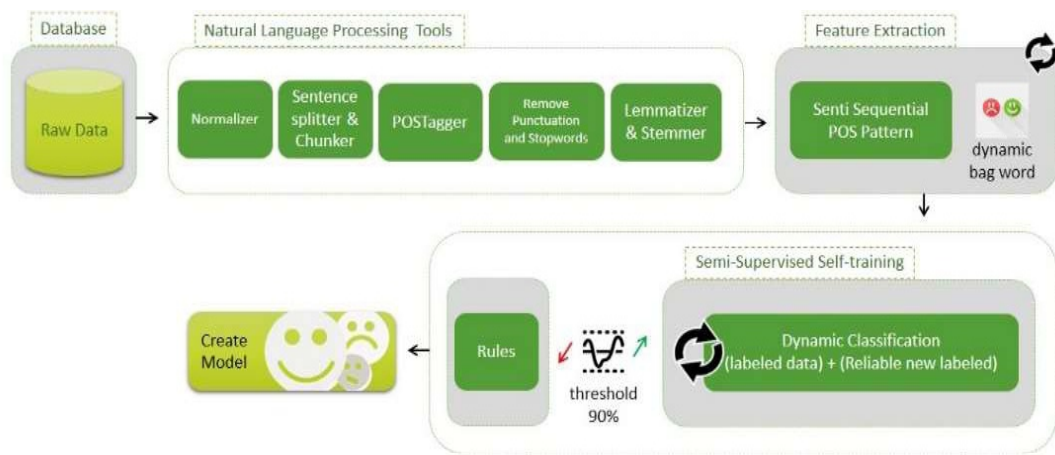
روش ارائه شده در این مقاله در سطح جمله و با رویکرد ترکیبی است که از روش های زیر بهره می برد:

(جدول-۳): خلاصه گام های روش پیشنهادی

(Table-3): Steps the proposed method

گام	هدف
آموزش	در این مرحله مجموعه داده آموزشی وارد برنامه شده و بدون پیش پردازش چندنگاشت ها (یک، دو و سه) حسی با شرایط خاص، ساخته می شوند و این واژه ها (به جز ایست واژه ها) و عبارات حسی ایجاد شده به لغت نامه وفقی پویا افزوده می شود؛ سپس پیش پردازش بر روی نظر جاری انجام می گیرد؛ سپس برداری شامل معادل عددی برچسب های تخصیص یافته، به عنوان ورودی های الگوریتم یادگیری ماشین تشکیل می شود؛ مدلی بر اساس مقادیر ورودی (بردارهای متناظر با برچسب جملات) و دسته های تعیین شده، ایجاد خواهد شد.

¹ <http://wtlab.um.ac.ir>



(شکل-۲): مدل مفهومی روش پیشنهادی
(Figure-2): Conceptual model of the proposed method

فرآیند یادگیری می‌تواند از نتایج پیش‌بینی خود هم به‌عنوان یک کمک در یادگیری استفاده کند. به همین دلیل این روش به نام خودآموزی^۵ و یا خودراه‌انداز^۶ نیز شناخته شده‌است.

Tain f from (X_t, Y_t)
Predict on $x \in X_u$
Add $(x, f(x))$ to labeled data
Repeat

(شبه‌کد-۱): روش نیمه‌نظارتی خودیادگیر
(Pseudocode-1): The Semi-supervised Self-training method

بعد از آموزش و راه‌اندازی دسته‌بند خودآموز، در مرحله آزمون، چنانچه معیار بیشترین شباهت^۷ - بر اساس محتمل‌ترین توالی بررسی‌شده در این مدل - به‌دست‌آمده از دسته‌بند، از یک مقدار آستانه بیشتر باشد (به‌نحوی که دسته‌بند طبقه‌بندی را با اطمینان بالایی انجام داده باشد) نظر موردبررسی و برچسب پیش‌بینی‌شده‌اش را می‌پذیریم. مقدار آستانه ۹۰٪ را بر اساس بررسی نرخ صحت‌های به‌دست‌آمده از اجرای روش پیشنهادی با آستانه‌های مختلف حاصل شد که در شکل (۳) قابل‌مشاهده است. با توجه به خودآموزبودن دسته‌بندمان، خروجی منتخب را به‌شرط یکسان‌بودن با نتیجه خروجی حاصل از قوانین به ورودی‌های دسته‌بند می‌افزاییم و دوباره دسته‌بند را آموزش خواهیم داد. همچنین برچسب‌های احساسی نظر جاری را مطابق آنچه در بخش ۳-۲-آمده است، به لغت‌نامه قطبیت‌دار اضافه می‌کنیم؛ اما چنانچه خروجی

۳-۱-۱- پیش‌پردازش

این مرحله شامل چندین گام مقدماتی است که قبل از ایجاد بردار (مرحله بعد) باید پیموده شوند. در این مرحله، جداسازی جملات، شناسایی کلمات، تبدیل جملات محاوره به رسمی (با استفاده از ابزار ارائه‌شده آزمایشگاه فناوری وب)، نرمال‌سازی، یکسان‌سازی نویسه‌ها و اصلاح نشانه‌گذاری‌ها انجام می‌شود.

۳-۱-۲- برچسب‌زنی احساسی اجزای کلام

در این مرحله (پس از حذف جداکننده‌ها، حروف ربط و ایست‌واژه‌ها) کلمات حاصل لم‌یابی و ریشه‌یابی می‌شوند و سپس به کمک برچسب‌زن اجزای کلام^۱ در کنار لغت‌نامه وقتی^۲ واژگان احساسی دارای برچسب قطبیت (که چگونگی ساخت آن‌ها در بخش ۳-۲-آمده است)، برچسب‌های حاوی حس را مشخص و به‌روزرسانی می‌کنیم. خروجی‌های حاصل از این مرحله، ورودی‌های دسته‌بند نیمه‌نظارتی‌مان را تشکیل خواهند داد.

۳-۱-۳- دسته‌بندی نیمه‌نظارتی و قوانین

در این مرحله به کمک یک دسته‌بند نیمه‌نظارتی که توسط داده‌های اولیه آموزش دیده شده‌است به پیش‌بینی نوع نظر، خواهیم پرداخت. یک دسته از روش‌های نیمه‌نظارتی^۳ روش‌های خودیادگیر^۴ است که به‌عنوان یکی از ایده‌دهندگان اولیه در [21] ارائه شد. ایده اصلی در این روش این است که

¹ Part-Of-Speech Tagging

² Adaptive Lexicon

³ Semi-Supervised

⁴ Self-training

⁵ Self-teaching

⁶ Bootstrapping

⁷ Maximum Likelihood

لغت نامه وفقی استفاده شده را شرح دهیم: در مرحله آموزش به کمک مجموعه داده برچسب خورده مثبت و منفی، سه ویژگی تک واژه (تک نگاشت)^۱، دو واژه (دونگاشت)^۲ و سه واژه (سه نگاشت)^۳ را با این شرط که به کلمات بازدارنده یا ایست واژه ها^۴ ختم نشود، استخراج می کنیم.

لازم به ذکر است در مرحله آموزش برای از دست ندادن کلمات یا چندنگاشت های محاوره ای افراد، از طرفی و از سوی دیگر، نظر به تنوع بسیار بالای کلمات محاوره ای در زبان فارسی، تبدیل محاوره به رسمی انجام نمی پذیرد؛ در مرحله آزمون نیز، در ابتدا عبارت ورودی، بدون انجام محاوره به رسمی، (کلمه به کلمه) نخست به صورت یک سه نگاشت در میان سه نگاشت های بررسی شده در مرحله آزمون بررسی می شود؛ در صورت عدم وجود به صورت دونگاشت بررسی شده و در صورت عدم وجود به صورت تک نگاشت بررسی می شود. اگر در هیچ کدام از این حالات چندنگاشت مورد نظر در میان چندنگاشت های موجود در مرحله آموزش یافت نشد، آنگاه تبدیل محاوره به رسمی انجام می پذیرد تا چنانچه شکل رسمی این چندنگاشت در فاز آزمون دیده شد، از قلم نیفتد. از آنجاکه ذاتاً بسیاری از چندنگاشت ها بدون ایست واژه ماهیت خود را از دست می دهند، در قسمت بررسی سه نگاشت و دونگاشت ها (که قبل از بررسی تک نگاشت ها انجام می پذیرد) ایست واژه ها هنوز در متن وجود دارند؛ اما پس از رسیدن به بررسی تک نگاشت ها ایست واژه ها (که خودشان نیز تک نگاشت هستند) نادیده گرفته می شوند.

(جدول ۵-): نمونه هایی از چندنگاشت های تولید شده در

لغت نامه وفقی

(Table-5): Examples of the generated n-grams in the adaptive dictionary

تک نگاشت	هوشمندانه، عالی، بی نظیریه، فوق العادس، مزخرف
دونگاشت	خوب هست، فاصله مناسب، ساده بود، کاهش کیفیت، راضی کننده، محشرهست، کم ندارد، رزولوشن پایین، Ram پایین، خیلی شیک
سه نگاشت	ارزشش را دارد، اصلا مشکلی ندارد، ولی رنگ صفحه، منحصر به فرد، سخت افزار قوی

همان طور که در جدول (۵) تعدادی از چندنگاشت های تولید شده را مشاهده می کنید، مواردی منحصر به فرد و به نوعی

¹ Unigram

² Bigram

³ Trigram

⁴ Stop words

دسته بند از مقدار آستانه کمتر باشد (دسته بند طبقه بندی را با اطمینان کافی انجام نداده باشد)، آنگاه نتیجه حاصل از قوانین و الگوهای تعریف شده در جدول (۴) برای خروجی پذیرفته خواهد شد.

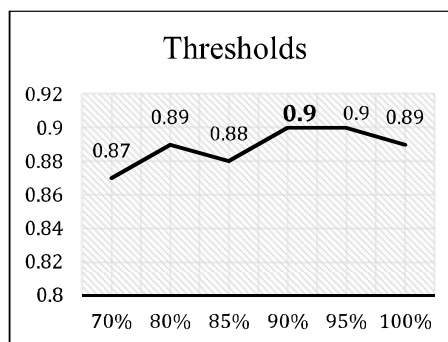
(جدول ۴-): قوانین و الگوهای حسی (برگرفته از [1] با اعمال

تغییراتی در شرط سوم)

(Table-4): Rules and emotional patterns (implied from [1] with modifications in the third rule)

عمل	شرط
به همان نتایج دسته بند اعتماد می شود.	چنانچه جمله حاوی برچسب حسی نباشد:
مثبت یا منفی بودن آن (ها) گزارش داده می شود.	اگر نظر حاوی یک برچسب حسی و یا چند برچسب حسی موافق باشد:
در صورت مشاهده نفی کننده (غیر فعل) نخستین برچسب حسی بعد آن معکوس خواهد شد؛ در صورت مشاهده فعل منفی قطبیت کلی معکوس خواهیم کرد؛ در سایر حالات بسته به نوع واژه و یا عبارت، قطبیت منفی و یا مثبت در نظر گرفته می شود؛ در نهایت مقدار کلی را گزارش خواهیم کرد.	اگر جمله حاوی چندین برچسب ناموافق باشد:

در ادامه این بخش، مدل زبانی پیشنهادی بر اساس لغت نامه وفقی و همچنین روش های دسته بند استفاده شده به طور مبسوط شرح داده خواهد شد.



(شکل ۳-): انتخاب بهترین آستانه بر اساس نرخ صحت

به دست آمده از آستانه های متفاوت

(Figure-3): Choose the best threshold based on the accuracy obtained from different thresholds

۲-۳- لغت نامه وفقی و مدل زبانی پیشنهادی

قبل از معرفی مدل زبانی پیشنهادی، لازم است روش ایجاد

مبتنی بر نوع اطلاعات مجموعه داده مورد استفاده ایجاد شده است که قطعاً بسیاری از این چندنگاشت‌ها و حتی کلمات در لغت‌نامه‌های معتبر وجود ندارند.

لذا مدل زبانی مورد نیاز برای حل مسئله به این صورت است که در ابتدا وضعیت‌ها را به صورت یک دوتایی (برچسب صرفی، احساس) در نظر می‌گیریم که برچسب صرفی یکی از حالات محدود فعل، اسم، صفت، قید ... بوده و احساس نیز یکی از سه حالت مثبت، منفی و خنثی است. تمامی دوتنگاشت و سه‌نگاشت‌ها برای پیش‌گیری از زیاد شدن تعداد وضعیت‌ها از برچسب صرفی‌ای به نام «چندنگاشت» بهره می‌جوئیم. درواقع، تحلیل هر چندنگاشت به‌طور مجزا (در مرحله آموزش و آزمون) تنها به دلیل شناسایی و حصول اطمینان نسبت به برچسب احساس آن است و چندنگاشت‌ها وضعیت‌های جداگانه‌ای ایجاد نمی‌کنند و سپس با بررسی چندنگاشت‌ها و کلمات در لغت‌نامه وقتی مشابه شبه‌کد ۲، برچسب را بر اساس نسبت جمع برداری منفی/مثبت‌ها بر جمع اسکالر آن‌ها، با علامتی مثبت، منفی و یا خنثی علامت‌گذاری می‌کنیم. در همین بخش منفی‌کننده‌ها^۱ بر اساس کلمه ماقبل مورد بررسی قرار می‌گیرند و در صورت وجودشان برچسب موجود در دادگان معکوس می‌شود. (افعال منفی که با حرف «ن» شروع می‌شوند، ولی ریشه آن با غیر «ن» است و واژه‌های «بی»، «بدون» و «نا»، منفی‌کننده امتیاز خواهند بود).

دسته نخست چندنگاشت‌هایی هستند که بسیار تکرار شده‌اند اما تعداد استفاده آن‌ها در جملات مثبت و منفی سربه‌سر بوده است. در این نوع چندنگاشت‌ها رابطه $\frac{pos-neg}{pos+neg}$ عددی کوچک منفی یا مثبت، خروجی می‌دهد. چند مثال از این دسته عبارت‌اند از فعل «است»، راستی»، «شاید» و ...

دسته دوم چندنگاشت‌هایی هستند که خیلی کم (به‌عنوان مثال یک‌بار) در کل دادگان تکرار شده‌اند و به‌طور تصادفی تکرارشان در یک یا چند جمله مثبت یا منفی بوده است. در این حالات رابطه $\frac{pos-neg}{pos+neg}$ عدد بزرگ (یک به ۱) خروجی می‌دهد که به نظر خیلی مثبت یا خیلی منفی می‌آید؛ اما این مثبت یا منفی بودن شدید، معنادار نیست چراکه به‌خاطر کمبود داده این اتفاق رخ داده‌است. مثال‌های این دسته عبارت‌اند از اسم شهر «اردکان»، کلمه «زادروز»، نام «داریوش» و ...

بر اساس منطق بالا بازه ۱- تا ۱ به پنج بخش تقسیم شد و بخش وسط (به‌خاطر شامل بودن دسته نخست از خنثی‌ها) و دو بخش کناری (به‌خاطر شامل بودن بخش دوم از خنثی‌ها) به‌عنوان خنثی در نظر گرفته شد و دو بخش دیگر برحسب منفی یا مثبت بودن برچسب منفی یا مثبت دریافت کردند؛ لذا مطابق جدول (۶) برچسب‌گذاری می‌کنیم:

(جدول-۶): بازه‌های تعیین شده جهت برچسب‌گذاری

چندنگاشت‌های موجود در لغت‌نامه وقتی

(Table-6): The determined ranges for labeling n-grams in the Adaptive Dictionary

برچسب	بازه	عنوان
خنثی	$[+0.6, +1]$	مقادیر مرزی با امتیاز خیلی مثبت
	$[-1, -0.6]$	مقادیر مرزی با امتیاز خیلی منفی
	$[0.2, +0.2]$	مقادیر نزدیک به صفر
مثبت	$[+0.2, +0.6]$	مقادیر میانی مثبت
منفی	$[-0.6, -0.2]$	مقادیر میانی منفی

بعد از ساخت لغت‌نامه وقتی با استفاده از داده‌های آموزش، خواهیم توانست بردار ویژگی داده‌های آزمون را بر اساس توالی برچسب‌های به‌وجودآمده تهیه کنیم.

۳-۳- روش‌های دسته‌بندی استفاده‌شده

مدل مخفی مارکوف^۲ در اواخر دهه ۱۹۶۰ میلادی معرفی شد. این مدل، درعمل گزینه مناسبی برای حل مسائل ترتیبی بوده است؛ ازهمین‌رو در مسائل پردازش زبان طبیعی نیز مورد استفاده قرار می‌گیرد. به‌عنوان نمونه در [19] یک توسعه

```

Input = myNGram, myNGramTag, previousWord
pos,neg=
findPolarityFromDics(myNGram);
score = (pos - neg) / (pos + neg);
if (previousWord in Negation)
    score *= -1;
if(score>0.2 and score<0.6)
    myNGramTag = "Pos"+myNGram;
else if(score<-0.2 and score>-0.6)
    myNGramTag = "Neg"+myNGram;
else
    myNGramTag = "Neu"+myNGram;
Output = myNGramTag
    
```

(شبه‌کد-۲): الگوریتم تعیین قطبیت چندنگاشت بر اساس تعداد

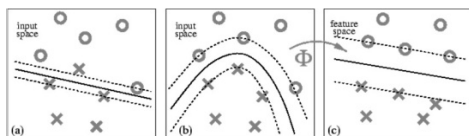
تکرار در لغت‌نامه وقتی

(Pseudocode-2): Algorithm of determining n-grams polarity, based on observations frequency in the Adaptive Dictionary

بر اساس تحلیلی که روی چندنگاشت‌ها صورت دادیم، به‌طورمعمول با دو دسته کلی از چندنگاشت‌های خنثی مواجه می‌شویم.

² Hidden Markov Models (HMM)

¹ Negations



(شکل-۵): شمایی از دسته بندی SVM مبتنی بر Kernel
(Figure-5): A schematic view of classification in Kernel-based SVM

۳-۴- معیارهای ارزیابی

جهت ارزیابی روش های پیشنهادی در این پژوهش، از چهار معیار دقت^۲، فراخوانی^۳، معیار F^4 و نرخ صحت^۵ بر مبنای ماتریس درهم ریختگی^۶ مطابق جدول (۷) استفاده شده است.

(جدول-۷): حالات ممکن صحت عمل یک دسته بند برای هر

طبقه TN درست منفی، TP درست مثبت، FN اشتباه منفی، FP اشتباه مثبت

(Table-7): Possible cases for output validity of a classifier in each class (TP: true positive, TN: true negative, FP: false positive, FN: false negative)

		دسته پیش بینی شده	
		مثبت	منفی
دسته واقعی	مثبت	True positive (TP)	False negative (FN)
	منفی	False positive (FP)	True negative (TN)

معیار دقت: تعداد مواردی مثبتی که به درستی پیش بینی شده اند، نسبت به تعداد کل موارد درست پیش بینی شده.

$$Precision = \frac{TP}{TP + FP} \quad (1)$$

معیار فراخوانی: تعداد مواردی مثبتی که به درستی پیش بینی شده اند از تعداد کل موارد مثبت.

$$Recall = \frac{TP}{TP + FN} \quad (2)$$

معیار F: این معیار در ارتباط با معیارهای دقت و فراخوانی به صورت میانگین وزن دار، از هر دو معیار، به کار می رود. (Recall=R, Precision=P)

$$FScore = \frac{2P \cdot R}{P + R} \quad (3)$$

نرخ صحت: میزان مقادیر درست پیش بینی شده نسبت به تعداد کل مجموعه داده، بیانگر نرخ صحت است؛ با توجه به اینکه این معیار تنها معیاری است (از بین موارد اشاره شده) که

² Precision

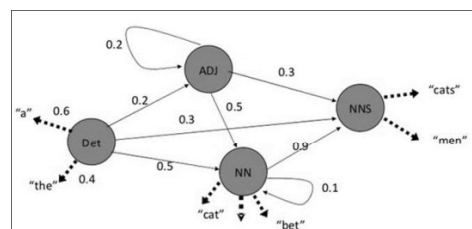
³ Recall

⁴ F- Measure (FScore)

⁵ Accuracy

⁶ Confusion Matrix

بر مدل مخفی مارکوف با استفاده از تخمین مرتبه دوم برای برچسب گذاری جملات مورد استفاده قرار گرفته است. یکی از مهم ترین الگوریتم های یادگیری بر اساس رخداد مشاهدات در مدل مخفی مارکوف، الگوریتم بام-ولش^۱ است [23]؛ این روش را به سادگی و با محاسبه احتمال رخداد پارامترها می توان تعریف کرد و یکی از ویژگی های مخصوص این الگوریتم این است که هم گرایی در آن تضمین شده است. این دسته بند در جهت حل مسائل پردازش زبان طبیعی که ارتباط مستقیم با توالی کلمات دارد از جمله برچسب زن اجزای کلام مورد استفاده قرار گرفته است [5].



(شکل-۴): نمونه ساده از کاربرد HMM در حل مسئله POS-

Tagger
(Figure-4): A simple example for HMM application in solving POS-Tagger problem

ماشین بردار پشتیبان در سال ۱۹۹۲ میلادی معرفی شد. این دسته بند نوع خاصی از مدل های خطی را می یابد که موجب بیشینه شدن تفکیک بین طبقات می شود. بردارهای پشتیبان درواقع نزدیک ترین نقاط به حاشیه ابرصفحه هستند. در مسائلی که داده ها به صورت خطی جداپذیر نباشند، یک Kernel (از نوع خطی Linear، چندجمله ای Polynomial، گوسی Gaussian و یا ...) برای آن در نظر گرفته می شود که در عمل نگاشتی را بین خط با آن ابرصفحه جداکننده برقرار کند؛ یا به بیان دیگر، داده ها به فضای با ابعاد بیشتر نگاشت پیدا می کنند تا بتوان آن ها را در این فضای جدید به صورت خطی جدا کرد؛ و با این شرایط SVM را، از این منظر، غیر خطی نیز می توان در نظر گرفت. به عنوان نمونه در تصویر زیر با پارامتر ϕ فضای اولیه داده ها به فضای ویژگی هایی نگاشت می شود که با همان SVM خطی می توان دسته بندی کرد.

در این پژوهش از دسته بند مدل مخفی مارکوف (به دلیل مناسب بودن آن برای مشاهدات ترتیبی و مدل زبانی پیشنهادی این پژوهش) به عنوان دسته بند اصلی و از دسته بند ماشین بردار پشتیبان (به دلیل بهترین عملکرد در بین سایر دسته بندها مورد استفاده در سایر پژوهش های نظر کاوی در زبان فارسی) بهره گرفته شده است.

¹ Baum-Welch

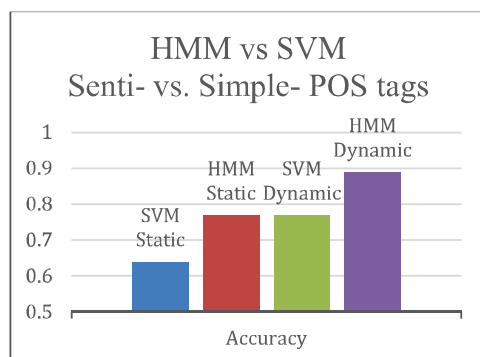
(جدول-۸): جزئیات مجموعه داده مورد استفاده

(Table-8): Details of the utilized dataset

تعداد نظرات مجموعه آموزش	1900
تعداد نظرات مجموعه آزمون	229
برچسب‌ها	مثبت و منفی
میانگین تعداد کلمات هر جمله	8.39
تعداد جملات مثبت	۱۶۸۴ (آموزش = ۱۵۰۷ و آزمون = ۱۷۷)
تعداد جملات منفی	۴۳۳ (آموزش = ۳۸۷ و آزمون = ۴۶)

۴-۲- نتایج

در ابتدا برای انتخاب دسته‌بند مناسب برای مدل متنی ارائه شده، دسته‌بند SVM را با دسته‌بند پیشنهادی (HMM) به کمک بردار متناظر با برچسب‌های مستخرج از دادگان احساسی ایستا واژگان فارسی دارای برچسب قطبیت تهیه شده در آزمایشگاه سامانه‌های هوشمند اطلاعات دانشگاه تهران توسط [8] و همچنین برچسب‌های مستخرج از دادگان احساسی پویا (وقفی) مورد ارزیابی قرار داده‌ایم؛ سپس با آزمایش بر روی دادگان نظرات مبنای، با تخصیص آموزش/آزمون به همان شیوه پژوهش مورد مقایسه^۳ نتایج پیش‌رو حاصل شد:



(شکل-۶): مقایسه عملکرد دسته‌بند SVM و HMM با به کارگیری

واژه‌نامه ایستا و پویا

(Figure-6): HMM classifier vs. SVM classifier & static dictionary vs. dynamic dictionary

همان‌گونه که در شکل (۶) قابل مشاهده است، برچسب‌های احساسی پویا هم در مورد HMM و هم در مورد SVM نسبت به برچسب‌های احساسی ایستا نتایج را به طور قابل ملاحظه‌ای بهبود داده‌اند که نشان‌دهنده کارایی برچسب‌های احساسی است. همچنین قابل مشاهده است که طبقه‌بند HMM در هر دو حالت، با کارایی بالاتری، نسبت به

^۳ در پژوهش مورد مقایسه ۹۰ درصد داده‌ها برای آموزش و ۱۰ درصد برای آزمون به کار بسته شده‌اند.

در آن علاوه بر مثبت‌های درست تشخیص داده شده، منفی‌های درست تشخیص داده شده نیز ملاک ارزیابی قرار می‌گیرند، این ملاک را ملاک اصلی مقایسه قرار داده‌ایم.

$$Accuracy = \frac{TP + TN}{TP + TN + FN + FP} \quad (۴)$$

سه ملاک بازیابی اطلاعاتی Precision، Recall و FScore در نظرکاوای جنبه‌گرا و غیرجنبه‌گرا دو مفهوم متفاوت بااهمیت‌های متفاوت را نشان می‌دهند؛ در نظرکاوای جنبه‌گرا این مفاهیم نشان‌گر این هستند که چه میزان جنبه‌هایی که نویسنده در مورد آن صحبت کرده است، به درستی بازیابی شده‌اند که همواره اطلاعات ارزشمندی است و حکایت از صحت تشخیص جنبه در الگوریتم نظرکاوای دارد؛ اما در نظرکاوای غیرجنبه‌گرا که این پژوهش نیز در این حوزه تعریف شده، این سه معیار میزان شناسایی درست متون با برچسب نظر مثبت است. از آنجاکه در دادگان مورد ارزیابی (دیجی کالا) شناسایی درست نظرات مثبت مزیتی نسبت به شناسایی درست نظرات منفی نداشتند، این معیار فقط در نمودار تشریح پایداری درج شده است تا پایستاربودن تمامی ملاک‌ها (شامل ملاک‌های یادشده) را نشان دهد.

ثبات و پایداری^۱: جهت نمایش قدرت روش پیشنهادی، پایداری آن را با در نظر گرفتن اندازه کوچک‌تری از مجموعه آموزش و به تناسب، اندازه بزرگ‌تر مجموعه آزمون مورد بررسی قرار می‌دهیم.

۴- نتایج تجربی

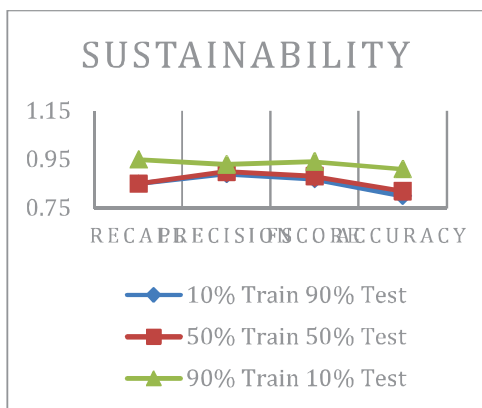
۴-۱- معرفی مجموعه داده

جهت ارزیابی روش ارائه شده، از آنجاکه این مقاله نسخه خودکار شده (بدون نیاز به دادگان احساسی وقفی خبره آورد) پژوهش [1] است، مجموعه داده آن مطالعه را که توسط نویسندگان مقاله ارائه شده بود، مجموعه داده مرجع این پژوهش قرار دادیم. دادگان مزبور بخشی از مجموعه داده گردآوری شده توسط آزمایشگاه فناوری وب دانشگاه فردوسی از پایگاه وب دیجی کالا^۲ است که با توجه به نوع نظردهنده (نظر خبره، کاربر فعال و کاربر عادی) شامل نظرات رسمی، نیمه رسمی و عامیانه است. جزئیاتی از این مجموعه داده به شرح زیر است:

^۱ Sustainability

^۲ <http://www.digikala.com>

نتایج حاصل در شکل (۸)، ثبات و پایداری الگوریتم با وجود حجم داده آموزشی کم مشهود است.



(شکل-۸): بررسی پایداری روش نیمه نظارتی پیشنهادی با کاهش نسبت مجموعه آموزش

(Figure-8): Sustainability of the proposed semi-supervised method with decreasing training data ratio

۵- نتیجه گیری و کارهای آینده

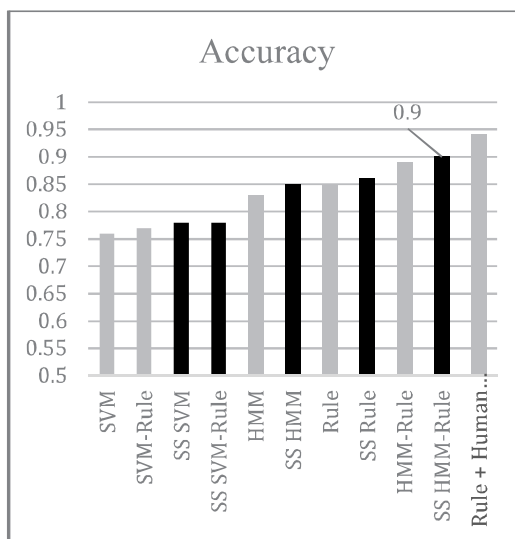
در این مقاله، یک مدل استخراج ویژگی برای نظرکاوی مبتنی بر دادگان احساسی حساس به مجموعه داده (وفقی) و سپس یک طبقه‌بند مناسب برای آن (مدل مخفی مارکوف خودآموز، HMM، در کنار یک طبقه‌بند مبتنی بر قانون به عنوان درجه اطمینان HMM) ارائه شد. مقایسه‌ها نمایانگر برتری دادگان احساسی وفقی نسبت به دادگان احساسی ایستا و همچنین برتری طبقه‌بند HMM نسبت به طبقه‌بند ماشین بردار پشتیبانی (SVM) و همچنین بهره‌وری فرآیند خودآموزی در HMM، در فرآیند نظرکاوی است. با توجه به پیچیدگی زمانی روش HMM، نظرکاوی ارائه شده در این مقاله مناسب نظرکاوی برای نظرات نسبتاً کوتاه (مثل نظرات کاربران در وبگاه‌ها، یا پست‌های کاربران در میکرو بلاگ‌ها) است.

نکته قابل توجه دیگر این پژوهش، این است که بی‌نیاز شدن به عنصر خبره انسانی مورد نیاز در لغتنامه وفقی مورد استفاده در مرزهای دانش، دقت نظرکاوی با اختلاف جزئی ۴٪ حفظ شده است که این در مقایسه با روش‌های نیمه نظارتی و نظارتی موجود یک بهبود قابل ملاحظه است.

در نسخه نهایی این پژوهش در حال انجام که به وسیله آن بر آن هستیم روشی حتی با کارایی بالاتر از دادگان احساسی خبره‌آورد دهیم، هم الگوریتم ایجاد دادگان احساسی و هم قانون‌های مورد استفاده در طبقه‌بندی فازی خواهند شد؛ به طوری که با پوشش عدم قطعیت موجود در سامانه بتوانیم فاصله عامل‌های کارایی روش ارائه شده را با روش خبره‌آورد پر کنیم و نیز جملات طولانی تقطیع جمله و به بررسی جداگانه

طبقه‌بند SVM فرآیند نظرکاوی را انجام داده است.

در آزمایش بعد، با برگزیدن HMM به همراه برجسب‌های احساسی اثر فرآیند خودآموزی را مورد مطالعه قرار می‌دهیم. شکل (۷) کارایی طبقه‌بند HMM را یکبار بدون فرآیند خودآموزی و یکبار به همراه خودآموزی نشان داده است. مشاهده می‌شود، فرآیند خودآموزی باعث رشد تمامی مؤلفه‌های کارایی گردیده است. یکی از مزایای الگوریتم ارائه شده این است که در مرحله آموزش درعمل دادگان تخصصی پویا به صورت وفقی (سازگار با مجموعه خاصمان) ایجاد می‌شود؛ برای بررسی این مزیت، دادگان حاضر را با مجموعه دادگان حساس معتبر ایستا، معرفی شده در [8] مورد بررسی و مقایسه قرار دادیم. نکته قابل ملاحظه دیگر این است که با وجود شرایط یکسان، اگرچه دقت روش پیشنهادی ما از روش مبتنی بر قانون و دادگان احساسی خبره آورد [1] کمتر است، اما از طرفی با این کاهش جزئی کارایی، عنصر پرهزینه خبره انسانی برای استخراج دادگان احساسی متناسب با مجموعه داده، از الگوریتم حذف شده است و از طرف دیگر این کاهش فقط به اندازه ۴٪ بوده که نتایج بسیار نزدیکی نسبت به حالت وجود خبره است.



(شکل-۷): مقایسه صحت نظرکاوی روش‌های مختلف

اعمال شده بر مجموعه داده مبنا (روش‌های نیمه نظارتی پررنگ تر رسم شده‌اند)

(Figure-7): Comparison of the opinion mining accuracy in different methods on the benchmark dataset (semi-supervised methods are highlighted)

در انتها، برای بررسی ثبات و پایداری روش نیمه نظارتی ارائه شده را بر روی مجموعه داده مبنا با اندازه‌هایی متفاوتی از مجموعه آموزش-آزمون مورد بررسی قرار دادیم که با توجه به

و پژوهشی کتابداری و اطلاع رسانی - آستان قدس رضوی، ۱۳۹۱.

- [3] H. Sotudeh and Z. Honarjoooyan, "A review on Persian challenges in digital paradigms, and their effect on efficiency of automatic text processing and information retrieval," Library and Information Science, 15 (4), Astan Quds Razavi. 2013
- [4] S. Alimardani and A. Aghaei, "Opinion Mining in Persian Language Using Supervised Algorithms," 2015.
- [5] A. Azimzadeh, M. M. Arab, and S. R. Quchani, "Persian part of speech tagger based on Hidden Markov Model," 9th JADT, 2008.
- [6] M. E. Basiri, A. R. Naghsh-Nilchi, and N. Ghassem-Aghace, "A Framework for Sentiment Analysis in Persian," 2014.
- [7] J. Bollen, H. Mao, and X. Zeng, "Twitter mood predicts the stock market," J. Comput. Sci., vol. 2, no. 1, pp. 1-8, 2011.
- [8] I. Dehdarbehbahani, A. Shakery, and H. Faili, "Semi-supervised word polarity identification in resource-lean languages," Neural Networks, vol. 58, pp. 50-59, 2014.
- [9] M. Gamon, "Sentiment classification on customer feedback data: noisy data, large feature vectors, and the role of linguistic analysis," in Proceedings of the 20th international conference on Computational Linguistics, 2004, p. 841.
- [10] V. Gupta and G. S. Lehal, "A survey of text mining techniques and applications," J. Emerg. Technol. web Intell., vol. 1, no. 1, pp. 60-76, 2009.
- [11] A. K. Jain and Y. Pandey, "Analysis and implementation of sentiment classification using lexical POS markers," Int. J., vol. 2, no. 1, 2013.
- [12] B. Liu, "Sentiment analysis and opinion mining," Synth. Lect. Hum. Lang. Technol., vol. 5, no. 1, pp. 1-167, 2012.
- [13] B. Liu, "Sentiment analysis: Mining opinions, sentiments, and emotions": Cambridge University Press, 2015.
- [14] J. Liu, Y. Cao, C.-Y. Lin, Y. Huang, and M. Zhou, "Low-Quality Product Review Detection in Opinion Summarization," in EMNLP-CoNLL, 2007, pp. 334-342.
- [15] R. Kumar and R. Vadlamani, "A survey on opinion mining and sentiment analysis: tasks, approaches and applications," Knowledge-Based Syst., vol. 89, pp. 14-46, 2015.
- [16] E. Sadikov, A. Parameswaran, and P. Venetis, "Blogs as predictors of movie success," 2009

آن‌ها پیشنهاد می‌شود و البته در این شرایط تنها باید برجسب‌های پراهمیت‌تر در نظر گرفته شوند تا از رشد زیاد لغت‌نامه پویا جلوگیری شود.

سپاس‌گزاری

در پایان از آقایان دکتر کاهانی و اصغری کمال سپاس دارم که ابزارهای پردازش زبان طبیعی آزمایشگاه فناوری وب دانشگاه فردوسی^۱ را در اختیارمان قرار دادند تا به کمک آن در کنار کتابخانه متن باز Accord.NET^۲ پیاده‌سازی مقاله را انجام دهیم؛ و همچنین از آقای دکتر اکبر زاده توتونچی و نیز آقای سید علی حسینی (دانشجوی دکترای ایشان) به دلیل یاری و راهنمایی‌های بی‌دریغشان و همچنین از خانم دکتر شاکری به دلیل در اختیار دادن واژگان قطبیت‌دار فارسی سپاس‌گزاریم.

6- References

۶- مراجع

- [۱] سید محمد اصغری نکاح، محسن کاهانی و احسان. عسگریان. «نظرکاوی با استفاده از برجسب‌های صرفی و معنایی و کشف روابط حسی جملات فارسی». بیستمین کنفرانس ملی سالانه انجمن کامپیوتر ایران. دانشگاه فردوسی مشهد. ۱۳۹۴
- [1] S.M. Asghari N, M. Kahani, and E. Askarian, "Opinion Mining by means of syntactic and semantic labels, and discovering emotional relations in Persian sentences". Computer Society of Iran (20th) Computer Conference (CSICC 2015). Ferdowsi University of Mashhad. 2015
- [۲] برهانی زرنندی، سمیه، علی اکبر نیک نفس، و مجید محمدی. "عقیده کاوی در نقد کالا با استفاده از شبکه واژگان احساسی"، دومین کنفرانس ملی مهندسی صنایع و سیستم ها، نجف آباد، دانشگاه آزاد اسلامی واحد نجف آباد، گروه مهندسی صنایع، ۱۳۹۲
- [2] S. Borhani Z., A.A. Niknafs, and M. Mohammadi, "Opinion mining in product reviews, using emotional vocabulary network" 2nd National Conference on Industrial Engineering & Systems (NIESC 2014). Najafabad branch, Islamic Azad University. 2014
- [۳] هاجر ستوده، زهره هنرجویان. "مروری بر دشواری‌های زبان فارسی در محیط دیجیتال و تأثیرات آنها بر اثر بخشی پردازش خودکار متن و بازیابی اطلاعات". فصلنامه علمی
- ^۱ ابزارهای پردازش متون زبان فارسی، آزمایشگاه فناوری وب دانشگاه فردوسی مشهد (۱۳۹۱-۳۹۴) (wtlab.um.ac.ir)
- ^۲ Accord.NET Framework, <http://accord-framework.net/>

- [29] S.Liu, W.Zhu, N.Xu, F.Li, X. Q.Cheng, Y.Liu, & Y.Wang, "Co-training and visualizing sentiment evolution for tweet events". In Proceedings of the 22nd International Conference on World Wide Web (pp. 105-106). ACM. 2013



محسن نجف زاده، کارشناسی خود را در رشته مهندسی کامپیوتر خرم افزار در سال ۱۳۸۷ و مدرک کارشناسی ارشد را در دانشگاه آزاد اسلامی واحد مشهد در رشته مهندسی کامپیوتر - هوش مصنوعی در سال ۱۳۹۵ به پایان رساند. زمینه پژوهشی وی پردازش زبان طبیعی، نظار کاوی و شبکه های عصبی است. نشانی رایانامه ایشان عبارت است از:

mohsen.najafzadeh@mshdiau.ac.ir



سعید راحتی قوچانی، در سال ۱۳۶۸ کارشناسی خود را در رشته الکترونیک دانشگاه تهران به پایان رساند و کارشناسی ارشد و دکترای خود را در رشته مخابرات به ترتیب در واحد تهران جنوب و واحد علوم و تحقیقات دانشگاه آزاد اسلامی در سال های ۱۳۷۲ و ۱۳۷۷ اخذ کرد. سپس با پایان دوره دکترا، به عنوان استادیار و از سال ۱۳۹۰ در سمت دانشیاری دانشگاه آزاد اسلامی واحد مشهد مشغول به خدمت شد. تاکنون بیش از یکصد مقاله در نشریات و کنفرانس ها از ایشان به چاپ رسیده است. گرایش پژوهشی ایشان پردازش زبان و گفتار و آموزش شبکه های عصبی و کاربرد آن در مدل سازی سامانه های زیستی است. نشانی رایانامه ایشان عبارت است از:

rahati@mshdiau.ac.ir



رضا قائمی، مدرک کارشناسی خود را در رشته مهندسی کامپیوتر در سال ۱۳۷۶ و همچنین مدرک کارشناسی ارشد را در رشته مهندسی کامپیوتر- نرم افزار در سال ۱۳۷۹ به اتمام رساند؛ سپس مدرک دکترای خود را در سال ۱۳۹۱ در رشته هوش مصنوعی از دانشگاه یوپایم مالزی اخذ کرد. وی هم اکنون عضو هیئت علمی و استادیار دانشگاه آزاد اسلامی واحد قوچان است. زمینه پژوهشی ایشان هوش محاسباتی و داده کاوی است. نشانی رایانامه ایشان عبارت است از:

r.ghaemi@iauu.ac.ir

- [17] M. Saraee and A. Bagheri, "Feature selection methods in Persian sentiment analysis," in Natural Language Processing and Information Systems, Springer, 2013, pp. 303-308.
- [18] M. Shams, A. Shakery, and H. Faili, "A non-parametric LDA-based induction method for sentiment analysis," in Artificial Intelligence and Signal Processing (AISP), 2012 16th CSI International Symposium on, 2012, pp. 216-221.
- [19] S. M. Thede and M. P. Harper, "A second-order hidden Markov model for part-of-speech tagging," in Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics, 1999, pp. 175-182.
- [20] A. Tumasjan, T. O. Sprenger, P. G. Sandner, and I. M. Welp, "Predicting Elections with Twitter: What 140 Characters Reveal about Political Sentiment," ICWSM, vol. 10, pp. 178-185, 2010.
- [21] H. Scudder, "Probability of error of some adaptive pattern-recognition machines," IEEE Transactions on Information Theory, vol. 11, pp. 363-371, 1965.
- [22] N. F. F. da Silva, L. F. Coletta, E. R. Hruschka, and E. R. Hruschka Jr, "Using unsupervised information to improve semi-supervised tweet sentiment classification," Information Sciences, vol. 355, pp. 348-365, 2016.
- [23] L. R. Welch, "Hidden Markov models and the Baum-Welch algorithm," IEEE Information Theory Society Newsletter, vol. 53, pp. 10-13, 2003.
- [24] M. Kang, J. Ahn, and K. Lee, "Opinion mining using ensemble text hidden Markov models for text classification." 2017.
- [25] N. F. F. D. Silva, L. F. Coletta, & E. R. Hruschka, "A survey and comparative study of tweet sentiment analysis via semi-supervised learning." ACM Computing Surveys (CSUR), 49(1), 15, 2016.
- [26] D. Rao, and D. Ravichandran, "Semi-supervised polarity lexicon induction" in Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics (pp. 675-682). Association for Computational Linguistics. 2009.
- [27] L. Becker, G. Erhart, D. Skiba, and V. Matula, "AVAYA: Sentiment Analysis on Twitter with Self-Training and Polarity Lexicon Expansion". SemEval@ NAACL-HLT, pp. 333-340, 2013.
- [28] S.Liu, F.Li, F.Li, X.Cheng, & H.Shen, "Adaptive co-training SVM for sentiment classification on tweets". In Proceedings of the 22nd International Conference on World Wide Web Information & Knowledge Management (pp. 2079-2088). ACM. 2013.

