

# حسن نگار : شبکه واژگان حسی فارسی



احسان عسکریان<sup>\*1</sup>، محسن کاهانی<sup>2</sup> و شهلا شریفی<sup>3</sup>

<sup>1</sup> و <sup>2</sup> آزمایشگاه فناوری وب، دانشکده مهندسی، دانشگاه فردوسی مشهد، مشهد، ایران

<sup>3</sup> گروه آموزشی زبان شناسی همگانی، دانشکده ادبیات و علوم انسانی، دانشگاه فردوسی مشهد، مشهد، ایران

## چکیده

وظیفه اصلی نظرکاوی استخراج و تشخیص حس مثبت یا منفی (رضایت مندی) افراد، از روی اطلاعات متنی است. نبود یک واژه نامه حسی فارسی عامل یکی از چالش های اصلی نظرکاوی در زبان فارسی است. در این مقاله روشی جدید برای تولید شبکه واژگان حسی فارسی (حسن نگار) با استفاده از منابع زبانی فارسی و انگلیسی ارائه می شود. همچنین پیکره نظرات فارسی ایجاد شده برای انجام پژوهش های نظرکاوی، معرفی خواهند شد. برای تولید حسن نگار ابتدا شبکه واژگان جامع زبان فارسی (فردوس نت) ساخته شده است؛ سپس میزان حس هر گروه هم معنی در شبکه واژگان حسی انگلیسی به کلمات متناظر آنها در حسن نگار (شبکه واژگان حسی فارسی) نگاشت می شود. در آزمایش های انجام شده، مشخص شد که حسن نگار دارای دقت ۰/۸۶ و نرخ بازیابی ۰/۷۵ است و می تواند به عنوان واژه نامه حسی مرجع برای زبان فارسی استفاده شود.

واژگان کلیدی: نظرکاوی، شبکه واژگان فردوس نت، واژه نامه حسی، ابزارهای پردازش متون زبان فارسی.

## HesNegar: Persian Sentiment WordNet

Ehsan Asgarian<sup>\*1</sup>, Mohsen Kahani<sup>2</sup> & Shahla Sharifi<sup>3</sup>

<sup>1,2</sup>Department of Computer Engineering, Ferdowsi University of Mashhad, Mashhad, Iran

<sup>3</sup>Department of Linguistics, Faculty of Letters and Humanities, Ferdowsi University of Mashhad, Mashhad, Iran

### Abstract

Awareness of others' opinions plays a crucial role in the decision making process performed by simple customers to top-level executives of manufacturing companies and various organizations. Today, with the advent of Web 2.0 and the expansion of social networks, a vast number of texts related to people's opinions have been created. However, exploring the enormous amount of documents, various opinion sources and opposing opinions about an entity have made the process of extracting and analyzing opinions very difficult. Hence, there is a need for methods to explore and summarize the existing opinions. Accordingly, there has recently been a new trend in natural language processing science called "opinion mining". The main purpose of opinion mining is to extract and detect people's positive or negative sentiments (sense of satisfaction) from text reviews. The absence of a comprehensive Persian sentiment lexicon is one of the main challenges of opinion mining in Persian.

In this paper, a new methodology for developing Persian Sentiment WordNet (HesNegar) is presented using various Persian and English resources. A corpus of Persian reviews developed for opinion mining studies are introduced. To develop HesNegar, a comprehensive Persian WordNet (FerdowsNet), with high recall and proper precision (based on Princeton WordNet), was first created. Then, the polarity of each synset in English SentiWordNet is mapped to the corresponding words in HesNegar. In the conducted tests, it was found that HesNegar has a precision score of 0.86 a recall score of 0.75 and it can be used as a comprehensive Persian

\* Corresponding author

\* نویسنده عهده دار مکاتبات

فصلنامه



سال ۱۳۹۷ شماره ۱ پیاپی ۳۵

۷۱

www.SID.ir

SentiWordNet. The findings and developments made in this study could prove useful in the advancement of opinion mining research in Persian and other similar languages, such as Urdu and Arabic.

**Keywords:** Opinion Mining, FerdowsNet (Persian WordNet), Sentiment Lexicon, Persian Text Processing Tools

جمله<sup>۵</sup> [3] یا قاب‌های معنایی<sup>۶</sup> [4]، ادامه دادند. معمولاً در این واحد، از تحلیل ذهنیت<sup>۷</sup> برای تمیز دادن جملات حاوی حس از جملات فاقد حس (شامل واقعیت مانند اخبار) استفاده می‌شود.

در سال‌های اخیر بسیاری از پژوهش‌های صورت گرفته در این زمینه به سمت زبان‌های غیرانگلیسی (به‌ویژه اسپانیایی، چینی، آلمانی، چکوسلوواکی<sup>۸</sup> و عربی)، معطوف شده‌است [5-8]. البته روش‌های نسبتاً جدیدی برای نظر کاوی چندزبانه<sup>۹</sup> در حال ایجاد هستند. در زمینه نظر کاوی چندزبانه بیشتر کارهای انجام‌شده به تحلیل حس در واحد متن (سند) و یا تحلیل ذهنیت و حس در واحد جمله پرداختند [9,10]. غالباً این روش‌ها با ایده‌های به‌نسبت ساده‌ای از مترجم‌های ماشینی به منظور بهره‌گیری از مجموعه اصطلاحات حاوی حس و سایر منابع و ابزارهای زبان انگلیسی برای زبان مورد نظر خود، استفاده می‌کنند [11-14]. ولی با توجه به تفاوت قواعد نحوی زبان‌ها، اصطلاحات حاوی حس و سایر پیچیدگی‌های زبانی، نتایج به‌دست‌آمده از این دسته روش‌ها دقت مناسبی برای استفاده در زبان‌های مختلف ندارند. در حال حاضر پژوهش‌های بسیار اندکی بر روی نظر کاوی در زبان فارسی صورت گرفته‌است. البته کمبود منابع (از قبیل پیکره استاندارد متون نظرات و مجموعه واژگان حاوی حس استاندارد) و ابزارهای استاندارد پردازش متن مناسب و دسترس و سایر پیچیدگی‌های نحوی و دستوری زبان فارسی مانع بزرگی برای انجام بررسی‌های نظر کاوی در زبان فارسی است.

ابزارهای پردازش زبان طبیعی و شبکه‌های واژگان<sup>۱۰</sup> (هستان‌شناسی ریخت‌شناسی بین کلمات<sup>۱۱</sup>) نقش مهمی در کاربردهای پردازش زبان طبیعی و بازیابی اطلاعات متنی

<sup>5</sup> Sentence level sentiment analysis

<sup>6</sup> Semantic frame

<sup>7</sup> Subjectivity Analysis

<sup>8</sup> Czech

<sup>9</sup> Multi-Lingual

<sup>10</sup> WordNet

<sup>11</sup> Lexical ontology of words

## ۱- مقدمه

آگاهی از تجربیات، نظرات و دیدگاه افراد نقش اساسی در فرایند تصمیم‌گیری مشتریان ساده تا مدیران سطح بالای شرکت‌های تولیدکننده و سازمان‌های مختلف دارد. امروزه با پیدایش تارنمای ۲/۰، حجم بسیاری از متون مربوط به نظرات افراد ایجاد شده‌است. ولی کاوش در حجم انبوه مستندات، منابع نظرسنجی مختلف و وجود نظرات مغایر درباره یک موجودیت، فرایند استخراج و پردازش نظرات را بسیار دشوار ساخته‌است. بنابراین نیاز به روش‌هایی برای کاوش و تحلیل نظرات موجود در وب احساس می‌شود. بدین منظور در دهه اخیر گرایش جدیدی در علم پردازش زبان‌های طبیعی<sup>۱</sup> به نام نظر کاوی<sup>۲</sup> ایجاد شده‌است. یکی از مهمترین وظایف این حوزه دسته‌بندی حسی مستندات بر اساس بار حسی<sup>۳</sup> مثبت یا منفی (میزان رضایت‌مندی) آنهاست.

اغلب پژوهش‌های نخستین در زمینه نظر کاوی، سعی در دسته‌بندی نظر یا حس کلی یک متن، در قالب دو دسته حس مثبت و منفی، داشتند [1]. در ادامه، پژوهشگران سعی در تعیین درجه رضایت (میزان کمی در بازه مشخص) متن نظرات (به جای دسته‌بندی قطبی یا دوحالتی) کردند [2]. برای انجام این دسته‌بندی‌ها اغلب از روش‌های نظارت‌شده<sup>۴</sup> که برچسب نمونه‌ها در آنها به‌صورت دستی مشخص شده‌بود، در حوزه‌های محصولات تجاری که نظرات به‌تصریح بیان شده‌بودند، استفاده می‌کردند. مشکل اصلی تحلیل حس در واحد متن، فرض یکسان بودن مضمون در همه متن یا متون جمع‌آوری شده است؛ درحالی‌که ممکن است بخش‌های مختلف متن (یا متون مختلف)، دارای مضمون‌های متفاوتی باشند. بنابراین قبل از تحلیل حس، لازم است مضمون بخش‌های مختلف شناسایی و جدا از هم بررسی شوند. درنتیجه، پژوهش‌گران نظر کاوی کار تحلیل حس را در واحد

<sup>1</sup> Natural Language Processing

<sup>2</sup> Opinion Mining

<sup>3</sup> polarity

<sup>4</sup> supervised

می‌شود. برای محاسبه شباهت معنایی به‌طور معمول از روش‌های زیر استفاده می‌شود:

۱- مبتنی بر شبکه واژگان یا سایر لغت‌نامه‌ها و دانش‌نامه‌ها [17]

۲- روابط وابستگی نحوی بین عبارت حاوی حس با کلمات موجود در واژه‌نامه حس

۳- هم‌رخدادی عبارات حاوی حس با کلمات موجود در فهرست نخستین کلمات حاوی حس (روش‌های یادگیر بدون ناظر) در درون پیکره‌های مختلف مستندات<sup>۲</sup> [18].

این رویکرد را زیرمجموعه‌ای از رویکرد تشخیصی حس مبتنی بر مجموعه لغات می‌توان به‌شمار آورد؛ با این تفاوت که فهرست کلمات حاوی حس برای مستندات ورودی (داده شده) تشکیل می‌شود و از این‌رو کارکرد بهتری برای تشخیص حس عبارات در حوزه‌های مختلف خواهد داشت.

با توجه به وابستگی زیاد روش‌های مختلف تحلیل احساسات به لغت‌نامه واژگان حس، در این بخش به توضیح روش‌های مختلف واژه‌نامه حس پرداخته می‌شود؛ به‌طور کلی، از سه رویکرد زیر برای تولید واژه‌نامه‌های حس استفاده می‌شود:

۱- مبتنی بر پیکره [1,19-21]

۲- مبتنی بر لغت‌نامه و پایگاه دانش [17,22-24]

۳- مبتنی بر روش‌های یادگیر باناظر [25-27]

در رویکرد ساخت لغت‌نامه حس با استفاده از روش‌های یادگیر باناظر، به داده‌های آموزشی دارای برچسب حس اولیه نیاز است. به‌دلیل مشکلات برچسب‌گذاری حس کلمات (تهیه پیکره آموزشی) گاهی از این رویکرد برای استخراج واژگان حس در دامنه محدود (خاص) نظرات استفاده می‌شود. اغلب از ترکیب این رویکرد با روش‌های مبتنی بر لغت‌نامه (پایگاه دانش) و یا پیکره نظرات (با قالب مشخص) استفاده می‌شود.

روش‌های مبتنی بر پیکره، از پیکره‌های متنی به‌نسبت بزرگ و از قواعد زبان‌شناسی استفاده می‌کنند. به‌طور معمول برای ایجاد واژه‌نامه حس وابسته به دامنه (موضوع) خاص استفاده می‌شود. البته با در نظر گرفتن پیکره‌های متنی بزرگ از این رویکرد برای تولید واژه‌نامه‌های حس عمومی (مستقل از دامنه) نیز می‌توان استفاده کرد. هاتز یواسیلوقلو و همکارش [19] از قواعد زبان‌شناسی برای تعیین بار حس صفت‌ها

<sup>۲</sup> مستندات نظرات کاربران، صفحات وب یا نتایج خروجی موتورهای جستجو

موجود در وب دارد. از این‌رو در این پژوهش، ابتدا مجموعه‌ای از ابزارهای پردازش متون زبان فارسی مورد نیاز ایجاد شده و شبکه واژگان جامع زبان فارسی (فردوس‌نت) معرفی می‌شوند. سپس با استفاده از شبکه واژگان ایجادشده، مفاهیم و عبارات حس استخراج‌شده برای زبان انگلیسی به زبان فارسی نگاشت داده می‌شوند. تهیه واژه‌نامه عبارات حاوی حس<sup>۱</sup> برای زبان طبیعی یکی از حوزه‌های پژوهشی در نظر کاوی به‌شمار می‌رود. علاوه‌براین، جهت ارزیابی روش‌های دسته‌بندی احساسات، پیکره متون نظرات مربوط به محصولات مختلف تجاری در زبان فارسی جمع‌آوری و بخشی از آن توسط انسان حاشیه‌نویسی (برچسب‌گذاری) حس شده است.

در بخش بعدی مقاله به مرور روش‌های دسته‌بندی احساسات پرداخته می‌شود. در بخش سوم ابزارهای ایجادشده برای پردازش متون زبان فارسی معرفی و به‌اختصار توضیح داده خواهند شد. در بخش چهارم، نحوه ایجاد شبکه واژگان فردوس‌نت و تولید شبکه واژگان حس فارسی (حس‌نگار) توضیح داده خواهد شد. در بخش آخر نیز نتایج ارزیابی کیفیت ابزارهای پردازش متن مختلف ایجادشده، شبکه‌های واژگان مختلف زبان فارسی و حس‌نگار گزارش خواهند شد.

## ۲- مروری بر کارهای گذشته

در چند سال گذشته، ساخت مجموعه لغات حاوی حس با بار حس مثبت و منفی، یکی از روش‌های مورد توجه پژوهش‌گران برای تشخیص حس جملات بوده است. به‌طور کلی روش‌های تحلیل احساسات را می‌توان به دو گروه ۱- روش‌های مبتنی بر واژه‌نامه حس و استفاده از دانش زمینه (یادگیری بدون ناظر یا شبه‌ناظر) و ۲- روش‌های یادگیری باناظر تقسیم‌بندی کرد. دقت روش‌های مبتنی بر واژه‌نامه حس به‌طور کامل وابسته به مجموعه لغات حاوی حس و وزن‌های از پیش تعیین شده است [15]. این روش‌ها به‌طور بدون ناظر و برای حوزه‌های عمومی قابل استفاده هستند [16]. در رویکرد دوم (دسته‌بندی حس متون) نیز از واژگان حس به‌عنوان یکی از ویژگی‌های مهم متن نظرات استفاده می‌شود.

از دیگر روش‌های تشخیص حس عبارات، استفاده از روش‌های محاسبه شباهت معنایی کلمات است. در این روش‌ها برای تشخیص حس نظرات به‌طور معمول از شباهت معنایی عبارات و فهرست کوچکی از کلمات حاوی حس نخستین استفاده

<sup>1</sup> Sentiment lexicon

استفاده کردند. بدین منظور، از حروف ربط<sup>۱</sup> بین صفت‌ها در پیکره متنی بزرگ برای دسته‌بندی حسی کلمات استفاده شد. ترنی [1] تعدادی از الگوهای زبانی برای استخراج عبارات حسی را مشخص کرد. همچنین از موتور جستجو و معیار پی‌ام‌آی<sup>۲</sup> برای محاسبه هم‌رخدادی و ارتباط معنایی عبارات حسی مختلف با کلمات حسی کلیدی مثبت و منفی از پیش تعیین شده (مانند کلمات “excellent” و “poor”)، استفاده شد. ولیکوویچ و همکاران [21] تمام چندکلمه‌ای‌های (n-grams) موجود را از میلیون‌ها صفحه وب به‌عنوان گره‌های یک گراف در نظر گرفتند. وزن یال بین این گره‌ها بر اساس فاصله کسینوسی میان بردارهای متناظر با آن عبارات محاسبه شد؛ سپس آنها با انتشار میزان بار حسی برخی عبارات حسی نخستین<sup>۳</sup> از طریق یال‌های گراف، بار حسی سایر عبارات را محاسبه کردند.

برخی از روش‌های استخراج واژگان حسی نیز وابسته به پیکره متنی نظرات با قالب مشخص هستند. کاجی و همکارش [20] ابتدا تکرار رخداد عبارات حسی در جملات مثبت (بخش نقاط قوت سایت) و جملات منفی (بخش نقاط ضعف سایت) را شمارش کردند. پس از آن، از معیار پی‌ام‌آی برای محاسبه بار حسی عبارات گزیده استفاده کردند. با ایده‌ای مشابه، نوفرستی و شمس‌فرد از یک پیکره از واژگان حسی مربوط به داروها با استفاده از کلمات کلیدی (موجودیت‌های اثرپذیر) در متن نظرات کمک گرفتند [28]. آنها از متن نظرات موجود در بخش فواید یا عوارض داروها (در منبع ورودی) برای تعیین بار حسی مثبت یا منفی استفاده کردند.

روش‌های مبتنی بر لغت‌نامه اغلب از شبکه‌واژگان برای تعیین روابط معنایی و محاسبه بار حسی کلمات استفاده می‌کنند. کمپس و همکاران [17] بر اساس روابط هم‌معنی شبکه‌واژگان یک گراف تشکیل دادند. آنها فرض کردند که کلمات هم‌معنی دارای بار حسی یکسانی هستند؛ سپس بر اساس فاصله کوتاه‌ترین مسیر بین کلمات (گره‌های گراف) میزان بار حسی سایر کلمات را محاسبه کردند. کانایاما و همکارش [29] نیز از روابط هم‌معنایی و تضاد در شبکه‌واژگان استفاده کردند. آنها بار حسی مجموعه کلمات دارای حس را با استفاده از یک روش بوت استرپ<sup>۴</sup> به سایر لغات هم‌جوار

انتشار دادند. تاکامورا و همکاران [22] با استفاده از روابط معنایی مختلف (از قبیل هم‌معنی، تضاد و روابط سلسله مراتبی) و هم‌رخدادی عبارات بخش توصیف<sup>۵</sup> در شبکه‌واژگان، شبکه‌ای معنایی از واژگان ایجاد و سپس با استفاده از یک مدل چرخشی<sup>۶</sup> بار حسی را بین کلمات مختلف بخش کردند. رائو و همکارش [24] نیز از شبکه‌واژگان و فرهنگ جامع کلمات<sup>۷</sup> در OpenOffice به‌عنوان منبع برای ساخت گراف کلمات استفاده کردند؛ سپس با بهره‌گیری از یک روش شبه‌ناظر برچسب‌های حسی را در گراف انتشار دادند. در بسیاری از مقالات [23,30-32] از روش‌های مبتنی بر گام تصادفی<sup>۸</sup> و رتبه‌دهی صفحات<sup>۹</sup> در گراف برای تعیین بار حسی (گرایش معنایی<sup>۱۰</sup>) کلمات استفاده شده‌است.

برای ایجاد شبکه‌واژگان حسی انگلیسی به نام سنتی‌وردنت<sup>۱۱</sup>، با یک الگوریتم یادگیر شبه‌ناظر در چهار مرحله به هر گروه از واژگان هم‌معنی در شبکه‌واژگان بار حسی مثبت، منفی و میزان غیرذهنی‌بودن<sup>۱۲</sup>، انتساب دادند [33, 34]. با توجه به استفاده از این منبع برای تولید واژه‌نامه حسی زبان فارسی، نحوه ایجاد شبکه‌واژگان حسی انگلیسی در بخش‌های بعد توضیح داده می‌شود. نیواروسکایا و همکارش [35] برای ایجاد لغات نام حسی به نام سنتی‌فول<sup>۱۳</sup> ابتدا هسته نخستین از کلمات حسی از پایگاه داده افکت<sup>۱۴</sup> [36] را استخراج و برای ابهام‌زدایی و آزمون درستی کلمات حسی نخستین از شبکه‌واژگان حسی انگلیسی (سنتی‌وردنت) استفاده کردند؛ سپس با روابط معنایی مختلف موجود در شبکه‌واژگان، مجموعه لغات سنتی‌فول گسترش داده شد. پوریا و همکاران [25] با استفاده از یک روش دسته‌بندی شش-کلاسه برچسب‌های پایگاه داده وردنت افکت<sup>۱۵</sup> [37] را به هر یک از مفاهیم پایگاه داده سنتیک‌نت<sup>۱۶</sup> [38] انتساب دادند.

5 WordNet Glosses

6 Spin Model

7 thesaurus

8 Random Walk

9 PageRank

10 Semantic orientation

11 SentiWordNet

12 objectivity

13 SentiFul

14 Affect

15 WordNet Affect

16 SenticNet

1 conjunctive

2 Pointwise Mutual Information

3 Sentiment seed words

4 boot-strapping method

خاص خود هستند. به‌عنوان مثال کلمات در مقوله اسم با هم دارای روابط معنایی‌ای مانند: هم‌معنایی<sup>۶</sup>، تضاد معنایی<sup>۷</sup>، رابطه شمول معنایی (دربرداشتن)<sup>۸</sup>، روابط سلسله‌مراتبی<sup>۹</sup> (جزء به کل<sup>۱۰</sup> و کل به جزء<sup>۱۱</sup>) و ... هستند. شبکه واژگان اغلب برای ابهام‌زدایی و تعیین شباهت معنایی در کاربردهای مختلف پردازش زبان طبیعی و بازیابی اطلاعات مانند ترجمه ماشینی، استخراج اطلاعات و خلاصه‌سازی و ... استفاده می‌شود. آخرین نسخه پی.دبلیو.ان<sup>۱۲</sup> به‌طور تقریبی شامل ۱۵۵۳۲۷ کلمه است که در قالب ۱۱۷۵۹۷ گروه هم‌معنی سازماندهی شدند. اخیراً در بعضی مقالات از شبکه واژگان برای استخراج واژگان حسی و ویژگی‌های موجودیت موردنظر استفاده شده است [23,44].

اکنون برای بیش از چهار زبان طبیعی در جهان شبکه واژگان ایجاد شده‌است که بین غالب آنها با پی.دبلیو.ان پیوند وجود دارد. در ایجاد شبکه واژگان برای سایر زبان‌ها به‌طور معمول از دو روش استفاده می‌شود:

- ۱- ساخت شبکه واژگان با استفاده از ترجمه گروه‌های هم‌معنی پی.دبلیو.ان مانند شبکه واژگان فارسی دانشگاه تهران [45] و آی‌دبلیو.ان دی<sup>۱۳</sup> [46] و
- ۲- استفاده شبکه واژگان از منابع زبان مقصد و روش‌های زبان‌شناسی ایجاد شده و برقراری پیوند بین گروه‌های هم‌معنی آن با پی.دبلیو.ان مانند فارسن ت [47] و یوروورلدنت<sup>۱۴</sup> [48].

### ۱-۳- شبکه‌های واژگان فارسی موجود

در زبان فارسی فعالیت‌های زیادی در جهت ساخت شبکه واژگان به‌صورت خودکار یا نیمه‌خودکار انجام شده‌است [45,47, 49-52] ولی تنها شبکه واژگان فارسن ت [47] و شبکه

در اغلب کارهای انجام‌شده در زمینه نظر کاوی در زبان فارسی، مجموعه‌ای از لغات حسی به‌صورت دستی جمع‌آوری و استفاده شده‌است [39,40]. همچنین در برخی از کارها [41,42] از ترجمه فارسی (توسط مترجم ساده انگلیسی به فارسی) منابع (واژه‌نامه‌های حسی) زبان انگلیسی استفاده شده‌است. ددهار بهبهانی و همکاران [32] در روشی جدید ساخت واژگان حسی فارسی را با استفاده از منابع زبان انگلیسی (لغات حسی نخستین و شبکه واژگان) پیشنهاد دادند. به این منظور آنها نخست فهرست اولیه لغات حسی انگلیسی (پیکره Micro-WNOp [43]) را به‌صورت دستی مشخص کرده، سپس با استفاده از روش گام تصادفی، بار (وزن) حسی سایر لغات موجود در گراف معنایی (شبکه واژگان) را تعیین کردند. با توجه به پراکندگی<sup>۱</sup> و ناقص بودن شبکه واژگان فارسی موجود آنها نمی‌توانستند روش خود را به‌طور مستقیم بر روی شبکه واژگان فارسی اجرا کنند. آنها سپس با استفاده از ارتباط (لینک) موجود بین گروه‌های هم‌معنی شبکه واژگان انگلیسی با شبکه واژگان فارسن ت (نسخه ۱) مجموعه‌ای از واژگان حسی فارسی (UTHS) ایجاد کردند. درنهایت، در مرحله پس‌پردازش با توجه به روابط هم‌معنایی و تضاد موجود در شبکه واژگان، بار حسی کلمات اصلاح شدند. با توجه به تعداد کم گروه‌های هم‌معنی در فارسن ت (نسبت به شبکه واژگان انگلیسی) سایر کلمات گروه‌های هم‌معنی شبکه واژگان انگلیسی با مترجم گوگل، ترجمه و وارد مجموعه لغات حسی فارسی شدند. در تمام مجموعه واژگان حسی فارسی ایجادشده شامل ۱۸۱۵ کلمه حسی مثبت و ۱۸۵۶ کلمه حسی منفی شامل نقش‌های صفت، اسم و فعل است.

### ۳- شبکه واژگان

شبکه واژگان دانشگاه پرینستون<sup>۲</sup> (PWN) یک پایگاه داده لغوی<sup>۳</sup> برای زبان انگلیسی است. شبکه واژگان حاوی لغات زبان طبیعی در قالب مجموعه‌های کلمات هم‌معنی<sup>۴</sup> (گروه هم‌معنی<sup>۵</sup>) است که در دسته‌هایی با توجه به نقش‌های نحوی‌ای مانند فعل، اسم، صفت و قید تقسیم‌بندی شده‌اند. گروه‌های هم‌معنی در هر مقوله دستوری دارای روابط معنایی

<sup>1</sup> Sparse

<sup>2</sup> Princeton WordNet

<sup>3</sup> Electronic Lexical Database

<sup>4</sup> synonymous sets

<sup>5</sup> synset

<sup>6</sup> synonymy

<sup>7</sup> antonymy

<sup>8</sup> meronymy

<sup>9</sup> Taxonomic

<sup>10</sup> hyponymy

<sup>11</sup> hypernymy

<sup>12</sup> WordNet 3.1 database statistics

<sup>13</sup> IWND

<sup>14</sup> EuroWordNet

واژگان فارسی دانشگاه تهران [45] منتشر شده و برای بهره‌برداری و مقایسه قابل استفاده هستند.

### ۱-۳-۱- فارسی‌نت

فارسی‌نت (شبکه واژگان عمومی زبان فارسی) یک پایگاه دانشی لغوی<sup>۱</sup> است که حاوی اطلاعات در مورد واژه‌ها و ترکیبات زبان (مفاهیم)، اجزای کلام آنها (POS) و روابط معنایی میان آنهاست. آخرین نسخه این پایگاه داده (فارسی‌نت نسخه ۲,۰) برای استفاده‌های پژوهشی در دسترس است.<sup>۲</sup> برای تولید این شبکه واژگان ایده به‌کاررفته در یوروورلدنت استفاده شد. به این ترتیب، نخست به‌صورت دستی هسته نخست شبکه واژگان فارسی ایجاد و سپس با استفاده از یک روش شبه‌ناظر در فرایندی بالا به پایین این شبکه واژگان تکمیل شد. هسته نخست فارسی‌نت به‌کمک ترجمه مفاهیم بالکاننت<sup>۳</sup> و برخی مفاهیم پرکاربرد زبان فارسی ایجاد شد. سپس با به‌کارگیری یک روش شبه‌ناظر و استفاده از منابع مختلف زبان فارسی و منابع دوزبانه (فارسی-انگلیسی) هسته نخست شبکه واژگان تکمیل شد [47].

### ۲-۳-۱- شبکه واژگان فارسی دانشگاه تهران

آخرین نسخه از شبکه واژگان فارسی تهیه‌شده در دانشگاه تهران [45] در سایت شبکه‌های واژگان چندزبانه<sup>۴</sup> قابل دریافت است. این شبکه واژگان با اجرای الگوریتم بدون ناظر EM<sup>۵</sup> و استفاده از یک پیکره متنی و شبکه واژگان انگلیسی (PWN) ایجاد شده‌است. آنان برای محاسبه احتمال نخست وجود هر کلمه را در هر گروه کلمات هم‌معنی از شبکه واژگان فارسی‌نت نسخه ۱,۰ استفاده و سپس سعی کردند که این احتمال را (برای هر کلمه) به کمک روش تکراری EM<sup>۶</sup> بیشینه کنند.

### ۲-۳-۲- فردوس‌نت

شبکه واژگان ایجادشده در این پژوهش فردوس‌نت نامیده شده است. هدف اصلی از ایجاد فردوس‌نت تولید یک شبکه واژگان با بازخوانی<sup>۷</sup> بالا (پوشش تقریباً کاملی از شبکه واژگان

پرینستون) با دقتی مناسب برای زبان فارسی است. برای ساخت شبکه واژگان فردوس‌نت، از منابع زبانی و پایگاه‌های دانش زیر استفاده شده‌است:

- شبکه واژگان دانشگاه پرینستون (PWN)
- لغت‌نامه‌های دوزبانه (انگلیسی-فارسی) مختلف
- مترجم متن گوگل
- پایگاه دانش ویکی‌پدیا و هستان‌شناسی یاگو<sup>۷</sup> [53]
- برای برقرار کردن پیوند بین ویکی‌پدیا و پی.دبلیو.ان
- پیکره‌های متنی فارسی (چند سایت خبری و صفحات ویکی‌پدیا فارسی)
- دانش‌نامه‌ها، دایره‌المعارف و لغت‌نامه‌های فارسی
- شبکه‌های واژگان فارسی فارسی‌نت و شبکه واژگان دانشگاه تهران

شبکه واژگان فردوس‌نت در طی فرایند نهم‌مرحله‌ای (مطابق شکل (۱)) برای تک‌تک گروه‌های هم‌معنی، تشکیل شده‌است:

**مرحله نخست:** تمام کلمات گروه‌های هم‌معنی با لغت‌نامه‌های<sup>۸</sup> دوزبانه (انگلیسی به فارسی) مختلف ترجمه می‌شوند.

**مرحله دوم:** به‌ازای هر گروه هم‌معنی (از شبکه واژگان انگلیسی) یک گراف دوبخشی تشکیل شده و مراحل دو الی نه انجام می‌شوند. در این گراف دوبخشی، به‌ازای هر کلمه انگلیسی یک گره در سمت چپ گراف ( $X_i$ ) و به‌ازای هر کلمه فارسی (لیست ترجمه‌ها) یک گره در سمت راست گراف ( $Y_i$ ) در نظر گرفته می‌شود؛ سپس به هر کلمه انگلیسی  $x_i$  و ترجمه فارسی آن  $y_j$  یالی به شکل  $(x_i, y_j)$  بین گره مربوط رسم می‌شود. وزن هر یال وابسته به تعداد تکرار کلمات در فهرست ترجمه (با لغت‌نامه‌های مختلف) و رتبه ترجمه<sup>۹</sup> آنهاست.

**مرحله سوم:** این مرحله شامل دو قسمت است:

- ۱- نخست از پایگاه‌های دانش ویکی‌پدیا و سایر شبکه‌های واژگان فارسی، کلمات فارسی مرتبط با گروه هم‌معنی جاری استخراج می‌شود. به این منظور، نخست با استفاده از هستان‌شناسی یاگو [53] مفاهیم مربوط به کلمات گروه هم‌معنی انتخاب‌شده در ویکی‌پدیا استخراج و

<sup>1</sup> lexical database

<sup>2</sup> <http://dadegan.ir/catalog/farsnet>

<sup>3</sup> BalkaNet

<sup>4</sup> <http://compling.hss.ntu.edu.sg/omw/>

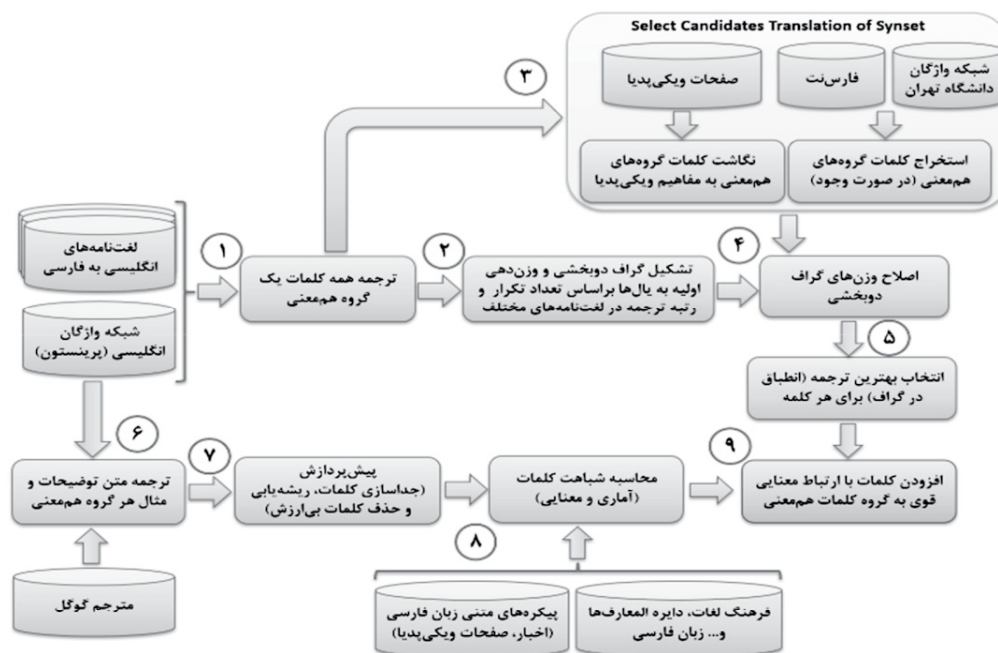
<sup>5</sup> Expectation Maximization

<sup>6</sup> Recall

<sup>7</sup> Yago

<sup>8</sup> Dictionary

<sup>9</sup> چون در اغلب لغت‌نامه‌های موجود کلمات ترجمه برحسب اولویت مرتب هستند.



(شکل-۱): معماری ایجاد گروه‌های هم معنی شبکه واژگان فردوس‌نت  
(Figure-1): The system architecture for construction of the synsets of FerdowsNet

ارتباطی آنها (با یکی از کلمات انگلیسی) بیشتر از متوسط وزن ارتباطی کلمات انتخاب شده (یال‌های انتخاب شده) باشد، به‌عنوان ترجمه گزیده دوم (هم مفهوم با کلمه انگلیسی)، انتخاب می‌شود.

**مرحله ششم:** توضیحات<sup>۴</sup> و مثال‌های مربوط به هر گروه هم معنی (در پی.دابلیوان) با مترجم گوگل ترجمه می‌شود.

**مرحله هفتم:** پیش برداش لازم بر روی متن ترجمه انجام شده و کلمات کلیدی متن (معادل با کلمات فارسی در سمت راست گراف دوبخشی) استخراج می‌شوند.

**مرحله هشتم:** با استفاده از پیگروه‌های معنی زبان فارسی (خبرگزاری‌های برخط مختلف [54-56] و محتویات صفحات ویکی پدیا فارسی) و همچنین فرهنگ واژگان مختلف فارسی، کلمات هم معنی و هم کاربرد (با استفاده از معیار پی‌ام‌آی [1,57]) با کلمات انتخاب شده در مرحله پنجم (مجموعه کلمات S) استخراج می‌شوند.

**مرحله نهم:** از میان کلمات استخراج شده در مرحله پیش، کلماتی که شباهت آنها با کلمات مجموعه S (که در مرحله پنجم انتخاب شدند) بیش از حد آستانه مشخصی باشد، به فهرست کلمات نهایی افزوده می‌شوند.

پس از ترجمه و گسترش کلمات گروه‌های هم معنی

معادل فارسی این مفاهیم (در صورت وجود صفحه فارسی ویکی پدیا برای آن مفهوم) استخراج می‌شوند.

۲- کلمات گروه‌های هم معنی معادل (با استفاده از پیوند بین شبکه‌های واژگان دانشگاه تهران و فارس‌نت با پی.دابلیوان در صورت وجود) استخراج می‌شوند.

**مرحله چهارم:** در این مرحله با استفاده از کلمات استخراج شده در مرحله قبل (از پایگاه‌های دانش ویکی پدیا و سایر شبکه‌های واژگان فارسی)، کلمات ترجمه گسترش (افزودن گره‌های جدید به سمت راست گراف دوبخشی)، یا وزن یال‌های مربوط به کلمات ترجمه موجود در گراف دوبخشی (تشکیل شده در مرحله دوم) تقویت (افزوده) می‌شوند.

**مرحله پنجم:** با استفاده از الگوریتم هانگاریس<sup>۱</sup> (مجاری)<sup>۲</sup> بهترین تطابق<sup>۳</sup> (مناسب‌ترین ترجمه‌های فارسی برای کلمات انگلیسی) از گراف دو بخشی وزن دار تشکیل شده استخراج می‌شوند. به این ترتیب فهرست کلمات نامزد S برای گروه هم معنی فعلی به دست می‌آید. همچنین در این مرحله، از بین کلمات فارسی انتخاب نشده، کلماتی که وزن

<sup>1</sup> Hungarian

<https://github.com/KevinStern/software-and-algorithm>

<sup>2</sup> [s/blob/master/src/main/java/blogspot/software\\_and\\_algorithm/s/stern\\_library/optimization/HungarianAlgorithm.java](https://github.com/KevinStern/software-and-algorithm/blob/master/src/main/java/blogspot/software_and_algorithm/s/stern_library/optimization/HungarianAlgorithm.java)

<sup>3</sup> Matching

<sup>4</sup> Glossary

کلمات در شبکه واژگان حسی انگلیسی (SentiWordNet) استفاده شده است.

روش پیشنهادی این مقاله بر پایه استفاده از بار حسی کلمات موجود در واژه‌نامه شبکه واژگان حسی انگلیسی نسخه ۳،۰ با نداشت و ارتباط بین فردوس‌نت و شبکه واژگان انگلیسی است.

#### ۱-۴- ابزارهای پردازش متن فارسی

زبان فارسی یکی از زبان‌های هندو-اروپایی<sup>۲</sup> است که بیش از صد میلیون نفر از مردم جهان با آن صحبت می‌کنند و به‌عنوان زبان رسمی سه کشور ایران، افغانستان (فارسی دری) و تاجیکستان (فارسی تاجیکی) است. به‌دلیل پیچیدگی‌های زبانی، منابع و مطالعات انجام‌شده در این زبان از دیدگاه محاسباتی پژوهش‌گران کمتر به آن توجه کرده‌اند [64,65]. استانداردهای، تفکیک متن به جملات، عبارات و کلمات و برچسب‌گذاری و حاشیه‌نویسی آنها تأثیر به‌سزایی بر روی پردازش و استخراج اطلاعات، دسته‌بندی یا دیگر کاربردهای پردازش زبان طبیعی دارد. با توجه به اهمیت مرحله پیش‌پردازش در تحلیل نظرات [66]؛ مجموعه‌ای از ابزارهای کاربردی پردازش متون زبان فارسی برای استفاده در شبکه واژگان حسی فارسی تهیه شده است.

#### ۲-۴- شبکه واژگان حسی انگلیسی

شبکه واژگان حسی انگلیسی [33] یکی از بهترین منابع موجود برای شناسایی کلمات حسی است که بر اساس تعیین میزان بار حسی هر کدام از گروه هم‌معنی در شبکه واژگان انگلیسی پرینستون (PWN) ایجاد شده است. شبکه واژگان حسی انگلیسی برای هر کدام از گروه‌های هم‌معنی میزان بار حسی منفی<sup>۳</sup>، مثبت<sup>۴</sup> و همچنین مقدار غیرحسی<sup>۵</sup> بودن (با توجه به مقدار حس مثبت و منفی) را با عددی بین صفر و یک مشخص می‌کند.

شبکه واژگان حسی انگلیسی نسخه ۱/۰<sup>۶</sup> [34] با یک الگوریتم یادگیر شبه‌ناظر در چهار مرحله تشکیل شده است:

۱- بار حسی مثبت و منفی تعدادی محدودی از گروه‌های

<sup>2</sup> Indo-European

<sup>3</sup> negativity

<sup>4</sup> positivity

<sup>5</sup> objectivity

<sup>6</sup> 1-(positivity score + negativity score)

<sup>7</sup> SentiWordNet v1.0

پی.دابلیوان، از روابط موجود بین گروه‌های هم‌معنی در این شبکه واژگان برای ارتباط بین گروه‌های هم‌معنی فردوس‌نت استفاده می‌شود.

در شبکه واژگان فردوس‌نت برای هر کلمه در گروه کلمات هم‌معنی مقداری به‌عنوان میزان اطمینان<sup>۱</sup> با نداشت وزن یال‌های مرتبط با آن کلمه در گراف دوبخشی کلمات (که در الگوریتم ساخت فردوس‌نت توضیح داده شد)، به عددی بین صفر و یک، محاسبه می‌شود.

#### ۴- تولید شبکه واژگان حسی

بدیهی است که تشخیص عبارات حسی و تعیین میزان حس آنها (کمی‌سازی میزان حس) بدون کمک انسان، برای ماشین غیرممکن است. بنابراین در روش‌های تحلیل حس، ابتدا فهرستی از عبارات حسی نخستین که به‌طور معمول دارای مقدار عددی تعیین‌کننده در میزان بار حسی هستند، با کمک اشخاص خبره تهیه شده و به‌عنوان ورودی به تحلیل‌گر حس متن داده می‌شود؛ سپس با الگوریتم‌های مختلفی فهرست نخستین عبارات حاوی حس بسط و تکمیل می‌شود و میزان و شدت حس نیز با توجه به برخی کلمات جمله از قبیل منفی‌کننده‌ها (معکوس‌کننده‌های حسی) تنظیم می‌شود. افعال اسنادی منفی (از قبیل نیست، نبود و نمی‌باشد) و برخی کلمات (مانند «به‌ندرت»، «عدم وجود» و «بدون») از مهم‌ترین معکوس‌کننده‌های حسی جملات در زبان فارسی هستند. همان‌طور که اشاره شد، تولید واژه‌نامه لغات حسی یکی از بخش‌های اساسی و مهم برای تشخیص حس و شدت آن است. برای ایجاد مجموعه واژگان حسی به‌طور عموم دو روش در مقالات استفاده شده است:

۱- توسعه یا ترجمه عبارات حسی از روی مجموعه واژگان حسی موجود [58-60]

۲- گسترش فهرست کلمات حسی نخستین با استفاده از شبکه واژگان، پیکره‌ای از متون نظرات و سایر منابع زبانی [33,61-63].

اغلب در رویکرد دوم، فهرست کلمات حسی نخستین به‌صورت دستی تهیه می‌شوند. در این مقاله با ایجاد یک شبکه واژگان کامل برای زبان فارسی و برقراری پیوند بین مفاهیم آن با شبکه واژگان از رویکرد نخست استفاده شده است. به عبارت دیگر، برای تولید مجموعه واژگان حسی فارسی از بار حسی

<sup>1</sup> Confidence



شبکه واژگان پرینستون به زبان فارسی، شبکه واژگان جامع (فردوسنت) ساخته شده است. در نهایت، با استفاده از فردوسنت، میزان بار حسنی محاسبه شده برای هر گروه هم معنی در شبکه واژگان حسنی انگلیسی به گروه های هم معنی متناظر با آن در حس نگار نگاشت می شود. پس در واقع با ابهام زدایی مفاهیم شبکه واژگان حسنی انگلیسی، یک شبکه واژگان حسنی برای زبان فارسی ایجاد شده است.

به این ترتیب واژه های حسنی زبان فارسی در چهار مقوله صفت (مانند توانا، زشت)، اسم (مانند پیش کسوت، درد)، قید (مانند عجولانه، به آرامی) و فعل (مانند عشق ورزیدن، خندیدن) دسته بندی می شوند.

از شبکه واژگان حسنی فارسی به عنوان یک واژه نامه حسنی مرجع برای زبان فارسی می توان استفاده کرد. واژه نامه حسنی فارسی به دست آمده از کلمات حسنی با میزان بار حسنی بیشتر از ۰/۵ (بدون در نظر گرفتن میزان اطمینان<sup>۴</sup>)، شامل ۶۰۶۳ واژه حسنی مثبت و ۹۴۴۱ واژه حسنی منفی است.

علاوه بر این، با توجه به وجود درجه اطمینان برای کلمات موجود در هر گروه کلمات هم معنی در فردوسنت، برای هر کلمه حسنی علاوه بر بار حسنی مثبت و منفی، میزان اطمینان (اعتبار) نیز خواهیم داشت. تعداد کلمات حسنی با نقش ها و میزان اطمینان مختلف در جدول (۱) گردآوری شده است. تأثیر میزان اطمینان بر دقت و بازخوانی شبکه واژگان فردوسنت و در نتیجه تأثیر آن بر واژه نامه حسنی در بخش ارزیابی نشان داده شده است.

## ۵- ارزیابی

برای ارزیابی واژگان حسنی استخراج شده، مجموعه نظرات وبگاه دیجی کالا<sup>۵</sup> با خزش گرو<sup>۶</sup> در سال ۱۳۹۲ جمع آوری شده است. دیجی کالا پنجمین سایت پربازدید ایران و بزرگترین فروشگاه اینترنتی<sup>۷</sup> در ایران و خاورمیانه است<sup>۸</sup> و با توجه به حجم بالای کاربران، از غنای نظرات تقریباً خوبی

<sup>۴</sup> در صورتی که مجموع بار حسنی مثبت و منفی کلمه بیشتر از ۰.۵ باشد، در این صورت کلمه حسنی و در غیر این صورت وزن غیر حسنی کلمه بیشتر است.

<sup>۵</sup> www.digikala.com

<sup>۶</sup> Web crawler

<sup>۷</sup> market leader in e-commerce

<sup>۸</sup> http://www.alexa.com/topsites/countries/IR

هم معنی اولیه به صورت دستی مشخص شده و این بار حسنی با توجه به روابط معنایی (هم معنی و متضاد) موجود در شبکه واژگان بر روی کلمات و گروه های هم معنی مرتبط بسط داده می شود.

۲- سپس تعدادی از گروه های هم معنی فاقد بار حسنی<sup>۱</sup> نیز مشخص می شوند (به صورت دستی برچسب گذاری می شوند) و توضیحات<sup>۲</sup> آنها به همراه گروه های هم معنی مشخص شده در مرحله پیش به عنوان داده های آموزشی برای استفاده در مرحله یادگیری یک روش دسته بندی باناظر به کار گرفته می شوند.

۳- با الگوریتم دسته بندی (با نرخ بازخوانی پایین و دقت بالا) سایر گروه های هم معنی برچسب گذاری حسنی<sup>۳</sup> می شوند.

۴- به منظور کاستن از خطای الگوریتم های دسته بندی، گروه های هم معنی آماده شده در مرحله دو برای آموزش چند الگوریتم دسته بندی استفاده می شوند. بعد از اجرای آنها (مرحله سوم)، در این مرحله نتایج آنها با یکدیگر ترکیب می شوند. بدین منظور به ازای هر گروه هم معنی کلمات، میانگین بار حسنی به دست آمده از روش های دسته بندی مختلف محاسبه شده و به عددی در بازه صفر الی یک نگاشت می شود.

سپس در نسخه سوم این شبکه واژگان حسنی [33] با استفاده از الگوریتم تکراری گام تصادفی بر روی گراف شبکه واژگان پی.دابلیوان، نتایج حاصل از نسخه پیش (SentiWordNet v1.0) را اصلاح کردند.

تاکنون شبکه واژگان حسنی انگلیسی به عنوان یک منبع واژگان حسنی مستقل از دامنه و موضوع در بسیاری از کاربردهای نظر کاوی استفاده شده است [44, 67, 68]. علاوه بر این، با توجه به رابطه این واژه نامه با پی.دابلیوان، از این منبع در بسیاری از کاربردهای نظر کاوی در زبان های دیگر نیز (با برقرار کردن پیوند بین شبکه های واژگان آن زبان ها با پی.دابلیوان) استفاده شده است [58, 61].

## ۳-۴- ساخت واژه نامه حسنی زبان فارسی

همان طور که اشاره شد، در این پژوهش برای تولید حس نگار ابتدا با استفاده از نگاشت مفاهیم (گروه های هم معنی) در

1 objective

2 The glosses of the synsets

<sup>۳</sup> برچسب های "neg"، "pos" و "obj"

(جدول-۳): ویژگی‌های انواع نظرات در پیکره دیجی کالا  
(Table-3): A variety of opinions available on the corpus

نوع بیان (نوع نوشتار)		ویژگی‌های مورد بررسی		میزان ذهنی <sup>۱</sup> و سلیقه‌ای بودن		تعداد متوسط نظر به ازای هر کالا		طول (تعداد کاراکتر) متوسط متن	
رسمی	انواع ویژگی‌ها	نسبتاً کم <sup>۲</sup>	تعداد زیاد	۱	۳/۷۹	۸۰۴۵	۱۷۴۳	۳۵۶	۹/۱۲
نیمه رسمی	تعداد زیاد	متوسط	تعداد زیاد						
عامیانه	تعداد محدود	زیاد	تعداد محدود						

### ۱-۵- ارزیابی شبکه‌ی واژگان فردوسنت و حس‌نگار

ابتدا شبکه‌ی واژگان فردوسنت به صورت کمی ارزیابی و با سایر شبکه‌های واژگان موجود در زبان فارسی مقایسه می‌شود. در جدول (۴) ویژگی‌های شبکه‌های واژگان فارسی مختلف گردآوری شده‌است.

برای ارزیابی کیفی فردوسنت و مقایسه آن با دیگر شبکه‌های واژگان موجود در زبان فارسی، ابتدا حدود هزار گروه هم‌معنی از شبکه‌ی واژگان انگلیسی (از هر یک از مقوله‌های اسمی، فعلی، صفت و قید، ۲۵۰ گروه هم‌معنی) به صورت تصادفی انتخاب و سپس کلمات گروه‌های هم‌معنی معادل آنها در شبکه‌های واژگان مختلف زبان فارسی استخراج و برای داوری به اشخاص خبره داده شدند؛ سپس ارزیابی کیفی کلمات درون گروه‌های هم‌معنی مختلف توسط داوران خبره در زمینه‌ی پردازش زبان طبیعی انجام شده‌است. برای ارزیابی دقیق کیفیت، معیارهای دقت و بازخوانی برای هر شبکه واژگان به صورت جداگانه محاسبه شده‌است.

بدین منظور ابتدا به کمک داوران، به ازای هر گروه کلمات هم‌معنی انگلیسی یک مجموعه مرجع از لغات فارسی معادل آن (مفهوم) به صورت  $S^* = \{s_1^*, s_2^*, s_3^*, \dots, s_n^*\}$  در نظر گرفته می‌شود. همچنین، مجموعه واژه‌های موجود در شبکه واژگان برای هر گروه هم‌معنی،  $S^{wn} = \{s_1^{wn}, s_2^{wn}, s_3^{wn}, \dots\}$

<sup>۱</sup> subjective

<sup>۲</sup> اغلب بیان ویژگی‌های کمی (objective)

(جدول-۱): تعداد کلمات حسی مثبت (#Pos) و منفی (#Neg) موجود در شبکه‌ی واژگان حسی فارسی (حس‌نگار) براساس محدودیت‌های مختلف

(Table-1): Positive (Pos#) and negative (Neg#) sentiment words in PSWN

اسم	صفت	قید	فعل	همه	محدودیت
Pos#	۱۱۱۹۸	۹۴۹۱	۲۳۳۴	۳۶۷۵	۲۶۶۹۸
Neg#	۱۳۱۷۵	۱۱۰۲۲	۶۹۸	۴۱۷۳	۲۹۰۶۸
Pos#	۱۷۴۹	۱۳۷۲	۴۰۸	۳۸۳	۳۹۱۲
Neg#	۲۱۰۰	۱۶۲۴	۱۱۹	۴۷۱	۴۳۱۴
Pos#	۳۶۶	۴۸۳	۲۴	۳۹	۹۱۲
Neg#	۵۵۳	۷۵۰	۱۵	۸۷	۱۴۰۵

برخوردار است. این مجموعه داده شامل نظرات متنی بیان شده درباره محصولات مختلف است.

تعداد کل نظرات پیکره جمع‌آوری شده شامل ۳۱۷۳۰ نظر در ده نوع کالای متفاوت است. داوران حدود ۳۰۸۰ متن نظرات را برای آموزش الگوریتم‌های یادگیری ماشین بانظر و ارزیابی روش‌های نظر کاوی، برچسب زده‌اند؛ ولی بقیه نظرات فاقد برچسب حسی هستند. برخی از ویژگی‌های این مجموعه داده در جدول (۲) فهرست شده‌اند.

این مجموعه داده دارای سه دسته نظر مختلف مربوط به «نقد خبره»، «نظرات کاربران فعال» و «نظرات کاربران عادی و مهمان» است. در جدول (۳) ویژگی هر دسته از نظرات و تفاوت‌های آنها قابل مشاهده است.

(جدول-۲): ویژگی‌های پیکره نظرات دیجیکالا  
(Table-2): Features of Digikala review corpus

دارای برچسب حسی	کل پیکره	تعداد انواع (گروه) محصولات
۱۰	۱۰	تعداد محصولات (کالاها)
۲۷۰	۷/۵۷۲	تعداد نظرات
۳/۰۸۰	۳۱/۷۳۰	



(جدول ۴-): ارزیابی کمی شبکه‌های واژگان زبان فارسی  
(Table-4): Quantitative assessment Persian WordNets

فردوس نت	شبکه واژگان فارسی (دانشگاه تهران)	فارس نت ۲	شبکه واژگان انگلیسی	تعداد کلمات (یکتا)
۱۰۰۰۶۲	۱۸۱۶۶	۳۱۲۳۰	۱۵۵۲۸۷	تعداد کلمات (یکتا)
۹۱۶۴۰	۱۷۷۵۹	۲۰۴۳۲	۱۱۷۶۵۹	تعداد گروه‌های هم‌معنی
۹۱۶۴۰	۱۷۷۵۹	۱۷۳۰۰	۱۱۷۶۵۹	تعداد گروه کلمات هم‌معنی لینک شده با شبکه واژگان انگلیسی
۲/۰۳۳	۱/۷۱۵	۱/۸۵۳	۱/۷۵۹	متوسط تعداد کلمه در هر گروه کلمات هم‌معنی
۶۰٪	۵۸٪	۵۶٪	۵۴٪	درصد گروه‌های هم‌معنی تک‌کلمه‌ای
۱/۸۸	۱/۵۷	۱/۱۶	۱/۵۱۴	تعداد متوسط اشتراک هر کلمه در گروه‌های هم‌معنی مختلف

لغات فارسی به‌طوری بود که برای داوران مشخص نبود که هر کلمه فارسی متعلق به کدام شبکه واژگان است. در نهایت تعدادی از داوران خبره (در زمینه مفاهیم پردازش زبان طبیعی و شبکه واژگان) کلمات اشتباه هر گروه هم‌معنی را مشخص کردند. کلمات اشتباه کلماتی هستند که معادل مفهوم بیان شده در گروه هم‌معنی انگلیسی نباشند. سپس با استفاده از برچسب‌های داوران میزان دقت و بازخوانی و F1-Measure برای هر یک از گروه‌های هم‌معنی محاسبه شده‌است. نتایج ارزیابی دقت، بازخوانی و مقدار F1-Measure شبکه‌های واژگان مختلف برای مجموعه هزار عددی از گروه‌های هم‌معنی در جدول (۵) نمایش داده شده‌است.

با توجه به وجود درجه اطمینان به‌ازای هر یک از کلمات موجود در گروه‌های هم‌معنی در شبکه واژگان فردوس نت، در این شبکه واژگان ارزیابی کیفی به‌ازای بازه‌های اطمینان مختلف محاسبه شده‌است. در جدول (۵)، منظور از *conf* درجه اطمینان کلمات درون هر گروه هم‌معنی در شبکه واژگان فردوس نت است. همچنین باید توجه شود که در شبکه‌های واژگان فارسی، گروه‌های هم‌معنی فارسی معادل با برخی از گروه‌های هم‌معنی موجود در شبکه واژگان انگلیسی وجود ندارد. به این منظور در جدول (۶) یک‌بار بازخوانی برای این گروه‌ها صفر و یک‌بار نیز این گروه‌ها در محاسبه میانگین بازخوانی (موجود) کلی در نظر گرفته نشده‌اند. به‌عنوان مثال، برای شبکه واژگان فارس نت، در مجموعه هزار عددی گروه‌های هم‌معنی انتخاب‌شده از شبکه واژگان انگلیسی برای ارزیابی، تنها تعداد ۱۴۷ (حدود ۱۵٪) گروه هم‌معنی معادل در شبکه واژگان فارس نت وجود داشتند. دقت کلمات درون گروه‌های هم‌معنی موجود حدود ۰/۸۹۸ و میزان بازخوانی کلی (از هزار گروه هم‌معنی) ۰/۱۰۱ است؛ ولی اگر مقدار بازخوانی تنها برای همان ۱۴۷ گروه هم‌معنی (گروه‌های هم‌معنی که معادل آنها در فارس نت موجود است) محاسبه شود، مقدار بازخوانی این شبکه واژگان حدود ۰/۶۹۵ خواهد بود.

برای ارزیابی کیفی شبکه واژگان حسی فارسی (حس نگار) و واژه‌نامه حسی UTIIS، داوران حدود ۱۵۰ کلمه حسی را از متن نظرات استخراج کردند. با رویکردی مشابه پیش، دقت و بازخوانی شبکه واژگان حسی، محاسبه و نتایج آن در جدول (۶) بیان شده‌است.

برای محاسبه بازخوانی مشخص شد که از بین ۱۵۰ عبارت حسی، ۱۱۳ مورد در بین لغات شبکه واژگان حسی (حس نگار) وجود داشتند. پس به عبارت دیگر، میزان بازخوانی شبکه واژگان حسی حدود ۰/۷۵ است. البته اغلب عبارات

$S_i^{wn}$  نامیده می‌شود. در نهایت دقت و بازخوانی به‌ازای هر یک از گروه‌های کلمات هم‌معنی به‌صورت زیر محاسبه شده و میانگین آنها به‌عنوان دقت و بازخوانی نهایی شبکه واژگان در نظر گرفته می‌شود.

$$Recall = \frac{|S^* \cap S^{wn}|}{|S^*|} \quad \text{رابطه ۱}$$

$$Precision = \frac{|S^* \cap S^{wn}|}{|S^{wn}|} \quad \text{رابطه ۲}$$

$$F_1 - Measure = 2 \cdot \frac{precision \cdot recall}{precision + recall} \quad \text{رابطه ۳}$$

عبارت  $|S^* \cap S^{wn}|$  بیان‌کننده تعداد کلمات درست انتخاب‌شده در گروه کلمات هم‌معنی موردنظر است. عبارت  $|S^*|$  تعداد کلمات فارسی معادل با مفهوم اصلی و  $|S^{wn}|$  تعداد کلمات در گروه هم‌معنی مورد نظر در شبکه واژگان فارسی است.

در مرحله بعد، کلمات و توضیحات هر گروه هم‌معنی در شبکه واژگان انگلیسی در کنار فهرست لغات فارسی (معادل با آن گروه هم‌معنی) به شکل یک فهرست آماده شد. فهرست

(جدول-۵): ارزیابی کیفی شبکه‌های واژگان زبان فارسی  
(Table-5): The qualitative assessment of Persian WordNets

شبکه‌های واژگان زبان فارسی	پوشش گروه‌های هم‌معنی شبکه واژگان انگلیسی	دقت	بازخوانی	F1-Measure	بازخوانی موجود	F1- Measure موجود
فارس‌نت ۲	۱۵٪	۰/۸۹۸	۰/۱۰۱	۰/۱۸۲	۰/۶۹۵	۰/۷۸۴
شبکه واژگان فارسی دانشگاه تهران	۱۱٪	۰/۷۹۷	۰/۰۵۷	۰/۱۰۶	۰/۵۳۱	۰/۶۳۷
فردوس‌نت ( $conf \geq 0.8$ )	۴۵٪	۰/۸۲۶	۰/۱۹۹	۰/۳۲۱	۰/۴۸۳	۰/۶۱
فردوس‌نت ( $conf \geq 0.4$ )	۵۶٪	۰/۸۰۱	۰/۳۱۰	۰/۴۴۷	۰/۶۲۷	۰/۷۰۳
فردوس‌نت (همه کلمات)	۸۲٪	۰/۷۴۰	۰/۶۰۳	۰/۶۶۵	۰/۹۰۱	۰/۸۱۳

این، برای افزایش کیفیت روش‌های تحلیل نظرات، مجموعه‌ای از ابزارهای پیش‌پردازش زبان فارسی ایجاد و توسعه داده شدند. همچنین در این راستا، پیکره نظرات محصولات تجاری برای نظرکاوی در زبان فارسی توسعه پیدا کرد؛ سپس ارزیابی دقیقی بر روی کیفیت ابزارهای پردازش متن و شبکه واژگان فردوس‌نت (در مقایسه با سایر شبکه‌های واژگان موجود برای زبان فارسی) انجام شد. در نهایت نرخ دقت و بازخوانی شبکه واژگان حسّی فارسی نیز تخمین زده شد.

این مقاله علاوه بر معرفی ابزارهای پردازش زبان فارسی، مجموعه واژگان و الگوهای حسّی مناسب و پیکره نظرکاوی فارسی، مبنای خوبی برای انتخاب و استفاده از ویژگی‌ها، روش‌های انتخاب ویژگی و روش‌های مختلف دسته‌بندی برای ادامه پژوهش‌های نظرکاوی در زبان فارسی (و سایر زبان‌های مشابه مانند اردو و عربی) برای پژوهش‌گران این حوزه می‌تواند باشد.

### سیاس‌گزاری

در این قسمت، لازم است از زحمات دانشجویان گروه زبان‌شناسی و اعضای آزمایشگاه فناوری تارنمای دانشگاه فردوسی مشهد که در برچسب‌زنی پیکره نظرات و ارزیابی ابزارهای پردازش متون مختلف مشارکت داشتند، سپاس‌گزاری شود.

حسی که در شبکه واژگان وجود ندارند، مربوط به کلمات عامیانه یا اشکالات املایی کلمات حسّی متن نظرانی هستند که توسط ابزارهای پیش‌پردازش متن تولید شده، تصحیح نشده‌اند. از میان ۱۱۳ کلمه موجود در شبکه واژگان حسّی، نوع حس (از نظر مثبت یا منفی بودن)، ۹۷ کلمه مطابق با حس بیان شده توسط داوران بود. پس دقت شبکه واژگان حسّی فارسی ایجاد شده، حدود ۰/۸۶ است. برای کلماتی که در چند گروه هم‌معنی حضور داشتند، متوسط بار حسّی آنها در گروه‌های مختلف در نظر گرفته شده است.

(جدول-۶): ارزیابی کیفی واژه‌نامه حسّی زبان فارسی  
(Table-6): Qualitative assessment of the Persian sentiment lexicons

شبکه واژگان حسّی فارسی (حس‌نگار)	دقت	بازخوانی	F1-Measure
شبکه واژگان حسّی فارسی (حس‌نگار)	۰/۸۶	۰/۷۵	۰/۸
واژه‌نامه حسّی UTHS [32]	۰/۸۷۷	۰/۶۵۳	۰/۷۵

بخشی از خطای شبکه واژگان حسّی فارسی مربوط به خطای موجود در شبکه واژگان فارسی فردوس‌نت است. علاوه بر آن بار حسّی مشخص شده برای کلمات در شبکه واژگان حسّی انگلیسی نیز دارای مقداری خطا است که این خطا در شبکه واژگان حسّی فارسی نیز انتشار می‌یابد و دقت آن را کاهش می‌دهد [69].

## 7-References

## ۷- مراجع

- [1] P. D. Turney, "Thumbs up or thumbs down?: semantic orientation applied to unsupervised classification of reviews," in *Proceedings of the 40th annual meeting on association for*

<sup>۱</sup> برای استفاده‌های پژوهشی (غیرتجاری) از پیکره و ابزارهای پردازش متن فارسی می‌توانید به تارنمای آزمایشگاه فناوری وب دانشگاه فردوسی (<http://wtlab.um.ac.ir>) مراجعه نمایید.

## ۶- نتیجه‌گیری

در این پژوهش ابزارهای پردازش متن، شبکه واژگان فردوس‌نت و به دنبال آن شبکه واژگان حسّی زبان فارسی ایجاد و معرفی شدند. برای تولید شبکه واژگان حسّی (واژه‌نامه حسّی) فارسی از روش نگاشت و استفاده از شبکه واژگان حسّی انگلیسی برای زبان فارسی، با استفاده از طراحی شبکه واژگان فردوس‌نت (و سایر منابع زبانی) پیشنهاد شد. علاوه بر

- [12] C. Banea, R. Mihalcea, and J. Wiebe, "Porting Multilingual Subjectivity Resources Across Languages," *IEEE Transactions on Affective Computing*, vol. 4, 2013.
- [13] A. Balahur and M. Turchi, "Comparative Experiments Using Supervised Learning and Machine Translation for Multilingual Sentiment Analysis," *Computer Speech & Language*, vol. 2, pp. 56-75, 2013.
- [14] M. Okada and K. Hashimoto, "Investigation of Preprocessing of Multilingual Online Reviews for Automatic Classification," in *Computer and Information Science (ICIS), 2012 IEEE/ACIS 11th International Conference on*, 2012, pp. 306-309.
- [15] X. Ding, B. Liu, and P. S. Yu, "A holistic lexicon-based approach to opinion mining," in *Proceedings of the international conference on Web search and web data mining*, 2008, pp. 231-240.
- [16] M. Taboada, J. Brooke, M. Tofiloski, K. Voll, and M. Stede, "Lexicon-based methods for sentiment analysis," *Computational linguistics*, vol. 37, pp. 267-307, 2011.
- [17] J. Kamps, M. Marx, R. J. Mokken, and M. De Rijke, "Using wordnet to measure semantic orientations of adjectives," in *Proceedings of the 4th International Conference on Language Resources and Evaluation (LREC 2004)*, 2004, pp. 1115-1118.
- [18] A. Fahrni and M. Klenner, "Old wine or warm beer: Target-specific sentiment analysis of adjectives," in *Proc. of the Symposium on Affective Language in Human and Machine, AISB*, 2008, pp. 60-63.
- [19] V. Hatzivassiloglou and K. R. McKeown, "Predicting the semantic orientation of adjectives," in *Proceedings of the eighth conference on European chapter of the Association for Computational Linguistics*, 1997, pp. 171-174.
- [20] N. Kaji and M. Kitsuregawa, "Building Lexicon for Sentiment Analysis from Massive Collection of HTML Documents," in *EMNLP-CoNLL*, 2007, pp. 1075-1083.
- [21] L. Velikovich, S. Blair-Goldensohn, K. Hannan, and R. McDonald, "The viability of web-derived polarity lexicons," in *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, 2010, pp. 777-785.
- [22] H. Takamura, T. Inui, and M. Okumura, "Extracting semantic orientations of words using spin model," in *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, 2005, pp. 133-140.
- [2] B. Pang and L. Lee, "Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales," presented at the Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics, 2005.
- [3] E. Riloff and J. Wiebe, "Learning extraction patterns for subjective expressions," presented at the Proceedings of the 2003 conference on Empirical methods in natural language processing, 2003.
- [4] S.-M. Kim and E. Hovy, "Extracting opinions, opinion holders, and topics expressed in online news media text," presented at the Proceedings of the Workshop on Sentiment and Subjectivity in Text, 2006.
- [5] C. O. Alm, "Subjective natural language problems: motivations, applications, characterizations, and implications," presented at the Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: short papers-Volume 2, 2011.
- [6] L. Barbosa and J. Feng, "Robust sentiment detection on twitter from biased and noisy data," presented at the Proceedings of the 23rd International Conference on Computational Linguistics: Posters, 2010.
- [7] M. Abdul-Mageed, M. Diab, and M. Korayem, "Subjectivity and sentiment analysis of modern standard Arabic," presented at the Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, 2011.
- [8] I. Habernal, T. Ptáček, and J. Steinberger, "Supervised sentiment analysis in Czech social media," *Information Processing & Management*, vol. 50, pp. 693-707, 2014.
- [9] H. Guo, H. Zhu, Z. Guo, X. Zhang, and Z. Su, "OpinionIt: a text mining system for cross-lingual opinion analysis," in *Proceedings of the 19th ACM international conference on Information and knowledge management*, 2010, pp. 1199-1208.
- [10] D. Gao, F. Wei, W. Li, X. Liu, and M. Zhou, "Cross-lingual Sentiment Lexicon Learning With Bilingual Word Graph Label Propagation," *Computational Linguistics*, vol. 41, pp. 21-40, 2015.
- [11] M.-T. Martín-Valdivia, E. Martínez-Cámara, J.-M. Perea-Ortega, and L. Alfonso Ureña-López, "Sentiment polarity detection in Spanish reviews combining supervised and unsupervised approaches," *Expert Systems with Applications*, vol. 40, pp. 3934-3942, 2012.

- [33] S. Baccianella, A. Esuli, and F. Sebastiani, "SentiWordNet 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining," in *LREC*, 2010, pp. 2200-2204.
- [34] A. Esuli and F. Sebastiani, "Sentiwordnet: A publicly available lexical resource for opinion mining," in *Proceedings of 5th International Conference on Language Resources and Evaluation (LREC)*, Genoa, 2006, pp. 417-422.
- [35] A. Neviarouskaya, H. Prendinger, and M. Ishizuka, "SentiFul: A lexicon for sentiment analysis," *IEEE Transactions on Affective Computing*, vol. 2, pp. 22-36, 2011.
- [36] A. Neviarouskaya, H. Prendinger, and M. Ishizuka, "Textual affect sensing for sociable and expressive online communication," in *International Conference on Affective Computing and Intelligent Interaction*, 2007, pp. 218-229.
- [37] C. Strapparava and A. Valitutti, "WordNet Affect: an Affective Extension of WordNet," in *Proceedings of the 4th International Conference on Language Resources and Evaluation (LREC)*, 2004, pp. 1083-1086.
- [38] E. Cambria, R. Speer, C. Havasi, and A. Hussain, "SenticNet: A Publicly Available Semantic Resource for Opinion Mining," in *AAAI fall symposium: commonsense knowledge*, 2010.
- [39] M. E. Basiri, A. R. Naghsh-Nilchi, and N. Ghassem-Aghaee, "A Framework for Sentiment Analysis in Persian," *Open Transactions on Information Processing*, vol. 1, pp. 1-14, 2014.
- [40] F. Amiri, S. Scerri, and M. H. Khodashahi, "Lexicon-based Sentiment Analysis for Persian Text," in *Recent Advances in Natural Language Processing*, 2015, pp. 9-16.
- [41] 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.
- [42] س.ع. مردانی و ع. آقایی، "راثة روش نظارتی برای نظرکاوی در زبان فارسی با استفاده از لغتنامه و الگوریتم SVM"، فصلنامه علمی-پژوهشی مدیریت فناوری اطلاعات، دوره ۷ (۲)، ۱۳۹۳، ۳۶۲-۳۴۵.
- [42] A. Mardani and S.A. Aghaie "A supervised method for opinion mining in Persian using lexicon and SVM algorithm", in *National Journal of Information Technology Management*, 2015(7), pp. 345-362.
- [43] S. Cerini, V. Compagnoni, A. Demontis, M. Formentelli, and G. Gandini, "Micro-WNOp: A gold standard for the evaluation of automatically compiled lexical resources for opinion mining,"
- [23] A. Esuli and F. Sebastiani, "Pageranking wordnet synsets: An application to opinion mining," presented at the Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics (ACL), Prague, Czech Republic, 2007.
- [24] 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, 2009, pp. 675-682.
- [25] S. Poria, A. Gelbukh, A. Hussain, N. Howard, D. Das, and S. Bandyopadhyay, "Enhanced SenticNet with affective labels for concept-based opinion mining," *IEEE Intelligent Systems*, vol. 28, pp. 31-38, 2013.
- [26] S. Gindl, A. Weichselbraun, and A. Scharl, "Extracting and Grounding Contextualized Sentiment Lexicons," 2013.
- [27] D. Tang, F. Wei, B. Qin, M. Zhou, and T. Liu, "Building Large-Scale Twitter-Specific Sentiment Lexicon: A Representation Learning Approach," in *the 25th International Conference on Computational Linguistics (COLING)*, 2014, pp. 172-182.
- [۲۸] س. نوفرستی، س. و. م. شمس فرد، "ساخت نیمه خودکار یک پیکره از نظرات غیرمستقیم در دامنه دارو و به کارگیری آن برای تعیین قطبیت نظرات". چاپ شده در مجله پردازش و علائم داده‌ها، دوره ۱۳ (۲)، سال ۱۳۹۵، ۴۹-۳۵.
- [28] S. Nofarsti and M. Shamsfard, "Automatic building a corpus and exploiting it for polarity classification of indirect opinions about drugs," in *Journal of Signal and Data Processing (JSDP)*, 2016; 13 (2), pp.35-49.
- [29] H. Kanayama and T. Nasukawa, "Fully automatic lexicon expansion for domain-oriented sentiment analysis," in *Proceedings of the 2006 conference on empirical methods in natural language processing*, 2006, pp. 355-363.
- [30] A. Hassan and D. Radev, "Identifying text polarity using random walks," in *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, 2010, pp. 395-403.
- [31] A. Hassan, A. Abu-Jbara, R. Jha, and D. Radev, "Identifying the semantic orientation of foreign words," in Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: short papers-Volume 2, 2011, pp. 592-597.
- [32] I. Dehdarbehbahani, A. Shakery, and H. Faili, "Semi-supervised word polarity identification in resource-lean languages," *Neural Networks*, vol. 58, pp. 50-59, 2014.



- [56] A. Balali, A. Rajabi, S. Ghassemi, M. Asadpour, and H. Faili, "Content diffusion prediction in social networks," in *5th Conference on Information and Knowledge Technology (IKT)*, 2013, pp. 467-471.
- [57] P. Turney, "Mining the web for synonyms: PMI-IR versus LSA on TOEFL," in *12th European Conference on Machine Learning (ECML 2001)*, Freiburg, Germany, 2001, pp. 491-502.
- [58] K. Denecke, "Using sentiwordnet for multilingual sentiment analysis," in *Data Engineering Workshop, 2008. ICDEW 2008. IEEE 24th International Conference on*, 2008, pp. 507-512.
- [59] C. M. Özsert and A. Özgür, "Word polarity detection using a multilingual approach," in *Computational Linguistics and Intelligent Text Processing*, ed: Springer, 2013, pp. 75-82.
- [60] J. Steinberger, M. Ebrahim, M. Ehrmann, A. Hurriyetoglu, M. Kabadjov, P. Lenkova, et al., "Creating sentiment dictionaries via triangulation," *Decision Support Systems*, vol. 53, pp. 689-694, 2012.
- [61] F. L. Cruz, J. A. Troyano, B. Pontes, and F. J. Ortega, "Building layered, multilingual sentiment lexicons at synset and lemma levels," *Expert Systems with Applications*, vol. 41, pp. 5984-5994, 2014.
- [62] F. H. Mahyoub, M. A. Siddiqui, and M. Y. Dahab, "Building an Arabic Sentiment Lexicon Using Semi-Supervised Learning," *Journal of King Saud University-Computer and Information Sciences*, vol. 26, pp. 417-424, 2014.
- [63] Y. Chen and S. Skiena, "Building sentiment lexicons for all major languages," in *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Short Papers)*, 2014, pp. 383-389.
- [64] M. Shamsfard, "Challenges and open problems in Persian text processing," *Proceedings of LTC*, vol. 11, 2011.
- [65] W. Feely, M. Manshadi, R. Frederking, and L. Levin, "The CMU METAL Farsi NLP Approach," in *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, 2014, pp. 4052-4055.
- [66] R. Duwairi and M. El-Orfali, "A study of the effects of preprocessing strategies on sentiment analysis for Arabic text," *Journal of Information Science*, vol. 40, pp. 501-513, 2014.
- [67] W. Chamlerwat, P. Bhattarakosol, T. Rungkasiri, and C. Haruechaiyasak, "Discovering Language resources and linguistic theory: Typology, second language acquisition, English linguistics," pp. 200-210, 2007.
- [44] A. Montejo-Ráez, E. Martínez-Cámara, M. T. Martín-Valdivia, and L. A. Ureña-López, "Ranked wordnet graph for sentiment polarity classification in twitter," *Computer Speech & Language*, vol. 28, pp. 93-107, 2014.
- [45] M. Montazery and H. Faili, "Automatic Persian wordnet construction," in *Proceedings of the 23rd International Conference on Computational Linguistics: Posters*, 2010, pp. 846-850.
- [46] K. N. Lam, F. A. Tarouti, and J. Kalita, "Automatically constructing Wordnet synsets," in *52nd Annual Meeting of the Association for Computational Linguistics (ACL 2014)*, Baltimore, USA, 2014.
- [47] M. Shamsfard, A. Hesabi, H. Fadaei, N. Mansoory, A. Famian, S. Bagherbeigi, et al., "Semi automatic development of farsnet; the persian wordnet," in *Proceedings of 5th Global WordNet Conference, Mumbai, India*, 2010.
- [48] P. Vossen, "A multilingual database with lexical semantic networks," *Computational Linguistics* vol. 25, pp. 628-630, 1998.
- [49] F. Keyvan, H. Borjian, M. Kasheff, and C. Fellbaum, "Developing persianet: The persian wordnet," in *3rd Global wordnet conference*, 2007, pp. 315-318.
- [50] A. Famian and D. Aghajaney, "Towards Building a WordNet for Persian Adjectives," *International Journal of Lexicography*, pp. 307-308, 2006.
- [51] M. Fadaee, H. Ghader, H. Faili, and A. Shakery, "Automatic WordNet Construction Using Markov Chain Monte Carlo," *Polibits*, pp. 13-22, 2013.
- [52] N. Taghizadeh and H. Faili, "Automatic Wordnet Development for Low-resource Languages using Cross-lingual WSD," *Journal of Artificial Intelligence Research*, vol. 56, pp. 61-87, 2016.
- [53] F. Mahdisoltani, J. Biega, and F. Suchanek, "YAGO3: A knowledge base from multilingual Wikipedias," in *7th Biennial Conference on Innovative Data Systems Research*, 2014.
- [54] A. AleAhmad, H. Amiri, E. Darrudi, M. Rahgozar, and F. Oroumchian, "Hamshahri: A standard Persian text collection," *Knowledge-Based Systems*, vol. 22, pp. 382-387, 2009.
- [55] H. Eghbalzadeh, B. Hosseini, S. Khadivi, and A. Khodabakhsh, "Persica: A Persian corpus for multi-purpose text mining and Natural language processing," in *Telecommunications (IST), 2012 Sixth International Symposium on*, 2012, pp. 1207-1214.

دانشگاه است. ایشان تاکنون مقالات متعدد در زمینه‌های مورد علاقه‌اش همچون کاربردشناسی زبان، ساخت واژه و رده‌شناسی زبان نگاشته که در مجلات علمی پژوهشی داخلی و خارجی یا مجموعه مقالات همایش‌ها منتشر شده‌اند و رساله‌ها و پایان‌نامه‌های متعددی را در دانشگاه فردوسی مشهد و سایر دانشگاه‌های کشور هدایت و راهنمایی کرده است. نشانی رایانامه ایشان عبارت است از:

[sh-sharifi@um.ac.ir](mailto:sh-sharifi@um.ac.ir)

Consumer Insight from Twitter via Sentiment Analysis," *J. UCS*, vol. 18, pp. 973-992, 2012.

[68] M.-T. Martín-Valdivia, E. Martínez-Cámara, J.-M. Perea-Ortega, and L. A. Ureña-López, "Sentiment polarity detection in Spanish reviews combining supervised and unsupervised approaches," *Expert Systems with Applications*, vol. 40, pp. 3934-3942, 2013.

[69] K. Denecke, "Are SentiWordNet scores suited for multi-domain sentiment classification?," presented at the Fourth International Conference on Digital Information Management, (ICDIM 2009), 2009.



**احسان عسکریان** دانشجوی دکترای

نرم‌افزار دانشگاه فردوسی است. وی مدرک کارشناسی ارشد خود را در گرایش نرم‌افزار از دانشگاه صنعتی شریف دریافت کرد. زمینه پژوهشی مورد علاقه‌ایشان نظر کاوی،

پردازش متن و داده‌کاوی است. ایشان از سال ۱۳۹۰ با عضویت در آزمایشگاه فناوری وب در پروژه‌های مختلف داده‌کاوی و پردازش متن مشارکت داشته و ابزارهای پایه‌ای متنوعی در رابطه با پردازش متن ایجاد کرده است. نشانی رایانامه ایشان عبارت است از:

[ehsan.asgarian@mail.um.ac.ir](mailto:ehsan.asgarian@mail.um.ac.ir)



**محسن کاهانی** استاد گروه مهندسی

کامپیوتر دانشگاه فردوسی مشهد و مدیر آزمایشگاه فناوری وب است. ایشان دکترای خود را در رشته مهندسی کامپیوتر از دانشگاه ولونگونگ استرالیا در سال ۱۳۷۷

اخذ کرده است. زمینه‌های پژوهشی مورد علاقه ایشان شامل وب‌معنایی، پردازش زبان طبیعی، سامانه‌های تصمیم‌یار و مهندسی نرم‌افزار است.

نشانی رایانامه ایشان عبارت است از:

[kahani@um.ac.ir](mailto:kahani@um.ac.ir)



**شهلا شریفی** مدرک کارشناسی ارشد

و دکترای خود را در رشته زبان‌شناسی همگانی از دانشگاه فردوسی مشهد دریافت کرده و در حال حاضر با درجه دانشیاری عضو هیئت علمی گروه زبان‌شناسی همین