

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

# دانشگاه یزد

پردیس فنی و مهندسی

گروه مهندسی کامپیوتر

پایان نامه برای دریافت درجه کارشناسی ارشد

نرم افزار

## تشخیص رویداد از وبسایت های خبری با استفاده از

### خوشه بندی

استادان راهنما: دکتر سجاد ظریف زاده و دکتر امیر جهانگرد رفسنجانی

پژوهش و نگارش: الهام رسولی

مهرماه ۱۳۹۷

کلیه‌ی حقوق مادی و معنوی مترتب بر نتایج مطالعات، ابتکارات و نوآوری‌های ناشی از تحقیق موضوع این پایان‌نامه/رساله متعلق به دانشگاه یزد است و هرگونه استفاده از نتایج علمی و عملی از این پایان‌نامه/رساله برای تولید دانش فنی، ثبت اختراع، ثبت اثر بدیع هنری، همچنین چاپ و تکثیر، نسخه‌برداری، ترجمه و اقتباس و ارائه مقاله در سمینارها و مجلات علمی از این پایان‌نامه/رساله منوط به موافقت کتبی دانشگاه یزد است.

## تقدیم به

### پدر و مادر عزیزم

و همه کسانی که درست اندیشیدن را به من آموختند.

## سپاس‌گزاری

سپاس خداوند یکتای عزتمندی که رحمت و دانش او در سراسر گیتی گسترده شده، آسمان‌ها و زمین همه از آن اوست و علم و دانش حقیقی را بر هر که بخواهد موهبت می‌فرماید. رحمت و لطف او را بی‌نهایت سپاس می‌گویم چرا که فهم و درک مطالب این پژوهش را بر من ارزانی داشت و مرا به این اصل رساند که علم و ایمان دو بال یک پروازند. توفیق تلاش به من داد و هر بار که خطا کردم فرصتی دوباره، تا با امید، تلاشی تازه را آغاز کنم. سپاس بیکران بر همدلی و همراهی و همگامی پدر و مادر دلسوز و مهربانم و با تقدیر و تشکر شایسته از استادان فرهیخته جناب آقای دکتر سجاد ظریف زاده و جناب آقای دکتر امیر جهانگرد رفسنجانی که راهنمایی این پایان‌نامه را به عهده گرفتند و در تمام مراحل مرا از راهنمایی‌های ارزشمند خود بهره‌مند ساختند. همچنین از همه‌ی استادان دوران تحصیل به ویژه جناب آقای دکتر علی‌محمد زارع بیدکی به خاطر تدریس فوق العاده‌ی ایشان و از در پایان همه‌ی دوستانم که در طی این مسیر با من همدل و همراه بودند کمال تشکر و قدردانی را دارم.

## چکیده

با افزایش تولید محتوا و اسناد الکترونیکی در وب نیاز به ابزارهای استخراج و پالایش اطلاعات بیشتر احساس می‌شود. چرا که یافتن اطلاعات مورد نظر از میان انبوه اطلاعات کاری دشوار و زمان‌بر است. از این رو، ابزارهای مختلفی جهت تسهیل دسترسی به اطلاعات مورد نظر کاربران ایجاد شده است. در همین راستا تشخیص و ردیابی رویداد جهت تشخیص جدیدترین و مهم‌ترین رویدادهای جهان واقعی مورد مطالعه‌ی پژوهشگران قرار گرفته است. از اطلاعات بدست آمده از این طریق می‌توان در کاربردهای مختلف نظیر تصمیم‌گیری‌های تجاری، پیش‌بینی نتایج انتخابات و پیش‌بینی شیوع بیماری استفاده کرد. به طور کلی، تشخیص رویداد به دو دسته‌ی تشخیص رویداد گذشته‌نگر و تشخیص رویداد برخط تقسیم می‌شود. در تشخیص رویداد گذشته‌نگر، هدف یافتن رویدادها در مجموعه‌ای از اسناد از پیش جمع‌آوری شده است. درحالیکه تشخیص رویداد برخط بر شناسایی رویدادهای جدید در جریان‌های خبری برخط تمرکز دارد. در این پایان‌نامه، به تشخیص رویداد برخط در مجموعه‌ی اسناد خبری با استفاده از روش‌های تشخیص جامعه که زیر مجموعه‌ای از روش‌های خوشه‌بندی هستند، پرداخته می‌شود. روش پیشنهادی که گراف وزن‌دار عبارات انفجاری (*WebKey*) نام دارد، شبکه‌ای از عبارات انفجاری ایجاد می‌کند که این شبکه براساس هم‌رخدادی عبارات اسناد در بازه‌ی زمانی کوتاه تشکیل می‌شود. سپس دو ویژگی جدید شامل کلیک‌های کاربران بر اسناد و بر عناوین خبری از دادگان استخراج و با به کارگیری آن‌ها در الگوریتم تشخیص جامعه، خوشه‌های عبارات توصیف‌کننده‌ی رویداد شناسایی می‌شوند. روش پیشنهادی در مقایسه با روش‌های معروف گذشته فراخوانی را ۱۹ درصد و دقت را ۲ برابر افزایش می‌دهد.

**کلیدواژه:** تشخیص و ردیابی موضوع، تشخیص رویداد، وب کاوی، تشخیص جامعه

# فهرست مطالب

ه	فهرست جداول
ز	فهرست تصاویر
ط	فهرست اختصارات
۱	۱ مقدمه
۲	۱-۱ پیشگفتار
۳	۲-۱ تاریخچه
۴	۳-۱ کاربردها
۴	۴-۱ چالش‌ها
۴	۱-۴-۱ وابستگی حوزه
۵	۲-۴-۱ محدودیت زمانی
۵	۳-۴-۱ صحت تشخیص
۵	۴-۴-۱ منابع داده‌ای متنوع
۶	۵-۴-۱ انبوه داده‌ها
۶	۵-۱ تعریف مساله و دستاوردها
۷	۶-۱ ساختار کلی پایان‌نامه
۹	۲ مفاهیم اولیه
۱۰	۱-۲ مقدمه
۱۰	۲-۲ رویداد چیست؟
۱۲	۳-۲ تشخیص گذشته‌نگر
۱۲	۴-۲ تشخیص برخط

۵-۲	مراحل اصلی	۱۳
۱-۵-۲	قطعه‌بندی نقل	۱۴
۲-۵-۲	تشخیص موضوع	۱۵
۳-۵-۲	ردیابی موضوع	۱۵
۴-۵-۲	تشخیص پیوند	۱۵

### ۳ پیشینه پژوهش ۱۷

۱-۳	مقدمه	۱۸
۲-۳	تشخیص رویداد	۱۸
۱-۲-۳	پیش‌پردازش داده‌ها	۱۸
۲-۲-۳	نمایش داده‌ها	۱۹
۳-۲-۳	سازماندهی داده‌ها	۱۹
۳-۳	ردیابی رویداد	۲۰
۴-۳	روش‌های تشخیص رویداد	۲۰
۱-۴-۳	روش‌های مبتنی بر نحوه‌ی مدل‌سازی دادگان	۲۰
۱-۱-۴-۳	روش‌های سند محور	۲۱
۲-۱-۴-۳	روش‌های ویژگی محور	۲۱
۳-۱-۴-۳	مدل‌های موضوع	۲۴
۵-۳	روش‌های متداول	۲۵
۱-۵-۳	روش‌های خوشه‌بندی	۲۵
۱-۱-۵-۳	خوشه‌بندی سلسله‌مراتبی	۲۵
۲-۱-۵-۳	خوشه‌بندی افزایشی	۲۶
۳-۱-۵-۳	خوشه‌بندی $K$ -میانگین	۲۷
۴-۱-۵-۳	خوشه‌بندی $K$ -میانه	۲۷
۵-۱-۵-۳	$K$ -میانگین++	۲۷

## ب



۲۸	۳-۵-۲ روش‌های مبتنی بر الگو
۲۹	۳-۵-۳ روش‌های مبتنی بر گراف
۳۱	۳-۵-۳-۱ الگوریتم‌های تقسیم‌بندی گراف
۳۱	۳-۵-۳-۲ خوشه‌بندی سلسله‌مراتبی
۳۱	۳-۵-۳-۳ خوشه‌بندی تفکیکی
۳۲	۳-۵-۳-۴ مرکزیت میانی
۳۶	۳-۶ شبکه‌های اجتماعی و میکرو بلاگ‌ها
۳۷	۳-۶-۱ تشخیص و ردیابی رویداد در توییتر
۳۸	۳-۷ دادگان پرکاربرد

#### ۴ روش پیشنهادی

۴۱	۴-۱ مقدمه
۴۲	۴-۲ پیش‌پردازش
۴۳	۴-۳ محاسبه‌ی ویژگی‌های اولیه
۴۴	۴-۴ تشخیص انفجار
۴۵	۴-۵ محاسبه‌ی ویژگی‌های انفجاری
۴۶	۴-۶ ساخت گراف
۴۷	۴-۷ نمونه‌برداری از گراف
۴۸	۴-۸ تشخیص رویداد

#### ۵ پیاده‌سازی و ارزیابی

۵۱	۵-۱ مقدمه
۵۲	۵-۲ دادگان
۵۲	۵-۲-۱ ویژگی‌های دادگان
۵۳	۵-۲-۲ مشخصات دادگان
۵۶	۵-۳ پیاده‌سازی

۵-۳-۱ مشخصات سیستم . . . . . ۵۶

۵-۳-۲ مشخصات برنامه . . . . . ۵۷

۵-۴ مدل ارزیابی . . . . . ۵۷

۵-۵ نتایج ارزیابی . . . . . ۶۰

۵-۵-۱ مرحله‌ی آموزش . . . . . ۶۰

۵-۵-۲ مرحله‌ی آزمون . . . . . ۶۷

۶ نتیجه‌گیری و پیشنهادها ۶۹

۶-۱ جمع‌بندی و نتیجه‌گیری . . . . . ۷۰

۶-۲ پیشنهادها . . . . . ۷۱

۷۴ واژه نامه فارسی به انگلیسی

۷۵ واژه نامه انگلیسی به فارسی

۷۷ منابع و مآخذ

# فهرست جداول

۱-۳	مقادیر آستانه در روش صیادی	۳۴
۲-۳	دادگان TDT	۳۹
۱-۵	مشخصات آماری دادگان	۵۴
۲-۵	جدول پایگاه‌های خبری برخط خزش شده توسط جویشرگر پارسی‌جو	۵۵
۳-۵	تعداد موجودیت نامدار در دادگان	۵۶
۴-۵	مقادیر اولیه برای شروع الگوریتم BKG	۶۲
۵-۵	جزئیات روش BKG برحسب گرم	۶۳
۶-۵	جزئیات روش BKG برحسب نمونه‌برداری	۶۴
۷-۵	مقادیر اولیه برای شروع الگوریتم WebKey	۶۵
۸-۵	معیارهای ارزیابی رویداد در روش WebKey	۶۵
۹-۵	معیارهای ارزیابی عبارت رویداد در روش WebKey	۶۵
۱۰-۵	تاثیر تقویت موجودیت نامدار بر معیارهای ارزیابی روش WebKey	۶۵
۱۲-۵	تاثیر نمونه‌برداری بر معیارهای ارزیابی روش WebKey	۶۶
۱۱-۵	مصرف زمان و حافظه با نمونه‌برداری و بدون نمونه‌برداری در روش WebKey	۶۶
۱۳-۵	مقایسه‌ی روش‌ها از نظر زمان و حافظه مصرفی	۶۸

# فهرست تصاویر

- ۱-۲ سلسله‌مراتب مفاهیم خبری در TDT . . . . . ۱۱
- ۲-۲ تشخیص رویدادهای جدید در جریان نقل‌های خبری . . . . . ۱۳
- ۱-۳ گرافی با سه جامعه . . . . . ۳۰
- ۱-۴ معماری کلی سیستم پیشنهادی (*WebKey*) . . . . . ۴۲
- ۲-۴ رویدادهای شناسایی شده توسط روش *WebKey* . . . . . ۵۰
- ۱-۵ تصویری از فایل خزش شده توسط جویشرگر پارسی جو . . . . . ۵۳
- ۲-۵ تصویری از دادگان مرجع درستی . . . . . ۵۹
- ۳-۵ نمودار دقت و فراخوانی برای روش BKG بر حسب  $\rho$  . . . . . ۶۲
- ۴-۵ نمودار دقت و فراخوانی برای روش BKG بر حسب گرم . . . . . ۶۳
- ۵-۵ نمودار دقت و فراخوانی برای روش BKG بر حسب  $\eta$  . . . . . ۶۳
- ۶-۵ نمودار دقت و فراخوانی برای روش BKG بر حسب نمونه‌برداری . . . . . ۶۴
- ۷-۵ نمودار دقت و فراخوانی برای روش *WebKey* بر حسب گرم . . . . . ۶۶
- ۸-۵ مقایسه‌ی روش پیشنهادی با سایر روش‌ها . . . . . ۶۸
- ۱-۶ نمودار دقت و فراخوانی برای روش BKG بر حسب گرم . . . . . ۷۳

# فهرست اختصارات

Area-based Detection Model .....	<b>ADM</b>
Application Programming Interface .....	<b>API</b>
Breath First Search .....	<b>BFS</b>
Defense Advanced Research Projects Agency .....	<b>DARPA</b>
First Story Detection .....	<b>FSD</b>
Group Average Clustering .....	<b>GAC</b>
Hierarchical agglomerative clustering .....	<b>HAC</b>
Latent Dirichlet Allocation .....	<b>LDA</b>
Latent Semantic Indexing .....	<b>LSI</b>
Node-based Detection Model .....	<b>NDM</b>
New Event Detection .....	<b>NED</b>
Natural Language Processing .....	<b>NLP</b>
Probabilistic Latent Semantic Indexing .....	<b>PLSI</b>
Retrospective Event Detection .....	<b>RED</b>
Singular value decomposition .....	<b>SVD</b>
Support Vector Machine .....	<b>SVM</b>
Topic Detection and Tracking .....	<b>TDT</b>
Term Frequency-Inverse Document Frequency .....	<b>TF-IDF</b>
Translingual Information Detection, Extraction, and Summarization .....	<b>TIDES</b>

# فصل ۱

## مقدمه

## ۱-۱ پیشگفتار

خبر علاوه بر نقش آگاهی‌دهنده از نظر تأثیری که بر زندگی افراد جامعه می‌گذارد از اهمیت بالایی برخوردار است. از خبر می‌توان برای پیش‌بینی وضعیت بازار بورس و اقتصاد، کمک در تصمیم‌گیری‌های سیاسی، اجتماعی، اقتصادی و غیره بهره برد. اخبار از گذشته تا به امروز به روش‌های مختلفی منتشر و در دسترس عموم قرار گرفته است. با بوجود آمدن پدیده‌ی اینترنت و بستر وب، تولید و انتشار اخبار به راحتی برای همه‌ی اقشار جامعه امری آسان و راحت شد. از طرفی، رشد روزافزون اینترنت و انتشار اخبار به صورت اسناد الکترونیکی در شبکه‌های اجتماعی و انواع پایگاه‌های خبری برخط سبب شده است تا پیگیری و یافتن اطلاعات مشخص در وب، کاری زمان‌بر و دشوار شود. همچنین تشخیص تمام موضوعات به‌روز توسط یک فرد، امری غیرممکن است. بنابراین پردازش این حجم عظیم از اطلاعات نیازمند روش‌ها و ابزارهای جدید است. این مسئله، چالش‌هایی در حوزه‌ی بازیابی اطلاعات<sup>۱</sup> ایجاد می‌کند.

بازیابی و استخراج خبر<sup>۲</sup> زیرمجموعه‌ای از بازیابی اطلاعات و همچنین استخراج اطلاعات<sup>۳</sup> است. این دو مفهوم دو حوزه‌ی مجزا هستند. بازیابی اطلاعات به بازیابی اسناد مرتبط به پرس‌وجوی کاربر گفته می‌شود، درحالی‌که استخراج اطلاعات به استخراج خودکار اطلاعات ساخت‌یافته<sup>۴</sup> از اطلاعات غیرساخت‌یافته<sup>۵</sup> یا نیمه‌ساخت‌یافته<sup>۶</sup> اطلاق می‌شود. بازیابی و استخراج خبر به چندین حوزه‌ی تحقیقاتی مجزا تقسیم می‌شود.

تشخیص و ردیابی موضوع (TDT<sup>۷</sup>) زیر مجموعه‌ای از بازیابی و استخراج اطلاعات است. انگیزه‌ی اولیه برای تحقیقات در این حوزه، فراهم کردن یک فناوری مرکزی برای ایجاد سامانه‌ای است که با نظارت بر اخبار منتشر شده، رخدادهای جالبی را که در جهان رخ می‌دهند، اعلام کند [۱]. از طرفی، شرکت‌ها و افراد به منظور توسعه‌ی الگوریتم‌های هوشمند، علاقه‌مند به استفاده از این دادگان هستند [۲]. بنابراین هدف اصلی از تشخیص و ردیابی موضوع، توسعه‌ی فناوری‌هایی است که بتوانند متون خبری را در انواع رسانه‌های خبری جست‌وجو، سازماندهی و ساختاردهی کنند [۱]. معمولاً در این حوزه از روش‌های متن‌کاوی<sup>۸</sup>، پردازش زبان

<sup>1</sup>Information retrieval

<sup>2</sup>News retrieval and extraction

<sup>3</sup>Information extraction

<sup>4</sup>Structured

<sup>5</sup>Unstructured

<sup>6</sup>Semi-structured

<sup>7</sup>Topic Detection and Tracking

<sup>8</sup>Text mining

طبیعی (NLP)<sup>۱</sup> و تکنیک‌های خوشه‌بندی<sup>۲</sup> استفاده می‌شود. موضوع، مفهوم کلی‌تری از رویداد است. یک موضوع می‌تواند شامل چندین رویداد باشد. رویداد به رخدادی گفته می‌شود که در زمان و مکان مشخصی روی می‌دهد. از این رو، در پژوهش‌های اخیر از عبارت تشخیص و ردیابی رویداد نیز استفاده می‌شود.

## ۲-۱ تاریخچه

تشخیص و ردیابی موضوع از یک مطالعه‌ی مقدماتی<sup>۳</sup> در سال ۱۹۹۶ شروع شد. این تحقیق توسط آژانس پروژه‌های تحقیقاتی پیشرفته‌ی وزارت دفاع دولت ایالات متحده‌ی آمریکا (DARPA)<sup>۴</sup> حمایت می‌شد و در واقع زیر مجموعه‌ای از برنامه‌ی TIDES<sup>۵</sup> در DARPA به شمار می‌رفت [۳، ۴]. سایر مشارکت‌کنندگان در این تحقیق علاوه بر DARPA، شامل دانشگاه کارنگی ملون<sup>۶</sup>، دراگون سیستمز<sup>۷</sup> و دانشگاه ماساچوست<sup>۸</sup> بودند. پژوهش فوق داده‌های خود را از رسانه‌های خبری<sup>۹</sup> و با استفاده از سامانه‌های تبدیل گفتار به متن<sup>۱۰</sup> که خبرهای تلویزیون و رادیو را به طور خودکار به متن تبدیل می‌کنند جمع‌آوری می‌کرد. درابتدا، پژوهش‌ها حول مفهوم کلی موضوع<sup>۱۱</sup> صورت می‌گرفت، اما در ادامه‌ی تحقیقات در سال ۱۹۹۸ مفهوم رویداد<sup>۱۲</sup> مطرح شد. هدف از تشخیص و ردیابی موضوع، شکستن متن به نقل<sup>۱۳</sup> جدید است تا با نظارت بر این نقل‌ها بتواند رویدادهایی که قبلاً دیده نشده‌اند شرا ناسایی کند و این نقل‌ها را به گروه‌هایی تقسیم کند که موضوع بحث یکسانی دارند [۱]. طی این تحقیقات چندین پیکره<sup>۱۴</sup> از خبرهای جمع‌آوری شده از انواع رسانه‌های خبری به زبان‌های مختلف تهیه شد که محققان در تحقیقات خود از آن‌ها استفاده می‌کردند. همچنین از این مجموعه دادگان<sup>۱۵</sup> برای مسابقات TDT استفاده می‌شد. با پیدایش شبکه‌های اجتماعی<sup>۱۶</sup> چالش‌ها و مفاهیم جدیدی

<sup>۱</sup>Natural Language Processing

<sup>۲</sup>Clustering techniques

<sup>۳</sup>Pilot study

<sup>۴</sup>Defense Advanced Research Projects Agency

<sup>۵</sup>Translingual Information Detection, Extraction, and Summarization

<sup>۶</sup>Carnegie Mellon University

<sup>۷</sup>Dragon Systems

<sup>۸</sup>University of Massachusetts at Amherst

<sup>۹</sup>Newswire

<sup>۱۰</sup>Speech-to-text

<sup>۱۱</sup>Topic

<sup>۱۲</sup>Event

<sup>۱۳</sup>Story

<sup>۱۴</sup>Corpus

<sup>۱۵</sup>Dataset

<sup>۱۶</sup>Social network



در حوزه‌ی تشخیص و ردیابی موضوع به وجود آمد. یکی از زیرمجموعه‌های جدید در این حوزه، تشخیص روندهای درحال‌ظهور<sup>۱</sup> است که توجه بسیاری از محققان را به خود جلب کرده است. روند درحال‌ظهور، به موضوعی گفته می‌شود که در طول زمان به علاقه‌مندی و سودمندی آن افزوده می‌شود [۵].

## ۳-۱ کاربردها

تشخیص رویداد امکان‌فهم، استخراج و خلاصه‌سازی خودکار پیشامدهای هر رویداد در زمینه‌های مختلف از جمله بیولوژیکی، امنیت، سلامت و اقتصاد را فراهم می‌سازد. با به کارگیری روش‌های TDT، می‌توان به اطلاعات بسیاری دسترسی یافت. از جمله می‌توان به افرادی که تحت تاثیر رویداد قرار می‌گیرند، زمان و مکان رویداد، میزان خسارت ناشی از رویداد و تاثیر آن بر محیط اطراف اشاره کرد. از جمله دیگر کاربردهای تشخیص رویداد می‌توان به تشخیص زودهنگام شیوع بیماری و کشف جرایم اشاره کرد [۶]. از طرفی می‌توان در سامانه‌های استخراج اطلاعات و تصمیم‌گیری نیز از روش‌های تشخیص رویداد استفاده کرد تا کاربران به سرعت از رویدادهای جهان واقعی مطلع شده و در جهت بهره‌وری از شرایط به نفع خود اقدام کنند. همچنین خروجی این نوع از سامانه‌ها می‌تواند در سامانه‌های توصیه‌گر<sup>۲</sup> نیز استفاده شود.

## ۴-۱ چالش‌ها

روش‌های تشخیص رویداد با چالش‌هایی روبرو هستند. از جمله می‌توان به موارد زیر اشاره کرد.

### ۱-۴-۱ وابستگی حوزه

روش و تکنیکی که در یک حوزه جواب می‌دهد ممکن است برای حوزه‌ی دیگر نتایج خوبی نداشته باشد و غیرقابل استفاده باشد. بنابراین بسته به شرایط و حوزه‌ی انتخابی اسناد باید از روش‌ها و تکنیک‌های متفاوتی

<sup>1</sup>Emerging trend detection

<sup>2</sup>Recommender system

جهت تشخیص رویداد استفاده کرد [۶، ۷]. با پیدایش شبکه‌های اجتماعی نیز چالش‌های بیشتری در این حوزه ایجاد شده است. با افزایش کاربران و محتوای تولید شده در این شبکه‌ها روز به روز بر اهمیت استفاده از داده‌های آن‌ها افزوده می‌شود. بنابراین نیاز به روش‌ها و ابزارهای جدید جهت استخراج اطلاعات از این داده‌ها احساس می‌شود.

#### ۲-۴-۱ محدودیت زمانی

رویدادها باید به درستی در محدوده‌ی زمانی مشخصی تشخیص داده شوند تا بتوان از حادثه‌ی احتمالی جلوگیری کرد و یا عکس‌العمل مناسب نشان داد. بسته به حوزه‌ی انتخابی این زمان می‌تواند از ثانیه تا دقیقه تغییر کند [۶، ۷].

#### ۳-۴-۱ صحت تشخیص

حوزه‌هایی که نیاز به تصمیم‌گیری‌های حیاتی دارند، نیازمند صحت<sup>۱</sup> تشخیص بالایی هستند. در چنین حوزه‌هایی مانند پزشکی و بانکداری، تشخیص اشتباه منجر به خسارات مالی و یا جانی می‌شود. بنابراین انتظار می‌رود که روش‌های تشخیص رویداد از نرخ مثبت حقیقی<sup>۲</sup> (تشخیص درست) بالا و نرخ مثبت کاذب<sup>۳</sup> (تشخیص نادرست) پایین برخوردار باشند [۶].

#### ۴-۴-۱ منابع داده‌ای متنوع

مجموعه‌ای عظیم از انواع داده‌ها در شبکه‌های اجتماعی و پایگاه‌های خبری برخط وجود دارد که این داده‌ها شامل اسناد متنی، تصاویر، ویدیو، صوت، داده‌های رابطه‌ای، رکوردهای چند متغیری و داده‌های فضایی و زمانی هستند. بنابراین باید داده‌ی مناسب با روش مطالعه از منابع انتخاب شود [۶].

---

<sup>۱</sup>Accuracy

<sup>۲</sup>True-positive rate

<sup>۳</sup>False-positive rate

حجم عظیم داده‌ها نیازمند الگوریتم‌های محاسباتی با قدرت بالا و فضای ذخیره‌سازی کلان برای ذخیره، دسترسی، پالایش و پردازش داده‌ها در محدوده‌ی زمانی مشخص است [۶].

## ۱-۵ تعریف مساله و دستاوردها

وب جهانی یکی از غنی‌ترین منابع اطلاعاتی است و روز به روز بر تعداد کاربران و محتوای تولید شده توسط آن‌ها افزوده می‌شود. بررسی انبوهی از اطلاعاتی که هر روزه در اینترنت منتشر می‌شود در عمل کاری زمان‌بر و غیرممکن است. موتورهای جستجو یکی از بهترین گزینه‌ها برای بازیابی و پالایش اطلاعات هستند. با این حال، نیاز به سامانه‌هایی که بتوانند به صورت خودکار مهم‌ترین رویدادهای جهان واقعی را کشف کنند، احساس می‌شود. اطلاعاتی که از چنین سامانه‌هایی بدست می‌آیند می‌توانند در زمینه‌های مختلف تجاری، پزشکی، اجتماعی، سیاسی، پژوهشی و غیره کاربرد داشته باشند. بنابراین نیاز به روش‌های جدیدی برای بازیابی و استخراج اطلاعات احساس می‌شود. یکی از مطالعات مهم در این حوزه تشخیص و ردیابی رویداد است که به فرایند شناسایی رویدادهای جهان واقعی توسط تحلیل و رصد رسانه‌های جمعی به ویژه اینترنت گفته می‌شود. دقت و زمان تشخیص رویداد دو چالش جدی برای این مسئله به حساب می‌آیند.

در این پایان‌نامه با ارائه‌ی روشی مبتنی بر گراف و به‌کارگیری روش تشخیص جامعه<sup>۱</sup> که نوعی روش خوشه‌بندی در گراف محسوب می‌شود، رویدادها در وب‌سایت‌های خبری شناسایی و خوشه‌بندی می‌شوند. روش‌های پیشین از گراف بدون وزن جهت تشخیص رویداد استفاده کرده‌اند، درحالی‌که در این پژوهش از دو ویژگی جدید شامل تعداد کلیک‌های کاربران بر اسناد و عناوین خبری به همراه فراوانی عبارات سند جهت وزن‌دهی گراف و تشخیص رویداد استفاده شده است. همه‌ی این ویژگی‌ها در بازه‌ی انفجاری عبارت محاسبه شده‌اند، بدین ترتیب ویژگی زمان نیز در این روش در نظر گرفته شده است. الگوریتم‌های تشخیص جامعه عموماً زمان اجرای بالایی دارند، در این پژوهش با استفاده از روش‌های نمونه‌برداری از گراف، سازوکاری جهت بهبود سرعت اجرا به کار گرفته شده است. بدین ترتیب رویدادهای خبری در طول زمان دسته‌بندی می‌شوند و

<sup>1</sup>Community detection

جست‌وجو در آن‌ها برای کاربر تسهیل خواهد شد. نتایج ارزیابی نشان می‌دهد که روش پیشنهادی فراخوانی را نسبت به روش‌های پیشین ۱۹ درصد بهبود می‌بخشد و دقت را تا ۲ برابر افزایش می‌دهد.

## ۶-۱ ساختار کلی پایان‌نامه

ساختار کلی پایان‌نامه در ادامه بدین شرح می‌باشد. در فصل دوم، مفاهیم اولیه‌ی حوزه‌ی تشخیص و ردیابی موضوع مطرح می‌شود. در فصل سوم مروری بر پیشینه‌ی پژوهش در این حوزه خواهیم داشت. فصل چهارم به شرح روش پیشنهادی اختصاص دارد. در فصل پنجم، درباره‌ی جزییات پیاده‌سازی و ارزیابی روش پیشنهادی بحث می‌شود و فصل ششم شامل نتیجه‌گیری و ارائه‌ی پیشنهادهایی برای کارهای آینده است.