

شیوه نامه برچسب گذاری

محدودیت های مجموعه داده های فارسی

مجموعه داده مورد استفاده در روش های مشابهت سنجی معنایی به صورت جفت جمله برچسب دار می باشند که با عددی در بازه ی (0 و 1) یا با یکی از اعضای مجموعه {0,1,2,3,4,5} میزان مشابهت دو جمله یا سند بیان می شود. با توجه به دانش ما، تاکنون تمامی پروژه های شباهت سنجی معنایی فارسی بر روی مجموعه داده های سرقت ادبی و با هدف کشف موارد تطابق متن صورت گرفته است. این نوع مجموعه داده ها شامل متون رسمی و آکادمیک و اسناد با طول بلند می باشند. تاکنون پژوهشی معنامحور بر روی متون کوتاه عامیانه شبکه های اجتماعی مانند داده های توئیت در زبان فارسی انجام نشده است. در نتیجه مجموعه داده ای برای این امر در حال حاضر وجود ندارد. در این پژوهش برای تست و ارزیابی مدل پیشنهادی، نیاز به مجموعه داده ای برچسب گذاری شده داریم، که با توجه به عدم وجود چنین مجموعه داده ای، برآن شدیم تا مجموعه داده ای از داده های توئیت فارسی به منظور تستِ مدل شباهت سنج معنایی ایجاد و در آزمایشات خود استفاده کنیم.

همانطور که پیش تر ذکر شد، مجموعه داده ی وظایف شباهت سنجی معنایی به صورت جفت-سند هم معنی با برچسبی به منظور نمایش امتیاز مشابهت دو سند است. یافتن توئیت های هم معنی نیازمند شناسایی پارافریزها می باشد. بسیاری از کارهای قبلی در متون انگلیسی در مورد شناسایی پارافریز بر روی یک مجموعه داده خاص، مجموعه پارافریز مایکروسافت [52]، که از مقالات خبری مشتق شده است، توسعه و ارزیابی شده است. اما این داده ها با داده های توئیت بسیار متفاوت است. در توئیت اخبار بسیار کوتاه درج می شوند و به دلیل محدودیت در تعداد کاراکتر، بیان مطالب توسط واژگان خاص انجام می شود. با شناسایی این واژگان مشترک بین توئیت های ارسالی در یک موضوع و یک زمان خاص، می توان توئیت های پارافریز را شناسایی کرد.

در کارهای گذشته در زبان های غیر فارسی مانند انگلیسی، عربی و ترکی [53] [54] [55]، روال کار به این صورت است که بعد از شناسایی توئیت های پارافریز، توسط چند حاشیه نویس خبره جفت توئیت ها از لحاظ معنایی

مقایسه می‌شوند و امتیازی به عنوان برچسب دریافت می‌کنند. درنهایت با توجه به رای حاشیه‌نویس‌ها برچسب میانگین شناسایی شده و به عنوان امتیاز هر جفت توثیت در مجموعه داده درج گردیده است.

جمع‌آوری و برچسب‌زنی مجموعه داده تست

در این بخش، مراحل گردآوری و برچسب‌گذاری مجموعه داده جمع‌آوری شده با هدف مشابهت سنجی معنایی داده‌های شبکه‌های اجتماعی، به ترتیب شرح داده می‌شود. لازم به ذکر است، مجموعه داده مورد استفاده این پژوهش مبتنی بر داده‌های جمع‌آوری شده فارسی توییتر در بستر سامانه ذکاوت پژوهشگاه ارتباطات و فناوری اطلاعات (مرکز تحقیقات مخابرات ایران) بوده و همچنین جهت پردازش زبان فارسی نیز از ابزارهای توسعه داده شده در این سامانه استفاده خواهد شد. به صورت کل پس از جمع‌آوری موضوعی داده‌ها در تاریخ‌های مشخص، با کمک روش بخش 5-3-1 پارافریزهای موجود شناسایی و برای عملیات برچسب‌زنی ذخیره و پالایش می‌شوند. دیتاست جمع‌آوری شده در مجموع 1123 جفت-توثیت را شامل می‌شود که توسط 4 حاشیه‌نویس مجرب برچسب‌زنی شده‌اند. برچسب‌نهایی هر نمونه از دیتاست، با میانگین‌گیری از برچسب این 4 نفر بدست می‌آید.

شناسایی جفت داده‌ی کاندید

در ابتدا نیاز است به این موضوع پرداخت که دیتاست‌های مربوط به وظیفه‌ی مشابهت سنجی معنایی، به صورت جفت (pair) می‌باشند. جفت متن به این صورت که دو متن از لحاظ معنایی به یکدیگر مشابهت دارند (دو متن پارافریز هم هستند) و این مشابهت با مقدار عددی به عنوان درجه تشابه تحت عنوان برچسب در دیتاست ذخیره شده‌اند. پس هر نمونه در دیتاست از یک جفت متن با برچسبی که نشان‌دهنده میزان مشابهت دو متن است، نمایش داده می‌شود.

برای ایجاد جفت های کاندید، ابتدا تمام داده ها پیش پردازش شدند و مواردی که به مقایسه‌ی معنایی متون کمکی نمی‌کردند (برای مثال stopwords، هشتگ های غیرمرتبط و ...) را حذف کردیم. در مرحله بعد، تمام جفت متون به عنوان جفت‌های کاندید در نظر گرفته شدند، سپس با توجه به معیارهای زیر فیلتر شدند. یک جفت کاندید حذف می شود اگر تنها یکی از شروط زیر برقرار باشد:

1. در یک جریان (trend) از شبکه ی اجتماعی مورد نظر نباشند.

2. فاصله زمانی درج دو متن بیشتر از یک روز باشد.

3. دو متن در کمتر از 5 لغت همپوشانی داشته باشند.

امتیازدهی شباهت

روش حاشیه نویسی ما از روشی پیروی می کند که در اکثر دیتاست های زبان انگلیسی برای شباهت سنجی معنایی استفاده می شود. به جای اینکه جفت متون را صرفاً به عنوان پارافریز یا غیرپارافریز برچسب گذاری کنیم، یک امتیاز تشابه معنایی دقیق تر به آن اختصاص داده می شود. ما دستورالعمل های ارائه شده برای وظیفه ی شباهت سنجی معنایی را دنبال کردیم و حاشیه نویس ها جفت متون را مستقیماً، طبق معیارهای نشان داده شده در جدول زیر نمره گذاری کردند.

نمرات تشابه معنایی برای جفت پارافریزهای کاندید

درجه شباهت معنایی	راهنما
5 - یکسان	کاملاً معادل هستند، زیرا آنها به یک معنا هستند.
4- نزدیک	به مقدار زیادی معادل هستند، اما برخی از جزئیات بی اهمیت متفاوت هستند.
3- مرتبط	تقریباً معادل هستند، اما برخی از اطلاعات مهم متفاوت است/از دست رفته است.

2- در یک زمینه	معادل نیستند، اما برخی جزئیات را به اشتراک می‌گذارند.
1- تا حدودی مرتبط	معادل نیستند، اما در یک موضوع هستند.
0- غیر مرتبط	در موضوعات مختلف هستند.

چهار گویشور بومی فارسی زبان، به عنوان حاشیه نویس استخدام شدند. ما مجموعه ای از دستورالعمل های حاشیه نویسی را برای توضیح روند امتیازدهی آماده کردیم. این دستورالعمل، نمرات شباهت را به همراه یک جفت مثال برای هر نمره و توضیح کوتاه معرفی کرده است. نمونه‌ها نیز از میان نمونه‌های متعدد انتخاب شدند.

به منظور استفاده بهتر و آزمایش‌های دقیق تر در آینده، نمرات اختصاص داده شده نیز به برچسب های باینری تبدیل شدند. ابتدا، نمرات هر حاشیه نویس با در نظر گرفتن جفت متون که به صورت یکسان (برچسب 5)، نزدیک (برچسب 4) و مرتبط (برچسب 3) به عنوان مثبت (یعنی پارافریز) و آنهایی که در یک زمینه (برچسب 2)، مرتبط (برچسب 1) و غیرمرتبط (برچسب 0) علامتگذاری شده بودند، برچسب منفی هستند (یعنی غیر پارافریز). تعداد تصمیمات مثبت و منفی برای هر نمونه ممکن است به صورت یک جفت خلاصه شود. به عنوان مثال، (1,3) نشان می‌دهد که تنها یک حاشیه نویس این جفت را به عنوان یک پارافریز برچسب گذاری کرده است، در حالی که سه نفر باقی مانده آن را به عنوان غیرپارافریز برچسب گذاری کرده اند. در جدول 5، معیارهای قضاوت باینری را نشان می‌دهیم که بر اساس تعداد پاسخ های حاشیه نویسان است. توجه داشته باشید که در مواردی که تعداد تصمیمات مثبت و منفی برابر باشد، یک جفت را "قابل بحث" (debatable) در نظر می‌گیریم. این رویکرد مشابه روش برچسب گذاری دیتاست TPC است که بر اساس توافق بین حاشیه نویسان نیز می‌باشد.

برای برچسب نهایی شباهت معنایی، ما میانگین نمرات را ارائه می‌دهیم. این محدوده بین 1.75 و 3.00 برای جفت های "قابل بحث" است، در حالی که جفت های پارافریز مثبت بالاتر از 3.00 و جفت های منفی کمتر از 1.75 امتیاز دارند.

معیارهای قضاوت باینری بر اساس تعداد پاسخ های حاشیه نویسی ها

تعداد پاسخ ها	قضاوت
(4,0); (3,1)	مثبت (1)
(0,4); (1,3)	منفی (0)
(2,2)	قابل بحث

برای مثال جدول شماره 6، سه جفت نمونه از داده ها را نشان می دهد. هر جفت با نمرات 4 حاشیه نویسی مختلف و میانگین نمرات شباهت نشان داده شده است. جفت "قابل بحث" توسط چهار حاشیه نویسی (4,2,3,0) نمره گذاری شده است و میانگین شباهت 2.25 است.

جفت های نمونه فرضی از دیتاست

جفت متن	امتیازات	پارافریز
تاجر سرشناس در مورد مافیای خودرو گفت: "اگر برخی چیزها را توضیح دهم، آنها نمی توانند در ایران بمانند" ن.ح. تاجر ایرانی در گفتگو با انتخاب 24 گفت: "مافیای خودرو بگذراند سکوت کنم، وگرنه جایی در این کشور ندارند"	(4,2,3,0) <i>Average = 2.25</i>	قابل بحث
مشخص شد که 13 خودرو از 15 خودرو که در قرعه کشی به یک فرد رسید، در نمایندگی های منطقه هشت سپرده شده است.	(3,4,5,5) <i>Average = 4.25</i>	مثبت

		13 شانس از 15 خودرو که به یک شخص رسید، در منطقه هشت است که از حساب دلالتان منطقه واریز شده است.
منفی	$(1,3,0,0)$ $Average = 1$	استان فارس با 52.9 درصد بیشترین سهم فروش وام مسکن در کل فروش مسکن را داشت. سهم اولین فروش از کل فروش خانه 45.4 درصد بوده است.

چند نمونه از برچسب دهی داده ها

برای درک بهتر و آشنایی بیشتر با نحوه ی لیبل دهی توسط حاشیه نویس ها، این بخش به ذکر چند مثال از دیتاست و توضیح مختصری در مورد آن می پردازد.

نمونه 1:

توئیت یک	توئیت دو	امتیاز شباهت
نیکلاس #مادرو رئیس جمهوری #ونزوئلا وارد #تهران شد.	نیکلاس مادرو، رئیس جمهور ونزوئلا وارد تهران شد	5

همانطور که در نمونه شماره یک دیده می شود، دو توئیت کاملاً از لحاظ معنایی و محتوایی یکسان می باشند و اندک تفاوتی از لحاظ لغوی در آنها مشاهده می شود. پس در نتیجه دو توئیت پارافریز یکدیگر بوده و با توجه به جدول شماره 1، با امتیاز شباهت 5 می توان آن ها را برچسب دهی کرد.

نمونه 2:

توئیت یک	توئیت دو	امتیاز شباهت
----------	----------	--------------

5	<p>#برانکو هم به درخواست #تیم_ملی ایران جواب منفی داد با وجود توافق فدراسیون های فوتبال #ایران و عمان و نامه نگاری هایی که برای برگزاری دیدار دوستانه انجام شده بود، اما این دیدار به دلیل مخالفت برانکو ایوانکوویچ منتفی شد.</p>	<p>دیدار دوستانه ایران - عمان به خاطر مخالفت برانکو منتفی شد با وجود توافق فدراسیون های فوتبال ایران و عمان و نامه نگاری هایی که برای برگزاری دیدار دوستانه انجام شده بود، اما این دیدار به دلیل مخالفت برانکو ایوانکوویچ منتفی شد</p>
---	---	--

از آنجایی که در نمونه شماره دو نسبت به شماره یک تفاوت های لغوی بیشتری دیده می شود، اما با بررسی معنایی دو توثیت موجود در این نمونه، مشاهده می شود که محتوای دو متن کاملاً یکسان است و تمامی معانی موجود در توثیت یک، در توثیت شماره دو نیز موجود است، از این رو امتیاز مشابهت 5 برای این جفت در نظر گرفته می شود.

نمونه 3:

توثیت یک	توثیت دو	امتیاز مشابهت
<p>در ایران خبرگزاری صدا و سیما می گوید نایب رئیس فدراسیون فوتبال به اتهام کلاهبرداری بازداشت شده. گفته شده این دستگیری با یک پرونده کلاهبرداری با استفاده از نوعی رمز ارز مرتبط است. گفت و گو با فرید اشرفیان، خبرنگار ورزشی:</p>	<p>نایب رئیس فدراسیون فوتبال بازداشت شد نایب رئیس فدراسیون فوتبال به دلیل ارتباط با پرونده #کلاهبرداری رمزارز جعلی کینگمانی بازداشت شد</p>	4

--	--	--

در نمونه شماره 3، دو توثیت از لحاظ معنا و محتوا به یک موضوع اشاره دارند، اما در توثیت یک به مواردی کم اهمیت دیگری پرداخته شده که به بیان اصل موضوع ارتباط چندانی ندارد. به همین دلیل امتیاز مشابهت 4 برای این جفت لحاظ می‌شود.

نمونه 4:

توثیت یک	توثیت دو	امتیاز مشابهت
مقام جهاد کشاورزی: ظرف یک هفته قیمت برنج را به نرخ مصوب برمی‌گردانیم	با افزایش قیمت #برنج ایرانی به کیلویی ۱۰۰ هزار تومان، جهاد کشاورزی از اجرای طرح تشدید مبارزه با احتکار برنج در ایران خبر داد و اعلام کرد: بر اساس این طرح، قیمت برنج ظرف یک هفته به نرخ مصوب باز خواهد گشت. قیمت‌های دستوری در چهار دهه گذشته جزیی از سیاست اقتصادی در #ایران بوده است	3

در نمونه 4، موضوع مورد بحث در هر دو توثیت به یک چیز اشاره می‌کند، اما برخی از اطلاعات مهم در توثیت یک وجود ندارد و یا اینکه متفاوت از توثیت دو است. به همین دلیل امتیاز مشابهت 3 در نظر گرفته می‌شود.

نمونه 5:

توثیت یک	توثیت دو	امتیاز مشابهت
----------	----------	---------------

2	<p>عکس روز - اولین گل #سردار_آزمون در</p> <p>#بوندسلیگا. آزمون که از ابتدا در ترکیب تیم</p> <p>لورکوزن حضور داشت، موفق شد در دقیقه ۱۸</p> <p>بازی مقابل گروتفورث گل دوم تیمش را به ثمر</p> <p>برساند.</p>	<p>سردار آزمون بهترین شد</p> <p>#سردار_آزمون با گرفتن نمره بسیار</p> <p>خوب 8.4 از وبسایت آماری</p> <p>هواسکورد بهترین بازیکن لورکوزن</p> <p>در بازی با گروتفورث شد</p>
---	---	---

در نمونه شماره 5، از لحاظ معنایی و محتوایی دو توثیت به دو موضوع متفاوت اشاره می‌کنند، اما در برخی جزئیات مشترک هستند، در نتیجه مرتبط نبوده اما می‌توانیم با استفاده از جدول شماره 1، امتیاز مشابهت 2 را به آنها اختصاص دهیم.

نمونه 6:

توثیت یک	توثیت دو	امتیاز مشابهت
<p>پیکر بی‌جانی که خانواده</p> <p>#عبدالباقی دارن بالای سرش اون</p> <p>نمایش مضحک و بی‌معنی رو بازی</p> <p>میکنن، پیکر عزیز یک خانواده است</p> <p>که احتمالاً در نزدیکی آوار #متروپل</p> <p>بی‌قرار و دل‌آشوب منتظر خبری</p>	<p>تست DNA و دم خروس #عبدالباقی ۱-</p> <p>تست DNA نیاز به یک تا سه هفته زمان برای</p> <p>کپی‌برداری و مقایسه دارد و در بهترین</p> <p>آزمایشگاههای دنیا حداقل ۳ روز زمان نیاز دارد و</p> <p>اعلام نتیجه ۲۸ ساعت بعد از پیداشون جسد</p> <p>تامل برانگیز است ۲-تست DNA مقایسه دو</p>	1

هستند؛ حقیقتاً قساوت و پستی این جماعت انتها نداره	DNA با هم است. روی DNA آدم‌ها اسمشون نوشته نشده	
--	--	--

دو توئیت نمونه شماره 6، از لحاظ معنایی پارافریز هم نبوده و ارتباطی با یکدیگر ندارند، اما در یک موضوع مشترک قرار دارند و به همین علت امتیاز مشابهت 1 برای آن‌ها لحاظ شده است.

نمونه 7:

توئیت یک	توئیت دو	امتیاز مشابهت
چرا و چرا و چرا #عبدالباقی در ساختمانی که مشکل دارد دفتر میسازد! شاید توهم داشته مقاوم است. یعنی اینقدر به خودش دروغ گفته که باورش شده.	«برادران لیلا» برنده جایزه «فیپرسی» #جشنواره_کن شد هیات داوران فدراسیون بین‌المللی منتقدان (فیپرسی) در جشنواره فیلم کن، جایزه بهترین فیلم بخش رقابتی اصلی را به فیلم «#برادران_لیلا» ساخته #سعید_روستایی اعطا کرد. روستایی جایزه‌اش را به مردم داغ‌دیده #آبادان تقدیم کرد_انتخاب	0

نمونه شماره 7، حاوی دو توثیت از دو موضوع متفاوت و کاملاً جدا از هم است که باید امتیاز مشابهت 0 را برای آن‌ها در نظر گرفت.

	tweet 1	tweet 2	نفر اول	نفر دوم	نفر سوم	نفر چهارم	average
1	ترجمه: یکی از هدیه‌هایی که ترامپ می‌تواند به کشور و جهان بدهد این است که در ۱۵ روز آخر ریاست جمهوری خود جنگی با ایران آغاز نکند. دشوار هست بگویم این کار عمداً از طرف دولت صورت می‌گیرد یا اینکه به طور ساده نتیجه بی‌کفایتی آنهاست. در هر صورت، آنها در حال انجام یک بازی خطرناک هستند. https://t.co/sDjySiyMSW	@Moghadam۹۶۳۱ (مشخص نیست شاید این صوت از طرف خودشون باشه) همدار هیلاری کلینتون: ترامپ، با ایران جنگ نکن! این یک بازی خطرناک است ♦ هیلاری کلینتون در پیام توئیتری نوشته: هدیه‌ای که ترامپ می‌تواند به این کشور و جهان بدهد این است که در ۱۵ روز باقی‌مانده از دولتش جنگی با ایران آغاز نکند. https://t.co/GhdvdyvYeb	4	3	5	4	4
2	#حسن_یزدانی استوری گذاشته و قبل از عذرخواهی بابت بخت دیشب به خانواده #مهساامینی تسلیت گفته. فهرمان و پهلوان قلب ما خودتی حسن آقا یزدانی	به نظر من حسن یزدانی یک قهرمان عالی و تکرار نشدنی هست ولی آقا حسن پهلوان نیست پهلوان باید پشت مردمش باشه لطفاً برداشت بد نکنید 🙏 #بختی https://t.co/rTHPtXrgg #حسن_یزدانی	2	4	4	3	3.25
3	«خدا روشکر با قتل پسر من به آرامش رسیدم/دیگر هیچ ناراحتی در زندگی ندارم/میدانستم پسر من در خانه ما رابطه‌ی جنسی برقرار میکند» اگه زیاد دلشو ندارید این رشتو رو نخوئید. بخشی از اعترافات پدر #بابک_خرمدین:	پدر #بابک_خرمدین: پسر دختران زیادی را به خانه می‌آورد و در اتاق با آنها تنها بود. ما فکر میکردیم با آنها رابطه جنسی برقرار میکنند. از رفتارهایش خسته شدیم و او را کشتیم. *منهم در هنگام خروج از اتاق با پسر دستهایش را رو به آسمان کرد و خدا را بابت	4	4	4	4	4

سه نمونه موجود از دیتاست گردآوری شده این پژوهش

شکل زیر سه نمونه از دیتاست گردآوری شده را نشان میدهد.

برای دریافت اطلاعات بیشتر از پروژه با ایمیل zekavat@itrc.ac.ir تماس بگیرید.