

خوشه بندی جملات فارسی مبتنی بر الگوریتم های هوش جمعی

مهدی بازقندی^۱، قمرناز تدین تبریزی^۲ و مجید وفایی جهان^۳

Mehdi_Bazghandi@ymail.com^۱

Tadayon@mshdiau.ac.ir^۲

VafaeiJahan@mshdiau.ac.ir^۳

چکیده

خوشه بندی یکی از مسائل مهمی است که امروزه بسیاری از محققین در زمینه های مختلف به آن پرداخته اند. تا کنون الگوریتم های کلاسیک زیادی در این زمینه ارائه شده است. که اغلب این روش ها دارای ناپایداری بوده و همچنین پارامترهای آن ها محدود به انتخاب کاربر می باشد. از کاربردهای خوشه بندی می توان به خوشه بندی متون و اسناد در موضوعات خلاصه سازی متون و بازیابی اطلاعات یاد کرد. در خوشه بندی جملات یک متن برای مشخص شدن جملات مشابه و نمی توان از روش مشابه آن (دسته بندی متون مشابه) استفاده کرد. بردارهایی به طول m و با مقادیر صفر بسیار زیاد پدید خواهد آمد. برای حل این مشکل، روشی جدید مبتنی بر PSO برای خوشه بندی جملات یک متن معرفی شده است. به طوریکه به جای استفاده از فاصله اقلیدسی و فاصله کسینوسی، از یک معیار جدید در محاسبه فاصله دو جمله استفاده شده است. معیاری که در آن، ارتباط معنایی کلمات با استفاده از ارتباطات آنها در متن در نظر گرفته می شود. همچنین تعیین تعداد خوشه های بهینه یکی دیگر از کارهای انجام شده در این مقاله است. برای ارزیابی یک مجموعه از خبرهای ورزشی فارسی انتخاب شده است. نتایج حاصل از ارزیابی روش پیشنهادی نشان می دهند که استفاده از خوشه بندی PSO معنایی، با تعیین تعداد خوشه های مطلوب، دقت بهتری را در خوشه بندی جملات در مقایسه با روش های دیگر، دارد.

کلمات کلیدی

خوشه بندی، الگوریتم PSO، بردارهای Context-Vector، شباهت معنایی

۱- مقدمه

یا در خوشه بندی سبد خرید مشتریان، فاصله بر اساس شباهت خرید تعیین می شود. لذا محاسبه فاصله بین دو داده در خوشه بندی بسیار مهم می باشد؛ زیرا کیفیت نتایج نهایی را دستخوش تغییر قرار خواهد داد. فاصله که همان معرف عدم تجانس است حرکت در فضای داده ها را می سازد و سبب ایجاد خوشه ها می گردد. با محاسبه فاصله بین دو داده می توان فهمید که چقدر این دو داده به هم نزدیک هستند و براین اساس در یک خوشه قرار داده می شود. توابع ریاضی مختلفی، برای محاسبه فاصله وجود دارند؛ فاصله اقلیدسی، فاصله همینگ و... با افزایش حجم منابع متنی موجود در وب چالشی که وجود داشته است آن است که چگونه کاربران می توانند به اطلاعات مورد نیاز خود در اسرع وقت و با دقت بالا دسترسی داشته باشند. خوشه بندی می تواند راه حلی برای حل این مسئله باشد. خوشه بندی می تواند نقش مهمی در انتخاب عنصرهای شایسته داشته باشد. با خوشه بندی اطلاعات می توان عنصرهای مشابه را در یک خوشه قرار داد و عنصر محوری را به عنوان نماینده خوشه انتخاب کرد. خوشه بندی در بسیاری از کاربردهای پردازش متن مانند خلاصه سازی، درک متن، ترجمه ماشینی و... استفاده می شود

پردازش داده، یکی از شاخص های بسیار مهم در دنیای اطلاعات است. خوشه بندی یکی از بهترین روش هایی است که برای کار با داده ها ارائه شده است. خوشه بندی قابلیت ورود به فضای داده و تشخیص ساختارش را امکان پذیر می نماید. لذا به عنوان یکی از ایده آل ترین مکانیزم ها، برای کار با دنیای عظیم داده ها محسوب می شود.

خوشه بندی، یافتن ساختاری در مجموعه ای از داده ها است که طبقه بندی نشده اند. به بیان دیگر می توان گفت که خوشه بندی قراردادن داده ها در گروه هایی است که اعضای هر گروه از زاویه خاصی شبیه یکدیگرند. در نتیجه شباهت بین داده های درون هر خوشه حداکثر و شباهت بین داده های درون خوشه های متفاوت حداقل می باشد. معیار شباهت در اینجا، فاصله بوده یعنی نمونه هایی که به یکدیگر نزدیکترند در یک خوشه قرار می گیرند. به عنوان نمونه در خوشه بندی اسناد دوری و یا نزدیکی داده ها متناسب با تعداد کلمه های مشترکی که در دو سند وجود دارد و



استخراج می کند. با توجه به اینکه در زبان فارسی با کمبود این ابزارها و منابع رو به رو هستیم. در این مقاله سعی شده است اسامی موجود در پیکره از قبل به صورت دستی در یک فایل ذخیره شوند. هدف این بخش آن است که ارتباط و میزان شباهت تمامی اسامی استخراج شده، تعیین شود و در ماتریس شباهت ذخیره شود. ما برای رسیدن به این هدف بردارهای Context Vector را معرفی می کنیم. برای هر کدام از واژه های اسمی موجود در متن یک بردار N بعدی تعیین می کنیم که شامل N کلمه با بیشترین ارتباط با واژه مورد نظر است. برای ای منظور ما باید بیشترین رخدادهای همزمان این واژه با سایر واژه ها را پیدا کنیم. احتمال شرطی متقارن (Symmetric Conditional Probability) روشی را معرفی می کند که می توان میزان ارتباط و همزمانی واژه ها را بدست آورد [۶] که معادله آن به صورت زیر است:

$$SCP(w_1, w_2) = \frac{p(w_1, w_2)}{p(w_1)p(w_2)} \quad (1)$$

که $p(w_1, w_2)$ احتمال رخداد همزمان دو واژه در یک همسایگی معین است. $p(w_1)$ و $p(w_2)$ به ترتیب احتمال رخداد واژه w_1 و w_2 به تنهایی را بیان می کنند. در پیاده سازی یک پنجره متنی به اندازه ۲۰ واژه لحاظ شده است. طبق اندازه پنجره، می توان واژه هایی که در همسایگی "باشگاه استقلال" آمده را مشاهده کرد. (جدول (۱))

جدول (۱). رخداد همزمان واژه ها با "باشگاه استقلال"

مدیر برنامه های مجتبی جباری گفت: جباری امروز قصد آشتی با مسوولان باشگاه استقلال را داشت اما، صحبت های کادرفنی این تیم همه چیز را به هم زد. خشایار محسنی در گفت و گو با خبرنگار ورزشی خبرگزاری دانشجویان ایران (ایسنا) گفت: زمانی که روز پنجشنبه به باشگاه استقلال رفتیم، برخورد بدی با جباری صورت گرفت. البته پس از این جلسه، مسوولان باشگاه استقلال از اقدامشان اظهار پشیمانی کردند و من نیز با جباری صحبت کردم. به همین دلیل این بازیکن قصد داشت امروز برای آشتی به باشگاه استقلال برود اما، صبح با من تماس گرفت و با اشاره به مصاحبه های حجازی و حاجیلو در مطبوعات، از تصمیم خود منصرف شد.

که در خوشه بندی آن ها معیار شباهت استفاده شده، نقش مهمی در کیفیت خوشه ها خواهد داشت. در این بخش به چند مورد از کارهای انجام گرفته در این زمینه اشاره می کنیم. یکی از عمومی ترین الگوریتم های خوشه بندی، الگوریتم K-Means می باشد این الگوریتم با داشتن ساده گی و خطی بودن، کارایی مناسبی در برخورد با داده های حجیم دارد [1]. در الگوریتم K-Means جواب نهایی وابستگی زیادی به انتخاب های اولیه مانند تعداد کلاسترها و مراکز نقاط دارند. و متناظر با این انتخاب ها راه حل های محلی خواهیم داشت [2]. در مقاله [3] از الگوریتم PSO برای خوشه بندی داده ها استفاده می کند این الگوریتم مشکل پاسخ های بهینه الگوریتم K-Means را برطرف می کند و با جستجوهای بیشتر پاسخ گلوبال را بدست خواهد آورد. در همین مقاله ترکیب الگوریتم های K-Means و PSO در نظر گرفته شده است و جواب بهتری نسبت به هر کدام از الگوریتم ها دارد. در [4] از الگوریتم PSO برای خوشه بندی اسناد استفاده کرده است اما معیار استفاده شده فاصله اقلیدسی است اما این معیار، فقط ظاهر کلمات موجود در متن را در نظر می گیرد و به ارتباط معنایی آن ها توجهی نمی کند. در Mead که یک سیستم خلاصه ساز مبتنی بر مرکز می باشد، با استفاده از خوشه بندی مبتنی بر مرکز، جملات مشابه یک متن را در یک خوشه قرار می دهد. [۵]. در زبان فارسی در مقاله [۹] با استفاده از خوشه بندی سلسله مراتبی و K-Means جملات مشابه را در یک خوشه قرار می دهد. در این مقاله سعی شده است ابتدا ارتباط معنایی کلمات با استفاده از روش های اماری تعیین شود. و سپس با استفاده از الگوریتم خوشه بندی PSO و بهبود یافته آن، جملات یک متن را خوشه بندی کرد.

۲- روش پیشنهادی

در این مقاله ابتدا با استفاده از بردارهای Context-Vector ارتباط معنایی کلمات تشخیص داده می شود. سپس الگوریتم PSO تشریح می شود. در ادامه آن الگوریتم خوشه بندی PSO پیشنهادی بیان می شود. و در انتها روش های پیشنهادی با یکدیگر مقایسه می شوند.

۲-۱- ارتباط معنایی واژه ها

در ابتدا نیاز داریم پیکره خود را مورد پردازش قرار دهیم و اسامی موجود در آن را به عنوان کلمات کاندید استخراج کنیم. در زبان انگلیسی می توان به ابزاری مانند TNT-Tager و SENTA اشاره کرد که اسامی یک جزئی یا دو جزئی را

			۰,۵۲	
رونالدو	الکلاسیکو	بارسلونا	میلان	قرارد
۰,۱۹	۰,۸۵	۰,۷۲	۰,۶۳	
بارسلونا	نیوکمپ	گواردیولا	اسپانیا	رتال مادرید
۰,۴۹	۰,۶۲	۰,۸۱	۰,۵۶	
علی دایی	سایپا	تیم ملی	اسیا	سرمربی
۰,۶۲	۰,۵۴	۰,۵۱	۰,۴۲	

با استفاده از رابطه شباهت (۳) (تشابه برداری) میزان شباهت دو واژه را در یک ماتریس ذخیره می کنیم این رابطه نشان می دهد که دو واژه تا چه اندازه از لحاظ معنایی به همدیگر مرتبط هستند. و چنانچه این اندازه از یک میزان استانه بیشتر باشد نشان می دهد که این دو واژه به یکدیگر مشابه هستند. و ماتریس شباهت را می سازد. این ماتریس برای یکبار ساخته می شود و ذخیره خواهد شد. جدول (۳).

جدول (۳): ماتریس شباهت

	علی دایی	اقای گل جهان	سید مهدی رحمتی	دروازه بان استقلال
علی دایی	۱	۰,۶۵	۰,۳۲	۰,۲۱
اقای گل جهان	۰,۶۵	۱	۰,۳۴	۰,۱۴
سید مهدی رحمتی	۰,۳۲	۰,۶۵	۱	۰,۷۱
دروازه بان استقلال	۰,۲۱	۰,۱۴	۰,۷۱	۱

اکنون با استفاده از ماتریس شباهت واژه ها، شباهت میان دو جمله به صورت زیر حاصل می شود [6].

$$Sim(s_i, s_j) = \sum_{k=1}^n \sum_{i=1}^m \inf osimBA(x_i, x_j) \quad (4)$$

به دو جمله زیر دقت کنیم:

(۱) علی دایی دروازه سید مهدی رحمتی را گشود.

(۲) برترین گلزن جهان، دروازه بان استقلال را مغلوب کرد.

بعد از استخراج بردارهای Context-Vector مربوط به اسامی می توان میزان مشابهت دو کلمه را در فضای برداری مقایسه کرد. و چنانچه دو بردار در فضای برداری بهم نزدیک باشند نشان می دهد که دو واژه به یکدیگر نزدیک می باشند و در یک متن مشابه بکار گرفته شده اند. و ارتباط آن ها تقریباً مشابهت است. در صورتی واژه هایی که در بردار Context-Vector یک واژه می آیند روابط متنوعی با آن خواهند داشت. اکنون برای بدست آوردن واژه هایی که بردارهای آن ها بهم شبیه است به صورت زیر عمل می کنیم. در این مرحله میزان شباهت واژه های اسمی را می بایم فرض کنیم که بردار زیر موجود است.

$$x_i = (x_{i,1}, x_{i,2}, \dots, x_{i,p}) \quad (2)$$

این بردار مربوط به واژه اسمی X می باشد. که میزان شباهت اطلاعاتی دو واژه $S_{i,j} = f(x_i, x_j)$ بصورت زیر محاسبه می شود [6].

(3)

$$\inf osimBA(x_i, x_j) = \frac{A_{i,j}}{B_i * B_j + A_{i,j}}$$

$$A_{i,j} = \sum_{k=1}^p \sum_{l=1}^p X_{i,k} * X_{j,l} * SCP(w_{i,k}, w_{j,l})$$

$$\forall i, B_i = \sqrt{\sum_{k=1}^p \sum_{l=1}^p X_{i,k} * X_{j,l} * SCP(w_{i,k}, w_{j,l})}$$

که در این رابطه $X_{i,j}$ مربوط به وزن واژه $W_{i,j}$ است. [4] برای ۱۰۰ واژه اسمی مهمتر میزان تشابه آنها را به کمک تابع $f(x_i, x_j)$ حساب می کنیم. و در یک ماتریس ذخیره می کنیم. به جدول (۲) نگاه کنید.

جدول (۲): Context Word

لیگ برتر	مدیرعامل سایپا	علی دایی	سرمربی	سایپا
	۰,۳۸	۰,۳۹	۰,۵۰	۰,۷۵
جباری	فتح اله زاده	تیم ملی	مدیرعامل	استقلال
		۰,۲۳	۰,۲۸	۰,۳۹

این دو جمله از لحاظ ظاهری هیچ شباهتی به یکدیگر ندارند. فاصله کسینوسی میزان شباهت آن ها را صفر فرض می کند. در صورتی که دو جمله بالا از لحاظ معنایی یکسان است. و رابطه (۲) میزان شباهت آن ها را عددی مخالف صفر فرض می کند. و نشان می دهد این دو جمله با یکدیگر ارتباط دارند.

به ماتریس جمله- کلمه دو جمله بالادقت کنید:

جدول (۴) ماتریس جمله- کلمه

	دروازه بان استقلال	برترین گلزن جهان	سید مهدی رحمتی	علی دایی
S1	0	0	1	1
S2	1	1	0	0

طبق رابطه کسینوسی:

$$\text{Cosine}(s_1, s_2) = \frac{\sum_{i=1}^n S_i * S_j}{\sqrt{(S_i)^2} * \sqrt{(S_j)^2}} \quad (5)$$

میزان تشابه دو جمله S_1 و S_2 طبق رابطه بالا صفر می شود. در صورتی که طبق رابطه (۴) میزان تشابه آن ها برابر ۰,۵۲ بدست خواهد آمد. بنابراین با این رابطه میزان تشابه دو جمله را حساب می کنیم و در یک ماتریس ذخیره می کنیم. ماتریسی شبیه زیر حاصل خواهد شد. جدول (۵)

خلاصه این روش به این صورت است: برای یک متن ورودی ابتدا برای هر جمله اسامی موجود در تکتک جملات استخراج می شوند. وزن هر کدام از آن ها بدست می آوریم. (معیار TF-ISF). سپس با استفاده از ماتریس شباهت و رابطه (۴) می توان میزان شباهت جملات را محاسبه کرد. و ماتریس شباهت جملات را بدست آورد.

جدول (۵): ماتریس شباهت جملات

	S_1	S_2	S_3
S_1	۱	۰,۶۵	۰,۲۳
S_2	۰,۶۵	۱	۰,۳۴
S_3	۰,۲۳	۰,۳۴	۱

۲-۲- الگوریتم PSO

الگوریتم اجتماع ذرات يك الگوریتم بهینه سازی تقلیدی از رفتارهای جوامع جانوری در پردازش دانش جامعه است. این الگوریتم از دو زمینه ریشه گرفته است. نخست زندگی مصنوعی (مانند دسته پرندگان، ماهی ها) و دوم محاسبات تکاملی [8]. مبنای توسعه الگوریتم PSO این است که جواب های ممکن در يك مسأله بهینه سازی به صورت پرندگان بدون حجم و خصوصیات کیفی در نظر گرفته می شوند که از آنان به عنوان ذرات یاد می شود، این پرندگان در يك فضای n بعدی پرواز کرده و مسیر حرکت خود در فضای جستجو را بر اساس تجارب گذشته خود و همسایگانشان تغییر می دهند. در دسته ای متشکل از n جز، موقعیت جزء iam تحت اثر يك بردار مکانی n بعدی مطابق معادله (۶) قرار دارد. این جزء همچنین دارای يك بردار سرعت به صورت معادله (۷) می باشد. بهترین موقعیت قبلی بدست آمده برای جز iam، با استفاده از معادله (۸) نمایش داده می شود.

$$x_i = (x_{i,1}, x_{i,2}, \dots, x_{i,n}) \quad (6)$$

$$v_i = (v_{i,1}, v_{i,2}, \dots, v_{i,1}) \quad (7)$$

$$p_i = (p_{i,1}, p_{i,2}, \dots, p_{i,n}) \quad (8)$$

در نهایت موقعیت جدید اجزای دسته با استفاده از معادلات زیر به دست می آید.

$$v_i(t+1) = \Omega v_i(t) + c_1 r_1 (p_1(t) - x_i(t)) + c_2 r_2 (p_g(t) - x_i(t))$$

$$x_i(t+1) = x_i(t) + v_i(t+1)$$

(9)

۲-۳- الگوریتم PSO باینری

هر ذره در این الگوریتم یک مقدار باینری دارد. و مقدار آن با یک تابع برازندگی سنجیده می شود. اما به جای استفاده از معادله سرعت، موقعیت جدید ذره به صورت زیر حساب می شود.

$$P_{ij}(t) \geq \frac{1}{1 + \exp(-v_{i,j}(t))} \rightarrow X_{ij}(t+1) = 0 \quad (10)$$

$$\text{else} \rightarrow X_{ij}(t+1) = 1$$

۲-۴- خوشه بندی معنایی PSO

۲-۴-۱- روش خوشه بندی پیشنهادی

الگوریتم های خوشه بندی خودکار، مجموعه ای از اشیاء را در تعدادی از خوشه قرار می دهند. به طوری که بهترین عدد برای تعداد خوشه ها توسط الگوریتم خوشه بندی مشخص می شود. در نهایت اشیاء موجود در هر خوشه باید بیشترین شباهت ممکن به یکدیگر را داشته باشند. در عین حال که با اشیاء دیگر خوشه ها فاصله داشته باشند. در روش خوشه بندی پیشنهادی، مجموعه ای از جملات $S = [s_1, s_2, s_3, \dots, s_n]$ وجود دارند به طوری که باید در مجموعه ای از خوشه ها $C = \{C_1, C_2, \dots, C_K\}$ بدون هیچ گونه همپوشانی، قرار گیرند. تعداد خوشه های بهینه (K) مشخص نمی باشد. این روش خوشه بندی دارای شرایط زیر می باشد:

(۱) دو کلاستر مختلف، جمله مشترکی ندارند.

$$\forall p \neq q: p, q \in \{C_1, C_2, \dots, C_K\}, C_p \cap C_q = \emptyset$$

(۲) هر جمله باید حتما در یک کلاستر قرار گیرد.

$$\cup C_p = D$$

(۳) هر کلاستر باید حداقل یک جمله را داشته باشد.

$$C_p \neq \emptyset \quad \forall p \in \{1, 2, 3, \dots, k\}$$

در الگوریتم PSO برای اینکه بتوانیم از روش خوشه بندی پیشنهادی استفاده کنیم، باید ساختار هر ذره طوری تعریف شود که سه شرط بالا را ارضا کند. در این ساختار تعریفی، طول هر ذره برابر تعداد جمله های متن تعریف می شود. به طوری که برای هر جمله باید عددی از مجموعه اعداد $\{1, 2, 3, \dots, K\}$ انتخاب شود. (شرط دوم). در این ساختار باید از تمام K عدد موجود استفاده کرد. (شرط سوم). هر جمله فقط یک عدد منحصر به فرد دارد. (شرط اول). این عدد در طول ذره منحصر به فرد نمی باشد. بنابراین به طور خلاصه، اعداد موجود در هر ذره (با توجه به تعداد خوشه ها) باید طور انتخاب شوند که از همه اعداد $\{1, 2, 3, \dots, K\}$ استفاده شود. در طول ذره ممکن است یک عدد که نشان دهنده شماره کلاستر می باشد، تکرار شود. اما باید از همه این K عدد استفاده شود. در این مقاله ساختار هر

۱	۱	۳	۲	۲	۱	۳	۴	۴
---	---	---	---	---	---	---	---	---

ذره مانند شکل (۱) تعریف می شود.

شکل (۱) ساختار یک ذره

روش پیشنهادی ما به اینصورت است که به جای استفاده از فاصله اقلیدسی دو بردار یا فاصله کسینوسی آنها، از فاصله معنایی آنها استفاده کنیم. ما در مراحل قبل، میزان شباهت معنایی دو جمله را با استفاده از بردار های متنی بدست آوردیم. در این روش باید دو معیار مهم در نظر گرفته شوند.

(۱) جملات باید طوری خوشه بندی شوند که جملات موجود در یک خوشه باید دارای بیشترین شباهت ممکن به یکدیگر باشند. (تابع درون خوشه ای (Intra Clustering))

(۲) جملات یک خوشه باید دارای فاصله زیادی با خوشه های دیگر داشته باشند. (تابع برون خوشه ای (Inter Clustering)) بنابراین تابعی که بتواند این معیارها را ارضاء کنند به صورت زیر تعریف می شود:

تابع درون خوشه ای: (Intra Clustering)

$$P = \left(\sum_{p=1}^k c_p \sum_{s_i, s_j} sim(s_i, s_j) \right) \rightarrow MAX \quad (11)$$

تابع برون خوشه ای (Inter Clustering)

$$R = \sum_{P=1}^{K-1} \sum_{Q=P+1}^K \sum_{s_i \in c_p} \sum_{s_j \in c_q} c_p c_q sim(s_i, s_j) \rightarrow MIN \quad (12)$$

توابع ترکیبی

$$H = P * R \rightarrow MAX \quad (13)$$

$$E = w_2 * p + (1 - w_2) * R \rightarrow MAX \quad (14)$$

$$F = \frac{2}{\frac{1}{P} + \frac{1}{R}} = \frac{2 * P * R}{P + R} \rightarrow MAX \quad (15)$$

تابع ترکیبی H، تابع حاصل ضرب دو تابع P و R است. تابع E تابع ترکیب وزنی دو تابع است. و تابع F میانگین هارمونیک دو تابع می باشد. تابع E نشاندهنده میزان مشارکت دو تابع P و R می باشد به طوری که اگر w_2 برابر صفر باشد تابع E نشاندهنده تابع R خواهد بود و چنانچه یک انتخاب شود نشاندهنده P خواهد بود. و اگر ۰.۵، انتخاب شود، هر کدام از توابع (P, R) سهم یکسانی را خواهند داشت. در این توابع K تعداد خوشه ها ست. C_p تعداد جملات موجود در خوشه است. $sim(s_i, s_j)$ میزان شباهت ذخیره شده در ماتریس شباهت جمله (جدول (۵)) است.

۳- تعداد خوشه های بهینه

شاخص های زیادی برای پیدا کردن تعداد خوشه های بهینه معرفی شده است که اغلب آن ها از دو پارامتر فشرده گی و تفکیک استفاده می کنند. که ما از تابع زیر استفاده می کنیم [7]:

$$V_{ram}(K) = \frac{comp(k)}{separ(k)} \quad (17)$$

هنگامی که تعداد خوشه ها زیاد می شود. میزان فشرده گی خوشه ها کاهش می یابد. و مقدار $V_{ram}(K)$ کاهش میابد. در نتیجه $V_{ram}(K+1) \leq V_{ram}(K)$.. ما برای رسیدن به تعداد خوشه های بهینه از روش سعی و تکرار استفاده می کنیم به طوریکه هدف مینیم کردن تعداد خوشه هلیبی است که در رابطه زیر صدق کند:

$$H = \frac{V_{ram}(K) - V_{ram}(K+1)}{V_{ram}(K)} \leq \varepsilon \quad (18)$$

که این رابطه تعداد خوشه های بهینه را تعیین می کند و ε یک استانه می باشد.

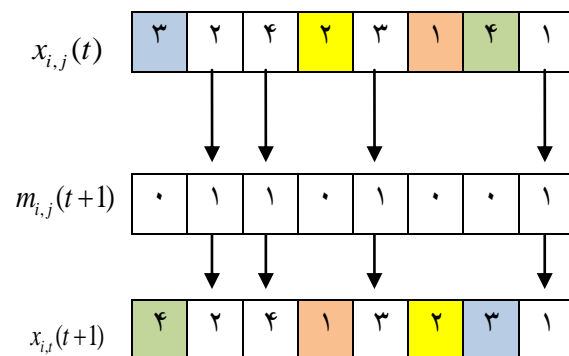
۴- نتایج

ما یک مجموعه خبر ورزشی مربوط به خبرگزاری دانشجویان ایران (ایسنا) را به عنوان پیکره متنی خود انتخاب کرده ایم. و تمامی آزمایشات و ارزیابی ها را روی این مجموعه داده انجام داده ایم. این مجموعه شامل ۸ مجموعه خبر ورزشی می باشد. دو روش ارزیابی برای سیستم های بازیابی اطلاعات وجود دارد ارزیابی مستقیم و ارزیابی معطوف به کاربرد. در ارزیابی مستقیم، نتیجه سیستم با نتیجه کار چند انسان بصورت مستقیم ارزیابی می شود. با توجه به اینکه این مقاله به خوشه بندی جملات فارسی پرداخته است و ابزاری برای ارزیابی آن وجود ندارد. از ارزیابی مستقیم استفاده خواهیم کرد. این روش ارزیابی مشابه روش استفاده شده در [5] می باشد. دو معیار مهم در خوشه بندی جملات عبارتند از سودمندی نسبی درون خوشه ای (CBRU) و اشتراک اطلاعات بین جمله ای (CSIS). در حالت ایده ال خوشه بندی مناسب است که CBRU را که بیانگر میزان ارتباط جملات موجود در یک خوشه است، ماکزیمم کند و معیار CSIS را کاهش دهد. در این ارزیابی ها از چند داور خواسته می شود که به جملات موجود در تمام کلاسترها نمره دهند و سپس کارایی کلاستر مورد نظر بدست خواهد آمد... (تمام)

این ذره نشان می دهد که جملات اول دوم و ششم در کلاستر اول. جملات سوم و هفتم در کلاستر سوم. جملات چهارم و پنجم در کلاستر دوم و جملات هشتم ونهم در کلاستر چهارم می باشد. تابع برازنده گی مورد استفاده در الگوریتم، یکی از توابع درون خوشه ای، برون خوشه ای و توابع ترکیبی بخش 2-4 می باشد. مسئله ای که در این الگوریتم باید به آن اشاره کرد این است که ذرات تولیدی نباید مقادیری خارج از سه شرط ذکر شده تولید کنند. اما این مسئله در الگوریتم PSO اجتناب ناپذیر است. که ناشی از وجود بردارهای سرعت ذرات می باشد. که ممکن است مقادیر ذرات تولیدی خارج از محدودیت های ذکر شده (حتی اعشاری) شوند. که برای حل این مسئله می توان از الگوریتم ژنتیک بهره برد و برداری به نام بردار جهش به صورت زیر استفاده کرد.

$$m_{i,j}(t+1) = \begin{cases} 1, & \text{if } Rand_j \leq sigm(v_{i,j}(t+1)) \\ 0, & \text{otherwise} \end{cases} \quad (16)$$

بردار $m_{i,j}(t+1)$ با تاثیر گذاشتن روی $x_{i,j}(t)$ ، آن را به بردار $x_{i,t}(t+1)$ تبدیل می کند. چنانچه ز امین عنصر از بردار $m_{i,j}(t+1)$ برابر یک باشد. عنصر ز ام بردار $x_{i,j}(t)$ بدون تغییر به ز امین عنصر بردار $x_{i,t}(t+1)$ منتقل می شود. و در صورتی که ز امین عنصر بردار $m_{i,j}(t+1)$ صفر باشد. عنصر ز ام بردار $m_{i,j}(t+1)$ جهش خواهد داشت. این جهش ناشی از تاثیر عملگر معکوس بروی آن می باشد شکل ۲۰ شکل گیری بردار $x_{i,t}(t+1)$ با استفاده از $m_{i,j}(t+1)$ را نشان می دهد.



شکل ۲) عملکرد بردار جهش

روش ها به ازای متن واحدی سنجیده می شود. (به شکل (۳) توجه کنید:

جدول (۵): متن خیر: دیدار دکتر احمدی نژاد از اردوی تیم ملی

[1] حضور رئیس جمهور اسلامی ایران در تمرینات تیم ملی فوتبال، فضای خاصی را به اردوی تیم ملی بخشید.
 [2] به گزارش خبرگزاری دانشجویان ایران (ایسنا)، به محض ورود رئیس جمهور بر سر تمرینات تیم ملی از درب ورودی ورزشگاه آزادی، تماشاگران تمرینات تیم ملی با شعار "انرژی هسته ای حق مسلم ماست" و "احمدی نژاد دوست داریم" به تشویق او پرداختند.
 [3] احمدی نژاد سپس با برانکو به گفت و گو پرداخت و نفس تلاش بازیکنان تیم ملی را ارزشمند و یک پیروزی خواند.
 [4] رئیس جمهور در ادامه با تک تک بازیکنان خوش و بش کرد و ابراهیم میرزاپور به گفت و گو و درد دل با رئیس جمهور پرداخت.
 [5] رییس جمهور در ادامه با تک تک بازیکنان خوش و بش کرد و ابراهیم میرزاپور به گفت و گو و درد دل با رییس جمهور پرداخت.
 [6] گفتنی است دکتر احمدی نژاد پیش از سخنانش، دقایقی را با ملی پوشان تمرین کرد و به زدن ضربات پنالتی پرداخت. رییس جمهور ۴ ضربه ی پنالتی به ابراهیم میرزاپور دروازه بان تیم ملی زد که ۳ تای آن تبدیل به گل شد و یکی از آنها را میرزاپور دفع کرد.
 [7] رئیس جمهور همچنین با حضور در رختکن تیم ملی و با پوشیدن گرمکن تیم ملی به زمین بازگشت.
 [8] دکتر احمدی نژاد هنگام ورود به زمین چمن به گرم کردن خود

همین روش نیاز است تعداد خوشه های بهینه مشخص شود. جدول ۵. با چهار خوشه در نظر گرفته شده است. تعداد خوشه های بهینه در جدول (7) تعیین شده است. (استانه 0.002).

جدول (۷): تعداد خوشه های بهینه در دروش P

K	1	2	3	4	5
H	0.007	0.006	0.003	0.002	0.001

اکنون خوشه بندی بدست آمده از روش قبل باید ارزیابی شود. به اینصورت کند که تعدادی داور به تمام جملات موجود در یک کلاستر امتیاز می دهند این امتیازات (۰ تا ۱۰) باید طوری به جملات نسبت داده شود که تعیین کند جمله مورد نظر تا چه اندازه به موضوع عمومی خوشه نزدیک است. به جدول (۸) توجه کنید. (خوشه بندی بدست آمده از مرحله قبل)

جدول (۸): ارزیابی یک کلاستر

داورها	داور اول	داور دوم	داور سوم	مجموع
داور اول	۱	۰,۴۳	۰,۳۵	۰,۵۹
داور دوم	۰,۴۳	۱	۰,۲۹	۰,۵۷
داور سوم	۰,۳۵	۰,۲۹	۱	۰,۵۴
				۰,۵۶

با توجه به اینکه از هر خوشه دو جمله با امتیاز بالا (متناسب با نرخ فشرده گی) توسط داورها انتخاب می شود. کارایی کلاستر اول توسط روش P برابر ۰,۵۶ است. بقیه کلاسترهای این روش توسط داورها ارزیابی می شوند. سپس میانگین آنها، کارایی کل خوشه بندی را نشان خواهند داد. در تمامی این روش ها پارامترهای الگوریتم PSO به اینصورت مقداردهی شدند: جمعیت اولیه ذرات برابر ۴۰ است. (N=40). تعداد تکرار ۲۵۰ (t=250) و پارامترهای شناختی و اجتماعی $c_1 = 2.5, c_2 = 4$ و وزن اینرسی برابر $w = 0.5$ است. تمام نتایج ثبت شده، حاصل اجرای میانگین ۱۰ بار اجرای هر روش می باشد در تمام ارزیابی ها در تمام این روش ها با تعداد خوشه های بهینه انجام می گیرد. بنابراین ابتدا باید تعداد خوشه های بهینه هر روش را محاسبه کرد. و در روش E فاکتور $w_2 = 0.4$ در نظر گرفته شده است. انتخاب این مقدار برای w_2 با استفاده از ارزیابی تاثیر آن روی مسئله کارایی صورت گرفته است [10]. برای ارزیابی نهایی از هشت مجموعه خبری استفاده شده است. به طوریکه نتایج ثبت شده در جدول (۹)

جدول (6). جواب همگرا در روش P

Iteration	gbest	خوشه بندی
1	2.0330	(1,1,2,3,4,1,2,3,2)
20	2.5123	(1,2,2,2,2,2,3,3,4)
40	2.7901	(1,1,1,1,3,3,2,2,1)
50	2.7901	(1,1,1,1,3,3,2,2,1)
100	3.2301	(1,1,1,1,3,3,2,2,1)
120	3.4578	(1,1,2,3,3,1,1,2,3)
150	3.4897	(1,1,3,2,2,3,1,2,2)
180	3.4897	(1,1,3,2,2,3,1,2,2)
200	4.7389	(1,1,1,1,2,3,3,2,2)
220	4.7389	(1,1,1,1,2,3,3,2,2)
250	4.7389	(1,1,1,1,2,3,3,2,2)

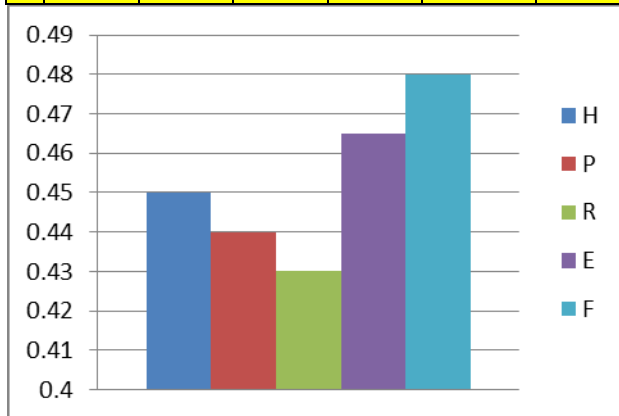
در جدول (6) همانطور که ملاحظه می شود پاسخ های گلوبال در تکرار ۲۰۰ به جواب نهایی همگرا شده است. (K=4) در

و (۱۰) کارایی روش های خوشه بندی را روی مجموعه های خبری (دیدار احمدی نژاد از اردوی تیم ملی و قهرمانی سایپا) نشان می دهد.

جدول (۹): مقایسه روش های خوشه بندی روی مجموعه خبر اول

کارایی خوشه بندی	P	R	H	E	F	
P	۰,۴۲۱	X	+0.119	+0.276	+0.266	+0.211
R	۰,۴۷۸	-0.135	X	+0.078	+0.096	+0.104
H	۰,۵۱۹	-0.232	-0.085	X	+0.018	+0.028
E	0.۵۲۹	-0.256	-0.628	-0.019	X	+0.009
F	۰,۵۳۴	-0.268	-0.117	-0.028	-0.009	X

کارایی خوشه بندی	P	R	H	E	F	
P	0.515	X	-0.040	+0.006	+0.013	+0.021
R	0.495	+0.038	X	+0.044	+0.051	+0.058
H	0.518	-0.006	-0.046	X	+0.008	+0.025
E	0.522	-0.013	-0.600	-0.007	X	+0.008
F	0.526	-0.021	-0.062	-0.015	-0.007	X



شکل (۳): مقایسه روش های خوشه بندی

در جدول (۹) برای تعیین کردن میزان بهبودی روش ها در مقایسه با یک روش خاص، از فرمول زیر استفاده می شود:

$$\frac{OurMethod - OtherMethod}{OtherMethod} * 100 \quad (19)$$

علامت مثبت نشان دهنده میزان بهبودی یک روش نسبت به روش خاص دیگر را نشان می دهد و علامت منفی نشان دهنده عکس آن است. به طور مثال کارایی روش H برابر ۰,۵۱۸ است و روش E و F در مقایسه با روش H، به ترتیب به اندازه ۰,۸ و ۲,۵ درصد عملکرد بهتری نسبت به آن دارند. جدول زیر عملکرد روش های مختلف خوشه بندی روی مجموعه خبری دو (قهرمانی سایپا را نشان می دهد).

جدول (۱۰): مقایسه روش های خوشه بندی روی مجموعه خبر دوم

با توجه به جدول (۱۰) همانطور که ملاحظه می شود روش خوشه بندی F عملکرد بهتری در مقایسه با روش های دیگر دارد. با اجرای تمام روش ها روی هشت مجموعه خبری نتایج نشان می دهد روش خوشه بندی میانگین هارمونیک (F)، کارایی بهتری در مقایسه با روش های دیگر دارد. در مقابل روش های P و R عملکرد ضعیفی نسبت به بقیه خواهند داشت. (میانگین روی هشت مجموعه خبر)

۵- نتیجه گیری

یک مسئله مهم در خلاصه سازی متون، استخراج موضوعات و جملاتی است که تا حدود زیادی بتوانند مفهوم متن اصلی را برسانند. در سال های اخیر مقالات زیادی به این موضوع پرداخته اند. و اغلب آن ها برای پیدا کردن شباهت میان متن ها و جملات از معیار فاصله کسینوسی و اقلیدسی استفاده کرده اند. در خوشه بندی جملات یک متن برای مشخص شدن جملات مشابه، نمی توان از روش مشابه (دسته بندی متون مشابه) استفاده کرد. بردارهایی به طول m و با مقادیر صفر بسیار زیاد پدید خواهد آمد. ما برای حل این مشکل در این مقاله با استفاده از بردارهای Context-Vector ارتباط معنایی واژه ها را یافتیم و سپس با تعریف معیارها و توابعی که سعی در ارضاء ویژگی های مسئله خوشه بندی داشتند، توانسیم توابع مناسب را استخراج کنیم. و برای بهینه شدن هر کدام از این توابع الگوریتم PSO را بکار برده ایم. ما در این مقاله چند تابع متناظر با ویژگی های خوشه بندی تعریف نمودیم. روش R بالا بودن میزان فشرده گی خوشه ها را در نظر می گیرد. در روش P بالا بودن

تعداد خوشه ها اهمیت می یابد. در بین روش های ترکیبی با توجه به نتایج حاصل از ارزیابی روش F (تابع میانگین هارمونیک) مناسبترین روش به لحاظ کارایی معرفی شد.

۶- مراجع

- [1] Ching-Yi Cheo , Fun Ye."Particle Swarm Optimization Algorithm and Its Application". Intemationai Conference on Networking, Sensing Control. Taiwi. Taiwan. March 21-23, 2004
- [2] Zalik, Krista Rizman . "An Efficient k'-means Clustering Algorithm."Pattern Reconition Letters, Vol. 29, I.9. Pag. 1385-1391. Elsevier 07/2008. to Clustering Analysis.
- [3] R. J. Kuo • M. J. Wang • T. W. Huang .” An application of particle swarm optimization algorithm to clustering analysis “ .Springer-Verlag 2009.
- [4] Ramanji Killani, K. Srinivasa Rao.“Effective Document Clustering with Particle Swarm Optimization”. Springer-Verlag Berlin Heidelberg 2010.
- [5] D. R. Radev, H. Jing, M. Stys, and D. Tam (2004), Centroid-based summarization of multiple documents. Information Processing and Management, 2004.
- [6] Dias, G. and Alves, E. (2005). Unsupervised Topic Segmentation Based on Word Co-occurrence and Multi-Word Units for Text Summarization. In Proceedings of the ELECTRA Workshop associated to 28th Annual International ACM SIGIR Conference, Salvador, Brazil, pp. 41-48.
- [7] R.M.Aliguliyev.”Automation Documents Summarization By Sentence Extraction:.
- [8] Eberhart R, Kennedy J (1995) A new optimizer using particle swarm theory. In: Proceedings of the sixth international symposium on micro machine and human science, pp 39–43

[۹] مشکئی محسن. انالویی مرتضی "بازیابی خبرهای مرتبط پیشین برای تولید خلاصه های پیشینه-خبر" اولین کنفرانس برق , دانشگاه علم و صنعت 1388

[۱۰] بازقندی , مهدی, تدین تبریزی, قمرناز. وفایی جهان مجید. "خلاصه سازی متون فارسی مبتنی بر خوشه بندی PSO". اولین کنفرانس بین المللی پردازش خط و زبان فارسی. دانشگاه سمنان, ۱۳۹۱.