

# مقایسه عملکرد اشتراک در مآخذ و اشتراک واژگان عناوین استنادها در خوشه‌بندی فازی پروانه‌های ثبت اختراع<sup>۱</sup>

آناهیتا کرمانی\*

کارشناسی ارشد کتابداری و اطلاع‌رسانی

نرگس نشاط<sup>۲</sup>

دانشیار،

سازمان اسناد و کتابخانه ملی ج.ا. ایران

عباس حری<sup>۳</sup>

استاد،

گروه کتابداری و اطلاع‌رسانی، دانشگاه تهران

دریافت: ۱۳۹۰/۰۹/۲۴ | پذیرش: ۱۳۹۱/۰۲/۱۳

فصلنامه علمی پژوهشی  
پژوهشگاه علوم و فناوری اطلاعات ایران  
شاپا(چاپی) ۸۲۲۳-۲۲۵۱  
شاپا(الکترونیکی) ۸۲۳۱-۲۲۵۱  
نمایه در SCOPUS و ISC  
<http://jipm.irandoc.ac.ir>  
دوره ۲۸ | شماره ۲ | صص ۴۱۱-۴۳۲  
زمستان ۱۳۹۱  
نوع مقاله: پژوهشی

**چکیده:** هدف پژوهش حاضر، ارزیابی و مقایسه خصیصه استناد در قالب اشتراک در مآخذ و خصیصه واژگان عناوین استنادها در خوشه‌بندی پروانه‌های ثبت اختراع است. بدین منظور، این بررسی به روش تجربی بر روی مجموعه‌ای شامل پروانه‌های ثبت اختراع یو.اس. مربوط به رده ۹۷۷/۷۷۴ از رده‌بندی یو.اس. انجام شد. ابتدا پروانه‌های ثبت اختراع با استفاده از خوشه‌بندی فازی C میانگین خوشه‌بندی شد و خوشه‌بندی‌های حاصل با استفاده از معیارهای ارزیابی مانعیت و جامعیت گسترش یافته بی کیوبد مورد ارزیابی قرار گرفت. یافته‌ها نشان داد که خوشه‌بندی پروانه‌های ثبت اختراع با استفاده از اشتراک در مآخذ، عملکرد بهتری را نسبت به خوشه‌بندی با استفاده از واژگان عناوین استنادها دارد و ساختار خوشه‌ای در سطوح گسترده‌تری از جامعیت در خوشه‌بندی با استفاده از اشتراک در مآخذ برقرار است.

**کلیدواژه‌ها:** پروانه‌های ثبت اختراع، خوشه‌بندی، اشتراک در مآخذ، واژگان عناوین استنادها، خوشه‌بندی فازی C میانگین، ارزیابی خوشه‌بندی، مانعیت گسترش یافته بی کیوبد، جامعیت گسترش یافته بی کیوبد.

. برگرفته از پایان‌نامه کارشناسی ارشد  
پدیدآور رابط با عنوان ارزیابی تأثیر استفاده از  
واژگان عناوین استنادی در مقایسه با استفاده از  
اشتراک استنادی در خوشه‌بندی پروانه‌های ثبت  
اختراع.

\*kermanianahita@gmail.com  
2. narges\_neshat@yahoo.com  
3. riwash@yahoo.com

## ۱. مقدمه

پروانه‌های ثبت اختراع یکی از منابع مهم اطلاعاتی در صنعت به‌شمار می‌رود که اطلاعات موجود در آنها می‌تواند به اهداف گوناگون، مورد استفاده و جستجو قرار گیرد. از جمله این اهداف می‌توان «تعیین قابلیت ثبت»<sup>۱</sup> اختراعات مورد ادعا و یا «رد اعتبار»<sup>۲</sup> اختراعات ثبت شده به‌واسطه یافتن اختراعی مشابه اشاره کرد ( Graf and Azzopardi 2008; Bonino, Ciaramella, and Corno 2010; Hideo, Azzopardi, and Vanderbauwhede 2010). اهمیت بازیابی دقیق پروانه‌های ثبت اختراع از آنجا ناشی می‌شود که گاه فقط بازیابی یک مورد مرتبط می‌تواند منجر به رد اعتبار ادعا(ها)ی اختراع دیگری شود. به‌عبارت دیگر، در نظام‌های پژوهشی (ماهوی) ثبت اختراع، ثبت یک اثر اختراع منوط به دارا بودن شرایطی از جمله تازگی و نو بودن<sup>۳</sup> و غیربدیهی بودن<sup>۴</sup> آن است (United States Patent and Trademark Office (USPTO) 2010). از جمله فعالیت‌هایی که برای بررسی چنین شرایطی صورت می‌گیرد، جستجو در آثار پیشین<sup>۵</sup> است. همچنین، ممکن است اختراعی به ثبت رسیده باشد، اما برای تأیید یا رد اعتبار آن نیاز به جستجو در آثار پیشین باشد.

خوشه‌بندی روشی است که برای گروه‌بندی مدارک مشابه مورد استفاده قرار می‌گیرد و شیوه‌ای مناسب در بازیابی (Kang et al. 2007) و ترسیم نقشه‌های علمی (Leydesdorff 1987) فراهم می‌کند. در این روش، مدارک در گروه‌هایی از پیش تعیین نشده به‌نام خوشه قرار می‌گیرند، به‌طوری که مدارک مشابه در کنار یکدیگر و مدارک و ماهیت‌های نامشابه دور از یکدیگر قرار گیرند. عوامل گوناگونی در خوشه‌بندی مدارک مؤثر است. یکی از این عوامل مربوط به نوع خصیصه‌ای<sup>۶</sup> است که به‌واسطه آن یک مدرک مورد بازنمایی قرار می‌گیرد (Dhillon, Kogan, and Nicholas 2003). به‌همین دلیل، انتخاب درست خصیصه برای مدارک در نتایج خوشه‌بندی آنها تأثیرگذار خواهد بود.

استناد یکی از پرکاربردترین خصیصه‌هایی است که برای بازنمایی پروانه‌های ثبت اختراع مورد استفاده قرار می‌گیرد. اهمیت استنادها در پروانه‌های ثبت اختراع از آن جهت است که از یک طرف از سوی افراد متخصصی به‌نام بازرس تعیین می‌شوند و از سوی دیگر، به‌دلیل مرتبط بودن با ادعا(های) اختراع انتخاب می‌شوند. این استنادها در واقع، گزارش‌هایی هستند که اختراعی بودن یک اثر را مورد تأیید قرار می‌دهند (United States Patent and Trademark Office

1. Patentability	2. invalidity	3. novelty
4. non-obviousness	5. prior art	6. attribute

(USPTO) 2010<sup>1</sup>. از این رو، می‌توانند به‌عنوان خصیصه‌ای مناسب برای بازنمون پروانه‌های ثبت ثبت اختراع مورد استفاده قرار گیرند.

استفاده از استناد در بازنمایی پروانه‌های ثبت اختراع می‌تواند به شیوه‌های گوناگون صورت پذیرد. یکی از این شیوه‌ها، استفاده از «استناد» در خوشه‌بندی است که بر مبنای «اشتراک در مآخذ» انجام می‌شود. یکی دیگر از راه‌ها، استفاده از «استناد» در خوشه‌بندی است که بر مبنای «واژگان عناوین استنادها» صورت می‌گیرد. بدیهی است هر روزه تعدادی اختراع به ثبت می‌رسد و مورد استناد نیز قرار می‌گیرد. این امر سبب افزایش حجم استنادها و در نتیجه عاملی تأثیرگذار در خوشه‌بندی آنها خواهد بود. به همین دلیل، پیامدهای این امر سبب می‌گردد که شکل‌های مختلف پردازش و خوشه‌بندی بر مبنای خصیصه استنادی مورد مذاقه قرار گیرد تا بهترین راه که با جامعیت و مانعیت بالاتری همراه است، انتخاب شود. استفاده از استناد در بازنمایی پروانه‌های ثبت اختراع می‌تواند به شیوه‌های گوناگون صورت پذیرد. از این رو، پژوهش حاضر بر آن است تا با ایجاد یک خوشه‌بندی مبنای پروانه‌های ثبت اختراع، به ارزیابی خوشه‌بندی از طریق واژگان عناوین استنادها و خوشه‌بندی از طریق اشتراک در مآخذ پردازد و در این راستا، به پرسش‌های زیر پاسخ دهد:

۱. میزان جامعیت و مانعیت گسترش یافته بی‌کیوبد در خوشه‌بندی پروانه‌های ثبت اختراع با استفاده از اشتراک در مآخذ چقدر است؟
۲. میزان جامعیت و مانعیت گسترش یافته بی‌کیوبد در خوشه‌بندی پروانه‌های ثبت اختراع با استفاده از واژگان عناوین استنادها چقدر است؟
۳. کدام یک از دو نوع خوشه‌بندی ایجادشده از عملکرد بهتری (با توجه به معیارهای یادشده) برخوردار است؟

## ۲. برخی تعاریف

### ۱-۲. خوشه‌بندی

خوشه‌بندی مدارک عبارت است از تقسیم مجموعه‌ای از مدارک به تعدادی طبقه نامشخص به نام «خوشه». هدف از خوشه‌بندی، قرار دادن مدارک در خوشه‌ها به گونه‌ای است که اعضای هر خوشه از یک سو با یکدیگر دارای بیشترین شباهت و از سوی دیگر، با اعضای سایر خوشه‌ها دارای کمترین شباهت باشند (Tan, Steinbach, and Kumar 2006).

۱. البته باید توجه داشت که نظام‌های ثبت اختراع در مراکز گوناگون متفاوت است و توصیف یادشده از استناد در مراکزی صدق می‌کند که در آن ثبت اختراع مستلزم بررسی ماهوی اختراع است.

برای خوشه‌بندی پروانه‌های ثبت اختراع (مدارک) لازم است نوعی خصیصه برای پروانه‌ها تعیین شود تا به واسطه آن پروانه‌ها مورد بازنمون قرار گیرند (Jain and Dubes 1988 cited in Jain, Murty, and Flynn 1999). پس از آن، میزان شباهت (نزدیکی) مدارک به یکدیگر براساس معیار شباهت، محاسبه می‌شود و مدارک مشابه در کنار یکدیگر و مدارک نامشابه دور از هم قرار خواهند گرفت.

گاه ممکن است ارزش خصیصه‌ها در مدارک گوناگون متفاوت باشد. برای قائل شدن این تفاوت می‌توان به هر خصیصه در هر مدرک وزنی متفاوت اختصاص داد. یکی از روش‌های وزن‌دهی به خصیصه‌ها، روش tf-idf است (tf-idf 2011). در این روش، وزن هر خصیصه از ضرب «بسامد حضور یک خصیصه در یک مدرک» (tf) در «لگاریتم تعداد همه مدارک بخش بر تعداد مدارک حاوی هر خصیصه» (idf) به دست می‌آید.

استفاده از فضای برداری، یکی از فنون بازنمایی مدارک است. در این روش، مدارک با استفاده از خصیصه‌های وزن داده شده (به عنوان مؤلفه‌های بردار هر مدرک) مورد بازنمایی قرار می‌گیرند. معیار شباهت کوسینوسی یکی از متداول‌ترین معیارها برای سنجش شباهت میان مدارک است (Larsen and Aone 1999 cited in Huang 2008). بنابر این معیار، اگر  $d_j$  و  $d_k$  دو مدرک باشند شباهت میان آنها (کوسینوس زاویه میان مدارک آنها) از فرمول زیر به دست می‌آید:

$$\text{sim}(d_j, d_k) = \frac{\vec{d}_j \cdot \vec{d}_k}{|\vec{d}_j| |\vec{d}_k|}$$

که در آن  $\vec{d}_j$ ،  $\vec{d}_k$  ضرب داخلی بردار دو مدرک در یکدیگر و  $|\vec{d}_j|$  و  $|\vec{d}_k|$  اندازه بردارهای  $\vec{d}_j$  و  $\vec{d}_k$  هستند. از آنجا که وزن هر خصیصه، مؤلفه‌های بردار هر مدرک را تشکیل می‌دهد، محاسبه شباهت بین مدارک را می‌توان به طریق زیر محاسبه کرد:

$$\text{sim}(d_j, d_k) = \frac{\sum_{i=1}^n w_{i,j} w_{i,k}}{\sqrt{\sum_{i=1}^n w_{i,j}^2} \sqrt{\sum_{i=1}^n w_{i,k}^2}}$$

در این صورت، فاصله میان مدارک برابر است با (Huang 2008):

$$\text{distance} = 1 - \text{similarity}$$

بنابراین، هرچه میزان شباهت بیشتر باشد، فاصله میان دو مدرک کمتر خواهد بود.

#### ۱-۱-۲. خوشه‌بندی فازی C میانگین<sup>۱</sup>

خوشه‌بندی مدارک از طریق الگوریتم‌های گوناگون امکان‌پذیر است. طراحی برخی از این الگوریتم‌ها به گونه‌ای است که یک مدرک فقط به یک خوشه تعلق خواهد گرفت. در مقابل، در برخی الگوریتم‌ها امکان تعلق یک مدرک به بیش از یک خوشه امکان‌پذیر است. گونه نخست، خوشه‌بندی ناهمپوشان و گونه دوم، خوشه‌بندی همپوشان نامیده می‌شود. خوشه‌بندی فازی C میانگین (C به معنی خوشه) نوعی خوشه‌بندی همپوشان است. استفاده از این الگوریتم در گروه‌بندی پروانه‌های ثبت اختراع این امکان را فراهم می‌کند که پروانه‌هایی که هم‌زمان به بیش از یک موضوع تعلق دارند بتوانند به بیش از یک خوشه نیز تعلق داشته باشند. در نتیجه، این الگوریتم نسبت به خوشه‌بندی‌های ناهمپوشان، تناسب بیشتری با ماهیت‌های مورد خوشه‌بندی (پروانه‌های ثبت اختراع) خواهد داشت.

برای ایجاد خوشه‌بندی فازی C میانگین مراحل زیر دنبال می‌شود (Wedding 2009):

۱. انتخاب تعداد دلخواه  $k$  خوشه برای مجموعه‌ای شامل  $N$  مدرک<sup>۲</sup>، به طوری که  $k < N$ ؛
- که در اینجا  $K$  برابر است با ۷۵ خوشه و  $N$  برابر با ۷۱۷ پروانه ثبت اختراع است.
۲. ایجاد یک نقطه مرکزی آغازین به‌ازای هر یک از  $k$  خوشه؛
۳. محاسبه فاصله هر یک از  $N$  مدرک تا هر یک از  $k$  خوشه؛
۴. تعیین یک عضویت فازی یا کسری در هر یک از  $k$  خوشه برای هر یک از  $N$  مدرک؛
۵. یافتن نقطه مرکزی جدیدی برای هر یک از  $k$  خوشه از طریق یافتن میانگین وزنی رکوردها؛ و
۶. تکرار مراحل ۳، ۴، و ۵ تا هنگامی که تغییری در عضویت خوشه‌ای به‌وجود نیاید (یا تا هنگام رسیدن به معیار همگرایی).

#### ۱-۱-۲-۱. تعیین درجه عضویت مدارک در هر یک از خوشه‌ها

درجه عضویت فازی، از محاسبه فاصله‌های مدارک از مرکز خوشه‌ها به‌دست می‌آید.

#### 1. Fuzzy C-means clustering

۲. در توصیف الگوریتم، از واژه data استفاده شده است، زیرا کاربرد خوشه‌بندی فقط به خوشه‌بندی مدارک محدود نیست. اما، از آنجا که در اینجا خوشه‌بندی برای مدارک استفاده شده واژه مدرک مورد استفاده قرار گرفته است.

هنگامی که مدرکی به مرکز یک خوشه نزدیک تر از مرکز دیگر خوشه‌ها باشد، درجه عضویت مدرک در آن خوشه بیشتر و هنگامی که مدرک دورتر باشد، درجه عضویت کمتری خواهد داشت. مجموع مقادیر عضویت فازی [یک مدرک] نیز برابر با یک است (Wedding 2009).

درجه عضویت از طریق فرمول زیر محاسبه می‌شود:

$$u_k = \frac{1}{\sum_{i=1}^j \left(\frac{d_k}{d_i}\right)^p}$$

که در آن  $k$  یکی از  $j$  خوشه،  $d$  فاصله، و  $p$  توان است. توان ( $p$ ) خود از فرمول زیر به دست

می‌آید:

$$p = \frac{2}{(m-1)}$$

که در آن  $m$  یک توان<sup>۱</sup> فازی است به طوری که  $1 < m < \infty$  است. مقدار  $m$  بر درجه فازی بودن خوشه‌ها تأثیرگذار است. هنگامی که  $m \rightarrow 1$  ( $m$  نزدیک به ۱ باشد)، درجه فازی شدن خوشه‌ها کمتر و هنگامی که  $m \rightarrow \infty$  ( $m$  به سمت بی‌نهایت میل کند)، خوشه‌ها به‌طور فزاینده‌ای فازی خواهند شد. برای مقادیر نزدیک به بی‌نهایت، مقدار عضویت‌های فازی یک مدرک در همه خوشه‌ها به برابری می‌گراید. به‌طور معمول، برای  $m$  مقدار ۲ در نظر گرفته می‌شود (Wedding 2009).

## ۲-۲. ارزیابی خوشه‌بندی‌های همپوشان: جامعیت و مانعیت گسترش یافته بی‌کیوبد<sup>۲</sup>

معیارهای مورد استفاده در ارزیابی خوشه‌بندی‌های همپوشان با خوشه‌بندی‌های ناهمپوشان تفاوت دارد (Amigo et al. 2008). جامعیت و مانعیت گسترش یافته بی‌کیوبد از معیارهای ارزیابی خوشه‌بندی‌های همپوشان هستند که از گسترش جامعیت و مانعیت بی‌کیوبد<sup>۳</sup> به دست آمده‌اند (Amigo et al. 2008). معیارهای بی‌کیوبد در حوزه‌هایی چون تعیین مدارک هم‌مرجع (در اینجا، هم‌خوشه<sup>۴</sup>)، بازیابی اطلاعات، و استخراج اطلاعات مورد استفاده قرار می‌گیرد (Bagga and Baldwin 1998).

1. exponent

2. extended BCubed precision and recall

۳. براساس ایمیل دریافت‌شده در ۶ فوریه ۲۰۱۲ از کوهان سوچای کارلس (Cohan Sujay Carlos)، بی‌کیوبد (B-Cubed)، به معنی سه  $b$ ، نام الگوریتم این معیارهاست که از اولین حرف نام ارائه‌دهندگان الگوریتم‌ها یعنی Biermann, Bagga, and Baldwin برگرفته شده است.

4. cross-document coreferencing

برای محاسبه جامعیت و مانعیت گسترش یافته بی کیوبد، ابتدا باید به محاسبه مانعیت و جامعیت چندگانه<sup>۱</sup> پرداخت. مانعیت چندگانه<sup>۲</sup> برابر است با مینیمم میان تعداد هم‌رخدادی دو مدرک در «خوشه‌بندی» و هم‌رخدادی آنها در «خوشه‌بندی مرجع» بخش بر تعداد هم‌رخدادی آنها در تعداد هم‌رخدادی در «خوشه‌بندی» همچنین جامعیت چندگانه<sup>۳</sup> برابر است با مینیمم تعداد هم‌رخدادی دو مدرک در «خوشه‌بندی» و تعداد هم‌رخدادی آنها در «خوشه‌بندی مرجع» بخش بر تعداد هم‌رخدادی آنها در «خوشه‌بندی مرجع». این تعاریف به بیان ریاضی به صورت زیر تعریف می‌شوند (Amigo et al. 2008):

$$\text{Multiplicity Precision}(e, e') = \frac{\text{Min}(|C(e) \cap C(e')|, |L(e) \cap L(e')|)}{|C(e) \cap C(e')|}$$

$$\text{Multiplicity Recall}(e, e') = \frac{\text{Min}(|C(e) \cap C(e')|, |L(e) \cap L(e')|)}{|L(e) \cap L(e')|}$$

که در آن  $e$  و  $e'$  نشان‌دهنده دو مدرک است.  $L(e)$  مجموعه خوشه‌های مرجع و  $C(e)$  مجموعه خوشه‌های مربوط به  $e$  است.  $|C(e) \cap C(e')|$  تعداد هم‌رخدادی  $e$  و  $e'$  در خوشه‌بندی و  $|L(e) \cap L(e')|$  تعداد هم‌رخدادی  $e$  و  $e'$  در خوشه‌بندی مرجع است. مقدار مانعیت گسترش یافته بی کیوبد برای یک مدرک برابر است با میانگین مانعیت چندگانه آن مدرک نسبت به مدارک با خوشه مرجع مشترک و مقدار نهایی مانعیت گسترش یافته بی کیوبد برابر با میانگین مانعیت گسترش یافته بی کیوبد برای هر مدرک است. به همین ترتیب، مقدار جامعیت گسترش یافته بی کیوبد نیز از فرایندی مشابه محاسبه می‌شود. توضیحات یادشده به زبان ریاضی در ادامه آمده است (Amigo et al. 2008):

$$\text{Precision BCubed} = \text{Avg}_e \left[ \text{Avg}_{e', C(e) \cap C(e') \neq \emptyset} [\text{Multiplicity precision}(e, e')] \right]$$

$$\text{Recall BCubed} = \text{Avg}_e \left[ \text{Avg}_{e', L(e) \cap L(e') \neq \emptyset} [\text{Multiplicity recall}(e, e')] \right]$$

1. multiplicity precision and recall

2. multiplicity precision

3. multiplicity recall

### ۳. پیشینه پژوهش

در میان آثار فارسی، اثری که در آن خوشه‌بندی پروانه‌های ثبت اختراع را با استفاده از عناصر استنادی مورد ارزیابی قرار داده باشد یافت نشد؛ با این حال، می‌توان به اثر داروغه (۱۳۸۰) اشاره کرد که در پایان‌نامه خود به بررسی میزان همپوشانی کلیدواژه‌های عنوان و عناوین مآخذ با توصیفگرهای نمایه‌سازی در پایان‌نامه‌های دکترای تخصصی روان‌پزشکی، زنان و زایمان، و قلب و عروق دانشگاه علوم پزشکی ایران پرداخته است. یافته‌های این پژوهش نشان داد که میزان همپوشانی سرجمع کلیدواژه‌های عنوان و عناوین مآخذ با توصیفگرهای نمایه‌سازی در هر رشته روان‌پزشکی، قلب و عروق، و زنان و زایمان بیش از میزان همپوشانی کلیدواژه‌های عنوان با توصیفگرهای نمایه‌سازی است.

کسلر طی مقاله‌ای مطالعه‌ای مقدماتی برای آزمون فرضیه اشتراک در مآخذ انجام داد. این مطالعه بر روی ۴۰ مدرک منتشرشده در ژوئن ۱۹۵۸ در مجموعه مقالات *آی.آر.ای*.<sup>۱</sup> با موضوع ترانزیستورها صورت گرفت. ارزیابی گروه‌های به‌وجودآمده از طریق بررسی موردی برخی مقالات و مقایسه موضوعی آنها با دیگر مقالات موجود در آن گروه، از سوی نویسنده صورت گرفته است. او در نهایت، نتیجه گرفت که می‌توان بدین شیوه به ایجاد گروه‌های کوچک‌تری از مقالات پرداخت به‌گونه‌ای که هر گروه زیرمجموعه درستی از گروه بزرگ‌تر باشد (Kessler 1963).

سالتون استفاده از استنادهای مآخذ در مقایسه با استفاده از واژگان در بازیابی مدارک مرتبط را مورد مقایسه قرار داده است. این پژوهش بر روی ۶۲ مدرک در حوزه زبان‌شناسی و ترجمه ماشینی صورت گرفته است. شواهد پژوهش به‌روشنی نشان نداده است که استنادها در نظام‌های بازیابی خودکار مدارک به‌طور معنی‌داری نقش مؤثری ایفا می‌کنند (Salton 1963).

سالتون عملکرد بازیابی را هنگام اضافه کردن اطلاعات استنادی به بردارهای تشکیل‌شده از واژه‌ها و مفاهیم و نیز استفاده از اطلاعات استنادی به‌تنهایی در ایجاد بردارها، نسبت به حالتی که از استناد استفاده نمی‌شود (استفاده از واژگان و مفاهیم)، مورد مقایسه قرار داده است. این پژوهش، بر روی چکیده ۲۰۰ مقاله در حوزه آنرویدینامیک به‌همراه ۴۲ پرسش اخذشده از پژوهشگران حوزه آنرویدینامیک صورت گرفته است. نتایج نشان داده است که اضافه کردن مفاهیم استنادی به توصیفگرهای موضوعی موجب افزایش میزان مانعیت تا بیش از ۱۰ درصد در آستانه‌ای از جامعیت خواهد شد. همچنین، خصیصه‌های استنادی عملکرد بسیار بهتری در سطوح پایین جامعیت (مانعیت محوری) دارند (Salton 1971).

1. Proceedings of the IRE



شاو به ارزیابی حضور ساختار خوشه‌ای به‌عنوان تابعی از جامعیت<sup>۱</sup> در مدارکی با موضوع فیروز کیستی<sup>۲</sup> پرداخته است. این ارزیابی برای چهار نوع بازنمایی موضوعی و دو نوع بازنمایی استنادی صورت گرفته است. نتایج نشان داد که برای تمامی انواع بازنمایی‌های مورد پژوهش، با کاهش جامعیت در بازنمایی، حضور ساختار خوشه‌ای کم خواهد شد. با وجود اینکه بسیاری از مدارک فاقد ارجاع یا استناد بوده‌اند، در بازنمایی استنادی حضور ساختار خوشه‌ای را در طیف گسترده‌تری از سطوح مختلف جامعیت نسبت به بازنمایی موضوعی نشان داد. هر دو نوع بازنمایی استنادی نشان‌دهنده شواهدی مبنی بر وجود ساختار خوشه‌ای در کمترین حالت جامعیت بود. ساختارهای اعمال‌شده در این مجموعه با استفاده از نمایه‌های موضوعی و استنادی، شرایط لازم را برای ایجاد خوشه‌بندی معنی‌دار فراهم کرده‌اند (Shaw 1990).

شاو میزان اثربخشی بازیابی مبتنی بر خوشه‌بندی را برای پنج نوع بازنمایی، مبتنی بر بازنمایی موضوعی و بازنمایی استنادی مورد بررسی قرار داده است. این آزمایش روی مجموعه‌ای شامل ۱۲۳۹ مدرک با موضوع فیروز کیستی صورت گرفت که با اصطلاح Cystic Fibrosis در کتابخانه ملی پزشکی (مدلاین) نمایه شده بودند. ترکیبات مورد بررسی شامل استفاده از سرعنوان موضوعی و منابع مورد استناد، سرعنوان موضوعی و منابع استنادکننده، سرعنوان موضوعی پزشکی و منابع مورد استناد، سرعنوان موضوعی پزشکی و منابع استنادکننده، و منابع مورد استناد و منابع استنادکننده بوده است. نتایج این پژوهش نشان داد که مقدار بهینه عملکرد تمامی ترکیب‌ها، قابل مقایسه و بیش از مقدار بهینه عملکرد در بازنمایی به‌صورت غیرترکیبی<sup>۳</sup> است. عملکرد بهینه، به‌طور ثابت، در سطوح پایینی از جامعیت محوری<sup>۴</sup> رخ داده است. در میان بازنمایی‌های ترکیبی جامعیت محور<sup>۵</sup>، مواردی که بازنمایی موضوعی را مورد استفاده قرار داده بوده‌اند، سطح عملکرد پایین‌تری داشته‌اند، نتایج بازیابی حاصل از ساختار تصادفی، قابل مقایسه با نتایج مشاهده‌شده بوده است. در میان بازنمایی‌های ترکیبی جامعیت محور، مواردی که منابع استنادکننده و استنادشونده را مورد استفاده قرار داده‌اند، نسبت به بازنمایی‌های حاصل از خصیصه‌های موضوعی، به‌طور قابل ملاحظه‌ای از اثربخشی بیشتری برخوردار بوده‌اند (Shaw 1991).

لای و وو در مقاله خود با عنوان «استفاده از رویکرد هم استنادی برای ایجاد نظامی نوین در رده‌بندی پروانه‌های ثبت اختراع» استفاده از رویکرد هم استنادی را به‌منظور رده‌بندی پروانه‌های ثبت اختراع پیشنهاد کرده‌اند. نتایج ارزیابی این مدل روی مجموعه‌ای آزمایشی شامل

- |  |                         |
|--|-------------------------|
| 1. Exhaustivity                        | 2. cystic fibrosis (CF) |
| 3. constituent representation          | 4. exhaustivity         |
| 5. exhaustive composite representation |                         |

پروانه‌های مرتبط با نانو تکنولوژی حاکی از مفید بودن اطلاعات استنادی در سازماندهی پروانه‌های ثبت اختراع بوده است (Lai and Wu 2005).

لی و همکاران در مقاله‌ای با عنوان «رده‌بندی خودکار پروانه‌های ثبت اختراع با استفاده از شبکه استنادی: مطالعه‌ای تجربی در فناوری نانو» استفاده از اطلاعات استنادی پروانه‌های ثبت اختراع به‌ویژه اطلاعات ساختار شبکه استنادی را در اختصاص خودکار رده‌ها به پروانه‌های ثبت اختراع مورد بررسی قرار داده‌اند. رویکرد مورد استفاده در پژوهش آنها مبتنی بر کرنل<sup>۱</sup> بوده و برای اخذ اطلاعات محتوایی و استنادی از تعریف توابع کرنل با استفاده از مجموعه‌ای از مدارک حاوی پروانه‌های ثبت اختراع در حوزه نانو تکنولوژی صورت گرفته است. نتایج این پژوهش حاکی از آن است که استفاده از ساختارهای شبکه استنادی به‌طور معنی‌داری عملکرد بهتری نسبت به حالتی ایجاد می‌کند که در آن از استناد استفاده نشده است (Li et al. 2007).

فوجی در مقاله‌ای با عنوان «تلفیق اطلاعات استنادی و محتوایی در بازیابی ثبت اختراع ان.تی.سی.آی.آر.۶» به ارزیابی بازیابی اطلاعات پروانه‌های ثبت اختراع با هدف رد اعتبار پروانه‌های ثبت اختراعات پرداخته است. مدل پیشنهادی او شامل تلفیق اطلاعات استنادی و متنی برای بازیابی پروانه‌های ثبت اختراع بوده است. در این مدل، ابتدا از بازیابی مبتنی بر متن و پس از آن، از اطلاعات استنادی استفاده است. در بازیابی مبتنی بر متن، از متن ادعای (های) موجود در پروانه‌های ثبت اختراع استفاده شده است. در نمره‌دهی استنادی مدارک دو شیوه در پیش گرفته شده است: حالتی که نمره استنادی در کل مدارک در نظر گرفته می‌شود و حالتی که نمره استنادی در N مدرک برتر حاصل از بازیابی مبتنی بر متن محاسبه می‌شود. در این مدل، به‌منظور نمره‌دهی ثبت اختراعات از نظر اطلاعات استنادی از روش پیچ‌رنک<sup>۲</sup> استفاده گردیده است. نتایج این پژوهش تجربی بر روی مجموعه آزمایشی پروانه‌های ثبت اختراع ان.تی.سی.آی.آر.۶ نشان داده است که تلفیق اطلاعات استنادی در حالت اول (محاسبه نمره استنادی بدون در نظر گرفتن N مدرک برتر) اثربخش‌تر از حالت دوم عمل کرده است. اما به‌طور کلی، تلفیق اطلاعات استنادی و متنی در بازیابی اثربخش بوده است (Fujii 2007).

تیوانا و هورویتز در «فایندسایت - یافتن خودکار پروانه‌های ثبت اختراع پیشین»، به ارائه الگوریتمی پرداخته‌اند که در آن از اسنادها برای یافتن پروانه‌های ثبت اختراع استفاده شده است. این الگوریتم بر روی نتایج بازیابی پیاده‌سازی می‌شود، به این معنی که ابتدا واژه‌ای جستجو و پس از آن الگوریتم بر روی مجموعه مدارک بازیابی شده، پردازش‌های بعدی را

1. Kernel

2. PageRank

انجام می‌دهد. ابتدا از طریق جستجوی کلیدواژه‌ای، مجموعه‌ای از پروانه‌های ثبت اختراع مشخص می‌شود. سپس، پروانه‌های ثبت اختراع مهم از طریق اسنادهای پیش‌آیند و پس‌آیند مشخص می‌گردد. استفاده از اسناد در الگوریتم پیشنهادی سبب بازیابی تعداد زیادی از منابع نامرتبط در کنار منابع مرتبط بوده است (Tiwana and Horowitz 2009).

زو و کرافت در مقاله خود با عنوان «تبدیل پروانه‌های ثبت اختراع به عبارت‌های جستجو برای آثار پیشین» به شناسایی بهترین ناحیه برای انتخاب خودکار عبارت جستجو پرداخته‌اند. نتایج پژوهش بر روی مجموعه حاوی پروانه‌های ثبت اختراع اداره یو.اس.بی.تی.<sup>۱</sup> حاکی از آن است که از میان نواحی عنوان، چکیده، خلاصه آثار پیشین<sup>۲</sup>، توصیف اشکال، شرح جزئیات، ادعاها، جستجوی آزاد و ادعاهای مقدم، ناحیه خلاصه آثار پیشین دارای بهترین تأثیر در بازیابی آثار پیشین بوده است. همچنین، در میان الگوریتم‌های مورد استفاده برای وزن‌دهی، الگوریتم tf دارای بهترین عملکرد و پس از آن، الگوریتم‌های tf-idf و bool قرار داشته‌اند (Xue and Croft 2009).

جستجو در آثار پیشین حاکی از آن است که موارد اندکی به بررسی نقش اسناد در سازماندهی پروانه‌های ثبت اختراع پرداخته‌اند و توجه آنها نیز بیشتر بر رده‌بندی خودکار یا بازیابی پروانه‌های ثبت اختراع بوده است. با توجه به آنکه اسناددهی در پروانه‌های ثبت اختراع از سوی متخصصان صورت می‌گیرد (United States Patent and Trademark Office (USPTO 2010)، احتمال وجود ارتباط میان اسنادها و اختراع ادعا شده بیشتر خواهد بود و در نتیجه به نظر می‌رسد که خوشه‌بندی پروانه‌های ثبت اختراع بر مبنای اسناد نتایج مطلوبی در برداشته باشد. به همین منظور، پژوهش حاضر بر آن است تا به ارزیابی خوشه‌بندی بر مبنای «اشتراک در مآخذ» و «واژگان عناوین اسنادها» پردازد.

#### ۴. روش پژوهش و ابزارهای گردآوری اطلاعات

در پژوهش حاضر که از نوع تجربی است، پس از استخراج و ذخیره اطلاعات در نرم‌افزار اکسل، ابزارهای زیر برای مراحل مختلف مورد استفاده قرار گرفت:

نرم‌افزار کد منبع باز ریپدماینر<sup>۳</sup> در وزن‌دهی به واژگان و تشکیل ماتریس‌های واژه-مدرک؛  
برنامه<sup>۴</sup> خوشه‌بندی فازی<sup>۵</sup> نوشته شده به زبان برنامه‌نویسی پرل<sup>۶</sup> برای ایجاد خوشه‌بندی فازی c

1. USPTO 2. background summary 3. <http://rapid-i.com/content/view/181/190/>  
4. <http://cpan.uwinnipeg.ca/htdocs/Algorithm-FuzzyCmeans/Algorithm/FuzzyCmeans.html>  
5. fuzzy c-means clustering 6. Perl

میانگین و برنامه<sup>۱</sup> ارزیابی کننده خوشه‌بندی براساس جامعیت و مانعیت گسترش یافته بی کیوبد که در کارگروه ویس<sup>۲</sup> معرفی و ارائه گردیده است. علاوه بر نرم‌افزارها و الگوریتم‌های پیش گفته، در تهیه نمودارها نرم‌افزار اکسل و اس.پی.اس.اس.<sup>۳</sup> مورد استفاده قرار گرفتند. جامعه آماری شامل ۷۱۷ پروانه ثبت اختراع یو.اس.<sup>۴</sup> مورد استناد در پروانه‌های ثبت اختراع یو.اس. متعلق به رده ۹۷۷/۷۷۴ (از رده‌بندی یو.اس.<sup>۵</sup>) بوده است. در رده‌بندی پروانه‌های ثبت اختراع یو.اس. رده ۹۷۷ رده‌ای اصلی و مربوط به موضوع نانو تکنولوژی است. ۹۷۷/۷۷۴ رده فرعی و زیرمجموعه‌ای<sup>۶</sup> از رده فرعی ۹۷۷/۷۷۳ است. رده ۹۷۷/۷۷۳ نیز زیرمجموعه‌ای از رده ۹۷۷/۷۰۰ است. رده ۹۷۷/۷۰۰ مربوط به نانو ساختارها و رده ۹۷۷/۷۷۳ شامل نانوذرات است (Class definition for class 9772011).

#### ۴-۱. نحوه ایجاد مجموعه مورد آزمایش<sup>۷</sup> با استفاده از تعیین خوشه‌بندی مرجع

ارزیابی خوشه‌بندی به‌طور معمول، با استفاده از مجموعه‌های آزمایشی صورت می‌گیرد. مجموعه آزمایشی، دربرگیرنده مدارکی است که نسبت به آنها قضاوت‌های ربط صورت گرفته است. در این مجموعه‌ها، موضوعاتی تعریف شده است که تعدادی از مدارک با آنها مرتبط است. خوشه‌بندی مرجع، گونه‌ای خوشه‌بندی است که در آن مدارک مرتبط به‌درستی در خوشه‌ها (و به‌عبارت دیگر، در کنار یکدیگر) قرار گرفته‌اند. بنابراین، پس از تطابق خوشه‌بندی پیشنهادی و خوشه‌بندی مرجع می‌توان به ارزیابی خوشه‌بندی پیشنهادی پرداخت. یکی از روش‌های ایجاد این مجموعه‌ها برای پروانه‌های ثبت اختراع، استفاده از روابط استنادی است. در این روش، هر پروانه ثبت اختراع مانند موضوعی است که استنادهای موجود در آن، مدارک مرتبط را مشخص می‌کند (Generation of a test collection based on citations [n.d.]). بنابراین، هر پروانه مانند خوشه‌ای حاوی مجموعه‌ای از مدارک مرتبط است. در این صورت، مجموعه‌ای از پروانه‌های ثبت اختراع خود می‌تواند به‌عنوان خوشه‌بندی مرجع در نظر گرفته شود.

در پژوهش حاضر، پروانه‌های ثبت اختراع رده ۹۷۷/۷۷۴ به‌عنوان خوشه‌بندی مرجع و استنادهای صورت گرفته در آنها به‌عنوان مجموعه مورد آزمایش در نظر گرفته شده است. همچنین، از آنجا که خوشه‌بندی پروانه‌های ثبت اختراع یو.اس. مورد نظر بوده است از میان منابع مورد استناد، فقط پروانه‌های ثبت اختراع یو.اس. در مجموعه قرار گرفته است.

1. www.evalita.it/sites/evalita.fbk.eu/files/doc2011/weps2007\_scorer\_1.1.tar.gz  
 2. WePS  
 3. SPSS  
 4. US patent  
 5. US patent classification  
 6. indented  
 7. test collection

برای ایجاد این مجموعه ابتدا تعداد ۱۰۴ پروانه ثبت اختراع مربوط به رده ۹۷۷/۷۷۴ به عنوان خوشه بندی مرجع برگزیده شد که برای به دست آوردن آنها عبارت ccl/977/774 در جستجوی پیشرفته پایگاه مربوط به پروانه های ثبت اختراع اداره یو.اس.پی.تی.ا. به نشانی <http://patft.uspto.gov/netahtml/PTO/search-adv.htm> مورد جستجو قرار گرفت. از میان ۱۰۴ خوشه (پروانه)، تعداد ۱۸ خوشه به دلیل یکسان بودن با پروانه های مورد استناد و تعداد ۱۱ خوشه نیز به دلیل عدم استناد به پروانه ثبت اختراع یو.اس.پی.تی.ا. از خوشه بندی مرجع حذف شد. در نهایت، تعداد ۷۵ خوشه به دست آمد. تعداد پروانه های ثبت اختراع یو.اس.پی.تی.ا. مورد استناد در این پروانه ها برابر با ۷۱۷ پروانه ثبت اختراع بوده است که در این پژوهش، به عنوان مجموعه آزمایشی مورد استفاده قرار گرفت.

#### ۲-۴. استخراج خصیصه ها

پس از مشخص شدن مجموعه آزمایشی، اقدام به استخراج «شماره ها» و «عناوین پروانه های ثبت اختراع مورد استناد» در هر پروانه ثبت اختراع (با استفاده از پایگاه های اطلاعاتی یو.اس.پی.تی.ا.، گوگل پتنتس<sup>۱</sup>، فری پتنتس آنلاین<sup>۲</sup>) و درج آنها در قالب جداولی در نرم افزار اکسل صورت گرفته است. در نهایت، دو جدول تشکیل شد: یک جدول برای ایجاد اشتراک در مآخذ و جدول دیگر برای استفاده از واژگان عناوین استنادها.

#### ۳-۴. ایجاد ماتریس خصیصه - پروانه ثبت اختراع

برای ایجاد ماتریس «استناد- پروانه ثبت اختراع» جدول مربوط به اشتراک در مآخذ به نرم افزار ریپیدماینر وارد شد و وزن دهی استنادها بر اساس روش  $tf-idf$  (2011) صورت گرفت. برای ایجاد ماتریس واژگان عناوین استنادی- پروانه ثبت اختراع نیز از جدول عناوین استنادی استفاده گردید. پیش از تشکیل این ماتریس به حذف واژگان موجود در سیاهه بازدارنده (با استفاده از سیاهه بازدارنده تعریف شده در این نرم افزار که با عنوان stopword (English) مشخص شده است) و ریشه یابی واژگان باقیمانده با استفاده از ریشه یاب پورتر<sup>۳</sup> (Porter 1980) موجود در نرم افزار ریپیدماینر پرداخته شد.

#### ۴-۴. خوشه بندی فازی پروانه های ثبت اختراع

پس از به دست آوردن خصیصه های وزن داده شده در نرم افزار ریپیدماینر، از برنامه

1. <http://www.google.com/patents>  
3. Porter stemmer

2. <http://www.freepatentsonline.com/search.html>

نوشته شده به زبان پرل برای ایجاد خوشه بندی فازی c میانگین استفاده شده است. از آنجا که خوشه بندی مرجع دارای ۷۵ خوشه بود، در تنظیمات خوشه بندی، تعداد خوشه برابر با ۷۵ خوشه (طی ۱۰۰ تکرار) تعیین شد.

خروجی این برنامه جدولی است که در آن میزان درجه عضویت هر پروانه ثبت اختراع را به هر یک از ۷۵ خوشه نشان می دهد.

جدول ۱. بخشی از ماتریس مربوط به عضویت فازی هر مدرک در هر خوشه براساس اشتراک استنادی

1	A	T1	T2	T3	T4	T5	T6
2	T2155658	0.0133	0.0133	0.0133	0.0133	0.0133	0.0133
3	T3203319	0.0133	0.0133	0.0133	0.0133	0.0133	0.0133
4	T3406228	0.0133	0.0133	0.0133	0.0133	0.0133	0.0133
5	T3410880	0.013	0.013	0.013	0.013	0.013	0.013
6	T3450673	0.0133	0.0133	0.0133	0.0133	0.0133	0.0133
7	T3485806	0.0133	0.0133	0.0133	0.0133	0.0133	0.0133
8	T3488327	0.0133	0.0133	0.0133	0.0133	0.0133	0.0133
9	T3488389	0.0133	0.0133	0.0133	0.0133	0.0133	0.0133
10	T3499032	0.0133	0.0133	0.0133	0.0133	0.0133	0.0133

جدول ۱ نشان دهنده عضویت های فازی پروانه های ثبت اختراع در خوشه بندی حاصل از اشتراک در مآخذ است. در ستون A، شماره پروانه های ثبت اختراع و در ستون های بعدی، نام هر خوشه قرار دارد.

جدول ۲. بخشی از ماتریس مربوط به عضویت فازی هر مدرک در هر خوشه براساس واژگان عناوین استنادی

1	A	T1	T2	T3	T4	T5	T6
2	T2155658	0.0133	0.0133	0.0133	0.0133	0.0133	0.0133
3	T3203319	0.0133	0.0133	0.0133	0.0133	0.0135	0.0133
4	T3406228	0.0133	0.0133	0.0131	0.0131	0.0132	0.0135
5	T3410880	0.013	0.0132	0.0131	0.013	0.013	0.0133
6	T3450673	0.0133	0.0133	0.0135	0.0133	0.0133	0.0134
7	T3485806	0.0128	0.0154	0.0127	0.0127	0.0128	0.0132
8	T3488327	0.0133	0.0133	0.0133	0.0133	0.0133	0.0133
9	T3488389	0.0133	0.0133	0.0133	0.0133	0.0133	0.0133
10	T3499032	0	0	0	0	0	0

در جدول ۲، عضویت های فازی پروانه های ثبت اختراع به هر خوشه در خوشه بندی مربوط به واژگان عناوین استنادها دیده می شود. در ستون A، شماره پروانه های ثبت اختراع و در ستون های بعدی، نام هر خوشه قرار گرفته است.

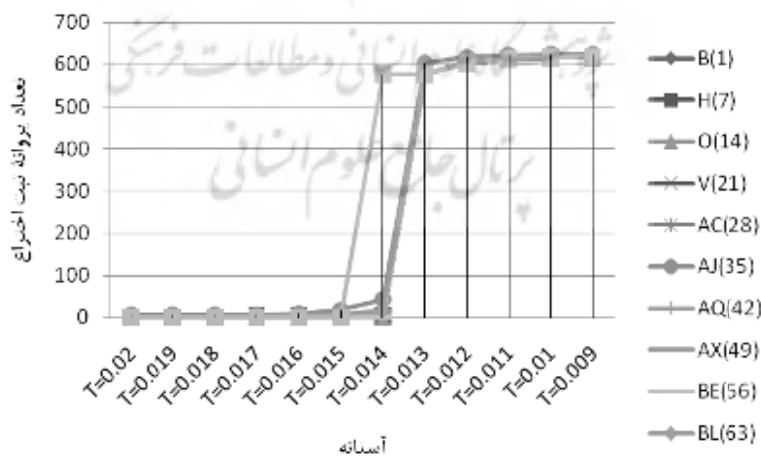
۵. یافته‌ها

آمار مربوط به درجه عضویت پروانه‌های ثبت اختراع در هر خوشه، در خوشه‌بندی با استفاده از اشتراک در مآخذ و واژگان عناوین استنادها در جدول ۳ آمده است.

جدول ۳. درجه عضویت در هر خوشه - اشتراک در مآخذ و واژگان عناوین استنادها

	درجه عضویت در خوشه (واژگان عناوین استنادها)	درجه عضویت در خوشه (اشتراک در مآخذ)
Min	0	0
Max	1	1
Average	0.013332	0.013314
Median	0.0128	0.0133
Mode	0	0.0133
Std	0.041732	0.04051

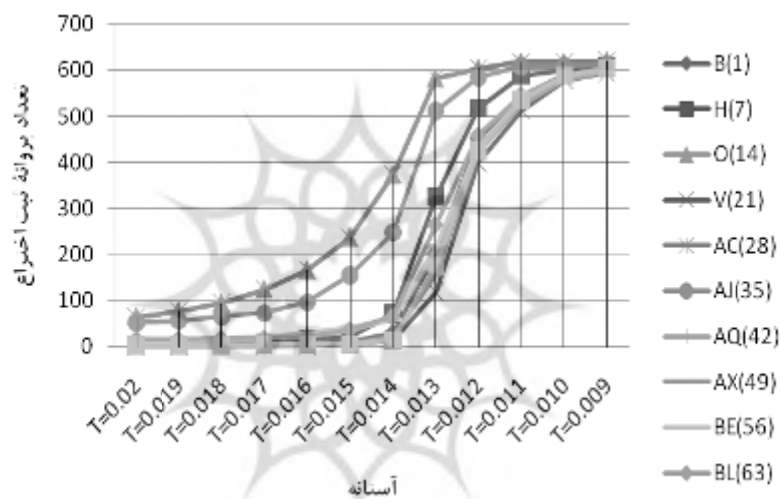
پس از مشخص شدن درجه‌های عضویت هر پروانه ثبت اختراع، در ۷۵ خوشه آستانه‌ای برای ارزیابی مانعیت محور و آستانه‌ای برای ارزیابی جامعیت محور تعیین شده است. این آستانه با استفاده از مطالعه مقدماتی بر روی تعدادی از خوشه‌ها تعیین شده است. برای تعیین آستانه مناسب، آستانه‌های موجود در بازه ۰/۰۲ تا ۰/۹۰۰ مورد بررسی قرار گرفته‌اند. این بازه به‌طور تجربی و با توجه به شاخص‌های مربوط به میانگین و میانه و نیز انحراف معیار در جدول ۵ و ۶ انتخاب گردید. نمودارهای مربوط به تغییرات آستانه و تغییرات تعداد اعضای هر خوشه برای تعداد خوشه انتخابی در ادامه ارائه شده است.



نمودار ۱. اشتراک در مآخذ- تغییرات تعداد اعضای هر خوشه با توجه به تغییرات آستانه عضویت

مطابق نمودار ۱ که مربوط به خوشه‌بندی مبتنی بر استاندارد (اشتراک در مآخذ) است، تغییرات تعداد اعضای هر خوشه، از آستانه ۰/۰۱۵ به ۰/۰۱۳ شدت گرفته (شیب منحنی) و پس از آستانه ۰/۰۱ تغییر چندانی نکرده است.

مطابق نمودار ۲ که مربوط به خوشه‌بندی با استفاده از عناوین استنادی است، تغییرات تعداد اعضای هر خوشه در این خوشه‌بندی در آستانه‌ای میان ۰/۰۱۲ و ۰/۰۱۴ شدت گرفته (شیب منحنی) و پس از آستانه ۰/۰۱ تغییر چندانی نکرده است.



نمودار ۲. واژگان عناوین استنادها- تغییرات اعضای هر خوشه با تغییر آستانه عضویت در هر خوشه

با توجه به نتایج حاصل از این مطالعه مقدماتی، در پژوهش حاضر در ارزیابی مانعیت محور از آستانه ۰/۰۱۴ و در ارزیابی جامعیت محور از آستانه ۰/۰۰۹ برای تعیین اعضای هر خوشه استفاده شده است و برای بررسی روند تغییرات از مانعیت محوری به سمت جامعیت محوری آستانه‌های ۰/۰۱۴، ۰/۰۱۳، ۰/۰۱۲، ۰/۰۱۱، ۰/۰۱۰، ۰/۰۰۹ و در تعیین اعضای هر خوشه مورد استفاده قرار گرفته است.

در میان آستانه‌های تعیین شده، ارزیابی صورت گرفته پس از اعمال آستانه ۰/۰۱۴ در خوشه‌بندی با استفاده از واژگان عناوین استنادها هیچ نتیجه‌ای دربر نداشت. این بدان معنی است که در این آستانه نقاطی که در خوشه‌بندی با یکدیگر هم‌رخداد شده‌اند، در خوشه‌بندی مرجع هم‌رخداد نشده‌اند و برعکس. این امر سبب می‌شود که مخرج کسر در فرمول مورد استفاده، هم



در محاسبه مانعیت گسترش یافته (به دلیل عدم هم‌رخدادی در خوشه‌بندی مرجع) و هم در محاسبه جامعیت گسترش یافته (به دلیل عدم هم‌رخدادی در خوشه‌بندی) صفر شود و کسری تعریف نشده (به لحاظ ریاضی) حاصل گردد. مقادیر جامعیت و مانعیت گسترش یافته بی‌کیوبد با تغییر آستانه‌ها در جدول ۴ و نمودار ۳ نمایان است.

**جدول ۴. مقادیر نهایی جامعیت و مانعیت گسترش یافته در خوشه‌بندی با استفاده از اشتراک در مآخذ و خوشه‌بندی با استفاده از واژگان عناوین استنادها در آستانه‌های اعمال شده**

آستانه	اشتراک در مآخذ		واژگان عناوین استنادها	
	م. گ. ب.	ج. گ. ب.	م. گ. ب.	ج. گ. ب.
0.014	1	0	-	-
0.013	1	0.7	1	0.26
0.012	1	0.73	1	0.52
0.011	1	0.75	1	0.63
0.010	1	0.75	1	0.70
0.009	1	0.76	1	0.73
میانگین	1	0.615	1	0.568
میانه	1	0.74	1	0.63

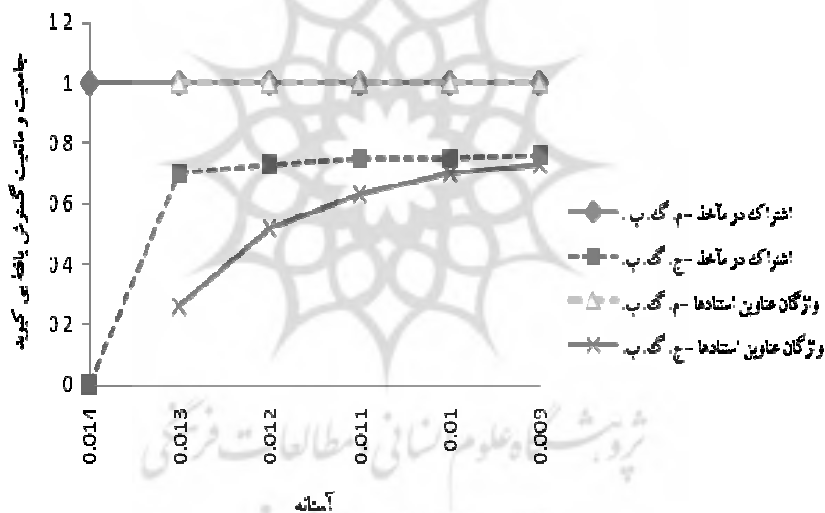
همان‌طور که جدول ۷ نشان می‌دهد، در آستانه ۰/۰۱۴ در خوشه‌بندی که براساس واژگان عناوین استنادی انجام شده است، جامعیت و مانعیت گسترش یافته نتیجه‌ای دربر نداشته است، بدین معنی که در این آستانه یکی از خوشه‌ها از هم پاشیده شده و شاید چند گره (نود) باقی مانده باشد. در همین آستانه، در خوشه‌بندی با استفاده از اشتراک استنادی مقدار جامعیت گسترش یافته صفر است، اما مانعیت آن عدد ۱ را نمایش می‌دهد. همچنین، به ترتیب هر چه از آستانه ۰/۰۱۴ به سمت پایین تر سیر می‌شود، یعنی در آستانه‌های ۰/۰۱۳، ۰/۰۱۲، ۰/۰۱۱، ۰/۰۱۰، و ۰/۰۰۹ مقدار جامعیت گسترش یافته در هر دو خوشه‌بندی افزایش می‌یابد.

در آستانه‌های تعیین شده، میانگین میزان مانعیت گسترش یافته بی‌کیوبد در خوشه‌بندی پروانه‌های ثبت اختراع با استفاده از اشتراک در مآخذ برابر با ۱ و میانگین میزان جامعیت آن برابر با ۰/۶۱۵ است (پاسخ به پرسش ۱).

میانگین میزان مانعیت گسترش یافته بی‌کیوبد در خوشه‌بندی پروانه‌های ثبت اختراع با

استفاده از واژگان عناوین استنادها برابر با ۱ و میانگین میزان جامعیت آن برابر با ۰/۵۶۸ است (پاسخ به پرسش ۲).

مطابق با جدول ۷، همان‌طور که مشهود است گرچه هرچه از آستانه تعریف‌شده کاسته می‌شود مقدار جامعیت گسترش یافته بی‌کیوبد در هر دو خوشه‌بندی افزایش می‌یابد، نرخ تغییرات در اشتراک استنادی کمتر از خوشه‌بندی واژگان عناوین استنادهاست. از سوی دیگر، مقایسه ستون‌های مربوط به جامعیت در هر دو خوشه‌بندی نشان می‌دهد که میزان جامعیت در اشتراک استنادی بیش از خوشه‌بندی با استفاده از واژگان عناوین استنادهاست، اما میزان مانعیت در هر دو برابر با یک است (پاسخ به پرسش ۳). نمودار ۳ نشان‌دهنده تغییرات جامعیت و مانعیت گسترش یافته برحسب تغییر در آستانه و حرکت از جامعیت محوری در خوشه‌بندی به سمت مانعیت محوری است.



نمودار ۳. تغییرات جامعیت و مانعیت گسترش یافته برحسب تغییرات آستانه

با توجه به نمودار ۳، هرچه از مانعیت محوری (آستانه ۰/۰۱۴) در خوشه‌بندی به سمت جامعیت محوری (آستانه ۰/۰۰۹) پیش می‌رویم، میزان جامعیت گسترش یافته بی‌کیوبد در هر دو خوشه‌بندی افزایش می‌یابد. با این حال، در تمامی آستانه‌های تعیین شده میزان مانعیت گسترش یافته مقداری ثابت و برابر با یک را به خود اختصاص داده است.

مقایسه منحنی جامعیت گسترش یافته (ج. گ. ب.) در خوشه‌بندی حاصل از اشتراک در مآخذ با منحنی خوشه‌بندی حاصل از واژگان عناوین استنادها نشان می‌دهد که در یک آستانه

واحد، میزان جامعیت گسترش یافته در خوشه‌بندی با استفاده از اشتراک در مآخذ مقدار بیشتری را نسبت به خوشه‌بندی با استفاده از واژگان عناوین اسنادها به خود اختصاص داده است. همچنین، یافته‌های نمودار ۳ حاکی از آن است که خوشه‌بندی با استفاده از اشتراک در مآخذ در طیف گسترده‌تری از آستانه قابل تعریف است (پاسخ به پرسش ۳). این بدان معنی است که ساختار خوشه‌ای در خوشه‌بندی با استفاده از اشتراک در مآخذ از انسجام بیشتری نسبت به استفاده از واژگان عناوین اسنادها برخوردار است.

## ۶. نتیجه‌گیری و بحث

یافته‌های پژوهش حاضر حاکی از آن است که خوشه‌بندی با استفاده از اشتراک در مآخذ در طیف گسترده‌تری از آستانه‌های تعیین شده (میزان جامعیت محوری) قابل تعریف است. این بدان معنی است که ساختار خوشه‌ای در خوشه‌بندی با استفاده از اشتراک در مآخذ از انسجام بیشتری نسبت به استفاده از واژگان عناوین اسنادها برخوردار است. پژوهش سالتون حاکی از عدم برتری خصیصه اسناد نسبت به خصیصه واژه در بازیابی مدارک بود (Salton 1963)، اما پژوهش حاضر حاکی از برتری خصیصه اسناد است. این امر ممکن است ناشی از شیوه ایجاد اسناد در پروانه‌های ثبت اختراع باشد. زیرا اسنادها از سوی بازرسان و در ارتباط با ادعاهای اختراع ایجاد می‌شوند و همین امر، احتمال حضور منابع مرتبط در میان اسنادها را افزایش می‌دهد. شاو نیز در پژوهش خود به این نتیجه رسید که بازنمایی اسنادی، حضور ساختار خوشه‌ای را در طیف گسترده‌تری از سطوح مختلف جامعیت نسبت به بازنمایی موضوعی نشان داده است و بازنمایی‌های اسنادی مورد استفاده در پژوهش او دارای ساختار خوشه‌ای در کمترین حالت جامعیت بوده است (Shaw 1990). سالتون نیز به این نتیجه رسید که اضافه کردن مفاهیم اسنادی به توصیفگرهای موضوعی موجب افزایش مانعیت تا بیش از ۱۰ درصد در آستانه‌ای از جامعیت خواهد شد. همچنین، خصیصه اسنادی عملکرد بسیار بهتری در سطوح پایین جامعیت محوری دارد (Salton 1971). شاو نیز به این نتیجه می‌رسد که در میان بازنمایی‌های ترکیبی جامعیت محور، مواردی که منابع اسنادکننده و اسنادشونده را مورد استفاده قرار داده‌اند نسبت به بازنمایی‌های حاصل از خصیصه‌های موضوعی، به‌طور قابل ملاحظه‌ای از اثربخشی بیشتری برخوردار بوده‌اند (Shaw 1991).

همچنین، نتایج پژوهش حاضر حاکی از برتری خصیصه اسناد به‌صورت یک واحد (به‌دلیل برتری خوشه‌بندی حاصل از اشتراک در مآخذ) نسبت به تجزیه آن به‌صورت خصیصه واژه است. این نتیجه تا حدودی همگام با نتایج پژوهش فوجی است که در آن نمره‌دهی

استنادی عملکرد بهتری نسبت به نمرده دهی برمبنای واژه داشته است (Fujii 2007). با این حال، تیوانا و هورویتز بیان داشته‌اند که استفاده از استناد در الگوریتم پیشنهادی آنها سبب بازیابی تعداد زیادی از منابع نامرتبط در کنار منابع مرتبط گشته است (Tiwana and Horowitz 2009). تجزیه استنادها به واژگان عناوین استنادها و اعمال پردازش‌هایی مانند حذف سیاه بازدارنده و ریشه‌یابی واژگان، منجر به کاهش تعداد خصیصه‌های مورد استفاده در عمل خوشه‌بندی خواهد شد. در نتایج به‌دست‌آمده تعداد خصیصه استناد (در اشتراک در مآخذ) حدود ۵۰۰۰ و تعداد خصیصه واژگان عناوین استنادها حدود ۳۰۰۰ واژه بوده است. این امر نشان از برتری واژگان عناوین استنادها نسبت به استفاده صرف از استنادها، در کاهش حجم مورد استفاده برای ذخیره خصیصه‌ها دارد. با توجه به اینکه میزان جامعیت در خوشه‌بندی با استفاده از واژگان عناوین استنادها در آستانه‌های جامعیت‌محور نزدیک به میزان جامعیت در اشتراک در مآخذ است، هرگاه جامعیت از ارزش کمتری برخوردار باشد، به‌نظر می‌رسد استفاده از این خصیصه گزینه‌ای مناسب در خوشه‌بندی پروانه‌های ثبت اختراع باشد. اما از سوی دیگر، جنبه حقوقی در ثبت اختراعات سبب شده است که نیازهای اطلاعاتی در جستجوی پروانه‌های ثبت اختراع نیازمند جامعیت بیشتر در نتایج جستجو باشد تا افرادی که در جهت بررسی قابلیت ثبت اثری به‌عنوان اختراع هستند بتوانند تمامی منابع مرتبط را مورد بررسی قرار دهند. از این رو، با توجه به نتایج پژوهش حاضر پیشنهاد می‌شود که روش اشتراک در مآخذ در طراحی نظام‌های بازیابی پروانه‌های ثبت اختراع مورد استفاده قرار گیرد.

## ۷. منابع

- داروغه، شیرین. ۱۳۸۰. بررسی میزان همپوشانی کلیدواژه‌های عنوان، عناوین مآخذ با توصیف‌گرهای نمایه‌سازی در پایان‌نامه‌های دکترای تخصصی روان‌پزشکی، زنان و زایمان و قلب و عروق دانشگاه علوم پزشکی ایران. پایان‌نامه کارشناسی ارشد، دانشکده مدیریت و اطلاع‌رسانی پزشکی، دانشگاه علوم پزشکی و خدمات بهداشتی درمانی ایران.
- Amigo, E., J. Gonzalo, J. Artiles, and F. Verdejo. 2008. A comparison of extrinsic clustering evaluation metrics based on formal constraints. *Information Retrieval* 12 (4): 461-486. <http://www.springerlink.com/content/812x06387152p045/fulltext.pdf> (accessed 3 Jun. 2011).
- Bagga, A. , and B. Baldwin. 1998. Algorithms for Scoring Coreference Chains. In *Proceedings of the Linguistic Coreference Workshop at the First International Conference on Language Resources and Evaluation (LREC'98)*, 563-566. <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.47.5848> (accessed 2 Feb. 2012).
- Bonino, D., A. Ciaramella, and F. Corno. 2010. Review of the state-of-the-art in patent information and forthcoming evolutions in intelligent patent informatics. *World Patent Information* 32 (1): 30-38.
- Class definition for class 977 - NANOTECHNOLOGY*. 2011. <http://www.uspto.gov/web/patents/classification/uspc977/defs977.htm#C977S774000> (accessed 04 Sep. 2011).

- Dhillon I., J. Kogan, and C. H. Nicholas. 2003. Feature selection and document clustering. In *A Comprehensive Survey of Text Mining*, 73-100. to be published by Springer-Verlag, M. Berry (Ed.). <http://www.cs.umbc.edu/csee/research/cadip/2002Symposium/kogan.pdf> (accessed 29 Jul. 2011).
- Fujii, A. 2007. Integrating content and citation information for the NTCIR-6 patent retrieval task. In *Proceedings of NTCIR-6 Workshop Meeting on Evaluation of Information Access Technologies: Information Retrieval, Question Answering and Cross-Lingual Information Access*, 377-380. Tokyo, Japan, 377-380. <http://research.nii.ac.jp/ntcir/ntcir-ws6/OnlineProceedings/NTCIR/76.pdf> (Accessed 22 April 2012).
- Generation of a test collection based on citations*. [n.d.]. <http://www.ir-facility.org/citation-based-test-collection> (accessed 04 Sep. 2011).
- Graf, E., and L. Azzopardi. 2008. A methodology for building a patent test collection for prior art search. In *Proceedings of the 2nd International Workshop on Evaluating Information Access (EVIA)*, December 16, 2008, Tokyo, Japan, 60-71. <http://research.nii.ac.jp/ntcir/workshop/OnlineProceedings7/pdf/EVIA2008/11-EVIA2008-GrafE.pdf> (accessed 27 Jan. 2010).
- Hideo, J., L. Azzopardi, and W. Vanderbauwhede. 2010. A survey of patent users. In *Proceeding of the third symposium on Information interaction in context, August 18-22*, 13-22. New Brunswick, New Jersey, USA: ACM. doi> 10.1145/1840784.1840789.
- Huang, A. 2008. Similarity measures for text document clustering. In *Proceedings of the Sixth New Zealand Computer Science Research Student Conference NZCSRSC 2008 Christchurch New Zealand*, 49-56. [http://nzcsrsc08.canterbury.ac.nz/site/proceedings/Individual\\_Papers/pg049\\_Similarity\\_Measures\\_for\\_Text\\_Document\\_Clustering.pdf](http://nzcsrsc08.canterbury.ac.nz/site/proceedings/Individual_Papers/pg049_Similarity_Measures_for_Text_Document_Clustering.pdf) (accessed 30 May 2011).
- Jain, A. K., and R. C. Dubes. 1988. Algorithms for clustering data. Prentice-Hall advanced reference series. NJ: Prentice-Hall, Inc., Upper Saddle River.
- Jain, A., M. Murty, and P. Flynn. 1999. Data clustering: A review, *ACM Computing Surveys* 31 (3):264 – 323. doi>10.1.1.18.2720.
- Kang, I., S. Na, J. Kim, and J. Lee. 2007. Cluster-based patent retrieval. *Information Processing and Management* 43 (5): 1173-1182.
- Kessler, M. M. 1963. An experimental study of bibliographic coupling between technical papers. *IEEE transactions on information theory* 9 (1): 49-51.
- Lai K. K., and S. H. J. Wu. 2005. Using the patent co-citation approach to establish a new patent classification system. *Information Processing and Management* 41 (2): 313-330.
- Larsen, B., and C. Aone. 1999. Fast and effective text mining using linear-time document clustering. In *Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 16-22. San Diego, California: ACM . doi>10.1145/312129.312186.
- Leydesdorff, L. 1987. Various Methods for the mapping of science. *Scientometrics* 11 (5-6): 295-324.
- Li, X., H. Chen, Z. Zhang, and J. Li. 2007. Automatic patent classification using citation network information: An experimental study in nanotechnology. In *JCDL'07: Proceedings of the 7th ACM/IEEE-CS joint conference on Digital libraries*, 419-427. New York, NY, USA. ACM. doi>10.1145/1255175.1255262.
- Porter, M. F. 1980. An algorithm for suffix stripping. *Program: Electronic Library and Information Systems* 14 (3): 130-137. [http://telemat.det.unifi.it/book/2001/wchange/download/stem\\_porter.html](http://telemat.det.unifi.it/book/2001/wchange/download/stem_porter.html) (accessed 20 Jun. 2011).
- Salton, G. 1963. Associative document retrieval techniques using bibliographic information. *Journal of the ACM* 10 (4): 440-457. [http://delivery.acm.org/10.1145/330000/321188/p440-salton.pdf?ip=217.218.83.9&CFID=36477701&CFTOKEN=92443976&\\_acm\\_=1313915650\\_9ed0bb61c9157adf2e126a72446f5244](http://delivery.acm.org/10.1145/330000/321188/p440-salton.pdf?ip=217.218.83.9&CFID=36477701&CFTOKEN=92443976&_acm_=1313915650_9ed0bb61c9157adf2e126a72446f5244) (accessed 22 Aug. 2011). doi> 10.1145/321186.321188.
- Salton, G. 1971. Automatic indexing using bibliographic citations. *Journal of documentation* 27 (2): 98-110.
- Shaw, W. M., Jr. 1990. Subject indexing and citation indexing. part I: Clustering structure in the cystic fibrosis document collection. *Information Processing and Management* 26 (6): 693-703.

- Shaw, W. M., Jr. 1991. Subject and citation indexing. Part II: The Optimal, cluster-based retrieval performance of composite representations. *Journal of The American Society For Information Science* 42 (9):676-664.
- tf-idf*. 2011. Wikipedia, the free encyclopedia. <http://en.wikipedia.org/wiki/Tf%E2%80%93idf> (accessed 1 Jun. 2011).
- Tan, P. N., M., Steinbach, and V. Kumar. 2006. Cluster analysis: Basic concept and algorithm. In *Introduction to Data Mining*. P. N. Tan, M. Steinbach and V. Kumar, 487-568. Boston, MA: Addison-Wesley Longman Publishing.
- Tiwana, S., and E. Horowitz. 2009. Find cite: automatically finding prior art patents. In *Proceedings of the 2nd international workshop on Patent information retrieval*, 37-40. [http://delivery.acm.org/10.1145/1660000/1651352/p37-tiwana.pdf?ip=217.218.83.9&CFID=39566988&CFTOKEN=46635715&\\_\\_acm\\_\\_=1315641819\\_3628c1ae1982ae3f1f1f9b5747a3f756](http://delivery.acm.org/10.1145/1660000/1651352/p37-tiwana.pdf?ip=217.218.83.9&CFID=39566988&CFTOKEN=46635715&__acm__=1315641819_3628c1ae1982ae3f1f1f9b5747a3f756) (accessed 3 Mar. 2011).
- United States Patent and Trademark Office (USPTO). 2010. Manual of Patent Examining Procedure (MPEP) (eighth edition). [http://www.uspto.gov/web/offices/pac/mpep/pdf\\_download\\_search\\_instructions.doc](http://www.uspto.gov/web/offices/pac/mpep/pdf_download_search_instructions.doc) (accessed 9 Jan. 2011).
- Wedding, D. K. 2009. Extending the Data Mining Software Packages SAS enterprise miner and spss clementine to handle fuzzy cluster membership: Implementation with examples (Master of Science Thesis in Data Mining), Central Connecticut State University. <http://web.ccsu.edu/datamining/Data%20Mining%20Theses/Don%20Wedding%20thesis.pdf> (accessed 13 May 2011).
- Xue, X., and W. B. Croft. 2009. Transforming patents into prior-art queries. In *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*, 808-809. Boston, MA, USA. ACM. doi>10.1145/1571941.1572139.

پژوهشگاه علوم انسانی و مطالعات فرهنگی  
پرتال جامع علوم انسانی

# Application of Bibliographic Coupling versus Cited Titles Words in Patent Fuzzy Clustering

**Anahita Kermani\***

Master in Library and Information Science

**Narges Neshat<sup>1</sup>**

Associated Professor in National Library and Archives of Iran

**Abbas Horri<sup>2</sup>**

Professor in Library and Information Science, Tehran University

Iranian Journal of  
**Information  
Processing &  
Management**

Iranian Research Institute  
For Science and Technology  
ISSN 2251-8223  
eISSN 2251-8231  
Indexed in LISA, SCOPUS & ISC  
Vol.28 | No.2 | pp: 411-432  
Winter 2013

**Abstract:** Attribute selection is one of the steps before patent clustering. Various attributes can be used for clustering. In this study, the effect of using citation and citation title words, respectively, in form of bibliographic coupling and citation title words sharing, were measured and compared with each other, as patent attributes. This study was done in an experimental method, on a collection of 717 US Patent cited in the patents belong to 977/774 subclass of US Patent Classification. Fuzzy C-means was used for patent clustering and extended BCubed precision and extended BCubed recall were used as evaluation measure. The results showed that the clustering produced by bibliographic coupling had better performance than clustering used citation title words and existence of cluster structure were in a wider range of exhaustivity than citation title words.

**Keywords:** patents, patent clustering, bibliographic coupling, citation title words, Fuzzy C-means clustering, clustering evaluation, BCubed Precision, BCubedRecall

\*Corresponding author: kermanianahita@gmail.com

1. narges\_neshat@yahoo.com

2. riwash@yahoo.com