

پیش‌بینی پیوند موضوعی پژوهش‌های علمی با استفاده از تحلیل شبکه‌ی مقالات علمی

رضا شکرچیان چالشتی

E-mail: rz.shekarchian@ut.ac.ir

گزارش ماهانه شماره (۳، ۴) - ۹۷/۸/۱۳

چکیده

یک راه برای تحلیل مجموعه‌ی اسناد دسته‌بندی نشده مدل‌سازی موضوعی می‌باشد. مدل‌سازی موضوعی می‌تواند کلماتی را که از نظر معنایی مشابه هستند کنار هم قرار دهد. با تحلیل اسناد مجموعه چندین موضوع به دست می‌آید. هر موضوع توزیعی بر روی کلمات می‌باشد و هر سند توزیعی از موضوعات می‌باشد. در مدل‌سازی موضوعی این توزیعات به دست می‌آیند و با مشخص شدن موضوعات می‌توان کلمات مشابه را در کنار هم دسته‌بندی کرد.

کلمات کلیدی: مدل‌سازی موضوعی

۱. معرفی LDA و PLSA

روش‌های مدل‌سازی احتمالاتی موضوعات، مجموعه‌ای از الگوریتم‌هایی هستند که هدف اصلی آن‌ها کشف ساختار نهان موضوعات در حجم وسیعی از اسناد می‌باشد. یکی از مطرح‌ترین و پایه‌ای‌ترین روش‌ها LDA^۱ [1] می‌باشد. LDA خود به منظور بهبود روش PLSA^۲ [2] پیشنهاد داده شد. این روش نیز خود بهبودی بر روش مشهور دیگری با نام LSA^۳ [3] می‌باشد.

LSA یکی از روش‌های پایه در مدل‌سازی موضوعی می‌باشد. ایده‌ی اصلی LSA این است که ماتریس سند-کلمه را به دو ماتریس سند-موضوع و موضوع-کلمه تبدیل کند. قدم اول ساخت ماتریس سند-کلمه می‌باشد. برای m سند و مجموعه کلمه‌ای به اندازه‌ی n کلمه می‌تواند ماتریسی $m \times n$ ساخت که در آن سطرها سندها و ستون آن کلمات می‌باشند.

یک روش ساده برای وزن‌دهی به درایه‌های ماتریس سند-کلمه تعداد دفعاتی می‌باشند که کلمه در سند ظاهر شده است. این روش اهمیت یک کلمه در سند را نادیده می‌گیرد. روش دیگر وزن‌دهی درایه‌های ماتریس با $tf \cdot idf$ کلمات می‌باشد. که به نوعی میزان اهمیت و یکتایی کلمات را در سند نشان می‌دهد. Tf تعداد دفعاتی است که کلمه در سند ظاهر شده و idf متناظر با عکس تعداد سندهایی در مجموعه است که کلمه در آن‌ها ظاهر شده است. این ماتریس تنک و با نویز بالا می‌باشد و بسیاری از درایه‌های آن اطلاعات مناسبی ندارند. به همین منظور به دنبال موضوعاتی می‌گردیم که ارتباط معنادارتری از کلمات و اسناد ارائه کنند LSA. از روش کاهش بعد ماتریس استفاده می‌کند و برای کاهش بعد از روش SVD بهره می‌گیرد. روش SVD ماتریس M را به ۳ ماتریس تبدیل می‌کند $M = U \cdot S \cdot V$:

^۱ Latent Dirichlet allocation

^۲ Probabilistic latent semantic analysis

^۳ Latent semantic analysis

که S ماتریس قطری مقادیر ویژه‌ی ماتریس M می‌باشد. از t مقدار ویژه‌ی اول برای ساخت S استفاده می‌شود که تعداد موضوعات را مشخص می‌کند U . ماتریس سند-موضوع و V ماتریس موضوع-کلمه می‌باشد. با استفاده از این ماتریس‌ها و روش‌های اندازه‌گیری‌ای همچون شباهت کسینوسی می‌توان شباهت اسناد، شباهت کلمات، شباهت کلمات با اسناد را اندازه‌گیری کرد.

از جمله مشکلات LSA نیاز آن به تعداد بالایی سند برای به دست آوردن نتایج دقیق می‌باشد. LDA یک مدل مولد^۱ می‌باشد. مدل مولد برای اسناد بر اساس یک سری قانون نمونه‌گیری احتمالاتی می‌باشد. این قانون‌ها مشخص می‌کنند کلمات اسناد چگونه ممکن است بر پایه‌ی متغیرهای نهان تولید شوند. پس از به دست آمدن مدل تولیدی مناسب، هدف پیدا کردن بهترین مجموعه از متغیرهای نهان می‌باشد که می‌تواند توصیف کننده‌ی مشاهدات باشند (مثلاً کلمات موجود در اسناد)، با فرض این که مدل به دست آمده داده‌ها را تولید کرده است. در واقع مجموعه‌ای از مشاهدات و اسناد موجود می‌باشد. می‌خواهیم مدل مولد را به گونه‌ای بسازیم که گویی این مدل تولید کننده‌ی مشاهدات بوده است. مدل‌های احتمالاتی موضوعی متنوعی وجود دارد. مانند [1], [2], [4]–[7]. ایده‌ی اصلی تمام این مدل‌ها یکسان می‌باشد. ایده این است که اسناد توزیعی از موضوعات می‌باشند. این مدل‌ها بیشتر در فرضیات آماری دارای تفاوت می‌باشند. فرض می‌کنیم $P(z)$ برای یک سند نشان دهنده‌ی توزیع بر روی تمام موضوعات z می‌باشد. $P(w|z)$ توزیع کلمات را بر روی موضوع z نشان می‌دهد و توزیع کلمه-موضوع^۲ معرفی می‌شود. به ازای هر سند تولید هر کلمه‌ی در دو مرحله صورت می‌گیرد. برای تولید هر کلمه‌ی w_i در یک سند، ابتدا یک نمونه‌گیری بر روی توزیع موضوعات صورت می‌گیرد و یک موضوع z انتخاب می‌شود. سپس یک کلمه از توزیع کلمه-موضوع $P(w|z)$ انتخاب می‌شود. از $P(z_i = j)$ برای نشان دادن احتمال انتخاب موضوع j برای کلمه‌ی i ام در نمونه‌گیری و از $P(w_i | z_i = j)$ برای احتمال کلمه‌ی w_i در موضوع j ام استفاده می‌شود. برای یک سند، احتمال تولید کلمات آن مطابق با فرمول ۱-۲ می‌باشد.

$$P(w_i) = \sum_{j=1}^T P(w_i | z_i = j) P(z_i = j) \quad (1-1)$$

که T در این جا تعداد موضوعات می‌باشد. فرمول ۱-۲ به طور خلاصه بیان می‌کند که احتمال تولید یک کلمه برای یک سند، برابر با احتمال تولید کلمه توسط موضوعات می‌باشد. به همین خاطر برای تمام موضوعات بررسی می‌کند سند چقدر به هر موضوع مرتبط است و برای هر موضوع، کلمه با چه احتمالی به موضوع تعلق می‌گیرد. این مدل از فرض تشکیل شدن هر سند از چندین موضوع استفاده می‌کند. نسبتی که هر سند از موضوعات دارد، با دیگر اسناد متفاوت است. یک سند ممکن است بیشتر راجع به دو موضوع اقتصاد و سیاست صحبت کرده باشد در حالی که سند دیگر ممکن است بیشتر راجع به سیاست و ورزش باشد. این خاصیت LDA می‌باشد که اسناد موجود، مجموعه‌ای یکسان از موضوعات را در بر می‌گیرند، ولی هر سند میزان تعلق متفاوتی به هر موضوع دارد. هدف مدل سازی موضوعی

^۱ generative model

^۲ Topic-word distribution

پیدا کردن موضوعاتی از روی مجموعه‌ی اسناد می‌باشد. اسناد به عنوان مشاهدات در نظر گرفته می‌شوند و در ابتدا وجود دارند. سه عنصر موضوعات، توزیع موضوعات بر روی هر سند و موضوعی که به هر کلمه‌ی سند نسبت داده می‌شود به عنوان دانش نهفته شناخته می‌شود که قرار است یادگیری شوند. همان‌طور که گفته شد مدلی که این دانش نهفته را به دست می‌آورد، یک مدل مولد می‌باشد. مدل مولد تلاش دارد تا پارامترهای مدل را به گونه‌ای به دست آورد که مشاهدات با بالاترین احتمال تولید شوند. هزینه‌ی اصلی مدل سازی موضوعی در قسمت استنتاج ساختار نهان موضوعات با استفاده از اسناد مشاهده شده می‌باشد. هدف استنتاج پیدا کردن توزیع شرطی متغیرهای نهان، مشروط به دیده شدن مشاهدات می‌باشد. به این احتمال شرطی، توزیع موخر^۱ نیز گفته می‌شود. می‌توان برای توصیف روش LDA از نمادگذاری زیر استفاده کرد. موضوعات $\beta_{1:k}$ می‌باشند که هر β_k توزیعی بر روی کلمات می‌باشد. k تعداد موضوعات می‌باشد. نسبتی که سند d از هر موضوع به خود اختصاص می‌دهد، با θ_d نمایش داده می‌شود و نسبت‌های موضوع^۲ خوانده می‌شود. $\theta_{d,k}$ نشان دهنده‌ی نسبتی است که موضوع k برای سند d دارد. تخصیص‌های موضوع^۳ برای سند d می‌باشند که $z_{d,n}$ تخصیص موضوع برای n امین کلمه در سند d می‌باشد. در نهایت کلمات دیده شده برای سند d ، می‌باشند که $w_{d,n}$ n امین کلمه برای سند d می‌باشد. با در نظر گرفتن نمادگذاری‌های گفته شده، فرآیند تولیدی برای LDA از توزیع توأم متغیرهای آشکار^۴ و متغیرهای نهان^۵ به صورت فرمول ۲-۲ می‌باشد.

$$p(\beta_{1:K}, \theta_{1:D}, z_{1:D}, w_{1:D}) = \prod_{i=1}^K p(\beta_i) \prod_{d=1}^D p(\theta_d) \left(\prod_{n=1}^N p(z_{d,n} | \theta_d) p(w_{d,n} | \beta_{1:K}, z_{d,n}) \right) \quad (1-2)$$

باید توجه داشت که در توزیع گفته شده یک سری وابستگی وجود دارد. به عنوان مثال تخصیص

موضوع $z_{d,n}$ وابسته به نسبت‌های موضوع θ_d برای هر موضوع می‌باشد. همین‌طور کلمه‌ی مشاهده شده $w_{d,n}$ وابسته به تخصیص موضوع $z_{d,n}$ و تمام موضوعات $\beta_{1:k}$ می‌باشد (عملاً کلمه با نگاه به موضوعی که $z_{d,n}$ به آن اشاره می‌کند و احتمالی که کلمه‌ی $w_{d,n}$ در آن موضوع دارد، به دست می‌آید). تمامی این وابستگی‌هاست که LDA را LDA کرده است. مدل‌های گرافیکی احتمالاتی^۶ این امکان را فراهم می‌کنند که برخی از توزیع‌های احتمالی را به زبان گرافیکی بیان کرد. شکل ۱ مدل گرافیکی برای LDA می‌باشد.

^۱ Posterior distribution

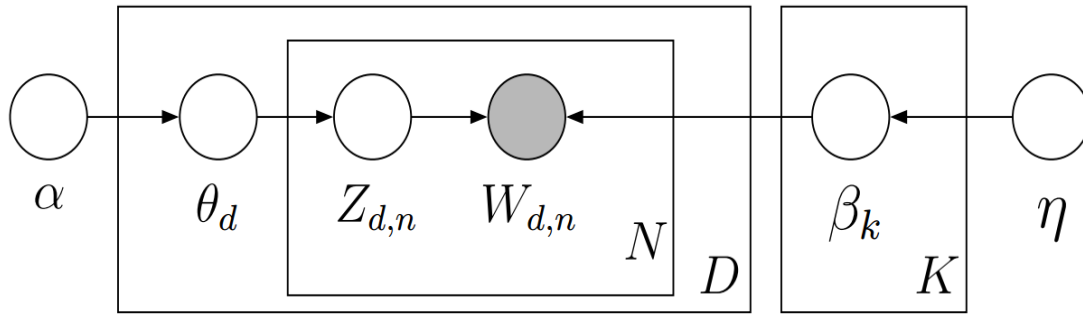
^۲ Topic proportions

^۳ Topic assignments

^۴ Observed variables

^۵ Hidden variables

^۶ Probabilistic graphical models



شکل (۱) مدل گرافیکی روش **LDA**. هر نود بیان کننده ی یک متغیر تصادفی می باشد و متناظر با نقشش در فرآیند مولد برچسب خورده است. نودهای نهان (نسبت های موضوع، تخصیص های موضوع و موضوع ها) به صورت ساده و نودهای آشکار (کلمات دیده شده در اسناد) به صورت سایه دار نمایش داده شده اند. مستطیل ها در شکل به عنوان صفحه^۱ شناخته می شوند و نشان دهنده ی تعداد تکرار می باشند. صفحه ی N نشان دهنده ی تعداد کلمات در هر سند و صفحه ی D نشان دهنده ی تعداد اسناد مجموعه می باشد.

LDA یکی از معتبرترین روش های موجود مدل سازی موضوعی اسناد می باشد. این روش خود بهبودی برای روش PLSA [2] می باشد. این روش بسیار شبیه به LDA می باشد. در PLSA متغیرهای پنهان موضوعات، $z_k \in \{z_1, \dots, z_K\}$ متناسب با تعداد اتفاق کلمات $w_j \in \{w_1, \dots, w_M\}$ در سند $d_i \in \{d_1, \dots, d_N\}$ می باشند. برای جفت سند و کلمه ی (d, w) داریم:

$$P(d_i, w_j) = P(d_i) \sum_{k=1}^K P(w_j | z_k) P(z_k | d_i) \quad (۱-۳)$$

$P(w_j | z_k)$ احتمال کلمه ی w_j در موضوع z_k می باشد و $P(z_k | d_i)$ احتمال موضوع z_k برای سند d_i می باشد. این پارامترها را می توان با بیشینه کردن لگاریتم راست نمای^۲ مجموعه ی C به دست آورد.

$$L(C) = \sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) \log \sum_{k=1}^K P(w_j | z_k) P(z_k | d_i) \quad (۱-۴)$$

پارامترهای مدل $\theta = \{P(z_k | d_i)\}$ و $\phi = \{P(w_j | z_k)\}$ می باشند که می توان آن ها را به استفاده از الگوریتم EM^۳ [8] تخمین زد.

مدل سازی موضوعی سابقه ی چندین ساله ای دارد. در ابتدا به موضوعات به صورت بسته ی لغات نگاه می شد. گسترش هایی بر این نگاه صورت گرفت. به دست آوردن موضوعات به صورت سلسله مراتبی به جای حالت تک لایه، به کارگیری موجودیت ها در کنار متن، استفاده از عبارات به جای لغات و استفاده از دانش داخلی و خارجی در مدل سازی از جمله ی این گسترش ها می باشند. هریک از این گسترش ها در ادامه بررسی خواهند شد.

^۱ Plate

^۲ Log likelihood

^۳ Expectation maximization

۲. استفاده از موجودیت ها، شبکه های اطلاعاتی^۱ و عبارات در گسترش مدل سازی موضوعی

روش های بسیاری برای بهبود LDA ارائه شدند. برخی از اطلاعات بیشتری از اسناد استفاده می کنند و مدل را پیچیده تر می کنند. به طور مثال در [9] از اطلاعات نویسندگان در کنار متن اسناد برای مدل سازی استفاده شده است. در مقاله ی [10] از موجودیت های موجود در اسناد برای بهبود مدل سازی کمک گرفته شده است. به طور مثال یک متن خبری می تواند در برگیرنده ی نام های اشخاص، محل ها، شرکت ها، مکان ها و ... باشد. این اطلاعات دانش بسیاری از سند در بر دارند و می توانند در مدل سازی بسیار کمک کننده باشند.

در برخی از روش ها از ارتباط میان اسناد شبکه های همگون و ناهمگونی ساخته اند و از این ارتباط برای مدل خود استفاده کرده اند. از جمله ی این روش ها می توان به [12], [11], [9] اشاره کرد. به طور مثال در [12] از اطلاعات مجله ای که مقاله در آن به چاپ رسیده است، نویسندگان مقاله و خود اسناد یک شبکه ی ناهمگون ساخته است. هر مقاله به یک مجله و چند نویسنده لینک دارد. هر نویسنده و مقاله نیز به چند مقاله لینک دارد و مرتبط است. از موضوعات به دست آمده برای هر سند برای مدل کردن موضوعات مربوط به نویسندگان و مجله ها استفاده می شود. از موضوعات به دست آمده برای مقالات و مجله ها نیز در مدل سازی موضوعی اسناد استفاده می شود. به نوعی در این شبکه موضوع هر نود در تاثیر پذیرفته و تاثیر گذار بر موضوعات همسایه ها می باشد و تاثیر موضوعات در شبکه ی ساخته شده از نودی به نود دیگر منتقل می شود. البته این مقاله از روش PLSA به عنوان مدل پایه برای مدل سازی موضوعی اسناد استفاده کرده است.

مقاله ی [11] از جمله تحقیقاتی است که بر روی مدل سازی موضوعات به صورت سلسله مراتبی کار کرده است. در این مقاله هدف ساختن سلسله مراتب موضوعات در شبکه ی اطلاعات ناهمگون می باشد. شبکه ی اطلاعات ناهمگون حاصل از اشیا از نوع های مختلفی می باشد که با یکدیگر در ارتباط می باشند. این نوع شبکه ی اطلاعات را می توان از بسیاری از منابع متداول همچون گزارشات شرکت های تجاری، مقالات منتشر شده ی علمی و رسانه های اجتماعی به دست آورد. ساختن سلسله مراتب مفاهیم با کیفیت بالا به منظور نمایش موضوعات در درشت دانه ی های مختلف می تواند در جستجو، بررسی و عبور از اطلاعات و شناسایی الگو کمک رسان باشد. در این مقاله الگوریتمی ارائه شده است که به صورت بازگشتی اقدام به ساخت سلسله مراتب موضوعات می کند. نوآوری و خلاقیتی که در مقاله به کار رفته این است که می تواند هم عبارات متنی و هم موجودیت های از نوع های مختلف را با یک روش ترکیبی از خوشه بندی و رتبه بندی بر روی شبکه داده های ناهمگون استفاده کند و سلسله مراتب موضوعات از نوع های مختلف را به دست آورد. در بسیاری از کارهایی که بر روی مدل های سلسله مراتبی موضوعات کار شده، استفاده از رابطه ی میان موجودیت ها با نوع های مختلف رایج نبوده است. در روشی که در این مقاله ارائه شده هم اطلاعات متن و هم ارتباطات میان موجودیت ها برای ساخت سلسله مراتب استفاده می شوند. نوآوری های به کار برده شده از این قرار می باشند. ساخت سلسله مراتب موضوعات به صورت بازگشتی به طوری که موضوعات به صورت لیست مرتبی از عبارات و موجودیت های از نوع های مختلف می باشند. از خوشه بندی و روش های رتبه بندی برای به دست آوردن زیر موضوع ها از موضوعات به صورت بازگشتی استفاده می شود.

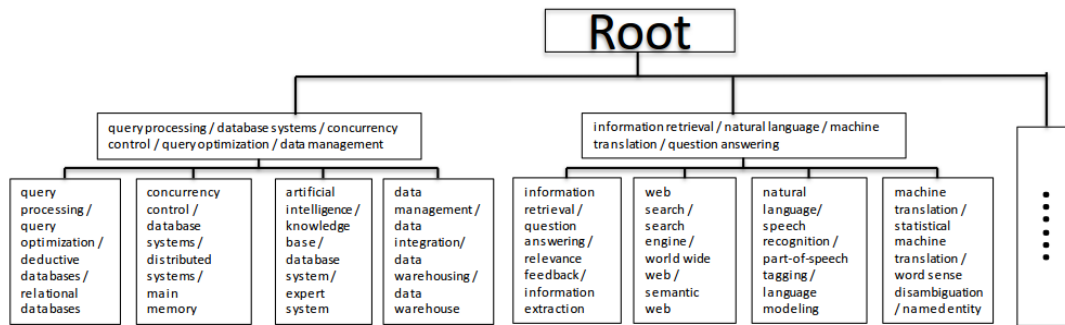
^۱ Information networks

این روش از مزایای روش [13] استفاده می‌کند که روشی برای تحلیل شبکه‌ها ناهمگون می‌باشد و برای به دست آوردن سلسله مراتب موضوعات بسیار مناسب می‌باشد. روشی که در این مقاله ارائه شده همچنین این قابلیت را دارد تا به صورت اتوماتیک اهمیت لینک‌های موجودیت‌ها با نوع‌های مختلف را مشخص کند. این مقاله به نوعی ادامه‌ی کار قبلی همین افراد است که روش [14] می‌باشد. این روش یک روش مبتنی بر آمار می‌باشد و سلسله مراتب موضوعات را تنها با تکیه بر محتوای متن انجام می‌دهد. در این روش هر موضوع به صورت سلسله مراتب و لیست مرتبی از عبارات با طول‌های مختلف نشان داده می‌شوند. این روش یک روش مبتنی بر عبارت می‌باشد و از این نگاه عبارتی به جای نگاه یک‌گرمی^۱ برای خوشه‌بندی و رتبه‌بندی موضوعات استفاده می‌کند. برای تخمین تکرار موضوعی عبارات از الگوریتم‌های کاوش الگوهای مکرر^۲ استفاده می‌کند و از این تکرار برای رتبه‌بندی موضوعی عبارات استفاده می‌کند. برای رتبه‌بندی عبارات موضوع، تابعی طراحی شده که چهار عامل را در کیفیت عبارات موضوع موثر می‌داند. پوشش^۳، خلوص^۴، کامل بودن^۵ و عبارت بودن^۶. پوشش بیان می‌کند که یک عبارت برای یک موضوع باید اسناد بسیاری را با آن موضوع پوشش دهد. به طور مثال «information retrieval» پوشش بهتری نسبت به «retrieval cross-language information» در موضوع بازیابی اطلاعات دارد. یک عبارت در یک موضوع خلوص دارد اگر تنها در اسناد با آن موضوع مکرر باشد و نه در اسنادی با موضوع دیگر. کامل بودن بیان می‌کند که یک موضوع کامل نباید زیر مجموعه‌ای از عبارت کامل بزرگ‌تر دیگری باشد. عبارت بودن هم نشان دهنده‌ی این است که کلمات عبارات باید به میزان کافی بیشتر از متوسط هم‌رخداد کلمات، با یکدیگر هم‌اتفاق باشند. در نهایت این روش از خوشه‌بندی بازگشتی برای ساخت سلسله مراتب استفاده می‌کند. سه شکل ۲ و ۳ و ۴ نتایج سه روش CATHY، NetClus و CATHYIN را نمایش می‌دهند. شکل‌ها از مقاله‌ی [11] می‌باشند. این سه روش مدل‌سازی موضوعی بر روی مقالات علوم کامپیوتر با استفاده از سه خصیصه‌ی کلمات، نویسندگان و مجله انجام شده‌اند.

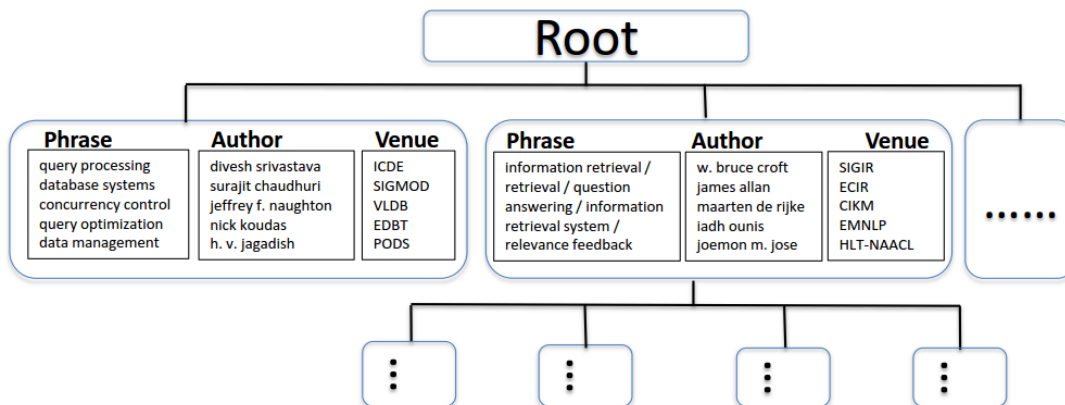
Term	Author	Venue	Term	Author	Venue	
data database queries system query	divesh srivastava jeffrey f. naughton christos faloutsos raghu ramakrishnan surajit chaudhuri	SIGMOD/ ICDE VLDB PODS CIKM	model retrieval learning information text	w. bruce croft chengxiang zhai james allan maarten de rijke c. lee giles	SIGIR ACL CIKM IJCAI AAAI

شکل ۲) روش مقاله‌ی [13]. خوشه‌های موجودیت‌های ناهمگون. هر مستطیل محیطی نشان دهنده‌ی یک خوشه می‌باشد. هر خوشه لیست مرتبی از کلمات و دو لیست مرتب از موجودیت‌ها می‌باشد.

^۱ unigram^۲ Frequent pattern mining^۳ coverage^۴ purity^۵ completeness^۶ phraseness



شکل ۳) روش مقاله‌ی [14]. سلسله مراتب موضوعی از متن‌ها. هر نود در شکل، لیست مرتبی از عبارات می‌باشد.



شکل ۴) روش مقاله [11]. سلسله مراتبی از موجودیت‌ها ناهمگون. هر نود در شکل لیست مرتبی از عبارات و دو لیست مرتب از موجودیت‌ها دارد.

دیده شد که چگونه از عبارات و موجودیت‌ها در مدل سازی و نمایش سازی موضوعات استفاده شد. از دیگر کارهایی که از عبارات برای مدل سازی موضوعی بهره گرفته، [15] می‌باشد. از قسمتی از ایده‌ی این مقاله که مربوط به استخراج عبارت‌ها^۱ می‌باشد، برای پیش برد این تحقیق بهره گرفته شده است. استفاده از عبارت‌ها برای بیان موضوعات بسیار به نظر بشر نزدیک‌تر از استفاده از کلمات می‌باشد و انسان بسیاری از مفاهیم را با دیدن چند کلمه در کنار هم متوجه می‌شود. بسیاری از مقالات از پس پردازش^۲ بر روی نتیجه‌ی مدل سازهای موضوعی مبتنی بر یک گرم استفاده می‌کنند تا آن‌ها را به صورت عبارت گونه نمایش دهند. برخی دیگر از کارها از الگوریتم‌های پیچیده‌ی مبتنی بر چند گرم^۳ بهره می‌گیرند. این روش‌ها یا کیفیت پایینی دارند و یا به خاطر پیچیدگی بسیار هزینه بر می‌باشند و مقیاس پذیر نمی‌باشند. روش این مقاله هم از نظر کیفیت و هم از نظر عملکرد بسیار کارا می‌باشد.

^۱ phrase

^۲ Post processing

^۳ N-gram

۳. استفاده از دانش در مدل سازی موضوعی

۳-۱. استفاده از دانش داخلی متن

مقاله‌ی [16] از جمله اولین کارهایی است که با تکیه بر متن و دانش برخاسته از متن، کیفیت مدل سازی موضوعی را بهبود داده است. هم چنین این مقاله روشی خودکار برای ارزیابی موضوعات پیشنهاد کرده است. فرمولی که در این مقاله برای ارزیابی ارائه شده بسیار شبیه به PMI^1 می باشد. میزان هم رخدادی^۲ دو به دوی M کلمه‌ی با بیشترین احتمال در هر موضوع موثر در خروجی معیار ارزیابی می باشد و انسجام موضوع^۳ را نشان می دهد. این مقاله روشی برای جلوگیری از ساخته شدن موضوعات ضعیف و با انسجام پایین دارد. با معیار ارائه شده، شناسایی موضوعات با کیفیت پایین، بدون نظر خبره و به صورت خودکار امکان پذیر شده است. در این مقاله، از هم رخدادی کلمات در مجموعه‌ی اسناد مستقیماً در مدل سازی موضوعها استفاده می کند. شاید در ابتدا این گونه به نظر آید که با اجازه ندادن به ظاهر شدن در یک موضوع به کلماتی که با هم دیگر هرگز هم رخداد نبوده اند، می توان به موضوعات با کیفیت دست یافت. ولی این طور نیست. طبق قانون توانی^۴، بیشتر کلمات به ندرت ظاهر می شوند و با بیشتر کلمات هم رخداد نیستند. حتی اگر به صورت معنایی با یکدیگر مرتبط باشند. در مدل ساز موضوعی ارائه شده تعداد رخداد کلمه‌ی w در موضوع t نه تنها احتمال دیده شدن دوباره‌ی کلمه را بالا می برد، بلکه باعث افزایش احتمال دیده شدن کلمات مرتبط نیز می شود. این روش جدید موضوعات اسناد را از طریق LDA به دست می آورد. با این تفاوت که مولفه‌ی مرsumی که در نمونه گیری موضوع-کلمه^۵ استفاده می شود و معروف به Polya urn می باشد را با Polya urn تعمیم یافته^۶ [17] جایگزین می کند. یک دنباله از نمونه های $i.i.d$ توزیع گسسته را می توان به دست آمده از تکرار انتخاب تصادفی یک توپ از ظرفی دانست. تعداد توپها از هر رنگ متناسب با احتمال آن توپ می باشد. پس از انتخاب هر توپ، جایگذاری انجام می شود. در Polya urn ساده هر توپ با توپی از همان رنگ جای گذاری می شود. هنگامی که یک توپ با رنگ مشخص از ظرفی برداشته می شود، توپ به همراه توپ جدیدی از همان رنگ به ظرف بازگردانده می شود. این گونه نمونه گیری باعث می شود احتمال انتخاب توپ رنگ w با هر بار انتخاب توپی با آن رنگ افزایش یابد. LDA از Polya urn ساده استفاده می کند. در روش Polya urn تعمیم یافته با انتخاب توپ رنگ w ، دو توپ از همان رنگ و A_{vw} توپ اضافی از هر رنگ $v \in \{1, \dots, W\}$ به ظرف بازگردانده می شوند. A_{vw} میزان تاثیر و ارتباط رنگها به یکدیگر را نشان می دهد. در مساله‌ی مدل سازی موضوعی کلمات در نقش رنگها و موضوعات در نقش ظرفها می باشند. توزیع کلمات یک موضوع متناسب با نسبت توپهای از هر رنگ در ظرف می باشد. در [16] از هم رخدادی

^۱ Point-wised Mutual Information

^۲ Co-occurrence

^۳ Topic coherence

^۴ Power-law

^۵ Topic-word

^۶ Generalized Polya urn

^۷ Independent and identically distributed

کلمات در ساخت ماتریس A_{vw} استفاده می شود.

۳-۲ استفاده از دانش دامنه های دیگر در مدل سازی موضوعی

یکی از روش هایی که از دانش دیگر دامنه ها برای مدل سازی مجموعه استفاده می کند روش AMC^1 [18] می باشد. یکی از ضعف های مدل سازی موضوعات کیفیت نامطلوب نتایج در صورت کم بودن تعداد اسناد مجموعه می باشد. الگو گرفتن از نحوه یادگیری انسان و نحوه تعامل آن با اطلاعات جدید می تواند در این زمینه کمک کننده باشد. یادگیری طولانی مدت 2 یکی از روش ها در یادگیری ماشین می باشد که AMC به نوعی از این الگو برای بهبود نتایج استفاده کرده است. انسان هنگام برخورد با یک رویداد جدید سعی می کند تا از آن چه در گذشته یاد گرفته است استفاده کند و نتیجه ی نحوه تعامل با رویداد جدید خود یک دانشی می شود برای رویدادهای بعدی. LDA و روش های گسترش داده ی آن از جمله روش های مر سوم و معروف برای مدل سازی موضوعی می باشند و نیاز به حجم زیادی از اسناد برای مدل سازی می باشند. مجموعه های سندی ای که شامل مقدار قابل توجهی سند باشند معمولاً اندک می باشند. به عنوان مثال نظرهای انجام گرفته بر روی کالاها را اگر در نظر بگیریم، اکثر کالاها بیشتر از ۱۰۰ نظر ندارند. اگر نظرهای مربوط به هر کالا را یک مجموعه در نظر بگیریم، با این تعداد از اسناد نمی توان موضوعات نهفته در نظرات را مدل سازی کرد. برای مقابله با این شرایط 3 روش می توان پیشنهاد داد.

- **ارائهی یک روش مدل سازی بهتر:** این روش زمانی می تواند تاثیر گذار باشد که تعداد زیادی از اسناد موجود باشد. چون که یادگیری مدل سازی موضوعی به صورت بدون ناظر 3 صورت می گیرد، اگر حجم اسناد کم باشد، اطلاعات کافی آماری قابل اطمینان وجود نخواهد داشت تا بتوان موضوعات منسجم به دست آورد. به نوعی نیاز به نظارت و یا اطلاعات خارج از اسناد داده شده نیاز می باشد.

- **درخواست از کاربر برای ارائهی دانش اولیه از دامنه:** یک نمونه از اطلاعات خارجی دانش اولیه کاربر از دامنه می باشد. به عنوان مثال کاربر می تواند دانش خود را به صورت ارتباط بایسته 4 و ارتباط نبایسته 5 وارد کند. ارتباط بایسته بیان می کند که دو کلمه باید در یک موضوع قرار بگیرند و ارتباط نبایسته بیان می کند که دو کلمه باید در دو موضوع متفاوت باشند. البته این نوع درخواست از کاربر می تواند مشکل زا باشد. چرا که مدل سازی غیر اتوماتیک می شود و همچنین کاربر نمی داند که چه دانشی برای چه قسمتی می تواند مفید واقع شود.

۴. ارزیابی کیفیت مدل ساز موضوع

^۱ Automatic Must-link Cannot-link

^۲ Lifelong learning

^۳ Unsupervised

^۴ Must-link

^۵ Cannot-link

روش‌های مدل‌سازی مجموعه‌ای از لغات که موضوع نام دارد را از مجموعه‌ی سندی استخراج می‌کنند. این کار به وسیله‌ی تکرارهای کلمات در سند انجام می‌شود. روش‌های ارزیابی انسجام^۱ موضوعات کمک می‌کنند تا موضوعات خوب از بد تمیز داده شوند و نشان می‌دهد که موضوعات به دست آمده تا چه حدی قابل فهم و معنادار می‌باشند. معیارهای ارزیابی انسجام بیشتر در حوزه‌ی پردازش زبان‌های طبیعی استفاده می‌شوند. در پردازش زبان‌های طبیعی موضوعات بیشتر به عنوان پیش‌پردازش اسناد مجموعه استفاده می‌شوند. موضوعات مجموعه و موضوعات هر سند، در خلاصه‌سازی اسناد [21]، روش‌های تشخیص معنای کلمات [22] و ترجمه‌ی ماشین [23] استفاده می‌شوند. در این تحقیق مسائل و دغدغه‌های مربوط به پردازش زبان مطرح نیست. در اینجا ملاک خوبی و بدی یک موضوع میزان نزدیکی آن به فهم بشری می‌باشد. چنگ و همکاران در [24] روشی پیشنهاد داده‌اند که مبتنی بر نفوذ^۲ کلمه می‌باشد. کلمات مزاحم^۳ به صورت تصادفی به موضوع‌ها وارد می‌شوند و از کاربران خواسته می‌شود که لغت نفوذی را پیدا کنند. روش نفوذ کلمات از این فرض استفاده می‌کند که کلمات مزاحم در یک موضوع منسجم بهتر قابل تشخیص می‌باشند. هرچه قدرت تشخیص بیشتر باشد، موضوع منسجم‌تر است. این روش به صورت غیر خودکار انجام می‌شود و این می‌تواند اصلی‌ترین نقطه‌ی ضعف این روش باشد. در روشی دیگر نیومن و همکاران [25] برای ارزیابی کیفیت موضوعات از کاربران خواستند که به موضوعات بر اساس کیفیت آن‌ها امتیاز دهند. در این روش هر موضوع با ده کلمه‌ی اول آن مشخص می‌شود. تلاش‌های بسیاری صورت گرفت تا روشی خودکار ارائه دهند که ارزیابی آن نزدیک به ارزیابی کاربران باشد. از ویکی‌پدیای انگلیسی و هم‌رخدادی کلمات موضوع در آن برای این ارزیابی خودکار استفاده شد. میمنو و همکاران [16] روش بسیار مشابهی ارائه کردند. تفاوت اصلی آن‌ها این است که به جای استفاده از اسناد ویکی‌پدیا برای نمونه‌گیری تعداد هم‌اتفاقی کلمات از اسناد خود مجموعه‌ی مدل شده برای نمونه‌گیری استفاده می‌کند. مشاهده شد که روش آن‌ها نزدیکی قابل توجهی به روشی که از امتیاز کاربران برای کیفیت‌سنجی موضوعات استفاده می‌کند، دارد. این مقاله علاوه بر ارائه‌ی یک معیار ارزیابی خودکار، روش جدیدی برای مدل‌سازی موضوعی ارائه کرده است که از هم‌رخدادی کلمات برای بهبود نتایج مدل‌سازی استفاده می‌کند. این روش خودکار نزدیکی قابل توجه به نظر کاربران در مورد کیفیت

Coherence^۱Intrusion^۲intruder^۳

موضوعات دارد. همین طور کیفیت به دست آمده از این روش نزدیکی قابل توجهی به روشی دارد که برای ارزیابی کیفیت موضوعات از نفوذ کلمات استفاده می کند.

۵. کاربردهای مدل سازی موضوعی

مدل سازهای موضوعی در بسیاری از زمینه ها وارد شده اند. از این میان می توان به مدل سازی زبان و انطباق مدل زبانی^۱ [27], [26], بازیابی اطلاعات [28] و [34]–[32], ابهام زدایی معنای کلمات [32], تحلیل شبکه های اجتماعی [34], [33] و کاوش نظرات [36], [35] اشاره کرد.

مقاله ی [37] به پیدا کردن انجمن ها^۲ و موضوعات با هم، در یک روش ترکیبی پرداخته است. بسیاری از اطلاعاتی که از کاربرها به دست می آید دارای متن می باشند. کاربرها با یکدیگر در ارتباط می باشند و در گراف روابط بین آنها یال وجود دارد. بسیاری از روش ها از لینک میان نودهای گراف برای مدل سازی موضوعی استفاده کرده اند. ولی در این مقاله شناسایی انجمن ها به کمک شناسایی موضوعات و شناسایی موضوعات به کمک شناسایی انجمن ها آمده اند. کاربرهایی که در یک انجمن قرار می گیرند موضوعات شبیه به هم زیادی دارند.

از دیگر کاربرهای مدل سازی موضوعی می توان به سیستم های توصیه گر اشاره کرد. مقاله ی [38] یکی از این کارها می باشد که یک توصیه گر برای مقالات علمی با کمک مدل سازی موضوعی طراحی کرده است. این مقاله روش های مرسوم توصیه گرها را با مدل سازی موضوعی ترکیب کرده و ساختار نهان کاربرها و موردها را به صورتی توصیف پذیر فراهم آورده است.

مقاله ی [39] تلفیقی از مدل سازی موضوعی و خلاصه سازی متن می باشد. موضوعات را به صورت سلسله مراتبی از مجموعه ی اسناد استخراج می کند. از الگوریتم های تحلیل شبکه استفاده می کند. از اسناد شبکه ای از تاثیر کلمات می سازد. کلمات خلاصه سازی که بیان کننده ی اصل موضوع هستند و بیشترین تاثیر را در شبکه دارند را می یابد. برای ساخت گراف تاثیر کلمات، کلمات نقش نودها را دارند. ارتباط بین دو کلمه بر اساس هم رخدادی آنها در کل مجموعه می باشد. گراف بیان کننده ی این است که اگر کلمه ی x دیده شد، چه میزان امکان دارد که به معنای کلمه ی y بدون دیدن آن برسیم. علاوه بر پیدا کردن کلمات موضوعی که بیشترین تاثیر را بر دیگر کلمات شبکه دارند،

^۱ Language model adaption

^۲ community

ارتباط میان موضوعات را نیز می یابد. کلمات خوشه بندی می شوند و سلسله مراتبی برای موضوعات ساخته می شود.

مدل سازی موضوعی می تواند به عنوان یک پایگاه دانش عمل کند و در کاربردهایی که مفاهیم جای کلمات را گرفته اند می تواند نقشی ایفا کند. این مفاهیم می توانند از پایگاه های دانشی هم چون وردنت و یا ویکی پدیا به دست آیند. مقاله ی [40] یکی از این مقالات می باشد. هم چنین می توان از موضوعات به دست آمده در ساخت پایگاه های دانش کمک گرفت. از پایگاه های دانش نیز می توان در ساخت موضوعات کمک گرفت و به جای نگاه کلمه ای به متن، نگاهی مفهومی با عناصر متن داشت و از این نگاه در مدل سازی موضوعی استفاده کرد. برخورد مفهومی با متن بیشتر در متون کوتاه حائز اهمیت می شوند. چرا که این متن ها کلمات کافی برای اعمال هویت خود ندارند. باید از همان چند کلمه و مفاهیم آن ها در تصمیم گیری ها (مثلا طبقه بندی) راجع به متن استفاده کرد. مقاله ی [41] از پایگاه دانشی احتمالاتی برای ساخت مفاهیم اسناد متن کوتاه استفاده می کند. پایگاه دانشی که استفاده می کند، Probase [42] و [43] نام دارد. این پایگاه دانش مفاهیم بسیار و غنی ای دارد که نزدیک به تصور بشری می باشند. پایگاه های دانشی هم چون وردنت [44]، ویکی پدیا [45] و Freebase [46] توسط خبره ها و نیروهای انسانی و یا تلاش های گروهی ساخته شده اند. تلاش های بسیاری صورت گرفته است که پایگاه های دانش به صورت خودکار ساخته شوند. KnowItAll [47]، TextRunner [48]، WikiTaxonomy [49] و YAGO [50] از جمله ی این تلاش ها می باشند. دو مشکل اصلی که بسیاری از پایگاه های دانش دارند، اندازه و جامعیت آن ها می باشد. به عبارتی پوشش و ریزدانگی مفاهیم را به خوبی ندارند. مثلا Freebase بعد از ۲۵ سال تلاش، ۲۰۰۰ دسته از مفاهیم دارد. Probase شامل میلیون ها مفهوم می باشد که دائما از میلیارد ها صفحه ی وب به دست می آید. روابطی که میان مفاهیم در Probase تعریف شده بسیار می باشد که شاید مهم ترین آن ها شباهت میان مفاهیم باشد.

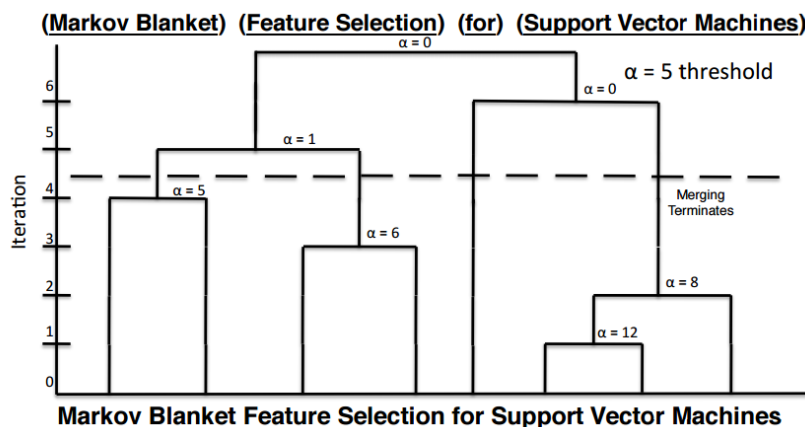
۶. نتیجه گیری

در این فصل نگاهی اجمالی به برخی از روش های موجود در مدل سازی موضوعی انداخته شد. تلاش ها و نگاه های مختلفی برای بهبود مدل سازی موضوعی وجود دارد. در این فصل از هر نگاهی نمونه ای معرفی شد که جدیدترین و بهترین نتایج اخیر را ارائه می کند. این نگاه از ساخت گراف از

اسناد و استفاده از آن برای مدل سازی تا استفاده از عبارتها، موجودیت، دانش های داخلی و خارجی و استفاده های مدل سازی موضوعی را شامل می شود.

یکی از راه های بهبود مدل سازی استفاده از چندگرم ها و عبارتها در مدل سازی موضوعی می باشد. یکی از کارهایی که از عبارات برای مدل سازی موضوعی بهره گرفته، [15] می باشد. از قسمتی از ایده ی این مقاله که مربوط به استخراج عبارتها می باشد، برای پیش برد این تحقیق بهره گرفته شده است. روش این مقاله هم از نظر کیفیت و هم از نظر عملکرد بسیار کارا می باشد و می تواند در مقیاس پذیری سیستم کمک کننده باشد. عبارتها در این مقاله چند ویژگی دارند. لیستی از عبارتها نشان دهنده ی یک موضوع منسجم هستند. عبارتهایی که به دست می آیند نزدیک به فهم بشری هستند. این روش موضوعات را با دقت و سرعت و پیچیدگی نزدیک به LDA به دست می آورد، با این تفاوت که به جای استفاده از یک گرم برای موضوعات، از عبارات برای مدل سازی و نمایش موضوعات استفاده می کند. این مقاله یک روش جدید کاوش عبارت برای تقسیم بندی متن به چندین عبارت یک یا چند کلمه ای و روشی جدید برای مدل سازی موضوعی بر روی بخش های ایجاد شده، ارائه کرده است. به دست آوردن عبارات در این مقاله به دو قسمت اصلی تقسیم می شود. پیدا کردن عبارات مکرر کاندید و تعداد تکرار آنها و در مرحله ی بعد ترکیب کلمات هر سند برای ساخت عبارات با کیفیت. قسمت کاوش عبارات مکرر را می توان در به دست آوردن تعداد تکرار کلمات پشت سر همی که از یک حداقلی بیشتر در کنار هم تکرار شده اند، خلاصه کرد. دو اصل می تواند در کارایی کاوش موثر باشد. اول

این که اگر عبارت P مکرر نباشد، هر عبارت بزرگتری که شامل P باشد نیز مکرر نیست. دوم هم این که اگر سندی هیچ عبارت مکرری به طول n نداشته باشد، آنگاه سند عبارت مکرری با طول بیشتر از n نخواهد داشت. نوآوری اصلی در کاوش عبارات، نحوه ساخت عبارات می باشد. عبارات از پایین به بالا ساخته می شوند. در هر چرخه به صورت حریصانه، مناسبترین جفت عبارت کاندید با هم ترکیب می شوند. این ترکیب از عبارات یک یا چند کلمه ای به دست می آید. از آن جا که تنها عباراتی که از تقسیم بندی سند به دست می آیند، معتبر می باشند، عباراتی که یک حداقل آستانه ای را دارند پذیرفته می شوند. شکل ۳-۸ نشان دهنده نحوه ساخت عبارات از پایین به بالا می باشد و برگرفته از مقاله ی [15] می باشد. تمام عبارات ترکیب شده مکرر می باشند.



ساخت پایین به بالا عبارات بر روی عنوان های مقالات علوم کامپیوتر

عبارت های پرتکرار معنادار به دست آمدند. این کلمات در کنار یکدیگر بسیار ظاهر شده اند. یک کلمه ای عبارت تکمیل کننده ی معنای کلمه ی همسایه می باشد. می توان مدعی شد قرار گرفتن کلمات یک عبارت در یک موضوع می تواند باعث بالا رفتن کیفیت موضوعات شود. این ایده مطرح می شود که از کلمات یک عبارت به عنوان ارتباط-بایسته^۱ در مدل سازی موضوعی به روش AMC استفاده شود. در این تحقیق از یک عبارت N کلمه ای، $N-1$ ارتباط-بایسته ساخته شد. قبلا ارتباط-بایسته ها از کلمات موضوعات به دست می آمدند. کلماتی که در موضوع های بسیاری با هم ظاهر می شوند این دانش را منتقل می کنند که این کلمات بهتر است در یک موضوع قرار گیرند. برای درک بهتر چگونگی عملکرد الگوریتم به بخش ۲-۴-۲ مراجعه شود.

دو اجرای متفاوت انجام شد. در هر دو اجرا مجموعه ی سندی به تعداد مشخصی خوشه شکسته می شود. از خوشه بندی موضوعی برای خوشه بندی استفاده می شود. برای مدل سازی موضوعی مجموعه از مدل سازی موضوعی خوشه ها استفاده می شود و موضوعات مجموعه در قالب موضوعات خوشه ها نمایش داده می شوند. در بخش ۳-۶-۱ نشان داده شد که برای استفاده از دانش خوشه ها چگونه می توان استفاده کرد تا یک خوشه مدل سازی موضوعی شود. باید خوشه ها در ابتدا با یک روش مدل سازی موضوعی، مدل سازی شوند. موضوعات خوشه ها به دست می آید. از موضوعات دانش در قالب ارتباط-بایسته ها و ارتباط-نبایسته ها استخراج می شوند. با کمک دانش به دست آمده از دیگر

^۱ Must-link

خوشه‌ها، خوشه‌ی مورد نظر مدل‌سازی موضوعی می‌شود. یک ایده در این‌جا استفاده از کلمات عبارت‌ها به عنوان ارتباط-بایسته در کنار ارتباط-بایسته‌های به دست آمده از موضوعات می‌باشد. در یک اجرا این ایده عملیاتی شد. فرض کنید قرار است خوشه‌ی N مدل‌سازی موضوعی شود. در فاز اول از موضوعات اولیه‌ی ساخته شده از تمام خوشه‌ها به جز خوشه‌ی N استفاده می‌شود و کلماتی که در کنار یکدیگر ارتباط-بایسته‌ها را شکل می‌دهند در یک فایل ثبت می‌شوند. در فاز بعد عبارت‌های پر تکرار و معنادار خوشه‌ی N با استفاده از روش مقاله‌ی [15] به دست می‌آیند. هر عبارت امتیازی دارد که متناسب با تکرار عبارت در مجموعه می‌باشد. عبارت‌ها بر اساس امتیازشان مرتب شده‌اند. N عبارت با امتیاز بیشتر برای استخراج ارتباط-بایسته استخراج می‌شوند. عدد N متناسب و تقریباً برابر با تعداد ارتباط-بایسته‌های استخراج شده از موضوعات خوشه‌ها می‌باشد. در نهایت از دو نوع ارتباط-بایسته استفاده شد و موضوعات جدید خوشه با استفاده از دانش حاصل از عبارت‌های خوشه و دانش حاصل از موضوعات دیگر خوشه‌ها به دست آمد.

فرض کنید قرار است موضوعات خوشه‌ی N با استفاده از دانش حاصل از ارتباط-بایسته‌ها به دست آورید. برای استخراج ارتباط-بایسته‌ها، موضوعات تمام دیگر خوشه‌ها باید به حافظه منتقل شوند تا دانش از آن‌ها استخراج شود. چرا که دانش، ارتباط دوتایی کلماتی است که در موضوعات خوشه‌ها بسیار با یکدیگر ظاهر شده‌اند. خوشه‌ی i به تمام دیگر خوشه‌ها وابسته می‌شود. این وابستگی مشکلاتی به همراه دارد. اول اینکه حافظه‌ی مصرفی را بالا می‌برد. چرا که اطلاعات موضوعات خوشه‌ها باید به حافظه منتقل شوند. هرچه تعداد موضوعات بیشتر باشد حجم اطلاعات موضوعات بیشتر می‌شود. نمونه‌ای از این اطلاعات، توزیع موضوعات برای هر سند می‌باشد. توزیع هر موضوع در یک سند عددی اعشاری می‌باشد. جمع احتمال توزیع‌ها ۱ می‌شود. چون تعداد موضوعات زیاد است، به غیر از اندک تعداد موضوعی که احتمال بالایی دارند، بقیه‌ی موضوع‌ها عددی کوچک می‌شوند که نیاز به اءشار بالا برای غیر صفر شدن آن‌ها ایجاد می‌شود. این وابستگی به تمام خوشه‌ها امکان اجرای هر خوشه بر روی یک ماشین را با مشکل مواجه می‌کند. این ایده مطرح می‌شود که جایگزینی برای استخراج ارتباط-بایسته‌ها یافت شود.

در اجرای دوم تلاش شد وابستگی یک خوشه از خوشه‌های دیگر برای استخراج دانش از بین برود. وابستگی در استخراج دانش بود. برای استخراج دانش برای یک خوشه باید موضوعات تمام دیگر خوشه‌ها در حافظه بارگزاری شوند. البته می‌توان با ایجاد محدودیت در انتخاب خوشه‌های مشابه تعداد خوشه‌هایی که برای استخراج دانش استفاده می‌شوند را کاهش داد. ولی باز هم حجم موضوعات هر خوشه می‌تواند یک سر بار باشد. تصمیم گرفته شد در اجرای دوم به جای استفاده از دانش حاصل از عبارت‌ها در کنار دانش حاصل از موضوعات دیگر خوشه‌ها برای بهبود مدل‌سازی موضوعی، تنها از دانش عبارت‌ها استفاده شود. با استفاده از این سیاست دانش هر خوشه وابسته به خود خوشه می‌شود. از عبارت‌های پرتکراری که در سندهای خوشه وجود دارند به عنوان دانش استفاده می‌شود. هر خوشه می‌تواند مستقل از دیگر خوشه‌ها مدل‌سازی موضوعی شود. این سیاست باعث می‌شود در مصرف منابعی همچون حافظه و پردازشگر صرفه‌جویی شود. هم‌چنین می‌توان خوشه‌ها را به صورت همزمان مدل‌سازی موضوعی کرد. در قسمت نتایج دقت دو اجرا مقایسه شده‌اند. اجرایی که دانش از عبارت‌ها در کنار دانش موضوعات استفاده می‌شود و زمانی که دانش تنها از دانش عبارت‌های خوشه می‌باشد.

- [1] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet Allocation," *J. Mach. Learn. Res.*, vol. 3, pp. 993–1022, Mar. 2003.
- [2] T. Hofmann, "Probabilistic Latent Semantic Indexing," in *Proceedings of the 22Nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, 1999, pp. 50–57.
- [3] S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, and R. Harshman, "Indexing by latent semantic analysis," *J. Am. Soc. Inf. Sci.*, vol. 41, no. 6, pp. 391–407, 1990.
- [4] M. Rosen-Zvi, T. Griffiths, M. Steyvers, and P. Smyth, "The Author-topic Model for Authors and Documents," in *Proceedings of the 20th Conference on Uncertainty in Artificial Intelligence*, 2004, pp. 487–494.
- [5] T. L. Griffiths and M. Steyvers, "A Probabilistic Approach to Semantic Representation." 2002.
- [6] T. L. Griffiths and M. Steyvers, "Prediction and Semantic Association," in *Advances in Neural Information Processing Systems*, 2003, p. 15.
- [7] D. Cohn and T. Hofmann, "The Missing Link - A Probabilistic Model of Document Content and Hypertext Connectivity." 2001.
- [8] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *J. R. Stat. Soc. Ser. B*, vol. 39, no. 1, pp. 1–38, 1977.
- [9] Y. Sun, J. Han, J. Gao, and Y. Yu, "iTopicModel: Information Network-Integrated Topic Modeling," in *Proceedings of the 2009 Ninth IEEE International Conference on Data Mining*, 2009, pp. 493–502.
- [10] H. Kim, Y. Sun, J. Hockenmaier, and J. Han, "ETM: Entity Topic Models for Mining Documents Associated with Entities.," *ICDM*, 2012.
- [11] C. Wang, M. Danilevsky, J. Liu, N. Desai, H. Ji, and J. Han, "Constructing Topical Hierarchies in Heterogeneous Information Networks," in *Data Mining (ICDM), 2013 IEEE 13th International Conference on*, 2013, pp. 767–776.
- [12] H. Deng, J. Han, B. Zhao, Y. Yu, and C. X. Lin, "Probabilistic Topic Models with Biased Propagation on Heterogeneous Information Networks," pp. 1271–1279, 2011.
- [13] Y. Sun, Y. Yu, and J. Han, "Ranking-based Clustering of Heterogeneous Information Networks with Star Network Schema," in *Proceedings of the 15th ACM SIGKDD*

- International Conference on Knowledge Discovery and Data Mining*, 2009, pp. 797–806.
- [14] C. Wang, M. Danilevsky, N. Desai, Y. Zhang, P. Nguyen, T. Taula, and J. Han, “A Phrase Mining Framework for Recursive Construction of a Topical Hierarchy,” in *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2013, pp. 437–445.
 - [15] A. El-Kishky, Y. Song, C. Wang, C. R. Voss, and J. Han, “Scalable Topical Phrase Mining from Text Corpora,” *CoRR*, vol. abs/1406.6, 2014.
 - [16] D. Mimno, H. M. Wallach, E. Talley, M. Leenders, and A. McCallum, “Optimizing Semantic Coherence in Topic Models,” in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2011, pp. 262–272.
 - [17] H. Mahmoud, *Polya Urn Models*, 1st ed. Chapman & Hall/CRC, 2008.
 - [18] Z. Chen and B. Liu, “Mining Topics in Documents: Standing on the Shoulders of Big Data,” in *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2014, pp. 1116–1125.
 - [19] Z. Chen, A. Mukherjee, B. Liu, M. Hsu, M. Castellanos, and R. Ghosh, “Discovering Coherent Topics Using General Knowledge,” in *Proceedings of the 22Nd ACM International Conference on Conference on Information & Knowledge Management*, 2013, pp. 209–218.
 - [20] Z. Chen, A. Mukherjee, B. Liu, M. Hsu, M. Castellanos, and R. Ghosh, “Leveraging Multi-domain Prior Knowledge in Topic Models,” in *Proceedings of the Twenty-Third International Joint Conference on Artificial Intelligence*, 2013, pp. 2071–2077.
 - [21] A. Haghighi and L. Vanderwende, “Exploring Content Models for Multi-document Summarization,” in *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, 2009, pp. 362–370.
 - [22] J. H. Lau, P. Cook, D. McCarthy, D. Newman, and T. Baldwin, “Word Sense Induction for Novel Sense Detection,” in *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, 2012, pp. 591–601.
 - [23] B. Zhao and E. P. Xing, “HM-BiTAM: Bilingual Topic Exploration, Word Alignment, and Translation,” in *Advances in Neural Information Processing Systems 20*, J. c. Platt, D. Koller, Y. Singer, and S. Roweis, Eds. Cambridge, MA: MIT Press, 2007, pp. 1689–1696.
 - [24] J. Chang, S. Gerrish, C. Wang, J. L. Boyd-graber, and D. M. Blei, “Reading Tea Leaves: How Humans Interpret Topic Models,” in *Advances in Neural Information Processing Systems 22*, Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta, Eds. 2009, pp. 288–296.

- [25] D. Newman, J. H. Lau, K. Grieser, and T. Baldwin, "Automatic Evaluation of Topic Coherence," in *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, 2010, pp. 100–108.
- [26] J. T. Chien and M. S. Wu, "Adaptive Bayesian Latent Semantic Analysis," *IEEE Trans. Audio. Speech. Lang. Processing*, vol. 16, no. 1, pp. 198–207, Jan. 2008.
- [27] M. S. Wu, H. S. Lee, and H. M. Wang, "Exploiting semantic associative information in topic modeling," in *Spoken Language Technology Workshop (SLT), 2010 IEEE*, 2010, pp. 384–388.
- [28] H. M. Wallach, "Topic Modeling: Beyond Bag-of-words," in *Proceedings of the 23rd International Conference on Machine Learning*, 2006, pp. 977–984.
- [29] X. Wang, A. McCallum, and X. Wei, "Topical N-Grams: Phrase and Topic Discovery, with an Application to Information Retrieval," in *Proceedings of the 2007 Seventh IEEE International Conference on Data Mining*, 2007, pp. 697–702.
- [30] X. Wei and W. B. Croft, "LDA-based Document Models for Ad-hoc Retrieval," in *Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2006, pp. 178–185.
- [31] A. Kotov, V. Rakesh, E. Agichtein, and C. K. Reddy, "Advances in Information Retrieval: 37th European Conference on IR Research, ECIR 2015, Vienna, Austria, March 29 - April 2, 2015. Proceedings," A. Hanbury, G. Kazai, A. Rauber, and N. Fuhr, Eds. Cham: Springer International Publishing, 2015, pp. 635–647.
- [32] L. Li, B. Roth, and C. Sporleder, "Topic Models for Word Sense Disambiguation and Token-based Idiom Detection," in *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, 2010, pp. 1138–1147.
- [33] H. Deng, J. Han, B. Zhao, Y. Yu, and C. X. Lin, "Probabilistic Topic Models with Biased Propagation on Heterogeneous Information Networks," in *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2011, pp. 1271–1279.
- [34] Q. Mei, D. Cai, D. Zhang, and C. Zhai, "Topic Modeling with Network Regularization," in *Proceedings of the 17th International Conference on World Wide Web*, 2008, pp. 101–110.
- [35] M. H. Alam and S. Lee, "Semantic Aspect Discovery for Online Reviews," in *Data Mining (ICDM), 2012 IEEE 12th International Conference on*, 2012, pp. 816–821.
- [36] F. Xianghua, L. Guo, G. Yanyan, and W. Zhiqiang, "Multi-aspect Sentiment Analysis for Chinese Online Social Reviews Based on Topic Modeling and HowNet Lexicon," *Know.-Based Syst.*, vol. 37, pp. 186–195, Jan. 2013.

- [37] Z. Yin, L. Cao, Q. Gu, and J. Han, “Latent Community Topic Analysis: Integration of Community Discovery with Topic Modeling,” *ACM Trans. Intell. Syst. Technol.*, vol. 3, no. 4, pp. 63:1–63:21, Sep. 2012.
- [38] C. Wang and D. M. Blei, “Collaborative Topic Modeling for Recommending Scientific Articles,” in *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2011, pp. 448–456.
- [39] C. Wang, X. Yu, Y. Li, C. Zhai, and J. Han, “Content Coverage Maximization on Word Networks for Hierarchical Topic Summarization,” in *Proceedings of the 22Nd ACM International Conference on Information & Knowledge Management*, 2013, pp. 249–258.
- [40] F. Wang, Z. Wang, Z. Li, and J.-R. Wen, “Concept-based Short Text Classification and Ranking,” in *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*, 2014, pp. 1069–1078.
- [41] Y. Song, H. Wang, Z. Wang, H. Li, and W. Chen, “Short Text Conceptualization Using a Probabilistic Knowledgebase,” in *Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence - Volume Volume Three*, 2011, pp. 2330–2336.
- [42] W. Wu, H. Li, H. Wang, and K. Zhu, “Towards a Probabilistic Taxonomy of Many Concepts,” Microsoft Technical Report, Mar. 2011.
- [43] W. Wu, H. Li, H. Wang, and K. Q. Zhu, “Probase: A Probabilistic Taxonomy for Text Understanding,” in *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*, 2012, pp. 481–492.
- [44] G. A. Miller, “WordNet: A Lexical Database for English,” *Commun. ACM*, vol. 38, no. 11, pp. 39–41, Nov. 1995.
- [45] D. B. Lenat and R. V Guha, *Building Large Knowledge-Based Systems; Representation and Inference in the Cyc Project*, 1st ed. Boston, MA, USA: Addison-Wesley Longman Publishing Co., Inc., 1989.
- [46] K. Bollacker, C. Evans, P. Paritosh, T. Sturge, and J. Taylor, “Freebase: A Collaboratively Created Graph Database for Structuring Human Knowledge,” in *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, 2008, pp. 1247–1250.
- [47] O. Etzioni, M. Cafarella, D. Downey, S. Kok, A.-M. Popescu, T. Shaked, S. Soderland, D. S. Weld, and A. Yates, “Web-scale Information Extraction in Knowitall: (Preliminary Results),” in *Proceedings of the 13th International Conference on World Wide Web*, 2004, pp. 100–110.
- [48] O. Etzioni, M. Banko, S. Soderland, and D. S. Weld, “Open Information Extraction from the Web,” *Commun. ACM*, vol. 51, no. 12, pp. 68–74, Dec. 2008.

- [49] S. P. Ponzetto and M. Strube, “Deriving a Large Scale Taxonomy from Wikipedia,” in *Proceedings of the 22Nd National Conference on Artificial Intelligence - Volume 2*, 2007, pp. 1440–1445.
- [50] F. M. Suchanek, G. Kasneci, and G. Weikum, “Yago: A Core of Semantic Knowledge,” in *Proceedings of the 16th International Conference on World Wide Web*, 2007, pp. 697–706.