

رتبه بندی لغات و آنتروپی کسری

مهری دهنوی، حسین^۱؛ آگاهی، حمزه^۲؛ مهری، علی^۱؛ تابش، مهسا^۱

^۱ گروه فیزیک، دانشگاه صنعتی نوشیروانی، خیابان شریعتی، ۷۱۱۶۷-۷۱۴۸، بابل

^۲ گروه ریاضی، دانشگاه صنعتی نوشیروانی، خیابان شریعتی، ۷۱۱۶۷-۷۱۴۸، بابل

چکیده

در این مقاله، روشی جدید در متن کاوی را با استفاده از آنتروپی کسری ارائه می‌دهیم. نتایج این روش برای استخراج نمایه‌ی کتاب آماری کسلا و برگر (۱۹۹۰) نشان می‌دهد که آنتروپی کسری، نسبت به آنتروپی متداول شانون، ابزار موفق‌تری در متن کاوی است.

Word Ranking and fractional entropy

Mehri-Dehnavi, Hossein¹; Agahi, hamzeh²; Mehri, Ali¹; Tabesh, Mahsa¹

¹ Department of Physics, Babol Noshirvani University of Technology, Shariati Ave., Babol, 47148-71167, Iran,

² Department of Mathematics, Babol Noshirvani University of Technology, Shariati Ave., Babol, 47148-71167, Iran,

In this paper, we propose a new method in text mining using the fractional entropy. Using this method for keyword extraction of the Statistical Inference book by Casella and Berger (1990) indicates that fractional entropy is more efficient than the traditional Shannon entropy in text mining.

PACS No. 05

$$S_{\alpha} = \sum_i \left\{ -\frac{p_i^{-\alpha}}{\Gamma(1-\alpha)} [\ln p_i + \psi(1) - \psi(1-\alpha)] \right\} p_i. \quad (2)$$

$\Gamma(x)$ و $\psi(x)$ به ترتیب نشان‌دهنده‌ی توابع گاما و دی‌گاما می‌باشند. آنتروپی کسری تعمیمی از آنتروپی شانون می‌باشد که در حالت خاص $\alpha=0$ بر آنتروپی شانون منطبق می‌شود. کاربردهایی از آنتروپی کسری در منابع [۷ و ۱] ارائه شده است.

در بخش بعدی، پس از مروری بر روش متداول متن‌کاوی بر اساس آنتروپی شانون، به ارائه مفهوم اهمیت تعمیم‌یافته کسری لغات، WI_{α} ، با استفاده از آنتروپی کسری، پرداخته و بر اساس آن به ارائه روش جدیدی برای متن‌کاوی بر اساس آنتروپی کسری خواهیم پرداخت. در بخش سوم روش ارائه شد خود برای استخراج لغات مهم کتاب آمار استنباطی استفاده می‌نماییم [۳].

۲ متن کاوی

احتمال در متن کاوی

مقدمه

یک مفهوم حائز اهمیت در آمار و احتمالات و فیزیک آماری آنتروپی شانون نام دارد. همان‌طور که از نام این مفهوم برداشت می‌شود، این مفهوم نشان‌دهنده‌ی میزان بی‌نظمی سیستم مورد مطالعه است. آنتروپی شانون به صورت زیر تعریف می‌شود

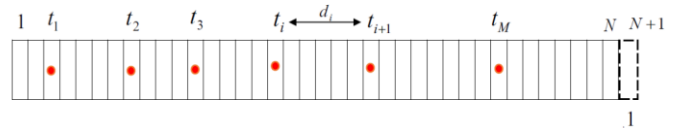
$$S = \sum_i p_i I(p_i) = -\sum_i p_i \ln p_i, \quad (1)$$

که $I(p_i) = -\ln p_i$ اطلاعات شانون و p_i احتمال حضور سیستم مورد مطالعه در حالت i ام می‌باشد، و در شرط $\sum_i p_i = 1$ صدق می‌کند.

یک شاخه‌ی مهم از ریاضیات که اخیراً مورد توجه محققین قرار گرفته حسابان با مشتقات و انتگرالهای کسری است [۲]. با استفاده از مفاهیم مشتق و انتگرال کسری، تعمیمی از آنتروپی شانون تحت عنوان آنتروپی کسری ارائه شد [۴]. آنتروپی کسری از مرتبه‌ی

$\alpha \in R$ به صورت زیر تعریف می‌شود

مطابق با شکل (۱) متنی به طول N را در نظر بگیرید. با این توصیف متن مورد بررسی ما می‌تواند شامل N لغت باشد. فرض بر این است که لغات متن مطابق شکل از ابتدا تا به انتها رتبه‌بندی شده‌اند. همچنین طبق معمول فرض دوره‌ای بودن متن را نیز به کار می‌بریم. کلمه‌ی (مورد نظر) W_1 را در نظر بگیرید که به تعداد (فرکانس) M بار در متن مورد مطالعه تکرار شده باشد. شماره و جایگاه‌های مختلف این کلمه در متن را مطابق با شکل به ترتیب با t_1, t_2, \dots, t_M نشان می‌دهیم. با توجه به توضیحات فوق فاصله‌ی کلمه (مورد نظر) W_1 شماره i با کلمه مشابه بعدی خود در متن برابر با $d_i = t_{i+1} - t_i$ خواهد بود.



شکل ۱: تصویر نمادین از یک کتاب با طول N کلمه و لغت مشخص W_1 با تعداد تکرار M مرتبه که در جایگاه‌های t_i قرار گرفته است. متن به طور دوره‌ای فرض شده است (یعنی جایگاه کلمه‌ی $N+1$ بر جایگاه کلمه‌ی اول منطبق است).

برای محاسبه‌ی فاصله‌ی آخرین کلمه‌ی متن با کلمه‌ی بعدی خود باید فرض چرخه‌ای را در نظر بگیریم. به عبارتی با فرض $t_{M+1} \rightarrow t_1 + N$ خواهیم داشت:

$$d_M = t_{M+1} - t_M = t_1 - t_M + N. \quad (۳)$$

واضح است که جمع فاصله‌های تمامی کلمات برابر با طول کل متن می‌باشد، به عبارتی $\sum_{i=1}^M t_i = N$. حال می‌توانیم برای کلمه مورد نظر W_1 قرار گرفته در جایگاه t_i ، یک احتمال به صورت $p_i = \frac{d_i}{N}$ ، تعریف نمود. واضح است که جمع احتمال‌های مربوط به یک لغت مورد نظر برابر با یک می‌باشد. روش‌های دیگری هم برای تعریف احتمال مربوط به یک کلمه در متن وجود دارد که تعدادی از آن‌ها در [۶ و ۵] مرور شده‌اند.

روش متداول رتبه‌بندی لغات با آنتروپی شنون

در روش متداول رتبه‌بندی لغات یک متن مطابق فرایند فوق برای هر کلمه‌ای یک مجموعه احتمالات $\{p_1, p_2, \dots, p_M\}$ را برای کلمه‌ی دلخواه W_j محاسبه می‌کنند، سپس با استفاده از رابطه‌ی (۱) آنتروپی شنون مربوط به آن کلمه را به دست می‌آورند. حال

اگر آن کلمه موردنظر یک لغت بی‌اهمیت مانند "است"، "در" و ... باشد، این کلمه تقریباً در تمامی متن به طور یکنواخت توزیع شده است. ولی اگر کلمه‌ی مربوطه کلمه‌ای مانند "احتمال"، "توزیع"، "همبستگی" و ... باشد، واضح است که این کلمه در بخش‌های خاصی از متن از قبیل بحث روی نتایج، و نمودارها و ... در مقایسه با قسمت مقدمات متن بیشتر تکرار می‌شود.

اگر در حالت ایده‌آل یک لغت بی‌اهمیت W_1 با تکرار M بار در متن به طور کاملاً یکنواخت توزیع شده باشد. واضح است که $d_i = N/M$ ، $p_i = d_i/N = \frac{1}{M}$ ، و مقدار آنتروپی شنون مربوطه برابر با $S_{uniform}^M = \sum_{i=1}^M \frac{1}{M} \ln M = \ln M$ خواهد بود. که بیشترین مقدار برای

یک کلمه با تعداد تکرار M عدد در متن است. ولی اگر کلمه‌ی موردنظر کلمه‌ای مهم از تعداد تکرار M مرتبه در متن باشد، واضح است که آنتروپی مربوطه از مقدار $\ln M$ کمتر خواهد شد. به همین دلیل هر چقدر مقدار آنتروپی یک لغت با مرتبه تکرار M از مقدار آنتروپی لغتی با همان فرکانس و توزیع یکنواخت ($S_{uniform}^M = \ln M$) متفاوت‌تر باشد این کلمه، کلمه‌ی مهم‌تری است. این روش، روش متداول رتبه‌بندی لغات با آنتروپی شنون است. پس به نظر می‌آید که اختلاف آنتروپی یک کلمه با فرکانس مشخص M و آنتروپی کلمه‌ای با توزیع یکنواخت و همان فرکانس معیار خوبی برای تشخیص اهمیت کلمه‌ی مربوطه باشد. این امر تا حدی درست می‌باشد! دلیل این امر نرمالیزه نبودن معیار اشاره شده است. برای نرمالیزه نمودن این معیار بهتر است اختلاف به دست آمده را بر مقدار آنتروپی لغت با توزیع یکنواخت از همان فرکانس، $S_{uniform}^M = \ln M$ ، تقسیم نماییم. این معیار را اهمیت لغت می‌نامیم.

تعریف ۱. اهمیت یک لغت با مرتبه تکرار M در یک متن به صورت زیر تعریف می‌شود

$$WI(W_1) = \frac{|S_{uniform}^M - S(W_1)|}{S_{uniform}^M}, \quad (۴)$$

که در آن $S(W_1)$ ، و $S_{uniform}^M = \ln M$ به ترتیب آنتروپی لغت W_1 با تکرار M ، مرتبه و آنتروپی یک لغت با توزیع یکنواخت همان مرتبه تکرار می‌باشند.

معیارهای ارزیابی روش‌ها در متن‌کاوی

در این بخش به مرور معیارهای متداول برای متن‌کاوی می‌پردازیم. متنی را در نظر بگیرید که لغات مهم آن از قبل نمایه شده باشند. تعداد لغات حاضر در نمایه این کتاب را N_{rel} فرض می‌کنیم. پس اجرای برنامه و امتیازبندی لغات طبق رابطه‌ی (۴)، به لیستی از لغات که بر حسب معیار اهمیت لغت مرتب شده‌اند خواهیم رسید. تعداد N_{ret} لغت اول به دست آمده در این لیست را در نظر بگیرید. از این تعداد، تعداد $N_{rel \cap ret}$ آن لغاتی هستند که نمایه متن موردنظر نیز وجود دارند. نسبت این دو عبارت را معیار "صحت"^۱ الگوریتم می‌نامیم:

$$P = \frac{N_{rel \cap ret}}{N_{ret}} \quad (5)$$

همچنین نسبت تعداد لغت مشترک به دست آمده در هر مرحله با نمایه، $N_{rel \cap ret}$ به تعداد کل لغات لیست نمایه را هم به عنوان معیار "بازیابی"^۲ تعریف می‌کنند:

$$R = \frac{N_{rel \cap ret}}{N_{rel}} \quad (6)$$

اگر N_{ret} به عنوان متغیر فرض کنیم و معیارهای صحت و بازیابی را بر حسب این متغیر رسم کنیم، می‌توانیم رفتار این توابع را بر حسب متغیر ذکر شده بررسی کنیم. با توجه به تعریف ارائه شده برای این معیارها، رفتار هر دو آن‌ها برای مقادیر کوچک متغیر N_{ret} رفتار نوسانی دارند. همچنین برای مقادیر بزرگ N_{ret} رفتار معیارهای بازیابی و صحت به ترتیب صعودی و نزولی هستند. معیار دیگری که میانگین هارمونیک دو معیار صحت و بازیابی می‌باشد را "معیار- F "^۳ می‌نامند:

$$F = \frac{2RP}{R+P} \quad (7)$$

رفتار کلی (به جز نقاط اولیه نوسانی) این معیار بر حسب متغیر N_{ret} ، با دو معیار بالا متفاوت است. این معیار ابتدا رفتار صعودی داشته و پس از نقطه‌ی بیشینه‌ی خود رفتار نزولی خواهد داشت.

برای هر الگوریتم داده‌کاوی و یا متن‌کاوی، مقدار بیشینه‌ی معیار- F عددی بخصوص برای متن مورد کاوش قرار گرفته است. الگوریتمی که مقدار بیشینه‌ی معیار- F آن، عدد بزرگ‌تری را برای متن خاصی به خود بگیرد، آن الگوریتم موفق‌تری در بررسی آن متن می‌باشد. اخیراً کاربردهای وسیعی از این آنتروپی در رتبه‌بندی کلمات و متن‌کاوی ارائه شده است [۵].

روش پیشنهادی برای رتبه‌بندی لغات با استفاده از آنتروپی کسری

در زیر بخش قبلی روشی را بیان نمودیم که برای هر لغت موردنظر W_1 می‌توان مجموعه احتمال تعریف نمود، و آنگاه با مجموعه احتمال به دست آمده و با استفاده از رابطه‌ی (۱) آنتروپی شنون مربوط به آن لغت را محاسبه نمود. همچنین نشان دادیم که لغات کم اهمیت آنتروپی بیشتری خواهند داشت. همانند روش اشاره شده در زیر بخش قبلی می‌توان احتمال‌های مربوط به یک لغت در یک متن را به دست آورد. پیشنهاد ما این است که برای به دست آوردن لغات مهم یک متن به جای آنتروپی شنون (۱) از آنتروپی کسری S^α ارائه شده در رابطه‌ی (۲) استفاده کنیم. نتایج حاصل از روش متن‌کاوی با آنتروپی کسری، برای حالت خاص $\alpha = 0$ ، برابر با نتایج به دست آمده برای متن‌کاوی با استفاده از آنتروپی شنون خواهد بود.

تعریف ۲. اهمیت تعمیم‌یافته کسری یک لغت با مرتبه تکرار M ، در یک متن به صورت زیر تعریف می‌شود

$$WI_\alpha(W_1) = \frac{|S_{\alpha, uniform}^M - S_\alpha(W_1)|}{S_{\alpha, uniform}^M}, \quad (8)$$

که در آن $S_\alpha(W_1)$ و $S_{\alpha, uniform}^M$ ، به ترتیب آنتروپی کسری مرتبه‌ی α لغت W_1 با تکرار M مرتبه، و آنتروپی کسری با مرتبه‌ی α یک لغت با همان مرتبه تکرار و توزیع یکنواخت، می‌باشند.

۳ نتایج مدل ارائه شده برای کتاب کسلا و برگر

در این بخش به بررسی روش پیشنهادی برای رتبه‌بندی کلمات در کتاب آماری کسلا و برگر خواهیم پرداخت. این مرجع شامل

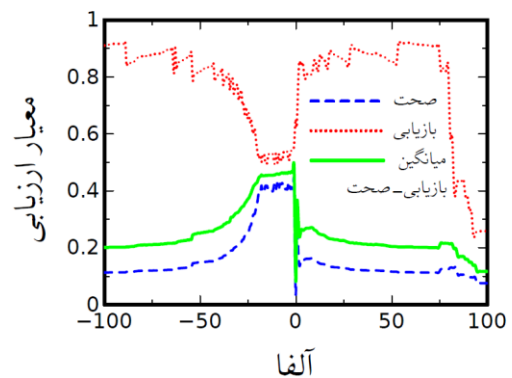
^۱ precision

^۲ Recall

^۳ F-measure

$N = 211133$ کلمه است. از این تعداد کلمه، تعداد کل $N_p = 11665$ کلمه متفاوت وجود دارد.

در این مقاله ابتدا طبق تعریف ۲ لغات کتاب را بر حسب اهمیت تعمیم یافته کسری لغات برای یک α مشخص رتبه بندی کردیم. برای یک مقدار مشخص α ، با استفاده از نمایه کتاب و روش پیشنهادی در بخش قبلی می توانیم نمودار صحت، بازیابی، و همچنین هارمونیک این دو معیار را بر حسب تعداد لغت ظاهر شده در لیست لغات مرتب شده بر اساس اهمیت لغات، به دست آوریم. سپس بیشینه مقدار معیار F را برای پارامتر کسری α مورد نظر به دست می آوریم. مقادیر صحت و بازیابی برای مربوط به نقطه ی بهینه ی معیار F را نیز می توانیم محاسبه کنیم. شکل ۲ نتایج به دست آمده، را بر حسب مقادیر مختلف پارامتر کسری را رسم نموده است. همان طور که در شکل مشاهده می نماییم، بیشترین مقدار میانگین هارمونیک بازیابی و صحت در $\alpha = -0.8$ رخ می دهد.



شکل ۲: بازیابی (خط چین آبی)، صحت (نقطه چین قرمز) و میانگین آن‌ها (خط سبز) برای نمایه های استخراج شده از کتاب آماری کسلا و برگر (۱۹۹۰) به کمک آنتروپی کسری با آلفاهای گوناگون. بیشترین مقدار میانگین هارمونیک بازیابی و صحت در $WI_{-0.8}$ رخ می دهد.

بحث و نتیجه گیری

با استفاده از آنتروپی کسری به معرفی اهمیت تعمیم یافته کسری لغات از مرتبه ی α ، WI_α پرداختیم، که در حالت خاص $\alpha = 0$ ، با اهمیت لغات به دست آمده از آنتروپی متداول شنون، WI برابر است. سپس با استفاده از این معیار به استخراج لغات با اهمیت مرجع [۳]، پرداختیم. نتایج حاصل نشان دادند که مقدار بهینه

پارامتر برای آنتروپی کسری برای متن کاوی کتاب مورد مطالعه برابر با $\alpha = -0.8$ می باشد.

جدول ۱: مقایسه کمی نتایج حاصل از آنتروپی کسری و آنتروپی شنون در استخراج نمایه برای کتاب آمار توصیفی کسلا-برگر. ردیف اول نشان دهنده نتایج به دست آمده با استفاده از اهمیت لغات تعمیم یافته $WI_{-0.8}$ ، ردیف دوم نتایج به دست آمده با استفاده از اهمیت لغات WI به دست آمدند.

میانگین هارمونیک بازیابی و صحت	صحت	بازیابی	جمعیت نسبی نمایه	
۰,۴۹	۰,۵۴	۰,۴۵	۰,۰۵	آنتروپی کسری
۰,۰۹	۰,۶۴	۰,۰۵	۰,۵۶	آنتروپی شنون

با توجه به بررسی های انجام شده در این پژوهش روش رتبه بندی کلمات پیشنهادی از روش رتبه بندی با استفاده از آنتروپی شنون، در جداسازی لغات مهم یک متن بهتر عمل می کند. همچنین جدول ۱ به صورت کمی نیز نشان که آنتروپی کسری، نسبت به آنتروپی متداول شنون، ابزار موفق تر و سریع تری در متن کاوی است.

مراجع

- [۱] Bagci, G. B. (2016). The third law of thermodynamics and the fractional entropies, *Physics Letters A*, 380(34), 2615-2618.
- [۲] Baeanu, D., Diethelm, K., Scalas, E., Trujillo, J.J. (2012). *Fractional Calculus*, world Scientific, Singapore.
- [۳] Casella, G. Berger, R. L. (1990). *Statistical Inference*, Wadsworth, California.
- [۴] Machado, J.T. (2014). Fractional order generalized information, *Entropy*, 16(4), 2350-2361.
- [۵] Mehri, A. Darooneh, A.H. (2011-a). Keyword extraction by nonextensivity measure, *Physical Review E* 83, 056106.
- [۶] Mehri, A. Darooneh, A.H. (2011-b). The role of entropy in word ranking, *Physica A*, 390, 3157-3163.
- [۷] Lopes, A. M., Machado, J.A.T. (2016). Integer and fractional-order entropy analysis of earthquake data series, *Nonlinear Dynamics*, 84 (1), 79-90.