

Комитет науки Министерства образования и науки
Республики Казахстан
РГП «Институт информационных и вычислительных технологий»
КН МОН РК



25 лет
Независимости
Республики Казахстан



25 лет
Институту
информационных и
вычислительных
технологий

МАТЕРИАЛЫ

Международной научной конференции
«Информатика и прикладная математика»
(*«Computer science and Applied Mathematics»*),
посвященной 25-летию Независимости Республики Казахстан и
25-летию Института информационных и
вычислительных технологий

Часть II

г. Алматы, 21-24 сентября 2016 года

Алматы
2016

Содержание

Секция 3. ТЕХНОЛОГИИ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА		7
<i>Kylyshkanov M.K., Uvaliyeva I.M., Belginova S.A.</i>	Informational and analytical system of the biochemical analysis of blood in children	8
<i>Мамбетаев О.А., Мұсабаев Р.Р., Тасболатұлы Н.</i>	Қазақ тілінің дыбыстық түрленімінің ерекшеліктері	20
<i>Несипханова А.Е.</i>	Жеке тұлғаны дауысы бойынша тану жүйесінің жұмысына сипаттама	28
<i>Солтангельдинова М., Мансурова М.Е., Бердібеков С.</i>	Гиперграф негізінде кілттік сөздерді шығару	35
<i>Амиргалиев Е.Н., Мустафин С.А., Муратханова Т.А.</i>	Вопросы определения движения робота на плоскости	44
<i>Барахнин В.Б., Кожемякина О.Ю., Забайкин А.В., Хаятова В.Д.</i>	Автоматизация процесса создания метрических справочников и конкордансов с использованием компьютерных алгоритмов анализа поэтических текстов на русском и казахском языках	47
<i>Кеншилов Ч., Амиргалиев Б., Арсланов М.</i>	Сравнение производительности слоев сверточных сетей для задачи визуального определения местоположения	55
<i>Койбагаров К.Ч., Мұсабаев Р.Р.</i>	Автоматическая классификация отзывов с использованием методов машинного обучения	62
<i>Мустафин С.А.</i>	Моделирование процесса развития закладочного материала	71
<i>Нұртазин А.Т., Хисамиеев З.Г.</i>	Автоматическое извлечение следствий из семантической структуры текста	75
<i>Түкеев У.А., Тұрганбаева А.</i>	Лексикон-фри стемминг для казахского языка	84
<i>Шарипбай А.А.</i>	Искусственный интеллект в Республике Казахстан	89
Секция 4. ИНФОРМАЦИОННАЯ БЕЗОПАСНОСТЬ И ЗАЩИТА ДАННЫХ		102
<i>Kostopoulos George</i>	An original numerical factorization algorithm	103

ГИПЕРГРАФ НЕГІЗІНДЕ КІЛТТІК СӨЗДЕРДІ ШЫГАРУ

Солтангельдинова М.К., Мансурова М.Е., Бердібеков С.

Әл-Фараби атындағы Қазақ Ұлттық Университеті, Қазақстан

Аннотация. Ұсынылып отырган жұмыста кілттік сөздерді шыгару алгоритмі, соның ішінде семантикалық байланысқан сөздерді шыгару технологиясы қарастырылған. Алгоритм графтың ерекше түрі – гиперграфқа негізделеді. Сонымен қатар, «Кездейсоқ кезу» әдісін қолданыла отырып, өндөлеттің құжаттың салмасы қарастырылады. Бұл семантикалық байланысқан сөздерді шыгаруда, яғни өндөлеттің ғылыми құжаттардың барлық ерекшеліктерін ескеруді қамтамасыз етеді. Нәтижесінде алынатын семантикалық байланысқан сөздер, фактографиілік ақпараттық жүйесін құруда жалғасын табады.

Kіріспе

Қазіргі таңда ақпарат көлемі өте үлкен жылдамдықпен ұлғаюда. Үлкен көлемді мәтіндерді оку және сондай үлкен массивті деректерден ақпарат іздеу тиімсіз болып келеді. Сол себепті деректерді өндөлеу жұмыстары алға шығып отыр. Мұндай деректерді өндөлеу ақпараттық-аналитикалық іздеу жүйелерінің негізі болып табылады. Деректерді өндөлеу, талдау бүтінгі таңда әр түрлі салаларда жиі қолданылады. Мәтіндерден қажетті белім шығара білу - талдау, өндөлеу жұмыстарының жаңа бағыты болып есептеледі. Бұл деректерді өндөлеу барысында семантикалық байланысқан сөздерді шыгару есебін өзекті етіп отыр. Шығарылатын сөздер/ сөз тіркестері кейін көп мақсатта қолданылуы мүмкін. Сондықтан да кілттік сөздер, сөз тіркестерін шыгарудың тиімді алгоритмін құру және зерттеу маңызды есептердің бірі.

Ұсынылып отырган алгоритмнің негізгі методикалық сипаттамалары жұмыстың теориялық мәселелер белімінде, ал практикалық жузеге асырылуы түсінілгенде тәжірибе белімінде сипатталынған.

Жұмыстың мақсаты – фактографиялық ақпараттық жүйесін құру мақсатында семантикалық байланысқан сөздерді шыгару. Осы мақсатқа жету негізінде келесі есептер алға қойылды:

1. Мәтіндерді морфологиялық талдаудан өткізу.
2. Мәтіндерді гиперграф түрінде модельдеу, «Кездейсоқ кезу» процесі негізінде Марков шынжыры әдісін қолдану.
3. Граф төбелерінің рангтерін есептеу, салтау.
4. Семантикалық сөздерді шыгару.

Жұмыстың тәжірибелі маңызы – деректерді өндеп, нәтижесінде кілттік сөздердің ишінуы кейін фактографиялық іздеу жүйелерін ұйымдастыруға арналған, құрылатын онтолияның пәндік аймағында сипатталатын, болмыстың мүмкін болатын мәндері ретінде қолданылуы.

I Мәтіннің морфологиялық белгілеу

Мәтіндерді морфологиялық белгілеу табиғи тілде мәтіндерді өндөудің ең маңызды қадамы. Мәтіндердің морфологиялық белгілеуі сөз формаларының белгілеуінен және әрбір сөз формасына лексемалық сипаттамалар мен оның грамматикалық белгілер жиынтығын меншіктеуден тұрады.

Мәтіндегі әрбір сөзге құрылатын морфологиялық ақпарат төрт «жолдан» немесе белгілеулер тобынан тұрады:

1. Сөз формасы жататын лексема (лексеманың «сөздік жазбасын» және қай сөз табы екендігін көрсетеді).

2. Лексеманың грамматикалық белгілер жиынын немесе сөз классификациялауши характеристикалары.

3. Сөз формасының грамматикалық белгілер жиыны немесе сөз түрлендіруші сипаттамалары (мысалы, зат есім септігі, етістік түрі яғни жекеше, көпше түрлері).

4. Грамматикалық форманың стандартты еместігі, орфографиялық бұрмалануы туралы ақпарат.

Мәтіндерді морфологиялық белгілеу бірнеше кезеңнен өтеді:

1. Парсинг. Мәтіндерді өндөудің бұл кезеңінде, сөздерге талдау жүргізіледі. Парсер әрбір сөзге анализ жасап, талдау нәтижесін шыгарады.

2. Фильтрация. Талдауды тазалаудан өткізеді және талдау нұсқаларын жояды.

3. Қолдан омонимияны шешу процессі.

2 Гиперграф негізінде кілттік сөздерді шыгару

Гиперграф – әрбір қабырғамен екі немесе одан да көп ішкі төбелер жиыны байланыса алатын жалпылама граф.

Айталық, бізде n документтен тұратын мәтін бар деп қарастырайық.

$$D = (d_1, d_2, \dots, d_n) \quad (1)$$

Мәтін анализі кезінде, документ ретінде, негізінде, жеке алынған сөйлемдер қарастырылады. Келесі белгілеулерді енгізейік:

$$W = (w_1, w_2, \dots, w_n) \quad (2)$$

- біздің мәтініміздің сөздігі.

Мәтінідер алдын ала морфологиялық талдаудан өткендіктен, әрбір сөздің морфологиялық сипаттамалары белгілі [2].

2.1 Мәтінді гиперграф түрінде модельдеу.

Талдау жүргізлетін мәтінді гиперграф түрінде модельдеу өте ыңғайлы әрі көп қолданылатын әдістің бірі. Талдау жасалатын мәтінімізді, $HG(V, E)$ гиперграф түрінде, V төбелер жиыны және E гиперқабырғалар жиыны ретінде көрсетуге болады. E гиперқабырға V төбелер ішкі жиыны болып келеді, мұндағы $\bigcup_{e \in E} e = V$. Гиперграфтағы

$v \in V$ төбелер ол мәтін сөздері, ал $e \in E$ гиперқабырғалар біздің мәтініміздің документтері.

$HG = R^{V \times E}$ арқылы біздің гиперграфтың араласын матрицасын белгілейік:

$$h(v, e) = \begin{cases} 1, & \text{егер } v \in e \\ 0, & \text{егер } v \notin e \end{cases} \quad (3)$$

$HG(V, E, w)$ - өлшенген гиперграф болсын. Мұндағы $w : E \rightarrow R^+$ - гиперқабырға салмағы.

Біздің жағдайда, гиперграфтың төбе және гиперқабырға дәрежесі келесі формуламен анықталады:

$$d(v) = \sum_{e \in E} w(e)h(v, e) \quad (4)$$

$$\delta(e) = \sum_{v \in V} h(v, e) = |e| \quad (5)$$

D_e және D_v гиперқабырғалар және төбелердің дәрежесін сипаттайтын диагональді матрица.

2.2 Гиперграф төбелерінің салмағын модельдеу

Мәтіндерді өндөуде әрбір сөз маңыздылығын елшеу, яғни сөздің мәтін TF-IDF өлшемі мәтіндер анализі және ақпараттық іздеу есептерінде жиі қолданылады. Мұндай өлшем көбінесе коллекция документтерін сандық вектор түрінде көрсетуде қолданыс табады, яғни қандай да бір сөздер жиынынан әрбір сөздің қолданылу маңыздылығын көрсетеді.

Біздің қарастырып отырған есебімізде, граф төбелерінің салмағы, яғни өнделетін мәтіндердің сөздерінің салмағын TF-IDF әдісіне негіздел, келесі формуламен анықтаймыз:

$$w(v_i)_{tf-idf} = \frac{tf(v_i)}{N_w} * \log \frac{N}{df(v_i)} \quad (6)$$

Мұндағы:

- $tf(v_i)$ - документте жиі кездесетін термин,
- N_w - документке кіретін барлық сөздер суммасы
- N - D документтер коллекциясындағы документтер саны
- $df(v_i)$ - v_i термині кездесетін D-дагы документтер жиыны.

Гиперқабырға салмағы ретінде, мысалы, документ танымалдылығын алуға болады. Осындай ерекшеліктер ескеріліп, келесі формула арқылы гиперқабырға салмағы есептеледі:

$$w(d_i) = \lambda R_{social}(d_i) + (1 - \lambda) R_{time}(d_i) \quad (7)$$

Мұндағы:

- λ - $0 < \lambda < 1$ аралығындағы тегістеу параметрі (сглаживаниес)
- R_{social} - әлеуметтік саптау (социальное ранжирование)
- R_{time} - уақыттық әсер ету өлшемі

Бұл параметрлер келесі формулалармен анықталады:

$$R_{time}(d_i) = Q^{(c-y_i)/24} \quad (8)$$

Бұл жердегі с және y_i ағымдағы уақыт және d_i күжатының жарияланған күні.

Оны 24-ке бөлу арқылы жарияланған уақыт пен ағымдағы уақыт арасындағы шырмашылықты күндер санынан көреміз. Q 0 мен 1 аралығындағы өшү параметрі болып табылады. Q -ды 0,5 деп алу қабылданған.

$$R_{social}(d_i) = \frac{s_i + 1}{\sum_e s_e + 1} \quad (9)$$

Мұндағы, s_i , d_i документі үшін социалды ерекшеліктер саны. $\sum_e s_e$ барлық документтегі социалды ерекшеліктер қосындысы.

Осы ерекшеліктердің бәрі қамтылып, есептеу барысында гиперқабырға салмағының формуласын анықтайды. Бұл гиперқабырға және төбелер салмағы кейін интуїдерде өз қолданысын табады.

2.3 «Кездейсоқ кезу» процесі. Марков тізбегі

Кездейсоқ кезу – дискретті уақыт мезеттерінде қадамдардың кездейсоқ өзгеру процесінің математикалық моделі. Сонымен қоса, әрбір қадамдағы өзгерулер алдындағылардан және уақыттан тәуелсіз деп үйіралады. Бұл модель әртүрлі сфераларда қолданылады, сондай – ақ, мұндай модель шынай процесстің айтартықтай женілдеуі болыш табылады.

Бірөлшемді дискретті кездейсоқ кезу – бұл дискретті уақытты $\{Y_n\}_{n \geq 0}$ кездейсоқ процесс және оның түрі келесідей:

$$Y_n = Y_0 + \sum_{i=1}^n X_i \quad (10)$$

Мұндагы:

- Y_0 – бастапқы күй.

- $X_i = \begin{cases} 1, & p_i \\ -1, & q_i = 1 - p_i \end{cases}, 0 < p_i < 1, i \in \mathbb{N}$

- $Y_0, X_i, i=1,2,\dots$ кездейсоқ шамалар бірге тәуелсіз.

Бірөлшемді дискретті кездейсоқ кезу бүтін күйлері бар *Марков тізбегі* болыш табылады. Оның бастапқы үлестірілімдері X_0 кездейсоқ шаманың ықтималдық функциясы арқылы беріледі, ал көшудің ықтималдықтар матрицасының түрі келесідей:

$$P \equiv (p_{ij})_{i,j \in \mathbb{Z}} = \begin{pmatrix} \dots & \dots & \dots & \dots \\ & q_{-1} & 0 & p_{-1} \\ & q_0 & 0 & p_0 \\ & q_1 & 0 & p_1 \\ & \ddots & \ddots & \ddots \end{pmatrix}$$

Яғни:

$$p_{i,i+1} \equiv P(X_{n+1} = i+1 | X_n = i) = p_i,$$

$$p_{i,i-1} \equiv P(X_{n+1} = i-1 | X_n = i) = q_i, \quad i \in \mathbb{Z},$$

$$p_{ij} \equiv P(X_{n+1} = j | X_n = i) = 0, \quad |i-j| \neq 1.$$

$\{X_n\}_{n \geq 0}$ дискретті кездейсоқ шамалар тізбегі қарапайым Марков тізбегі деңгейлік аталады, егер: $P(X_{n+1} = i_{n+1} | X_n = i_n, X_{n-1} = i_{n-1}, \dots, X_0 = i_0) = P(X_{n+1} = i_{n+1} | X_n = i_n)$. Осылайша, қарапайым жағдайда, Марков тізбегінің келесі күйлерінің шартты үлестірілімдері ағымдағы күйге ғана тәуелді және алдыңғы күйлердің барлығына тәуелсіз.

$\{X_n\}$ кездейсоқ шаманың мәндер облысы – тізбек күйлерінің кеңістігі деңгейлік аталады, ал n - қадам нөмірі.

$P(n)$ матрицасы n - ші қадамдағы көшулер ықтималдығының матрицасы деңгейлік аталады:

$$P_y(n) = P(X_{n+1} = j | X_n = i) \quad (11)$$

Ал вектор $p = (p_1, p_2, \dots, p_n)^T$, мұндагы $p_i = P(X_0 = i)$ – Марков тізбегінің бастапқы үлестірілімі. Көшулер ықтималдығының матрицасы стохастикалық болатындығы айқын, яғни:

$$\sum_j P_{ij}(n) = 1, \forall n \in N \quad (12)$$

Марков тізбегі бірынғай деп аталады, егер көшулер ықтималдығының матрицасы қадам нөмірінен тәуелсіз болса:

$$P_{ij}(n) = P_{ij}, \forall n \in N \quad (13)$$

Гиперграф төбелерін саптау үшін кездейсоқ кезу процесsein гиперграфқа жалпылаймыз. Кездейсоқ кезу процесsei графтағы төбелер арасындағы ауысулар, яғни Әрбір t қадамды дискретті уақыттан кейін берілген төбеден көрші төбеге ауысу болып табылады. Біз төбелерді $\{s_1, s_2, \dots, s_n\}$ күйлер жиыны ретінде қарастыра аламыз, ал ауысулар осы күйлердің ақырлы M Марков тізбегі болады. Ауысу ықтималдығы $P(u, v) = \Pr[ob(s_{t+1} = v | s_t = u)]$ түрінде есептеледі, бұл M Марков тізбегінің v -да $t+1$ уақытта болатындығын, t уақытта u -да болғандығын білдіреді. Біздің жағдайда, Марков тізбегі бірынғай болып табылады, ауысу ықтималдығы t уақыттан тәуелсіз. Әрбір u төбесі үшін $\sum_v P(u, v) = 1$.

M тек бір ғана ауысу үшін есептелген ықтималдықтары бар бірынғай болып табылады. Барлық жүріс үшін ауысу матрицасын $P \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$ есептеуге болады. P ауысулар матрицасы сәйкес сондай ықтималдықтармен төбелер арасындағы серфер іс әрекетінің кездейсоқ түрде орын ауыстыруын көрсететін төбелер арасында ауысуларды толығымен қамтиды.

Қарапайым графтарда, кездейсоқ кезу процесsei түсінікті болып келеді, тек қана ықтималдығы анықталған мақсатты төбеге баратын қабырға таңдалады. Сонда да, гиперграф үшін бұл жағдай біраз өзгеше орын алады. Бұл гиперграф құрылымының өзгешелігінде. Мысалы, гиперграфта, гиперқабырғаның екіден көп $\delta(e) \geq 2$ төбелер шүктелері болуы мүмкін.

Гиперграфта кездейсоқ кезу процесsein жалпылау үшін, біз гиперқабырғада бір бірімен инцидентті болатын екі төбе арасындағы ауысады кезу түрінде модельдеміз. Анықтама, кездейсоқ кезу процесsei бір қадамның орнында екі қадамдық процес: бірінші, кездейсоқ серфер ағымдағы u төбесімен инцидентті болатын гиперқабырғасын таңдайды. Екінші, серфер $u, v \in e$ шартын қанағаттандыратын, таңдалған гиперграфтан мақсатты v төбесін таңдайды.

Кездейсоқ кезу гиперграфта жалпылама аталады, ал қарапайым графта кездейсоқ кезудің ерекше жағдайы бар. Ол яғни, қабырғаның бір ғана төбесінің болуы, ал гиперграфта көп төбелер арасынан таңдай аламыз. Гиперграфпен кездейсоқ кезу процесsein Марков шынжыры арқылы анықтайтын болсақ, мұнда төбелер жиыны түйлөр жиынын құрайтын болады. Әрбір t уақытты қадам сайын, серфер инцидентті гиперқабырғада басқа төбелерге орын ауыстырады.

Бұл жұмыста, біз кездейсоқ кезудің өлшемен гиперграфтағы жалпы анықтамасын анықтауға тырысамыз, яғни гиперқабырғамен қоса төбелердің де салмағы ескерілетін болады. Мұндай жағдайда, кездейсоқ кезу процесsei гиперқабырға және төбелер салмағын қолдану арқылы кеңеje түседі. Төбе салмағын барлық инцидентті гиперқабырғалар арқылы анықтаймыз, ол ерекшеліктер векторы:

$$\vec{v}_w = \{w(v_{e1}), w(v_{e2}), \dots, w(v_{d(v)})\} \quad (14)$$

Яғни, бізде v төбесі кездесетін әрбір e гиперқабырғасы үшін төбелердің әртүрлі салмағы бар. Кездейсоқ кезудің болжамалы процессін келесі түрде сипаттаймыз: и тобесінен бастап серфер гиперқабырға салмағына $w(e)$ пропорционалды болатын, и төбесімен инцидентті болатын e гиперқабырғасын таңдайды. Одан кейін серфер, дәл солай, гиперқабырғада төбе салмағына пропорционалды v төбесі таңдалады, яғни ағымдағы гиперқабырғада қарастырылып отырган салмақ.

Өлшенген гиперграфтың инцидентті матрицасын $H_w \in R^{|V| \times |E|}$ келесідей анықтайық:

$$h_w(v, e) = \begin{cases} w(v_e), & \text{if } v \in e \\ 0, & \text{if } v \notin e \end{cases} \quad (15)$$

Осылайша, гиперқабырға дәрежесін қайтадан анықтасақ:

$$\delta(e_w) = \sum_{v \in V} h_w(v, e) \quad (16)$$

Жоғарыда көрсетілген формулаларды қолдана отырып, P көшу (ауысу) матрицасын есептеуге болады:

$$P(u, v) = \sum_{e \in E} w(e) \frac{h(u, e)}{\sum_{\hat{e} \in \hat{E}(u)} w(\hat{e})} \frac{h_w(v, e)}{\sum_{\hat{v} \in \hat{V}(e)} h_w(\hat{v}, e)} \quad (17)$$

Немесе матрицалық белгілеу:

$$P = D_v^{-1} H W_e D_{ve}^{-1} H_w^T \quad (18)$$

Мұндагы:

- D_v - төбенің өлшенген дәрежесінің диагональді матрицасы;
- H - гиперграф төбелерінің инцидентті матрицасы;
- W_e - гиперқабырға салмақтарының диагональді матрицасы;
- D_{ve} - гиперқабырғаның өлшенген дәрежесінің диагональді матрицасы;
- H_w - өлшенген графтың инцидентті матрицасы.

Мұнда P көшу матрицасы стохастикалық болып табылады және әрбір жолдар қосындысы 1-ге тең [1].

P көшу матрицасын есептегеннен кейін, кездейсоқ кезудің π стационарлы үслетірілімін түсіне қажеттігі түндейді. Стационарлы үлестірілім $\vec{v}_0 \in R^{|V| \times 1}$ векторынан бастап есептелуі мүмкін. Ол $1/|V|$ қосындылары 1-ге тең ықтималдықар. Ен алдымен, P^T көшу матрицасын \vec{v}_0 баған векторына көбейту арқылы $\vec{v}_0 = P^T * \vec{v}_0$ алынады. Осылайша, \vec{v} векторы өзгеруін тоқтатқанға дейін итерацияны жалғастырамыз. \vec{v} ықтималдықты үлестіру векторын көше матрицасына көбейту, келесі үлестірім қадамын $\vec{x} = P^T * \vec{v}$ береді. x_i ағымдағы i төбесінде болу ықтималдылығы болсын. Онда $x_i = \sum_j p_{ij} v_j$, v_j серфердің j түйінінде ал болуының ықтималдылығы, және p_{ij} j -дан i -ға көшу ықтималдылығы.

нәртүрлі
ймыз: и
иатын, и
шфер, дәл
ды, яғни

келесідей

(15)

і векторының үлестірілім ықтималдылығы өзгеруі п қадамнан кейін тоқтайды, егер кездесоқ кезу эргодикалық болса. Кездесоқ кезу эргодикалық болады, егер келесі шарттар орындалса:

1 Кез келген екі $s_i, s_j \in M$ күйлері үшін тізбек аударылмайтын (неприводимый) болса, және олар келесі шартты қанагаттандыруы тиіс: $P(s_i, s_j) > 0$.

2 Тізбек апериодикалық болса, яғни әрбір күйдің ЕҮОБ $\{t : P_t(s_i, s_j) > 0\}$ ол 1.

Аударылмаушылық және апериодикалық қасиеттерді тұрақтандыру үшін PageRank алгоритмін қолдануға болады [1].

(16)

3 PageRank алгоритмі
PageRank алгоритмі - бұл сілтемелік саптау алгоритмдерінің бірі. Бұл алгоритм көбінесе гиперсілтемелермен байланысқан документтер колекцияларына қолданылады.

(17)

Бұл алгоритмнің кездесоқ кезу процесіне қолданудағы негізгі идея, бұл алгоритм телепортация идеясын қолданады, яғни кездесоқ кезу процесін қайтадан жібереді. Ал бұл жағдай кезедесоқ кезудің эргодикалық қасиеттері үшін тиімді. Телепортация өте ықтималдықты өшу факторы (damping factor) арқылы көрсетіледі. Бұл сондай –ак графты аударылмайтындай жасауға кепіл береді, себебі кездесоқ кезудің кез келген басқа бір түйіндерге телепортация жасау ықтималдылығы бар.

$$\vec{v}_{i+1} = \alpha P^T \vec{v}_i + (1 - \alpha) * 1/n \quad (19)$$

Мұндагы:

- α – өшу факторы, оны 0,85 деп аламыз.

- n - графтагы түйіндер саны.

- $\alpha P^T v_i$ кездесоқ кезу бір инцидентті гиперқабырғадан таңдалынатынын білдіреді.

4 Шаблондар әдісі

Жалпы тұрғыда кілттік сөздерді шыгару әдістері өте көп. Ондай әдістердің қатарына N-gram, TF-IDF, TextRank, POS chunking және т.б. Соның ішінде POS chunking әдісі салыстырмалы түрде жақсы нәтиже көрсетеді.

POS chunking әдісі негізінде, морфологиялық шаблондарға негізделген. Бұл кілттік сөздерді шыгару кезінде, мәтін сөздерінің қатаң сөз таптары мен оның характеристикаларын бақылауды талап етеді. POS моделін аңдатылған деректерден шыгаруға болады. Кілттік сөздерді шыгару процесі кезінде, ең алдымен, мәтіндерге талдау жүргізіліп, әрбір сөздің морфологиялық сипаттамалары алынады. Сол сипаттамаларды қолдану арқылы, сәйкес шаблондар құрылып, кілттік сөздерді шыгару шағуға асырылады.

POS chunking әдісінің артықшылығы, деректерден мүмкіндігінше мағынасы бар кілттік сөздерді шыгаруга тырысады, яғни, ол мағынасыз кілттік сөздер санын қысқартады.

5 Есептеуіш тәжірибе

Жоғарыда нұсқалған алгоритмге тәжірибе жүргізу барысында Леонид Орловтың «Как создать электронный магазин в интернет» оку құралы қолданылды. Бұл кітап толықтай 384 бетті қамтиды. Тәжірибе барысында ең алдымен, кітап толықтай морфологиялық талдаудан өткізілді. Бағдарлама Java тілінде жазылғандықтан,

морфологиялық талдау жасауда Java тіліне негізделген Apache Lucene кітапханасы пайдаланылды. Талдаудан өткізілгеннен кейін құжатты қамтитын барлық сөздер сипаттамасы жеке түрде сакталынып отырды. Бұл сипаттамалар семантикалық сөздер / сөз тіркестерін шығаруда, шаблондар әдісін қолдану кезінде қажет болып табылады.

Kіріс ақпарат өндеуден өткізілген соң, ендігі сол деректермен жұмыс жасалынды. Өндөлген деректерді гиперграфқа негіздей, қажетті есептеулерді жүргіздік, яғни, ол құжаттың өзінің салмағын, құжат мәтініне кіретін әрбір сөздің салмағын анықтау, байланысу ықтималдықтарын есептеу. Бұл жағдайда, гиперграфтың олшенген және өлшенбеген түрлері де қарастырылды.

Алгоритмін ең маңызды бөлігі, сөздер арасындағы байланысты анықтау, сөздер рангін есептеу, «Кездесік кезу» моделіне негізделді. Кездесік кезу моделі қарапайым Марков тізбегінде көрсетілді. Яғни, тізбектің күйлері ретінде граф төбелері, қарастырып отырған сөздер жиыны алынды. Семантикалық байланысты анықтау кездесік түрде, бір күйден екінші күйге көшу арқылы көшу ықтималдықтарын есептеу көмегімен анықталды. Көшу ықтималдығы, құжат мәтініне кіретін әрбір сөздің барлық сөздермен байланысу ықтималдығын көрсетеді. Үйкималдық жоғары болған сайын, байланысу мүмкіндігі де жоғары.

Кездесік кезу процессын кейін стационарлы үлестірім мәнін есептеу жүргізілді. Бұл мән PageRank алгоритміне негізделді, және де бұл есептеу бізге сөздің рангін береді. Сол ранг арқылы кілттік сөздерді анықтауга болады. Ранг жоғары болған сөздер, құжат аясында жиі кездесетін кілттік сөздер болу мүмкіндігі бар. Ранг есептеу кезінде (19) формулада көрсетілген, а өшү факторы мәні ретінде 0,89 мәнін алу келісілді. Осылайша, әрбір сөздің рангтері анықталды.

Кілттік сөздер / сөз тіркестерін шығару мақсатына қол жеткізуде шаблондар әдісі қолданылды. Шаблондар әдісі негізінде өнделетін тілдің құрылымына байланысты болып табылады, және де сол бойынша шаблондар анықталады. Біздің жағдайда, орыс тілді құжаттар қарастырылғандықтан, сәйкес шаблондарды анықтап алдық. Ол шаблондар келесідей болды:

- [прилагательное + существительное];
- [причастие + существительное];
- [существительное + существительное, Род.падеж];
- [существительное + существительное, Твор.падеж];
- [числительное + существительное, Род. падеж + числительное];
- [существительное, Им.падеж].

Бұл шаблондарда көрсетілген сипаттамалар, өнделетін құжатқа ең бастапқыда жүргізілген морфологиялық сипаттамалардан алынады.

Төменде осы алгоритм нәтижесінде альпиган кілттік сөздер / сөз тіркестері бөлігі көрсетілген:

электронный бизнес, электронная коммерция, данные клиента, высокопроизводительная система, финансовая информация, неоплаченные заказы, первый взнос, конкретный платеж, расчетный банк, другой перевод, сформированный заказ, сеть, графические элементы, курьерская доставка, цифровая подпись, гипертекстовая связь, коммуникация, бюджет, компьютерную сеть и т.д.

Пай	1.
for Short	2.
osts2014)/	3.
онтологии	4.
дународні	Asian Fed
бірск, 6-8	5.
систем ин	IEEE 10th
турирован	6.
	In M. Colli
	7.
	decomposi
	tational Li
	8.]
	ing. In CII
	knowledge
	9.]
	ированнс
	10.]
извлечени	
университ	

Пайдаланылған әдебиеттер тізімі

1. Bellaachia A., AlDhelaan M. HGRANK: A Hypergraphbased Keyphrase Extraction for Short Documents in Dynamic Genre. // Making Sense of Microposts (# Microposts2014)/http://ceur-ws.org/Vol-1141/paper_06.pdf (available: 01.06.2016)
2. Пастушков И.С., Барахнин В.Б. Алгоритм автоматизированного наполнения онтологии фактографической поисковой системы//Всероссийская конференция с международным участием "Знания - Онтологии - Теории" (ЗОНТ-2015), с.87-95, Новосибирск, 6-8 октября 2015 г.
3. Барахнин В. Б., Федотов А.М., Шокин Ю.И. Технология создания программных систем информационного обеспечения научной деятельности, работающих со слабоструктурированными документами // Вычислительные технологии. – 2010 – Т. 15. – N 6.
4. Bougouin A., Boudin F., Daille B. Topicrank. Graph-based topic ranking for keyphrase extraction. In Proceedings of the Sixth IJCNLP, Nagoya, Japan, October 2013. Asian Federation of Natural Language Processing. – P. 543–551.
5. Gao Y., Liu J., and Ma P. The hot keyphrase extraction based on tf*pdf. In The 2011 IEEE 10th TrustCom, pages 1524-1528, 2011.
6. Hulth A. Improved automatic keyword extraction given more linguistic knowledge. In M. Collins and M. Steedman, editors, Proceedings of the 2003 EMNLP, 2003. – P. 216-223.
7. Liu Z., Huang W., Zheng Y., and Sun M. Automatic keyphrase extraction via topic decomposition. In Proceedings of the 2010 EMNLP, pages 366–376. Association for Computational Linguistics, October 2010.
8. Kushal S. Dave, Varma V. Pattern based keyword extraction for contextual advertising. In CIKM '10 Proceedings of the 19th ACM international conference on Information and knowledge management, pages 1885-18887 New York, USA, 2010.
9. Бритков В.Б., Булычев А.В. Методы анализа больших объемов слабоструктурированной информации. – 2010.
10. Шереметьева С.О., Осминин П.Г. Методы и модели автоматического извлечения ключевых слов. – Вестник ЮжноУральского государственного университета, – №1/12. – 2015.