

УДК 004.89

Лиман К.С. м.н.с.

Моделі тематик в застосуванні до зняття семантичної неоднозначності

Задача визначення значення слова, з множини можливих значень – онтології, відома як зняття семантичної неоднозначності (WSD, Word Sense Disambiguation). Значення конкретного слова у тексті залежить від його контексту. Значною перевагою використання онтологій, є те, що в них описано відношення між поняттями.

В тематичних моделях, таких як LDA (Latent Dirichlet Allocation), документ представляється як вектор пропорцій тематик, а кожна тематика як розподіл над словами. Але для потреб WSD необхідно мати розподіл над значеннями слів, як наприклад в LDAWN (розширення LDA на WordNet). Але LDAWN не враховує значення семантичної схожості і пов'язаності між концептами. Тому було запропоновано тематичну модель пов'язаних концептів, одним з параметрів якої є взаємне «розташування» концептів задане матрицею пов'язаності.

Ключові слова: семантична неоднозначність, тематичні моделі, семантична схожість.

¹ Міжнародний науково-навчальний центр інформаційних технологій і систем НАН України та МОН України, 03680 м.Київ Глушкова 40, e-mail: lyman.kostiantyn@gmail.com

Lyman K.S. researcher

Topic models application to word sense disambiguty

The problem of choosing a sense for a particular word from fixed set – ontology, is known as Word Sense Disambiguation (WSD). The main advantage of ontologies is that they contain information about semantic relationships among senses, although meaning of a word in a text depends on its context.

Within topic models, such as LDA (Latent Dirichlet Allocation), a document is presented as a vector of topic proportions, and each topic is a distribution over words. But for WSD purposes the distribution over senses is required, as for example in LDAWN (an extension of LDA over WordNet). Unfortunately, LDAWN does not taking into account a semantic similarity and relatedness estimations between concepts. That's why a topic model of related concepts was introduced, which has a matrix of relatedness as one of its parameters.

Key Words: word sense disambiguation, topic models, semantic similarity.

¹ International Research and Training Center for Information Technologies and Systems 03680, Kyiv, Glushkova st. 40, e-mail: lyman.kostiantyn@gmail.com

Статтю представив чл.-кор. НАНУ, д.ф.-м.н., проф. Анісімов А.В.

Зняття семантичної неоднозначності (Word Sense Disambiguation, WSD) є класичною задачею комп'ютерної лінгвістики. Вона полягає в тому, що комп'ютерна система повинна вибрати з запропонованих правильне значення слова в

даному контексті. Питання про те якою має бути множина цих значень, наскільки повною та детальною, і те як вона представляється є, напевне, одними з найпроблематичніших у сфері комп'ютерного семантичного аналізу та

інжинірингу знань. Проблема в тому, що люди, хоча і розуміють один одного, не зважаючи на різний мовленнєвий досвід, все таки не завжди відчувають різницю між різними значеннями слова¹. То чи варто очікувати, що комп'ютер зможе краще? Виділяють два основні підходи до визначення значення слів [1]. Перший, назовемо його *дистрибутивним*, базується на факті, що смисл слова визначається його контекстом. Відомий приклад – вірш Jabberwocky Льюїса Керрола².

Інший підхід використовує задані наперед множини можливих значень слів, такі як словники, онтології чи бази знань, створені експертами які і визначають структуру та рівень деталізації семантики слів. Онтології зазвичай, не покривають всі можливі смисли та їх варіації, вони обмежуються або спеціальним доменом знань або деякою частиною загальноживаної лексики (як наприклад WORDNET [8]). Назвемо цей підхід *знання-залежним* (knowledge based). Використання таких джерел структурованої семантики має ту перевагу, що прописані взаємозв'язки між поняттями відображають певні знання експерта про предметну область, але їх створення та підтримання в актуальному стані є трудомісткою задачею.

Задача автоматичного визначення всіх можливих смислів кожного слова відома як Word Sense Induction, WSI [6]. Вона тісно пов'язана з WSD, але відрізняється від неї: в WSD ми знаємо множину понять і нашою задачею є обрати такий елемент цієї множини, що найкраще пасує до контексту. Натомість WSI представляється як задача розділення контексту цільового слова на кластери, кожен з яких відповідає певному смислу цього слова. Приклади таких систем описані в [10] та [6].

Моделі тематик (наприклад, Latent Dirichlet Allocation [4]) є генеративними імовірнісними моделями колекції документів D , тобто вони описують імовірнісний процес генерації спостережуваної колекції документів (задаючи розподіл $p(x/y)$, де x – спостереження, а y –

приховані змінні): кожен документ є вектором пропорцій у просторі тематик отриманим з певного тематичного розподілу; тематики при цьому задаються певним розподілом у просторі всіх слів. Інформація закладена у тому яким чином слова сумісно зустрічаються дають змогу знайти такі параметри моделі, що отримані тематичні розподіли можна інтерпретувати саме як тематики. В останні роки з'являються багато робіт покликаних об'єднати WSD та моделі тематик. В цій статті буде розглянуто та проаналізовано одну з таких моделей - LDAWN [5], в якій тематичний розподіл включає розподіл над прихованими змінними, що відповідають синсетам WordNet. В [7], наприклад, розподіл над синсетами не залежить від тематики і слова генеруються або з тематики, або з синсета.

Більш детальний огляд сучасного стану справ в області зняття семантичної неоднозначності, а також огляд систем, підходів та результатів можна знайти в наступних оглядових роботах [12], [13] та збірнику [1].

LDAWN: Тематично-залежне випадкове блукання по онтології

В LDAWN тематика описується не поліноміальним розподілом над елементами словника, як в LDA, а випадковим блуканням по WordNet – WORDNET-WALK. Основною структурною одиницею WordNet є синсет – множина синонімів, об'єднаних певним значенням. Кожен синсет пов'язаний з іншим певними семантичними та лексичними відношеннями.

WORDNET-WALK це імовірнісний процес генерації слів, що використовує для цього гіпонімічні зв'язки³ між синсетами у WordNet. Згідно [5], уявімо агента, що починає рух з синсету s_0 entity – кореня ієрархії іменників (організована у вигляді таксономії). Агент обирає наступний вузол в своєму випадковому блуканні серед гіпонімів поточного концепту. Дана процедура повторюється поки не буде досягнуто листового вузла, який відповідає певному слову. Вибір наступного синсету в блуканні визначається відповідним розподілом над множиною синсетів, яка є множиною прихованих випадкових змінних, розподіл яких впливає на

¹ Особливо враховуючи, те що погоджуються люди між собою теж не так часто: в експерименті з двома анотаторами, які виконували розмітку текстів значеннями з WORDNET, згода була досягнута лише в 72.5% випадках [13]

² У перекладі Інни Коваль (2004):

Шов печір. Яштпорки слибкі
В трамиці дзигали якраз
Свистали свердлов'ї бідкі
І звинки мрюкали весь час.

³ Гіпоніми певного поняття є конкретизація ми цього поняття

процес генерації послідовності слів
 $\rho(w^d | z^d, \Lambda^d)$.

Для позначення змінних, які приймають значення зі скінченної множини, будемо використовувати нотацію фіктивних змінних:

$$b \in A: b = A_s \Leftrightarrow b = a \Leftrightarrow b = (0_1, \dots, 1_a, \dots, 0_A) \quad (1)$$

де $A = |A|$. Також замість $b \in A$ іноді будемо писати $b \in 1:A$, при цьому запис $b = 1:A$ означає, що b послідовно приймає всі значення з A . Введемо наступні позначення:

- $w^d = (w_1^d, \dots, w_{N_d}^d)$ – вектор випадкових величин, що приймають значення зі словника V , і представляють документ розміром у N_d слів, $V = |V|$;

- $z^d = (z_1^d, \dots, z_{N_d}^d)$ – вектор випадкових величин, значення яких вказують на одну з обраних для кожного слова тематик $z_i^d \in K$, $K = |K|$;

- $\Lambda^d = (\Lambda_1^d, \dots, \Lambda_{N_d}^d)$ – вектор шляхів від кореня онтології до синсета з якого було згенеровано відповідне слово. Кожен шлях Λ_i^d отримано з відповідного тематичного блукання, обраного згідно значення z_i^d .

Множину всіх синсетів позначимо S , $S = |S|$.

LDAWN провадить наступний генеративний процес:

- 1) Для кожної тематики $t = 1:K$:
 - а. Для кожного синсету $s \in S$, обираємо розподіл переміщень $\beta_s^t \sim \text{Dir}^{c(s)}(\tau_s)$ – імовірність переходу від s до i -того гіпоніма c_i^s .
- 2) Для кожного документу $w^d \in D$, $d = 1:D$:
 - а. Обираємо тематичний розподіл $\theta^d \sim \text{Dir}^K(\alpha)$;
 - б. Для кожного слова w_u^d , $u = 1:N_d$:
 - і. Обираємо тематику $z_u^d \sim \text{Mult}^K(1, \theta^d)$;

- ii. Будуємо шлях Λ_n^d який починається в кореневому вузлі λ_0 .

- iii. Серед нащадків λ_i :

1. Обираємо наступний вузол у блуканні

$$\lambda_{i+1} \sim \text{Mult}^{c(\lambda_i)}\left(1, \beta_{\lambda_i}^{z_n^d}\right);$$

2. Якщо λ_{i+1} є листковим вузлом, то генеруємо відповідне слово. Інакше повторюємо.

де $\text{Dir}^X(\alpha)$ – X -вимірний розподіл Діріхле,

$\text{Mult}^X(1, \theta)$ – X -значний поліноміальний розподіл з однією спробою, а $c(s)$ – множина нащадків (гіпонімів) синсету s .

Прихованими змінними цієї моделі є параметри K різних тематичних випадкових блукань WORDNET-WALK $\beta_{1:K}^t$, вектор пропорції тематик у документі θ^d , вектор змінних z^d , що визначають тематику для кожного слова, а також шлях у онтології для кожного слова Λ^d .

Для даної колекції документів $w_{1:D}$ повна апостеріорна імовірність прихованих параметрів:

$$\rho(\beta_{1:D}, z_{1:D}, \theta_{1:D}, \Lambda_{1:D} | w_{1:D}, \alpha, \tau) \propto \prod_{t=1}^K \rho(\beta_t | \tau) \prod_{d=1}^D \rho(\theta_d | \alpha) \times \quad (2).$$

$$\times \prod_{u=1}^{N_d} \rho(w_u^d | \Lambda_u^d) \rho(\Lambda_u^d | \beta_{z_u^d}^d) \rho(z_u^d | \theta_d)$$

Тепер задачу визначення значення слова можна поставити наступним чином: яка імовірність, що слово w_u^d було отримано зі шляху, що закінчується синсетом s . Тобто потрібно отримати апостеріорний розподіл прихованих змінних.

Інтуїтивне припущення, на якому базується LDAWN, полягає в тому, що слова однієї тематики семантично близькі і тому будуть поділяти шляхи в WORDNET. При спостереженні певного синсета s (збільшується імовірність спостерігати далі синсети, чії геренеруючі шляхи мають спільні вузли з аналогічним шляхом для s , тобто мають спільних прашчурів.

Нехай апіорне значення $\tau_{s,c}$ імовірності переходу від синсету s до його нащадка c є пропорційним загальній кількості спостережуваних слів, які є нащадками c [5]. В цьому випадку отримаємо імовірнісну інтерпретацію поняття ІНФОРМАЦІЙНОГО КОНТЕНТУ [14].

$$IC(s) = -\ln p(s) \quad (3)$$

Де $p(s)$ імовірність зустріти s або одного з його нащадків. Цю величину можна оцінити наступним чином

$$\hat{p}(s) = \frac{|W \in d : \text{cover}(s, W)|}{N_d} \quad (4)$$

тобто яка частина слів документу D покривається синсетом s (слово W покривається синсетом s , $\text{cover}(W, s) = \text{True}$, якщо W входить до s або будь-якого його нащадка).

Апостеріорний розподіл для LDAWN

Обчислення апостеріорного розподілу прихованих параметрів пов'язано з інтегруванням по всіх можливих значеннях цих параметрів (для обчислення граничного розподілу спостережуваних даних). Тому на практиці використовуються методи Monte Carlo Markov Chains (MCMC). Одним з них є Гіббсівське семплювання⁴. Для записати умовну імовірність одного елементу вектору випадкових величин, від фіксованих значень решти елементів, таким чином одночасно переміщуючись тільки в одному напрямку.

В LDAWN для кожного слова W_u^d в кожному документі семплюється призначення тематики Z_u^d і тематично-залежний шлях через WordNet Λ_u^d . Для обох змінних потрібно обчислити умовний розподіл від спостережуваних даних. Призначення тематики та шляху семплюється спільно, при умовній залежності від значень прихованих змінних обраних для всіх інших слів. Запишемо апостеріорний розподіл тематик і шляхів:

$$\begin{aligned} p(Z_u^d = t, \Lambda_u^d = \pi | Z_{-u}^d, \Lambda_{-u}^d, \tau, \alpha) = \\ = \underbrace{p(Z_u^d = t | Z_{-u}^d, \alpha)}_{\text{тематика}} \underbrace{p(\Lambda_u^d = \pi | Z_u^d = t, \Lambda_{-u}^d, \tau)}_{\text{шлях}} \end{aligned} \quad (5)$$

Позначимо $n_{-u}^{d,t}$ кількість слів відмінних від u , яким було призначено тематику t в документі d , якому слово W_u належить. Також, нехай $b_{-u}^{t,s,c}$ позначає кількість переходів в випадковому блуканні t -тої тематики від синсету s до його нащадка c , не враховуючи при цьому шлях π пов'язаний зі словом u , тобто кількість переходів від s до c у шляхах $\Lambda_{-u}^{d,t}$. Розглянемо перший множник (імовірність тематики) спільного розподілу:

$$\begin{aligned} p(Z_u^d = t | Z_{-u}^d, \alpha) = \\ = \int_{\theta} p_{Mult}(Z_u^d = t | \theta) p_{Dir}(\theta | Z_{-u}^d, \alpha) d\theta = \\ = \int_{\theta} \frac{\Gamma(\sum_{k=1}^K [n_{-u}^{d,k} + \alpha_k])}{\prod_{k=1}^K \Gamma(n_{-u}^{d,k} + \alpha_k)} \prod_{k=1}^K \theta_k^{n_{-u}^{d,k} + \alpha_k - 1} d\theta = \quad (6) \\ = \frac{\Gamma(\sum_{k=1}^K (n_{-u}^{d,k} + \alpha_k))}{\prod_{k=1}^K \Gamma(n_{-u}^{d,k} + \alpha_k)} \int_{\theta} \prod_{k=1}^K \theta_k^{n_{-u}^{d,k} + \alpha_k - 1} d\theta \end{aligned}$$

Нехай $\tilde{n}_{d,k} = n_{-u}^{d,k}$ для $k \neq t$ і $\tilde{n}_{d,k} = n_{-u}^{d,k} + 1$ для $k = t$. Тоді, враховуючи рівність $\int \prod_{k=1}^K \theta_k^{\alpha_k - 1} = \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^K \alpha_k)}$ отриману з визначення розподілу Діріхле і властивість Гамма-функції $\Gamma(\alpha + 1) = \alpha \Gamma(\alpha)$, отримуємо вираз (7) для імовірності тематики.

$$\begin{aligned} p(Z_u^d = t | Z_{-u}^d, \alpha) = \\ = \frac{\Gamma(\sum_{k=1}^K [n_{-u}^{d,k} + \alpha_k])}{\prod_{k=1}^K \Gamma(n_{-u}^{d,k} + \alpha_k)} \times \frac{\prod_{k=1}^K \Gamma(\tilde{n}_{d,k} + \alpha_k)}{\Gamma(\sum_{k=1}^K [\tilde{n}_{d,k} + \alpha_k])} = \quad (7) \\ = \frac{n_{-u}^{d,t} + \alpha_t}{\sum_{k=1}^K [n_{-u}^{d,k} + \alpha_k]} \end{aligned}$$

Аналогічним чином, розкривається другий множник. Зазначимо, що в цьому випадку доречніше було б писати $\Lambda_{-\pi}^d$ замість Λ_{-u}^d , щоб

⁴ Семплювання – процес отримання зразків деякої випадкової величини.

підкреслити те, що беремо всі шляхи окрім π – шляху, що веде від вершини онтології до слова W_u , але для однорідності нотації продовжимо використовувати Λ_{-u}^d , пам'ятаючи про альтернативне написання. Запишемо (8) – імовірність шляху $\pi = \{\lambda_1, \lambda_2, \dots, \lambda_L\}$, що дорівнює добутку по всім переходам шляху, $|\pi| = L, \lambda_j \in S$:

$$\rho(\Lambda_u^d = \pi | z_u^d = t, \Lambda_{-u}^d, \tau) = \prod_{l=0}^{L-1} \frac{b_{-u}^{t, \lambda_l, \lambda_{l+1}} + \tau_{\lambda_l, \lambda_{l+1}}}{\sum_{c=1}^{C_{\lambda_l}} [b_{-u}^{t, \lambda_l, c} + \tau_{\lambda_l, c}]} \quad (8)$$

Отже, отримаємо спільний розподіл призначення тематики і вибору шляху за умови зафіксованих тематик і шляхів для решти слів документу:

$$\rho(z_d^u = t, \Lambda_d^u = \pi | z_d^{-u}, \Lambda_d^{-u}, \tau, \alpha) \propto \frac{n_{-u}^{d,t} + \alpha_t}{\sum_{k=1}^K [n_{-u}^{d,k} + \alpha_k]} \left(\prod_{l=0}^{L-1} \frac{b_{-u}^{t, \lambda_l, \lambda_{l+1}} + \tau_{\lambda_l, \lambda_{l+1}}}{\sum_{c=1}^{C_{\lambda_l}} [b_{-u}^{t, \lambda_l, \lambda_{l+1}} + \tau_{\lambda_l, c}]} \right) \quad (9)$$

Схожість та пов'язаність

Як було сказано вище, для гіперпараметрів τ апріорними значеннями обираються значення їх інформаційного контенту, що у [14] було використано для визначення семантичної схожості синсетів наступним чином:

$$sim_{RES}(s, c) = IC(LCS(s, c))$$

де $LCS(s, c)$ – найближчий спільний родовий вузол синсетів s та c , тобто такий предок s та c , що має мінімальну кількість вузлів у шляхах від нього до s та c , відповідно. Інший приклад [10]:

$$sim_{LIN}(s, c) = \frac{2 \times IC(LCS(s, c))}{IC(s) + IC(c)}$$

Leacock Claudia та Chodorow Martin в [8] застосували ідею, що чим більша відстань між двома поняттями на графі онтології, тим менш вони семантично схожі:

$$sim_{LCH}(s, c) = -\log \frac{minPath(s, c)}{2D}$$

де D – діаметр таксономії, а $minPath(s, c)$ – довжина найкоротшого шляху між s та c .

Третій підхід до визначення семантичної схожості було запропоновано у [9], який базується на простій ідеї: чим більше спільних слів у словниковому описі двох понять, тим більше вони семантично схожі. Деякі результати порівняння цих методів і суджень людей про схожість слів можна знайти в [2].

Центральною проблемою з оцінкою семантичної схожості двох понять чи слів є визначення того, що ми розуміємо під схожістю. В літературі часто виділяють два поняття: схожості і пов'язаності. Два концепти можуть бути просто семантично пов'язаними, наприклад «двигун» і «коробка передач», тобто зустрівши одне з цих понять, зустріч іншого не буде несподіванкою. Але при цьому не можна сказати, що «двигун» і «коробка передач» у чомусь схожі, принаймні не більше ніж «двигун» і «мотор». Остання пара слів, взагалі кажучи є синонімами і це нам вказує на основну властивість схожості: схожі сутності в деякому контексті взаємозамінні. Зауважимо, для наведених мір «схожості» важко визначити чим вони є насправді мірами схожості чи пов'язаності, тому будемо розглядати вихід цих методів як оцінку семантичної пов'язаності.

Уявімо документ не як набір слів $W^d = (W_1^d, \dots, W_N^d)$, $W_i^d \in V$, а як набір синсетів $c^d = (c_1^d, \dots, c_N^d)$, $c_i^d \in S$. Тобто, у випадку LDAWN генеративний процес зупиняється, не генеруючи слів.

Тоді, виходячи з попередніх аргументів, семантична пов'язаність двох синсетів S_m та S_h пропорційна імовірності їх одночасної появи у певному контексті. Тобто, імовірності того, що випадкові величини c_u^d та c_v^d одночасно приймуть значення S_m та S_h в документі d . При цьому розглядаємо лише одну певну тематику t :

$$\begin{aligned} rel_{d,t}(S_m, S_h) &= p(c_u^d = S_m, c_v^d = S_h | c_{-u,v}^d) = \quad (10) \\ &= p_{LDAWN}(\Lambda_u^d = \pi_m | \Lambda_{-u}^d, z_u = t) \times \\ &\quad \times p_{LDAWN}(\Lambda_v^d = \pi_h | \Lambda_{-u,v}^d, z_v = t) = \\ &= \prod_{j=0}^{L_m-1} \rho_t(\lambda_{j+1}^m | \lambda_j^m, \Lambda_{-u}^d) \prod_{j=0}^{L_h-1} \rho_t(\lambda_{j+1}^h | \lambda_j^h, \Lambda_{-u,v}^d) \end{aligned}$$

де, згідно з результатом (8) для апостеріорного розподілу шляхів Λ_u^d :

$$\prod_{j=0}^{L_m-1} \rho_t(\lambda_{j+1}^m | \lambda_j^m, \mathbf{A}_{-u}^d) = \prod_{j=0}^{L_m-1} \frac{b_{-u,v}^{t,\lambda_j^m, \lambda_{j+1}^m} + \tau_{\lambda_j^m, \lambda_{j+1}^m}}{\sum_{c=1}^c \lambda_j^m \left[b_{-u,v}^{t,\lambda_j^m, c} + \tau_{\lambda_j^m, c} \right]} \quad (11)$$

Другий множник (10) розкладається аналогічно. Введемо наступні позначення:

$$\pi_{m,h} = \pi_m \cap \pi_h = \{ \lambda_0 = s_0, \dots, \lambda_{m,h} = LCS(s_m, s_h) \} \text{ та}$$

$$\pi_{m \setminus h} = \pi_m \setminus \pi_h = \{ \lambda_{m \setminus h} = LCS(s_m, s_h), \dots, \lambda_{L_m} = s_m \},$$

$\pi_{h \setminus m}$ аналогічно. Тоді:

$$rel_{d,t}(s_m, s_h) = \rho(\pi_{m,h} | \mathbf{A}_{-u,v}^d)^2 \rho(\pi_{m \setminus h} | \mathbf{A}_{-u}^d) \rho(\pi_{h \setminus m} | \mathbf{A}_{-u,v}^d) \quad (12)$$

Відмітимо, що для отримання граничної пов'язаності s_m та s_h потрібно просумувати по всіх можливих шляхах до цих синсетів та всіх контекстах $d=1:D$.

Тематична модель пов'язаних концептів

Як було показано вище, апостеріорна імовірність вибору певного синсету $\rho_t(c_u^d = s_m \Leftrightarrow \Lambda_u^d = \pi_m | \mathbf{A}_{-u}^d)$ в LDAWN, а також сумісний розподіл двох синсетів $\rho_t(c_u^d = s_m, c_v^d = s_h | \mathbf{A}_{-u,v}^d) = rel_t(s_m, s_h)$ залежить лише від кількості спостережуваних переходів у тематичному випадковому блуканні між синсетами, що входять до шляхів π_m та π_h .

Наявність чи відсутність сиблінгів синсету s (вузлів таксономії того ж рівня, що і s) не впливає на його розподіл. Тому розглянемо наступну генеративну тематичну модель, в якій буде можливо врахувати інформацію про «відстані» між синсетами, що обчислюються за допомогою мір схожості таких як LIN чи LESK.

Розіб'ємо множину синсетів S на рівні: $S_\delta = \{s \in S : depth(s) = \delta, \delta = 0:\Delta\}$

Генеративний процес:

1) Для кожної тематики $t \in 1:K$

а. Для кожного рівня $\delta=0:\Delta$ обираємо розподіл для множини випадкових величин, що відповідають синсетам, $s_\delta^t \sim N(s_\delta | \mu_\delta^t(\mu_{\delta-1}^t), \Sigma_\delta)$, де Σ_δ – матриця коваріацій, чисе апіорне значення відповідає rel_δ – матриці попарних схожостей синсетів рівня δ , а μ_δ – вектор середніх значень синсетів рівня δ , що в загальному випадку залежить від $\mu_{\delta-1}$: чим більше $\mu_{\delta-1,s}$ тим більші $\mu_{\delta,c}$, де $c \in$ нащадком $s : c \in c(s)$.

2) Для кожного документу $w^d \in D$, $d=1:D$:

а. Обираємо тематичний розподіл $\theta^d \sim \text{Dir}^K(\alpha)$.

б. Обираємо розподіл рівнів $G^d \sim \text{Dir}^\Delta(\gamma)$.

с. Для кожного слова w_u^d , $u=1:N_d$:

i. Обираємо тематику $z_u^d \sim \text{Mult}^K(1, \theta^d)$

ii. Обираємо рівень $g_u^d \sim \text{Mult}^\Delta(1, G^d)$

iii. Обираємо синсет $c_u^d \sim \text{Mult}^S\left(1, \sigma\left(s_{g_u^d}^{z_u^d} | \mu_{g_u^d}^{z_u^d}, \Sigma_{g_u^d}^{z_u^d}\right)\right)$

iv. Генеруємо слово w_u^d з обраного синсету c_u^d .

де $\sigma_i(\alpha) = \frac{\exp(\alpha_i)}{\sum_j \exp(\alpha_j)}$ – softmax функція, яка

приймає значення в інтервалі $[0;1]$. Тобто імовірність, що u -те слово буде згенеровано з m -того синсету S_m :

$$\rho(c_u^d = s_m | \mu_\delta^t, \Sigma_\delta^t) = \sigma_m(s_\delta^t | \mu_\delta^t, \Sigma_\delta^t) \quad (13)$$

Повна апостеріорна імовірність прихованих параметрів:

$$\rho(\theta^d, G^d, g^d, z^d, c^d, s^d, \mu^d | w^d, \Sigma, \alpha, \gamma) = \frac{\rho(w^d | \theta^d, G^d, g^d, z^d, c^d, \mu, \Sigma, \alpha, \gamma)}{\int_{\theta, G, g, s, c, z} \rho(w^d | \theta^d, G^d, g^d, z^d, c^d, \mu, \Sigma, \alpha, \gamma)} \quad (14)$$

Запишемо функцію правдоподібності для документу \mathbf{w}^d :

$$\begin{aligned} & p(\mathbf{w}^d | \boldsymbol{\theta}^d, \mathbf{G}^d, \mathbf{g}^d, \mathbf{z}^d, \mathbf{c}^d, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\alpha}, \boldsymbol{\gamma}) = \\ & = p(\boldsymbol{\theta}^d | \boldsymbol{\alpha}) p(\mathbf{G}^d | \boldsymbol{\gamma}) \times \\ & \times \prod_{t=1}^K p(\mathbf{s}^t | \boldsymbol{\mu}^t, \boldsymbol{\Sigma}^t) \prod_{u=1}^{N_d} p(w_u^d | c_u^d) \times \\ & \times p(c_u^d | \sigma(s_{g_u^d}^z)) p(g_u^d | \mathbf{G}^d) p(z_u^d | \boldsymbol{\theta}^d) \end{aligned} \quad (15)$$

Вивід апостеріорного розподілу для рівнів \mathbf{g}^d і пропорції рівнів \mathbf{G}^d , схожий на такий для тематик \mathbf{z}^d та пропорцій тематик у документі $\boldsymbol{\theta}^d$, бо вони мають однакову структуру розподілу: $p_{\text{Mult}}(\mathbf{y} | \boldsymbol{\theta}) p_{\text{Dir}}(\boldsymbol{\theta} | \boldsymbol{\alpha})$. Для пропорції рівнів легко отримати оцінку максимальної правдоподібності $\hat{G}_\delta = \frac{n_\delta}{N_d}$, де n_δ – кількість спостережуваних синсетів рівня δ .

Розглянемо вивід для розподілу синсетів при фіксованій тематиці t . Нехай ми спостерігаємо документ у вигляді синсетів $\mathbf{c}^d = (c_1^d, \dots, c_{N_d^t}^d)$, $c_i \in S: c_i = S_m \Leftrightarrow c_i = (0, \dots, 1, \dots, 0_S)$; верхній індекс t в N_d^t використано, щоб підкреслити, що розглядаються лише ті синсети для яких було обрано тематику t , але далі опустимо цей індекс, щоб спростити запис. Згідно опису моделі

$$\begin{aligned} & \mathbf{c}^d \sim \text{Mult}(N_d, \sigma(\mathbf{s})) = \\ & = \frac{N_d!}{c^{(1)}! c^{(2)}! \dots c^{(s)}!} p_1^{c^{(1)}} \cdot p_2^{c^{(2)}} \cdot \dots \cdot p_s^{c^{(s)}} \end{aligned} \quad (16)$$

де $p_m = \sigma(s_m)$, а $c^{(m)}$ – кількість таких c_u^d , що $c_u^d = s_m$. Пригадаємо, що $\mathbf{s} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ і тоді апостеріорний розподіл синсетів:

$$p(\mathbf{s} | \mathbf{c}) = \frac{p(\mathbf{c} | \mathbf{s}) p(\mathbf{s} | \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\int_{\mathbf{s}} p(\mathbf{c} | \mathbf{s}) p(\mathbf{s} | \boldsymbol{\mu}, \boldsymbol{\Sigma})} \quad (17)$$

Цей інтеграл не обчислюється аналітично, із-за входження softmax-функції, тому апроксимуємо його методом Лапласа

нормальним розподілом $q(\mathbf{s})$. Для цього запишемо логарифмічну функцію правдоподібності:

$$\begin{aligned} & L = -\ln p(\mathbf{s} | \mathbf{c}) = \\ & = \frac{1}{2} (\mathbf{s} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{s} - \boldsymbol{\mu}) - \\ & - \sum_{u=1}^{N_d} \sum_{m=1}^s c_u^m \ln \sigma(s_m) + \text{const} \end{aligned} \quad (18)$$

Тепер розглянемо першу і другу похідну другого доданку:

$$\begin{aligned} & \nabla_{s_j} \left(\sum_{m=1}^s c^{(m)} \ln \sigma(s_m) \right) = \\ & = - \sum_{m=1}^s c^{(m)} (I_{m,j} - \sigma(s_j)) \\ & = \begin{cases} N_d \sigma(s_j), & m \neq j \\ N_d \sigma(s_j) - c^{(m)}, & m = j \end{cases} \end{aligned} \quad (19)$$

$$\begin{aligned} & \nabla_{s_k} \nabla_{s_j} L = \nabla_{s_k} \left(- \sum_{m=1}^s c^{(m)} (I_{m,j} - \sigma(s_j)) \right) = \\ & = N_d \nabla_{s_k} \sigma(s_j) = N_d \sigma'(s_j) (I_{j,k} - \sigma(s_k)) \end{aligned} \quad (20)$$

Тоді:

$$\nabla \nabla L = \boldsymbol{\Sigma}^{-1} + N_d (\sigma(\mathbf{s}) \mathbf{I} - \sigma(\mathbf{s}) \sigma^\top(\mathbf{s})) \quad (21)$$

Тепер можемо записати апроксимуючий розподіл $q(\mathbf{s})$ [3]:

$$q(\mathbf{s}) = N(\mathbf{s} | \boldsymbol{\mu}_{\text{MAP}}, \boldsymbol{\Sigma}_N) \quad (22)$$

де $\boldsymbol{\Sigma}_N^{-1} = \nabla \nabla L$, а $\boldsymbol{\mu}_{\text{MAP}}$ максимальна апостеріорна оцінка, яку можна отримати з ∇L .

Розглянемо тепер дві випадкові величини c_u^d та c_v^d , що відповідають u -тому та v -тому слову і мають значення s_m та s_h відповідно. Нехай для них було обрано однакову тематику $z_u^d = t$, $z_v^d = t$, і обрано рівні $g_u^d = \delta_u$ та $g_v^d = \delta_v$ відповідно. Запишемо вираз для пов'язаності синсетів s_m та s_h , тобто сумісний розподіл

$p_t(c_u^d = s_m, c_v^d = s_h | \mu^d, \Sigma)$, де μ^d – апостеріорне значення середніх для розподілу синсетів, $s^d \sim N(\mu^d, \Sigma)$ після спостереження частини документу $W_{-u,v}^d$, тобто контексту для c_u^d я u та v :

$$rel_{d,t}(s_m, s_h) = p_t(c_u^d = s_m, c_v^d = s_h | \mu^d, \Sigma) \approx \sigma(s_m)\sigma(s_h)q(s^d | \mu^d, \Sigma) \quad (23)$$

$$\text{де } (\Sigma^d)^{-1} = \Sigma^{-1} + (N_d - 2)(\sigma(s)I - \sigma(s)\sigma^T(s))$$

Висновок

Таким чином, в запропонованій моделі дистрибутивна пов'язаність прямо залежить від значення структурної пов'язаної – матриці Σ .

Список використаних джерел

- [1] *E. Agirre, P. Edmonds* Word Sense Disambiguation: Algorithms and Applications. - Dordrecht, Netherlands : Springer, 2006. - p. 388.
- [2] *Anisimov A.V., Lyman K.S., Marchenko O.O.* The computational methods for the semantic proximity measures of natural language words // Artificial Intelligence (Штучний інтелект). - 2010. - 3. - pp. 170-176.
- [3] *Christopher M. Bishop* Pattern recognition and Machine learning. - : Springer, 2006.
- [4] *D. Blei, A. Ng, M. Jordan* Latent dirichlet allocation // Journal of Machine Learning Research. - : JMLR.org, 2003. - сс. 993-1022.
- [5] *J. Boyd-Graber, D. Blei, X. Zhu* A topic model for word sense disambiguation // Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL). - 2007. - pp. 1024-1033.
- [6] *S. Brody, M. Lapata* Bayesian word sense induction // Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics. - : Association for Computational Linguistics, 2009. - срп. 103-111.
- [7] *C. Chemudugunta et al.* Modeling Documents by Combining Semantic Concepts with Unsupervised Statistical Learning // The Semantic Web - ISWC 2008 - Berlin : Springer, 2008. - Vol. 5318.
- [8] *C. Fellbaum ed.* WordNet: an electronic lexical database.. - Cambridge, MA : The MIT Press, 1998.
- [9] *M. Lesk* Automatic sense disambiguation using machine readable dictionaries: How to tell a pine cone from a ice cream cone. // Proceedings of SIGDOC '86.. - 1986. - срп. 24-26.
- [10] *D. Lin* Automatic retrieval and clustering of similar words // Proceedings of COLING-ACL 98. - Montreal, Canada :, 1998.
- [11] *D. McCarthy* Word Sense Disambiguation: An Overview // Language and Linguistics compass. - : Wiley Online Library, 2009 - 2 : Vol. 3. - срп. 537-558.
- [12] *R. Navigli* Word Sense Disambiguation: A Survey // ACM Computing Surveys (CSUR). - : ACM, 2009. - 2 : Vol. 41. - pp. Article 10.
- [13] *B. Snyder, M. Palmer* The English all-words task // Proceedings of the ACL senseval-3 workshop. - Barcelona, Spain : ACL, 2004. - срп. 41-43.
- [14] *P. Resnik* Using Information Content to Evaluate Semantic Similarity in a Taxonomy // Proceedings of the 14th International Joint Conference on Artificial Intelligence. - 1995. - pp. 448-453.

Надійшла до редколегії 12.02.2013