

МНОГОФАКТОРНЫЙ МЕТОД ПОСТРОЕНИЯ ТЕМАТИЧЕСКОГО ПРЕДСТАВЛЕНИЯ НОВОСТНОГО КЛАСТЕРА

Аспирант ВМК
Алексеев Алексей

ОБРАБОТКА ПОТОКОВ НОВОСТЕЙ

2

- ❑ Новостные сервисы (30-40 тыс. документов в день)
- ❑ Кластеризация новостей на одну тему – новостной кластер (вхождения слов)
 - Удаление дубликатов
 - Рубрикация по тематическим рубрикам
 - Автоматическое аннотирование
 - Определение новизны
 - Извлечение информации
- ❑ Многие операции выполняются на основе пословного представления

ПРОБЛЕМЫ ПОСЛОВНОГО ПРЕДСТАВЛЕНИЯ

3

- ❑ Одна сущность названа посредством цепочки слов (многословным выражением)
- ❑ В кластере используется много разных наименований одной и той же сущности
- ❑ Авиабазы США в Киргизии:
 - база Манас, авиабаза Манас, Манас,
 - база в международном аэропорту Манас,
 - база США, американская авиабаза
- ❑ Проблемы:
 - Определение границ кластера
 - Автоматическое порождение аннотации
 - Определение новизны информации
 - Выделение подкластеров и др.

ПОСТАНОВКА ЗАДАЧИ - 1

4

- ❑ Кластер документов должен соответствовать ситуации или совокупности связанных ситуаций (основная тема кластера)
- ❑ В ситуации есть набор участников, которые в кластере:
 - Могут быть выражены не только словами, но и словосочетаниями
 - Каждый участник может выражаться совокупностью разных выражений
- ❑ Тестирование показывает, что качество кластеров оценивается около **70%** (Добров Б.В, Павлов А.С. 2010), т.е. основные участники в среднем представлены правильно
- ❑ **Предположение:** более точное выявление основных участников ситуации поможет более правильно осуществлять разные операции с кластерами

ПОСТАНОВКА ЗАДАЧИ - 2

5

❑ **Входные данные:**

Кластер тематически близких документов, собранных на основе пословной модели представления документов

❑ **Задача:**

Построить тематическое представление входного новостного кластера – совокупность основных участников ситуации и отношений между ними

❑ **Применение:**

- Добавление в пословное представление документа выделенных структур
- Использование усложненного представления документов кластера для улучшения качества операций с кластером

ПЛАН ДОКЛАДА

6

- **Природа вариативности и обзор существующих методов**
- ❑ Тематическое представление текста
- ❑ Метод выделения структур, описывающих участников ситуации
- ❑ Тестирование подхода в автоматическом аннотировании
- ❑ Заключение

ПРИРОДА ВОЗНИКНОВЕНИЯ ВАРИАТИВНОСТИ - 1

7

□ Цель использования:

- **Референция** (отнесенность языкового выражения к одному и тому же объекту действительности)

3 февраля президент Киргизии Курманбек Бакиев заявил о решении правительства прекратить деятельность авиабазы на территории республики... Президент не стал скрывать, что экономические резоны стали главной причиной побудившей правительство страны принять такое решение.

- **Перефразирование** (изменение текста без изменения смысла - рерайтинг)

Судьбу авиабазы США в "Манасе" решил парламент Киргизии. Парламент Киргизии в четверг примет окончательное решение о судьбе авиабазы США.

ПРИРОДА ВОЗНИКНОВЕНИЯ ВАРИАТИВНОСТИ - 2

8

□ Привязка к контексту:

➤ *Общеизвестно (Киргизия – Кыргызстан)*

Правительство Киргизии передало для ратификации в законодательный орган... Парламент Кыргызстана в четверг примет окончательное решение о судьбе авиабазы США... стали главной причиной побудившей правительство страны принять такое решение

➤ *Выводится из контекста*

В декабре 2006 года 46-летний водитель топливозаправщика киргизской фирмы, занимающейся обслуживанием аэропорта "Манас", Александр Иванов, был расстрелян в упор охранником авиабазы Закари Хатфилдом на КПП при въезде на перрон аэропорта"... Американский военный, несмотря на неоднократные требования киргизского МИДа, также был тайно вывезен с территории страны и до сих пор не предстал перед судом.

УСЛОЖНЕНИЕ ПОСЛОВНОГО ПРЕДСТАВЛЕНИЯ

ЛЕКСИЧЕСКИЕ ЦЕПОЧКИ - 1

9

- Использование информации об объектах и связях между ними, описанной в определенных ресурсах

- **Популярные подходы:**

- *Английский язык: **Wordnet***

Barzilay R., Elhadad M., 1999

- *Русский язык: **PyTез***

Лукашевич Н.В., Добров Б.В., 2009

Лексические цепочки в форме **тематических узлов** - к одному выделенному элементу относятся все другие элементы лексической цепочки

УСЛОЖНЕНИЕ ПОСЛОВНОГО ПРЕДСТАВЛЕНИЯ

ЛЕКСИЧЕСКИЕ ЦЕПОЧКИ - 2

10

Пример основных тематических узлов на основе РунТез:

ТЕННИСНЫЙ КОРТ	14
ТЕННИС	12
АВСТРИЙЦЫ	12
АВСТРИЯ	6
КИПРИОТЫ	16
КИПР	11
ХОРВАТЫ	10
СЕТ (ПАРТИЯ В ТЕННИСЕ)	6
ИГРОВАЯ ПАРТИЯ	5
ЧЕТВЕРТЬФИНАЛ	10
ПОЛУФИНАЛ	29
ПОЛУФИНАЛИСТ	2

МАТЧ	12
СПОРТИВНЫЙ ФИНАЛ	36
СПОРТИВНОЕ СОРЕВНОВАНИЕ	54
СПОРТ	8
СПОРТСМЕН	2
ФИНАЛИСТ	1
ЮЖНЫЙ, МИХАИЛ	23
РОССИЯНЕ	12
РОССИЙСКАЯ ФЕДЕРАЦИЯ	10
ТЕННИСИСТ	6
ЗАГРЕБ	70
ХОРВАТИЯ	36

УСЛОЖНЕНИЕ ПОСЛОВНОГО ПРЕДСТАВЛЕНИЯ РЕФЕРЕНЦИАЛЬНЫЕ ЦЕПОЧКИ

11

- ❑ Выявление различных упоминаний (отсылок) к одной и той же **реальной** сущности
- ❑ Учет различных видов отсылок, в том числе местоимений

Пример:

Президент РФ положительно оценил отчет премьера в Госдуме

За выступлением премьера в Думе наблюдал Владимир Путин. Он просмотрел большую часть выступления Дмитрия Медведева и в целом дал позитивную оценку отчёту главы правительства, сообщил журналистам пресс-секретарь Президента Дмитрий Песков.

- ❑ Узкая специализация: выявление строго-определенного спектра сущностей – люди, организации, места и т.д.

Latent Dirichlet Allocation (LDA)

University of California, USA, 2003
David M. Blei, Andrew Y. Ng, Michael I. Jordan

12

- ❑ Вероятностная модель, разбивающая слова на тематические группы
- ❑ Двухуровневое представление:
 - Распределение Topics-VS-Documents (TvsD)
 - Распределение Words-VS-Topics (WvsT)
- ❑ Наблюдаемый документ – автоматический семплинг из указанных распределений:
 - I. Выбор топика из TvsD
 - II. Выбор слова из WvsT
- ❑ **Проблемы:**
 - Вероятностный результат
 - Сложность интерпретации результатов работы алгоритма

Latent Dirichlet Allocation (LDA)

ПРИМЕР

13

“Arts”

“Budgets”

“Children”

“Education”

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services,” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

ПЛАН ДОКЛАДА

14

- ❑ Природа вариативности и обзор существующих методов
 - **Тематическое представление текста**
- ❑ Метод выделения структур, описывающих участников ситуации
- ❑ Тестирование подхода в автоматическом аннотировании
- ❑ Заключение

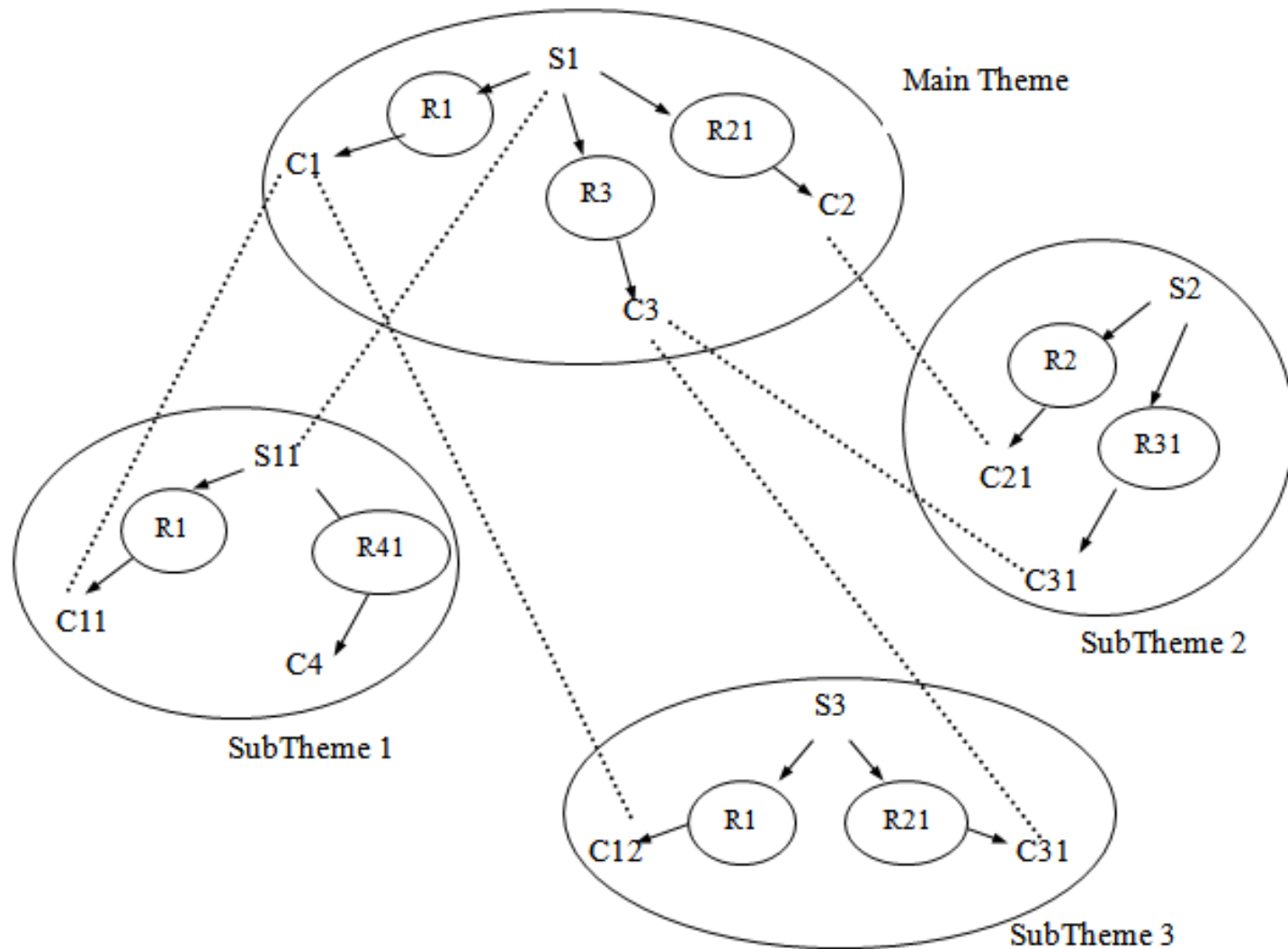
ГЛОБАЛЬНАЯ СВЯЗНОСТЬ ТЕКСТА

15

- ❑ **Van Dijk и гипотеза глобальной связности (1985)**
- ❑ Связный текст имеет одну главную тему и эта тема может быть выражена как пропозиция
- ❑ Тема целого текста раскрывается в тексте посредством локальных тем
- ❑ Каждое предложение текста соответствует некоторой теме текста
- ❑ Механизм глобальной связности позволяет контролировать локальные связи и переходы

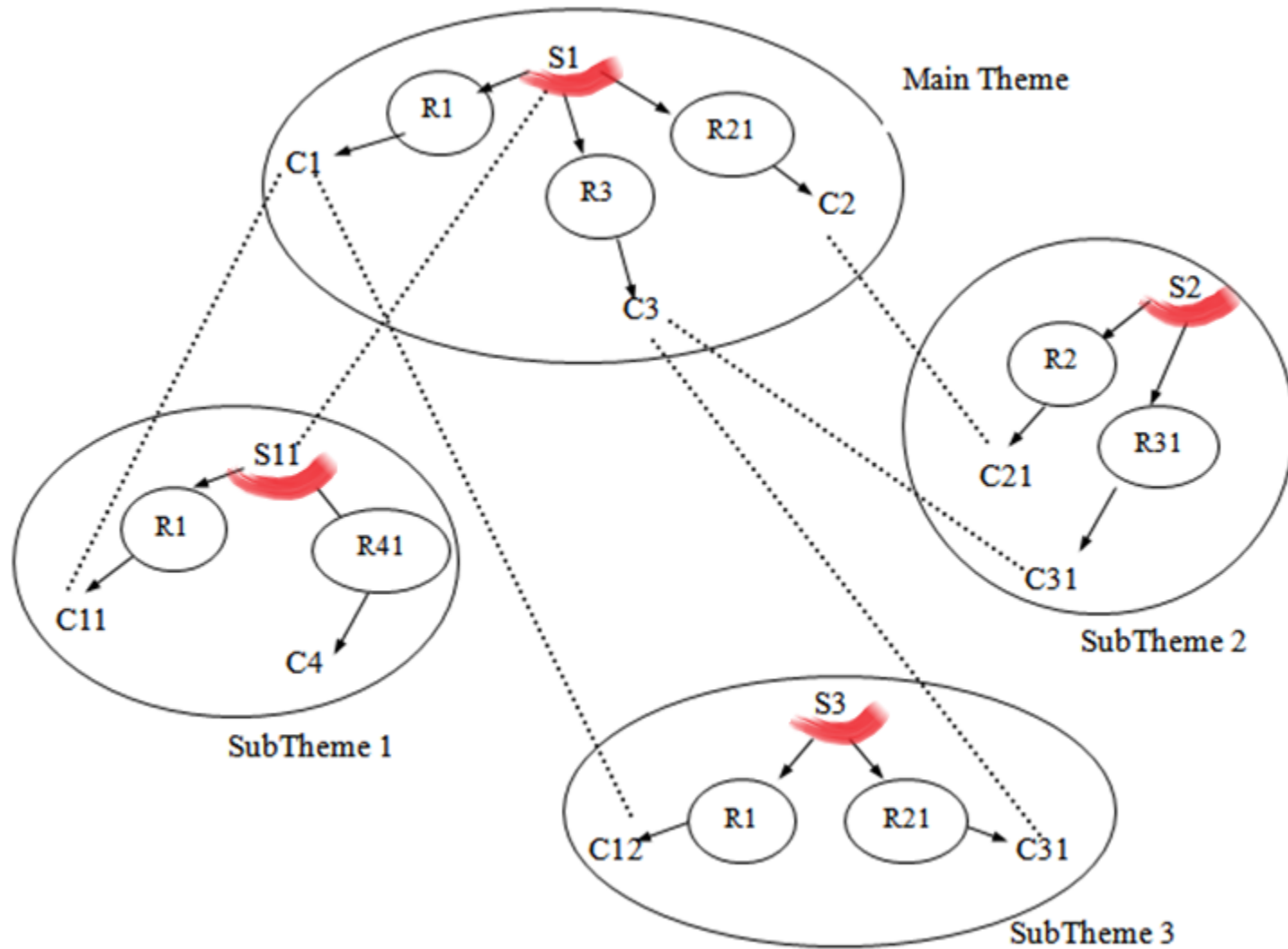
ГЛАВНАЯ ТЕМА И ЕЕ УТОЧНЕНИЕ В КОНКРЕТНЫХ ПРЕДЛОЖЕНИЯХ

16



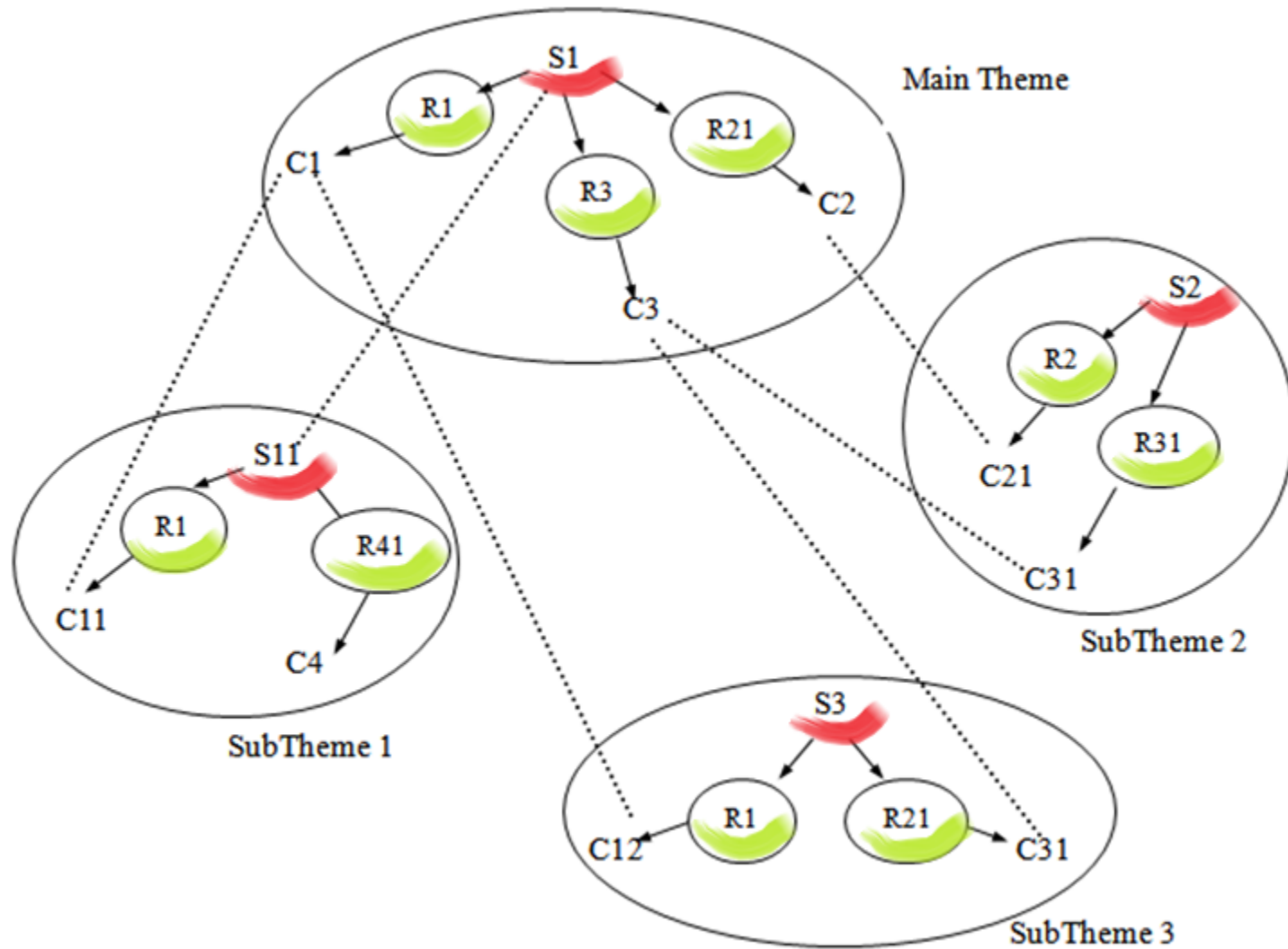
ГЛАВНАЯ ТЕМА И ЕЕ УТОЧНЕНИЕ В КОНКРЕТНЫХ ПРЕДЛОЖЕНИЯХ

17



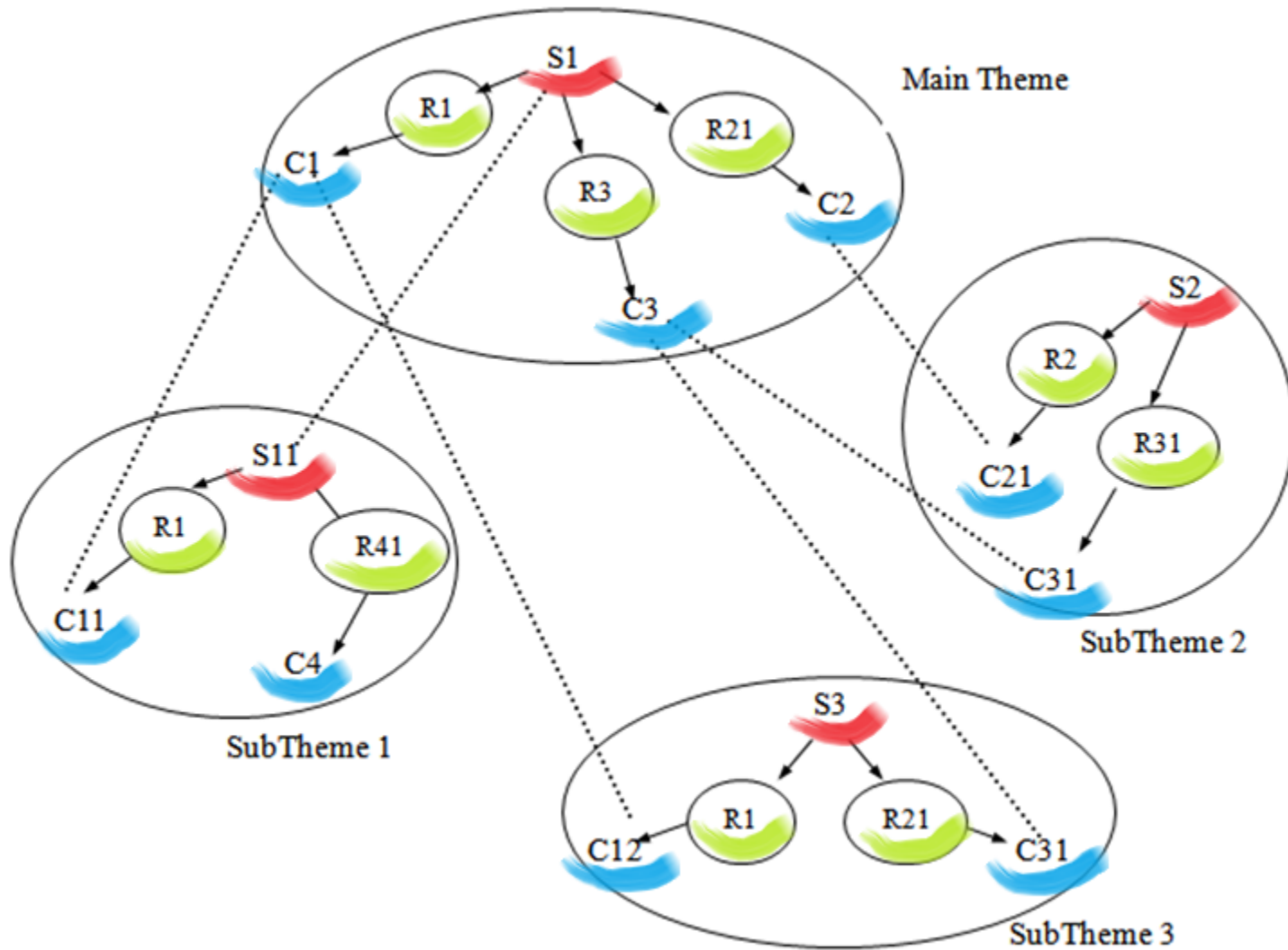
ГЛАВНАЯ ТЕМА И ЕЕ УТОЧНЕНИЕ В КОНКРЕТНЫХ ПРЕДЛОЖЕНИЯХ

18



ГЛАВНАЯ ТЕМА И ЕЕ УТОЧНЕНИЕ В КОНКРЕТНЫХ ПРЕДЛОЖЕНИЯХ

19



ЛЕКСИЧЕСКАЯ СВЯЗНОСТЬ vs. ГЛОБАЛЬНАЯ СВЯЗНОСТЬ

20

- ❑ Связный текст обладает лексической связностью: лексические и семантические повторы
- ❑ Лексическая связность – инструмент глобальной связности
- ❑ Чем больше две сущности упоминаются в одних и тех же предложениях текста, тем более важно отношение между ними для содержания текста
- ❑ Если сущности часто упоминаются, но мало встречаются в одних и тех же предложениях текста, то возможно они связаны между собой по смыслу (синоним, род-вид, часть-целое)
- ❑ Элементы тематического представления - имеют внутреннюю структуру тематического узла (ТУ): главное выражение и относящиеся к нему объекты

ПРОВЕРКА ГИПОТЕЗЫ ВСТРЕЧАЕМОСТИ В ПРЕДЛОЖЕНИЯХ - 1

21

- ❑ Проверка предположений была произведена с помощью Тезауруса русского языка РуТез
- ❑ В качестве правильных примеров (связанных языковых выражений) рассматривались выражения имеющие связь по Тезаурусу
- ❑ Различные типы связи рассматривались отдельно (*синонимия; род-вид; часть-целое*)
- ❑ Две группы по частям речи:

СУЩ. + СУЩ.

#

ПРИЛ. + СУЩ.

- ❑ Для каждой пары объектов вычислялись количество вхождений в одни и те же предложения (***Fsegm***) и в соседние (***Fsent***)

ПРОВЕРКА ГИПОТЕЗЫ ВСТРЕЧАЕМОСТИ В ПРЕДЛОЖЕНИЯХ - 2

22

Тип связи	Fsegm / Fsent	Число пар
Синонимы (СУЩ + СУЩ)	0.309	31
Синонимы (ПРИЛ + СУЩ)	0.491	53
Род – Вид (СУЩ + СУЩ)	1.130	88
Род – Вид (ПРИЛ + СУЩ)	1.471	28
Часть – Целое (СУЩ + СУЩ)	0.779	58
Часть – Целое (ПРИЛ + СУЩ)	1.580	29
Без связи по Тезаурусу	1.440	21483

Вывод: превышение Fsent по отношению к Fsegm свидетельствует о связи между объектами

НОВОСТНЫЕ КЛАСТЕРЫ И СВОЙСТВА СВЯЗНОГО ТЕКСТА

23

- ❑ Кластер – не является связным текстом,
 - но имеет тему кластера
 - статистические особенности усиливаются
- ❑ Извлечение
 - Многословных выражений,
 - Структур, описывающих основных участников ситуации (тематических узлов)
- ❑ Пример: Новостной кластер от 03.02.2007
 - Тема: Отставка президента алмазодобывающей компании АК «Алроса» Александра Ничипорука
 - 12 новостных документов

ПЛАН ДОКЛАДА

24

- ❑ Природа вариативности и обзор существующих методов
- ❑ Тематическое представление текста
 - **Метод выделения структур, описывающих участников ситуации**
- ❑ Тестирование подхода в автоматическом аннотировании
- ❑ Заключение

ПРИНЦИПЫ ФОРМИРОВАНИЯ ТЕМАТИЧЕСКИХ УЗЛОВ

25

- ❑ Структуры, соответствующие участникам ситуации, собираются в форме тематических узлов:
 - центр и элементы ("похожие" на центр слова)
 - элементы могут принадлежать нескольким (двум) тематическим узлам
- ❑ Элементы тематического узла могут быть "похожи" на центр по нескольким параметрам
- ❑ Если два выражения часто употребляются в одном предложении, они не могут оказаться в одном тематическом узле - ограничивающий фактор
- ❑ Предварительная сборка многословных выражений

ИЗВЛЕЧЕНИЕ МНОГОСЛОВНЫХ ВЫРАЖЕНИЙ

26

- ❑ Рассматривается как необходимый этап для построения тематических узлов
- ❑ Превышение встречаемости слов непосредственно рядом друг с другом по сравнению с отдельной встречаемостью во фрагментах предложений:

$$\text{Near} > 2 * (\text{Av} + \text{NotN})$$

- ❑ Дополнительные ограничения по частотности встречаемости слов рядом друг с другом
- ❑ **Примеры:** *Президент Алроса; президент компании; владелец компании; контрольный пакет акций; акции компании; пакет акций*

АЛГОРИТМ ПОСТРОЕНИЯ ТЕМАТИЧЕСКОГО ПРЕДСТАВЛЕНИЯ

27

- ❑ Разные признаки сходства:
 - **Контекстно-независимые признаки** (формальное сходство, тезаурус)
 - **Контекстно-зависимые признаки** («строгие» контексты, статистика контекстов)
- ❑ 4 типа контекстов: через глагол (**AV**); рядом (**Near**); не рядом (**NN**); в соседних предложениях (**NS**)
- ❑ Итеративный алгоритм. Одна итерация = объединение одной пары объектов
- ❑ **Управляющее правило:**
$$NS \geq 2 * (AV + Near + NN)$$

КОНТЕКСТНО-НЕЗАВИСИМЫЕ ПРИЗНАКИ

28

А. Формальное сходство объектов (BS)

- Простая метрика – одинаковые начала слов
- Пример: *Руководитель* – *Руководство*
- Вес признака: [0, 1]

$$BS = 1.0 - 0.1 * (\{ \text{Число_слов_с_отличными_началами} \})$$

В. Наличие связи в тезаурусе (TS)

- Рассмотрение цепочек отношений между объектами (использование вывода связей)
- Пример: *Руководитель* – *Глава*
- Вес признака: [0, 1]

$$TS = 1.0 - 0.2 * (\{ \text{Длина_пути_по_отношениям_тезауруса} \})$$

КОНТЕКСТНО-ЗАВИСИМЫЕ ПРИЗНАКИ

29

A. Частота появления в соседних предложениях (NSF)

$$NSF = \text{Min} \left[1, \frac{C}{\text{Avg}(C)} \right] \quad C = NS - 2 * (AcrossVerb + Near + NotNear)$$

B. Анализ «строгих» контекстов (SC)

- 4-граммы: по 2 слова влево и вправо
- Значение: $[0, 1]$ (относительно «лучшей» пары)

C. Скалярные произведения контекстов (SPS)

- Косинусная мера угла между векторами контекстов
- Сумма по характеристикам внутри предложения (AV + Near + NotNear)

ПРИМЕР РАНЖИРОВАНИЯ ПАР

30

<div>Features</div> <div>Pairs</div>	Context-independent		Context-dependent			SCORE
	BS	TS	NSF	SC	SPS	
Президент России – Президент РФ	0.90	1.00	0.00	0.50	0.68	2.74
Инвестгруппа– Инвестиционная группа	0.90	1.00	0.40	0.00	0.63	2.42
ГМК Норильский никель – Норильский никель	1.00	1.00	0.40	0.00	0.21	2.31
Российская Федерация – Россия	0.90	1.00	0.00	0.00	0.51	2.15
Отставка – Отставка с должности	0.90	1.00	0.40	0.00	0.00	2.10

ПРИМЕР СБОРКИ УЗЛА

31

- ❑ Итерация 7: (*Отставка*) \leftarrow (*Отставка с должности*)
- ❑ Итерация 33: (*Отставка, Отставка с должности*) \leftarrow (*Уход в отставку*)
- ❑ Итерация 44: (*Отставка, Отставка с должности, Уход в отставку*) \leftarrow (*Отставка президента*)
- ❑ Итерация 61: (*Уход с поста*) \leftarrow (*Уход в отставку*)
- ❑ Итерация 62: (*Отставка, Отставка с должности, Уход в отставку, Отставка президента*) \leftarrow (*Уход с поста, Уход в отставку*)
- ❑ Итерация 102: (*Отставка, Отставка с должности, Уход в отставку, Отставка президента, Уход с поста*) \leftarrow (*Пост*)
- ❑ Итерация 103: (*Пост, Отставка, Отставка с должности, Уход в отставку, Отставка президента, Уход с поста*) \leftarrow (*Должность*)
- ❑ Итерация 104: (*Пост, Отставка, Отставка с должности, Уход в отставку, Отставка президента, Уход с поста, Должность*) \leftarrow (*Уход*)

РЕЗУЛЬТАТЫ РАБОТЫ АЛГОРИТМА НА КЛАСТЕРЕ ПРИМЕРЕ

32

- ❑ **Пост**: уход с поста; должность; уход; отставка; отставка с должности; уход в отставку; отставка президента
- ❑ **Алроса**: президент Алроса; АК Алроса
- ❑ **Компания**: акция компании; владелец компании; объединение компаний; акция; акционер компании; владелец; пакет акций; состав владельцев; контрольный пакет акций; контрольный пакет; владение
- ❑ **Ничипорук**: Александр Ничипорук
- ❑ **Якутия**: президент Якутии; якутский; якутский президент

ПЛАН ДОКЛАДА

33

- ❑ Природа вариативности и обзор существующих методов
- ❑ Тематическое представление текста
- ❑ Метод выделения структур, описывающих участников ситуации
 - **Тестирование подхода в автоматическом аннотировании**
- ❑ Заключение

ПРИМЕНЕНИЕ ТЕМАТИЧЕСКОГО ПРЕДСТАВЛЕНИЯ

34

- Тематическое представление – дополнительная информация об устройстве новостного кластера

Слова → Объекты (слова + MB) → Тематические узлы

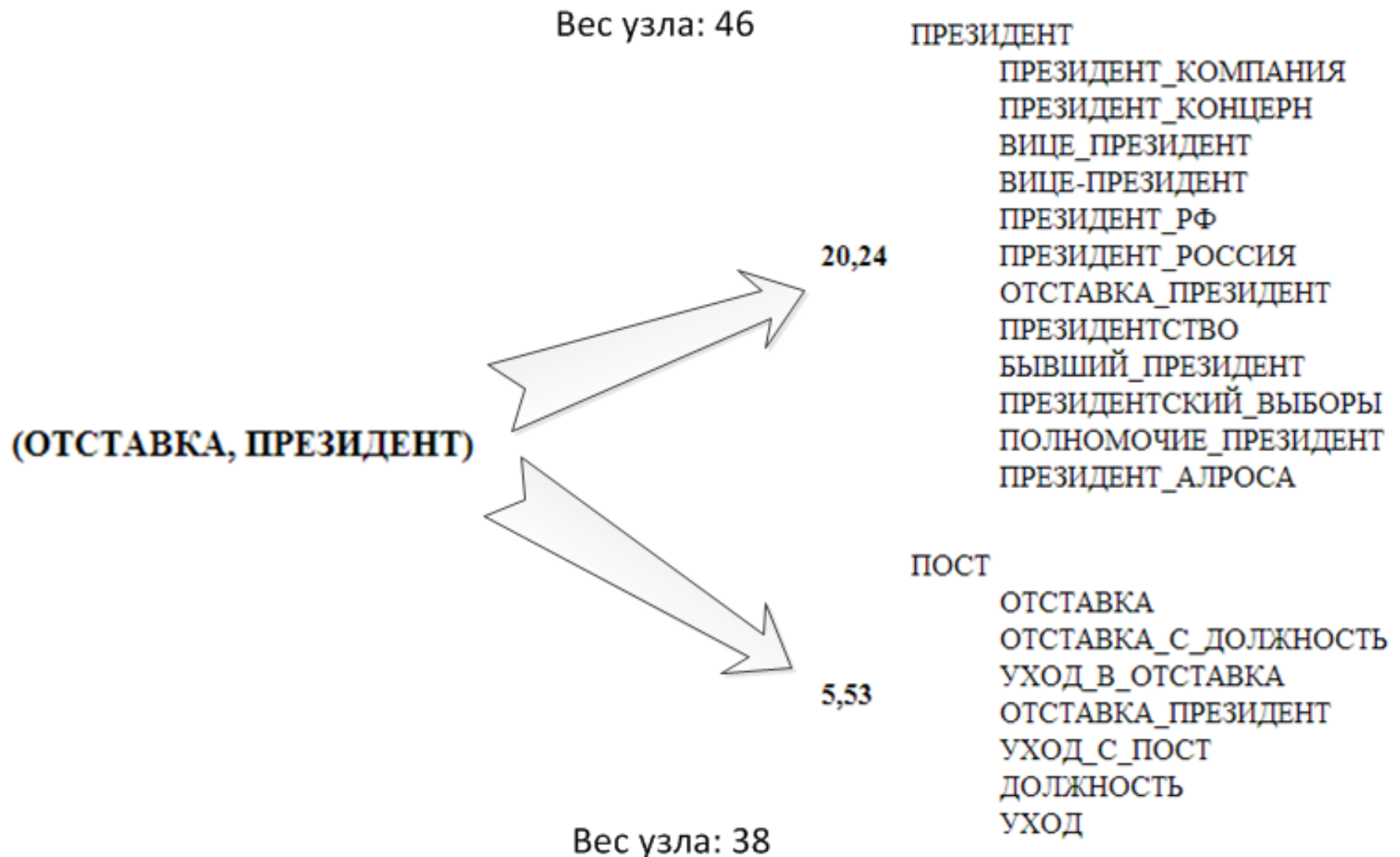
- Вес тематического узла:

$$weight(TU) = \sum_{ЭЛЕМЕНТ_ТУ_i \in TU} freq(ЭЛЕМЕНТ_ТУ_i)$$

- Вес объекта разбивается на веса в соответствующих тематических узлах (пропорционально силе связи с центром узла)

ПРИМЕНЕНИЕ ТЕМАТИЧЕСКОГО ПРЕДСТАВЛЕНИЯ - ПРИМЕР

35



ТЕМАТИЧЕСКОЕ ПРЕДСТАВЛЕНИЕ И АННОТИРОВАНИЕ

36

- Известные метода аннотирования:
 - Maximal Marginal Relevance (MMR)
 - SumBasic
 - Аннотирование на основе тем-го представления РуТез
- Специализированные методы аннотирования
 - Учет основных тематических узлов
 - Учет основных связей тематических узлов
- Учет **IDF**:

$$MWE_IDF = \ln \left(\frac{DOC_COUNT}{\left[\prod_{w_i \in MWE} Freq(w_i) \right] / [DOC_COUNT]^{N-1}} \right)$$

Maximal Marginal Relevance (MMR)

37

- ❑ Известный метод для запрос-ориентированного аннотирования (Carbonell, Goldstein, 1998)
- ❑ Итеративный метод
- ❑ Ранжирование предложений-кандидатов:
 - ✓ Максимизировать сходство с запросом
 - ✓ Минимизировать сходство с уже отобранными в аннотацию предложениями

Пусть: Q – запрос к системе, S – множество предложений кандидатов, s – рассматриваемое предложение кандидат, E – множество выбранных предложений. Тогда:

$$MMR = \arg \max_{s \in S} \left[\lambda \cdot Sim_1(s, Q) - (1 - \lambda) \cdot \max_{s_j \in E} Sim_2(s, s_j) \right]$$

- ❑ **Добавление тематического представления:** вектора кластера и предложений строятся по тематическим узлам

SumBasic

38

- ❑ Метод аннотирования, основанный на частотных характеристиках слов (Nenkova, Vanderwende, 2005)
- ❑ **Основная идея:** наиболее частотные слова исходного кластера с большей вероятностью должны оказаться в аннотации кластера:
$$p(w_i) = \frac{n}{N}$$

n – число вхождений слова
 N – общее число токенов

- ❑ Итеративный метод. На каждой итерации происходит расчет вероятностей слов, отбирается предложение с максимальной средней вероятностью слов:

$$weight(S_j) = \sum_{w_i \in S_j} \frac{p(w_i)}{|\{w_i \mid w_i \in S_j\}|}$$

- ❑ После отбора предложения происходит пересчет вероятностей для слов из отобранного предложения

SumBasic

+ ТЕМАТИЧЕСКОЕ ПРЕДСТАВЛЕНИЕ

39

$$p_{\text{new}}(w_i) = p_{\text{old}}(w_i) * p_{\text{old}}(w_i)$$

- ❑ Переход от вероятностей отдельных слов к вероятностям тематических узлов:

$$p(w_i) = \{p(T_1), p(T_2), \dots p(T_K)\} \quad \mathbf{K} - \text{число ТУ}$$

- ❑ В расчете веса предложения участвует тематический узел с бОльшей вероятностью:

$$\text{weight}(w_i) = \max (\{p(T_1), p(T_2), \dots p(T_K)\})$$

- ❑ Понижение веса для отобранного предложения происходит для тех же узлов, что участвовали при расчете веса предложения

АННОТИРОВАНИЕ НА ОСНОВЕ ТЕМАТИЧЕСКОГО ПРЕДСТАВЛЕНИЯ ПО ТЕЗАУРУСУ РУТЕЗ

40

- «Автоматическое аннотирование новостных кластеров на основе тематического представления» (Лукашевич Н.В., Добров Б.В, 2009)
- Пусть: ТУ ТЕЗ – тематический узел на основе тезауруса РуТез
ИС – именованная сущность
* OLD – упомянутый ТУ_ТЕЗ или ИС
* NEW – новый ТУ_ТЕЗ или ИС
- Тогда на каждой итерации алгоритма выбирается предложение:
 1. $S_i \supset \{ТУ_ТЕЗ_OLD \vee ИС_OLD\}$
 2. $S_i \supset \{ТУ_ТЕЗ_NEW \vee ИС_NEW\}$
 3. $\max(\sum_{ТУ_ТЕЗ \in S_i} weight(ТУ_ТЕЗ) + \sum_{ИС \in S_i} weight(ИС))$

СОБСТВЕННЫЕ АЛГОРИТМЫ АННОТИРОВАНИЯ НА ОСНОВЕ ТЕМАТИЧЕСКОГО ПРЕДСТАВЛЕНИЯ

41

Пусть: ТУ – тематический узел на основе предложенного алгоритма
ТУ REL – пара ТУ, * NEW – новый ТУ или ТУ_REL

Предлагается два алгоритма аннотирования:

I. Отбор по тематическим узлам (*OurSummary_Nodes*)

$$\text{➤ } s_i \Rightarrow \max \left(\begin{array}{c} \text{desc weight}(TY_NEW_j) \\ \sum_{TY_NEW_j \in s_i, i=1..3} \text{weight}(TY_NEW_j) \end{array} \right)$$

II. Отбор по связям ТУ (*OurSummary_Relations*)

➤ $\text{weight}(TY_REL)$ – число вхождений в одни и те же предложения

$$s_i \Rightarrow \max_{TY_REL_NEW_j \in Cluster} \left(\text{weight}(TY_REL_NEW_j) \right)$$

$$\text{➤ } s_i \Rightarrow \max \left(\sum_{TY_REL_NEW_j \in s_i} \text{weight}(TY_REL_NEW_j) \right)$$

АЛГОРИТМ ТЕСТИРОВАНИЯ

42

- ❑ Оценка результата нетривиальна: высокая степень субъективности и низкая согласованность экспертов
- ❑ 11 новостных кластеров, 2-4 ручные аннотации к каждому
- ❑ Оценка содержания аннотаций методом ROUGE
- ❑ Различные алгоритмы аннотирования в различных модификациях:
 - Maximal Marginal Relevance (MMR, 4 модификации)
 - SumBasic (2 модификации)
 - Собственный алгоритм аннотирования (4 модификации)
 - Аннотирование на основе тем-го представления PyTез
- ❑ Модификации алгоритмов как на основе пословной модели, так и на основе тематического представления

ПРИМЕР АННОТАЦИИ

43

Президент "АЛРОСА" Александр Ничипорук уходит со своего поста.

1. Отношения между нынешним и бывшим президентами АЛРОСА (господин Штыров возглавлял компанию до того как стал президентом Якутии) не сложились с момента прихода господина Ничипорука два года назад.

2. Александр Ничипорук покидает пост президента компании АЛРОСА, на которую приходится четверть мировой добычи алмазов.

3. Акционерами "АЛРОСА" являются Росимущество - 37 % акций, Министерство по управлению госимуществом Якутии - 32 %, физические и юридические лица - 23 %.

ПРИМЕР АННОТАЦИИ

Президент "АЛРОСА" Александр Ничипорук уходит со своего поста.

1. Отношения между нынешним и бывшим президентами АЛРОСА (господин Штыров возглавлял компанию до того как стал президентом Якутии) не сложились с момента прихода господина Ничипорука два года назад.

2. Александр Ничипорук покидает пост президента компании АЛРОСА, на которую приходится четверть мировой добычи алмазов.

3. Акционерами "АЛРОСА" являются Росимущество - 37 % акций, Министерство по управлению госимуществом Якутии - 32 %, физические и юридические лица - 23 %.

ПРИМЕР АННОТАЦИИ

45

Президент "АЛРОСА" Александр Ничипорук уходит со своего поста.

1. Отношения между нынешним и бывшим президентами АЛРОСА (господин Штыров возглавлял компанию до того как стал президентом Якутии) не сложились с момента прихода господина Ничипорука два года назад.

2. Александр Ничипорук покидает пост президента компании АЛРОСА, на которую приходится четверть мировой добычи алмазов.

3. Акционерами "АЛРОСА" являются Росимущество - 37 % акций, Министерство по управлению госимуществом Якутии - 32 %, физические и юридические лица - 23 %.

ПРИМЕР АННОТАЦИИ

46

Президент "АЛРОСА" Александр Ничипорук уходит со своего поста.

1. Отношения между нынешним и бывшим президентами АЛРОСА (господин Штыров возглавлял компанию до того как стал президентом Якутии) не сложились с момента прихода господина Ничипорука два года назад.

2. Александр Ничипорук покидает пост президента компании АЛРОСА, на которую приходится четверть мировой добычи алмазов.

3. Акционерами "АЛРОСА" являются Росимущество - 37 % акций, Министерство по управлению госимуществом Якутии - 32 %, физические и юридические лица - 23 %.

ПРИМЕР АННОТАЦИИ

47

Президент "АЛРОСА" Александр Ничипорук уходит со своего поста.

1. Отношения между нынешним и бывшим президентами АЛРОСА (господин Штыров возглавлял компанию до того как стал президентом Якутии) не сложились с момента прихода господина Ничипорука два года назад.

2. Александр Ничипорук покидает пост президента компании АЛРОСА, на которую приходится четверть мировой добычи алмазов.

3. Акционерами "АЛРОСА" являются Росимущество - 37 % акций, Министерство по управлению госимуществом Якутии - 32 %, физические и юридические лица - 23 %.

ПРИМЕР АННОТАЦИИ

48

Президент "АЛРОСА" Александр Ничипорук уходит со своего поста.

1. Отношения между нынешним и бывшим президентами АЛРОСА (господин Штыров возглавлял компанию до того как стал президентом Якутии) не сложились с момента прихода господина Ничипорука два года назад.

2. Александр Ничипорук покидает пост президента компании АЛРОСА, на которую приходится четверть мировой добычи алмазов.

3. Акционерами "АЛРОСА" являются Росимущество - 37 % акций, Министерство по управлению госимуществом Якутии - 32 %, физические и юридические лица - 23 %.

МЕТОД ROUGE

49

- ❑ ROUGE - Recall-Oriented Understudy for Gisting Evaluation (Chin-Yew Lin, 2004)
- ❑ **Основная идея:** сопоставление порожденных аннотаций с ручными (составленными экспертами)
- ❑ Различные метрики сравнения: N-граммы, Longest Common Substring и их модификации

$$ROUGE-N = \frac{\sum_{S \in \{ReferenceSummaries\}} \sum_{gram_n \in S} Count_{match}(gram_n)}{\sum_{S \in \{ReferenceSummaries\}} \sum_{gram_n \in S} Count(gram_n)}$$

РЕЗУЛЬТАТЫ

50

Метод	1	2	L	S	SU	Avg
MMR + Groups	0,58538 (1)	0,37498 (2)	0,55918 (2)	0,31262 (1)	0,32351 (1)	1,4
MMR_WithIDF	0,57996 (3)	0,3833 (1)	0,55891 (3)	0,29818 (2)	0,31003 (2)	2,2
MMR_WithoutIDF	0,58128 (2)	0,35748 (3)	0,56248 (1)	0,27601 (3)	0,28815 (3)	2,4
OurSummary_Nodes	0,54595 (5)	0,32606 (5)	0,52076 (5)	0,27397 (4)	0,285 (4)	4,6
OurSummary_Nodes_WithIDF	0,55526 (4)	0,32374 (6)	0,52925 (4)	0,26181 (7)	0,2738 (7)	5,6
OurSummary_Relations	0,54307 (6)	0,30604 (7)	0,51626 (6)	0,26375 (5)	0,27526 (5)	5,8
ThematicLines	0,52831 (7)	0,3319 (4)	0,50912 (7)	0,26312 (6)	0,27431 (6)	6,0
SumBasic + Groups	0,51884 (9)	0,24591 (10)	0,48974 (9)	0,25057 (8)	0,26142 (8)	8,8
MMR_WithIDF + Groups	0,45376 (10)	0,27439 (8)	0,43607 (10)	0,23926 (9)	0,24899 (9)	9,2
SumBasic	0,52753 (8)	0,23993 (11)	0,49326 (8)	0,23434 (10)	0,24628 (10)	9,4
OurSummary_Relations_WithIDF	0,45172 (11)	0,25309 (9)	0,4333 (11)	0,17712 (11)	0,1882 (11)	10,6

РЕЗУЛЬТАТЫ

51

Метод	1	2	L	S	SU	Avg
MMR + Groups	0,58538 (1)	0,37498 (2)	0,55918 (2)	0,31262 (1)	0,32351 (1)	1,4
MMR_WithIDF	0,57996 (3)	0,3833 (1)	0,55891 (3)	0,29818 (2)	0,31003 (2)	2,2
MMR_WithoutIDF	0,58128 (2)	0,35748 (3)	0,56248 (1)	0,27601 (3)	0,28815 (3)	2,4
OurSummary_Nodes	0,54595 (5)	0,32606 (5)	0,52076 (5)	0,27397 (4)	0,285 (4)	4,6
OurSummary_Nodes_WithIDF	0,55526 (4)	0,32374 (6)	0,52925 (4)	0,26181 (7)	0,2738 (7)	5,6
OurSummary_Relations	0,54307 (6)	0,30604 (7)	0,51626 (6)	0,26375 (5)	0,27526 (5)	5,8
ThematicLines	0,52831 (7)	0,3319 (4)	0,50912 (7)	0,26312 (6)	0,27431 (6)	6,0
SumBasic + Groups	0,51884 (9)	0,24591 (10)	0,48974 (9)	0,25057 (8)	0,26142 (8)	8,8
MMR_WithIDF + Groups	0,45376 (10)	0,27439 (8)	0,43607 (10)	0,23926 (9)	0,24899 (9)	9,2
SumBasic	0,52753 (8)	0,23993 (11)	0,49326 (8)	0,23434 (10)	0,24628 (10)	9,4
OurSummary_Relations WithIDF	0,45172 (11)	0,25309 (9)	0,4333 (11)	0,17712 (11)	0,1882 (11)	10,6

РЕЗУЛЬТАТЫ

52

Метод	1	2	L	S	SU	Avg
MMR + Groups	0,58538 (1)	0,37498 (2)	0,55918 (2)	0,31262 (1)	0,32351 (1)	1,4
MMR_WithIDF	0,57996 (3)	0,3833 (1)	0,55891 (3)	0,29818 (2)	0,31003 (2)	2,2
MMR_WithoutIDF	0,58128 (2)	0,35748 (3)	0,56248 (1)	0,27601 (3)	0,28815 (3)	2,4
OurSummary_Nodes	0,54595 (5)	0,32606 (5)	0,52076 (5)	0,27397 (4)	0,285 (4)	4,6
OurSummary_Nodes_WithIDF	0,55526 (4)	0,32374 (6)	0,52925 (4)	0,26181 (7)	0,2738 (7)	5,6
OurSummary_Relations	0,54307 (6)	0,30604 (7)	0,51626 (6)	0,26375 (5)	0,27526 (5)	5,8
ThematicLines	0,52831 (7)	0,3319 (4)	0,50912 (7)	0,26312 (6)	0,27431 (6)	6,0
SumBasic + Groups	0,51884 (9)	0,24591 (10)	0,48974 (9)	0,25057 (8)	0,26142 (8)	8,8
MMR_WithIDF + Groups	0,45376 (10)	0,27439 (8)	0,43607 (10)	0,23926 (9)	0,24899 (9)	9,2
SumBasic	0,52753 (8)	0,23993 (11)	0,49326 (8)	0,23434 (10)	0,24628 (10)	9,4
OurSummary_Relations_WithIDF	0,45172 (11)	0,25309 (9)	0,4333 (11)	0,17712 (11)	0,1882 (11)	10,6

ЗАКЛЮЧЕНИЕ

53

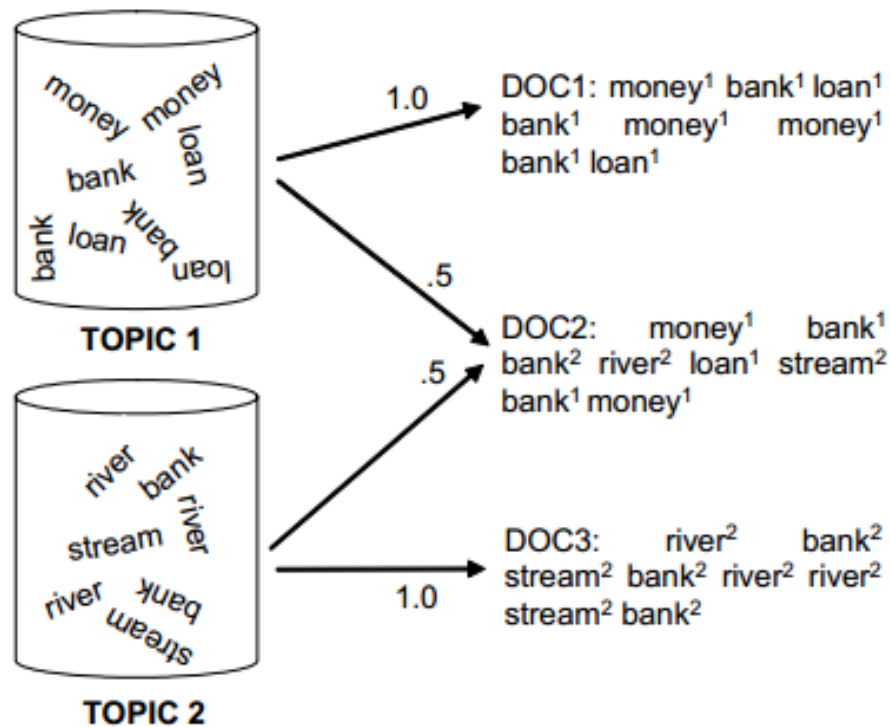
- ❑ Предложена и реализована модель автоматического построения тематического представления новостного кластера, основанная на совокупности нескольких факторов
- ❑ Новым при построении тематических является разделяющий фактор: частая встречаемость выражений в одних и тех же предложениях текста
- ❑ Предложены методы интеграции тематического представления в алгоритмы автоматического аннотирования новостного кластера
- ❑ Показано улучшение качества работы алгоритмов аннотирования на основе тематического представления

ДОПОЛНИТЕЛЬНЫЕ СЛАЙДЫ

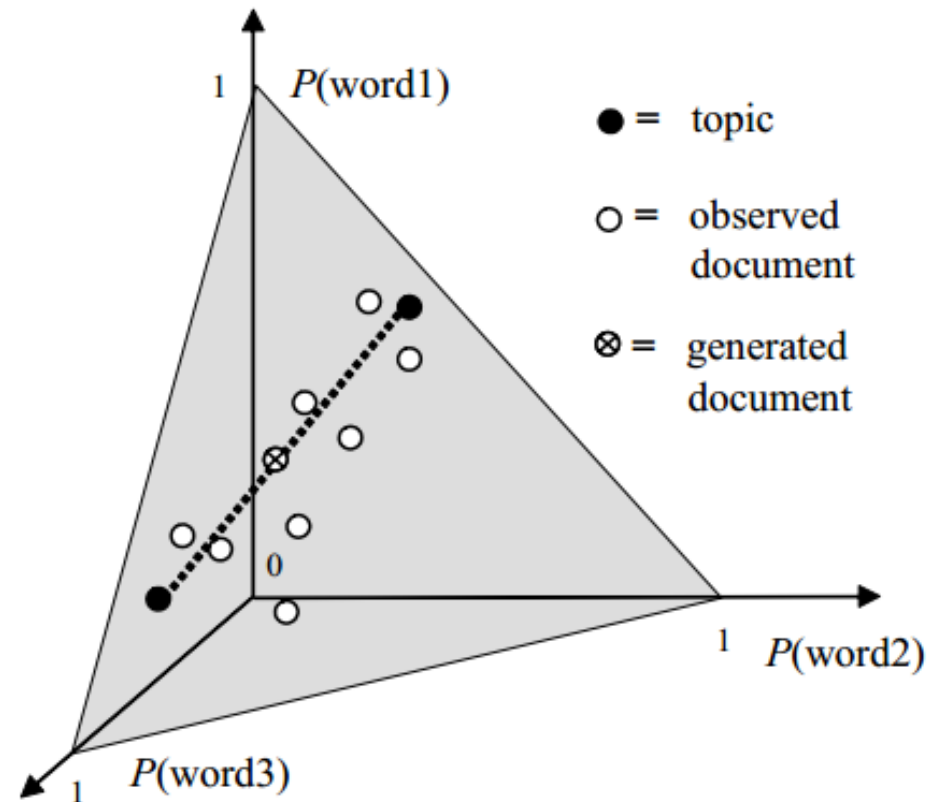
Latent Dirichlet Allocation (LDA)

University of California, USA, 2003
David M. Blei, Andrew Y. Ng, Michael I. Jordan

55



Пример автоматического семплинга



Геометрическая интерпретация LDA

РАЗЛИЧИЯ С МОДЕЛЬЮ LDA

56

LDA	НАША МОДЕЛЬ
1. Совместная встречаемость => описание ситуации, сферы деятельности	1. Выделение различного именования основных участников
2. Результат разбиения – непрерывный (ненулевые вероятности)	2. Дискретное (строгое) разбиение
3. Не использует никаких лингвистических знаний	3. Ориентирована на аккумуляцию всех типов информации
4. Нет поддержки многословных выражений	4. Изначальное построение модели с учетом многословных выражений