

ОСНОВЫ ОБРАБОТКИ ТЕКСТОВ

Лекция №8 Разрешение кореферентности

Лектор: н.с. ИСП РАН Сысоев Андрей Анатольевич

7 ноября 2018

План

- ▶ Кореферентность vs анафоричность.
- ▶ Подходы к разрешению кореферентности:
 - ▶ модель "пара упоминаний" (mention-pair model);
 - ▶ модель "сущность - упоминание" (entity-mention model);
 - ▶ подходы на основе ранжирования (ranking models).
- ▶ Поиск упоминаний в тексте.
- ▶ Оценка качества.

Разрешение кореферентности

Упоминание - последовательность слов текста, обозначающая некоторую сущность. Выражение с выделенным главным словом.

Разрешение кореферентности: найти в тексте все упоминания, относящиеся к одной и той же сущности.

сущность = кореферентная цепочка = кластер упоминаний

Шариков злобно покосился на профессора, а тот отправил ему косой взгляд. Через десять минут Шариков уехал в цирк. Филлип Филлипович остался один в своем кабинете. Он начал шагами мерять комнату.

Разрешение кореферентности

Упоминание - последовательность слов текста, обозначающая некоторую сущность. Выражение с выделенным главным словом.

Разрешение кореферентности: найти в тексте все упоминания, относящиеся к одной и той же сущности.

сущность = кореферентная цепочка = кластер упоминаний

Шариков злобно покосился на профессора, а тот отправил ему косой взгляд. Через десять минут Шариков уехал в цирк. Филлип Филлипович остался один в своем кабинете. Он начал шагами мерять комнату.

Разрешение анафоры

Анафор - упоминание, которое соотносится с предшествующим упоминанием (**антецедентом**), заменяя его или указывая на него. Чаще всего в роли анафора выступает местоимение.

Шариков злобно покосился на профессора, а тот отправил ему косой взгляд. Через десять минут Шариков уехал в цирк. Филлип Филлипович остался один в своем кабинете. Он начал шагами мерять комнату.

Разрешение анафоры

Анафор - упоминание, которое соотносится с предшествующим упоминанием (**антецедентом**), заменяя его или указывая на него. Чаще всего в роли анафора выступает местоимение.

Шариков злобно покосился на **профессора**, а **тот** отправил **ему** косой взгляд. Через десять минут Шариков уехал в цирк. **Филлип Филлипович** остался один в **своем кабинете**. **Он** начал шагами мерять **комнату**.

Корреферентность vs анафоричность

— Кот, — сообразил **Борменталь** и выскочил за ним вслед.

...

— Где он?! **Иван Арнольдович**, успокойте, ради бога, пациентов в приемной!

Борменталь, Иван Арнольдович

Корреферентность vs анафоричность

— Кот, — сообразил **Борменталь** и выскочил за ним вслед.

...

— Где он?! **Иван Арнольдович**, успокойте, ради бога, пациентов в приемной!

Борменталь, Иван Арнольдович

- ▶ корреферентность есть
- ▶ анафоричности нет

Корреферентность vs анафоричность

Я вошел в квартиру: стены белые, потолки высокие.

Кореферентность vs анафоричность

Я вошел в квартиру: стены белые, потолки высокие.

Бриджинг – это вид анафоры, при котором анафорически связанные элементы не кореферентны.

Разрешение кореферентности

Найти в тексте все упоминания, относящиеся к одной и той же сущности.

План

- ▶ Кореферентность vs анафоричность.
- ▶ Подходы к разрешению кореферентности.
- ▶ Поиск упоминаний в тексте.
- ▶ Оценка качества.

Классификация

- ▶ модель "пара упоминаний" (mention-pair model)
- ▶ модель "сущность - упоминание" (entity-mention model)
- ▶ ранжирующие модели (ranking models)

Модель "пара упоминаний"

- ▶ Для каждой пары упоминаний определить, входят ли они в одну кореферентную цепочку.
- ▶ Использовать классифицированные пары упоминаний для построения финальных кореферентных цепочек.

Модель "пара упоминаний"

Пример

Шариков злобно покосился на профессора, а тот отправил ему косой взгляд. Через десять минут Шариков уехал в цирк. Филлип Филлипович остался один в своем кабинете. Он начал шагами мерять комнату.

Модель "пара упоминаний"

Пример

Шариков злобно покосился на профессора, а тот отправил ему косой взгляд. Через десять минут Шариков уехал в цирк. Филлип Филлипович остался один в своем кабинете. Он начал шагами мерять комнату.

Модель "пара упоминаний"

Пример

Шариков злобно покосился на **профессора**, а тот отправил ему косой взгляд. Через десять минут Шариков уехал в цирк. Филлип Филлипович остался один в своем кабинете. Он начал шагами мерять комнату.

Модель "пара упоминаний"

Пример

Шариков злобно покосился на профессора, а тот отправил ему косой взгляд. Через десять минут Шариков уехал в цирк. Филлип Филлипович остался один в своем кабинете. Он начал шагами мерять комнату.

профессора ← тот

Модель "пара упоминаний"

Пример

Шариков злобно покосился на **профессора**, а **тот** отправил **ему** косой взгляд. Через десять минут Шариков уехал в цирк. Филлип Филлипович остался один в своем кабинете. Он начал шагами мерять комнату.

профессора ← тот
Шариков ← ему

Модель "пара упоминаний"

Пример

Шариков злобно покосился на профессора, а тот отправил ему косой взгляд. Через десять минут Шариков уехал в цирк. Филлип Филлипович остался один в своем кабинете. Он начал шагами мерять комнату.

профессора ← тот
Шариков ← ему
...

{Шариков, ему, Шариков}

{професора, тот, Филлип Филлипович, Он}

{своем кабинете, комнату}

Модель "пара упоминаний"

Подзадачи

- ▶ классификатор пар упоминаний
 - ▶ признаки
 - ▶ алгоритм машинного обучения
- ▶ отбор пар упоминаний
- ▶ объединение пар упоминаний в кореферентные цепочки

Классификатор пар упоминаний

Классификатор для каждой пары упоминаний должен выдать:

- ▶ метку: да / нет
- ▶ (опционально) уверенность

Классификатор пар упоминаний

Признаки

- ▶ **базовые** (для главных слов упоминаний и контекста):
 - ▶ токен
 - ▶ лемма
 - ▶ часть речи
 - ▶ граммемы (род, число, одушевленность)
- ▶ **тип именованной сущности**
- ▶ **согласованность**
 - ▶ род
 - ▶ число
 - ▶ одушевленность
 - ▶ оба упоминания - местоимения
 - ▶ тип именованной сущности

Классификатор пар упоминаний

Признаки

- ▶ **позиционные:** расстояние между упоминаниями в
 - ▶ токенах
 - ▶ упоминаниях
 - ▶ предложениях
 - ▶ параграфах
- ▶ **структурные:**
 - ▶ размер в токенах
 - ▶ взаимодействие с другими упоминаниями:
 - ▶ включение
 - ▶ вложенность

его кабинет

Классификатор пар упоминаний

Признаки

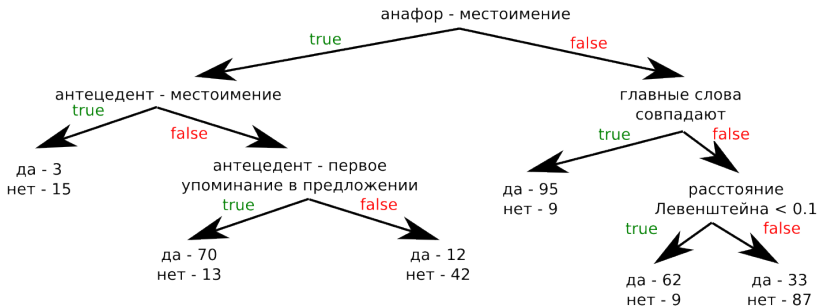
- ▶ **согласованность текстовых представлений:**
 - ▶ одно упоминание является частью другого:
профессор Преображенский - профессор
 - ▶ аббревиатуры:
нэп - новая экономическая политика
домком - домовый комитет
 - ▶ лексикографическая близость: по буквенным n-граммам, расстояние Левенштейна
 - ▶ совпадение главных слов упоминаний
- ▶ **синтаксические:**
 - ▶ роль в предложении
 - ▶ принадлежность одной клаузе (простое предложение в рамках сложного)
 - ▶ наличие общего родителя в синтаксическом дереве
- ▶ **семантические / синонимические:** косинус угла между векторами главных слов упоминаний

Классификатор пар упоминаний

Алгоритм машинного обучения

- ▶ SVM
- ▶ дерево решений
- ▶ случайный лес (random forest)

Дерево решений



Случайный лес

- ▶ строится набор из N деревьев решений
- ▶ при построении дерева / узла дерева выбор происходит не из всех признаков, а лишь из некоторого случайного подмножества

Построение составных признаков

метка	ант_пер	анаф_пер	лекс_близ >0.5	одуш	ед
да	F	F	T	F	F
нет	T	T	F	F	T
да	T	F	F	F	T
нет	F	F	F	T	F
да	T	F	T	F	T
нет	T	T	F	T	F
да	F	T	T	F	F

M. Segond, C. Borgelt

Item Set Mining Based on Cover Similarity

Построение составных признаков

метка	ант_пер	анаф_пер	лекс_близ >0.5	одуш	ед
да	F	F	T	F	F
нет	T	T	F	F	T
да	T	F	F	F	T
нет	F	F	F	T	F
да	T	F	T	F	T
нет	T	T	F	T	F
да	F	T	T	F	F

Построение составных признаков

метка	ант_пер	анаф_пер	лекс_близ > 0.5	одуш	ед
да	F	F	T	F	F
нет	T	T	F	F	T
да	T	F	F	F	T
нет	F	F	F	T	F
да	T	F	T	F	T
нет	T	T	F	T	F
да	F	T	T	F	F

Можно добавить новый признак:

$$\text{ант_пер} \wedge \text{анаф_пер} \wedge \neg(\text{лекс_близ} > 0.5)$$

Модель "пара упоминаний"

Подзадачи

- ▶ классификатор пар упоминаний
 - ▶ признаки
 - ▶ алгоритм машинного обучения
- ▶ отбор пар упоминаний
- ▶ объединение пар упоминаний в кореферентные цепочки

Отбор пар упоминаний

На этапе обучения

Для каждого упоминания:

- ▶ ближайшее к нему слева упоминание из одной с ним кореферентной цепочки - положительный пример
- ▶ упоминания между ними - отрицательные примеры

Отбор пар упоминаний

На этапе обучения

Для каждого упоминания:

- ▶ ближайшее к нему слева упоминание из одной с ним кореферентной цепочки - положительный пример
- ▶ упоминания между ними - отрицательные примеры

Шариков злобно покосился на профессора, а тот отправил ему косой взгляд. Через десять минут Шариков уехал в цирк. Филлип Филлипович остался один в своем кабинете. Он начал шагами мерять комнату.

Отбор пар упоминаний

На этапе обучения

Для каждого упоминания:

- ▶ ближайшее к нему слева упоминание из одной с ним кореферентной цепочки - положительный пример
- ▶ упоминания между ними - отрицательные примеры

Шариков злобно покосился на профессора, а **тот** отправил **ему** косой взгляд. Через десять минут **Шариков** уехал в цирк. **Филлип Филлипович** остался один в своем кабинете. Он начал шагами мерять комнату.

- ▶ **тот** — Филлип Филлипович
- ▶ **ему** — Филлип Филлипович
- ▶ **Шариков** — Филлип Филлипович

Отбор пар упоминаний

На этапе применения

- ▶ перебор всех возможных пар

Отбор пар упоминаний

На этапе применения

- ▶ перебор всех возможных пар - **неэффективно!**
- ▶ для каждого упоминания анализируются **N** упоминаний слева. Выбор **N**:
 - ▶ эвристически
 - ▶ подбор по валидационному набору
 - ▶ по расстояниям в тренировочном наборе

Модель "пара упоминаний"

Подзадачи

- ▶ классификатор пар упоминаний
 - ▶ признаки
 - ▶ алгоритм машинного обучения
- ▶ отбор пар упоминаний
- ▶ объединение пар упоминаний в кореллированные цепочки

Объединение пар упоминаний в коррелентные цепочки

Все коррелентные пары

- ▶ использовать только пары, классифицированные как коррелентные
- ▶ отфильтровать по порогу уверенности
- ▶ преобразовать пары в кластеры по полученным компонентам связности

Объединение пар упоминаний в коррелентные цепочки

Все коррелентные пары

[0.98] Шариков - Шариков

[0.95] профессора - Филлип Филлипович

[0.91] Шариков - ему

[0.85] профессора - Он

[0.80] порог отсечения

[0.78] профессора - тот

Коррелентные цепочки:

{Шариков, Шариков, ему}

{профессора, Филлип Филлипович, Он}

Объединение пар упоминаний в коррелентные цепочки

Ближайшее коррелентное упоминание

Для каждого упоминания в качестве антецедента выбирается ближайшее к нему слева, классифицированное как коррелентное.

Объединение пар упоминаний в корреферентные цепочки

Ближайшее корреферентное упоминание

Для каждого упоминания в качестве antecedента выбирается ближайшее к нему слева, классифицированное как корреферентное.

[Шариков, профессора, тот] — ему

[Шариков, Филлип Филлипович, своем кабинете] — Он

Объединение пар упоминаний в коррелентные цепочки

Ближайшее коррелентное упоминание

Для каждого упоминания в качестве антецедента выбирается ближайшее к нему слева, классифицированное как коррелентное.

[Шариков, профессора, тот] — ему

[Шариков, Филлип Филлипович, своем кабинете] — Он

Объединение пар упоминаний в коррелентные цепочки

Лучшее коррелентное упоминание

Для каждого упоминания в качестве antecedenta выбирается упоминание, классифицированное как коррелентное с большей уверенностью.

Объединение пар упоминаний в коррелентные цепочки

Лучшее коррелентное упоминание

Для каждого упоминания в качестве antecedента выбирается упоминание, классифицированное как коррелентное с большей уверенностью.

[Шариков_[0.8], профессора_[0.7], тот] — ему

[Шариков_[0.6], Филлип Филлипович_[0.9], своем кабинете] — Он

Объединение пар упоминаний в коррелентные цепочки

Лучшее коррелентное упоминание

Для каждого упоминания в качестве antecedента выбирается упоминание, классифицированное как коррелентное с большей уверенностью.

[Шариков_[0.8], профессора_[0.7], тот] — ему

[Шариков_[0.6], Филлип Филлипович_[0.9], своем кабинете] — Он

Объединение пар упоминаний в коррелентные цепочки

Easy-First Mention Pair

- ▶ Каждое упоминание образует свой собственный кластер.
- ▶ Упорядочить все пары упоминаний (положительные и отрицательные) по убыванию уверенности классификатора.
- ▶ Для каждой пары
 - ▶ Если пара отрицательная и не противоречит текущей структуре кластеров, то запомнить ее в качестве ограничения. Иначе - проигнорировать.
 - ▶ Если пара положительная и не противоречит текущим ограничениям, то объединить соответствующие кластера в один. Иначе - проигнорировать.

Объединение пар упоминаний в коррелентные цепочки

Easy-First Mention Pair

Классифицированные пары

Ограничения

профессора - Филлип Филиппович

тот - ему

ему - Шариков

профессора - тот

тот - Филлип Филиппович

ему - Филлип Филиппович

Кластеры

{профессора}, {Филлип Филиппович}, {тот}, {ему}, {Шариков}

Объединение пар упоминаний в коррелирующие цепочки

Easy-First Mention Pair

Классифицированные пары

Ограничения

тот - ему

ему - Шариков

профессора - тот

тот - Филлип Филиппович

ему - Филлип Филиппович

Кластеры

{профессора, Филлип Филиппович}, {тот}, {ему}, {Шариков}

Объединение пар упоминаний в коррелирующие цепочки

Easy-First Mention Pair

Классифицированные пары

Ограничения

ему - Шариков
профессора - тот
тот - Филлип Филиппович
ему - Филлип Филиппович

тот - ему

Кластеры

{профессора, Филлип Филиппович}, {тот}, {ему}, {Шариков}

Объединение пар упоминаний в коррелирующие цепочки

Easy-First Mention Pair

Классифицированные пары

Ограничения

профессора - тот
тот - Филлип Филиппович
ему - Филлип Филиппович

тот - ему

Кластеры

{профессора, Филлип Филиппович}, {тот}, {ему, Шариков}

Объединение пар упоминаний в коррелирующие цепочки

Easy-First Mention Pair

Классифицированные пары

Ограничения

тот - Филлип Филиппович

ему - Филлип Филиппович

тот - ему

Кластеры

{профессора, Филлип Филиппович, тот}, {ему, Шариков}

Объединение пар упоминаний в коррелентные цепочки

Easy-First Mention Pair

Классифицированные пары

Ограничения

тот - ему

ему - Филлип Филиппович

Кластеры

{профессора, Филлип Филиппович, тот}, {ему, Шариков}

Объединение пар упоминаний в коррелентные цепочки

Easy-First Mention Pair

Классифицированные пары

Ограничения

тот - ему

Кластеры

{профессора, Филлип Филиппович, тот}, {ему, Шариков}

Объединение пар упоминаний в коррелирующие цепочки

Общие алгоритмы кластеризации

- ▶ Преобразовать список пар в граф (возможно, взвешенный)
- ▶ Применить алгоритм кластеризации графа:
 - ▶ Chinese whispers
 - ▶ Markov clustering
 - ▶ кластеризация на основе алгебраической связности

C. Biemann

Chinese Whispers: An Efficient Graph Clustering Algorithm and Its Application to Natural Language Processing Problems

S. van Dongen

Graph Clustering by Flow Simulation

Классификация

- ▶ модель "пара упоминаний" (mention-pair model)
- ▶ модель "сущность - упоминание" (entity-mention model)
- ▶ ранжирующие модели (ranking models)

Модель "сущность - упоминание"

Модель "сущность - упоминание"

- ▶ Текст анализируется слева направо.
- ▶ Для каждого упоминания нужно определить:
 - ▶ принадлежит ли оно одной из уже построенных кореферентных цепочек
 - ▶ или начинает новую цепочку

Модель "сущность - упоминание"

Пример

Шариков злобно покосился на профессора, а тот отправил ему косой взгляд. Через десять минут Шариков уехал в цирк. Филлип Филлипович остался один в своем кабинете. Он начал шагами мерять комнату.

{Шариков, ему, Шариков} - Филлип Филлипович

{профессора, тот} - Филлип Филлипович

Модель "сущность - упоминание"

Пример

Шариков злобно покосился на профессора, а тот отправил ему косой взгляд. Через десять минут Шариков уехал в цирк. Филлип Филлипович остался один в своем кабинете. Он начал шагами мерять комнату.

{Шариков, ему, Шариков} - своем кабинете

{профессора, тот, Филлип Филлипович} - своем кабинете

Модель "сущность - упоминание"

Пример

Шариков злобно покосился на профессора, а тот отправил ему косой взгляд. Через десять минут Шариков уехал в цирк. Филлип Филлипович остался один в своем кабинете. Он начал шагами мерять комнату.

Модель "сущность - упоминание"

Признаки

Признаки для пары <корреляционная цепочка, упоминание> получаются из признаков для пар упоминаний с помощью кванторов и агрегирующих функций:

- ▶ ALL
- ▶ ANY
- ▶ MOST
- ▶ MIN
- ▶ MAX
- ▶ AVG
- ▶ ...

Ранжирующие модели

Для каждого упоминания в явном виде происходит ранжирование возможных antecedентов / кореферентных цепочек.

Примеры моделей ранжирования:

- ▶ модель турнира
- ▶ RankSVM

Модель турнира

Бинарный классификатор для тройки:

- ▶ анафор
- ▶ антецедент-кандидат₁
- ▶ антецедент-кандидат₂

Метка указывает на то, какой антецедент правильный.

Для анафора перебираем все возможные пары антецедентов. В качестве правильного выбираем тот антецедент, который победил большее количество раз.

RankSVM

x - анафор

y^+ - правильный antecedent

$y_1^-, y_2^-, \dots, y_n^-$ - неправильные antecedенты

$\varphi(x, y)$ - вектор признаков для пары анафор-antecedent $\langle x, y \rangle$

Нужно подобрать вектор w :

$$(w, \varphi(x, y^+)) > (w, \varphi(x, y_i^-)), \text{ для } i = 1..n$$

RankSVM

x - анафор

y^+ - правильный antecedent

$y_1^-, y_2^-, \dots, y_n^-$ - неправильные antecedенты

$\varphi(x, y)$ - вектор признаков для пары анафор-антецедент $\langle x, y \rangle$

Нужно подобрать вектор w :

$$(w, \varphi(x, y^+)) > (w, \varphi(x, y_i^-)), \text{ для } i = 1..n$$

То есть

$$(w, \varphi(x, y^+) - \varphi(x, y_i^-)) > 0, \text{ для } i = 1..n$$

RankSVM

x - анафор

y^+ - правильный antecedent

$y_1^-, y_2^-, \dots, y_n^-$ - неправильные antecedенты

$\varphi(x, y)$ - вектор признаков для пары анафор-антецедент $\langle x, y \rangle$

Нужно подобрать вектор w :

$$(w, \varphi(x, y^+)) > (w, \varphi(x, y_i^-)), \text{ для } i = 1..n$$

То есть

$$(w, \varphi(x, y^+) - \varphi(x, y_i^-)) > 0, \text{ для } i = 1..n$$

Обучаем SVM на примерах:

$$\varphi(x, y^+) - \varphi(x, y_i^-) \mapsto 1, \text{ для } i = 1..n$$

$$\varphi(x, y_i^-) - \varphi(x, y^+) \mapsto -1, \text{ для } i = 1..n$$

RankSVM

x - анафор

y^+ - правильный antecedent

$y_1^-, y_2^-, \dots, y_n^-$ - неправильные antecedенты

$\varphi(x, y)$ - вектор признаков для пары анафор-антецедент $\langle x, y \rangle$

Нужно подобрать вектор w :

$$(w, \varphi(x, y^+)) > (w, \varphi(x, y_i^-)), \text{ для } i = 1..n$$

То есть

$$(w, \varphi(x, y^+) - \varphi(x, y_i^-)) > 0, \text{ для } i = 1..n$$

Обучаем SVM на примерах:

$$\varphi(x, y^+) - \varphi(x, y_i^-) \mapsto 1, \text{ для } i = 1..n$$

$$\varphi(x, y_i^-) - \varphi(x, y^+) \mapsto -1, \text{ для } i = 1..n$$

Применение:

$$y^* = \underset{y}{\operatorname{argmax}}(w, \varphi(x, y))$$

План

- ▶ Кореферентность vs анафоричность.
- ▶ Подходы к разрешению кореферентности.
- ▶ Поиск упоминаний в тексте.
- ▶ Оценка качества.

Поиск упоминаний в тексте

Где брать упоминания (для разрешения кореферентности)?

- ▶ из существующей разметки (**в реальной жизни их не будет!**)
- ▶ получать с помощью синтаксического парсера (**много лишних!**)
- ▶ обучить алгоритм определения упоминаний

Поиск упоминаний в тексте

Алгоритм определения упоминаний:

- ▶ поиск главных слов
- ▶ "расширение" главных слов до упоминаний

Поиск главных слов

Задача бинарной классификации токенов.

Алгоритм машинного обучения - например, SVM.

Признаки:

- ▶ базовые:
 - ▶ лемма
 - ▶ часть речи
 - ▶ граммемы (род, число, одушевленность)
- ▶ позиция в предложении
- ▶ синтаксические связи
- ▶ *базовые* для контекста и родителя в синтаксическом дереве
- ▶ количество детей в синтаксическом дереве
- ▶ частота токена в документе

"Расширение" главных слов до упоминаний

Итерационное "наращивание" главного слова влево / вправо до полного упоминания.

Задача бинарной классификации пары токенов: добавлять ли текущий токен к упоминанию с соответствующим главным словом.

Алгоритм машинного обучения - например, SVM.

Признаки:

- ▶ базовые для анализируемого токена и главного слова:
 - ▶ лемма
 - ▶ часть речи
- ▶ направление
- ▶ расстояние
- ▶ позиция в предложении
- ▶ анализируемый токен и главное слово - части именованной сущности
- ▶ наличие связей между анализируемым токеном и главным словом в синтаксическом дереве

"Расширение" главных слов до упоминаний

Пример

Негодяй в грязном колпаке — повар **столовой** Нормального питания служащих — плеснул кипятком и обварил мне левый бок.

"Расширение" главных слов до упоминаний

Пример

Негодяй в грязном колпаке — повар столовой Нормального питания служащих — плеснул кипятком и обварил мне левый бок.

"Расширение" главных слов до упоминаний

Пример

Негодяй в грязном колпаке — повар столовой Нормального питания
служащих — плеснул кипятком и обварил мне левый бок.

"Расширение" главных слов до упоминаний

Пример

Негодяй в грязном колпаке — повар столовой Нормального питания
служащих — плеснул кипятком и обварил мне левый бок.

"Расширение" главных слов до упоминаний

Пример

Негодяй в грязном колпаке — повар столовой Нормального питания служащих — плеснул кипятком и обварил мне левый бок.

"Расширение" главных слов до упоминаний

Пример

Негодяй в грязном колпаке — повар столовой Нормального питания
служащих — плеснул кипятком и обварил мне левый бок.

План

- ▶ Кореферентность vs анафоричность.
- ▶ Подходы к разрешению кореферентности.
- ▶ Поиск упоминаний в тексте.
- ▶ Оценка качества.

Оценка качества

Поиск главных слов и упоминаний

$$\text{точность} = \frac{\text{количество правильно определенных}}{\text{всего найденных системой}}$$

$$\text{полнота} = \frac{\text{количество правильно определенных}}{\text{всего правильных}}$$

$$F1 = \frac{2 \times \text{точность} \times \text{полнота}}{\text{точность} + \text{полнота}}$$

Оценка качества

Разрешение кореферентности

Точность, полнота, F1 по версиям:

- ▶ MUC (Message Understanding Conference)
- ▶ B-cubed
- ▶ CEAF (Constrained Entity-Alignment F-Measure)

$$F1 = \frac{2 \times \text{точность} \times \text{полнота}}{\text{точность} + \text{полнота}}$$

S - правильная кореферентная цепочка

$R = R_1, R_2, \dots, R_m$ - кореферентные цепочки, вычисленные системой

$P_R(S)$ - разбиение цепочки S в соответствии с ответами системы

R_1, R_2, \dots, R_m , пересекающимися с S

Каждое упоминание, не входящее в ответы системы, образует кластер из одного элемента.

MUC

S - правильная кореферентная цепочка

$R = R_1, R_2, \dots, R_m$ - кореферентные цепочки, вычисленные системой

$P_R(S)$ - разбиение цепочки S в соответствии с ответами системы

R_1, R_2, \dots, R_m , пересекающимися с S

Каждое упоминание, не входящее в ответы системы, образует кластер из одного элемента.

Если $S = \{A, B, C, D\}$, $R = \{\{A, B\}\}$, то $P_R(S) = \{\{A, B\}, \{C\}, \{D\}\}$

S - правильная кореферентная цепочка

$R = R_1, R_2, \dots, R_m$ - кореферентные цепочки, вычисленные системой

$P_R(S)$ - разбиение цепочки S в соответствии с ответами системы

R_1, R_2, \dots, R_m , пересекающимися с S

Каждое упоминание, не входящее в ответы системы, образует кластер из одного элемента.

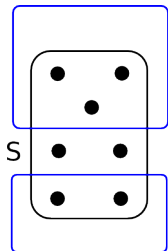
$c(S) = |S| - 1$ - минимальное количество связей, необходимых для объединения элементов из S в кореферентную цепочку

$m_R(S) = |P_R(S)| - 1$ - количество связей, которые надо добавить в $P_R(S)$ для получения S

$$\text{полнота}_S = \frac{c(S) - m_R(S)}{c(S)} = \frac{|S| - |P_R(S)|}{|S| - 1}$$

MUC

Пример



$$\text{полнота}_S = \frac{7 - 4}{7 - 1} = \frac{1}{2}$$

$$\text{полнота} = \frac{\sum_S (|S| - |P_R(S)|)}{\sum_S (|S| - 1)}$$

$$\text{полнота} = \frac{\sum_S (|S| - |P_R(S)|)}{\sum_S (|S| - 1)}$$

Для вычисления точности используются те же формулы, но правильные ответы и ответы системы "меняются местами".

B-cubed

M - множество всех упоминаний m

$S(m)$ - правильная цепочка, содержащая упоминание m

$R(m)$ - цепочка, выданная системой, содержащая упоминание m

Для каждого упоминания m :

$$\text{точность}_m = \frac{|R(m) \cap S(m)|}{|R(m)|}$$

$$\text{полнота}_m = \frac{|R(m) \cap S(m)|}{|S(m)|}$$

B-cubed

M - множество всех упоминаний m

$S(m)$ - правильная цепочка, содержащая упоминание m

$R(m)$ - цепочка, выданная системой, содержащая упоминание m

Для каждого упоминания m :

$$\text{точность}_m = \frac{|R(m) \cap S(m)|}{|R(m)|}$$

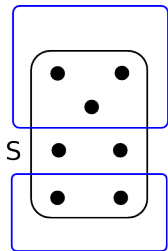
$$\text{полнота}_m = \frac{|R(m) \cap S(m)|}{|S(m)|}$$

$$\text{точность} = \frac{1}{|M|} \sum_{m \in M} \text{точность}_m$$

$$\text{полнота} = \frac{1}{|M|} \sum_{m \in M} \text{полнота}_m$$

B-cubed

Пример



$$\text{полнота} = \frac{1}{7} \left(\frac{3}{7} + \frac{3}{7} + \frac{3}{7} + \frac{1}{7} + \frac{1}{7} + \frac{2}{7} + \frac{2}{7} \right) = \frac{1}{7} \times \frac{15}{7} = \frac{15}{49}$$

$$\text{точность} = \frac{1}{7} \left(\frac{3}{3} + \frac{3}{3} + \frac{3}{3} + \frac{1}{1} + \frac{1}{1} + \frac{2}{2} + \frac{2}{2} \right) = \frac{1}{7} \times 7 = 1.0$$

CEAF

$\mathcal{S} = \{S_1, S_2, \dots, S_n\}$ - правильные кореферентные цепочки

$\mathcal{R} = \{R_1, R_2, \dots, R_m\}$ - кореферентные цепочки, вычисленные системой

$\varphi(S, R) \geq 0$ - оценка близости множеств S и R

$\varphi(S, R) = 0$, если $S = \emptyset$ или $R = \emptyset$. Далее можно считать, что $n = m$ - меньшее множество дополняем с помощью \emptyset до большего.

$g : \mathcal{S} \mapsto \mathcal{R}$ - взаимнооднозначное соответствие между \mathcal{S} и \mathcal{R}

$\Phi(g) = \sum_{i=1}^n \varphi(S_i, g(S_i))$ - суммарная близость

$g^* = \operatorname{argmax}_g \Phi(g)$

$$\text{точность}_{CEAF} = \frac{\Phi(g^*)}{\sum_{i=1}^n \varphi(R_i, R_i)}$$

$$\text{полнота}_{CEAF} = \frac{\Phi(g^*)}{\sum_{i=1}^n \varphi(S_i, S_i)}$$

X. Luo

On Coreference Resolution Performance Metrics

CEAF

CEAF_{упоминание} (CEAF_m):

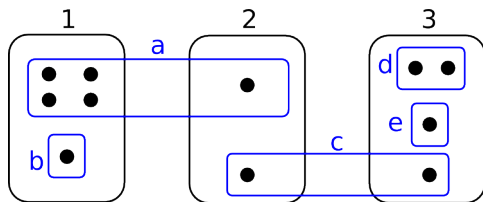
$$\varphi_m(S, R) = |S \cap R|$$

CEAF_{сущность} (CEAF_e):

$$\varphi_e(S, R) = \frac{2|S \cap R|}{|S| + |R|}$$

CEAF

Пример



	1	2	3	\emptyset	\emptyset
a	4	1	0	0	0
b	1	0	0	0	0
c	0	1	1	0	0
d	0	0	2	0	0
e	0	0	1	0	0

$$\text{точность}_{CEAF_m} = \text{полнота}_{CEAF_m} = \frac{7}{11}$$

Следующая лекция

Другие задачи обработки текстов

лектор: Майоров Владимир Дмитриевич