Лабораторная работа №3

Кластеризация текстов по тематикам

Цель работы: научиться реализовывать на выбранном языке программирования алгоритмы кластеризации текстов по тематикам на основе латентно-семантического анализа.

Краткие теоретические сведения

Латентно-семантический анализ отображает документы и отдельные слова в так называемое «семантическое пространство», в котором и производятся все дальнейшие сравнения. При этом делаются следующие предположения:

- 1) Документы это просто набор слов. Порядок слов в документах игнорируется. Важно только то, сколько раз то или иное слово встречается в документе.
- 2) Семантическое значение документа определяется набором слов, которые как правило идут вместе. Например, в биржевых сводках, часто встречаются слова: «фонд», «акция», «доллар»
- 3) Каждое слово имеет единственное значение. Это, безусловно, сильное упрощение, но именно оно делает проблему разрешимой.

Пример

Для примера возьмем несколько заголовков из различных новостей. На первом шаге из этих заголовков исключаем, так называемые, стопсимволы. Это слова которые встречаются в каждом тексте и не несут в себе смысловой нагрузки, это, прежде всего, все союзы, частицы, предлоги и множество других слов. Полный список использованных стопсимволов можно посмотреть в Приложении 1.

Далее необходимо выполнить операцию стемминга. Она не является обязательной, если набор текстов достаточно большой, или если тексты на английском языке, в силу того, что количество вариаций той или иной словоформы в английском языке существенно меньше чем в русском. В нашем же случае, пропускать этот шаг не стоит т.к. это приведет к существенной деградации результатов.

Дальше исключаем слова встречающиеся в единственном экземпляре. Это тоже необязательный шаг, он не влияет на конечный результат, но сильно упрощает математические вычисления. В итоге у нас остались, так называемые, индексируемые слова, они выделены жирным шрифтом:

- 1. Британская полиция знает о местонахождении основателя WikiLeaks
- 2. В суде США начинается процесс против россиянина, рассылавшего спам
- 3. Церемонию вручения Нобелевской премии мира бойкотируют 19 стран
- 4. В Великобритании арестован основатель сайта Wikileaks Джулиан Ассандж
- 5. Украина игнорирует церемонию вручения Нобелевской премии
- 6. Шведский **су**д отказался рассматривать апелляцию **основате**ля **Wikileaks**
- 7. НАТО и **США** разработали планы обороны **стран** Балтии **прот**ив России
- 8. Полиция Великобритании нашла основателя WikiLeaks, но, не арестовала
- 9.В Стокгольме и Осло сегодня состоится вручение Нобелевских премий

Латентно семантический анализ

На первом шаге требуется составить частотную матрицу индексируемых слов. В этой матрице строки соответствуют индексированным словам, а столбцы — документам. В каждой ячейке матрицы указано какое количество раз слово встречается в соответствующем документе (рис. 1).

	m1	T2	m2	Εл	ms.	m۶	7	Ιrο	mα
	1 1	14	13	14	17	10	1/	10	13
wikileaks	1	0	0	1	0	1	0	1	0
арестова	0	0	0	1	0	0	0	1	0
великобритан	0	0	0	1	0	0	0	1	0
вручен	0	0	1	0	1	0	0	0	1
нобелевск	0	0	1	0	1	0	0	0	1
основател	1	0	0	1	0	1	0	1	0
полиц	1	0	0	0	0	0	0	1	0
прем	0	0	1	0	1	0	0	0	1
прот	0	1	0	0	0	0	1	0	0
стран	0	0	1	0	0	0	1	0	0
суд	0	1	0	0	0	1	0	0	0
сша	0	1	0	0	0	0	1	0	0
церемон	0	0	1	0	1	0	0	0	0

Рисунок 1. Частотная матрица индексируемых слов

Следующим шагом необходимо провести сингулярное разложение полученной матрицы (рис. 2). Т.е. исходную матрицу М мы представляем в виде:

M = U*W*Vt

где U и Vt — ортогональные матрицы, а W — диагональная матрица. Причем диагональные элементы матрицы W упорядочены в порядке убывания. Диагональные элементы матрицы W называются сингулярными числами.

wikileaks	0.57	-0.01	0.01	-0.2	0.13	0.16	-0.16	-0.25	-0.64																			
арестова	0.34	-0	0.07	0.41	-0.42	-0.02	0.1	0.17	0.01	2	410	0	lo	h	lo.	0	lo.	To.	Τ	1	T2	T3	T4	T5	T6	T7	T8	T9
великобритан	0.34	-0	0.07	0.41	-0.42	-0.02	0.1	0.17	-0.01	3.4	3.		0	0	0	0	0	-10	0.	43	0.05	0.01	0.54	0	0.37	0.01	0.63	0
вручен	Û	0.52	0.07	-0.06	-0.08	-0.15	-0.17	0.02	-0.07	ŀ	0	2.27	20	0	0	0	0	-	-0		0.02	0.65	-0.01	0.59	-0	0.09	-0.01	0.47
нобелевск	0	0.52	0.07	-0.06	-0.08	-0.15	-0.17	0.02	0.32		0	0.6	-	90	0	0	0	0	0.	03	-0.7	-0.04	0.06	0.1	-0.16	-0.67	0.09	0.09
основател	0.57	-0.01	0.01	-0.2	0.13	0.16	-0.16	-0.25	0.64]. E	0	0	0.9		90	0	0	١,	-0	.22	-0.24	0.15	0.28	-0.11	-0.68	0.44	0.33	-0.1
полиц	0.31	-0	0.05	0.07	0.57	-0.6	0.29	0.37	-0	1	0	0	0	0.1	0.98	0	0	-	0.	69	-0.32	0.22	-0.45	-0.12	-0.03	0.27	-0.02	-0.1
прем	0	0.52	0.07	-0.06	-0.08	-0.15	-0.17	0.02	-0.25		0	6	0	-		0.7	10	-	-0	.27	-0.34	0.44	0.29	-0.13	0.45	0.12	-0.31	-0.4
прот	0.02	0.03	-0.6	0.13	-0.05	-0.22	0	-0.25	0	16	0	6	0	0	0	0.7	0.4	20	-0	.03	0.3	0.14	-0.17	0.44	-0.15	-0.3	0.24	-0.7
стран	0.01	0.22	-0.3	0.39	0.41	0.56	-0.22	0.4	-0	1 0	0	n	0	0	0	0	0.4	20	-0	.3	0.12	0.4	-0.39	-0.53	0.12	-0.23	0.46	0.13
суд	0.12	0.01	-0.3	3-0.62	-0.3	0.12	0.21	0.55	-0	10	Į0	M	M	V	ĮV.	V	Į.	Į.	0.	35	0.35	0.35	0.35	-0.35	-0.35	-0.35	-0.35	0
сша	0.02	0.03	-0.6	0.13	-0.05	-0.22	0	-0.25	0																			
перемон	0	0.38	0.03	0.02	0.08	0.31	0.82	-0.29	0	1																		

Рисунок 2. Сингулярное разложение частотной матрицы индексируемых слов

Согласно простым правилам произведения матриц, видно, что столбцы и строки соответствующие меньшим сингулярным значениям дают наименьший вклад в итоговое произведение. Например, мы можем отбросить последние столбцы матрицы U и последние строки матрицы V^{\star} t, оставив только первые 2. Важно, что при этом гарантируется, оптимальность полученного произведения. Разложение такого вида называют двумерным сингулярным разложением (рис. 3):

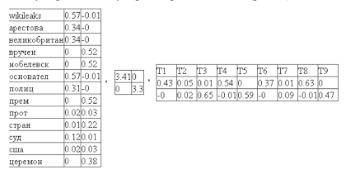


Рисунок 3. Двумерное сингулярное разложение

На практике, конечно, количество групп будет намного больше, пространство будет не двумерным а многомерным, но сама идея остается той же. Можно определять местоположения слов и статей в нашем пространстве и использовать эту информацию для, например, определения тематики статьи.

В полученном факторном пространстве документы и термины концентрируются областями, имеющими общий семантический и латентный смысл.

Применительно к кластеризации, получаемые области и есть кластеры. С помощью математических преобразований можно определить центры кластеров.

Задание к лабораторной работе

Написать программу на выбранном языке программирования, реализующую описанный выше алгоритм для кластеризации заголовков по темам. Количество заголовков задать не более 10 (оптимально от 7 до 10), количество кластеров не более 3. Программа должна запрашивать заголовки (они могут храниться в файлах). После получения двумерного сингулярного разложения осуществить прямую кластеризацию самостоятельно выбранным и изученным методом (иерархические алгоритмы, поиск k-средних и т.д.). Результатом работы программы должен быть файл, содержащий заголовки, разбитые по кластерам.

Контрольные вопросы

Стоп-символы русского языка

Стоп символы, они же стоп-слова, это слова, встречающиеся практически во всех текстах и не несущие специальной смысловой нагрузки. В русском языке, к стоп символам относятся предлоги, суффиксы, причастия, междометия, частицы и т.п. Неполный список стопслов представлен ниже:

_	еще	него	сказать
a	ж	нее	CO
без	же	ней	совсем
более	жизнь	нельзя	так
больше	38	нет	такой
будет	зачем	ни	там
будто	здесь	нибудь	тебя
бы	И	никогда	тем
был	ИЗ	ним	теперь
была	из-за	них	то
были	или	ничего	тогда
было	ИМ	но	тогда
быть	иногда	ну	тоже
В	их	0	только
вам	К	об	TOM
вас	кажется	один	TOT
вдруг	как	ОН	три
ведь	какая	она	TYT
во	какой	ОНИ	ты
ВОТ	когда	ОПЯТЬ	y
впрочем	конечно	ОТ	уж
все	которого	перед	уже
всегда	которые	по	хорошо
всего	кто	под	хоть
всех	куда	после	чего
всю	ли	ПОТОМ	человек
ВЫ	лучше	потому	чем
Γ	между	почти	через
где	меня	при	что
говорил	мне	про	чтоб
да	МНОГО	раз	чтобы
даже	может	разве	чуть
два	МОЖНО	c	ЭТИ
для	мой	сам	этого
до	МОЯ	свое	этой

МЫ	СВОЮ	ЭТОМ
на	себе	ЭТОТ
над	себя	эту
надо	сегодня	Я
наконец	сейчас	
нас	сказал	
не	сказала	
	на над надо наконец нас	на себе над себя надо сегодня наконец сейчас нас сказал

Помимо указанных слов имеет смысл еще фильтровать цифры, отдельные буквы и знаки препинания.