

# Иерархические вероятностные тематические модели

Надежда Чиркова <sup>1</sup>

Научный руководитель: Воронцов К. В. <sup>2</sup>

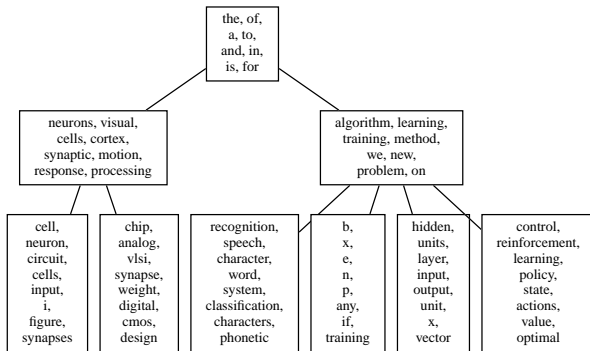
<sup>1</sup>ВМК МГУ <sup>2</sup>ВЦ РАН

Конференция «Ломоносов 2015», 15.04.2015

# План выступления

- 1 Постановка задачи
- 2 Обучение модели
- 3 Виды регуляризаторов
- 4 Результаты

Иерархическая тематическая модель — это многодольный граф тем.  
Каждая тема — это набор **слов**, авторов, **документов**, ... + множество **подтем**



- ✓ Тематические иерархии легко интерпретируемы и понятны человеку
- ✓ Иерархическая структура упрощает навигацию по коллекции документов

*David M Blei, Hierarchical Topic Models and the Nested Chinese Restaurant Process, 2010*

## Постановка задачи

$D$  - коллекция текстовых документов

$W$  - множество терминов

**Дана** коллекция текстовых документов:

$n_{dw}$  - матрица частот слов в документах (мешок слов):

$$F_{dw} = p(w|d) = \frac{n_{dw}}{n_d}$$

**Построить** модель:

построить  $L$  уровней иерархии

## Постановка задачи: построение первого уровня

**Дано:**  $F_{dw} = p(w|d) = \frac{n_{dw}}{n_d}$

**Построить** модель:

$$p(w|d) = \sum_{t \in T} p(w|t)p(t|d) = \sum_{t \in T} \varphi_{wt}\theta_{td} \Leftrightarrow F = \Phi\Theta$$

с параметрами  $\Phi = \{\varphi_{wt}\}_{W \times T}$  и  $\Theta = \{\theta_{td}\}_{T \times D}$ :

$\varphi_{wt} = p(w|t)$  — распределение слов в теме  $t$ ;

$\theta_{td} = p(t|d)$  — распределение тем в документе  $d$ .

**Критерий:** максимизация логарифма правдоподобия:

$$\mathcal{L}(\Phi, \Theta) = \sum_{d,w} n_{dw} \ln \sum_t \varphi_{wt}\theta_{td} + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

## Постановка задачи: построение нижних уровней иерархии

Дано:  $F_{dw} = p(w|d) = \frac{n_{dw}}{n_d}$

Построить модель:

$$p(w|d) = \sum_{t \in T} \sum_{s \in S} p(w|s) p(s|t) p(t|d) = \sum_{t \in T} \sum_{s \in S} \varphi_{ws} \psi_{st} \pi_{td} \Leftrightarrow F = \Phi \Psi \Pi$$

с параметрами  $\Phi = \{\varphi_{ws}\}_{W \times T}$ ,  $\Psi = \{\psi_{st}\}_{S \times T}$ :

$\varphi_{ws} = p(w|s)$  — распределение слов в дочерней теме  $s$ ;

$\psi_{st} = p(s|t)$  — распределение дочерних тем в родительской теме  $t$ .

Распределение родительских тем в документе  $d$   $\pi_{td} = p(t|d)$  фиксировано (вычисляется на предыдущем уровне)

**Критерий:** максимизация логарифма правдоподобия:

$$\mathcal{L}(\Phi, \Theta) = \sum_{d,w} n_{dw} \ln \sum_{t \in T} \sum_{s \in S} \varphi_{ws} \psi_{st} \pi_{td} + R(\Phi, \Psi, \Pi) \rightarrow \max_{\Phi, \Psi}$$

# Обучение модели первого уровня

ARTM – регуляризованный PLSA (Probabilistic Latent Semantic Analysis)

Основан на итерационном алгоритме:

## EM-алгоритм

$$\text{E-шаг: } p(t|d, w) = \frac{p(w|t)p(t|d)}{\sum_{t' \in T} p(w|t')p(t'|d)} = \frac{\varphi_{wt}\theta_{td}}{\sum_{t' \in T} \varphi_{wt'}\theta_{t'd}};$$

$$\text{M-шаг: } \varphi_{wt} \propto \left( n_{wt} + \varphi_{wt} \frac{\partial R}{\partial \varphi_{wt}} \right)_+; \quad n_{wt} = \sum_{d \in D} n_{dw} p(t|d, w);$$

$$\theta_{td} \propto \left( n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right)_+; \quad n_{td} = \sum_{w \in W} n_{dw} p(t|d, w).$$

# Обучение моделей нижних уровней

Обобщение ARTM на случай трехматричного разложения  
Итерационный алгоритм:

## EM-алгоритм

$$\mathbf{E}\text{-шаг: } p(t, s|d, w) = \frac{p(w|s)p(s|t)p(t|d)}{\sum_{t' \in T} \sum_{s' \in S} p(w|s')p(s'|t')p(t'|d)} = \frac{\varphi_{ws}\psi_{st}\pi_{td}}{\sum_{s' \in S} \sum_{t' \in T} \varphi_{ws'}\psi_{s't'}\pi_{t'd}};$$

$$\mathbf{M}\text{-шаг: } \varphi_{ws} \propto \left( n_{ws} + \varphi_{ws} \frac{\partial R}{\partial \varphi_{ws}} \right)_+ ; n_{ws} = \sum_{d \in D} \sum_{t \in T} n_{dw} p(t, s|d, w)$$
$$\psi_{st} \propto \left( n_{st} + \psi_{st} \frac{\partial R}{\partial \psi_{st}} \right)_+ ; n_{st} = \sum_{d \in D} \sum_{w \in W} n_{dw} p(t, s|d, w) .$$



## Требования к тематическим распределениям слов:

- На каждом уровне выделяем множества предметных и фоновых тем (по одной на каждую предметную тему родительского уровня); в корневой вершине фоновая тема – это слова общей лексики языка и коллекции, на следующих уровнях – слова общей лексики родительских тем;
- Предметные темы декоррелированы между собой и с фоновыми темами;
- Предметные и фоновые темы разрежены, т.е. содержат небольшое количество слов.

Формализация требований в виде регуляризатора матрицы  $\Phi$ :

$$R(\Phi) = -\tau_1 \sum_{w \in W} \sum_{s \in Sub} \beta_w \ln \varphi_{ws} - \tau_2 \sum_{w \in W} \sum_{s \in Bcg} \beta_w \ln \varphi_{ws} - \tau_3 \sum_{s, s' \in S} \sum_{w \in W} \varphi_{ws} \varphi_{ws'}$$

## Возможные конфигурации матрицы $\Psi$ :

- Если  $\Psi$  плотная, то мы имеем полный многодольный граф, если разреженная – граф становится деревом;
- Разреживание  $\Psi$  позволяет модели самой распределять потомков между родителями, определять их количество, а также разрешает множественное наследование потомков;
- Обнуляя некоторые значения в начальной инициализации  $\Psi$ , мы наоборот задаем структуру графа извне.

## Регуляризаторы матрицы $\Psi$ :

$$R(\Psi) = -\tau_1 \sum_{s \in S} \sum_{t \in T} \ln \psi_{st} + \tau_2 \sum_{s \in Bcg} \sum_{t \in T} \ln \psi_{st}$$

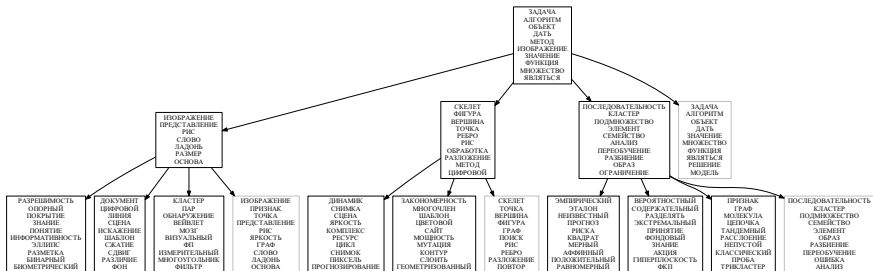
Эксперименты проводятся на коллекции документов конференций «Математические методы распознавания образов» и «Интеллектуализация обработки информации».

$|D| = 850$ ,  $|W| = 13031$

Веб-навигатор:

[MMRONavigator.vv.si](http://MMRONavigator.vv.si)



# Фрагмент иерархии



root - Chromium

Друзья Надежды | Диалоги | conference-presen | beamer conferenc | beamer - conferen | root

mmrnavigator.vv.si/root.html

# тематический навигатор

Главная | Корневая тема | О проекте

## Тема root (уровень 0)



Подтема	Слова(13031)	Авторы(1016)	Статьи(850)
Подтема 1.0 изображение обработка фотамет	задача алгоритм объект	Кельманов А. В. Моттль В. В. Красотина О. В.	
Подтема 1.1 система признак текст слово знание	дать метод изображение значение	Пытьев Ю. П. Моттль В. В. Дорофеек А. А. Воронцов К. В.	
Подтема 1.2 скелет система точка фигура	функция множество являются решение	Дюкова Е. В. Панкратов А. Н. Ветров Д. П. Середин О. С.	
Подтема 1.3 признак набор описание вешина	модель оценка работа точка признак класс	Дорофеек Ю. А. Кельманов А. В. Чехович Ю. В. Кропотов Д. А. Хамидуллин С. А. Федотов Н. Г.	<ul style="list-style-type: none"> <li>Внедрение системы Антиплагиат в Российской государственной библиотеке</li> <li>Применение методов распознавания образов в системе управления коллекциями графических документов</li> <li>Формализация задачи распознавания последовательности состояний сложного источника</li> <li>Прикладные технологии распознавания количественных характеристики растительности по цифровым многоспектральным и гиперспектральным аэрокосмическим изображениям</li> <li>Построение параметрического портрета динамической системы на основе синдромальных представлений</li> <li>Автоматизированная система мониторинга финансовых рынков Check4Trick</li> <li>Байесовский подход к задаче обучения распознаванию образов в нестационарной генеральной совокупности</li> <li>СРС-методы как методы интеллектуального анализа данных при исследовании реальных хаотических процессов</li> <li>Анализ техники живописи по изображениям в задачах атрибуции. Обзор</li> <li>Параметрическое семейство гранично-скелетных моделей формы</li> <li>Новая технология численного исследования динамических систем методами распознавания образов</li> <li>Вычислительная методика обработки и интерпретации</li> </ul>
Подтема 1.4 граф покрытие описание вешина	система результат вектор	Дьячкова А. Г. Ивахненко А. А. Сулимова В. В. Середин О. С.	
Подтема bcg 1.5	параметр случай распознавание мощь	Дедус Ф. Ф. Стрижков В. В. Чуличков А. И. Татарчук А. И.	

Menu [Системный мони... [compute\_p\_s.py (... [Terminal root - Chromium Chirkova\_Iomonos... [конференция МФ... en Вт., 14 апр., 23:59

1,2 - Chromium

Друзья Надежды x Диалоги x conference-presen x beamer conferenc x beamer - conferen x 1,2

mmronavigator.vv.si/1,2.html

# тематический навигатор

Главная | Корневая тема | О проекте

## Тема 1,2 (уровень 1)



Надтема root	Слова(2988)	Авторы(1016)	Статьи(164)
<p>Подтема 2.0</p> <p>ансамбль точечный источник технология</p> <p>Подтема 2.2</p> <p>область преобразование</p> <p>Подтема 2.3</p> <p>преобразование символ метрика</p> <p>Подтема 2.8</p> <p>классификатор функционал контур</p> <p>Подтема 2.9</p> <p>набор ансамбль равный довт</p>	<p>скелет система точка фигура закономерность траектория движение уровень длина сегмент окрестности описание эталон исследование локальный ряд фрагмент рис состояние построение примитив аксиома поток граница разметка</p>	<p>Кельманов А. В. Моттль В. В. Красотина О. В. Пытьев Ю. П. Моттль В. В. Дорофеюк А. Воронцов К. В. Доюкова Е. В. Панаратов А. Н. Ветров Д. П. Середин О. С. Дорофеюк Ю. А. Кельманов А. В. Чехович Ю. В. Кропотов Д. А. Хамидуллин С. А. Федотов Н. Г. Дьячков А. Г. Ивахненко А. А. Сулимова В. В. Середин О. С. Дедус Ф. Ф. Стрижов В. В. Чуличков А. И. Татарчук А. И.</p>	<ul style="list-style-type: none"> <li>• Параметрическое семейство гранично-скелетных моделей формы</li> <li>• Методы и средства распознавания сетевых атак с помощью нейросетевых и иммунолепечных технологий</li> <li>• Методы и алгоритмы распознавания объектов сельских поселений на цифровой карте</li> <li>• Построение обобщенных скелетов многоугольных бинарных фигур с многоугольными выпуклыми структурирующими элементами</li> <li>• О некоторых аспектах интеллектуального анализа пучков временных рядов</li> <li>• О свойствах задач и алгоритмов разметки элементов точечных конфигураций</li> <li>• Многолистая многоугольная фигура и ее скелет</li> <li>• Критерии локальной разрешимости и регулярности как инструмент исследования морфологии аминокислотных последовательностей</li> <li>• Система верификации владельца карманного компьютера по фотопортрету</li> <li>• Процедура оптимального ларного выравнивания для сравнения скелетных графов, заданных цепочками примитивов</li> <li>• О классификационном подходе к имитационному моделированию транспортных потоков</li> <li>• Формализация жестикуляции с помощью направленного</li> </ul>

Menu [Системный мони... [compute\_p\_s.py (... [Terminal] 1,2 - Chromium Chirkova\_Iomonos... [конференция МФ... en 15 апр., 00:00

2,2 - Chromium

Друзья Надежды x Диалоги x conference-presen x beamer conferenc x beamer - conferen x 2,2

mmronavigator.vv.si/2,2.html

# тематический навигатор

Главная | Корневая тема | О проекте

## Тема 2,2 (уровень 2)

Надтема 1.1

Надтема 1.2

### Слова(8736)

область преобразование высокий показать рассматривать конечный кривая требовать порядок пример схема фиксировать функционал определять контур описывать ребро набор пиксель форма собственный способ получить код характеристика

### Авторы(1016)

Кельманов А. В.  
Моттль В. В.  
Красотина О. В.  
Пытьев Ю. П.  
Моттль В. В.  
Дорофеюк А. А.  
Воронцов К. В.  
Доикова Е. В.  
Панаратов А. Н.  
Ветров Д. П.  
Середин О. С.  
Дорофеюк Ю. А.  
Чехович Ю. В.  
Кельманов А. В.  
Кропотов Д. А.  
Хамидуллин С. А.  
Федотов Н. Г.  
Дьячков А. Г.  
Ивахненко А. А.  
Сулимова В. В.  
Середин О. С.  
Дедус Ф. Ф.  
Стрижов В. В.  
Чуличков А. И.  
Татарчук А. И.

### Статьи(297)



- Параметрическое семейство гранично-скелетных моделей формы
- Методы и средства распознавания сетевых атак с помощью нейросетевых и иммуносетевых технологий
- Методы и алгоритмы распознавания объектов сельских поселений на цифровой карте
- Построение обобщенных скелетов многоугольных бинарных фигур с многоугольными выпуклыми структурирующими элементами
- О некоторых аспектах интеллектуального анализа пучков временных рядов
- О свойствах задач и алгоритмов разметки элементов точечных конфигураций
- Многолистая многоугольная фигура и ее скелет
- Критерии локальной разрешимости и регулярности как инструмент исследования морфологии аминокислотных последовательностей
- Система верификации владельца карманного компьютера по фотопортрету
- Процедура оптимального лярного выравнивания для сравнения скелетных графов, заданных цепочками примитивов
- О классификационном подходе к имитационному моделированию транспортных потоков
- Формализация жестики, пальмы с локтевым запястьем

Menu [Системный мони... [compute\_p\_s.py (... [Terminal] 2,2 - Chromium Chirkova\_Iomonos... [конференция МФ... en 15 апр., 00:01

1,0 - Chromium

Друзья Надежды x Диалоги x conference-presen x beamer conferenc x beamer - conferen x 1,0

mmronavigator.vv.si/1,0.html

# тематический навигатор

Главная | Корневая тема | О проекте

## Тема 1,0 (уровень 1)

Надтема root	Слова(3836)	Авторы(1016)	Статьи(189)
<p>Подтема 2.4</p> <ul style="list-style-type: none"> <li>виртуальный итоговый дискретный</li> </ul> <p>Подтема 2.6</p> <ul style="list-style-type: none"> <li>признак форма кадр контур ось сдвиг</li> </ul> <p>Подтема 2.8</p> <ul style="list-style-type: none"> <li>классификатор функционал контур</li> </ul> <p>Подтема 2.10</p> <ul style="list-style-type: none"> <li>преобразование область плоскость</li> </ul> <p>Подтема 2.11</p> <ul style="list-style-type: none"> <li>область плоскость фоома поедлагаться</li> </ul>	<p>изображение обработка фрагмент рис система цифровой палец ладонь область спектр повтор строка текст спектральный длина участок блок точка баз периодичность скрыть сегментация человек днк документ</p>	<p>Кельманов А. В. Мотль В. В. Красотина О. В. Пытьев Ю. П. Мотль В. В. Дорофеюк А. А. Воронцов К. В. Доюкова Е. В. Панаратов А. Н. Ветров Д. П. Середин О. С. Дорофеюк Ю. А. Кельманов А. В. Чехович Ю. В. Кропотов Д. А. Хамидуллин С. А. Федотов Н. Г. Дьячков А. Г. Иващенко А. А. Сулимова В. В. Середин О. С. Дедус Ф. Ф. Стрижов В. В. Чуличков А. И. Татарчук А. И.</p>	<ul style="list-style-type: none"> <li>Внедрение системы Антиплагиат в Российской государственной библиотеке</li> <li>Система Антиплагиат.РФ: задачи, проблемы, результаты, перспективы</li> <li>Комплекс и технологии тематической обработки данных дистанционного зондирования Земли</li> <li>Распознавание скрытой периодичности в геномах модельных организмов</li> <li>Вычислительные методы обработки и интерпретации многоспектральных и гиперспектральных аэрокосмических изображений</li> <li>Распознавание скрытой периодичности в кодирующих последовательностях ДНК</li> <li>Прикладные технологии распознавания количественных характеристик растительности по цифровым многоспектральным и гиперспектральным аэрокосмическим изображениям</li> <li>Об опыте анализа текстов с помощью системы Антиплагиат.РФ</li> <li>Спектральный метод дифференцирования функций в задаче поиска мегасателлитных tandemных повторов</li> <li>Комплекс программно-информационных средств оперативного дешифрирования космических изображений</li> <li>Распознавание природно-техногенных объектов по данным гиперспектральных систем аэрокосмического зондирования</li> </ul>



Menu [Системный мони... [compute\_p\_s.py (... [Terminal 1,0 - Chromium Chirkova\_Iomonos... [конференция МФ... en Cp., 15 апр., 00:01



2,10 - Chromium

Друзья Надежды x Диалоги x conference-presen x beamer conferenc x beamer - conferen x 2,10 x

mmronavigator.vv.si/2,10.html

MMPO   **тематический навигатор**

Главная | Корневая тема | О проекте

## Тема 2,10 (уровень 2)

Надтема 1.0

Надтема 1.2

Надтема 1.4

### Слова(9545)

преобразование  
область  
плоскость  
набор  
форма  
характеристика  
профиль  
эффективный  
равный  
ребро  
приближение  
граница  
ограничение  
структурный  
определять  
группа  
способ  
дальнейший  
поле  
рамка  
двумерный  
дискретный  
путь  
исследовать  
ось

### Авторы(1016)

Кельманов А. В.  
Моттль В. В.  
Красотина О. В.  
Пытьев Ю. П.  
Моттль В. В.  
Дорофеюк А. А.  
Воронцов К. В.  
Доикова Е. В.  
Панаратов А. Н.  
Ветров Д. П.  
Середин О. С.  
Дорофеюк Ю. А.  
Кельманов А. В.  
Чехович Ю. В.  
Кропотов Д. А.  
Хамидуллин С. А.  
Федотов Н. Г.  
Дьячков А. Г.  
Ивахненко А. А.  
Сулимова В. В.  
Середин О. С.  
Дедус Ф. Ф.  
Стрижков В. В.  
Чуличков А. И.  
Татарчук А. И.

### Статьи(449)

- Параметрическое семейство гранично-скелетных моделей формы
- Методы и средства распознавания сетевых атак с помощью нейросетевых и иммуноключевых технологий
- Методы и алгоритмы распознавания объектов сельских поселений на цифровой карте
- Построение обобщенных скелетов многоугольных бинарных фигур с многоугольными выпуклыми структурирующими элементами
- О некоторых аспектах интеллектуального анализа пучков временных рядов
- О свойствах задач и алгоритмов разметки элементов точечных конфигураций
- Многолистая многоугольная фигура и ее скелет
- Критерии локальной разрешимости и регулярности как инструмент исследования морфологии аминокислотных последовательностей
- Система верификации владельца карманного компьютера по фотопортрету
- Процедура оптимального парного выравнивания для сравнения скелетных графов, заданных цепочками примитивов
- Новая технология численного исследования динамических систем методами распознавания образов
- О плоскостных разбиениях плоскости и имитационных...

Menu [Системный мони... [compute\_p\_s.py (... [Terminal] 2,10 - Chromium Chirkova\_Iomonos... [конференция МФ... en 15 апр., 00:02

# Результаты и перспективы

## Результаты:

- предложен способ построения тематических иерархий с помощью двух- и трехматричных разложений матрицы частот слов;
- разработана пилотная версия тематического навигатора ММРО-ИОИ.

## Дальнейшие исследования:

- оценивание и улучшение качества иерархии;
- более детальное исследование влияния различных стратегий регуляризации на качество иерархии;
- частичное обучение (использование экспертной разметки);
- автоматическое именоване тем.

*Zavitsanos E., Paliouras G., Vouros G. A.* Non-parametric estimation of topic hierarchies from texts with hierarchical Dirichlet processes // *Journal of Machine Learning Research*, 2011. Vol. 12. Pp. 2749–2775.

*Wang C., Danilevsky M., Desai N., Zhang Y., Nguyen P., Taula T., Han J.* A Phrase Mining Framework for Recursive Construction of a Topical Hierarchy // *KDD '13, Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2013. Pp. 437–445.

*Wang C., Liu X., Song Y., Han J.* Scalable and robust construction of topical hierarchies // *arXiv:1403.3460v1 [cs.LG]* 13 Mar 2014.

*Vorontsov K. V., Potapenko A. A.* Tutorial on Probabilistic Topic Modeling: Additive Regularization for Stochastic Matrix Factorization // *AIST'2014, Analysis of Images, Social networks and Texts*. Springer International Publishing Switzerland. *Communications in Computer and Information Science (CCIS)*. 2014, Vol. 436. Pp. 29–46.