

LDA тематическое моделирование с использованием pyro

План

- Задача
- Тренировочный и тестовый сет
- LDA
- Pyro
- Схема процесса
- Примеры
- Результаты

Задача

- Подготовить тренировочный и тестовый наборы данных из webis-ambient15 - <https://webis.de/data/webis-ambient-15.html>
- Вычислить LDA-модель корпуса webis-ambient15
- Проецировать обычные текстовые документы в пространство LDA
- Проецировать темы в пространство LDA
- Реализовать в pyro - <https://docs.pyro.ai/en/stable/>

Webis-Ambient-15: Документы

- 3.100 <http://www.dog-tracker.com/breeds/beagle.html> Beagle, puppies, dog, breeder, stud service, standard Beagle breed info and standard, with free ads for puppies for sale , stud service and a lot more. ... Known as the "singing Beagle," he has a sweet hunting ...
- 4.1 <http://en.wikipedia.org/wiki/Bronx> The Bronx - Wikipedia, the free encyclopedia ... five boroughs were independent cities, the Bronx would rank as the ninth most ... The name refers to the Bronx River, and rivers are commonly referred to with the ...
- 5.99 <http://www.amazon.ca/s?ie=UTF8&keywords=cain&index=books&page=1> Amazon.ca: cain: Books ... Decision Making, conflict by Jim Cain and Barry Jolliff (Paperback - Nov 1997) ... From America's #1 Food Magazine by Anne Chappell Cain (Hardcover - Sep 15 2006) ...
- 6.93 <http://www.camelsaust.com.au/> Camels Australia Export - Central Australia - Northern Territory - Australia ... Camel Industry ... sustainable development of the camel industry through the use, ... estimated the present feral camel population in the Northern ...
- 7.65 <http://www.coralseavillastobago.com/> Coral Sea Villas Tobago Coral Sea Villas, Bon Accord, Tobago. Caribbean villas, 5 minutes from the airport and close to beaches, shops, restaurants
- ...

Webis-Ambient-15: Темы

- 3 Beagle
- 4 Bronx
- 5 Cain
- 6 Camel
- 7 Coral Sea
- 8 Cube
- 9 Eos
- 10 Excalibur
- 32 Purple Haze
- 33 Raam
- 34 Rhea
- 35 Scorpion
- 36 The Little Mermaid
- 37 Tortuga
- 38 Urania
- 39 Wink
- 40 Xanadu
- ...

LDA: процесс

Topics

gene 0.04
dna 0.02
genetic 0.01
...

life 0.02
evolve 0.01
organism 0.01
...

brain 0.04
neuron 0.02
nerve 0.01
...

data 0.02
number 0.02
computer 0.01
...

Documents

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson at Uppsala University in Sweden. "We arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

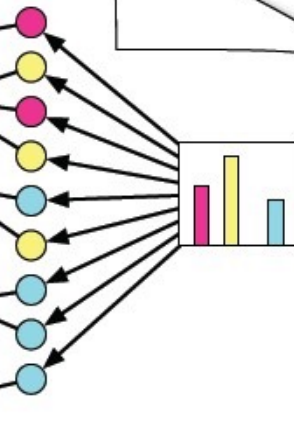


* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments



LDA: модель

1. Маша любит зеленые и синие вещи.
2. Петя нравятся красные и розовые вещи.
3. Даша любит прямоугольники и треугольники.
4. Дима нравятся круги.
5. Паша любит зеленые треугольники.

Распределение по предложениям:

Предложение 1 и 2: 100% цветов темы

Предложение 3 и 4: 100% тематических форм

Предложение 5: 50% цвета темы, 50% формы темы

Topic "Colors": { 20% зеленый. 20% синий. 20% красный. 20% розовый. }

Тема "Формы": { 30% кругов. 30% прямоугольников. 30% треугольников. }

LDA: процедура

Семплирование Гиббса:

Набор документов с предполагаемым количеством тем T

Случайно присваиваем тему каждому слову в каждом документе

Для каждого документа d Для каждого слова w

Вычислить $p(t | d)$ // долю слов в документе d , которые относятся к теме t

Вычислить $p(w | t)$ // доля назначений на тему t среди всех документов, которые приходятся на данное слово w .

Присвоить w новую тему t с вероятностью $p(\text{тема } t | \text{документ } d) * p(\text{слово } w | \text{тема } t)$.

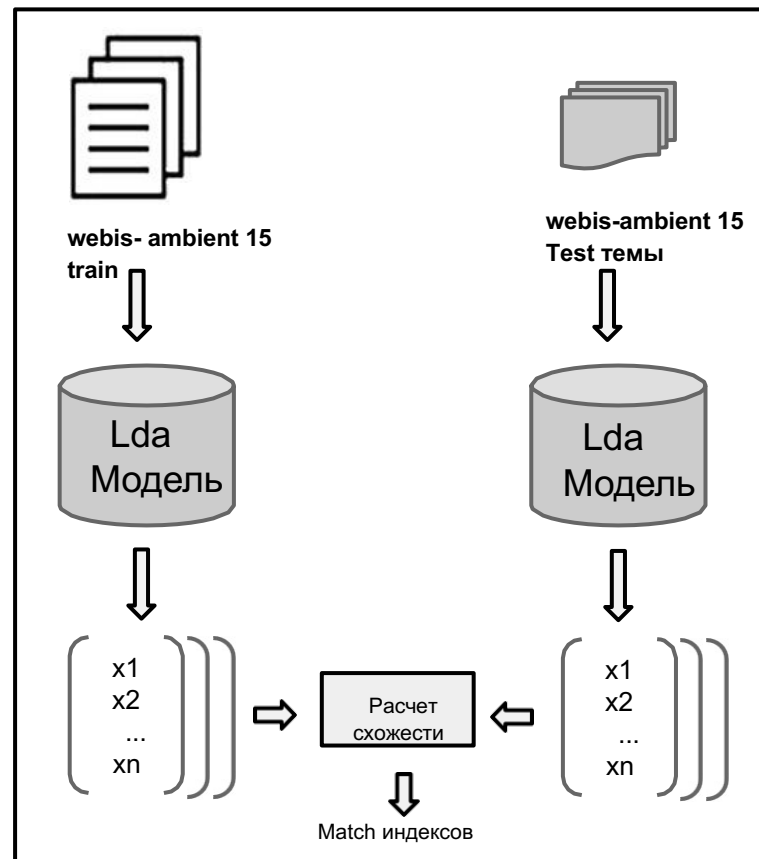
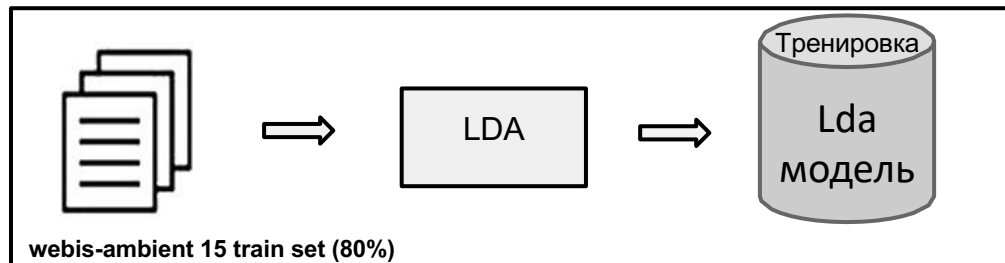
Pyro: модель

- Библиотека тематического моделирования, реализованная на языке программирования Python
- Использует фреймворк pytorch
- Подаем на вход данные как массив идентификаторов слов в форме [num_words_per_doc, num_documents]
- Используем вывод локальных тематических переменных, достигаемый с помощью многослойного перцептрона
- Используем распределения Гамма и Дирихле
- Задействуем Adam оптимизатор и gradient clipping

Pyro: параметры модели

- num-topics
- num-words
- num-docs
- num-words-per-doc
- num-steps
- layer-sizes
- learning-rate
- batch-size

Схема процесса



Примеры тем

Пример 1

$20.5 * \text{query} + 10.6 * \text{queries} + 3.1 * \text{index} + 1.3 * \text{database} + 1.0 * \text{answer} + 0.9 * \text{indexing} + 0.7 * \text{databases} + 0.6 * \text{answers} + 0.6 * \text{indexes} + 0.6 * \text{querying}$

Пример 2

$13.0 * \text{server} + 8.5 * \text{client} + 5.3 * \text{request} + 4.1 * \text{requests} + 3.5 * \text{clients} + 2.0 * \text{session} + 1.7 * \text{proxy} + 1.5 * \text{remote} + 1.3 * \text{response} + 1.0 * \text{service}$

Пример 3

$19.5 * \text{pp} + 10.3 * \text{fig} + 8.7 * \text{proc} + 5.2 * \text{vol} + 3.5 * \text{ii} + 2.3 * \text{conf} + 1.9 * \text{iii} + 1.8 * \text{he} + 1.6 * \text{unit} + 1.5 * \text{int}$

Результаты

Topic x1 - 82.61%

$0.042 \cdot \text{translation} + 0.035 \cdot \text{word} + 0.023 \cdot \text{english} + 0.015 \cdot \text{alignment} +$
 $0.013 \cdot \text{phrase} + 0.013 \cdot \text{chinese} + 0.011 \cdot \text{sentence} + 0.009 \cdot \text{mt} + 0.009 \cdot \text{corpus} +$
 $0.009 \cdot \text{target}$

Topic x2 - 78.05%

$0.020 \cdot \text{word} + 0.013 \cdot \text{corpus} + 0.013 \cdot \text{sentence} + 0.010 \cdot \text{lexical} + 0.010 \cdot \text{verb} +$
 $0.009 \cdot \text{noun} + 0.009$
 $\cdot \text{semantic} + 0.009 \cdot \text{sentences} + 0.009 \cdot \text{syntactic} + 0.008 \cdot \text{linguistics}$

Topic x3 – 80.11%

$27 : 0.068 \cdot \text{items} + 0.058 \cdot \text{item} + 0.026 \cdot \text{trust} + 0.024 \cdot \text{recommendation} +$
 $0.021 \cdot \text{ratings} + 0.019 \cdot \text{rating} + 0.018 \cdot \text{profile} + 0.015 \cdot \text{recommendations} +$
 $0.014 \cdot \text{recommender} + 0.014 \cdot \text{collaborative}$