

# ОТ НАИВНОГО БАЙЕСА ДО ТЕМАТИЧЕСКОГО МОДЕЛИРОВАНИЯ II

---

Сергей Николенко

НИУ ВШЭ — Санкт-Петербург

13 октября 2017 г.

---

*Random facts:*

- 13 октября 54 г. жена императора Клавдия отравила его грибами, и на престол взошел Нерон
- 13 октября 1884 г. Гринвич был утвержден как место прохождения нулевого меридиана
- 13 октября 1917 г. в португальском городе Фатима случилось чудо: явление Девы Марии и чудесное превращение Солнца, официально признанные чудом католической церковью

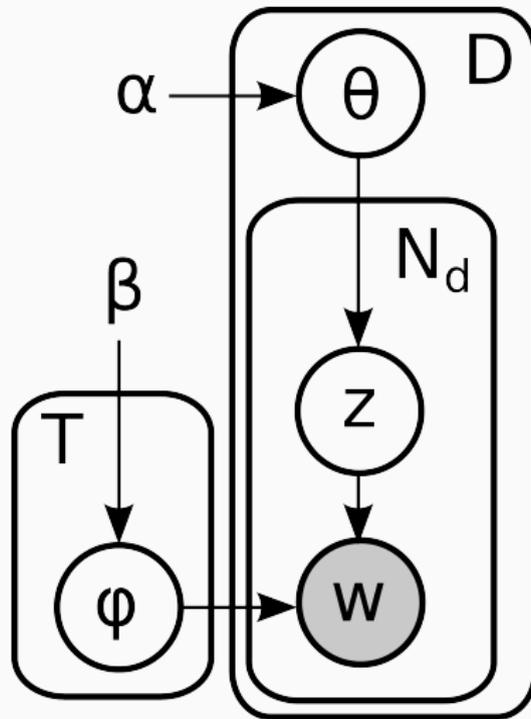
LDA

---

- Развитие идей pLSA – LDA (Latent Dirichlet Allocation).
- Это фактически байесовский вариант pLSA, сейчас нарисуем картинку, добавим априорные распределения и посмотрим, как работают наши методы приближённого вывода.
- Задача та же: смоделировать большую коллекцию текстов (например, для information retrieval или классификации).

- У одного документа может быть несколько тем. Давайте построим иерархическую байесовскую модель:
  - на первом уровне – смесь, компоненты которой соответствуют «темам»;
  - на втором уровне – мультиномиальная переменная с априорным распределением Дирихле, которое задаёт «распределение тем» в документе.

- Если формально: слова берутся из словаря  $\{1, \dots, V\}$ ; слово – это вектор  $w$ ,  $w_i \in \{0, 1\}$ , где ровно одна компонента равна 1.
- Документ – последовательность из  $N$  слов  $\mathbf{w}$ . Нам дан корпус из  $M$  документов  $\mathcal{D} = \{\mathbf{w}_d \mid d = 1..M\}$ .
- Генеративная модель LDA выглядит так:
  - выбрать  $\theta \sim \text{Di}(\alpha)$ ;
  - для каждого из  $N$  слов  $w_n$ :
    - выбрать тему  $z_n \sim \text{Mult}(\theta)$ ;
    - выбрать слово  $w_n \sim p(w_n \mid z_n, \beta)$  по мультиномиальному распределению.



# LDA: ЧТО ПОЛУЧАЕТСЯ [BLEI, 2012]

## 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

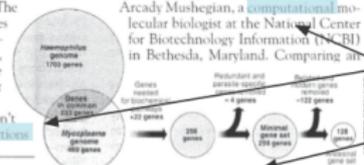
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here, two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known 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

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

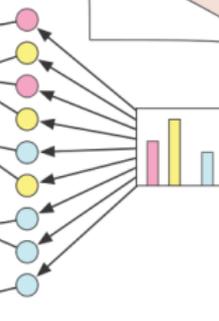
SCIENCE • VOL. 272 • 24 MAY 1996

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson at the University in Sweden. He also 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



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

## Topic proportions and assignments



- Два основных подхода к выводу в сложных вероятностных моделях, в том числе LDA:
  - *вариационные приближения*: рассмотрим более простое семейство распределений с новыми параметрами и найдём в нём наилучшее приближение к неизвестному распределению;
  - *сэмплирование*: будем набрасывать точки из сложного распределения, не считая его явно, а запуская марковскую цепь под графиком распределения (частный случай – сэмплирование по Гиббсу).
- Сэмплирование по Гиббсу обычно проще расширить на новые модификации LDA, но вариационный подход быстрее и часто стабильнее.

- Рассмотрим задачу байесовского вывода, т.е. оценки апостериорного распределения  $\theta$  и  $z$  после нового документа:

$$p(\theta, \mathbf{z} \mid \mathbf{w}, \alpha, \beta) = \frac{p(\theta, \mathbf{z}, \mathbf{w} \mid \alpha, \beta)}{p(\mathbf{w} \mid \alpha, \beta)}.$$

- Правдоподобие набора слов  $\mathbf{w}$  оценивается как

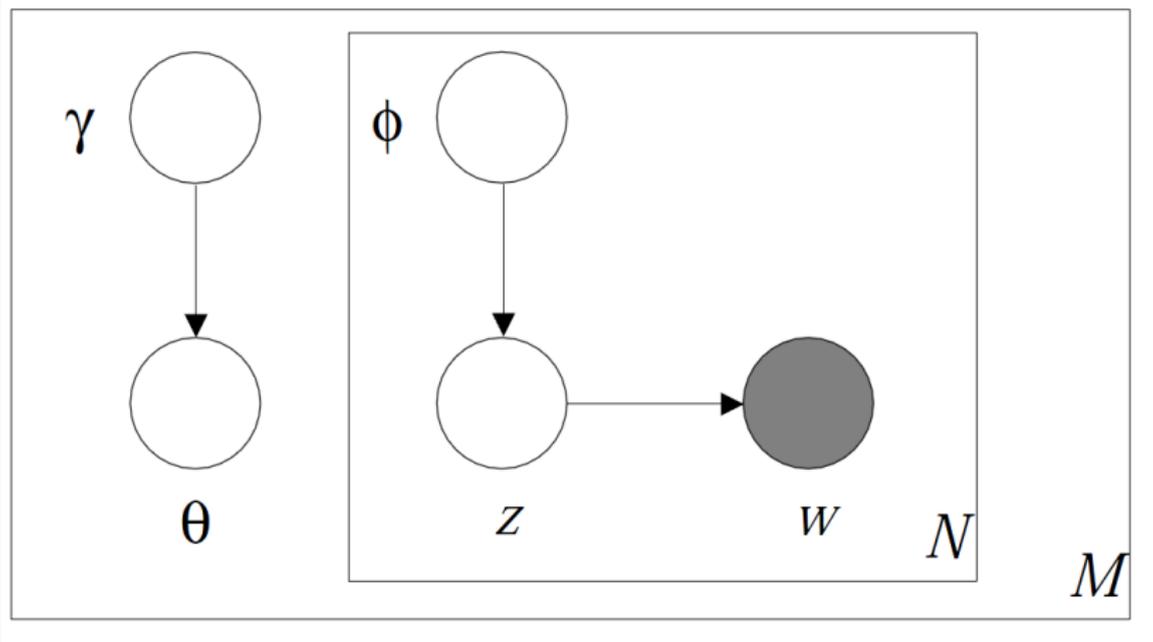
$$p(\mathbf{w} \mid \alpha, \beta) = \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \int \left[ \prod_{i=1}^k \theta_i^{\alpha_i - 1} \right] \left[ \prod_{n=1}^N \sum_{i=1}^k \prod_{j=1}^V (\theta_i \beta_{ij})^{w_n^j} \right] d\theta,$$

и это трудно посчитать, потому что  $\theta$  и  $\beta$  путаются друг с другом.

- Вариационное приближение – рассмотрим семейство распределений

$$q(\theta, z | \mathbf{w}, \gamma, \phi) = p(\theta | \mathbf{w}, \gamma) \prod_{n=1}^N p(z_n | \mathbf{w}, \phi_n).$$

- Тут всё расщепляется, и мы добавили вариационные параметры  $\gamma$  (Дирихле) и  $\phi$  (мультиномиальный).
- Заметим, что параметры для каждого документа могут быть свои – всё условно по  $\mathbf{w}$ .



- Теперь можно искать минимум KL-расстояния:

$$(\gamma^*, \phi^*) = \arg \min_{(\gamma, \phi)} \text{KL}(q(\theta, z | \mathbf{w}, \gamma\phi) \| p(\theta, \mathbf{z} | \mathbf{w}, \alpha, \beta)).$$

- Для этого сначала воспользуемся уже известной оценкой из неравенства Йенсена:

$$\begin{aligned} \log p(\mathbf{w} | \alpha, \beta) &= \log \int_{\theta} \sum_{\mathbf{z}} p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta) d\theta = \\ &= \log \int_{\theta} \sum_{\mathbf{z}} \frac{p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta) q(\theta, \mathbf{z})}{q(\theta, \mathbf{z})} d\theta \geq \\ &\geq E_q [\log p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta)] - E_q [\log q(\theta, \mathbf{z})] =: \mathcal{L}(\gamma, \phi; \alpha, \beta). \end{aligned}$$

- Распишем произведения:

$$\mathcal{L}(\gamma, \phi; \alpha, \beta) = E_q [p(\theta | \alpha)] + E_q [p(\mathbf{z} | \theta)] + E_q [p(\mathbf{w} | \mathbf{z}, \beta)] - E_q [\log q(\theta)] - E_q [\log q(\mathbf{z})].$$

- Свойство распределения Дирихле: если  $X \sim \text{Di}(\alpha)$ , то

$$E[\log(X_i)] = \Psi(\alpha_i) - \Psi\left(\sum_i \alpha_i\right),$$

где  $\Psi(x) = \frac{\Gamma'(x)}{\Gamma(x)}$  – дигамма-функция.

- Теперь можно выписать каждый из пяти членов.

- Теперь осталось только брать частные производные этого выражения.
- Сначала максимизируем его по  $\phi_{ni}$  (вероятность того, что  $n$ -е слово было порождено темой  $i$ ); надо добавить  $\lambda$ -множители Лагранжа, т.к.  $\sum_{j=1}^k \phi_{nj} = 1$ .
- В итоге получится:

$$\phi_{ni} \propto \beta_{iv} e^{\Psi(\gamma_i) - \Psi(\sum_{j=1}^k \gamma_j)},$$

где  $v$  – номер того самого слова, т.е. единственная компонента  $w_n^v = 1$ .

- Потом максимизируем по  $\gamma_i$ ,  $i$ -й компоненте апостериорного Дирихле-параметра.
- Получится

$$\gamma_i = \alpha_i + \sum_{n=1}^N \phi_{ni}.$$

- Соответственно, для вывода нужно просто пересчитывать  $\phi_{ni}$  и  $\gamma_i$  друг через друга, пока оценка не сойдётся.

- Теперь давайте попробуем оценить параметры  $\alpha$  и  $\beta$  по корпусу документов  $\mathcal{D}$ .
- Мы хотим найти  $\alpha$  и  $\beta$ , которые максимизируют

$$\ell(\alpha, \beta) = \sum_{d=1}^M \log p(\mathbf{w}_d | \alpha, \beta).$$

- Подсчитать  $p(\mathbf{w}_d | \alpha, \beta)$  мы не можем, но у нас есть нижняя оценка  $\mathcal{L}(\gamma, \phi; \alpha, \beta)$ , т.к.

$$\begin{aligned} p(\mathbf{w}_d | \alpha, \beta) &= \\ &= \mathcal{L}(\gamma, \phi; \alpha, \beta) + \text{KL}(q(\theta, z | \mathbf{w}_d, \gamma\phi) \| p(\theta, \mathbf{z} | \mathbf{w}_d, \alpha, \beta)). \end{aligned}$$

- EM-алгоритм:
  1. найти параметры  $\{\gamma_d, \phi_d \mid d \in \mathcal{D}\}$ , которые оптимизируют оценку (как выше);
  2. зафиксировать их и оптимизировать оценку по  $\alpha$  и  $\beta$ .

- Для  $\beta$  это тоже делается нехитро:

$$\beta_{ij} \propto \sum_{d=1}^M \sum_{n=1}^{N_d} \phi_{dni} w_n^j.$$

- Для  $\alpha_i$  получается система уравнений, которую можно решить методом Ньютона.

- В базовой модели LDA сэмплирование по Гиббсу после несложных преобразований сводится к так называемому *сжато* сэмплированию по Гиббсу (collapsed Gibbs sampling), где переменные  $z_w$  итеративно сэмплируются по следующему распределению:

$$p(z_w = t \mid \mathbf{z}_{-w}, \mathbf{w}, \alpha, \beta) \propto q(z_w, t, \mathbf{z}_{-w}, \mathbf{w}, \alpha, \beta) =$$

$$\frac{n_{-w,t}^{(d)} + \alpha}{\sum_{t' \in T} (n_{-w,t'}^{(d)} + \alpha)} \frac{n_{-w,t}^{(w)} + \beta}{\sum_{w' \in W} (n_{-w,t}^{(w')} + \beta)},$$

где  $n_{-w,t}^{(d)}$  – число слов в документе  $d$ , выбранных по теме  $t$ , а  $n_{-w,t}^{(w)}$  – число раз, которое слово  $w$  было порождено из темы  $t$ , не считая текущего значения  $z_w$ ; заметим, что оба этих счётчика зависят от других переменных  $\mathbf{z}_{-w}$ .

- Из сэмплов затем можно оценить переменные модели

$$\theta_{d,t} = \frac{n_{-w,t}^{(d)} + \alpha}{\sum_{t' \in T} (n_{-w,t'}^{(d)} + \alpha)},$$

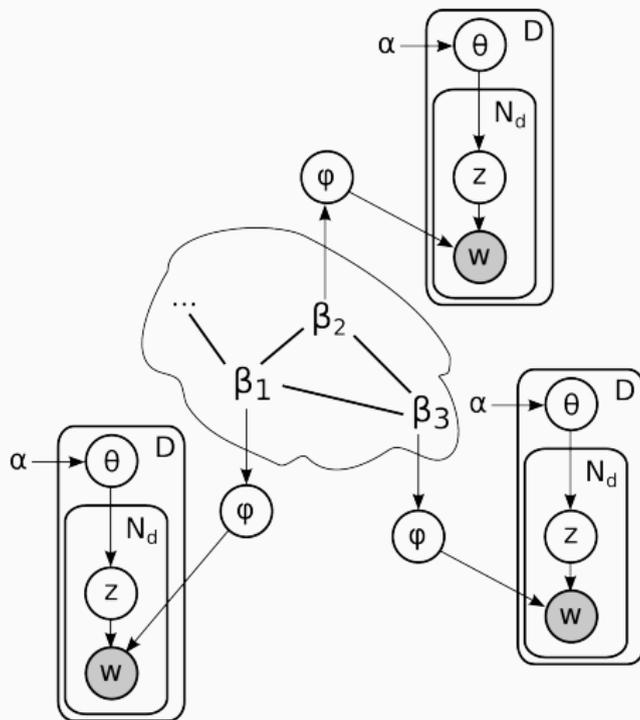
$$\phi_{w,t} = \frac{n_{-w,t}^{(w)} + \beta}{\sum_{w' \in W} (n_{-w,t}^{(w')} + \beta)},$$

где  $\phi_{w,t}$  – вероятность получить слово  $w$  в теме  $t$ , а  $\theta_{d,t}$  – вероятность получить тему  $t$  в документе  $d$ .

- В последние десять лет эта модель стала основой для множества различных расширений.
- Каждое из этих расширений содержит либо вариационный алгоритм вывода, либо алгоритм сэмплирования по Гиббсу для модели, которая, основываясь на LDA, включает в себя ещё и какую-либо дополнительную информацию или дополнительные предполагаемые зависимости.
- Обычно – или дополнительная структура на темах, или дополнительная информация.

- В базовой модели LDA распределения слов по темам независимы и никак не скоррелированы; однако на самом деле, конечно, некоторые темы ближе друг к другу, многие темы делят между собой слова.
- Коррелированные тематические модели (correlated topic models, СТМ); отличие от базового LDA здесь в том, что используется логистическое нормальное распределение вместо распределения Дирихле; логистическое нормальное распределение более выразительно, оно может моделировать корреляции между темами.
- Предлагается алгоритм вывода, основанный на вариационном приближении.

- Марковские тематические модели (Markov topic models, MTM): марковские случайные поля для моделирования взаимоотношений между темами в разных частях датасета (разных корпусах текстов).
- MTM состоит из нескольких копий гиперпараметров  $\beta_i$  в LDA, описывающих параметры разных корпусов с одними и теми же темами. Гиперпараметры  $\beta_i$  связаны между собой в марковском случайном поле (Markov random field, MRF).
- В результате тексты из  $i$ -го корпуса порождаются как в обычном LDA, используя соответствующее  $\beta_i$ .
- В свою очередь,  $\beta_i$  подчиняются априорным ограничениям, которые позволяют «делить» темы между корпусами, задавать «фоновые» темы, присутствующие во всех корпусах, накладывать ограничения на взаимоотношения между темами и т.д.



- Реляционная тематическая модель (relational topic model, RTM) – иерархическая модель, в которой отражён граф структуры сети документов.
- Генеративный процесс в RTM работает так:
  - сгенерировать документы из обычной модели LDA;
  - для каждой пары документов  $d_1, d_2$  выбрать бинарную переменную  $y_{12}$ , отражающую наличие связи между  $d_1$  и  $d_2$ :

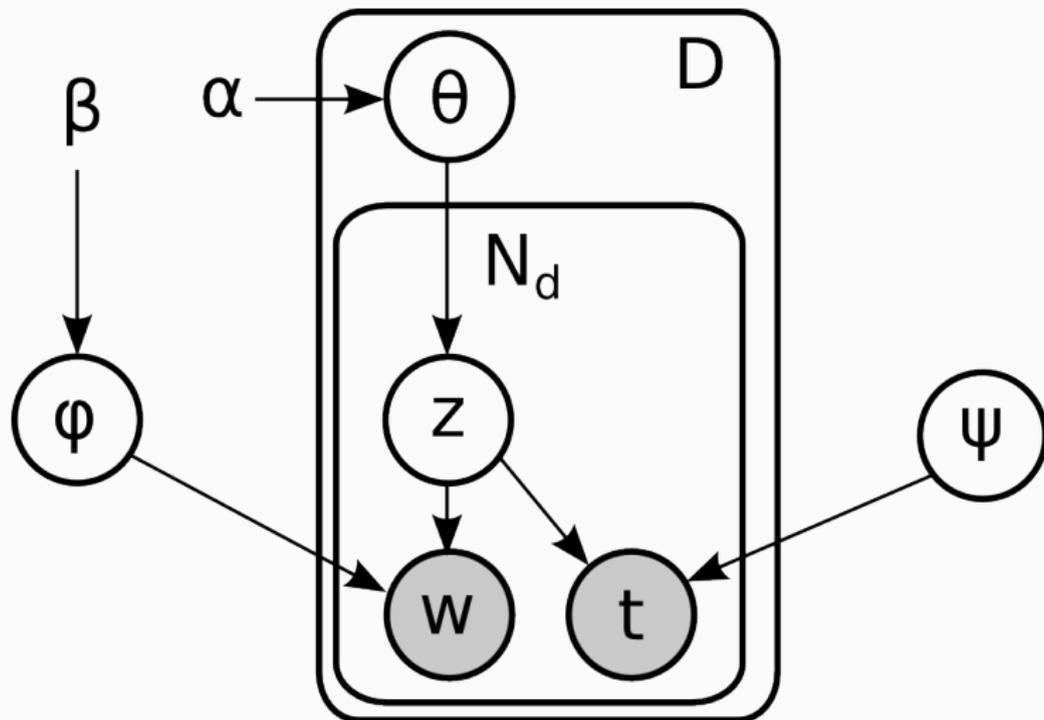
$$y_{12} \mid \mathbf{z}_{d_1}, \mathbf{z}_{d_2} \sim \psi(\cdot \mid \mathbf{z}_{d_1}, \mathbf{z}_{d_2}, \eta).$$

- В качестве  $\psi$  берутся разные сигмоидальные функции; разработан алгоритм вывода, основанный на вариационном приближении.

- Ряд важных расширений LDA касается учёта трендов, т.е. изменений в распределениях тем, происходящих со временем.
- Цель – учёт времени, анализ «горячих» тем, анализ того, какие темы быстро становятся «горячими» и столь же быстро затухают, а какие проходят «красной нитью» через весь исследуемый временной интервал.

- В модели TOT (Topics over Time) время предполагается непрерывным, и модель дополняется бета-распределениями, порождающими временные метки (timestamps) для каждого слова.
- Генеративная модель модели Topics over Time такова:
  - для каждой темы  $z = 1..T$  выбрать мультиномиальное распределение  $\phi_z$  из априорного распределения Дирихле  $\beta$ ;
  - для каждого документа  $d$  выбрать мультиномиальное распределение  $\theta_d$  из априорного распределения Дирихле  $\alpha$ , затем для каждого слова  $w_{di} \in d$ :
    - выбрать тему  $z_{di}$  из  $\theta_d$ ;
    - выбрать слово  $w_{di}$  из распределения  $\phi_{z_{di}}$ ;
    - выбрать время  $t_{di}$  из бета-распределения  $\psi_{z_{di}}$ .

- Основная идея заключается в том, что каждой теме соответствует её бета-распределение  $\psi_z$ , т.е. каждая тема локализована во времени (сильнее или слабее, в зависимости от параметров  $\psi_z$ ).
- Таким образом можно как обучить глобальные темы, которые всегда присутствуют, так и подхватить тему, которая вызвала сильный краткий всплеск, а затем пропала из виду; разница будет в том, что дисперсия  $\psi_z$  будет в первом случае меньше, чем во втором.



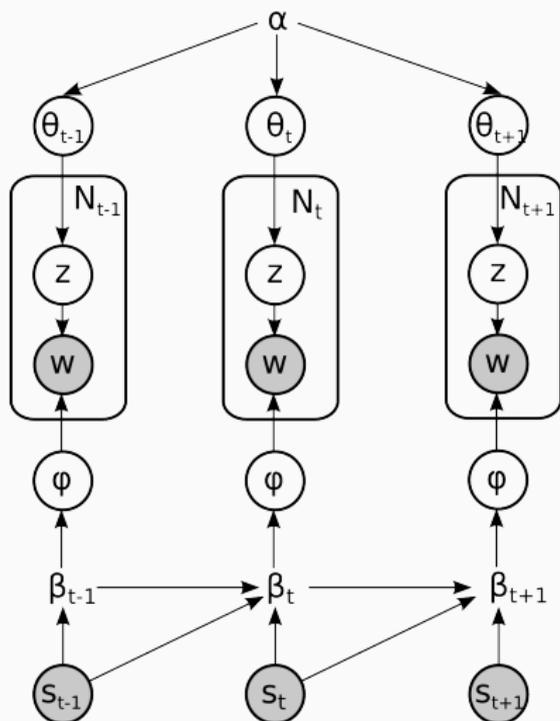
- *Динамические тематические модели* представляют временную эволюцию тем через эволюцию их гиперпараметров  $\alpha$  и/или  $\beta$ .
- Бывают дискретные ([d]DTM), в которых время дискретно, и непрерывные, где эволюция гиперпараметра  $\beta$  ( $\alpha$  здесь предполагается постоянным) моделируется посредством броуновского движения: для двух документов  $i$  и  $j$  ( $j$  позже  $i$ ) верно, что

$$\beta_{j,k,w} \mid \beta_{i,k,w}, s_i, s_j \sim \mathcal{N}(\beta_{i,k,w}, v\Delta_{s_i, s_j}),$$

где  $s_i$  и  $s_j$  – это отметки времени (timestamps) документов  $i$  и  $j$ ,  $\Delta(s_i, s_j)$  – интервал времени, прошедший между ними,  $v$  – параметр модели.

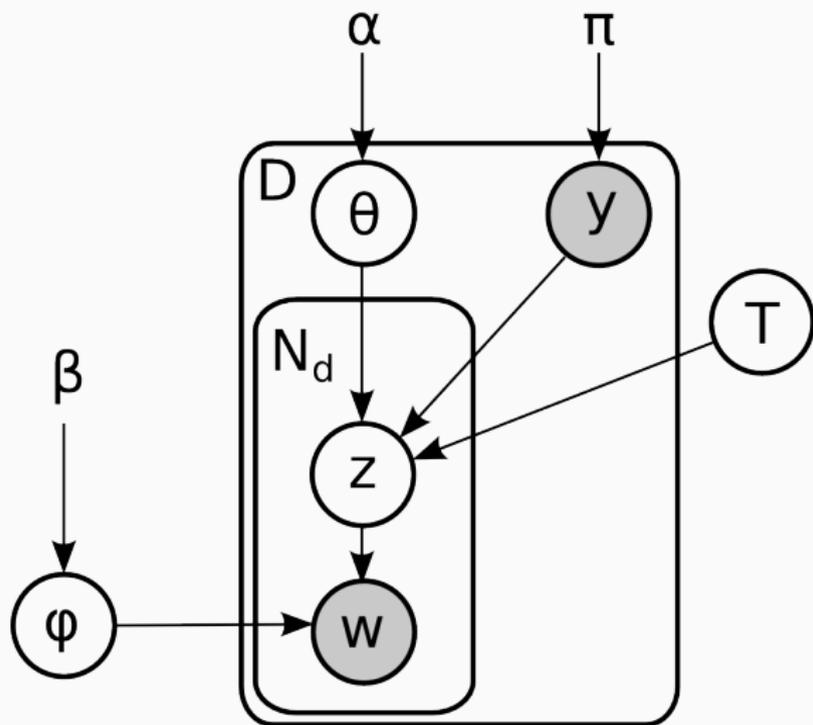
- В остальном генеративный процесс остаётся неизменным.

# НЕПРЕРЫВНАЯ ДИНАМИЧЕСКАЯ ТЕМАТИЧЕСКАЯ МОДЕЛЬ (CDTM)

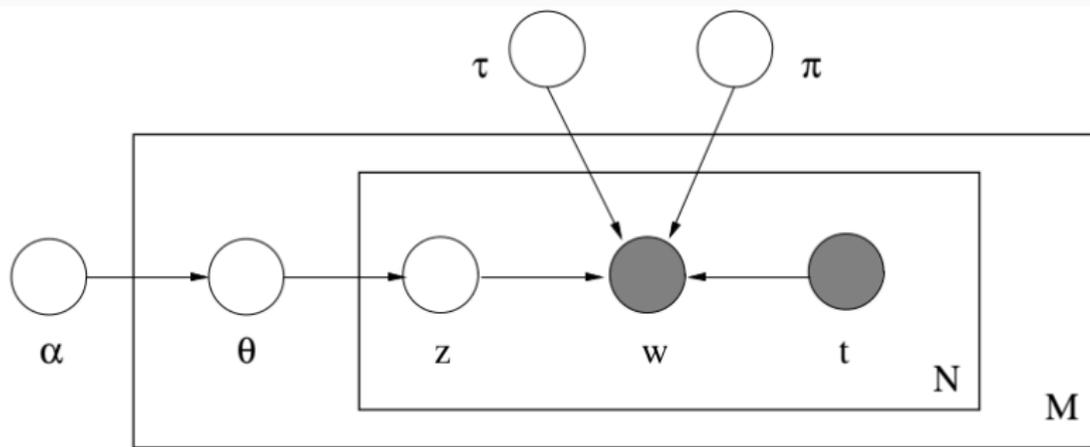


- Supervised LDA: документы снабжены дополнительной информацией, дополнительной переменной отклика (обычно известной).
- Распределение отклика моделируется обобщённой линейной моделью (распределением из экспоненциального семейства), параметры которой связаны с полученным в документе распределением тем.
- Т.е. в генеративную модель добавляется ещё один шаг: после того как темы всех слов известны,
  - сгенерировать переменную-отклик  $y \sim \text{glm}(\mathbf{z}, \eta, \delta)$ , где  $\mathbf{z}$  – распределение тем в документе, а  $\eta$  и  $\delta$  – другие параметры  $\text{glm}$ .
- К примеру, в контексте рекомендательных систем дополнительный отклик может быть реакцией пользователя.

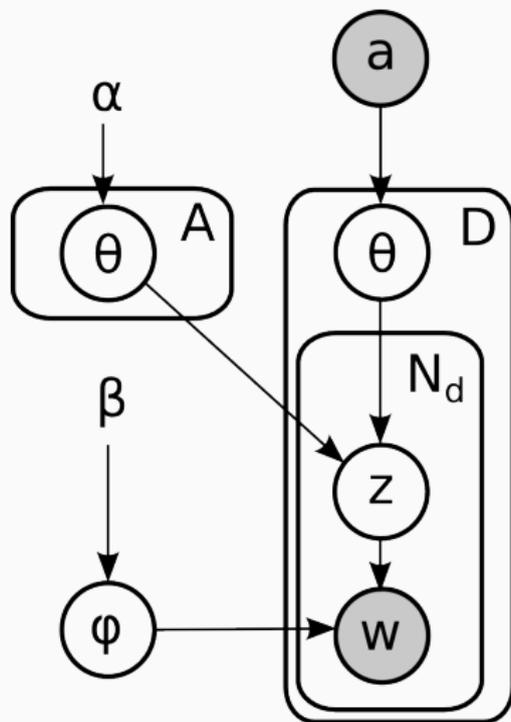
- Дискриминативное LDA (DiscLDA), другое расширение модели LDA для документов, снабжённых категориальной переменной  $y$ , которая в дальнейшем станет предметом для классификации.
- Для каждой метки класса  $y$  в модели DiscLDA вводится линейное преобразование  $T^y : \mathbb{R}^K \rightarrow \mathbb{R}_+^L$ , которое преобразует  $K$ -мерное распределение Дирихле  $\theta$  в смесь  $L$ -мерных распределений Дирихле  $T^y\theta$ .
- В генеративной модели меняется только шаг порождения темы документа  $z$ : вместо того чтобы выбирать  $z$  по распределению  $\theta$ , сгенерированному для данного документа,
  - сгенерировать тему  $z$  по распределению  $T^y\theta$ , где  $T^y$  – преобразование, соответствующее метке данного документа  $y$ .



- TagLDA: слова имеют теги, т.е. документ не является единым мешком слов, а состоит из нескольких мешков, и в разных мешках слова отличаются друг от друга.
- Например, у страницы может быть название – слова из названия важнее для определения темы, чем просто из текста. Или, например, теги к странице, поставленные человеком – опять же, это слова гораздо более важные, чем слова из текста.
- Математически разница в том, что теперь распределения слов в темах – это не просто мультиномиальные дискретные распределения, они факторизованы на распределение слово-тема и распределение слово-тег.



- Author-Topic modeling: кроме собственно текстов, присутствуют их авторы; или автор тоже представляется как распределение на темах, на которые он пишет, или тексты одного автора даже на разные темы будут похожи.
- Базовая генеративная модель Author-Topic model (остальное как в базовом LDA):
  - для каждого слова  $w$ :
    - выбираем автора  $x$  для этого слова из множества авторов документа  $a_d$ ;
    - выбираем тему из распределения на темах, соответствующего автору  $x$ ;
    - выбираем слово из распределения слов, соответствующего этой теме.



- Алгоритм сэмплирования, соответствующий такой модели, является вариантом сжатого сэмплирования по Гиббсу:

$$p(z_w = t, x_w = a \mid \mathbf{z}_{-w}, \mathbf{x}_{-w}, \mathbf{w}, \alpha, \beta) \propto$$

$$\propto \frac{n_{-a,t}^{(a)} + \alpha}{\sum_{t' \in T} (n_{-w,t'}^{(a)} + \alpha)} \frac{n_{-w,t}^{(w)} + \beta}{\sum_{w' \in W} (n_{-w,t}^{(w')} + \beta)},$$

где  $n_{-a,t}^{(a)}$  – то, сколько раз автору  $a$  соответствовала тема  $t$ , не считая текущего значения  $x_w$ , а  $n_{-w,t}^{(w)}$  – число раз, которое слово  $w$  было порождено из темы  $t$ , не считая текущего значения  $z_w$ ; заметим, что оба этих счётчика зависят от других переменных  $\mathbf{z}_{-w}, \mathbf{x}_{-w}$ .

1. Основы обработки текстов: какие задачи нужно уметь решать.
2. Классический метод категоризации: наивный байесовский классификатор.
3. Обобщаем наивный байес: кластеризация EM-алгоритмом.
4. Тематическое моделирование: pLSA, LDA, расширения LDA.

Спасибо за внимание!