

Категоризация текстов и модель LDA

Сергей Николенко

Казанский Федеральный Университет, 2014

Outline

- 1 Категоризация текстов
 - Naive Bayes
 - Latent Dirichlet allocation

Категоризация текстов

- Классическая задача машинного обучения и information retrieval – категоризация текстов.
- Дан набор текстов, разделённый на категории. Нужно обучить модель и потом уметь категоризовать новые тексты.
- Атрибуты a_1, a_2, \dots, a_n – это слова, v – тема текста (или атрибут вроде «спам / не спам»).
- Bag-of-words model: забываем про порядок слов, составляем словарь. Теперь документ – это вектор, показывающий, сколько раз каждое слово из словаря в нём встречается.

Naive Bayes

- Заметим, что даже это – сильно упрощённый взгляд: для слов ещё довольно-таки важен порядок, в котором они идут...
- Но и это ещё не всё: получается, что $p(a_1, a_2, \dots, a_n | x = v)$ – это вероятность *в точности такого набора слов* в сообщениях на разные темы. Очевидно, такой статистики взять неоткуда.
- Значит, надо дальше делать упрощающие предположения.
- Наивный байесовский классификатор – самая простая такая модель: давайте предположим, что все слова в словаре условно независимы при условии данной категории.

Naive Bayes

- Иначе говоря:

$$p(a_1, a_2, \dots, a_n | x = v) = p(a_1 | x = v) p(a_2 | x = v) \dots p(a_n | x = v).$$

- Итак, наивный байесовский классификатор выбирает v как

$$v_{NB}(a_1, a_2, \dots, a_n) = \arg \max_{v \in V} p(x = v) \prod_{i=1}^n p(a_i | x = v).$$

- В парадигме классификации текстов мы предполагаем, что разные слова в тексте на одну и ту же тему появляются независимо друг от друга. Однако, несмотря на такие бредовые предположения, naive Bayes на практике работает очень даже неплохо (и этому есть разумные объяснения).

Многомерная модель

- В деталях реализации наивного байесовского классификатора прячется небольшой дьяволёнок.
- Сейчас мы рассмотрим два разных подхода к naive Bayes, которые дают разные результаты: мультиномиальный (multinomial) и многомерный (multivariate).

Многомерная модель

- В многомерной модели документ – это вектор бинарных атрибутов, показывающих, встретилось ли в документе то или иное слово.
- Когда мы подсчитываем правдоподобие документа, мы перемножаем вероятности того, что встретилось каждое слово из документа и вероятности того, что не встретилось каждое (словарное) слово, которое не встретилось.
- Получается модель многомерных испытаний Бернулли. Наивное предположение в том, что события «встретилось ли слово» предполагаются независимыми.

Многомерная модель

- Математически: пусть $V = \{w_t\}_{t=1}^{|V|}$ – словарь. Тогда документ d_i – это вектор длины $|V|$, состоящий из битов B_{it} ; $B_{it} = 1$ iff слово w_t встречается в документе d_i .
- Правдоподобие принадлежности d_i классу c_j :

$$p(d_i | c_j) = \prod_{t=1}^{|V|} (B_{it}p(w_t | c_j) + (1 - B_{it})(1 - p(w_t | c_j))).$$

- Для обучения такого классификатора нужно обучить вероятности $p(w_t | c_j)$.

Многомерная модель

- Обучение – дело нехитрое: пусть дан набор документов $D = \{d_i\}_{i=1}^{|D|}$, которые уже распределены по классам c_j (возможно, даже вероятностно распределены), дан словарь $V = \{w_t\}_{t=1}^{|V|}$, и мы знаем биты B_{it} (знаем документы).
- Тогда можно подсчитать оптимальные оценки вероятностей того, что то или иное слово встречается в том или ином классе (при помощи лапласовой оценки):

$$p(w_t | c_j) = \frac{1 + \sum_{i=1}^{|D|} B_{it} p(c_j | d_i)}{2 + \sum_{i=1}^{|D|} p(c_j | d_i)}.$$

Многомерная модель

- Априорные вероятности классов можно подсчитать как $p(c_j) = \frac{1}{|D|} \sum_{i=1}^{|D|} p(c_j | d_i)$.
- Тогда классификация будет происходить как

$$\begin{aligned}
 c &= \arg \max_j p(c_j) p(d_i | c_j) = \\
 &= \arg \max_j \left(\frac{1}{|D|} \sum_{i=1}^{|D|} p(c_j | d_i) \right) \prod_{t=1}^{|V|} (B_{it} p(w_t | c_j) + (1 - B_{it})(1 - p(w_t | c_j))) = \\
 &= \arg \max_j \left(\log \left(\sum_{i=1}^{|D|} p(c_j | d_i) \right) + \sum_{t=1}^{|V|} \log (B_{it} p(w_t | c_j) + (1 - B_{it})(1 - p(w_t | c_j))) \right)
 \end{aligned}$$

Мультиномиальная модель

- В мультиномиальной модели документ – это последовательность событий. Каждое событие – это случайный выбор одного слова из того самого «bag of words».
- Когда мы подсчитываем правдоподобие документа, мы перемножаем вероятности того, что мы достали из мешка те самые слова, которые встретились в документе. Наивное предположение в том, что мы достаём из мешка разные слова независимо друг от друга.
- Получается мультиномиальная генеративная модель, которая учитывает количество повторений каждого слова, но не учитывает, каких слов *нет* в документе.

Мультиномиальная модель

- Математически: пусть $V = \{w_t\}_{t=1}^{|V|}$ – словарь. Тогда документ d_i – это вектор длины $|d_i|$, состоящий из слов, каждое из которых «вынуто» из словаря с вероятностью $p(w_t | c_j)$.
- Правдоподобие принадлежности d_i классу c_j :

$$p(d_i | c_j) = p(|d_i|) |d_i|! \prod_{t=1}^{|V|} \frac{1}{N_{it}!} p(w_t | c_j)^{N_{it}},$$

где N_{it} – количество вхождений w_t в d_i .

- Для обучения такого классификатора тоже нужно обучить вероятности $p(w_t | c_j)$.

Мультиномиальная модель

- Обучение: пусть дан набор документов $D = \{d_i\}_{i=1}^{|D|}$, которые уже распределены по классам c_j (возможно, даже вероятностно распределены), дан словарь $V = \{w_t\}_{t=1}^{|V|}$, и мы знаем вхождения N_{it} .
- Тогда можно подсчитать апостериорные оценки вероятностей того, что то или иное слово встречается в том или ином классе (не забываем сглаживание – правило Лапласа):

$$p(w_t | c_j) = \frac{1 + \sum_{i=1}^{|D|} N_{it} p(c_j | d_i)}{|V| + \sum_{s=1}^{|V|} \sum_{i=1}^{|D|} N_{is} p(c_j | d_i)}.$$

Мультиномиальная модель

- Априорные вероятности классов можно подсчитать как $p(c_j) = \frac{1}{|D|} \sum_{i=1}^{|D|} p(c_j | d_i)$.
- Тогда классификация будет происходить как

$$\begin{aligned}
 c &= \arg \max_j p(c_j) p(d_i | c_j) = \\
 &= \arg \max_j \left(\frac{1}{|D|} \sum_{i=1}^{|D|} p(c_j | d_i) \right) p(|d_i|) |d_i|! \prod_{t=1}^{|V|} \frac{1}{N_{it}!} p(w_t | c_j)^{N_{it}} = \\
 &= \arg \max_j \left(\log \left(\sum_{i=1}^{|D|} p(c_j | d_i) \right) + \sum_{t=1}^{|V|} N_{it} \log p(w_t | c_j) \right).
 \end{aligned}$$

LDA

- Более сложная модель – LDA (Latent Dirichlet Allocation).
- Задача: смоделировать большую коллекцию текстов (например, для information retrieval или классификации).
- Мы знаем наивный подход: скрытая переменная – тема, слова получаются из темы независимо по дискретному распределению.
- Аналогично работают и подходы, основанные на кластеризации.
- Давайте чуть усложним.

LDA

- Очевидно, что у одного документа может быть несколько тем; подходы, которые кластеризуют документы по темам, никак этого не учитывают.
- 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.

LDA

- 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 w . Нам дан корпус из M документов $\mathcal{D} = \{w_d \mid d = 1..M\}$.
- Генеративная модель LDA выглядит так.
 1. Выбрать $N \sim p(N \mid \xi)$.
 2. Выбрать $\theta \sim \text{Di}(\alpha)$.
 3. Для каждого из N слов w_n :
 1. выбрать тему $z_n \sim \text{Mult}(\theta)$;
 2. выбрать слово $w_n \sim p(w_n \mid z_n, \beta)$ по мультиномиальному распределению.

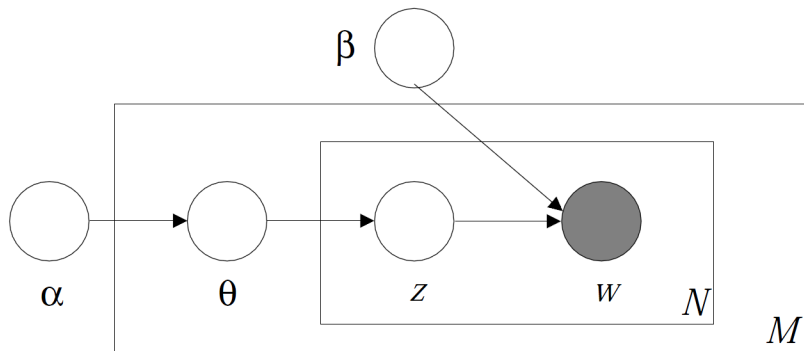
LDA

- Мы пока для простоты фиксируем число тем k , считаем, что β – это просто набор параметров $\beta_{ij} = p(w^j = 1 | z^i = 1)$, которые нужно оценить, и не беспокоимся о распределении на N .
- Совместное распределение тогда выглядит так:

$$p(\theta, \mathbf{z}, \mathbf{w}, N | \alpha, \beta) = p(N | \xi) p(\theta | \alpha) \prod_{n=1}^N p(z_n | \theta) p(w_n | z_n, \beta).$$

- В отличие от обычной кластеризации с априорным распределением Дирихле, мы тут не выбираем кластер один раз, а затем накидываем слова из этого кластера, а для каждого слова выбираем по распределению θ , по какой теме оно будет набросано.

LDA: графическая модель



Вывод в LDA

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

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

- Правдоподобие набора слов \mathbf{w} оценивается как

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

и это трудно посчитать, потому что θ и β путаются друг с другом.

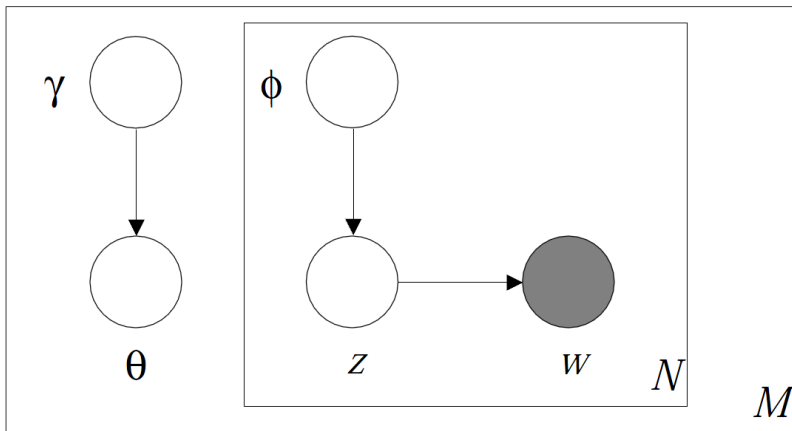
Вывод в LDA

- 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 γ (Дирихле) и ϕ (мультиномиальный).
- Заметим, что параметры для каждого документа могут быть свои – всё условно по \mathbf{w} .

LDA: вариационное приближение



LDA: вариационный вывод

- Теперь можно искать минимум KL-расстояния:

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

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

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

LDA: вариационный вывод

- 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.

LDA: вариационный вывод

$$\begin{aligned}
\mathcal{L}(\gamma, \phi; \alpha, \beta) &= \log \Gamma\left(\sum_{i=1}^k \alpha_i\right) - \sum_{i=1}^k \log \Gamma(\alpha_i) + \sum_{i=1}^k (\alpha_i - 1) \left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^k \gamma_j\right) \right] + \\
&+ \sum_{n=1}^N \sum_{i=1}^k \phi_{ni} \left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^k \gamma_j\right) \right] + \\
&+ \sum_{n=1}^N \sum_{i=1}^k \sum_{j=1}^V w_n^j \phi_{ni} \log \beta_{ij} - \\
&- \log \Gamma\left(\sum_{i=1}^k \gamma_i\right) + \sum_{i=1}^k \log \Gamma(\gamma_i) - \sum_{i=1}^k (\gamma_i - 1) \left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^k \gamma_j\right) \right] - \\
&- \sum_{n=1}^N \sum_{i=1}^k \phi_{ni} \log \phi_{ni}.
\end{aligned}$$

LDA: вариационный вывод

- Теперь осталось только брать частные производные этого выражения.
- Сначала максимизируем его по ϕ_{ni} (вероятность того, что n -е слово было порождено темой i); надо добавить λ -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 $w_n^v = 1$.

LDA: вариационный вывод

- Потом максимизируем по γ_i , i -й компоненте апостериорного Дирихле-параметра.
- Получится

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

- Соответственно, для вывода нужно просто пересчитывать ϕ_{ni} и γ_i друг через друга, пока оценка не сойдётся.

LDA: оценка параметров

- Теперь давайте попробуем оценить параметры α и β по корпусу документов \mathcal{D} .
- Мы хотим найти α и β , 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.

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

LDA: оценка параметров

- 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 α и β .

LDA: оценка параметров

- Для β это тоже делается нехитро:

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

- Для α_i получается система уравнений, которую можно решить методом Ньютона.

LDA: сэмплирование по Гиббсу

- В базовой модели LDA сэмплирование по Гиббсу после несложных преобразований сводится к так называемому *сжатому 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} .

LDA: сэмплирование по Гиббсу

- 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 .

- Сэмплирование по Гиббсу обычно проще расширить на новые модификации LDA, но вариационный подход быстрее и часто стабильнее.

Варианты и расширения модели LDA

- 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.

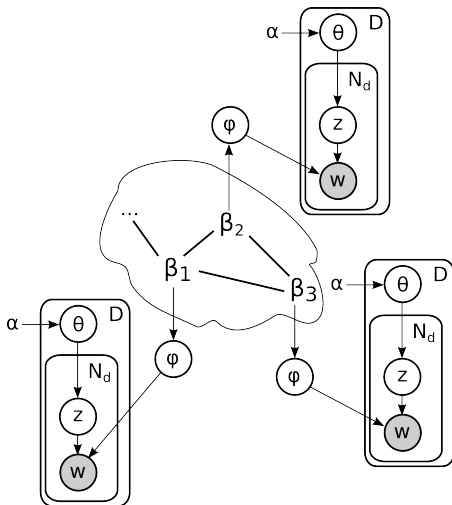
Коррелированные тематические модели

- 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

- 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 β_i в LDA, описывающих параметры разных корпусов с одними и теми же темами. Гиперпараметры β_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 β_i .
- В свою очередь, β_i подчиняются априорным ограничениям, которые позволяют «делить» темы между корпусами, задавать «фоновые» темы, присутствующие во всех корпусах, накладывать ограничения на

Марковские тематические модели



Реляционная тематическая модель

- 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 :

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

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

Модели, учитывающие время

- 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.

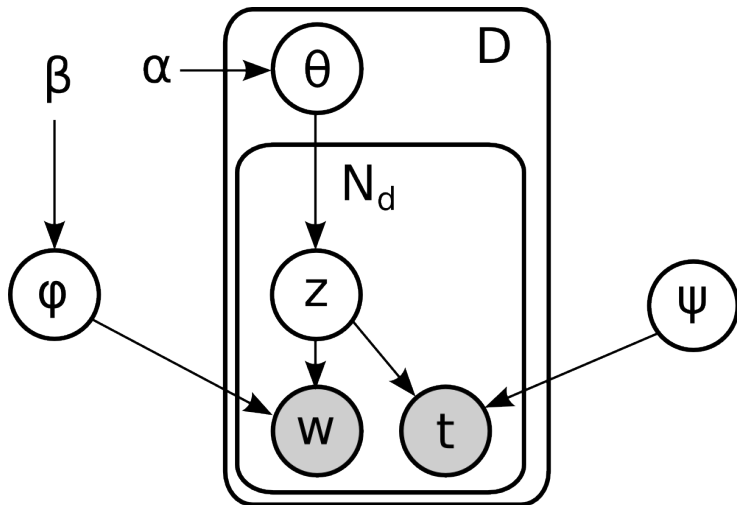
Topics over Time

- В модели ТОТ (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 ϕ_z из априорного распределения Дирихле β ;
 - для каждого документа d выбрать мультиномиальное распределение θ_d из априорного распределения Дирихле α , затем для каждого слова $w_{di} \in d$:
 - выбрать тему z_{di} из θ_d ;
 - выбрать слово w_{di} из распределения $\phi_{z_{di}}$;
 - выбрать время t_{di} из бета-распределения $\psi_{z_{di}}$.

Topics over Time

- Основная идея заключается в том, что каждой теме соответствует её бета-распределение ψ_z , т.е. каждая тема локализована во времени (сильнее или слабее, в зависимости от параметров ψ_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 ψ_z будет в первом случае меньше, чем во втором.

Topics over Time



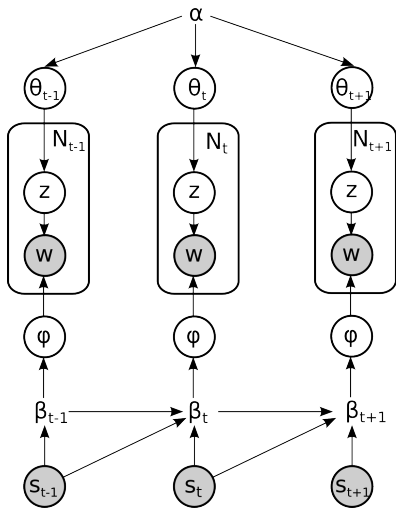
Динамические тематические модели

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

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

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

- В остальном генеративный процесс остаётся неизменным.

Непрерывная динамическая тематическая модель
(cDTM)

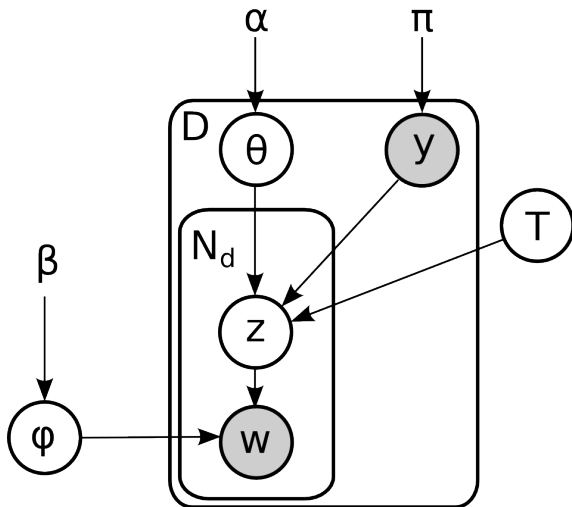
Supervised LDA

- 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 η и δ – другие параметры glm.
- К примеру, в контексте рекомендательных систем дополнительный отклик может быть реакцией пользователя.

DiscLDA

- 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 θ 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 θ , сгенерированному для данного документа,
 - сгенерировать тему z по распределению $T^y\theta$, где T^y – преобразование, соответствующее метке данного документа y .

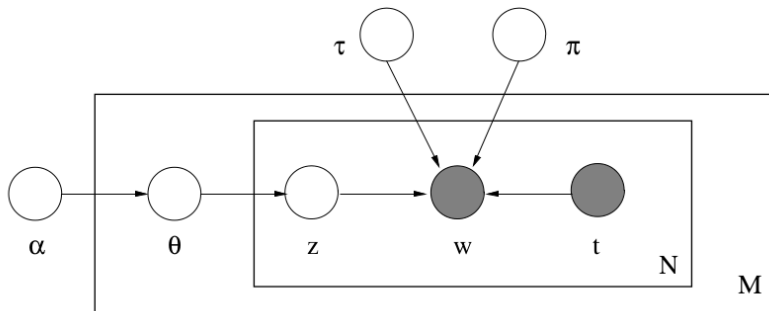
DiscLDA



TagLDA

- 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.

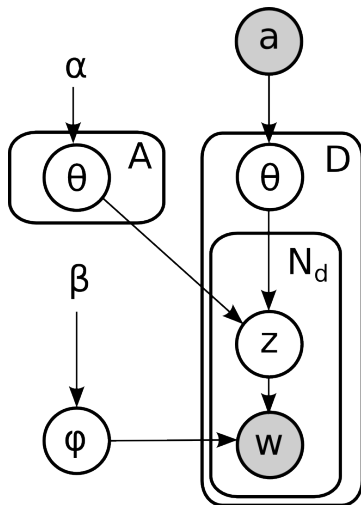
TagLDA



Author-Topic model

- 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 \mathbf{a}_d ;
 - выбираем тему из распределения на темах, соответствующего автору x ;
 - выбираем слово из распределения слов, соответствующего этой теме.

Author-Topic model



Author-Topic model

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

$$p(z_w = t, x_w = a \mid \mathbf{z}_{-w}, \mathbf{x}_{-w}, \mathbf{w}, \alpha, \beta) \propto \frac{n_{-a,t}^{(a)} + \alpha}{\sum_{t' \in T} (n_{-w,t'}^{(a)} + \alpha)} \frac{n_{-w,t}^{(w)} + \beta}{\sum_{w' \in W} (n_{-w,t}^{(w')} + \beta)},$$

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} .

Thank you!

Спасибо за внимание!