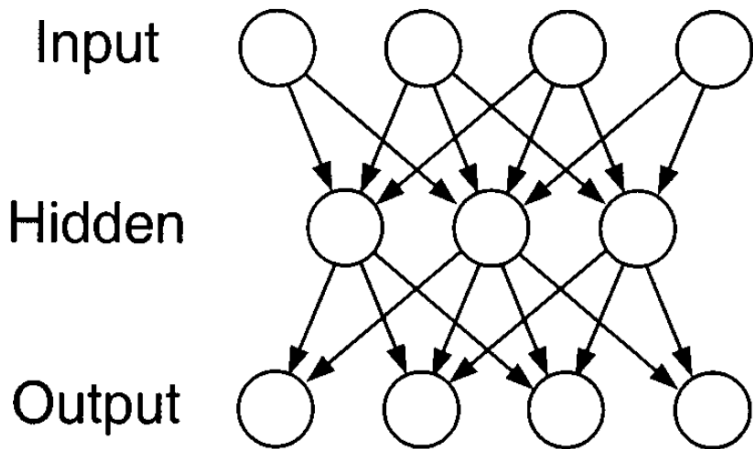


Нейронные сети и LDA

Сергей Николенко

Академический Университет, весенний семестр 2011

Outline



- В узле используем сигмоид-функцию (sigmoid belief networks):

$$p(S_i = 1 | S_{\text{pa}(i)}) = 1 / \left(1 + e^{-\sum_{j \in \text{pa}(i)} \theta_{ij} S_j - \theta_{i0}} \right).$$

- Это можно переписать в одну формулу для обоих значений S_i :

$$p(S_i | S_{\text{pa}(i)}) = \frac{e^{(\sum_{j \in \text{pa}(i)} \theta_{ij} S_j + \theta_{i0}) S_i}}{1 + \left(1 + e^{\sum_{j \in \text{pa}(i)} \theta_{ij} S_j + \theta_{i0}} \right)}, \text{ и}$$

$$p(S | \theta) = \prod_i \left[\frac{e^{(\sum_{j \in \text{pa}(i)} \theta_{ij} S_j + \theta_{i0}) S_i}}{1 + \left(1 + e^{\sum_{j \in \text{pa}(i)} \theta_{ij} S_j + \theta_{i0}} \right)} \right].$$

- Мы хотим подсчитать условные вероятности в этом совместном распределении.

- Мы снова применяем метод самосогласованного поля, факторизуем полностью:

$$q(H | E, \mu) = \prod_{i \in H} \mu_i^{S_i} (1 - \mu_i)^{1 - S_i}.$$

- Вычисляем КЛ точно так же, как в случае машин Больцмана, отличие только в $q \ln p$; в итоге получаем:

$$\begin{aligned} \ln p(E | \theta) \geq & \sum_{i < j} \theta_{ij} \mu_i \mu_j + \sum_i \theta_{i0}^c \mu_i - \\ & - \sum_i \left\langle \ln \left[1 + e^{\sum_{j \in \text{pa}(i)} \theta_{ij} S_j + \theta_{i0}} \right] \right\rangle - \\ & - \sum_i [\mu_i \ln \mu_i + (1 - \mu_i) \ln(1 - \mu_i)], \end{aligned}$$

где $\langle \cdot \rangle$ – ожидание по распределению q .

- Но есть проблема: среднее от $\ln(1 + e^{z_i})$, где $z_i = \sum_{j \in \text{pa}(i)} \theta_{ij} S_j + \theta_{i0}$, не вычислить точно.

- Для этого введём ещё одно вариационное приближение с параметрами ξ_j , из неравенства Йенсена:

$$\begin{aligned}\langle \ln(1 + e^{z_i}) \rangle &= \langle \ln \left[e^{\xi_i z_i} e^{-\xi_i z_i} (1 + e^{z_i}) \right] \rangle = \\ \xi_i \langle z_i \rangle + \langle \ln \left[e^{-\xi_i z_i} + e^{(1-\xi_i)z_i} \right] \rangle &\leq \xi_i \langle z_i \rangle + \ln \langle e^{-\xi_i z_i} + e^{(1-\xi_i)z_i} \rangle.\end{aligned}$$

- Можно доказать (мы не будем), что в пределе (при большом числе родителей у узла) ξ_j – это ожидание $\sigma(z_j) = 1/(1 + e^{-z_j})$.

- Для фиксированных ξ_i , дифференцируя KL по μ_i , получаем:

$$\mu_i = \sigma \left(\sum_j \theta_{ij} \mu_j + \theta_{i0} + \sum_j \theta_{ij} (\mu_j - \xi_j) + \sum_j K_{ji} \right),$$

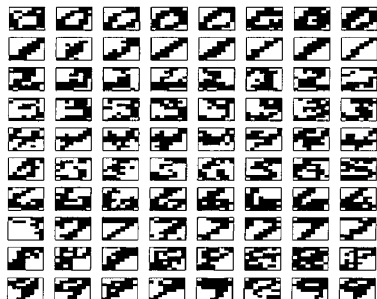
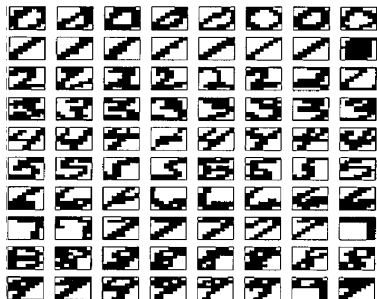
где $K_{ji} = \partial (-\ln \langle e^{-\xi_i z_i} + e^{(1-\xi_i) z_i} \rangle) / \partial \mu_i$.

Упражнение. Докажите это.

- Эту производную можно подсчитать из весов узла i , его детей j и других родителей его детей.
- Можно получить уравнения для апдейта ξ_i : [Saul, Jaakkola, Jordan, 1996], [Saul, Jordan, 1999].
- И, наконец, можно получить уравнения для апдейта θ_{ij} при обучении:

$$\Delta \theta_{ij} \propto (\mu_i - \xi_i) \mu_j - \theta_{ij} \xi_i (1 - \xi_i) \mu_i (1 - \mu_j).$$

Пример: распознавание цифр



- Ещё один конкретный пример применения – 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.

- 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.

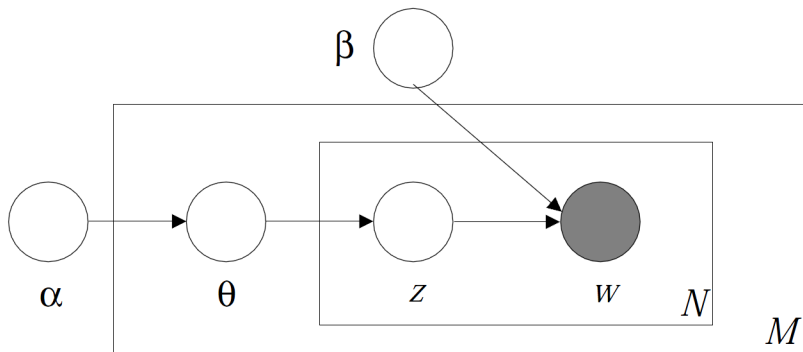
- Если формально: слова берутся из словаря $\{1, \dots, V\}$; слово – это вектор w , $w_i \in \{0, 1\}$, где ровно одна компонента равна 1.
- Документ – последовательность из N слов w . Нам дан корпус из M документов $\mathcal{D} = \{w_d \mid d = 1..M\}$.
- Генеративная модель LDA выглядит так.
 1. Выбрать $N \sim p(N \mid \xi)$.
 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.

- 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



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

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

- Мы здесь рассматриваем α и β как известные константы; по идее, конечно, их тоже нужно оценивать, но пока для простоты так.
- Знаменатель – правдоподобие – оценивается как

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

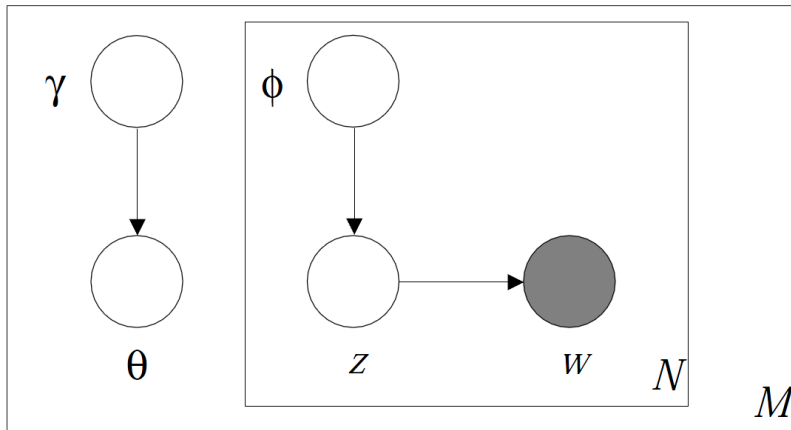
и это трудно посчитать, потому что θ и β путаются друг с другом.

- 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



- 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.

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

- 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 $w_n^v = 1$.

- 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.

- 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.

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

- 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 .

- 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.

Thank you!

Спасибо за внимание!