

Латентная семантическая модель для представления СМЫСЛОВ МНОГОЗНАЧНЫХ СЛОВ

Дмитрий Кондрашкин,
научный руководитель: к.ф.-м.н. Ветров Д. П.

26 февраля 2015 г.

Skip-gram model

- ▶ По слову w предсказывается слово из контекста v :

$$p(v | w, \theta) = \frac{\exp\{\text{In}_w^T \text{Out}_v\}}{\sum_{t=1}^V \exp\{\text{In}_w^T \text{Out}_t\}},$$

где $\theta = \{\text{In}_v, \text{Out}_v\}_{v=1}^V$.

- ▶ Каждому слову v соответствуют “входные” и “выходные” представления $\text{In}_v, \text{Out}_v \in \mathbb{R}^D$.

(Mikolov 2013) Distributed representations of words and phrases and their compositionality.

Skip-gram model

- ▶ Входной текст $\mathbf{o} = \{o_1, \dots, o_N\}$.
- ▶ Обучающий объект (x_i, \mathbf{y}_i) :
 - ▶ $x_i = o_i$,
 - ▶ $\mathbf{y}_i = \{o_{i-C/2}, \dots, o_{i-1}, o_{i+1}, \dots, o_{i+C/2}\}$.
- ▶ Правдоподобие:

$$\prod_{i=1}^N \prod_{j=1}^C p(y_{ij} | x_i, \boldsymbol{\theta}).$$

- ▶ Максимизировать будем лог-правдоподобие

$$L(X, Y, \boldsymbol{\theta}) = \sum_{i=1}^N \sum_{j=1}^C \log p(y_{ij} | x_i, \boldsymbol{\theta}).$$

Стохастическая оптимизация

- ▶ Типичная задача в машинном обучении

$$L(X, Y, \boldsymbol{\theta}) = \sum_{i=1}^N l(x_i, y_i, \boldsymbol{\theta}) + \lambda R(\boldsymbol{\theta}) \rightarrow \min_{\boldsymbol{\theta}}$$

- ▶ Стохастический градиентный спуск:

$$\begin{aligned}\boldsymbol{\theta}^{(t+1)} &= \boldsymbol{\theta}^{(t)} - \eta_t g(\boldsymbol{\theta}^{(t)}), \\ \nabla L(X, Y, \boldsymbol{\theta}) &\approx g(\boldsymbol{\theta}).\end{aligned}$$

- ▶ Условия сходимости:

$$\begin{aligned}\mathbb{E} g(\boldsymbol{\theta}) &= \nabla L(X, Y, \boldsymbol{\theta}), \\ \sum_t \eta_t &= \infty, \quad \sum_t \eta_t^2 < \infty.\end{aligned}$$

- ▶ Как правило, $g(\boldsymbol{\theta}) = N \nabla l(x_i, y_i, \boldsymbol{\theta}) + \lambda \nabla R(\boldsymbol{\theta})$.

Стохастическая оптимизация

- ▶ В нашем случае:

$$L(X, Y, \boldsymbol{\theta}) = \sum_{i=1}^N \underbrace{\sum_{j=1}^C \log p(y_{ij} | x_i, \boldsymbol{\theta})}_{l(x_i, \mathbf{y}_i, \boldsymbol{\theta})} \rightarrow \max_{\boldsymbol{\theta}}.$$

- ▶ $\nabla L(X, Y, \boldsymbol{\theta}) \approx N \sum_{j=1}^C \nabla \log p(y_{ij} | x_i, \boldsymbol{\theta})$.

Hierarchical softmax

$$p(v | w, \theta) = \frac{\exp\{\text{In}_w^\top \text{Out}_v\}}{\sum_{t=1}^V \exp\{\text{In}_w^\top \text{Out}_t\}}$$

- ▶ Размер словаря $V \approx 10^5$, линейная сложность для подсчета функции и градиента — очень долго.
- ▶ Нужна функция, такая что
 - ▶ $p(v | w, \theta) > 0$ и $\sum_{v=1}^V p(v | w, \theta) = 1$,
 - ▶ Считалась быстрее, чем за $O(V)$.

Hierarchical softmax

$$p(v | w, \theta) = \prod_{n \in \text{path}(v)} \sigma(\text{ch}(n) \text{In}_w^T \text{Out}_n).$$

- ▶ Бинарное дерево: каждому листу соответствует слово из словаря.
- ▶ Теперь “выходные” представления соответствуют не словам, а внутренним вершинам в дереве.
- ▶ $\text{path}(v)$ — номера вершин на пути из корня в лист, соответствующий слову v .
- ▶ $\text{ch}(n)$ — $+1$ или -1 в зависимости от того, что следующая вершина на пути — это правый или левый сын n .
- ▶ Используем, что $\sigma(x) + \sigma(-x) = 1$.

Linguistic regularities

x	$\operatorname{argmax}_w \cos(w, x)$
Berlin-Germany+Russia	Moscow
Obama-USA+Russia	Putin
king-man+woman	queen

Многозначные слова

- ▶ Проблемы:
 - ▶ Смешивание смыслов.
 - ▶ Доминирование наиболее частотного смысла.
- ▶ Ближайшие соседи по косинусному расстоянию:
 - ▶ для слова apple: macintosh, iigs, ipad, ibook;
 - ▶ для слова python: perl, php, **molurus**, c++, ..., **monty**;
- ▶ Как учесть то, что слова могут иметь больше одного смысла?

Наивная многосмысловая модель

- ▶ Введем скрытую переменную z — номер смысла, тогда:

$$p(v \mid z = k, w, \boldsymbol{\theta}) = \frac{\exp\{\text{In}_{w,k}^T \text{Out}_v\}}{\sum_{t=1}^V \exp\{\text{In}_{w,k}^T \text{Out}_t\}},$$

- ▶ Теперь каждому слову соответствуют K “входных” векторов.
- ▶ Наблюдаемые данные (x_i, y_i) , скрытые переменные z_i (неполная разметка!).
- ▶ Полное правдоподобие:

$$p(Y, Z \mid X, \boldsymbol{\theta}) = \prod_{i=1}^N p(z_i) \prod_{j=1}^C p(y_{ij} \mid z_i, x_i, \boldsymbol{\theta}).$$

Обучение

- ▶ Будем максимизировать логарифм неполного правдоподобия:

$$\log p(Y | X, \theta) = \log \sum_Z p(Y, Z | X, \theta) \rightarrow \max_{\theta}.$$

- ▶ Его можно представить в виде:

$$\log p(Y | X, \theta) = \mathcal{L}(q(Z), \theta) + \text{KL}(q(Z) \| p(Z | X, Y, \theta)),$$

где \mathcal{L} вариационная нижняя оценка:

$$\mathcal{L}(q(Z), \theta) = \mathbb{E}_{q(Z)} \log \left[\frac{p(Y, Z | X, \theta)}{q(Z)} \right].$$

- ▶ Нижняя, потому что $\text{KL}(q \| p) \geq 0$.

Обучение

- ▶ Перейдем к задаче максимизации нижней оценки:

$$\mathcal{L}(q(Z), \boldsymbol{\theta}) \rightarrow \max_{q(Z), \boldsymbol{\theta}}.$$

- ▶ Будем искать $q(Z)$ в виде $\prod_{i=1}^N q(z_i)$.
- ▶ Перепишем нижнюю оценку:

$$\mathcal{L}(q(Z), \boldsymbol{\theta}) = \sum_{i=1}^N \mathbb{E}_{q(z_i)} \left(\log p(z_i) + \sum_{j=1}^C \log p(y_{ij} | z_i, x_i, \boldsymbol{\theta}) \right) - \sum_{i=1}^N \mathbb{E}_{q(z_i)} \log q(z_i).$$

EM-алгоритм

- ▶ E-шаг:

$$q(z_i = k) = p(z_i = k | \mathbf{y}_i, x_i, \boldsymbol{\theta}^{old}) \propto \exp \left\{ \sum_{j=1}^C \log p(y_{ij} | k, x_i, \boldsymbol{\theta}^{old}) \right\}$$

для $k = 1, \dots, K$ и $i = 1, \dots, N$.

- ▶ M-шаг:

$$\mathcal{L}(q(Z), \boldsymbol{\theta}) \rightarrow \max_{\boldsymbol{\theta}}$$

или

$$\sum_{i=1}^N \sum_{j=1}^C \mathbb{E}_{q(z_i)} \log p(y_{ij} | z_i, x_i, \boldsymbol{\theta}) \rightarrow \max_{\boldsymbol{\theta}}.$$

EM-алгоритм

- ▶ Как найти $\mathbb{E}_{q(z_i)} \log p(y_{ij} | z_i, x_i, \theta)$?
- ▶ Стандартный прием:

$$p(y_{ij} | z_i, x_i, \theta) = \prod_{k=1}^K p(y_{ij} | k, x_i, \theta)^{\mathbb{I}[z_i=k]}.$$

- ▶ В результате:

$$\begin{aligned} \mathbb{E}_{q(z_i)} \log p(y_{ij} | z_i, x_i, \theta) &= \sum_{k=1}^K \mathbb{E}_{q(z_i)} \mathbb{I}[z_i = k] \log p(y_{ij} | k, x_i, \theta) = \\ &= \sum_{k=1}^K q(z_i = k) \log p(y_{ij} | k, x_i, \theta). \end{aligned}$$

Стохастический EM-алгоритм

- ▶ Для i -го обучающего объекта (x_i, y_i) выполняем E-шаг для нахождения распределения $q(z_i)$.
- ▶ При выполнении M-шага делаем шаг по стохастическому градиенту нижней оценки:

$$\nabla \mathcal{L}(q(Z), \theta) \approx N \sum_{j=1}^C \sum_{k=1}^K q(z_i = k) \nabla \log p(y_{ij} | k, x_i, \theta).$$

- ▶ Сравним с градиентом Skip-gram:

$$\nabla L(X, Y, \theta) \approx N \sum_{j=1}^C \nabla \log p(y_{ij} | x_i, \theta).$$

Проблемы наивного подхода

- ▶ Предположение о том, что все слова имеют K смыслов нереалистично.
- ▶ Как автоматически подбирать число смыслов для каждого слова?
- ▶ Можно использовать процесс Дирихле.
- ▶ См. нашу работу:

Breaking Sticks and Ambiguities with Adaptive Skip-gram.

Bartunov S., Kondrashkin D., Osokin A., Vetrov D.

<http://arxiv.org/abs/1502.07257>

Примеры найденных смыслов

Смысл 1	Смысл 2
almond	macintosh
cherry	iifx
plum	iigs
apricot	computers
orange	kaupro

Для слова apple.

Смысл 1	Смысл 2	Смысл 3
monty	perl	molurus
spamalot	php	pythons
cantsin	java	peafowl
zirkus	c++	tortoise
circus	objective-c	snake

Для слова python.

Разрешение лексической многозначности

Определение смысла слова по контексту:

$$\operatorname{argmax}_k p(z = k \mid \mathbf{y}, x, \boldsymbol{\theta}).$$

Примеры для слова apple

- ▶ $p(z \mid (\text{tasty, sweet}), \text{apple}) = [0.99998, 0.00002]$,
- ▶ $p(z \mid (\text{announce, today}), \text{apple}) = [0.121476, 0.878524]$.

Примеры для слова python

- ▶ $p(z \mid (\text{interesting, show}), \text{python}) = [0.950317, 0.0327315, 0.0169517]$,
- ▶ $p(z \mid (\text{code}), \text{python}) = [0.0219413, 0.976705, 0.00135406]$,
- ▶ $p(z \mid (\text{dangerous, animal}), \text{python}) = [0.262747, 0.00012, 0.737133]$.

Основные идеи

- ▶ Модель Skig-gram и иерархический soft-max.
- ▶ Стохастическая оптимизация.
- ▶ Введение скрытых переменных.
- ▶ Непараметрический байес.

Библиография

- ▶ (Mikolov 2013) Distributed representations of words and phrases and their compositionality. Mikolov T., Sutskever I., Chen K., Corrado G., Dean J. Advances in Neural Information Processing Systems.
- ▶ (Hoffman 2013) Stochastic variational inference. Hoffman M., Blei D., Wang C., Paisley J. The Journal of Machine Learning Research.