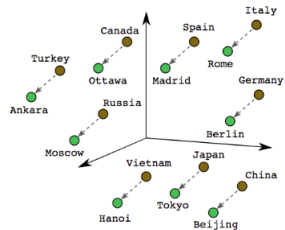
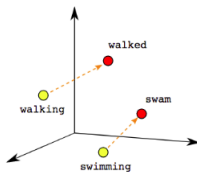
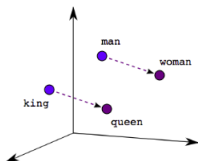


Векторные представления объектов

Виктор Китов

v.v.kitov@yandex.ru



Содержание

- 1 Векторное представление слов
- 2 Word2vec
- 3 Регулярности в пространстве представлений
- 4 Настройка skip-gram
- 5 Методы на основе матрицы встречаемости
- 6 Представления параграфов
- 7 Сиамская сеть

Стандартное представление слов

- Обозначим V = размер словаря.
- Стандартные представления слов $x \in \mathbb{R}^V$:
 - $x_w = \mathbb{I}[w \text{ встретился в документе}]$
 - $x_w = TF_w = \#[w \text{ встретился в документе}]$
 - $x_w = TF_w IDF_w, IDF_w = \frac{N}{N_w}$
 - N - # документов
 - N_w - # документов, содержащих w хотя бы раз.
- V велико, поэтому нужно компактное представление (word embedding) $x \in \mathbb{R}^K, K \ll V$:
 - меньше входов => меньше параметров => ниже переобучение
 - возможность учитывать семантическое сходство/различие
 - например, синонимы "автомобиль" и "машина"

Интерпретируемые векторные представления слов

- Можно из слов извлекать интерпретируемые признаки:
 - x^1 : часть речи
 - x^2 : род (м/ж/ср - для существительных)
 - x^3 : время (пр/наст/буд - для глаголов)
 - x^4 : \mathbb{I} [начинается с заглавной буквы]
 - x^5 : $\#$ букв
 - x^6 : категория: машинное обучение, физика, биология, ...
 - x^7 : подкатегория: обучение с учителем, без учителя, частичное обучение, ...
 - ...
- Необходимо придумывать признаки под задачу, производить разметку.
- Легче работать с неинтерпретируемыми признаками, но которые извлекаются автоматически.

Неинтерпретируемые представления слов

- Хотим, чтобы семантически близким словам соответствовали близкие представления.
- Дистрибутивная гипотеза (distributional hypothesis): слова близки по смыслу \Leftrightarrow они часто встречаются совместно
- "точность бустинга", "бустинг дал точность", "ниже точность, по сравнению с бустингом"
 - "точность" и "бустинг" связаны!
- Типичная размерность векторного представления $\in [300, 500]$.

Представления фраз

Можно обрабатывать фразы как отдельные "слова".

- Коллокации (неслучайно часто встречающиеся слова):

$$(w_i, w_j)\text{-коллокация} \iff \frac{p(w_i w_j) - \delta}{p(w_i)p(w_j)}$$

Представления фраз

Можно обрабатывать фразы как отдельные "слова".

- Коллокации (неслучайно часто встречающиеся слова):

$$(w_i, w_j)\text{-коллокация} \iff \frac{p(w_i w_j) - \delta}{p(w_i)p(w_j)}$$

$> threshold$. δ - параметр, снижающий значимость редко встречающихся слов.

Содержание

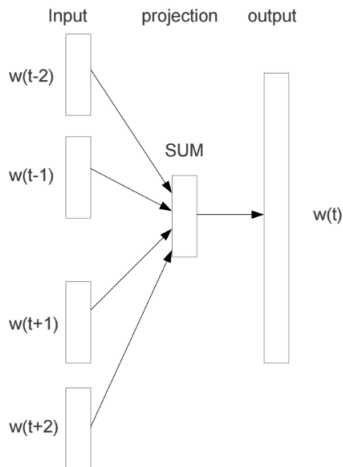
- 1 Векторное представление слов
- 2 Word2vec
- 3 Регулярности в пространстве представлений
- 4 Настройка skip-gram
- 5 Методы на основе матрицы совстречаемости
- 6 Представления параграфов
- 7 Сиамская сеть

Word2vec

- Для каждого w оценим:
 - целевое представление слова v_w
 - контекстное представление слова \tilde{v}_w
 - впоследствии можно не использовать, усреднить или конкатенировать с целевым представлением

CBOW: идея

Continuous bag of words (CBOW): предсказываем центральное слово по контексту.



CBOW: модель

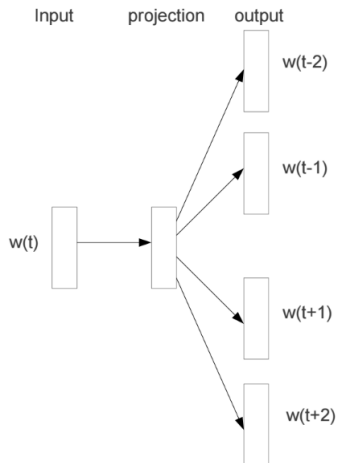
$$\frac{1}{T} \sum_{t=1}^T \ln p(w_t | w_{t-c}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+c}) \rightarrow \max_{\theta}$$

где $\tilde{v}_{context} = \sum_{-c \leq i \leq c, i \neq 0} \tilde{v}_{w_{t+i}}$ и

$$p(w_t | w_{t-c}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+c}) = \frac{\exp(\tilde{v}_{context}^T v_{w_t})}{\sum_{w=1}^V \exp(\tilde{v}_{context}^T v_w)}$$

Skip-gram: идея

Skip-gram: предсказываем контекст по центральному слову:



Skip-gram: модель

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq i \leq c, i \neq 0} \ln p(w_{t+i} | w_t) \rightarrow \max_{\theta}$$

$$p(w_{t+i} | w_t) = \frac{\exp(\tilde{v}_{w_t}^T v_{w_{t+i}})}{\sum_{w=1}^V \exp(\tilde{v}_{w_t}^T v_w)}$$

Проблема: знаменатель вычисляется за $O(V)$.

Комментарии

- Можем извлекать представления для др. объектов из последовательностей.
 - символы, биграммы, триграммы символов (см. *FastText*), предложения
 - нуклеотиды в ДНК последовательности
 - сервисы, заказанные клиентом компании
- Можем использовать ансамбли представлений
 - сумма, среднее, конкатенация

Содержание

- 1 Векторное представление слов
- 2 Word2vec
- 3 Регулярности в пространстве представлений**
- 4 Настройка skip-gram
- 5 Методы на основе матрицы встречаемости
- 6 Представления параграфов
- 7 Сиамская сеть

Похожие слова по представлению¹

- Ближайшие соседи слова пространстве эмбедингов - слова, похожие по смыслу (корпус GoogleNews, cosine-sim):
 - student -> teacher, faculty, school, university
 - car -> truck, jeep, vehicle
 - country -> nation, continent, region

¹http://epsilon-it.utu.fi/wv_demo/

Формы слов

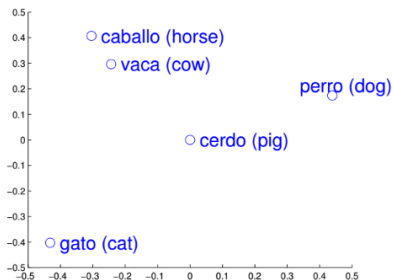
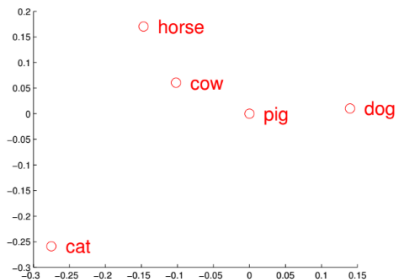
Одинаковые слова в разных формах образуют похожие структуры:



Представления могут помочь строить др. формы новых и редких слов.

Слова на разных языках

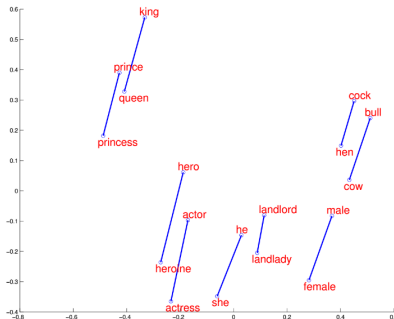
Слова на разных языках группируются похожим образом:



Представления слов могут помочь в переводе на др. язык.

Семантическая регулярность

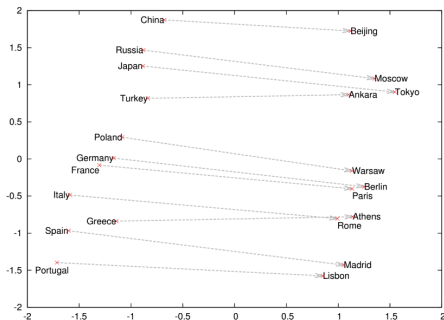
Слова, связанные семантически определенным образом группируются единообразно:



(prince-princess)+queen \approx king. Может помочь в системе автоматических ответов на вопросы.

Семантическая регулярность

Слова, связанные семантически определенным образом группируются единообразно:



(Beijing-China)+Russia \approx Moscow! Может помочь в системе автоматических ответов на вопросы.

Содержание

- 1 Векторное представление слов
- 2 Word2vec
- 3 Регулярности в пространстве представлений
- 4 Настройка skip-gram**
- 5 Методы на основе матрицы совстречаемости
- 6 Представления параграфов
- 7 Сиамская сеть

Skip-gram: модель

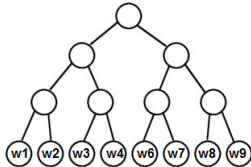
$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq i \leq c, i \neq 0} \ln p(w_{t+i} | w_t) \rightarrow \max_{\theta}$$

$$p(w_{t+i} | w_t) = \frac{\exp(\tilde{v}_{w_t}^T v_{w_{t+i}})}{\sum_{w=1}^V \exp(\tilde{v}_{w_t}^T v_w)}$$

Проблема: знаменатель вычисляется за $O(V)$.

Иерархический SoftMax

- Для каждого слова контекста w_t :
 - строим бинарное дерево, все слова языка - листья.
 - для каждого узла j определяем вероятности пойти вправо $\sigma(\tilde{v}_j^T v_{w_t})$ и влево $\sigma(-\tilde{v}_j^T v_{w_t})$
- $p(w|w_t)$ = произведение вероятностей дойти до него.
 - сложность вычисления $O(V) \rightarrow O(\log_2 V)$



- Эффективнее работает не сбалансированное дерево, а дерево Хаффмана.
 - более частотным словам - более короткие пути.
- Теперь контекст - набор $\{\tilde{v}_j\}_j$ для всех вершин дерева.

Негативное сэмплирование

- Негативное сэмплирование (negative sampling)² - аппроксимация максимизация правдоподобия.
- Для каждой реальной (позитивной) пары (w_t, w_{t+i}) сэмплируем K негативных случайно $(w_t, w_{j(1)}), \dots (w_t, w_{j(K)})$.

$$\underbrace{\ln \left(\frac{1}{1 + e^{-\tilde{v}_{w_t}^T v_{w_{t+i}}}} \right)}_{\sigma(+\tilde{v}_{w_t}^T v_{w_{t+i}})} + \sum_{k=1}^K \underbrace{\ln \left(\frac{1}{1 + e^{+\tilde{v}_{w_t}^T v_{w_{t+i}}}} \right)}_{\sigma(-\tilde{v}_{w_t}^T v_{w_{t+i}})} \rightarrow \max_{\tilde{v}, v}$$

- $K \sim 2-5$. $p(w_{j(k)}) \propto p(w)^{3/4}$.

²Distributed Representations of Words and Phrases

fastText³

- Работает как skip-gram (предсказываем слова контекста по центральному слову)
- Раньше совместимость была $\tilde{w}_t^T w_{t+i}$
- В fastText $\tilde{w}'_t := \tilde{w}_t + \sum_{p \in n\text{-grams}(w_t)} \text{part}_p$, w_{t+i} - прежняя.
- Новая совместимость

$$\tilde{w}_t^T w_{t+i} + \sum_{p \in n\text{-grams}(w_t)} \text{part}_p^T w_{t+i}$$

- Пример триграмм: "person" -> "<<person>", "<pe", "per", "ers", "rso", "son", "on>"
 - предлагается использовать все n-граммы, $3 \leq n \leq 6$.
- Работает для слов вне словаря, для слов с опечатками за счет $\sum_{p \in n\text{-grams}(w)} \text{part}_p$.

³<https://arxiv.org/pdf/1607.04606.pdf>, код и данные: fasttext.cc

Содержание

- 1 Векторное представление слов
- 2 Word2vec
- 3 Регулярности в пространстве представлений
- 4 Настройка skip-gram
- 5 Методы на основе матрицы совстречаемости**
- 6 Представления параграфов
- 7 Сиамская сеть

Матрица совстречаемости слов

- $X \in \mathbb{R}^{V \times V}$ - матрица со-встречаемости слов (word co-occurrence matrix)
- $X_{ij} = \#\{\text{слово } j \text{ встретилось в контексте слова } i\}$.
- Пример для контекста ± 1 слово:

I like deep learning. I like NLP. I enjoy flying.

counts	I	like	enjoy	deep	learning	NLP	flying	.
I	0	2	1	0	0	0	0	0
like	2	0	0	1	0	1	0	0
enjoy	1	0	0	0	0	0	1	0
deep	0	1	0	0	1	0	0	0
learning	0	0	0	1	0	0	0	1
NLP	0	1	0	0	0	0	0	1
flying	0	0	1	0	0	0	0	1
.	0	0	0	0	1	1	1	0

Разложение матрицы совстречаемости

- Hyperspace Analogue to Language (HAL)⁴: эмбединги из низкорангового разложения
 - напр. строки U либо столбы V^T из SVD.
 - эмбединги доминируются частыми словами!

⁴Lund and Burgess, 1996.

⁵Bullinaria and Levy, 2007

⁶<https://aclanthology.org/D14-1162.pdf>

Разложение матрицы совстречаемости

- Hyperspace Analogue to Language (HAL)⁴: эмбединги из низкорангового разложения
 - напр. строки U либо столбы V^T из SVD.
 - эмбединги доминируются частыми словами!
- Модификация⁵: счётчик совстречаемости \rightarrow PPMI

$$PPMI(w_1, w_2) = \max\{0, PMI(w_1, w_2)\} = \max\{0, \ln \frac{P(w_1, w_2)}{P(w_1)P(w_2)}\}$$

⁴Lund and Burgess, 1996.

⁵Bullinaria and Levy, 2007

⁶<https://aclanthology.org/D14-1162.pdf>

Разложение матрицы совстречаемости

- Hyperspace Analogue to Language (HAL)⁴: эмбединги из низкорангового разложения
 - напр. строки U либо столбы V^T из SVD.
 - эмбединги доминируются частыми словами!
- Модификация⁵: счётчик совстречаемости \rightarrow PPMI

$$PPMI(w_1, w_2) = \max\{0, PMI(w_1, w_2)\} = \max\{0, \ln \frac{P(w_1, w_2)}{P(w_1)P(w_2)}\}$$

- GloVe⁶: матр. факторизация $\log X \approx W^T \tilde{W} + B + \tilde{B}$

$$\sum_{i,j=1}^V f(X_{ij}) \left(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij} \right)^2 \rightarrow \min_{w, \tilde{w}, b, \tilde{b}}$$

$f(X_{ij})$ некоторая \uparrow функция весов, $f(0) = 0$.

⁴Lund and Burgess, 1996.

⁵Bullinaria and Levy, 2007

⁶<https://aclanthology.org/D14-1162.pdf>

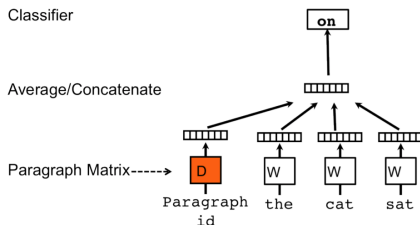
Содержание

- 1 Векторное представление слов
- 2 Word2vec
- 3 Регулярности в пространстве представлений
- 4 Настройка skip-gram
- 5 Методы на основе матрицы совстречаемости
- 6 Представления параграфов**
- 7 Сиамская сеть

Представления параграфов - мотивация

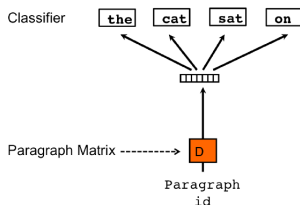
- Необходимо получить векторные представления параграфов (документов, предложений,...).
- Простой подход: усреднить слова, входящие в параграф.
 - или взвешенно усреднить, учитывая частоту встречаемости слов и их тематику.
- Точнее работает непосредственное представление самих параграфов.

Paragraph vector: модель PV-DM



- Во время обучения делим документы на параграфы. Каждому параграфу -> векторное представление.
- Оценивается CBOW, контекст: представление слов и параграфов.
- Можно усреднять или конкатенировать контексты слов и параграфа.
- Называется *Distributed Memory Model of Paragraph Vectors (PV-DM)*.

Paragraph vector: модель PV-DBOW



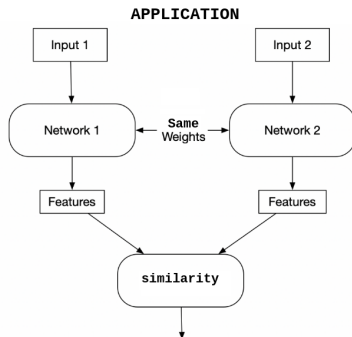
- Во время обучения делим документы на параграфы. Каждому параграфу -> векторное представление.
- Оценивается skip-gram: предсказываются случайные слова параграфа по представлению параграфа.
 - контекст: только представления параграфов
- Называется *Distributed Bag of Words version of Paragraph Vector (PV-DBOW)*

Содержание

- 1 Векторное представление слов
- 2 Word2vec
- 3 Регулярности в пространстве представлений
- 4 Настройка skip-gram
- 5 Методы на основе матрицы совстречаемости
- 6 Представления параграфов
- 7 Сиамская сеть**

Сиамская сеть

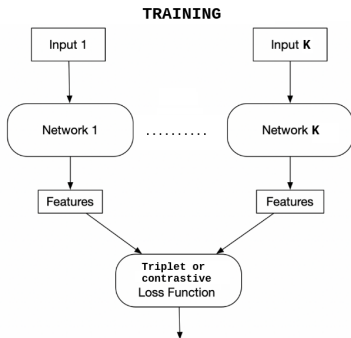
- Сиамская сеть (siamese network) генерирует 2 представления произвольных объектов.
 - объекты могут быть разных доменов.
- Обучение: похожие объекты \Rightarrow похожие представления.
 - похожесть: $\|\cdot\|_2^2$, $\langle \cdot, \cdot \rangle$, \cos – sim



Примеры приложений

- Классификация:
 - вход: 2 объекта (обучение) или тестовый объект (применение)
 - выход: класс на основе близости к центроиду класса или по K-NN
- Поисковая система
 - вход: документ и поисковый запрос
 - возможен и поиск по изображению
 - выход: степень релевантности документа запросу
- Обнаружение перефразирования:
 - вход: 2 предложения
 - выход: насколько они близки по смыслу
- Проверка подписи
 - вход: сканы 2х подписей
 - выход: степень их принадлежности одному человеку

Обучение



- Идея функции потерь:
 - представления похожих объектов должны быть близки
 - представления различных объектов должны быть далеки

Функции потерь⁷

Контрастные потери (contrastive loss):

- обучение на случайных парах объектов x_i, x_j

$$\mathbb{I}[y_i = y_j] \|f_{\theta}(x_i) - f_{\theta}(x_j)\|^2 + \mathbb{I}[y_i \neq y_j] \max\{0, \alpha - \|f_{\theta}(x_i) - f_{\theta}(x_j)\|\}^2$$

⁷ Обзор более продвинутых ф-ций потерь

Функции потерь⁷

Контрастные потери (contrastive loss):

- обучение на случайных парах объектов x_i, x_j

$$\mathbb{I}[y_i = y_j] \|f_\theta(x_i) - f_\theta(x_j)\|^2 + \mathbb{I}[y_i \neq y_j] \max\{0, \alpha - \|f_\theta(x_i) - f_\theta(x_j)\|\}^2$$

Тройные потери (triplet loss):

- обучение на случайных тройках x, x^+, x^- .
- x - опорный объект (anchor)
- x^+ - похожий на x (например, того же класса)
- x^- - не похожий на x (например, др. класса)
- $\alpha > 0$ - гиперпараметр

$$\mathcal{L}(x, x^+, x^-) = \max\left\{\|f_\theta(x) - f_\theta(x^+)\|^2 - \|f_\theta(x) - f_\theta(x^-)\|^2 + \alpha; 0\right\}$$

⁷Обзор более продвинутых ф-ций потерь

Функции потерь⁷

Контрастные потери (contrastive loss):

- обучение на случайных парах объектов x_i, x_j

$$\mathbb{I}[y_i = y_j] \|f_\theta(x_i) - f_\theta(x_j)\|^2 + \mathbb{I}[y_i \neq y_j] \max\{0, \alpha - \|f_\theta(x_i) - f_\theta(x_j)\|\}^2$$

Тройные потери (triplet loss):

- обучение на случайных тройках x, x^+, x^- .
- x - опорный объект (anchor)
- x^+ - похожий на x (например, того же класса)
- x^- - не похожий на x (например, др. класса)
- $\alpha > 0$ - гиперпараметр

$$\mathcal{L}(x, x^+, x^-) = \max\left\{\|f_\theta(x) - f_\theta(x^+)\|^2 - \|f_\theta(x) - f_\theta(x^-)\|^2 + \alpha; 0\right\}$$

Помимо сиамских сетей потери используются для metric learning $\rho_\theta(x, x')$.

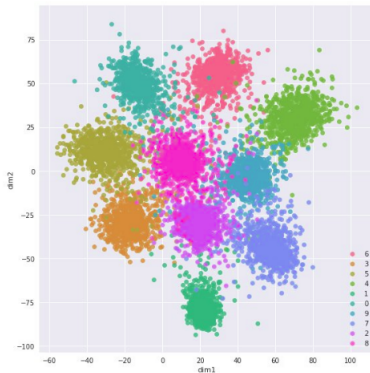
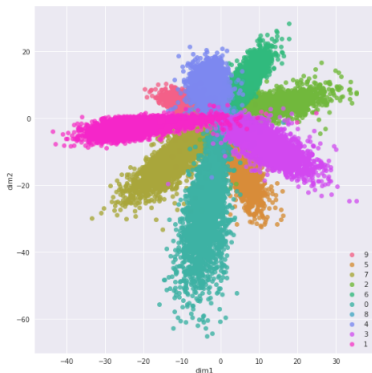
⁷Обзор более продвинутых ф-ций потерь

Сиамская сеть и классификация

- Классификация
 - выучивает "что представляет каждый класс".
 - выдает степени соответствия x каждому классу.
- Сиамская сеть
 - выучивает "что отличает классы друг от друга".
 - выдает расстояния от x до каждого класса.
 - более устойчива к дисбалансу классов и редким классам (*one shot learning*)
 - при обучении каждый класс учитывается поровну
 - модель выучивает признаки, по которым можно судить о сходстве классов на частотных классах, потом сразу подхватывает их для редких.
 - извлекает больше информации из выборки
 - обучение не на объектах, а на парах и тройках объектов.
 - хороша в ансамбле с классификатором (\uparrow разнообразие)

Представления объектов: классификация и сиамская сеть

Представления объектов: классификация и сиамская сеть для MNIST:



Заключение

- **Представления слов** отображают слова в компактные векторные представления.
 - может применяться
 - к биграммам, триграммам, коллокациям.
 - к символам - удобно для новых слов
 - к любым объектам из посл-тей (нуклеотиды в ДНК и др.)
- **Представления параграфов** отображают параграфы в векторные представления.
 - работают лучше, чем усреднение слов параграфа
- Представления можно находить для целевой или связанной задачи (language modeling, transfer learning)
- **Сиамская сеть** оценивает похожесть пар объектов.
 - применения: классификация (особенно one shot learning), детекция перефразирования, нахождение похожих изображений, ...