



Эмбеддинги. Word2vec,
fasttext



Представления слов

Слова представляются в виде one-hot вектора:

умный	$[0, \dots, 1, \dots, 0, \dots, 0, \dots, 0]$
сообразительный	$[0, \dots, 0, \dots, 1, \dots, 0, \dots, 0]$
собака	$[0, \dots, 0, \dots, 0, \dots, 1, \dots, 0]$

Недостатки:

- Большая размерность векторов;
- Все вектора одинаково похожи;
- Сложно кодировать дополнительную информацию.



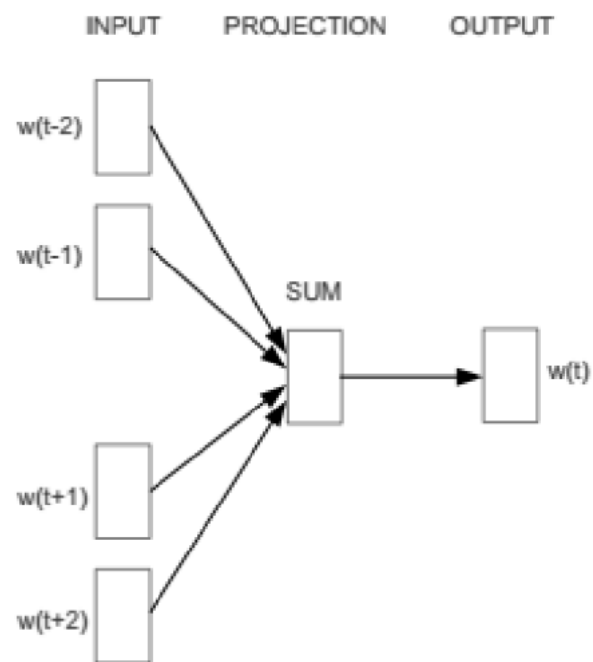
Гипотеза компактности

Ближние (похожие) слова встречаются в похожем контексте.

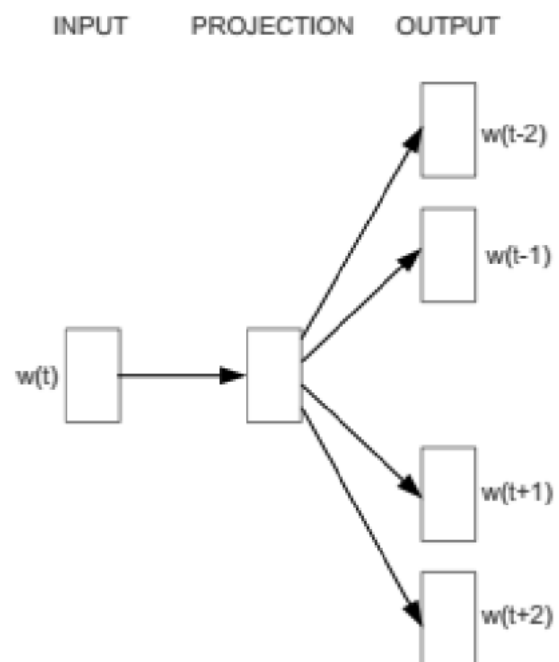
Идея: у близких слов (синонимов, родственных слов) и представления получают близкие.



Word2Vec



CBOW

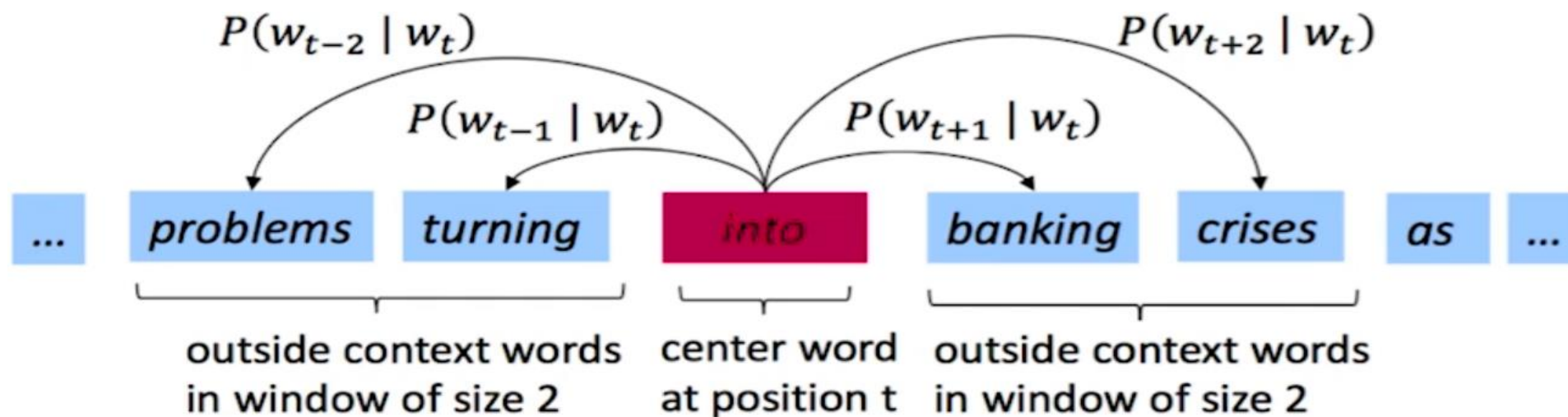


Skip-gram



Skip-gram

Зная контекст, мы хотим максимизировать вероятность центрального слова. Перемещаемся скользящим окном.



Skip-gram

Вероятность встретить слово w_0 рядом со словом w_1 :

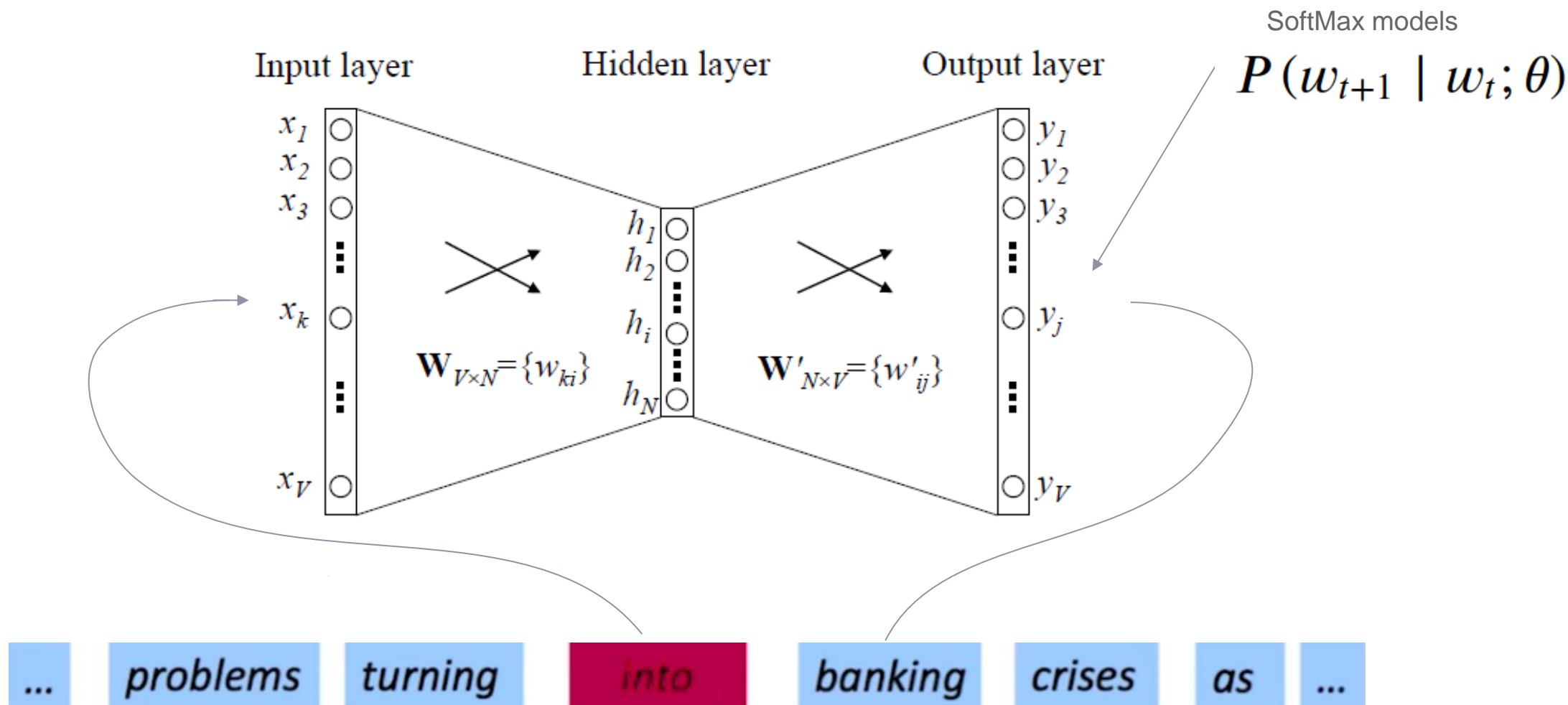
$$p(w_0|w_1) = \frac{\exp(\langle v'_{w_0}, v_{w_1} \rangle)}{\sum_{w \in W} \exp(\langle v'_w, v_{w_1} \rangle)}$$

Функционал для текста $T=(w_1 w_2 \dots w_n)$:

$$\sum_{i=1}^n \sum_{\substack{-c \leq j \leq c \\ j \neq 0}} \log p(w_{i+j}|w_i) \rightarrow \max$$



Skip-gram



Skip-gram

Простейшая модель классификации — линейный классификатор, применяемый к вектору контекста.

$$\textit{softmax}(\mathbf{x})_i = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)} = p_i$$

Для обучения минимизируем функцию потерь.



Skip-gram

Вероятность встретить слово w_o рядом со словом w_i :

$$p(w_o|w_i) = \frac{\exp(\langle v'_{w_o}, v_{w_i} \rangle)}{\sum_{w \in W} \exp(\langle v'_w, v_{w_i} \rangle)}$$

Считать знаменатель крайне затратно. Значит, и производные считать тоже долго.



Negative sampling

$$p(w_o | w_I) = \log \sigma(\langle v'_{w_o}, v_{w_I} \rangle) + \sum_{i=1}^k \log \sigma(-\langle v'_{w_i}, v_{w_I} \rangle)$$

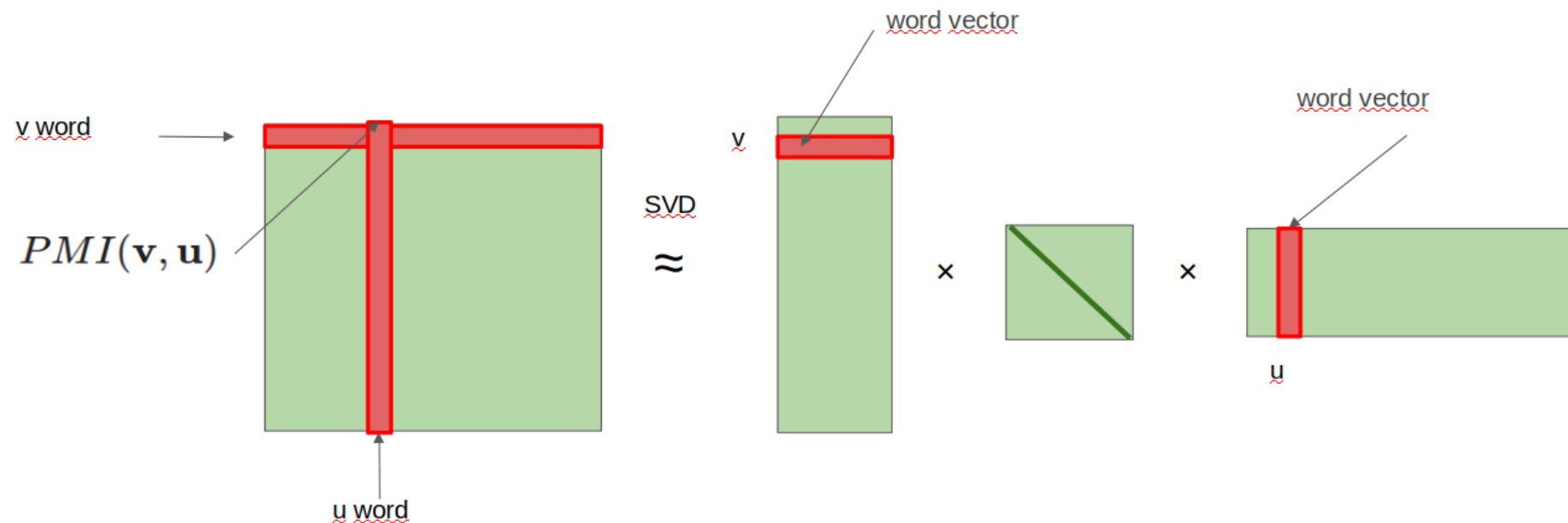
w_i — случайно выбранные слова. Слово w генерируется с вероятностью $P(w)$ — шумовое распределение.

$$P(w) = \frac{U(w)^{\frac{3}{4}}}{\sum_{v \in W} U(v)^{\frac{3}{4}}}, U(v) \text{ — частота слова } v \text{ в корпусе текстов.}$$



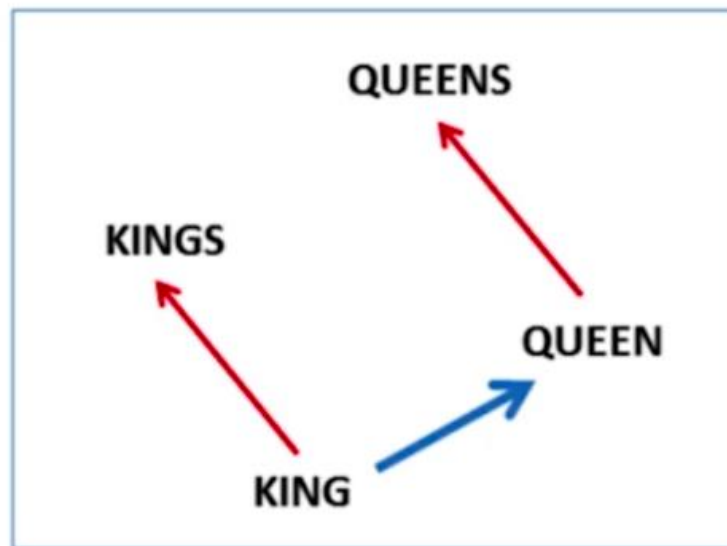
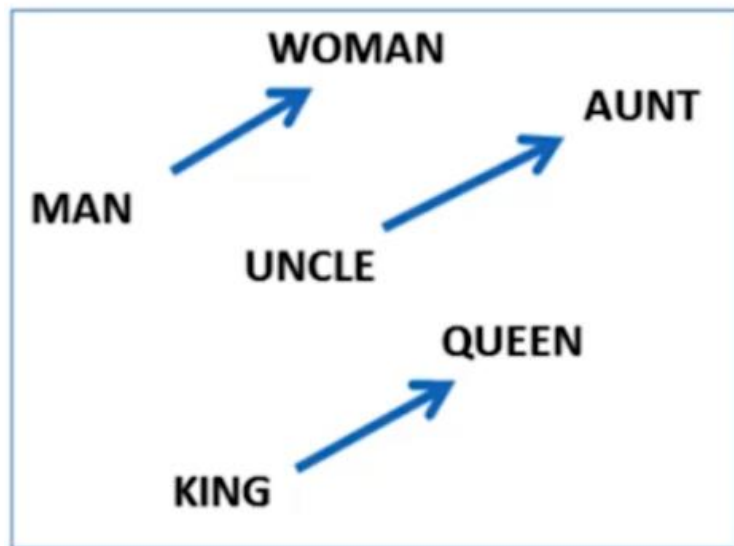
Word2vec vs SVD

Word2Vec with negative sampling \approx matrix factorization



Word2vec

$$\underline{v(\text{king})} - \underline{v(\text{man})} + \underline{v(\text{woman})} \approx \underline{v(\text{queen})}$$



Fasttext

Разделите слово на мешок из n -граммов: *apple* = $\langle ap, ppl, ple, le \rangle$ (BPE).

Вычислите вектор для каждого n -грамма.

Вектор для слова = сумма <вектора слова, векторов для n -граммов слов>.



Ваши вопросы?

Классические подходы к обработке естественного языка, урок 3

