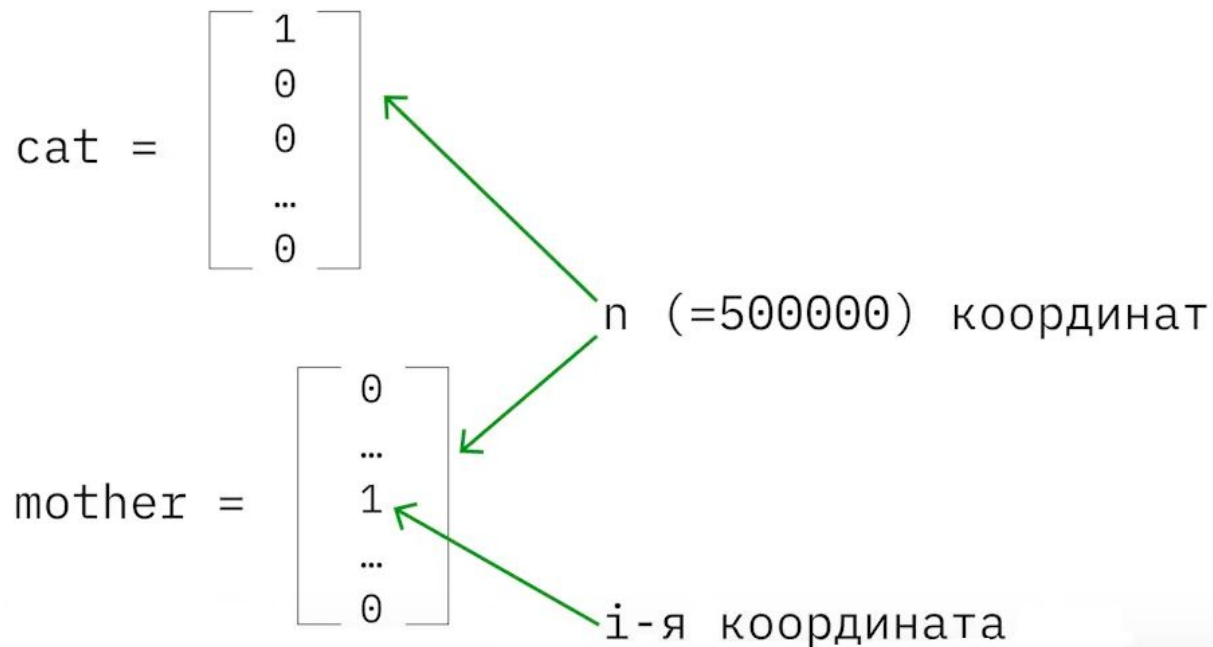


Word Embeddings

Простой способ: one-hot encoding

One-hot вектор размерности длины словаря
(например, 500.000)



Bag-of-Words

Сумма one-hot векторов слов

my dog is on the table



my cat dog is now on table the

the dog is on the table



are cat dog is now on table the

TF-IDF

$P(w, d, n_{dw}) = (N_w/N)^{n_{dw}}$ - вероятность встретить n_{dw} раз слово w в документе d

$$-\log P(w, d, n_{dw}) = n_{dw} \cdot \log(N/N_w) = TF(w, d) \cdot IDF(w)$$

$TF(w, d) = n_{dw}$ - term frequency;

$IDF(w) = \log(N/N_w)$ - inverted document frequency;

Context embeddings

Let's take into account words meanings in some way:

$v(\text{word}_i)[j] = \text{count}(\text{co-occurrences } \text{word}_i \text{ with } \text{word}_j \text{ in dataset})$

$v(\text{word}_1) = [12, 1, 0, 10, 5, \dots]$
 ↑ ↑ ↑ ↑ ↑
 horse ride wheel roof hair breed

$v(\text{word}_2) = [20, 10, 0, 0, 1, \dots]$
 ↑ ↑ ↑ ↑ ↑
 car ride wheel roof hair breed

Pointwise mutual information

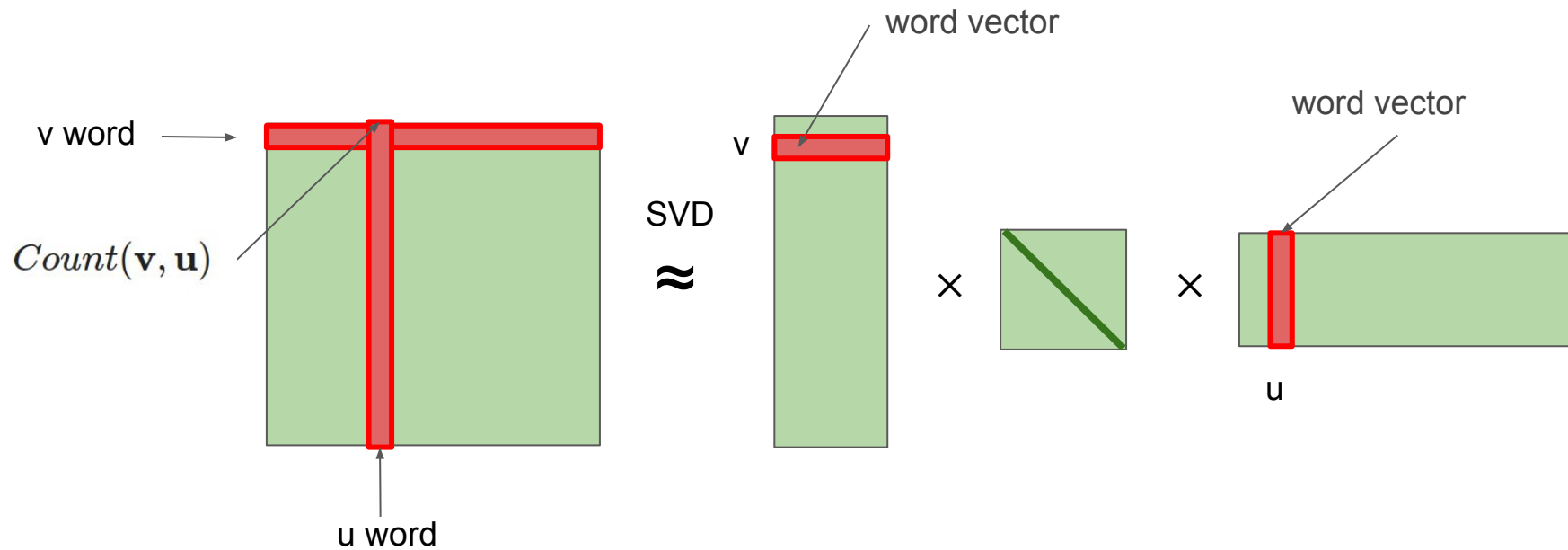
Скользящее окно фиксированной длины:

n_{uv} - встречаемость слова u и v вместе

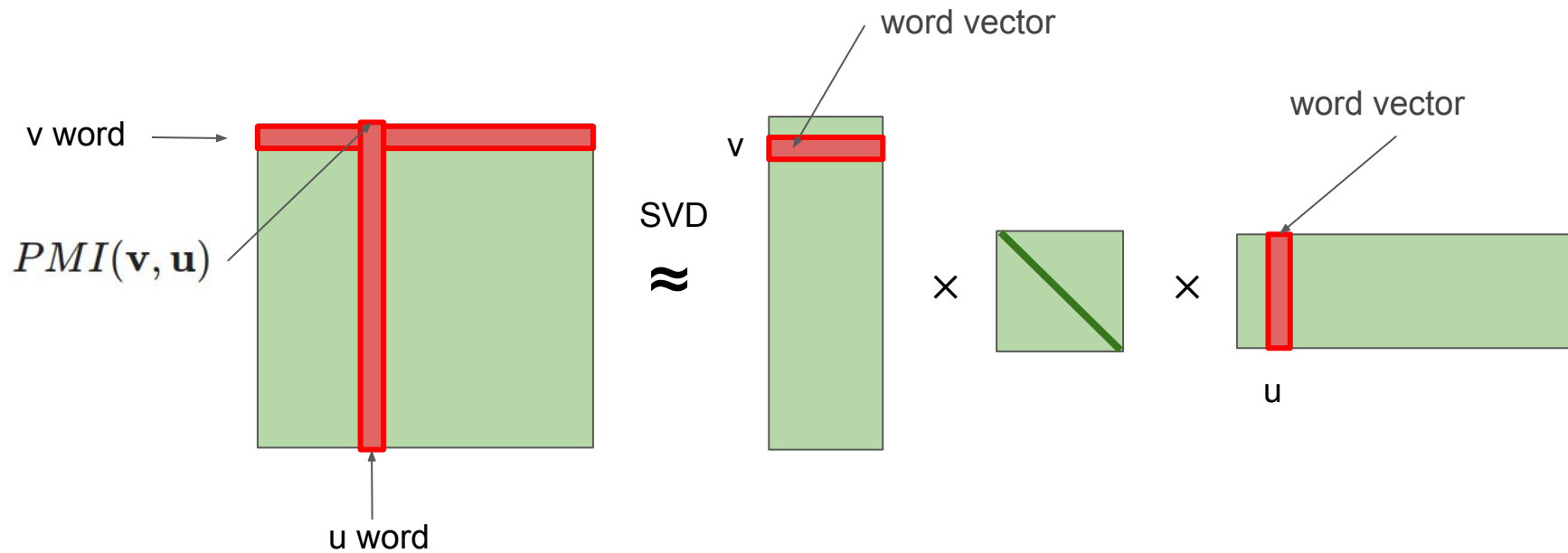
$$PMI = \log \frac{p(u,v)}{p(u)p(v)} = \log \frac{n_{uv}n}{n_u n_v}$$

$$pPMI = \max(0, PMI)$$

Co-Occurrence Counts



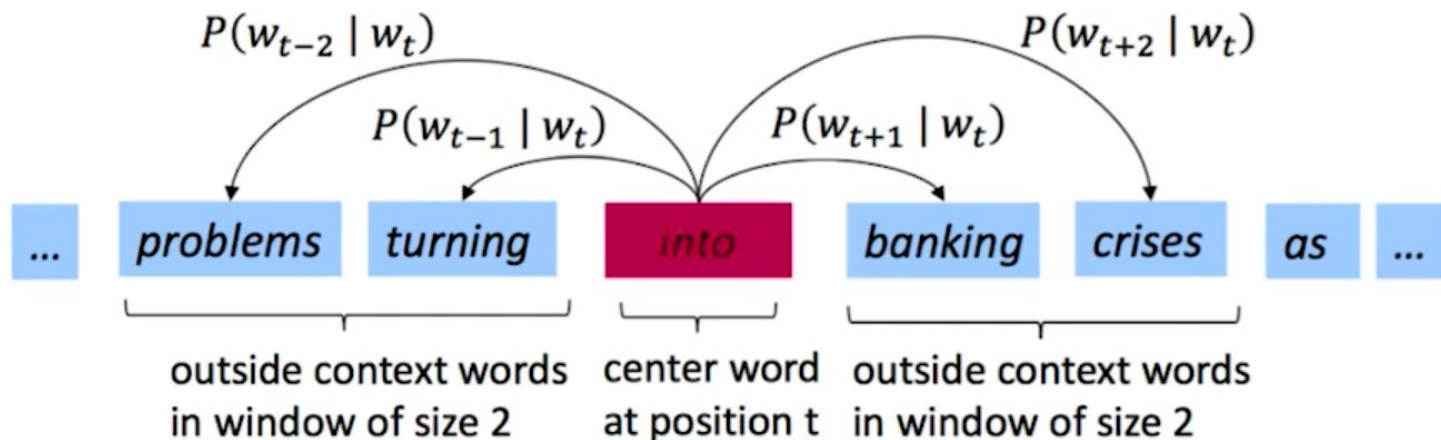
Pointwise mutual information



Word2Vec

We want to maximize probabilities of seeing a **surrounding** word based on **center** words.

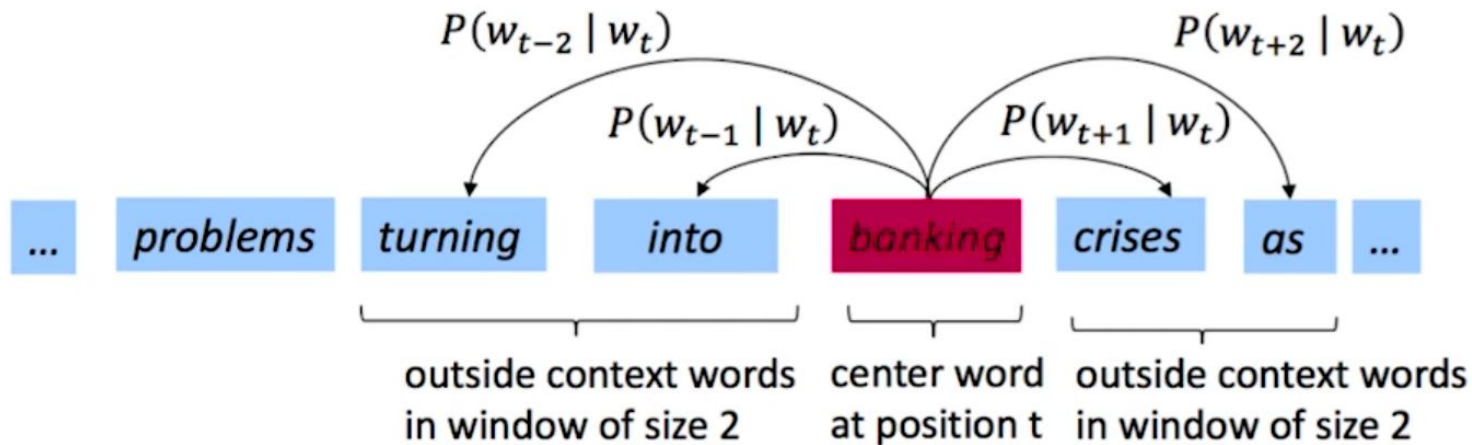
Going through text corpus by sliding windows.



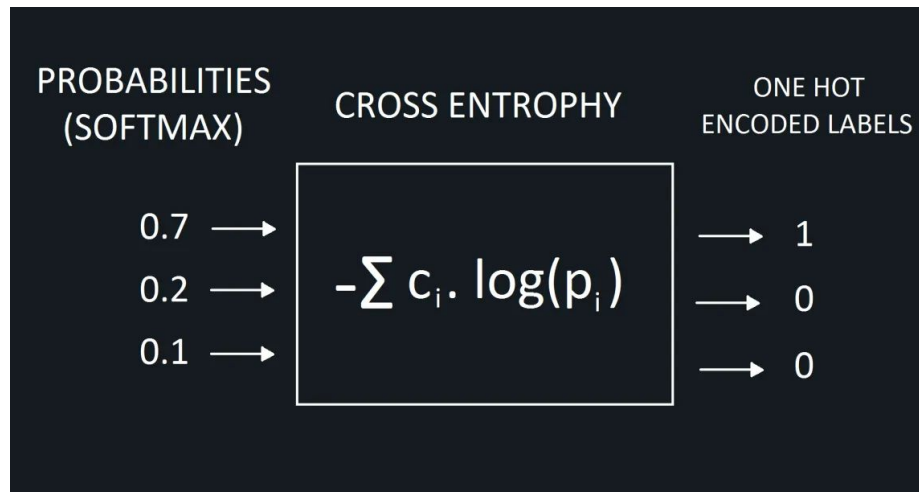
Word2Vec

We want to maximize probabilities of seeing a **surrounding** word based on **center** words.

Going through text corpus by sliding windows.



Cross-Entropy



$$H(q, p) = - \sum_x q(x) \log p(x)$$

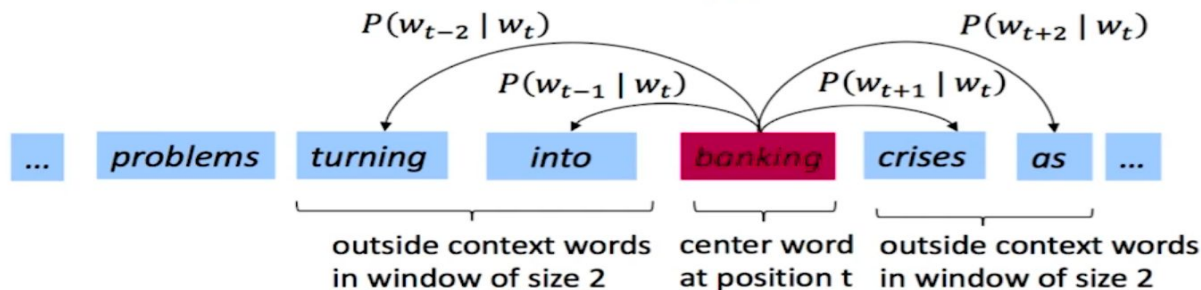
Word2Vec

Target function:

$$L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j} | w_t; \theta)$$

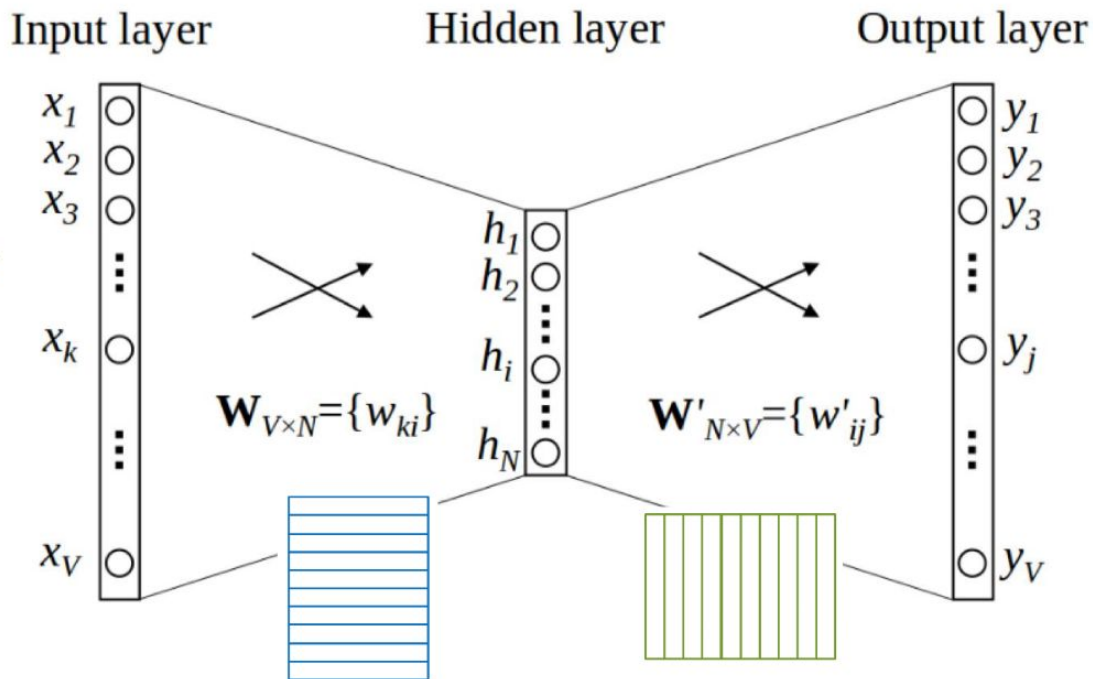
Or in log-likelihood point of view:

$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t; \theta)$$

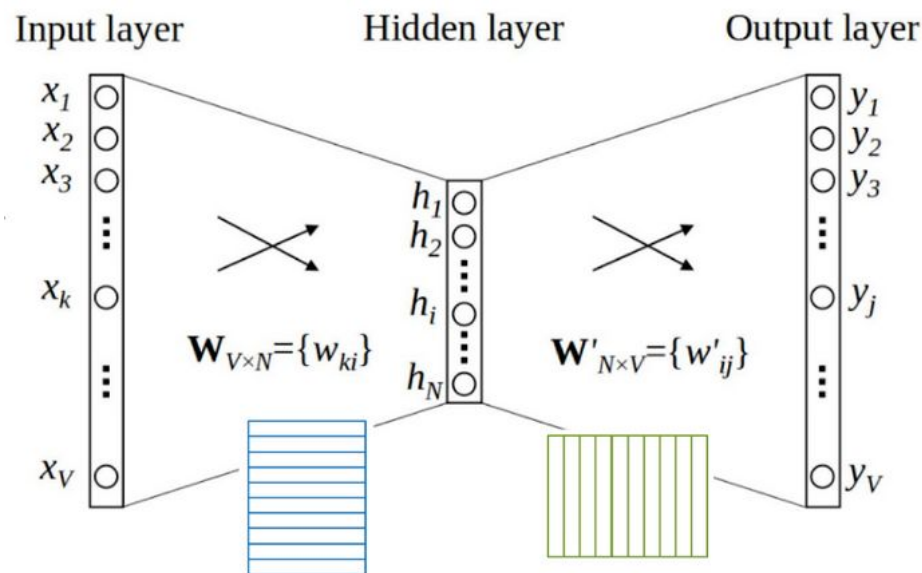


Word2Vec

$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix} \in \mathbb{R}^{2NV}$$



Word2Vec



$$\mathbf{W} - V \times N$$

$$\mathbf{W}' - N \times V$$

$$1. \mathbf{W}^T \cdot \mathbf{x} = h \implies (N \times V) \cdot (V \times 1) = N \times 1$$

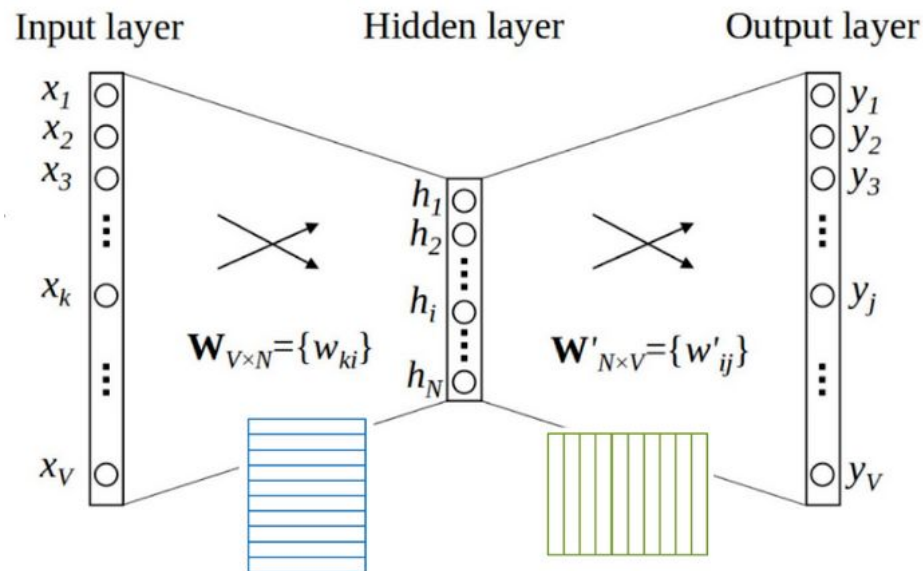
$$\mathbf{W}^T = [w_1^T w_2^T \cdots w_V^T] \implies \sum_{i=1}^V w_i^T \mathbf{x}_i = w_k^T = h$$

Word2Vec

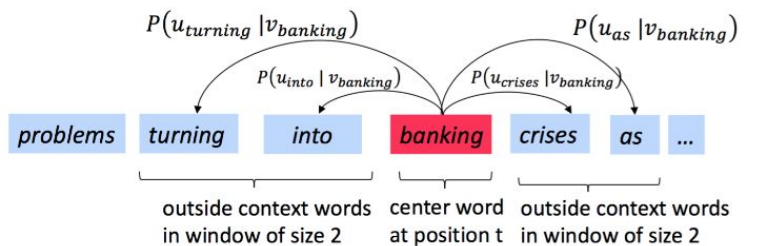
$$\mathbf{W} - V \times N$$
$$\mathbf{W}' - N \times V$$

$$2. \mathbf{W}'^T \cdot h = y \implies (V \times N) \cdot (N \times 1) = V \times 1$$

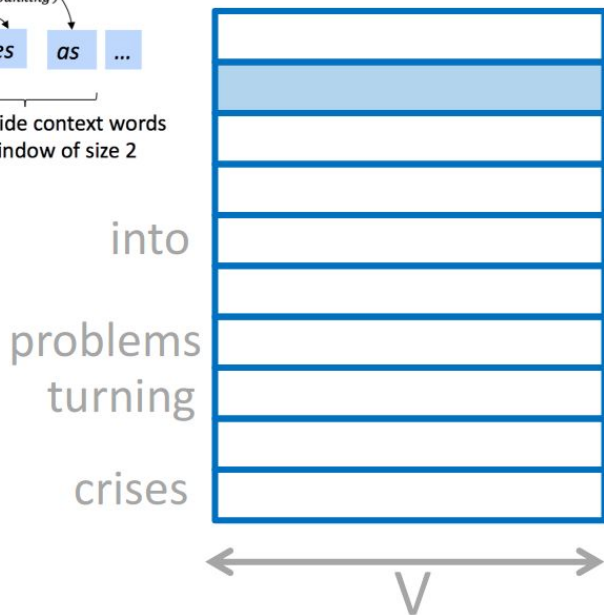
$$\mathbf{W}' = [w'_1 w'_2 \cdots w'_V] \implies y_j = (\mathbf{W}'^T \cdot h)_j = (w'_j)^T w_k^T = \langle w_k, w'_j \rangle$$



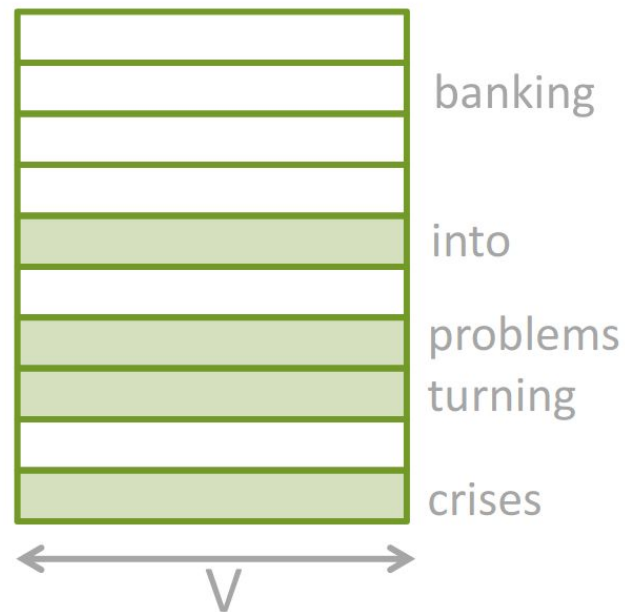
Word2Vec



When it is a
center word

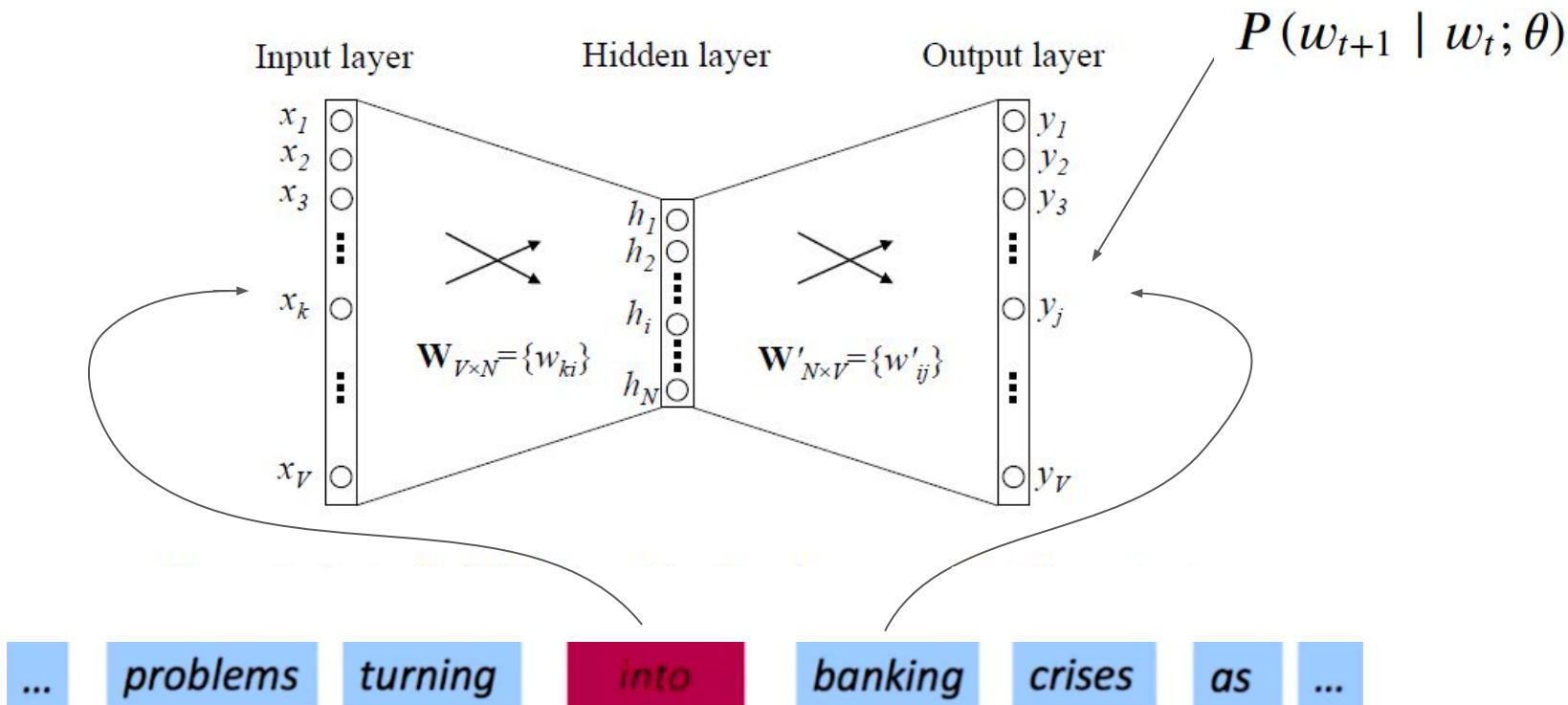


When it is a
context word



Word2Vec

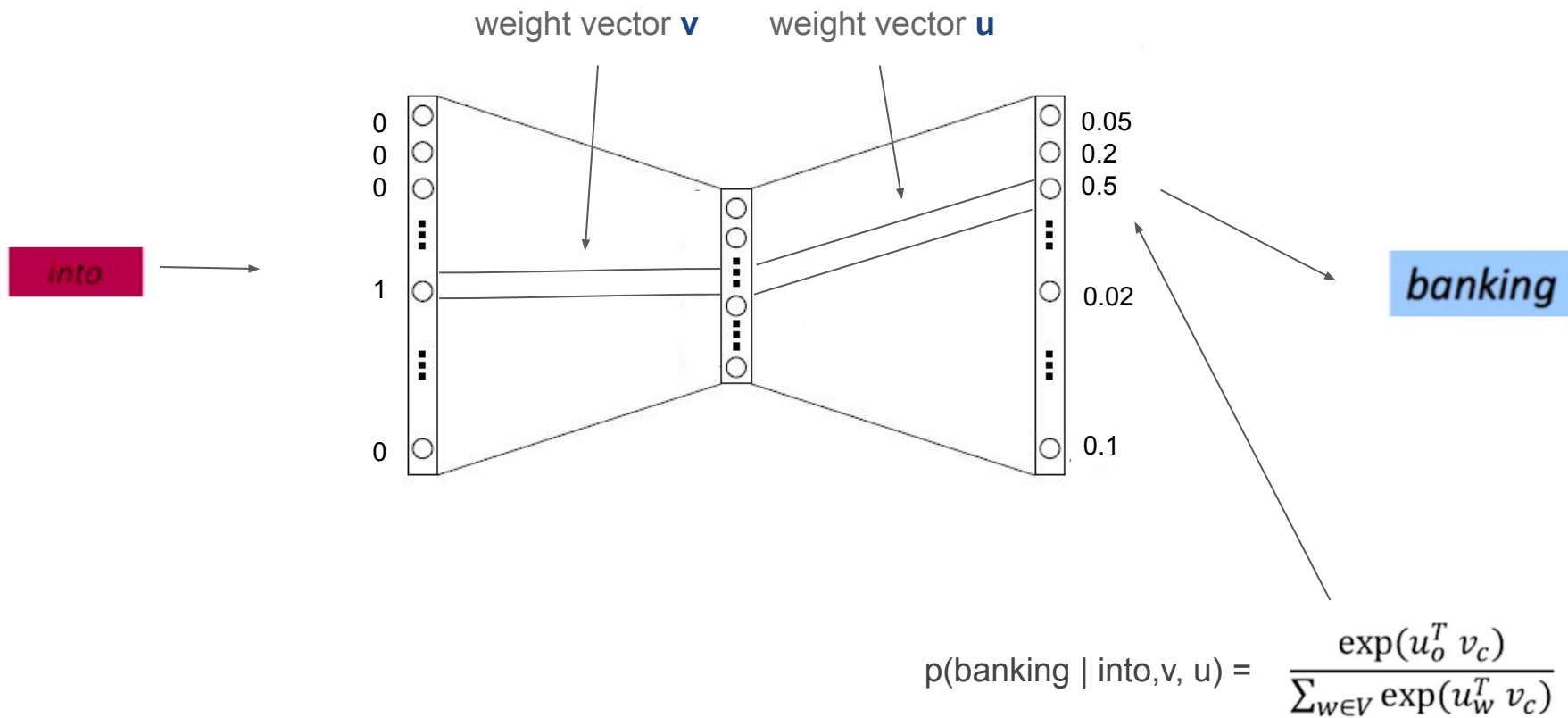
SoftMax models



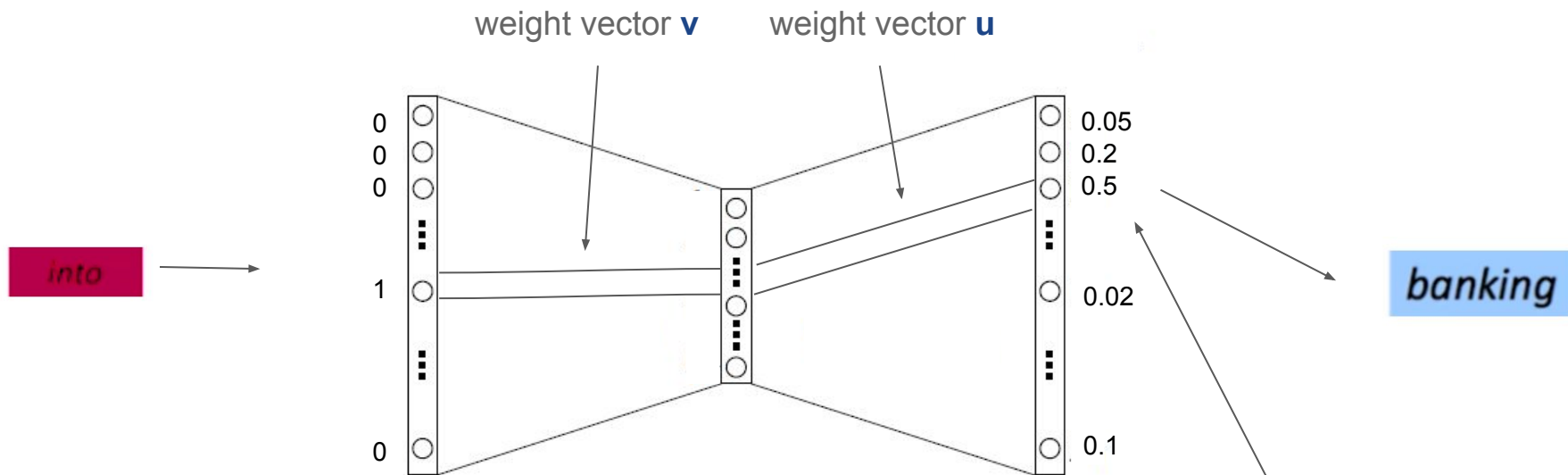
Softmax

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

Word2Vec



Word2Vec



- u is a context word vector
- v is a center word vector

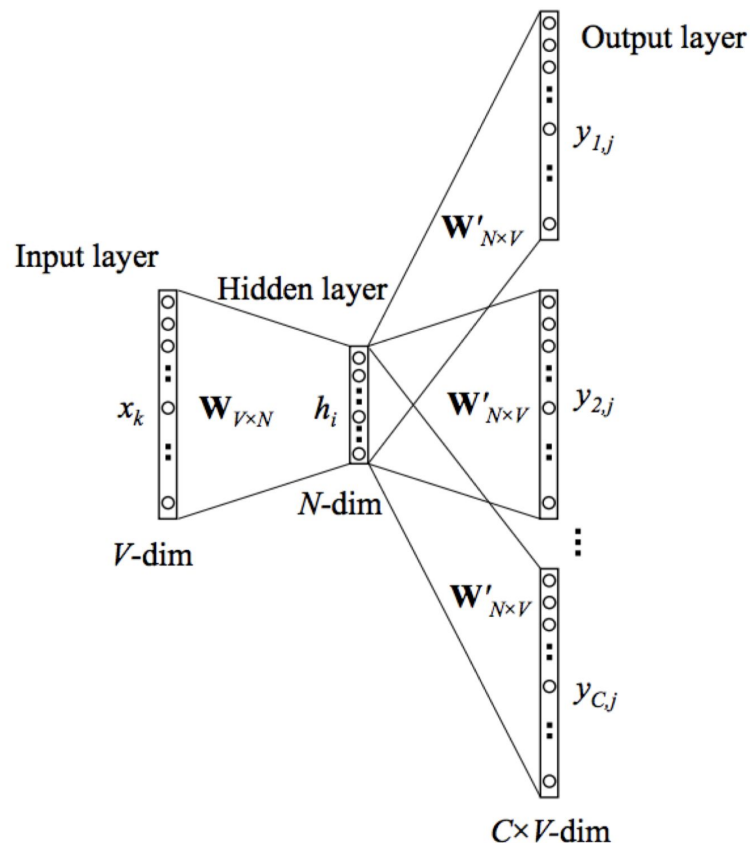
“Scalar product between me and my neighbour must be as big as possible”

$$p(\text{banking} \mid \text{into}, v, u) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

Word2Vec: Skip-Gram

Maximize probabilities of seeing a
surrounding word based on center words.

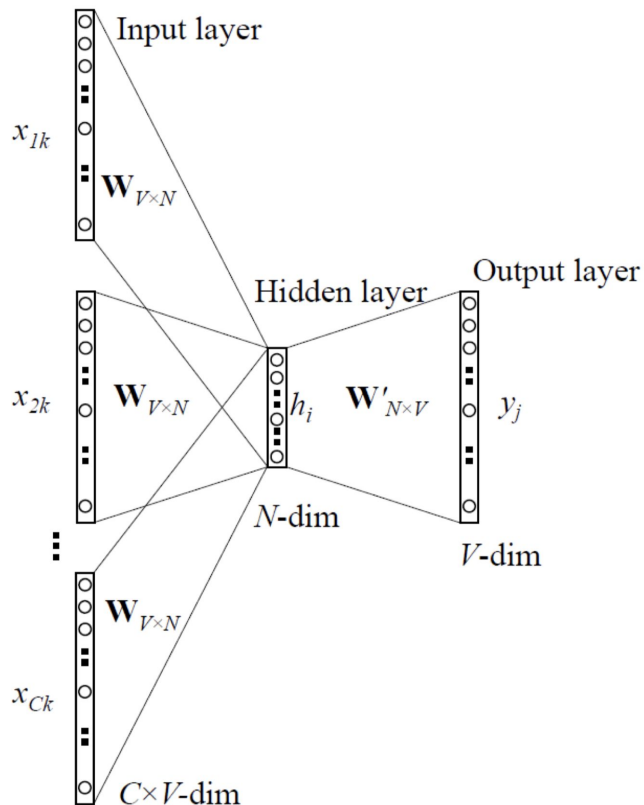
$$L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j} | w_t; \theta)$$



Word2Vec: CBOW

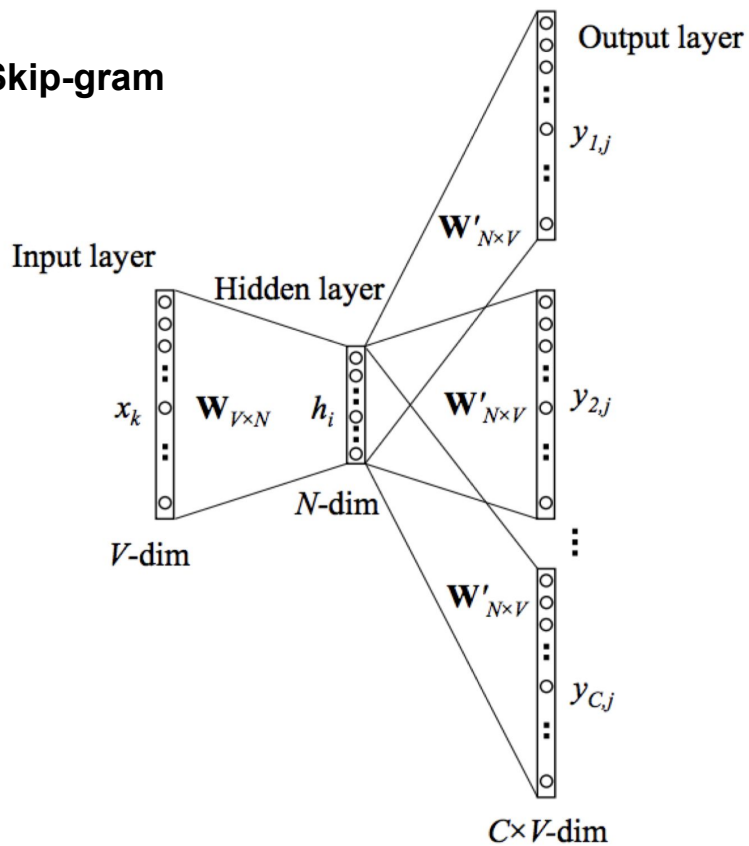
Maximize probabilities of seeing a **center** word based on **surrounding** words.

$$L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_t | w_{t+j}; \theta)$$

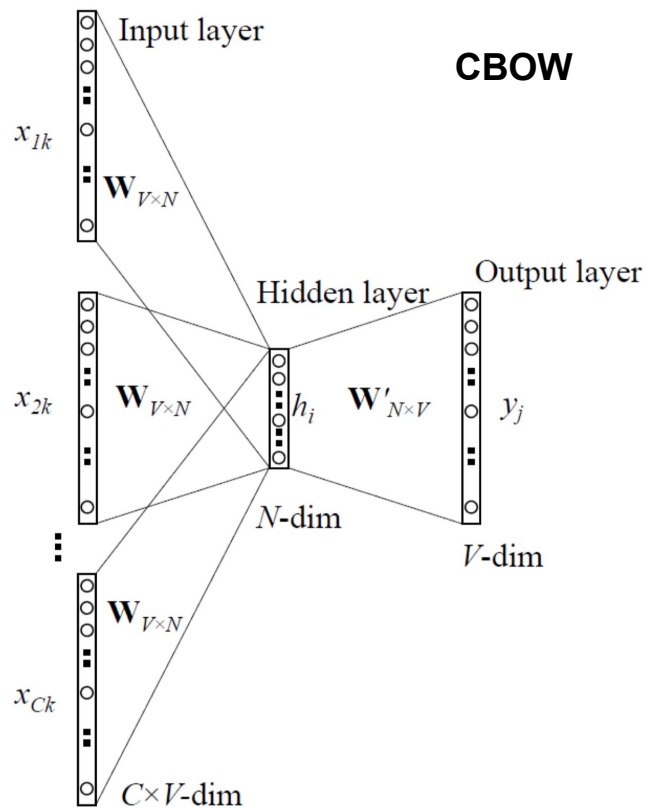


Word2Vec

Skip-gram



CBOW



Word2Vec

Problem:

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} \quad \leftarrow \text{Still big sum to compute}$$

Word2Vec

Problem:

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} \quad \leftarrow \text{Still big sum to compute}$$

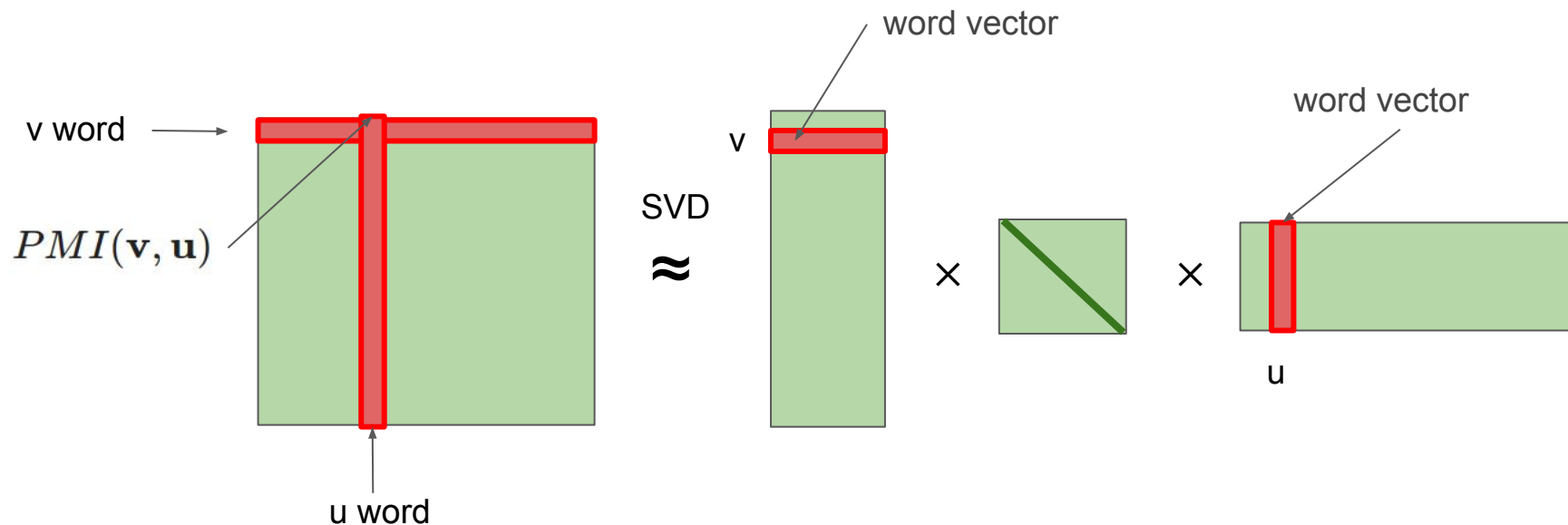
Possible solution:

- Negative sampling

$$P(w_i) = \frac{f(w_i)}{\sum_{j=0}^n (f(w_j))}$$

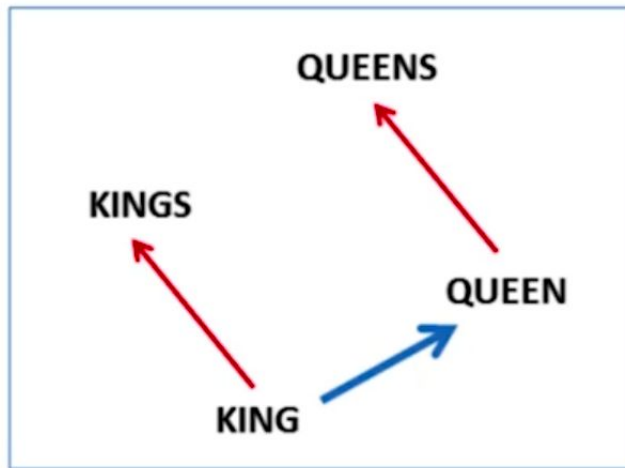
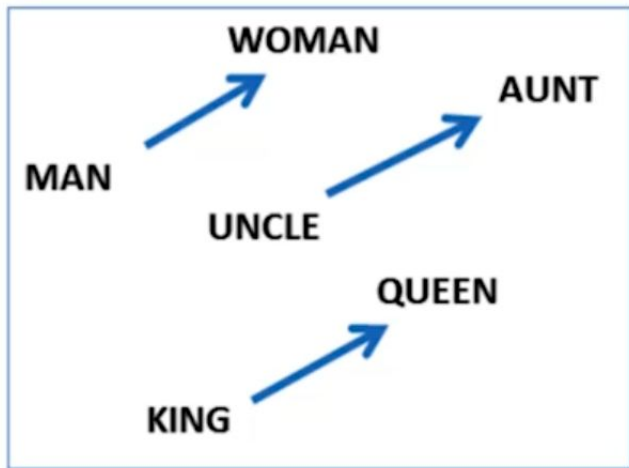
Word2Vec vs SVD

Word2Vec with negative sampling \approx matrix factorization - [Link](#)



Word2Vec

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



Word2Vec

Target function:

$$L(\theta) = \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log P(w_{t+j} \mid w_t, \theta)$$

$$L_{\log}(\theta) = \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log P(w_{t+j} \mid w_t, \theta) = \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log \frac{\exp(u_{t+j}^T v_t)}{\sum_{w \in V} \exp(u_w^T v_t)}$$

$$L_{\log}(\theta) = \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} (u_{t+j}^T v_t - \log \sum_{w \in V} \exp(u_w^T v_t))$$

GloVe

Before training count occurrences of pairs $[\text{word}_i, \text{word}_j]$ in corpus

Compute probabilities: $P_{ij} = \frac{\text{Count}(v_i, v_j)}{\text{Count}(v_i)}, \text{Count}(v_i) = \sum_k \text{Count}(v_i, v_k)$

Objective function:

$$J(\theta) = \frac{1}{2} \sum_{i,j=1}^W \boxed{f(P_{ij})} (u_i^T v_j - \log P_{ij})^2$$

Discount factor for rare words

FastText

- Divide word into bag of n-grams: apple = <ap, ppl, ple, le>
- Compute vector for each n-gram
- Vector for a word = sum of vectors for word n-grams
- Hash table
- Some n-grams have the same vector

FastText

- Divide word into bag of n-grams: apple = <ap, ppl, ple, le> (**BPE**)
- Compute vector for each n-gram
- Vector for a word = sum of vectors for word n-grams
- Hash table
- Some n-grams have the same vector

Advantages:

- Reasonable embeddings for rare words and words with mistakes
- Model is the same as before, we can even use model trained on words to train it further on n-grams!

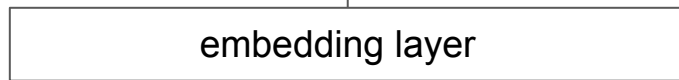
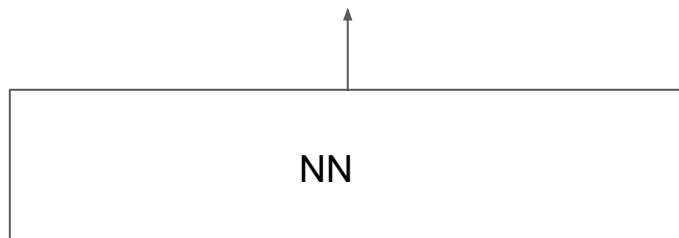
Quality of embeddings

- Оценка семантической близости между словами
 - поиск близких слов
 - поиск лишних слов
- Поиск аналогий
- Решение почти любой NLP задачи
 - ранжирование
 - классификация

Classification task

When you have small
text data for your task

Output (class,
e.g. genre, or
emotion)



do something, e.g.
average them all
or take sentence
embedding

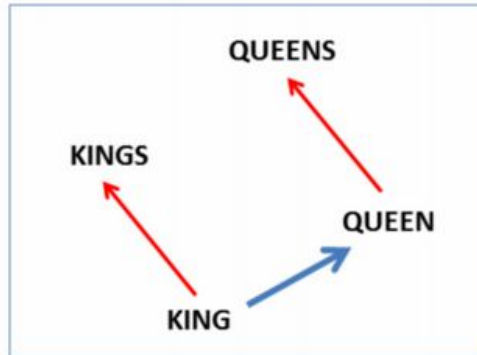
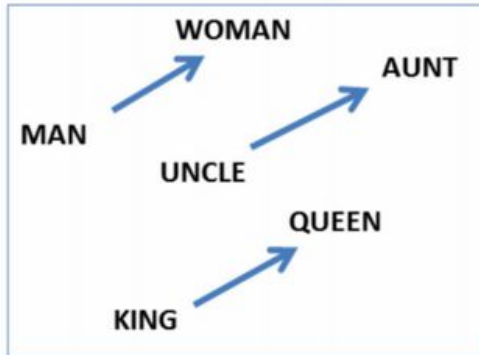
I have a cat UNK

pre-trained

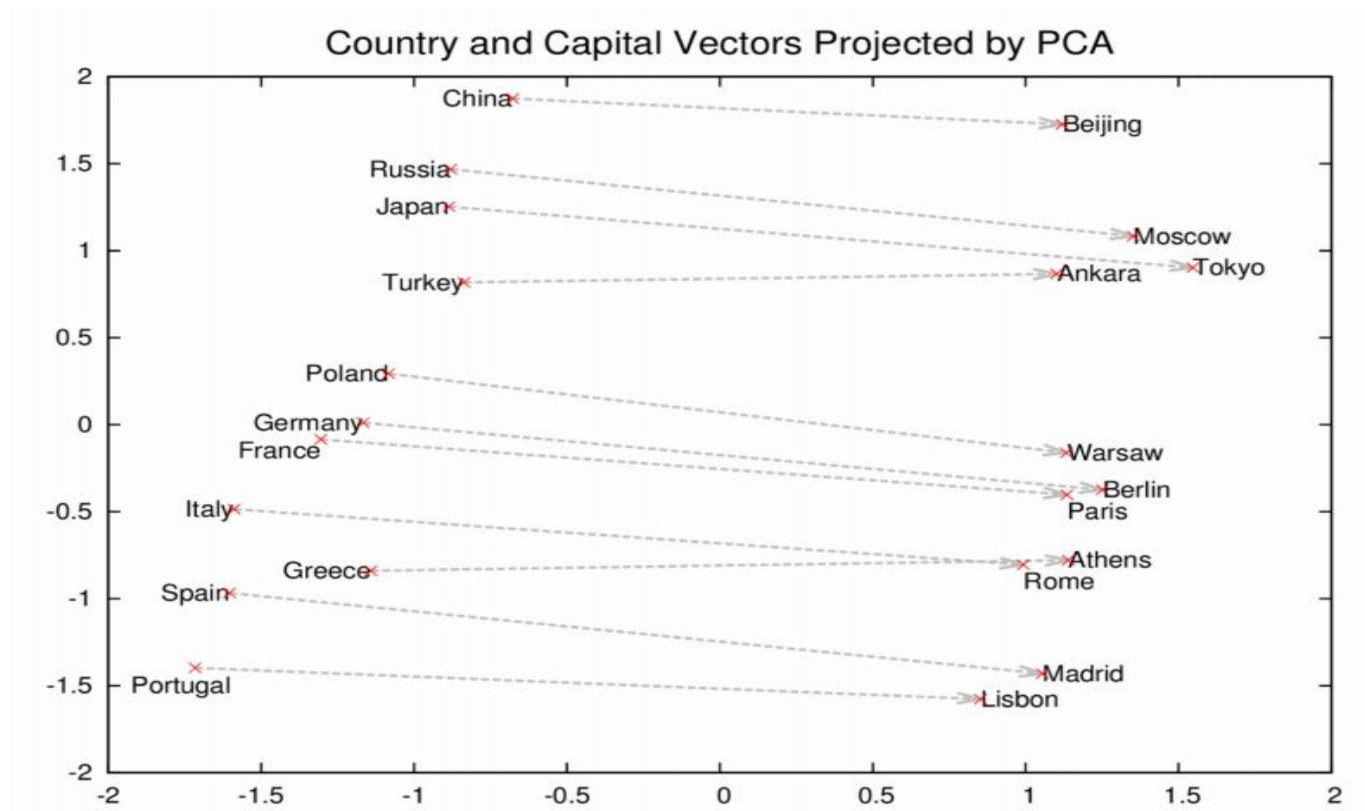
Поиск аналогий

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

syntactic: $v(\text{kings}) - v(\text{king}) + v(\text{queen}) \approx v(\text{queens})$

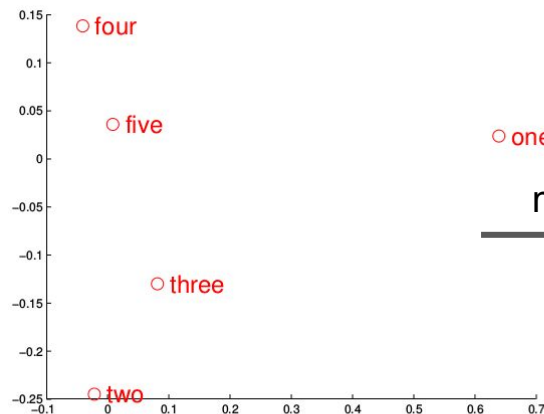


Поиск аналогий

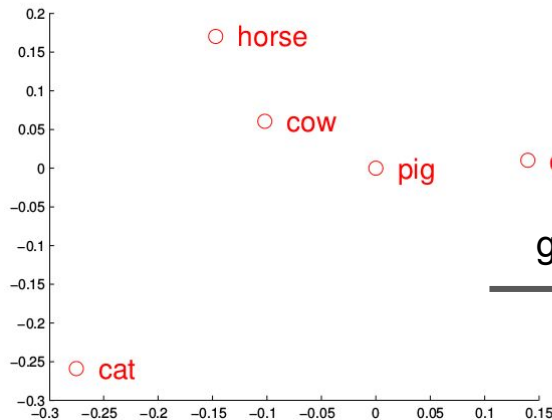
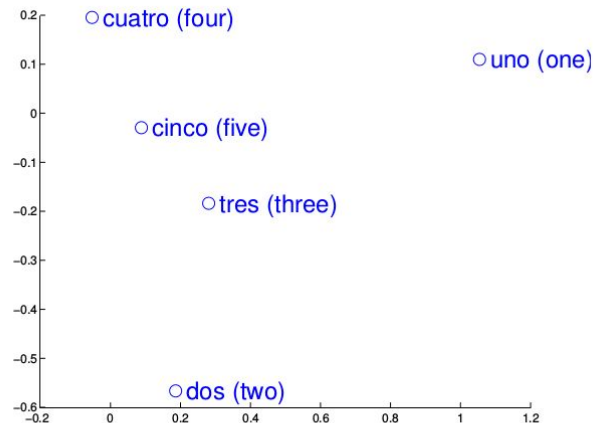


Language similarities

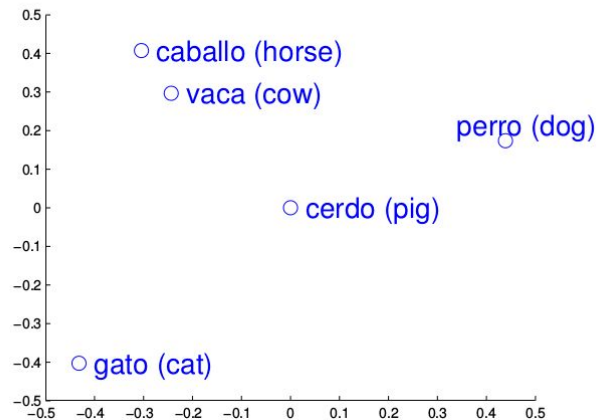
- train embeddings for english
- find mapping $f()$ from english to spanish
- get new english word -- use $f()$ to compute translation!



map

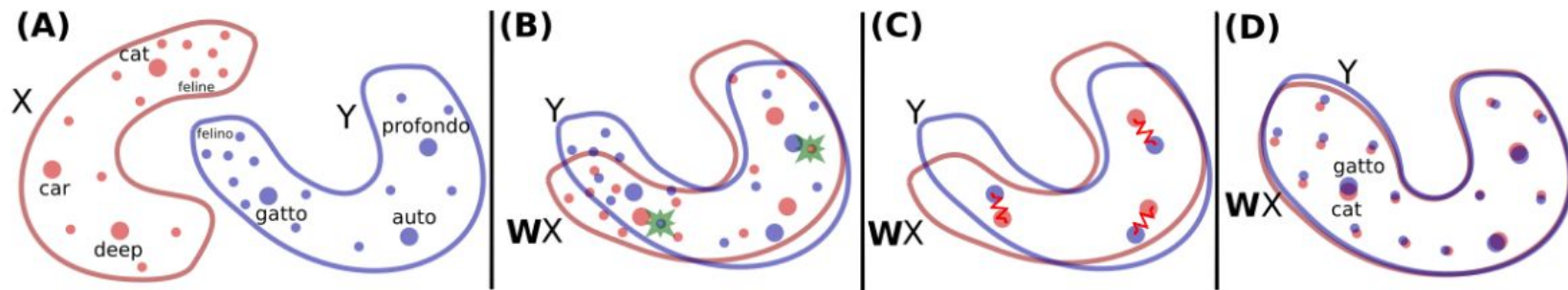


get!



Language similarities

$$W^* = \underset{W \in M_d(\mathbb{R})}{\operatorname{argmin}} \|WX - Y\|_F$$



Finally...

Okay, that's great, but why do we actually need embeddings?

Timeline

