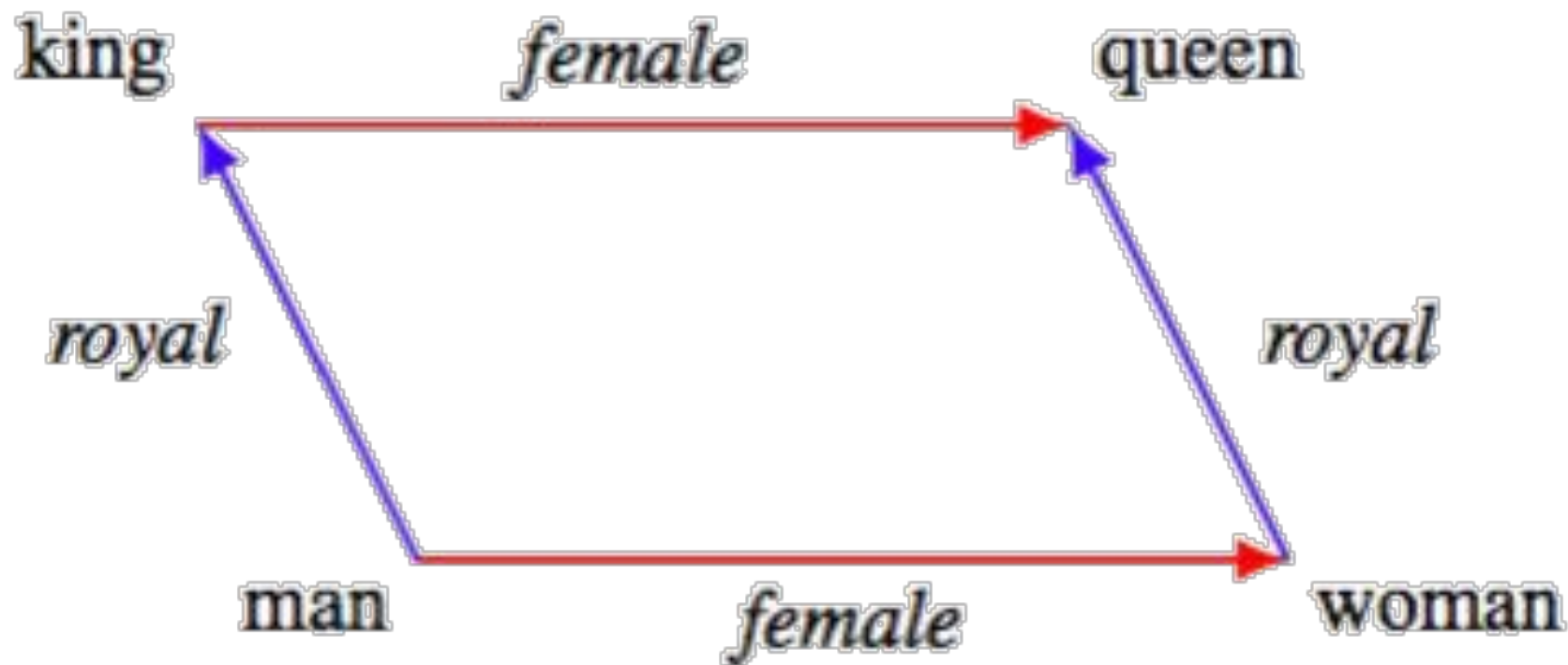


Distributed Representations of Words and Phrases and their Compositionality

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado,
Jeffrey Dean

Paper report by:
Nikita Ivlev, Lev Leontev, Aleksandr Kariakin

Introduction



Skip-gram architecture

Skipgrams

Step - 1

The product is really good The product is wonderful The product is awful

Step - 3

1	The	1	0	0	0	0	0	0
2	product	0	1	0	0	0	0	0
3	is	0	0	1	0	0	0	0
4	really	0	0	0	1	0	0	0
5	wonderful	0	0	0	0	1	0	0
6	good	0	0	0	0	0	1	0
7	awful	0	0	0	0	0	0	1



Skip-gram architecture

Skipgrams

Step - 1

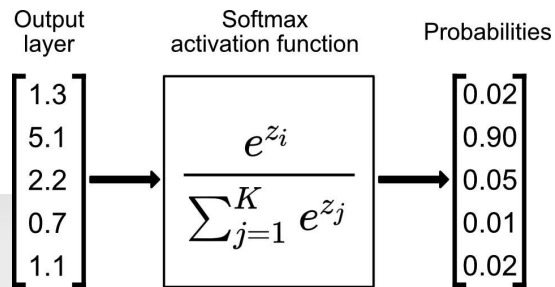
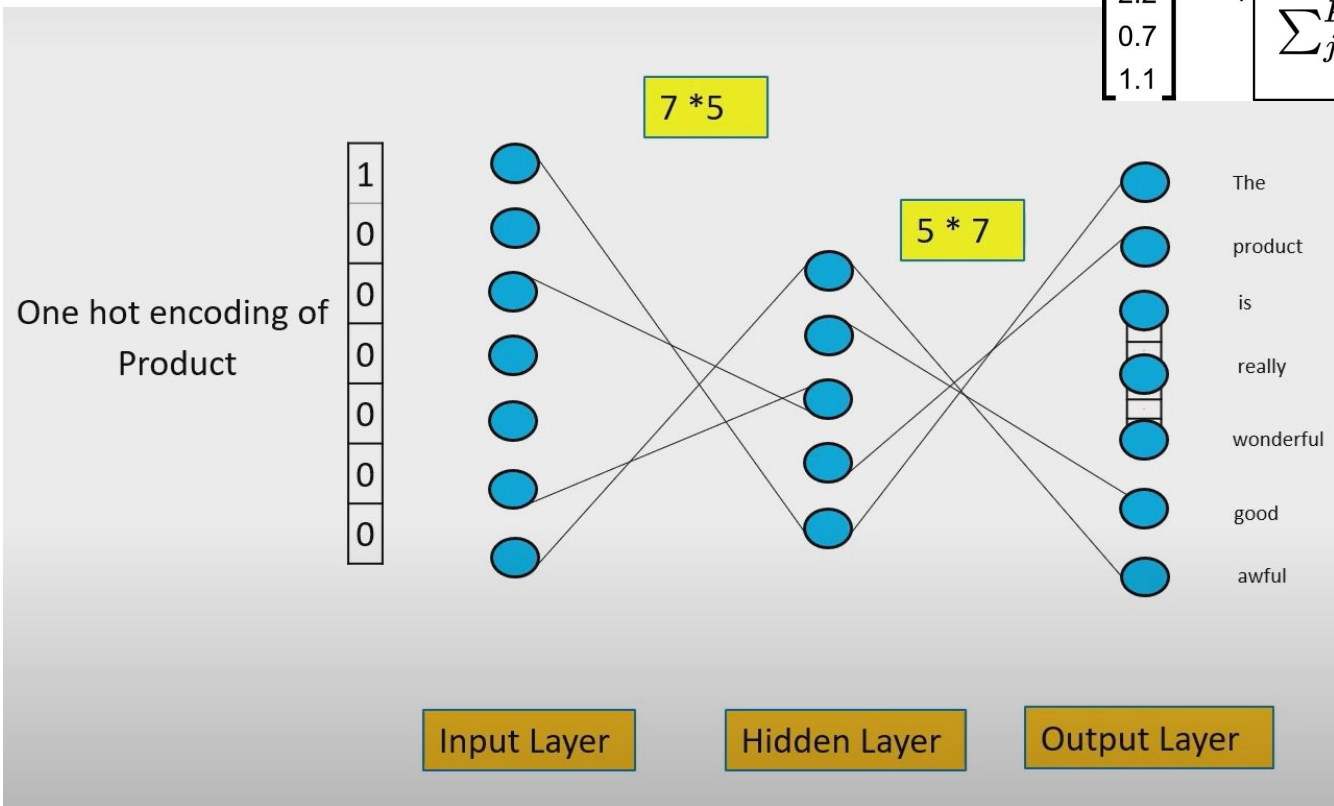
The product is really good The product is wonderful The product is awful

Step - 4

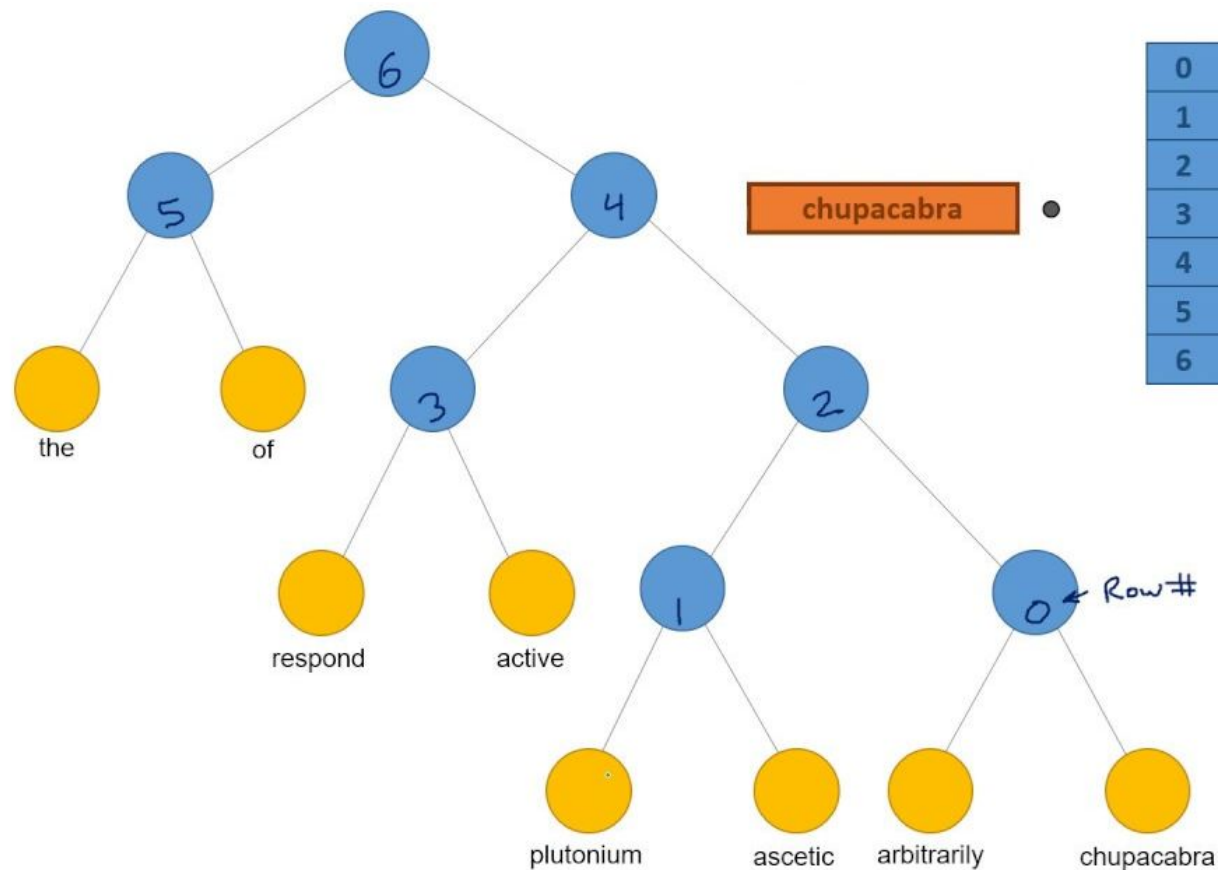
Input Words	Target Word
product	The
product	is



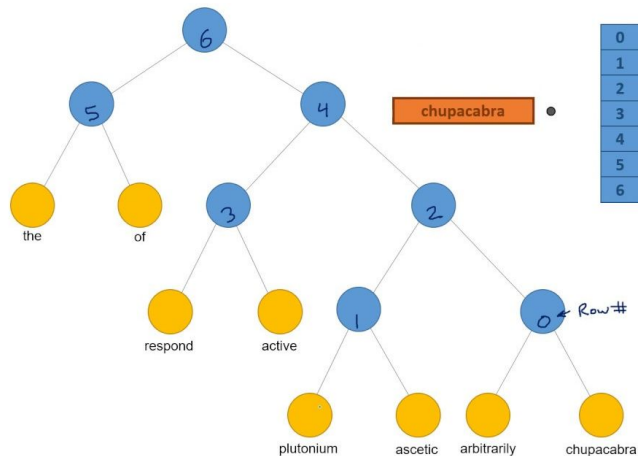
Skip-gram architecture



Hierarchical Softmax



Hierarchical Softmax



$$p(w|w_I) = \prod_{j=1}^{L(w)-1} \sigma \left(\mathbb{I}[n(w, j+1) = \text{ch}(n(w, j))] \cdot v'_{n(w, j)}{}^\top v_{w_I} \right)$$

$$\sum_{w=1}^W p(w|w_I) = 1 \quad \mathcal{O}(W^2) \rightarrow \mathcal{O}(W \log(W))$$

Negative Sampling

Computed V times
for all vocabs

$$\left\{ \begin{bmatrix} p(w_1|w^{(t)}) \\ p(w_2|w^{(t)}) \\ p(w_3|w^{(t)}) \\ \vdots \\ p(w_V|w^{(t)}) \end{bmatrix} \right\} = \frac{\exp(W_{\text{output}} \cdot h)}{\sum_{i=1}^V \exp(W_{\text{output}_{(i)}} \cdot h)} \in \mathbb{R}^V$$

Complexity = $O(V + V) \approx O(V)$
where V is very large

V computations are needed to
get normalization factor

Subsampling of Frequent Words

“in”, “the”, “a” — ?

To counter the imbalance between the rare and frequent words, we used a simple subsampling approach: each word w_i in the training set is discarded with probability computed by the formula

$$P(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}} \quad (5)$$

Learning Phrases

$$\text{score}(w_i, w_j) = \frac{\text{count}(w_i w_j) - \delta}{\text{count}(w_i) \times \text{count}(w_j)}$$

NBA Teams			
Detroit Oakland	Detroit Pistons Golden State Warriors	Toronto Memphis	Toronto Raptors Memphis Grizzlies
Airlines			
Austria Belgium	Austrian Airlines Brussels Airlines	Spain Greece	Spainair Aegean Airlines
Company executives			
Steve Ballmer Samuel J. Palmisano	Microsoft IBM	Larry Page Werner Vogels	Google Amazon

Phrase Skip-Gram Results

Czech + currency	Vietnam + capital	German + airlines	Russian + river	French + actress
koruna <u>Check crown</u> Polish zolty CTK	Hanoi Ho Chi Minh City Viet Nam Vietnamese	airline Lufthansa carrier Lufthansa flag carrier Lufthansa Lufthansa	Moscow Volga River upriver Russia	Juliette Binoche Vanessa Paradis Charlotte Gainsbourg Cecile De

Table 5: Vector compositionality using element-wise addition. Four closest tokens to the sum of two vectors are shown, using the best Skip-gram model.

Comparison to Published Word Representations

Model (training time)	Redmond	Havel	ninjutsu	graffiti	capitulate
Collobert (50d) (2 months)	conyers lubbock keene	plauen dzerzhinsky osterreich	reiki kohona karate	cheesecake gossip dioramas	abdicate accede rearm
Turian (200d) (few weeks)	McCarthy Alston Cousins	Jewell Arzu Ovitz	- - -	gunfire emotion impunity	- - -
Mnih (100d) (7 days)	Podhurst Harlang Agarwal	Pontiff Pinochet Rodionov	- - -	anaesthetics monkeys Jews	Mavericks planning hesitated
Skip-Phrase (1000d, 1 day)	Redmond Wash. Redmond Washington Microsoft	Vaclav Havel president Vaclav Havel Velvet Revolution	ninja martial arts swordsmanship	spray paint grafitti taggers	capitulation capitulated capitulating

Table 6: Examples of the closest tokens given various well known models and the Skip-gram model trained on phrases using over 30 billion training words. An empty cell means that the word was not in the vocabulary.