

Natural Language Processing Language Models

Sudeshna Sarkar

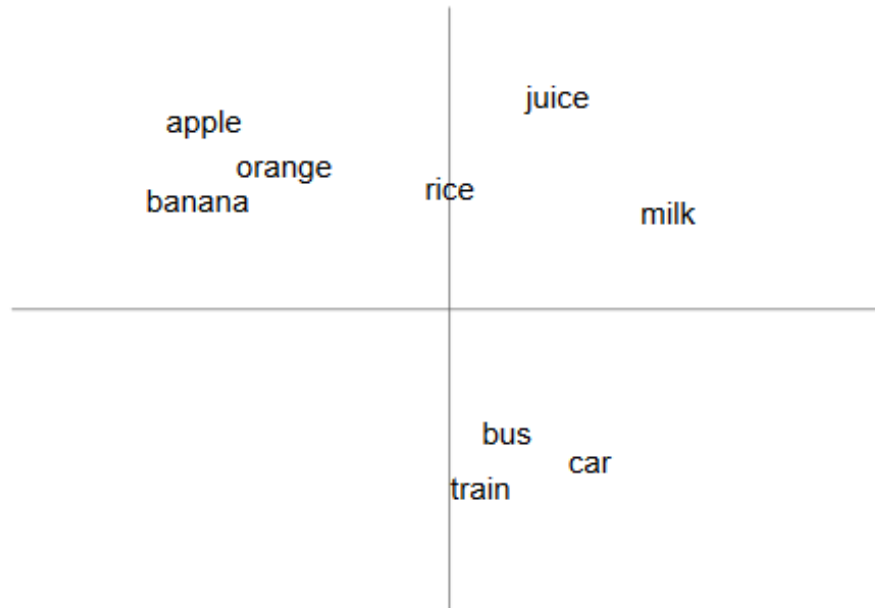
24 July 2019

Neural language model

$$P(w_t | w_{t-n}, \dots, w_{t-1}) = \frac{C(w_{t-n}, \dots, w_t)}{C(w_{t-n}, \dots, w_{t-1})}$$
$$= f_{\theta}(w_{t-n}, \dots, w_{t-1})$$

- Parametric estimator
- We need numerical representation of words

“A word is known by the company it keeps”



Word Representation

- Continuous Representation: based on context
- Distributional hypothesis

You can get a lot of value by representing a word by means of its neighbors

“You shall know a word by the company it keeps”

(J. R. Firth 1957: 11)

One of the most successful ideas of modern NLP

government debt problems turning into banking crises as has happened in
saying that Europe needs unified banking regulation to replace the hodgepodge

↖ These words will represent *banking* ↗

Representation

- We need effective representation of :

- Words, Sentences, Text

1: Use existing thesauri or ontologies like WordNet

Drawbacks:

- Manual
- Not context specific

2: Use co-occurrences for word similarity. Drawbacks:

- Quadratic space needed
- Relative position and order of words not considered

co-occurrence

1. Use co-occurrences for word similarity.
2. Singular Value Decomposition (SVD) on co-occurrence matrix
 - \hat{X} is the best rank k approximation to X , in terms of least squares

- $m = n =$
size of vocabulary

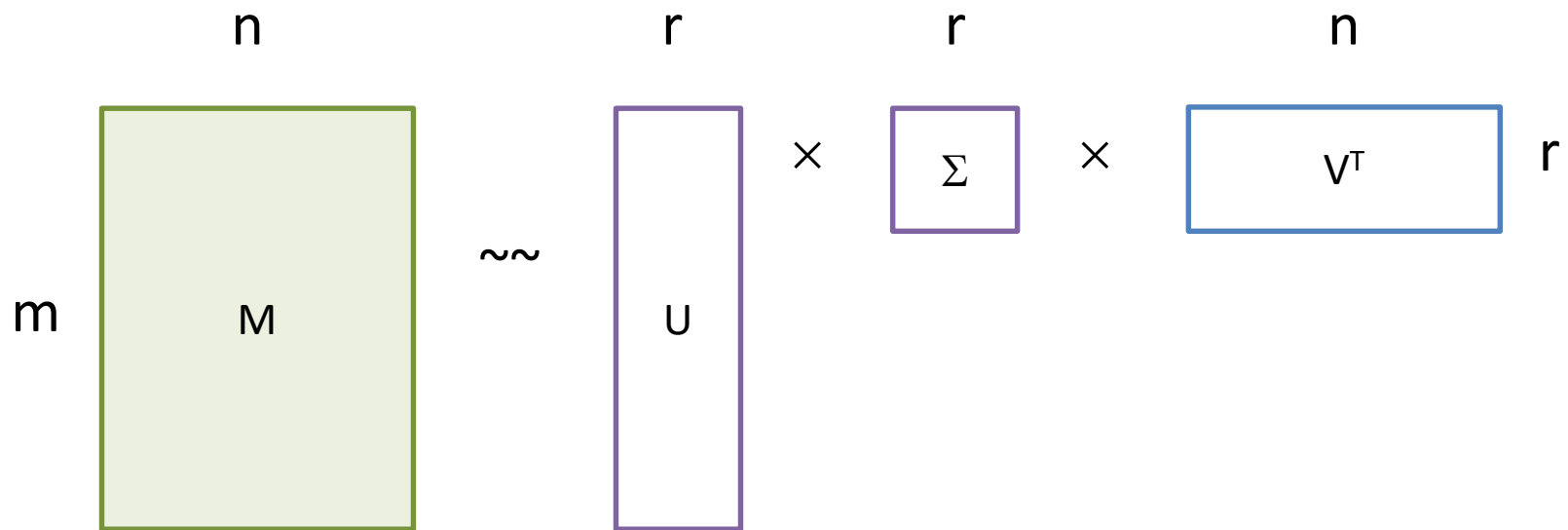
$$\begin{array}{ccccc}
 & m & & r & r & m \\
 n & \boxed{} & = & n \boxed{\begin{array}{c} | \\ | \\ | \\ | \end{array}} & r \boxed{\begin{array}{ccc} s_1 & s_2 & 0 \\ & s_3 & \ddots \\ 0 & & s_r \end{array}} & r \boxed{\begin{array}{c} \text{---} V_1 \text{---} \\ \text{---} V_2 \text{---} \\ \text{---} V_3 \text{---} \\ \vdots \end{array}} \\
 & X & & U & S & V^T \\
 \\
 & m & & k & k & m \\
 n & \boxed{} & = & n \boxed{\begin{array}{c} | \\ | \\ | \\ | \end{array}} & k \boxed{\begin{array}{ccc} s_1 & s_2 & 0 \\ & s_3 & \ddots \\ 0 & & s_k \end{array}} & k \boxed{\begin{array}{c} \text{---} V_1 \text{---} \\ \text{---} V_2 \text{---} \\ \text{---} V_3 \text{---} \\ \vdots \end{array}} \\
 & \hat{X} & & \hat{U} & \hat{S} & \hat{V}^T
 \end{array}$$

- \hat{S} is the same matrix as S except that it contains only the top largest singular values

SVD

- A decomposition of any matrix into a product of three matrices.
- From this decomposition, you can choose any number r of intermediate concepts (latent factors) in a way that minimizes the RMSE error given that value of r .
- The rank of a matrix is the maximum number of rows (or equivalently columns) that are linearly independent.
- If a matrix has rank r , then it can be decomposed exactly into matrices whose shared dimension is r .
- Vectors are orthogonal if their dot product is 0.
- An orthonormal basis is a set of unit vectors any two of which are orthogonal.

Form of SVD



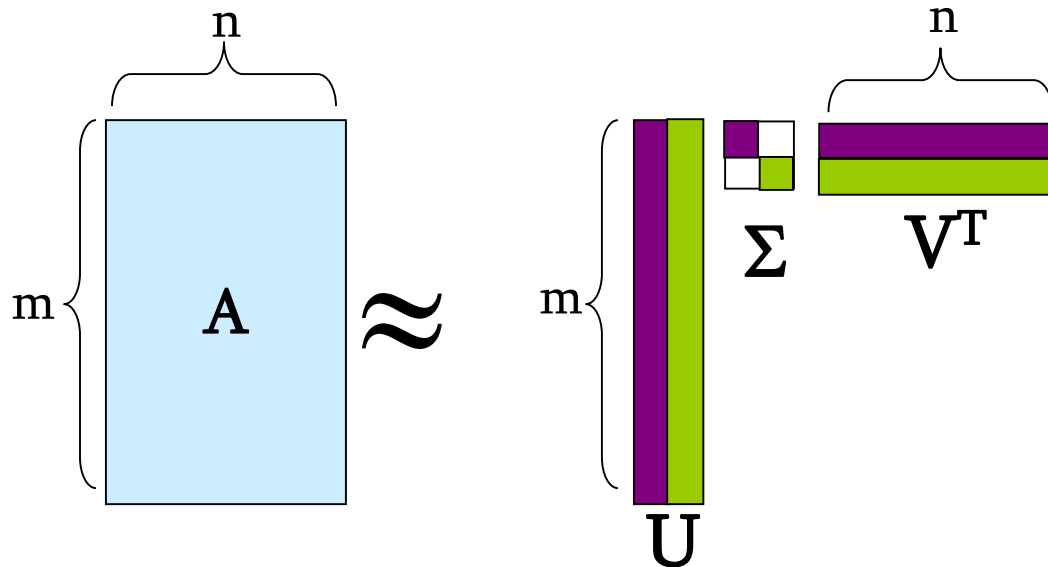
Special conditions:

U and V are column-orthonormal
(so V^T has orthonormal rows)
 Σ is a diagonal matrix
The values of Σ along the diagonal
are called the *singular values*.

- It is possible to decompose M **exactly**, if r is the rank of M .
- But we may make r much smaller, by setting to 0 the smallest singular values.
 - making the corresponding columns of U and V useless

Linkage Among Components of U , V , Σ

$$\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_i \sigma_i \mathbf{u}_i \circ \mathbf{v}_i^T$$



word2vec approach to represent the meaning of word

- Represent each word with a low-dimensional vector
- Word similarity = vector similarity
- Key idea: Predict surrounding words of every word
- Faster and can easily incorporate a new sentence/document or add a word to the vocabulary

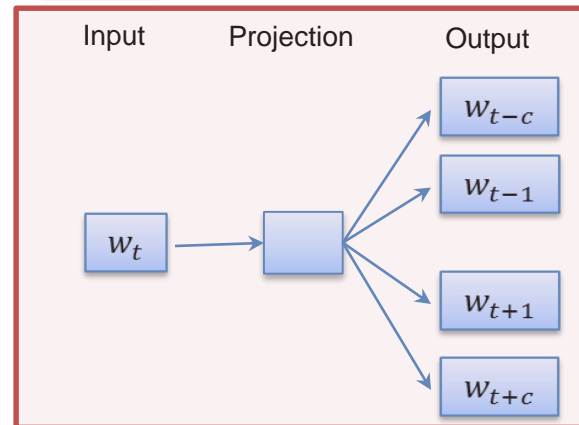
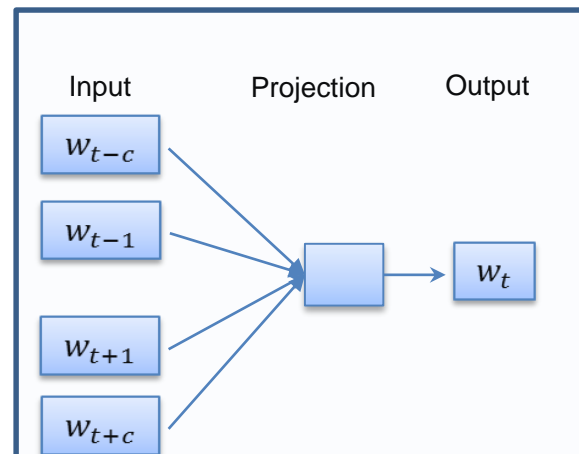
Word2vec

- Representation of words

“Similar words have similar contexts”

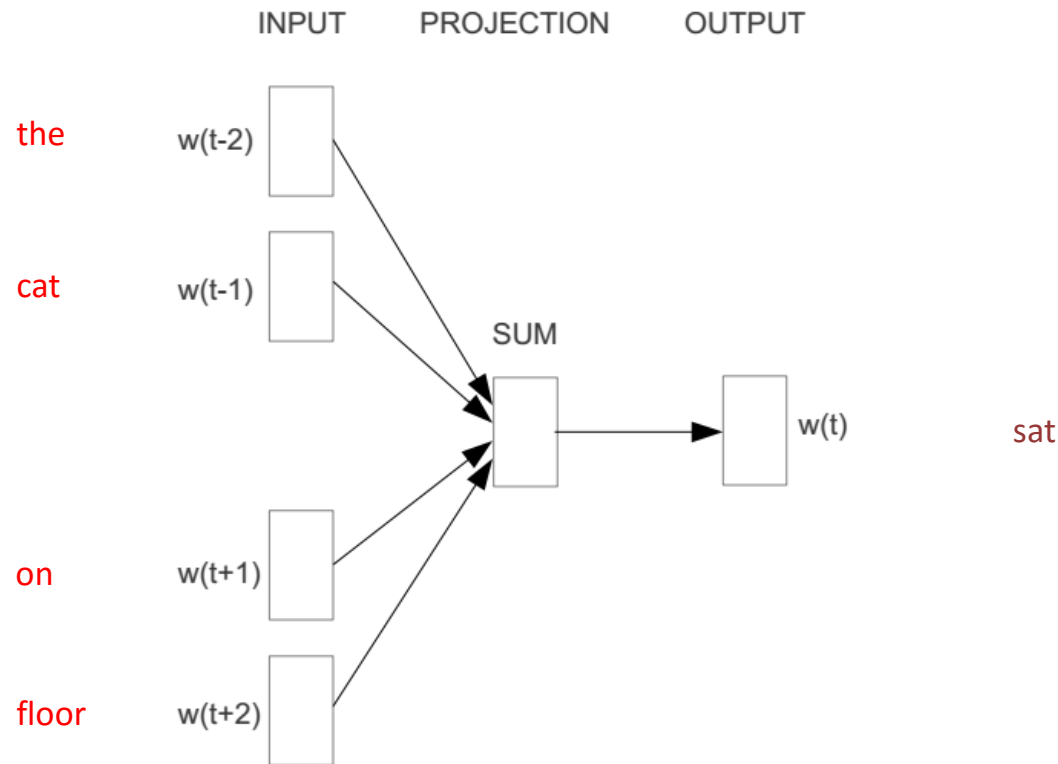
1. CBOW: $P(\text{Word}|\text{Context})$

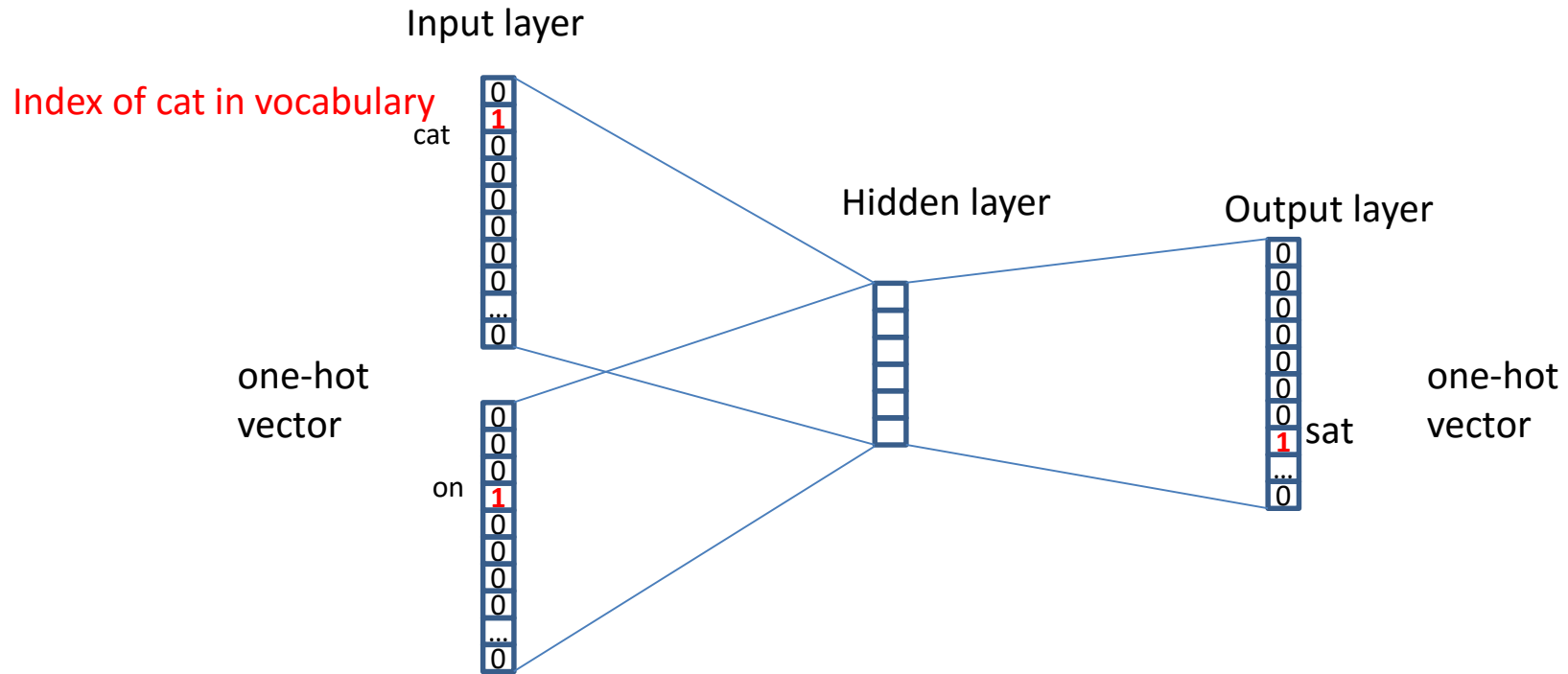
2. Skipgram: $P(\text{Context}|\text{Word})$

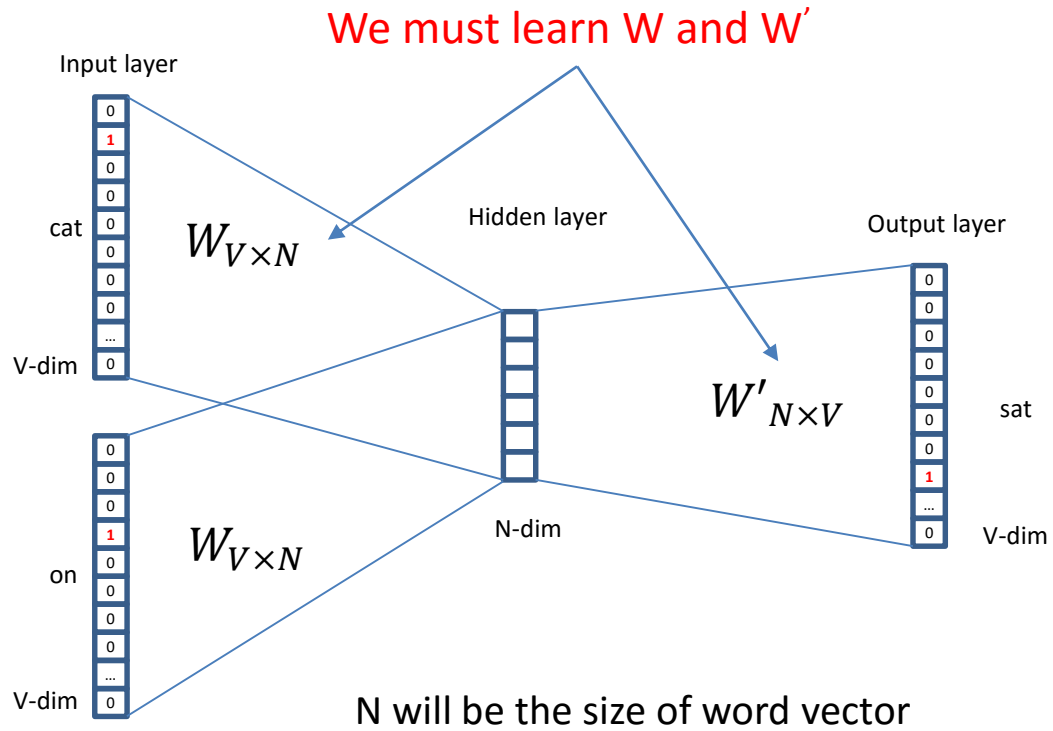


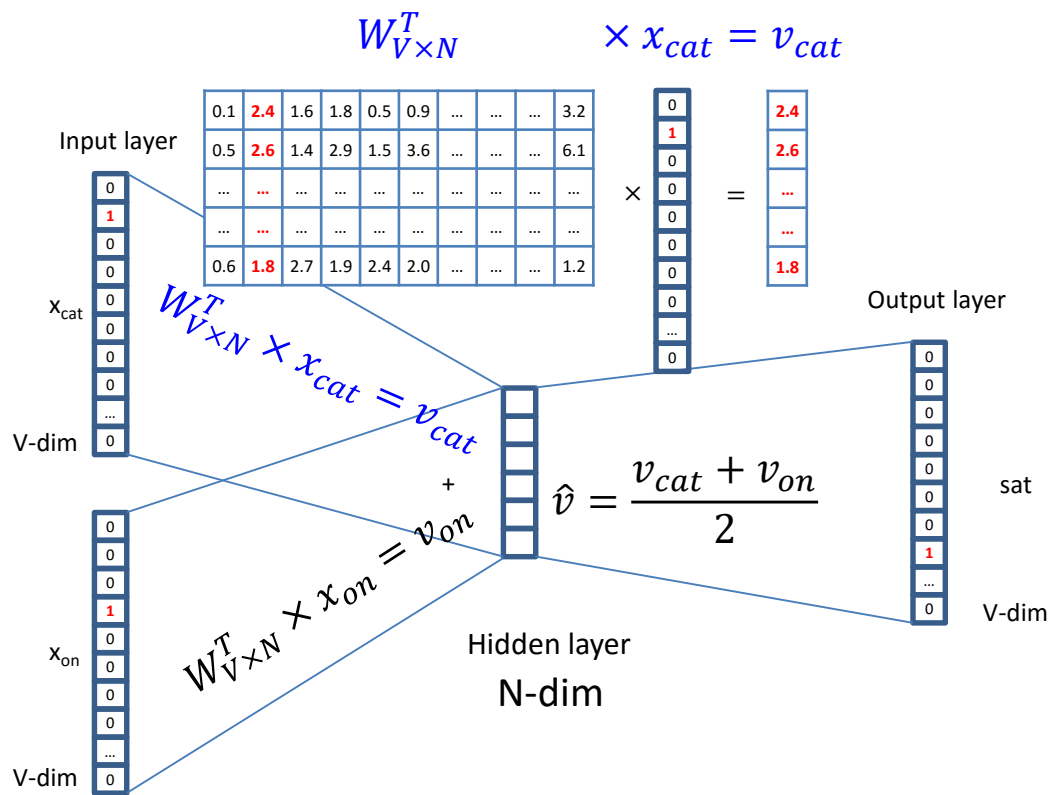
Word2vec – Continuous Bag of Word

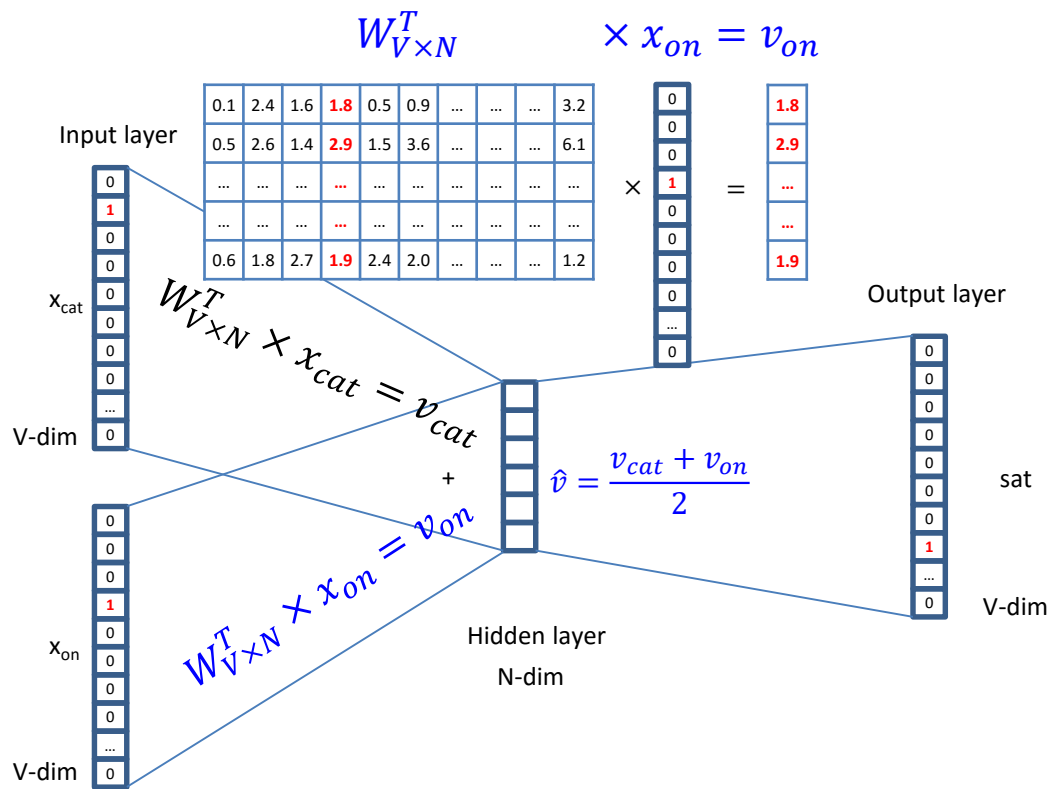
- E.g. “The cat sat on floor”
 - Window size = 2

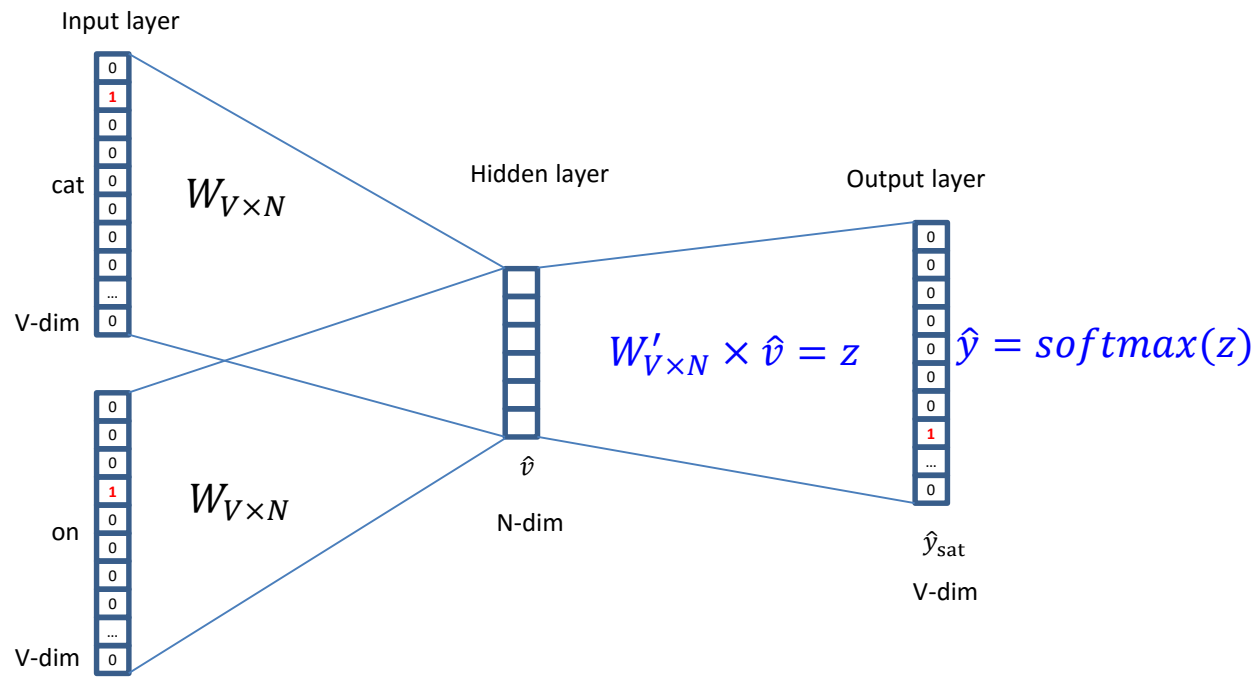


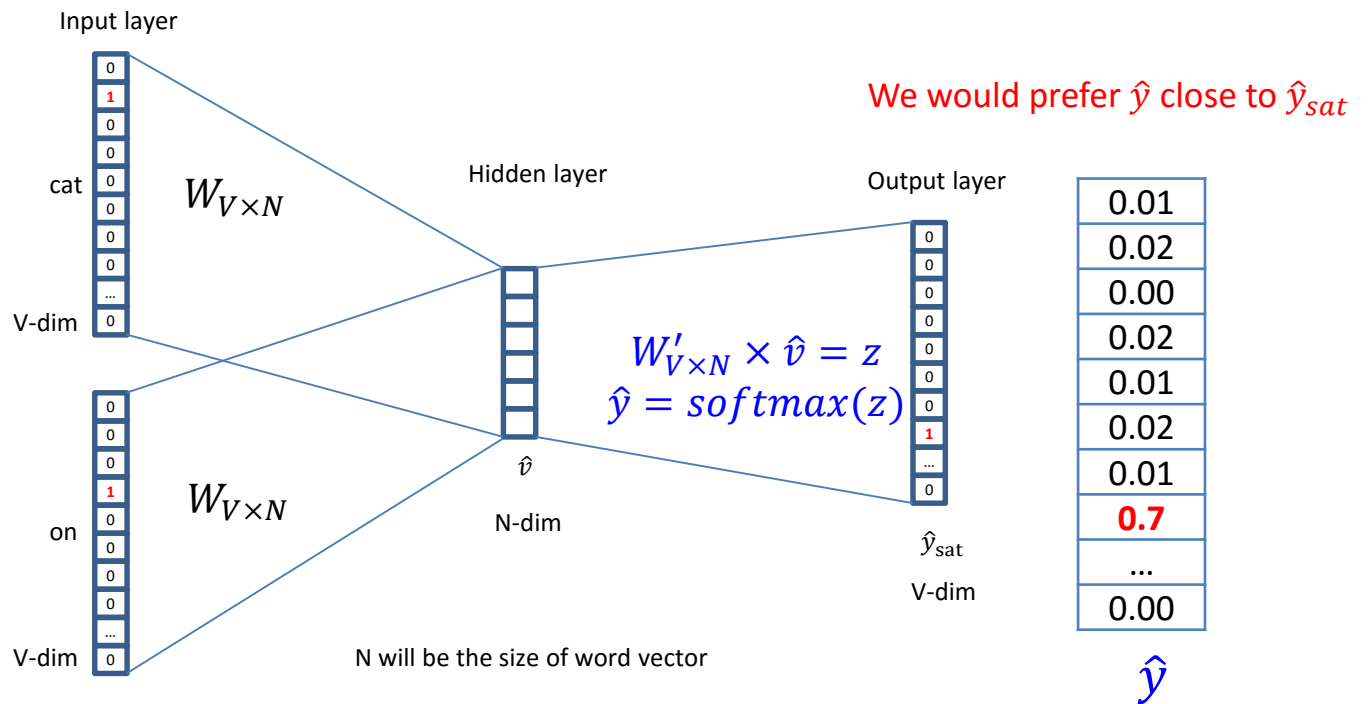


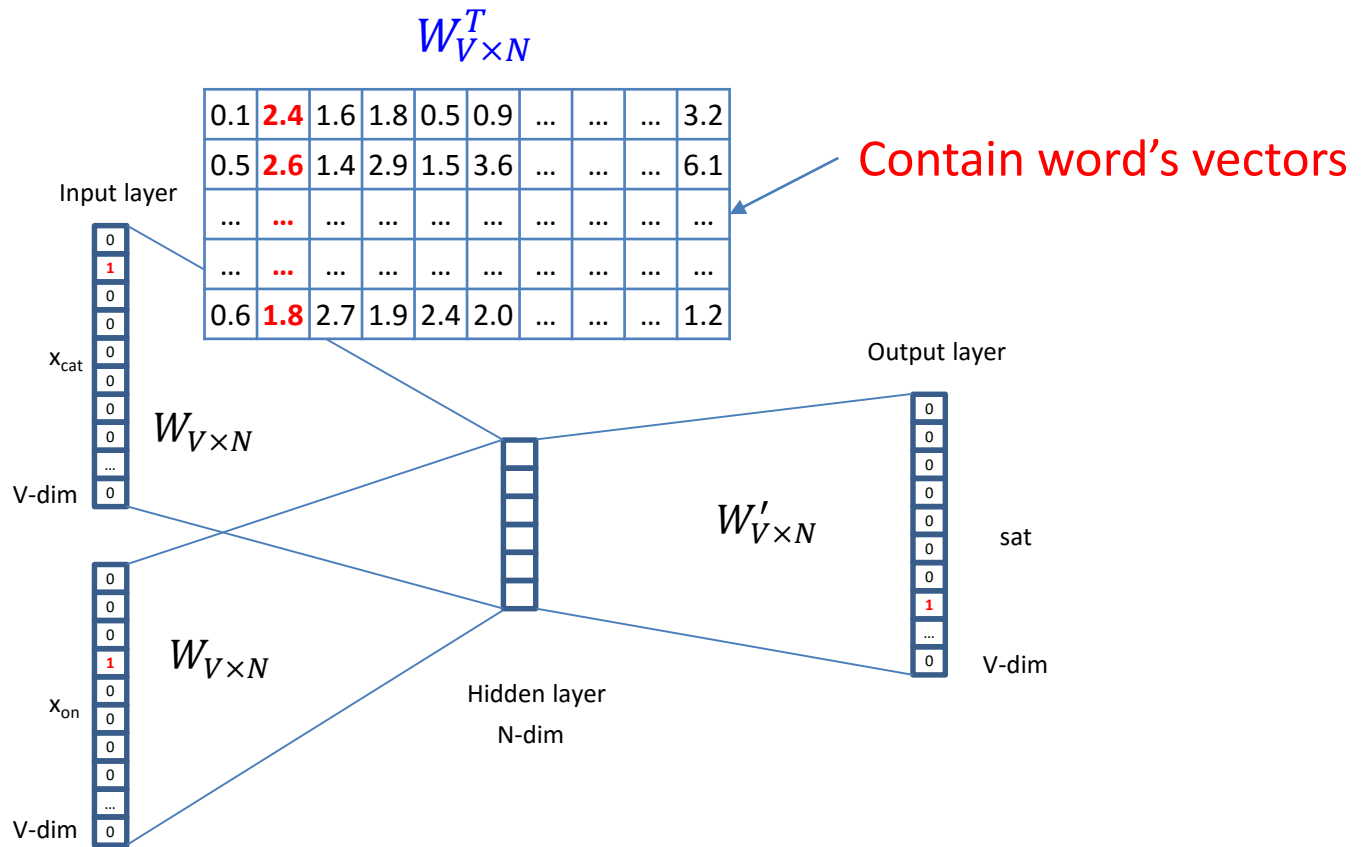












We can consider either W or W' as the word's representation. Or even take the average.

Some interesting results

Word Analogies

Test for linear relationships, examined by Mikolov et al. (2014)

a:b :: c:?



$$d = \arg \max_x \frac{(w_b - w_a + w_c)^T w_x}{||w_b - w_a + w_c||}$$

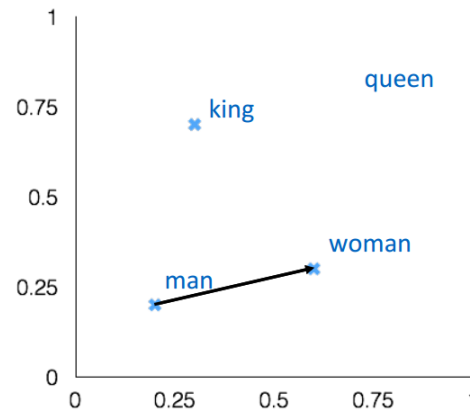
man:woman :: king:?

+ king [0.30 0.70]

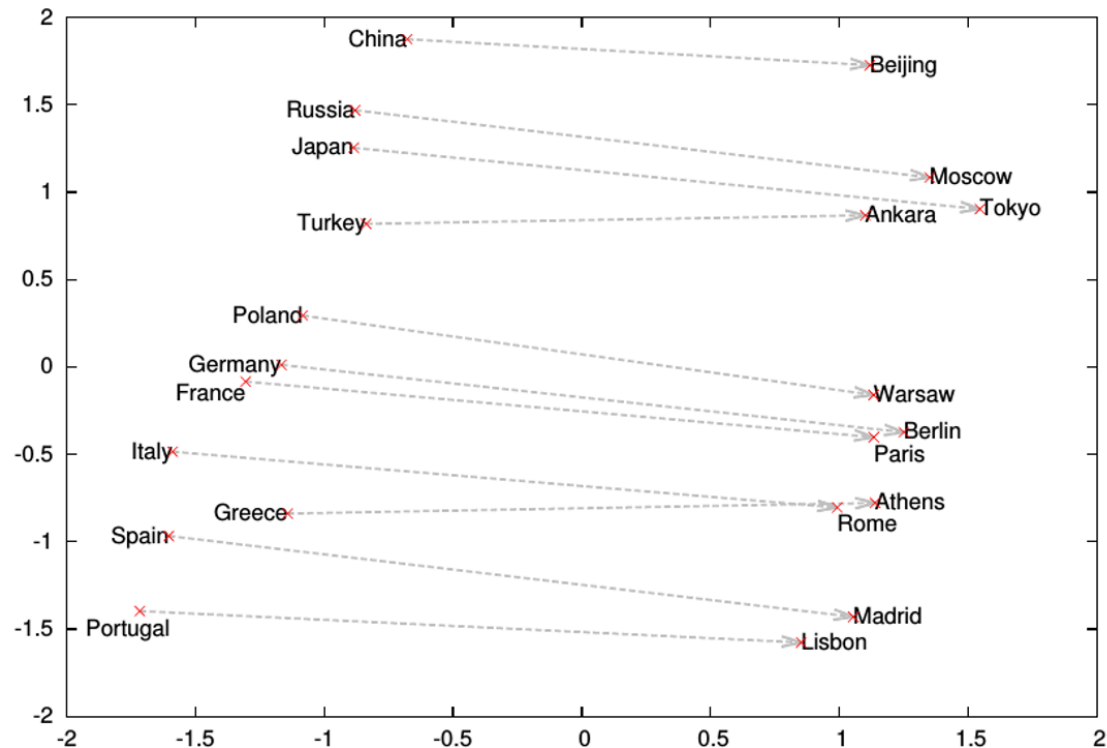
- man [0.20 0.20]

+ woman [0.60 0.30]

queen [0.70 0.80]



Word analogies



Word2Vec Objective

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log p(w_j | c)$$

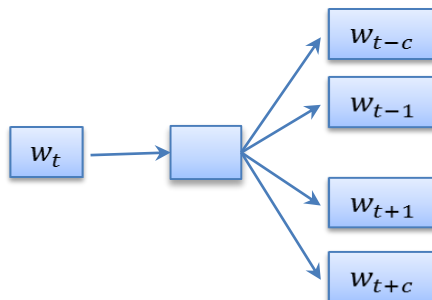
Word2Vec Objective

$$p(w_j|c) = \frac{\exp(w_j^T c)}{\sum_{i=1}^N \exp(w_i^T c)}$$

Skipgram Model

- Input: Central word w_t
- Output: Words in its context: w_{con}
 $\{w_{t-c}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+c}\}$
- Each input word represented by a 1-hot encoding of size V

INPUT PROJECTION OUTPUT



Source Text:

Deep Learning attempts to learn multiple levels of representation from data.

Input output pairs :

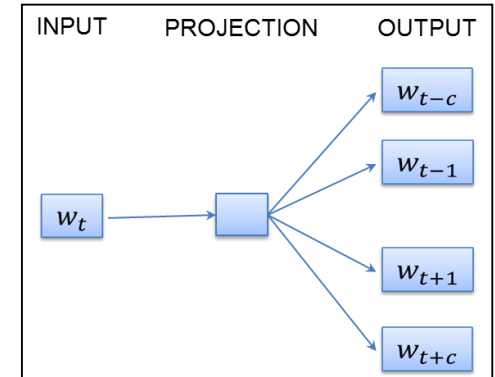
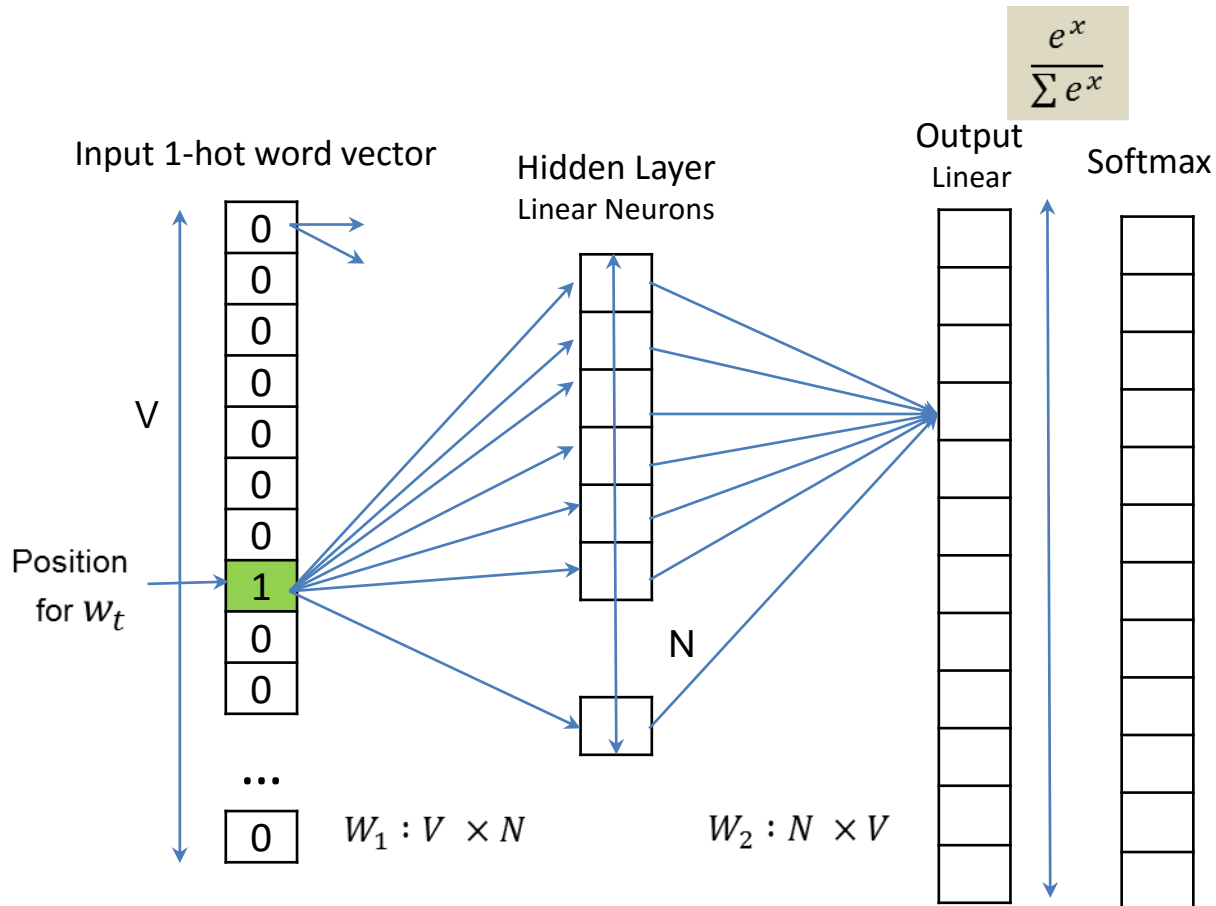
Positive samples:

(representation, levels)
(representation, of)
(representation, from)
(representation, data)

Negative samples:

(representation, x)
[x: all other words except the 4 positive]

Skipgram Model



Probability that the word in a context position is w_i

Positive sampling
Negative sampling

Skipgram: Loss function

- Maximize
$$\frac{1}{T} \sum_{t=1}^T \sum_{\text{context}} \log p(w_{\text{context}} | w_t)$$

$p(w_{\text{con}} | w_t)$ is the output of softmax classifier

$$p(w_{\text{con}} | w_t) = \frac{\exp(v'_{w_{\text{con}}} \cdot v_{w_t})}{\sum_{w=1}^W \exp(v'_w \cdot v_{w_t})}$$

Let the model parameters be θ . The solution is given by

$$\begin{aligned} & \underset{\theta}{\operatorname{argmax}} \sum_{(w_t, w_c) \in D} \log p(w_c | w_t; \theta) \\ &= \sum_{(w_t, w_c) \in D} \left(\log e^{(v'_{w_c} \cdot v_{w_t})} - \log \sum_x e^{(v'_x \cdot v_{w_t})} \right) \end{aligned}$$

Time $O(V)$

V : vocabulary size

Improve Efficiency

1. Hierarchical softmax:

$O(\log V)$

Skipgram: Loss function

- Maximize $\frac{1}{T} \sum_{t=1}^T \sum_{\text{context}} \log p(w_{\text{context}} | w_t)$

$p(w_{\text{con}} | w_t)$ is the output of softmax classifier

$$p(w_{\text{con}} | w_t) = \frac{\exp(v'_{w_{\text{con}}} \cdot v_{w_t})}{\sum_{w=1}^W \exp(v'_w \cdot v_{w_t})}$$

Let the model parameters be θ . The solution is given by

$$\begin{aligned} & \underset{\theta}{\operatorname{argmax}} \sum_{(w_t, w_c) \in D} \log p(w_{\text{con}} | w_t; \theta) \\ &= \sum_{(w_t, w_c) \in D} \left(\log e^{(v'_{w_{\text{con}}} \cdot v_{w_t})} - \log \sum_x e^{(v'_x \cdot v_{w_t})} \right) \end{aligned}$$

2. Negative sampling: Sample instead of taking all contexts into account

$$\sum_{(w_t, w_c) \in D} \left(\log \sigma(v'_{w_{\text{con}}} \cdot v_{w_t}) + \sum_{i=1}^k \mathbb{E}_{w_i \sim P_n(w)} [\log \sigma(-v'_{w_i} \cdot v_{w_t})] \right)$$

Subsampling of frequent words

Skip-Grams with Negative Sampling (SGNS)

Marco saw a furry little **wampimuk** hiding in the tree.

“word2vec Explained...”
Goldberg & Levy, arXiv 2014

Skip-Grams with Negative Sampling (SGNS)


Marco saw a furry little wampimuk hiding in the tree.

words

wampimuk
wampimuk
wampimuk
wampimuk
...

contexts

furry
little
hiding
in
...



D (data)

“word2vec Explained...”
Goldberg & Levy, arXiv
2014

Skip-Grams with Negative Sampling (SGNS)

Maximize: $\sigma(\vec{w} \cdot \vec{c})$

- c was **observed** with w

words

wampimuk

wampimuk

wampimuk

wampimuk

contexts

furry

little

hiding

in

Minimize: $\sigma(\vec{w} \cdot \vec{c}')$

- c' was **hallucinated** with w

words

wampimuk

wampimuk

wampimuk

wampimuk

contexts

Australia

cyber

the

1985