

# Vector Representation of Text

Deep Learning for NLP

# To compare pieces of text

- We need effective representation of :
  - Words
  - Sentences
  - Text
- Approach 1: Use existing thesauri or ontologies like WordNet and Snomed CT (for medical). Drawbacks:
  - Manual
  - Not context specific
- Approach 2: Use co-occurrences for word similarity. Drawbacks:
  - Quadratic space needed
  - Relative position and order of words not considered

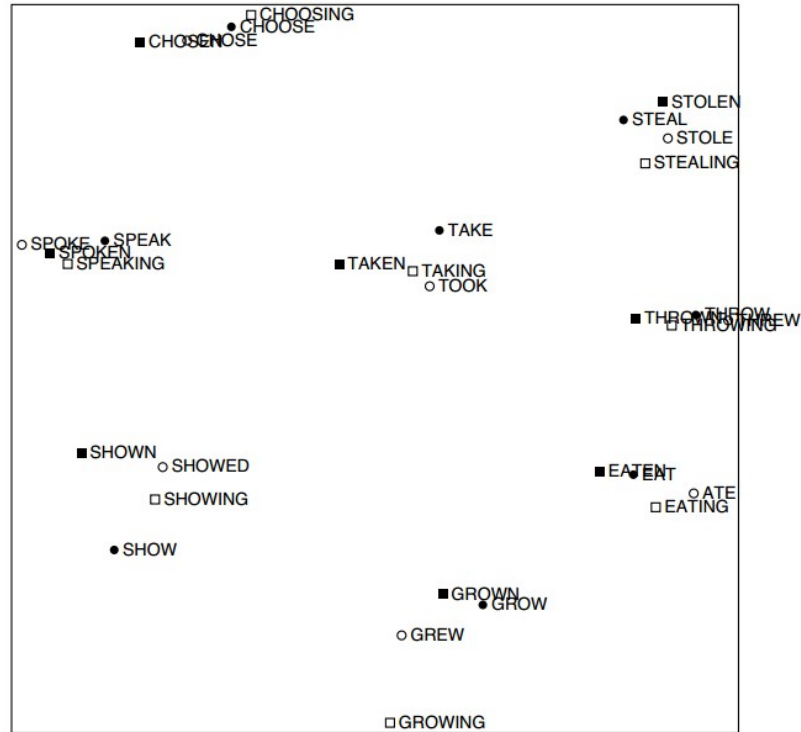
# Approach 3: low dimensional vectors

- Store only “important” information in fixed, low dimensional vector.
- Singular Value Decomposition (SVD) on co-occurrence matrix
  - is the best rank  $k$  approximation to  $X$ , in terms of least squares
  - Motel = [0.286, 0.792, -0.177, -0.107, 0.109, -0.542, 0.349, 0.271]
- $m = n$  = size of vocabulary
- is the same matrix as  $S$  except that it contains only the top largest singular values

$$\begin{array}{c}
 \begin{array}{ccccc}
 & m & & r & r & m \\
 & \boxed{\phantom{X}} & = & \boxed{\phantom{U}} & \boxed{\phantom{S}} & \boxed{\phantom{V^T}} \\
 n & & & n & r & r \\
 & X & & U & S & V^T
 \end{array} \\
 \\
 \begin{array}{ccccc}
 & m & & k & k & m \\
 & \boxed{\phantom{\hat{X}}} & = & \boxed{\phantom{\hat{U}}} & \boxed{\phantom{\hat{S}}} & \boxed{\phantom{\hat{V}^T}} \\
 n & & & n & k & k \\
 & \hat{X} & & \hat{U} & \hat{S} & \hat{V}^T
 \end{array}
 \end{array}$$

# Approach 3: low dimensional vectors

- An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence



# Problems with SVD

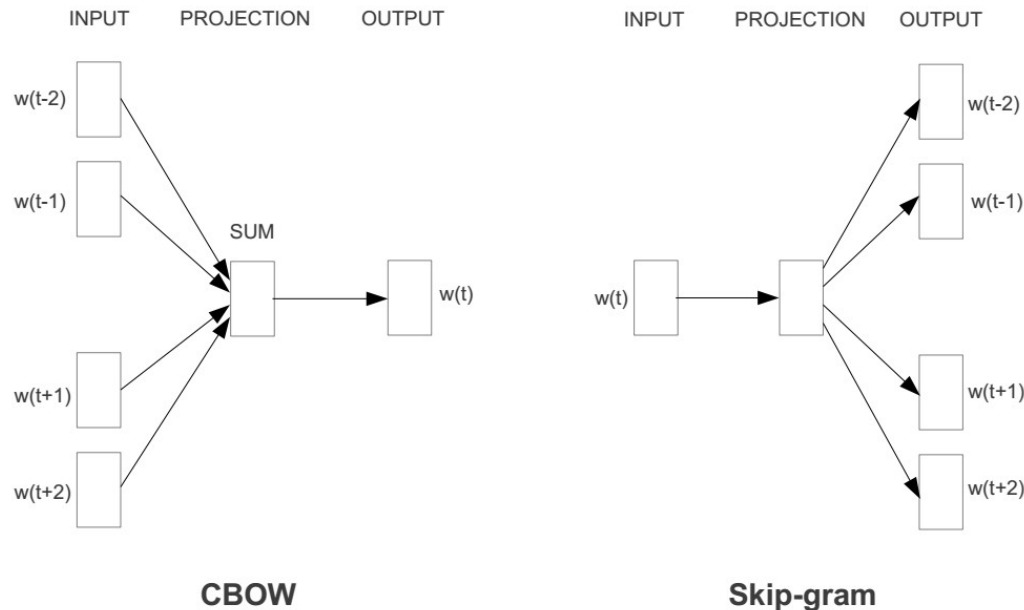
- Computational cost scales quadratically for  $n \times m$  matrix:  $O(mn^2)$  flops (when  $n < m$ )
- Hard to incorporate new words or documents
- Does not consider order of words

# **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

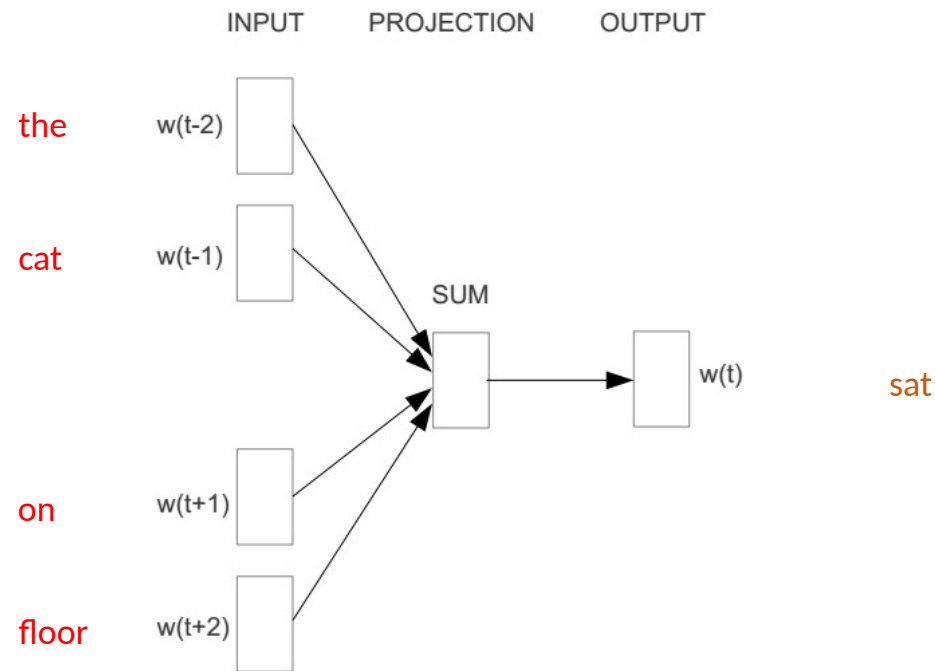
# Represent the meaning of word – word2vec

- 2 basic neural network models:
  - Continuous Bag of Word (CBOW): use a window of word to predict the middle word
  - Skip-gram (SG): use a word to predict the surrounding ones in window.

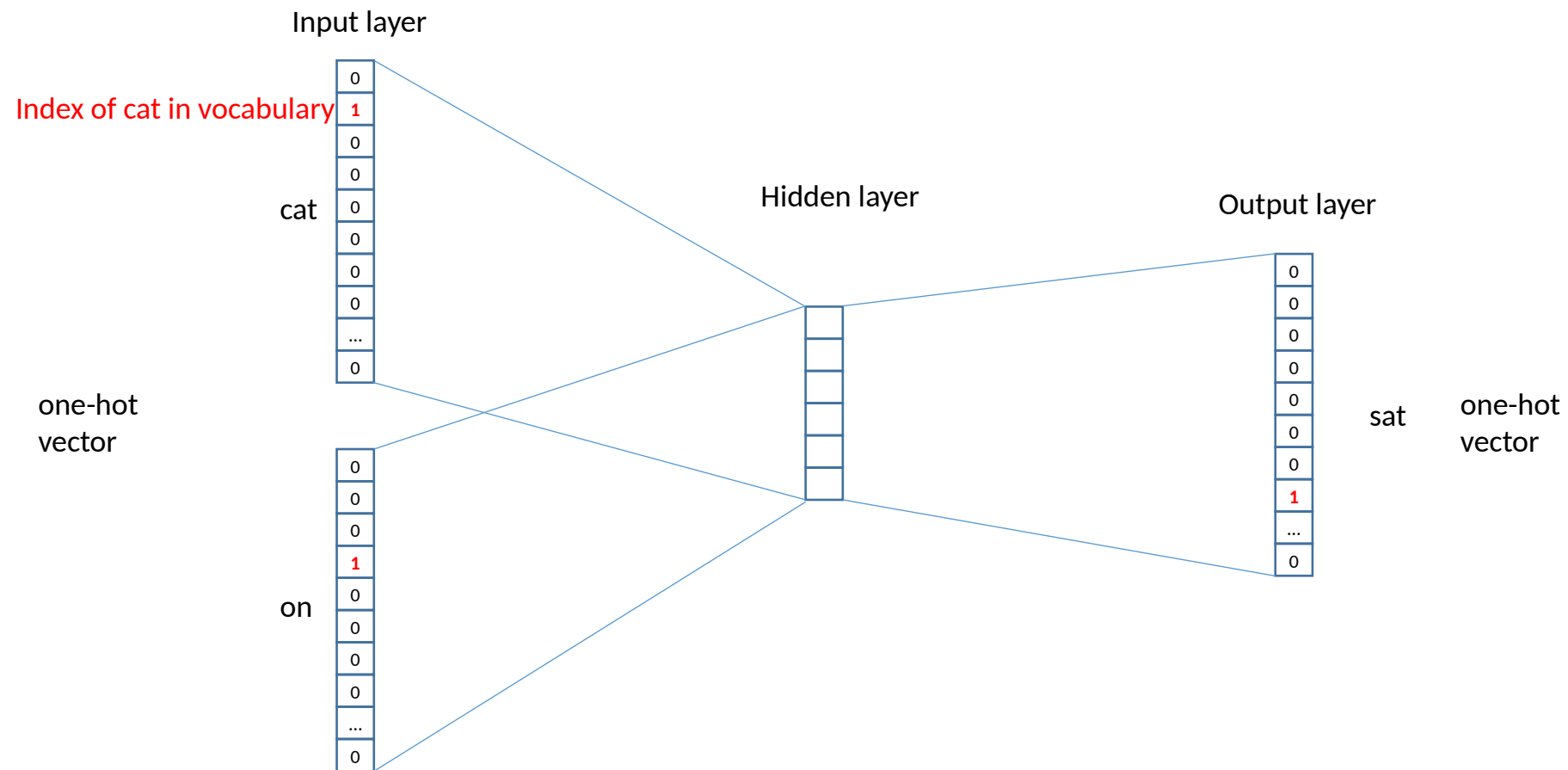


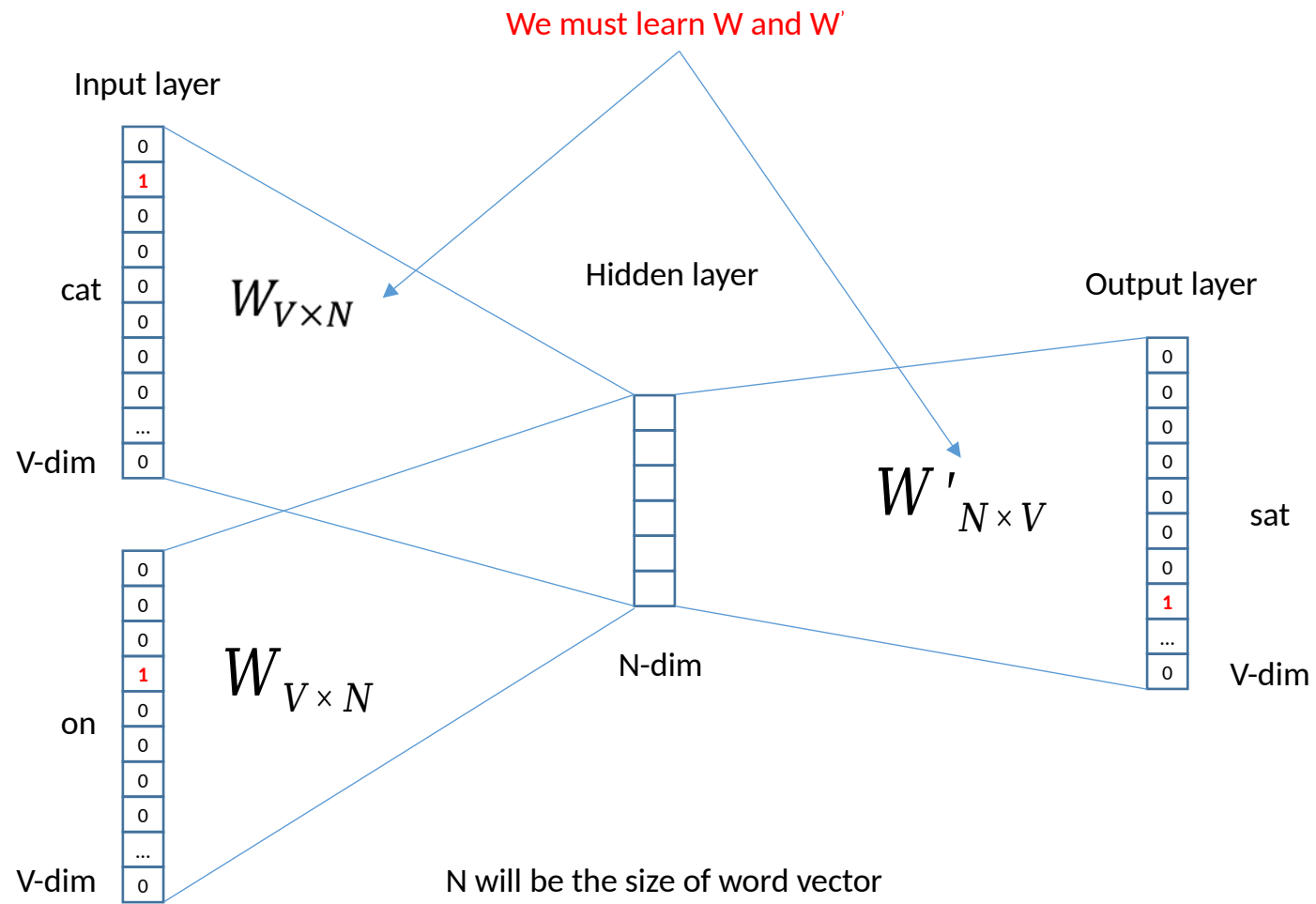
# Word2vec – Continuous Bag of Word

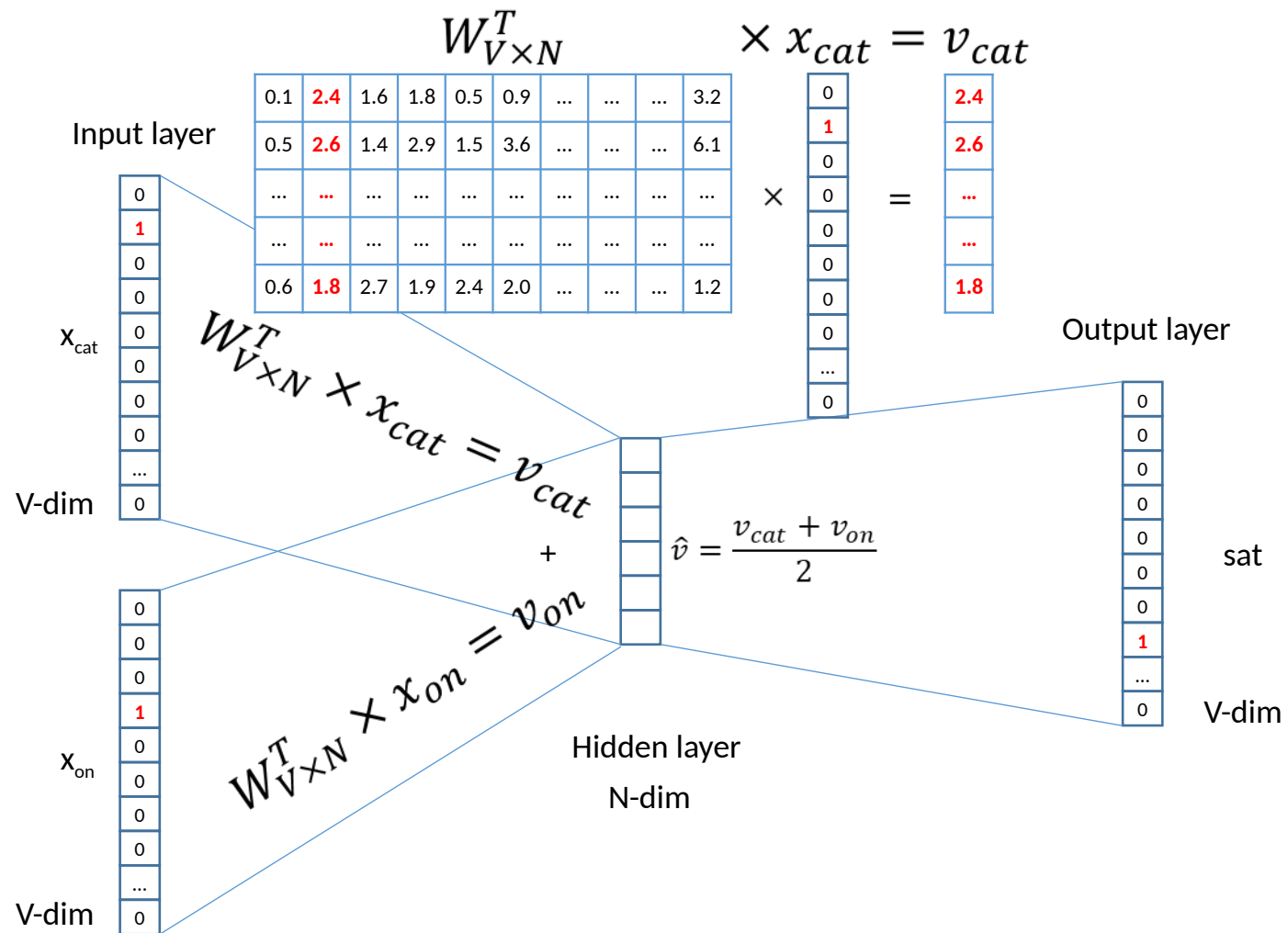
- E.g. “The cat sat on floor”
  - Window size = 2

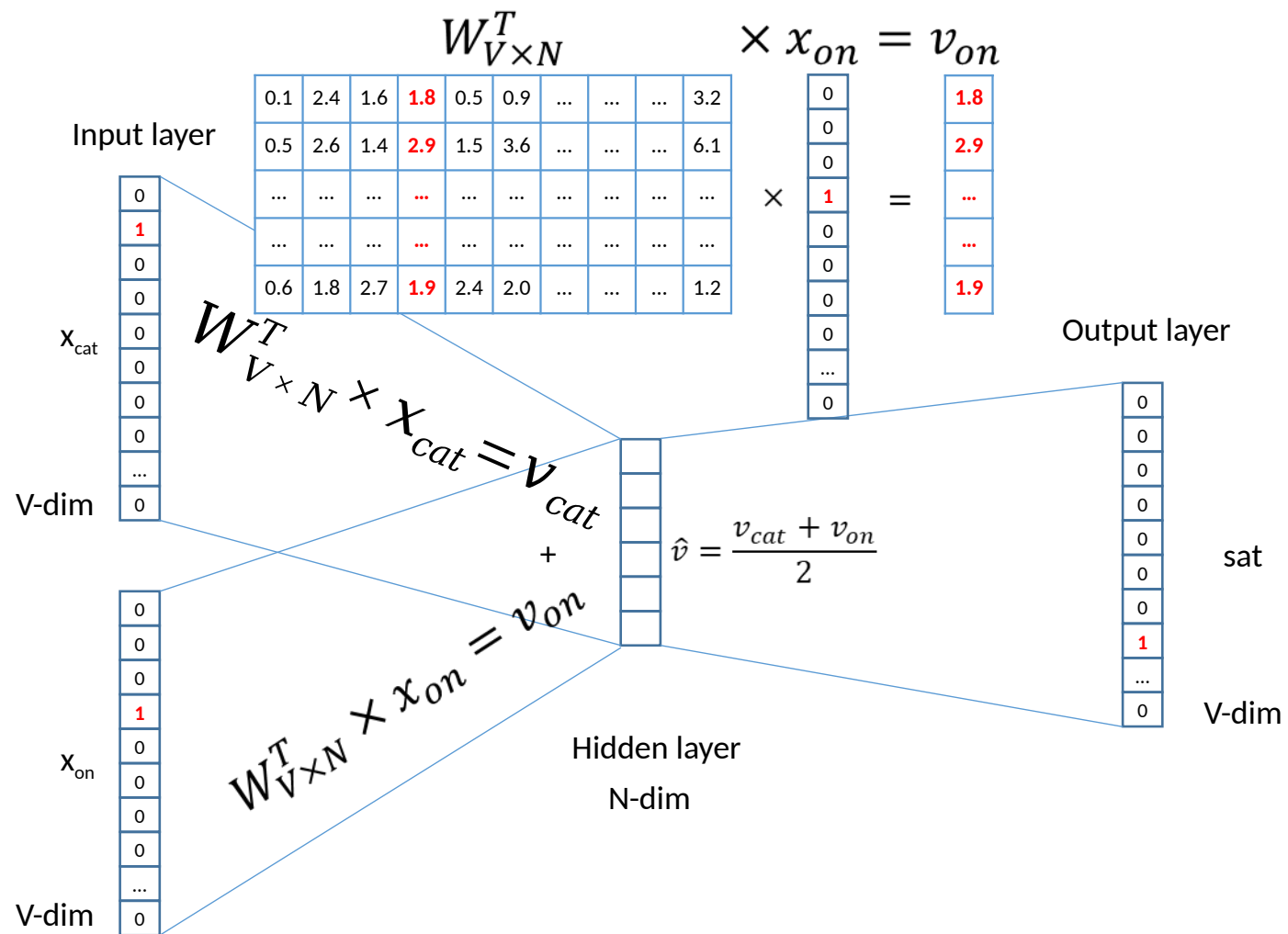


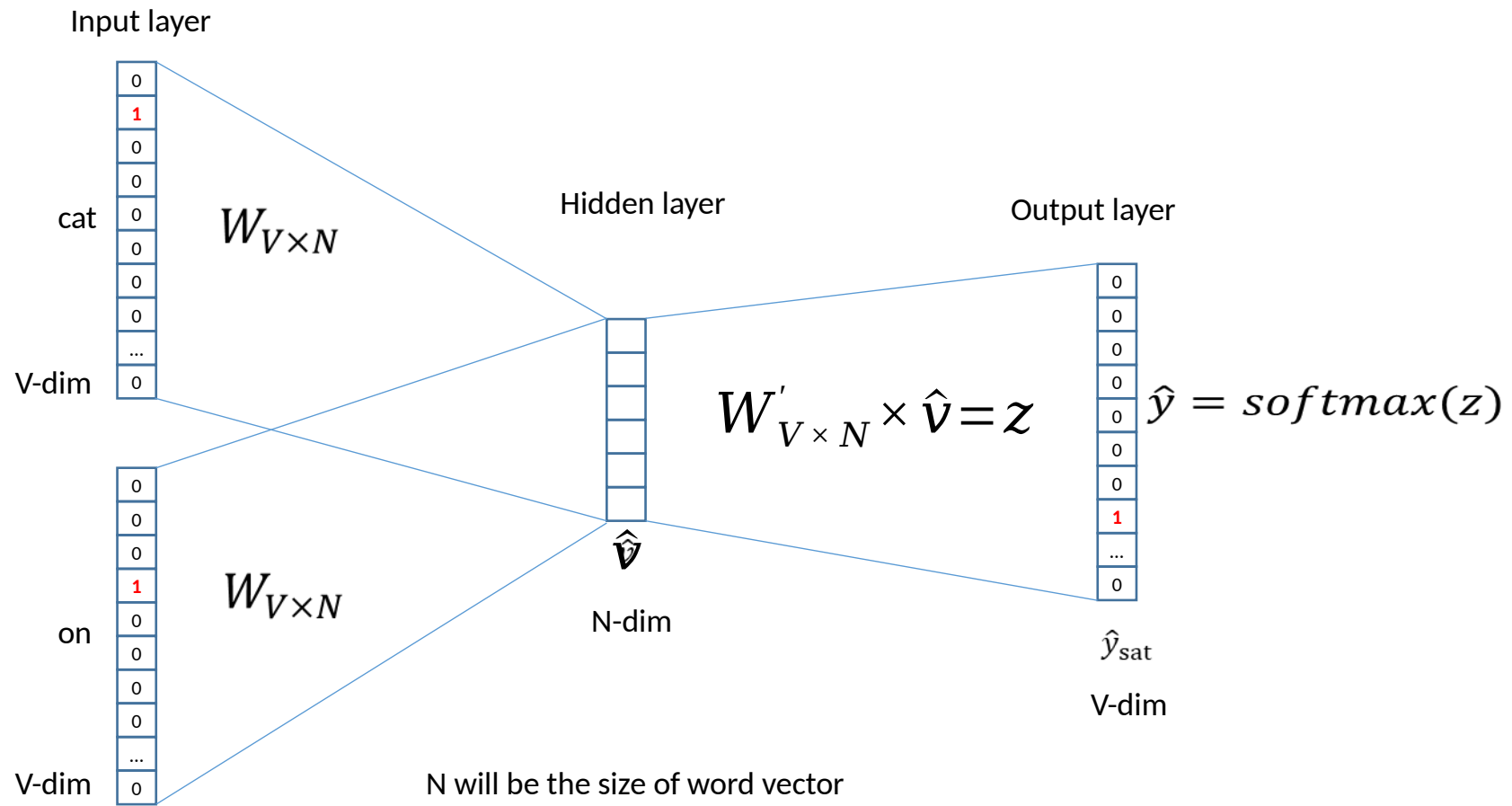


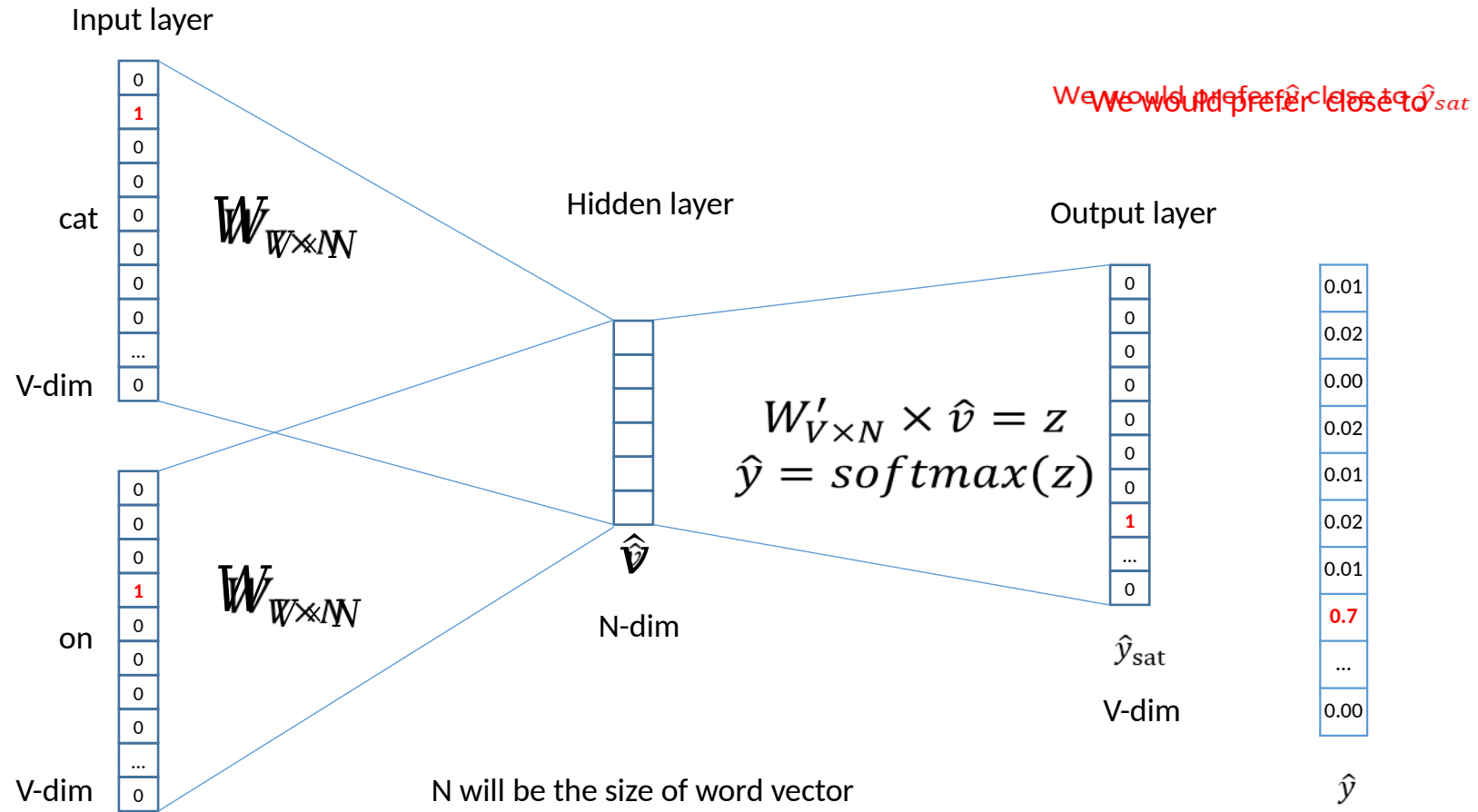


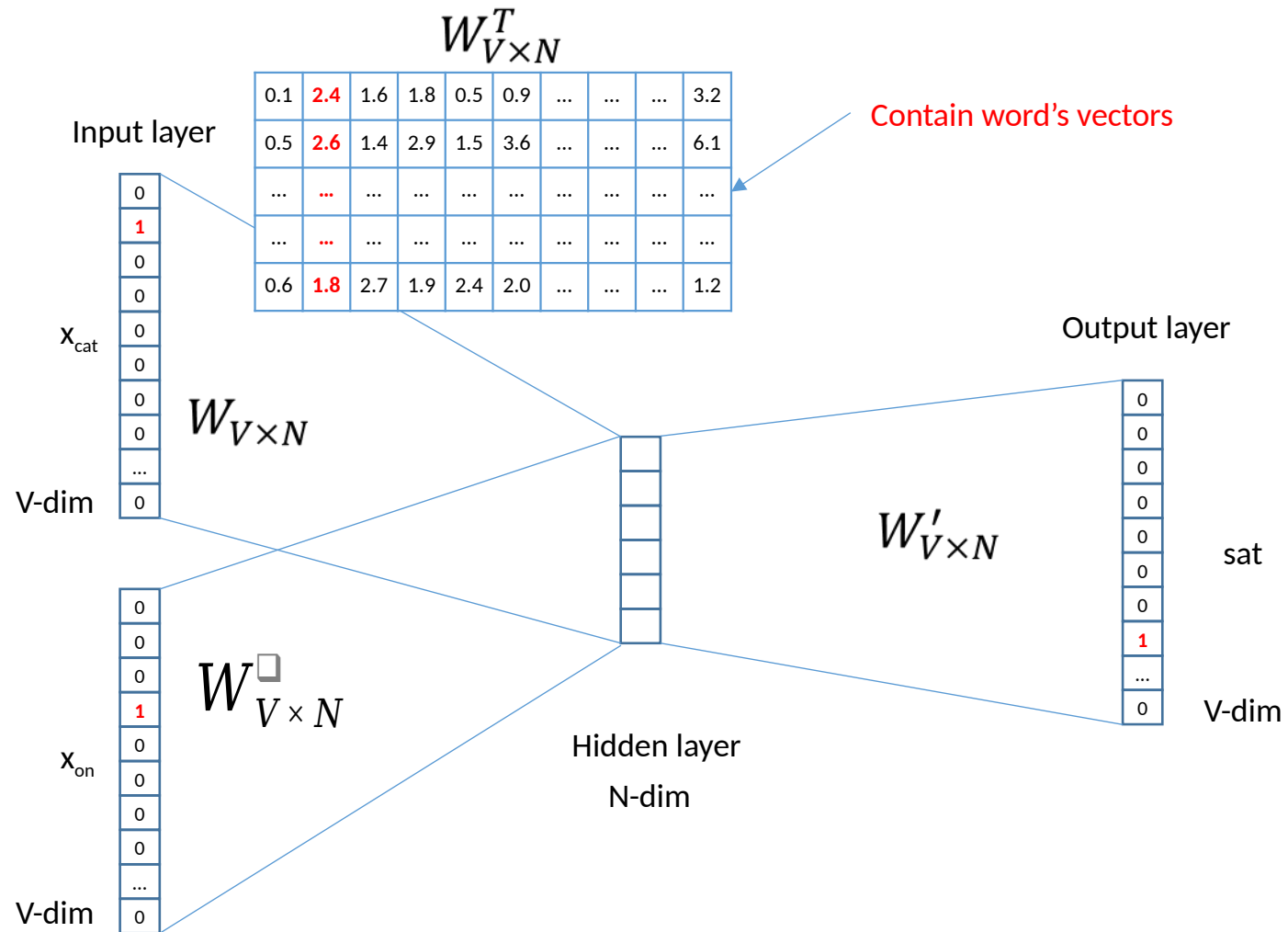












We can consider either  $W$  or  $W'$  as the word's representation. Or even take the average.

# 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:?

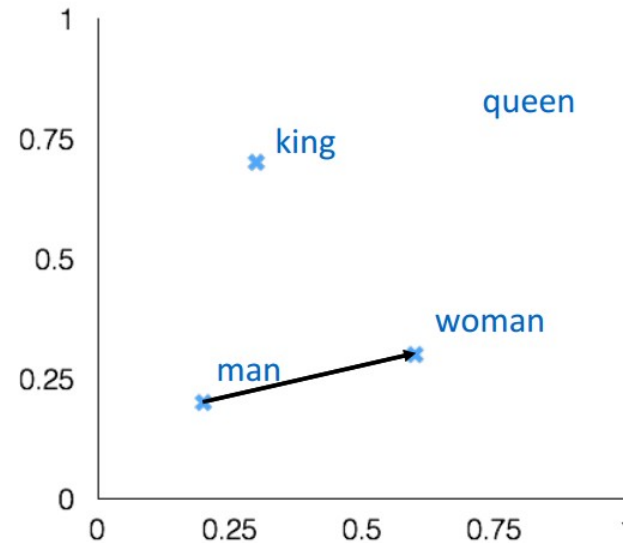
+ king [ 0.30 0.70 ]

- man [ 0.20 0.20 ]

+ woman [ 0.60 0.30 ]

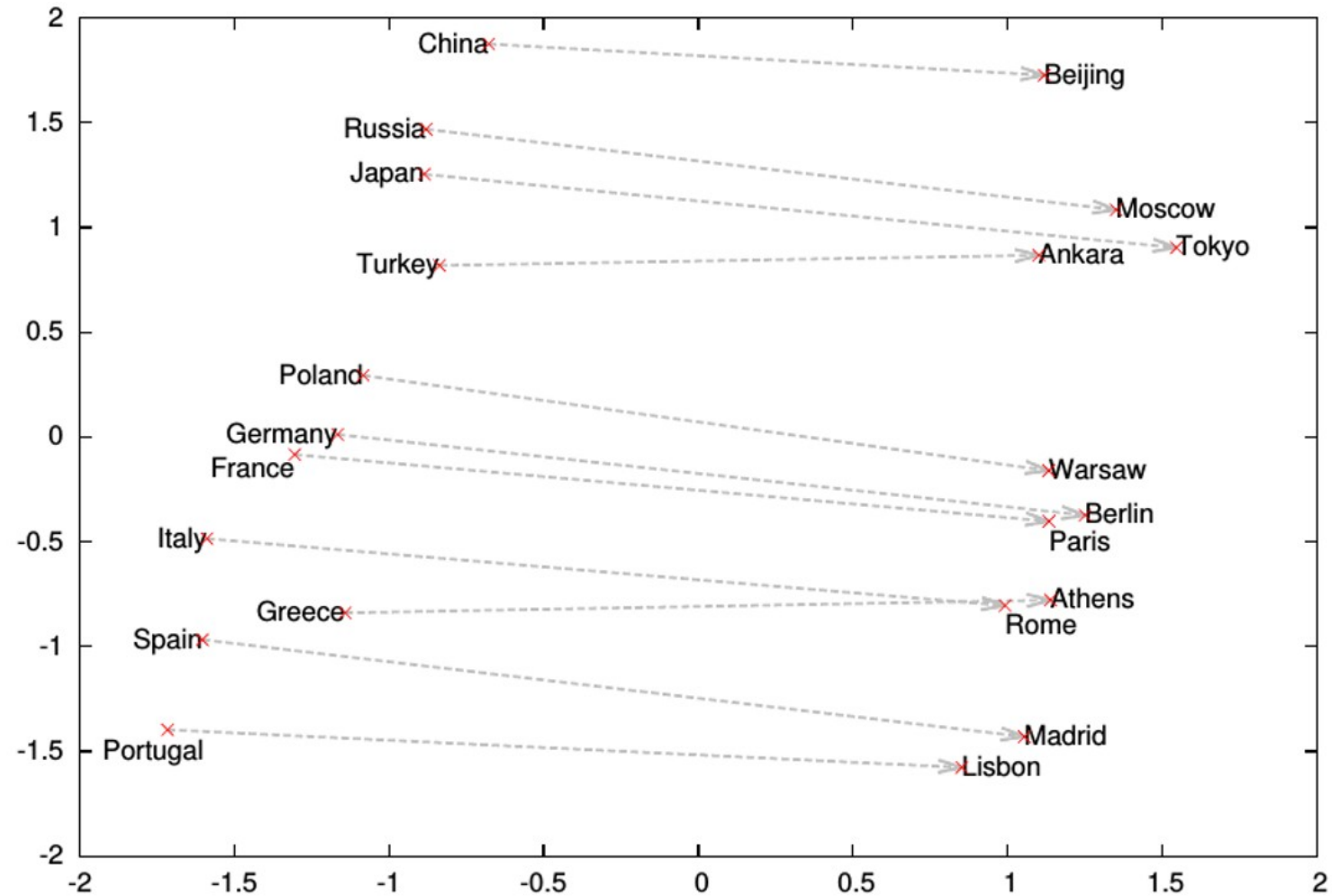
---

queen [ 0.70 0.80 ]



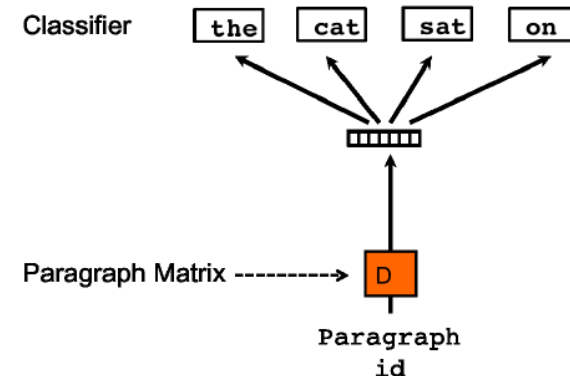
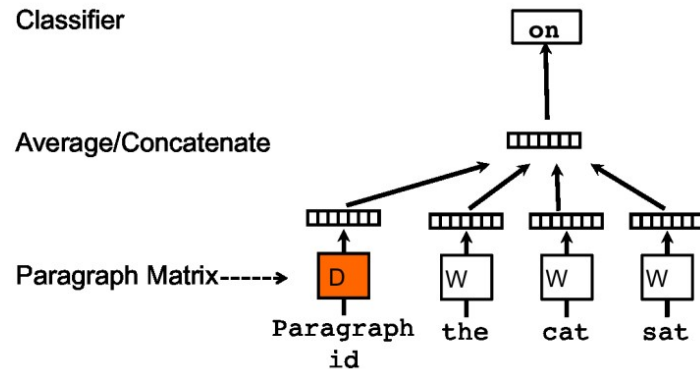


# Word analogies



# Represent the meaning of **sentence/text**

- Simple approach: take avg of the word2vecs of its words
- Another approach: Paragraph vector (2014, Quoc Le, Mikolov)
  - Extend word2vec to text level
  - Also two models: add paragraph vector as the input



# Applications

- Search, e.g., query expansion
- Sentiment analysis
- Classification
- Clustering

# Resources

- Stanford CS224d: Deep Learning for NLP
  - <http://cs224d.stanford.edu/index.html>
  - The best
- “word2vec Parameter Learning Explained”, Xin Rong
  - <https://ronxin.github.io/wevi/>
- Word2Vec Tutorial - The Skip-Gram Model
  - <http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/>
- Improvements and pre-trained models for word2vec:
  - <https://nlp.stanford.edu/projects/glove/>
  - <https://fasttext.cc/> (by Facebook)