

ANN AND DEEP LEARNING (CS636)

BY: Nidhi S. Periwal,
Teaching Assistant,
COED, SVNIT, Surat

Application of DL-> NLP

There are a lot of foods but I like meat best. I think that everyone always eats meat and meat is popular food. Meat is sold everywhere. For example: Supermarket, Market. Everybody always buys meat for the big party because it is delicious and cheap. Meat has a lot of proteins, if everyone eats a lot of meat they are intelligent and strong that's the reason. Why everyone eats meat. I like meat but sometimes I don't eat meat because I have to change different food each day otherwise. I will become fatter if I always eat only meat. However, I like meat best.

Bag-of-words

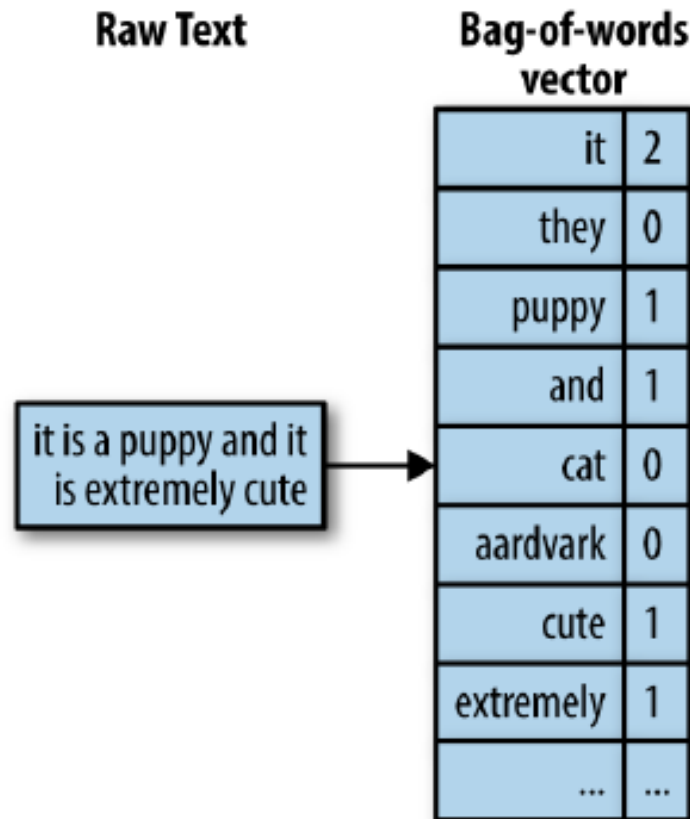


Fig. Turning raw text into a bag-of-words representation

Bag-of-Words

- *EG:*

- Similarly for other lines:
- "it was the worst of times" = [1, 1, 1, 0, 1, 1, 1, 0, 0, 0]
- "it was the age of wisdom" = [1, 1, 1, 0, 1, 0, 0, 1, 1, 0]
- "it was the age of foolishness" = [1, 1, 1, 0, 1, 0, 0, 1, 0, 1]
- These docs can also be represented as **document-term matrix**:

	it	was	the	best	of	times	worst	age	wisdom	foolishness
Doc 1	1	1	1	0	1	1	1	0	0	0
Doc 2	1	1	1	0	1	0	0	1	1	0
Doc 3	1	1	1	0	1	0	0	1	0	1

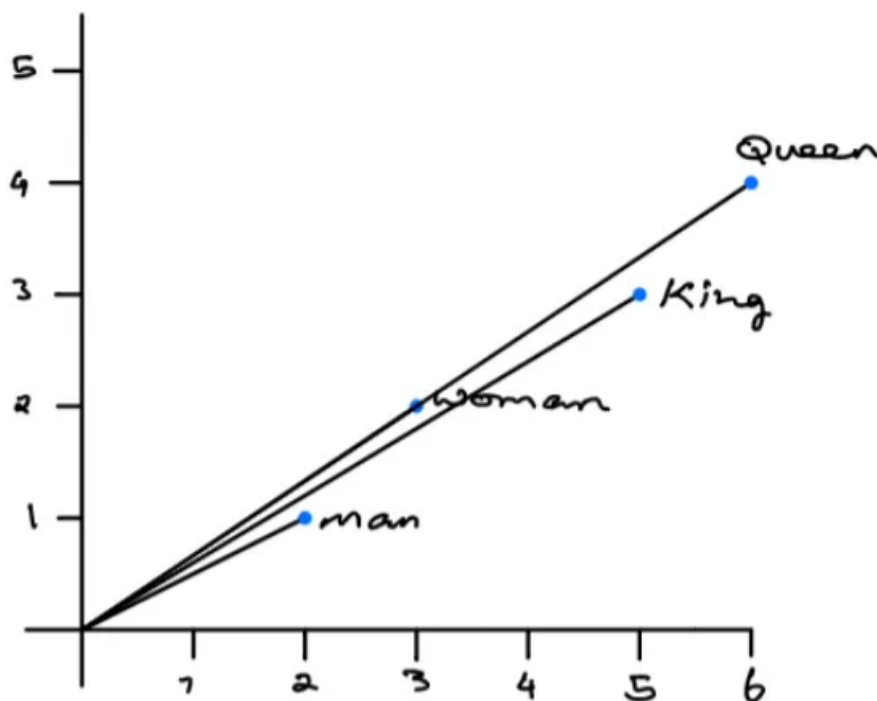
Table 1 : An example document-term matrix

Word Embedding

- Word Embedding is a technique where **individual words are transformed into a numerical representation of the word (a vector)**.
- Where each word is mapped to one vector, this vector is then learned in a way which resembles a neural network.
- **The vectors try to capture various characteristics of that word with regard to the overall text.**
- **These characteristics can include the semantic relationship of the word, definitions, context, etc.**
- **With these numerical representations we can identify similarity or dissimilarity between words.**

- Eg: 2 dimensional embedding vector of "king" - the 2 dimensional embedding vector of "man" + the 2 dimensional embedding vector of "woman" yielded a vector which is very close to the embedding vector of "queen".

King	-	Man	+	Woman	=	Queen
[5,3]	-	[2,1]	+	[3, 2]	=	[6,4]



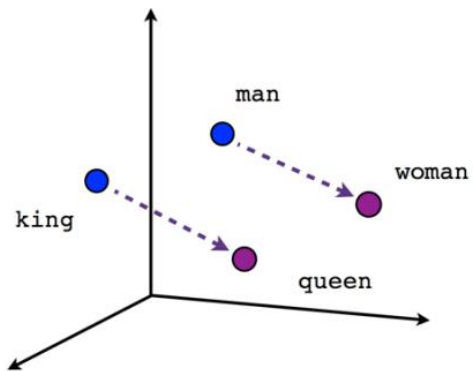
Word Embedding

- A word embedding is an approach to provide a dense vector **representation of words that capture something about their meaning.**
- Each word is represented by a point in the embedding space and these points are learned and moved around **based on the words that surround the target word.**
- Word Embeddings are the texts converted into numbers and there may be different numerical representations of the same text
- It is one of the most popular representation of document vocabulary.
- It is capable of capturing context of a word in a document, **semantic and syntactic similarity**, relation with other words, etc.

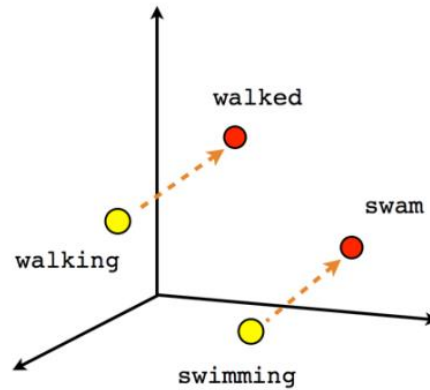
Word Embedding

- Need?
 - Machine Learning algorithms and Deep Learning Architectures are incapable of processing *strings* or *plain text* in their raw form.
 - They require numbers as inputs to perform any sort of job, be it classification, regression etc.

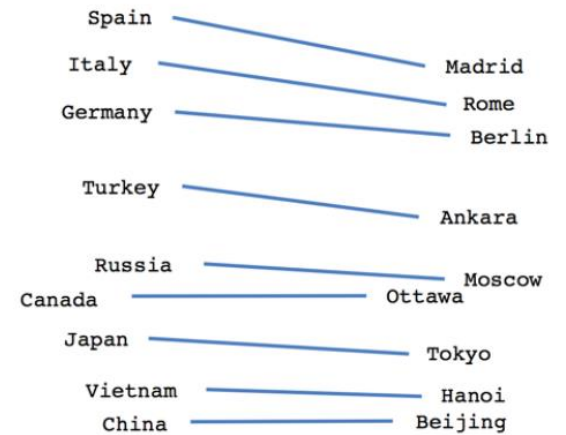
Word Embedding



Male-Female



Verb tense



Country-Capital

Word Embedding v/s Bag of words (BoW)

BoW	Word Embedding
<ul style="list-style-type: none">• Each word in the vocabulary is represented by one bit position in a HUGE vector	<ul style="list-style-type: none">• Stores each word in as a point in space, where it is represented by a vector of fixed number of dimensions
<ul style="list-style-type: none">• For example, if we have a vocabulary of 10000 words, and “SVNIT” is the 4th word in the dictionary, it would be represented by: 0 0 0 1 0 0 0 0 0 0	<ul style="list-style-type: none">• For example, “SVNIT” might be represented as : [0.4, -0.11, 0.55, 0.3 . . . 0.1, 0.02]
<ul style="list-style-type: none">• It uses word counts and frequencies that result in large and sparse vectors (mostly 0 values) that describe documents but not the meaning of the words.	<ul style="list-style-type: none">• Semantic meaning is retained of the document.

Word Embedding

- Different types of Word Embeddings
 - Frequency based Embedding
 - Count Vectors
 - TF-IDF
 - Co-Occurrence Matrix
 - Prediction based Embedding
 - CBOW
 - Skip-Gram

Word Embedding

- The word, which has been focused on to learn its representation is called **center word** and the words around it are called **context words**
- **Context can be** anything – a surrounding n-gram, a randomly sampled set of words from a fixed size window around the word

 : Center Word

 : Context Word

c=0 The cute  jumps over the lazy dog.

c=1 The    over the lazy dog.

c=2      the lazy dog.

One hot encoding

- Simplest embedding of text data
- One hot encoding is a vector representation of words in a “vocabulary

		cat	mat	on	sat	the
the =>		0	0	0	0	1
cat =>		1	0	0	0	0
sat =>		0	0	0	1	0
...						

Tf-Idf

- Tf-idf is a simple twist on the bag-of-words approach. It stands for ***term frequency–inverse document frequency***.
- It is a statistical measure **used to evaluate how important a word is to a document** in a collection or corpus.
- The tf-idf weight is composed by two terms:
 - **TF: Term Frequency**
 - **IDF: Inverse Document Frequency,**

Tf-Idf

– TF: Term Frequency

- **How frequently a term** occurs in a document.
- $TF(t) = (\text{Number of times term } t \text{ appears in a document}) / (\text{Total number of terms in the document})$.

– IDF: Inverse Document Frequency

- **measures how important a term** is
- While computing TF, all terms are considered equally important.
- However it is known that certain terms, such as "is", "of", and "that", may appear a lot of times but have little importance.
- Thus we need to weigh down the frequent terms while scale up the rare ones, by computing the following:
- $IDF(t) = \log(\text{Total number of documents} / \text{Number of documents with term } t \text{ in it})$.

Tf-Idf

– Hence, **Tf-Idf = tf*idf**

– Eg:

- Consider a document containing 100 words wherein the word *cat* appears 3 times.
- **tf for *cat* = (3 / 100) = 0.03.**
- Now, assume we have **10 million documents** and the word *cat* appears in **one thousand of these**.
- **idf = $\log(10,000,000 / 1,000) = 4$.**
- **Tf-idf : $0.03 * 4 = 0.12$.**

Tf-Idf

- Common words like ‘is’, ‘the’, ‘a’ etc. tend to **appear quite frequently in comparison** to the words which are important to a document.
- For example, a document **A on Lionel Messi is going to contain more occurrences of the word “Messi”** in comparison to other documents.
- But common words like **“the” etc. are also going to be present in higher frequency** in almost every document
- TF-IDF works by penalising these common words by assigning them lower weights **while giving importance to words like Messi in a particular document.**

Co-Occurrence Matrix

- Co-occurrence matrix records the co-occurrence count of the context words w.r.t. a center word in the specific window.
- Similar words tend to occur together and will have similar context
- Eg: Sun rises in East (Window size 1)

	Sun	Rises	In	East
Sun	0	1	0	0
Rises	1	0	1	0
In	0	1	0	1
East	0	0	1	0

- It preserves the semantic relationship between words i.e man and woman tend to be closer than man and apple.
- However, It requires huge memory to store the co-occurrence matrix

Word2Vec

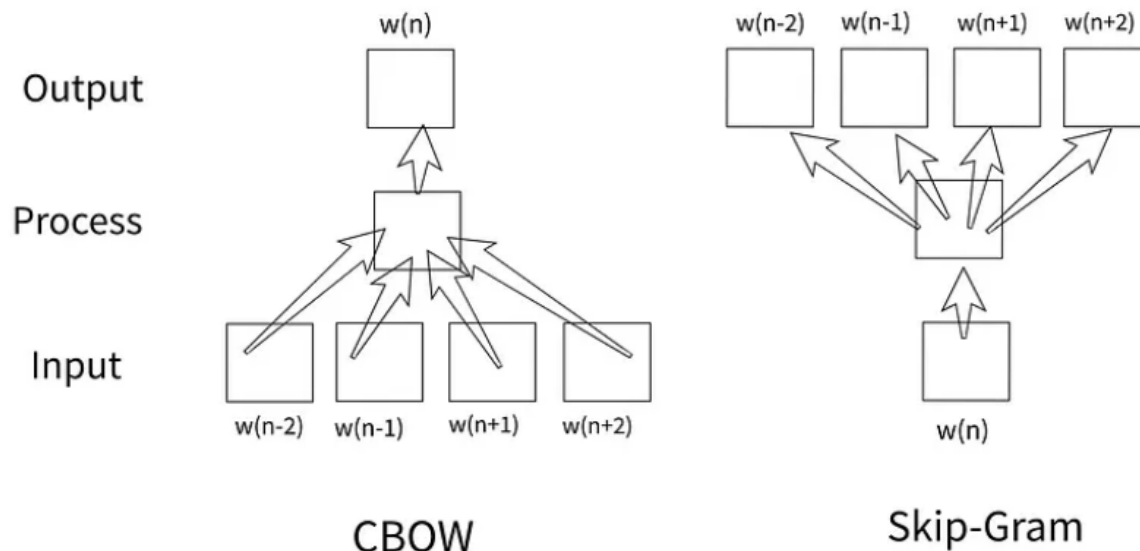
- **Word2Vec** model is used for **learning vector representations of words called “word embeddings”**
- Word2vec is algorithm for **learning a word embedding from a text corpus.**
- It learns the **similarity of word meaning from simple information.**
- The idea is based on the assumption that **the meaning of a word is affected by the words around it.**

- **Word2vec** is a technique for natural language processing (NLP) published in 2013. The word2vec algorithm uses a neural network model to learn word associations from a large corpus of text.
- Word2vec is a group of related models that are used to produce word embeddings. **These models are shallow, two-layer neural networks that are trained to reconstruct linguistic contexts of words.**
- **Word2vec takes as its input a large corpus of text and produces a vector space, typically of several hundred dimensions, with each unique word in the corpus being assigned a corresponding vector in the space.**
- Unlabeled data is trained via artificial neural networks to create the Word2Vec model that generates word vectors.

- Word2vec can utilize either of two model architectures to produce these distributed representations of words: continuous bag-of-words (CBOW) or continuous skip-gram.

Word2Vec

- Two 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.

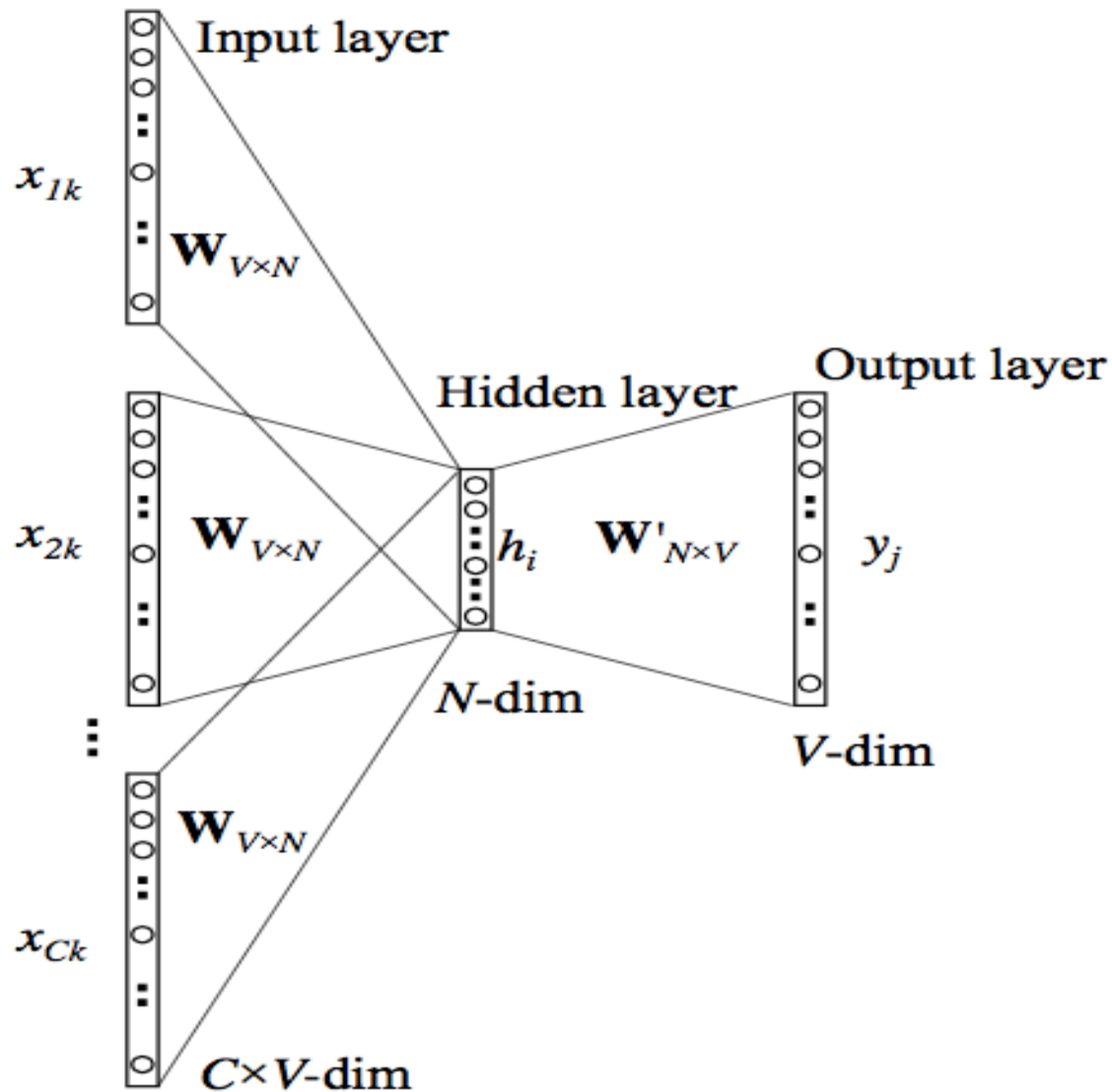


CBOW

- It takes the **context of each word as the input** and tries to predict the word corresponding to the context.
- The CBOW model takes a **window of surrounding words** as input and tries to predict the target word in the center of the window.
- The model is trained on a large text dataset and learns to predict the target word based on the patterns it observes in the input data.
- Here, we try to predict centre word by summing vectors of surrounding words.

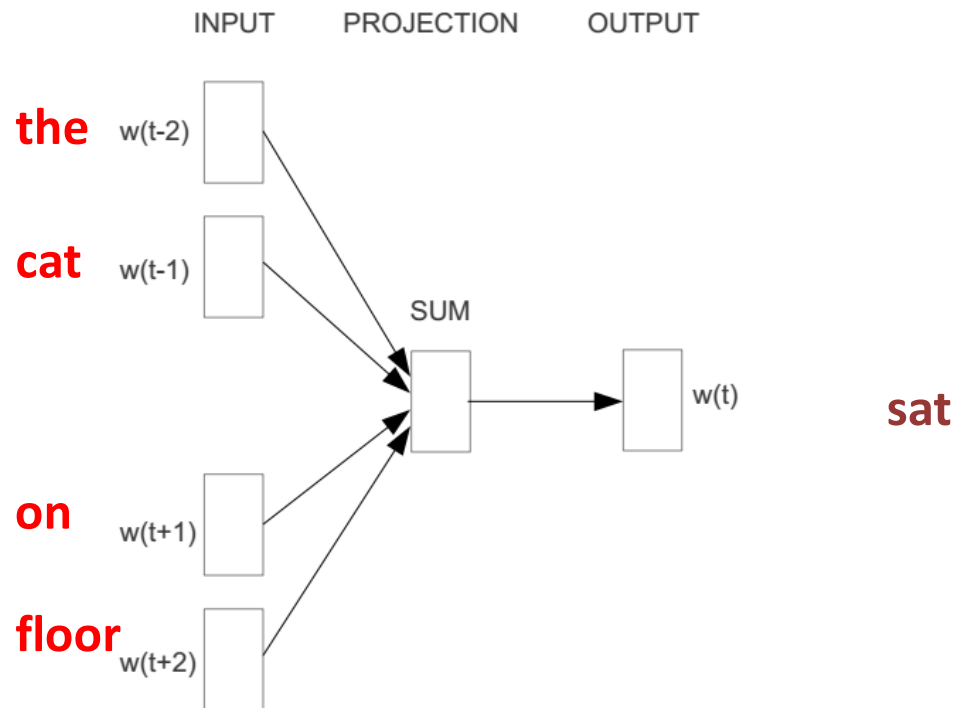
CBOW

- The architecture for multiple context words is shown in fig below:



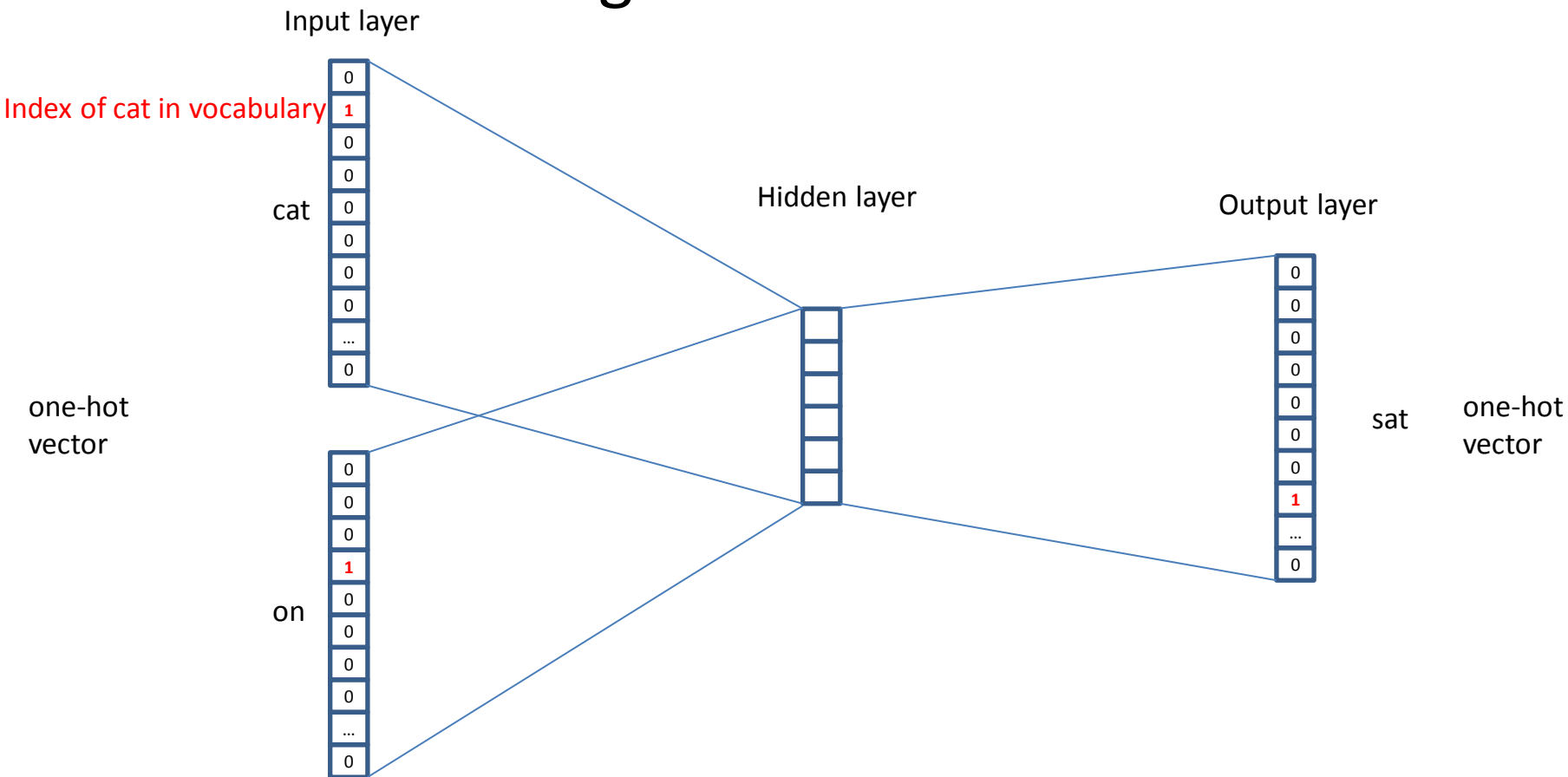
Word2Vec

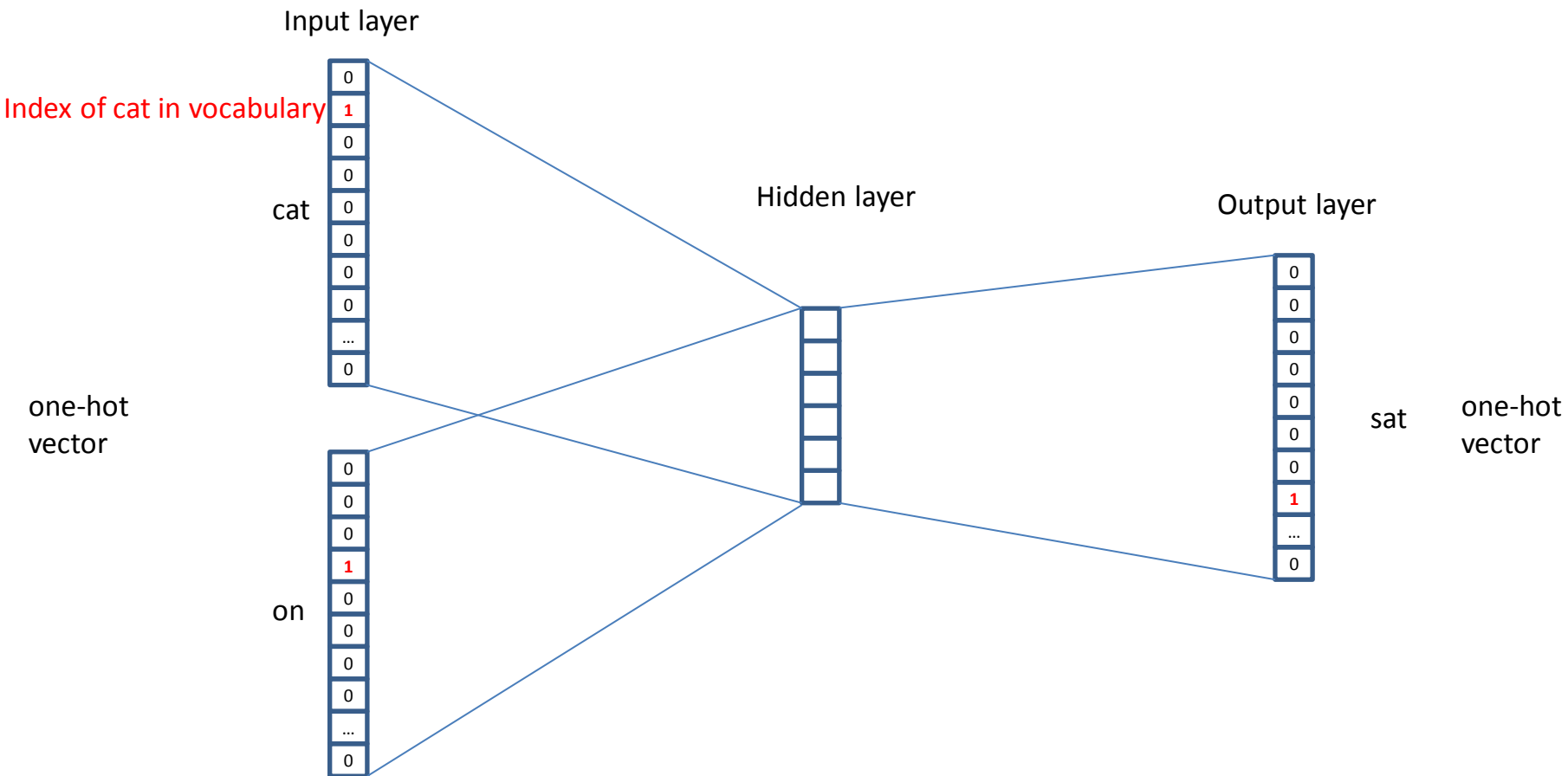
- Continuous Bag of Word (CBOW)
 - E.g. “The cat sat on floor” : Window size = 2



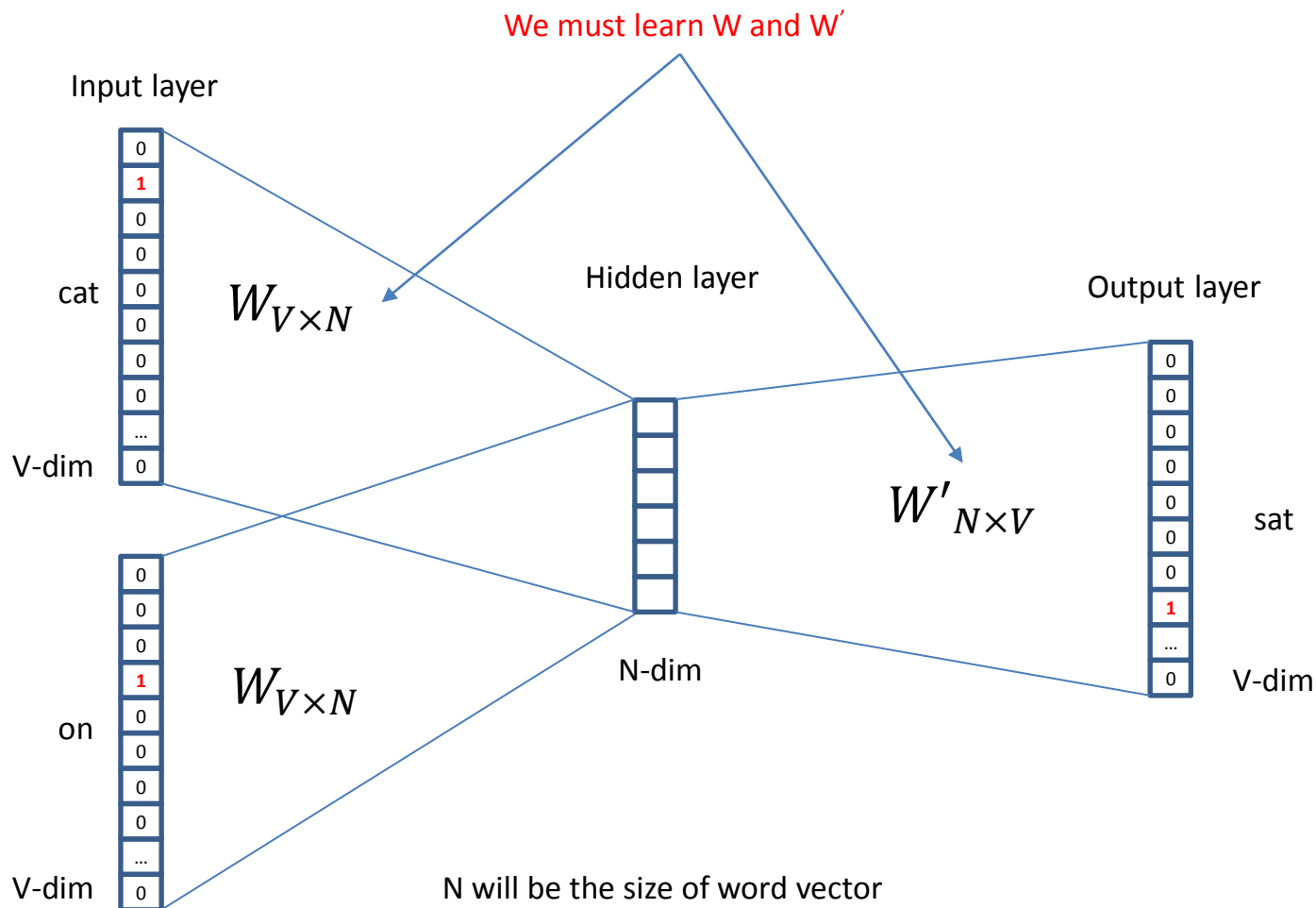
Word2Vec

- Continuous Bag of Word

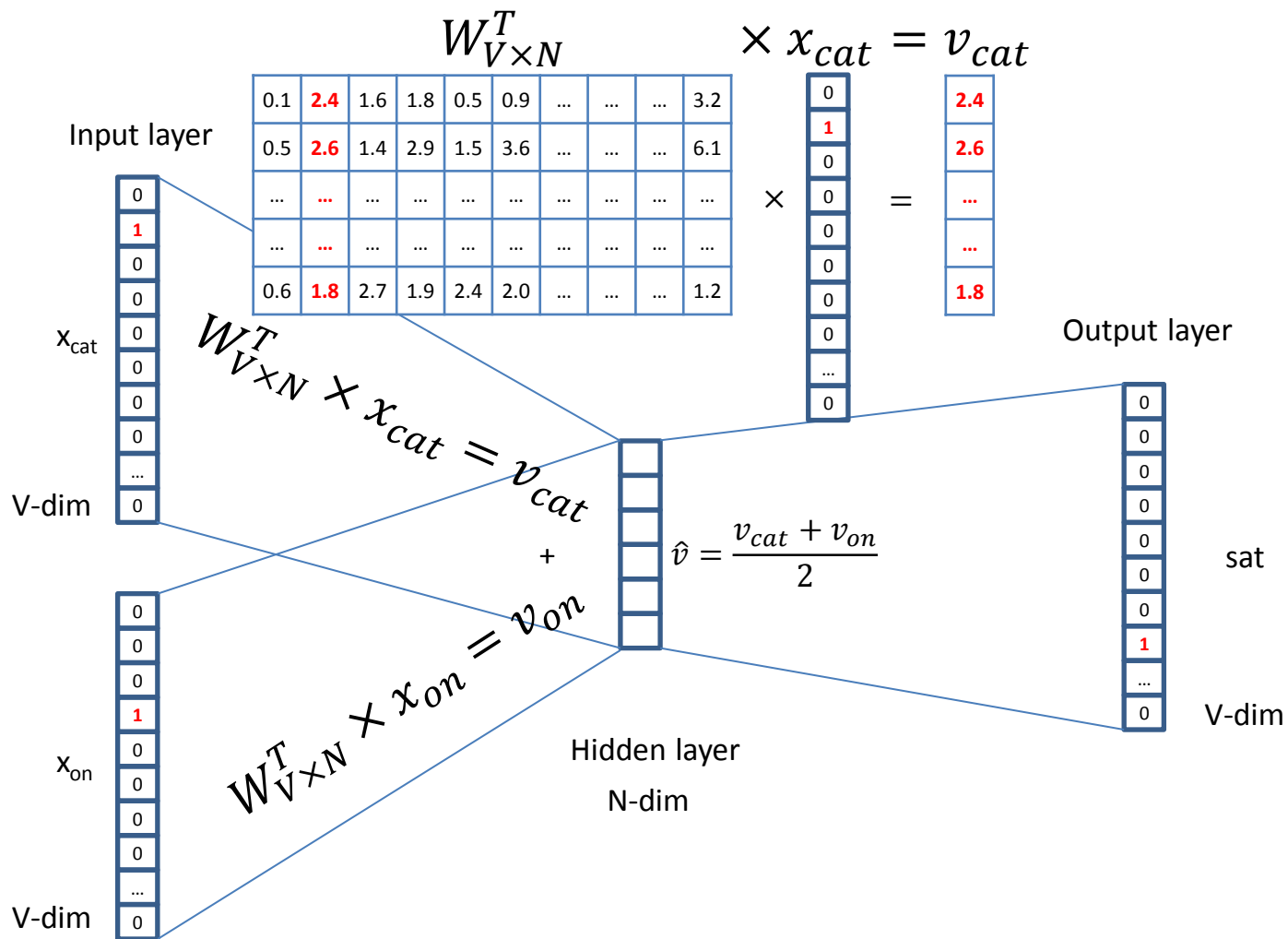


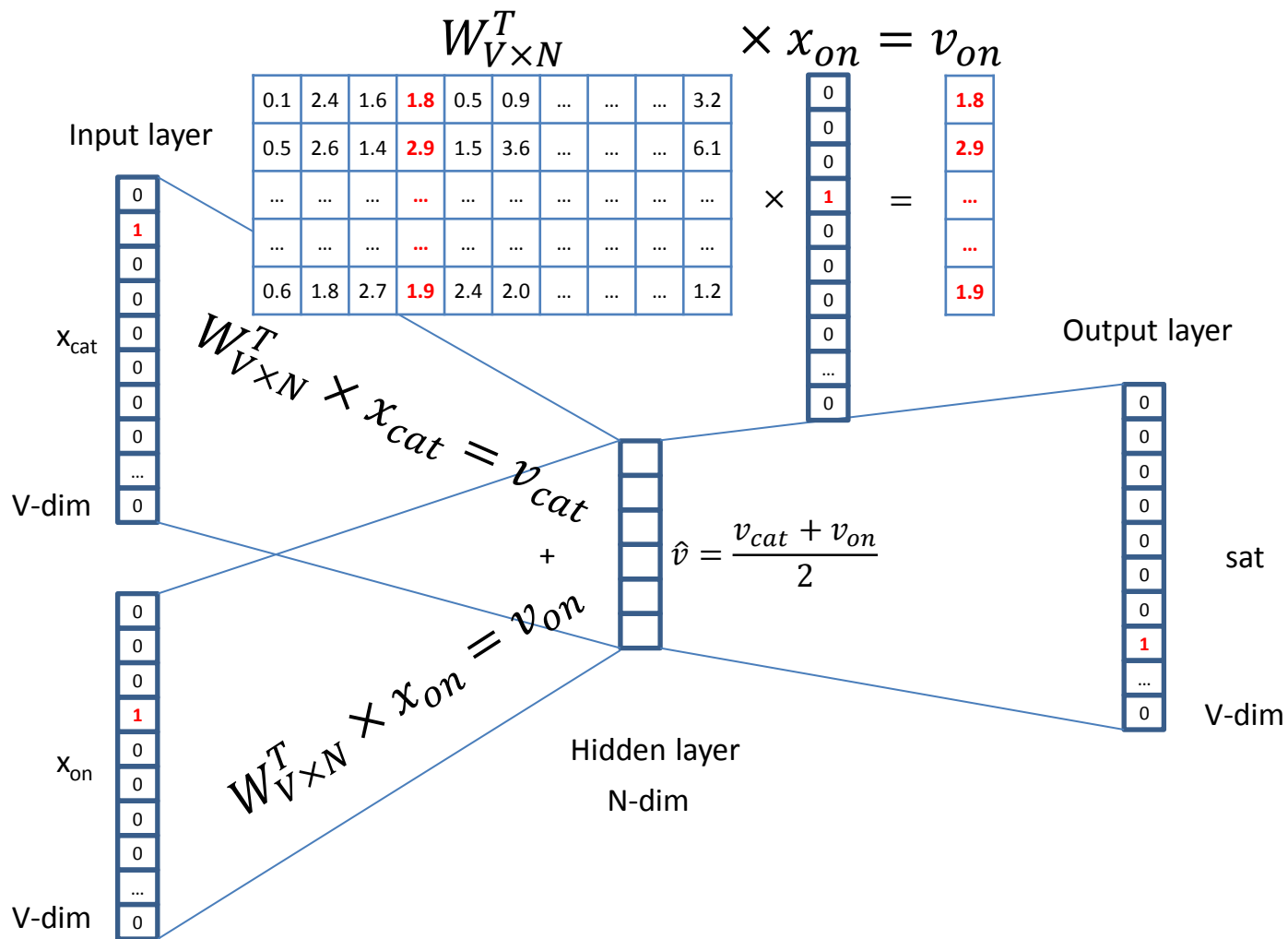


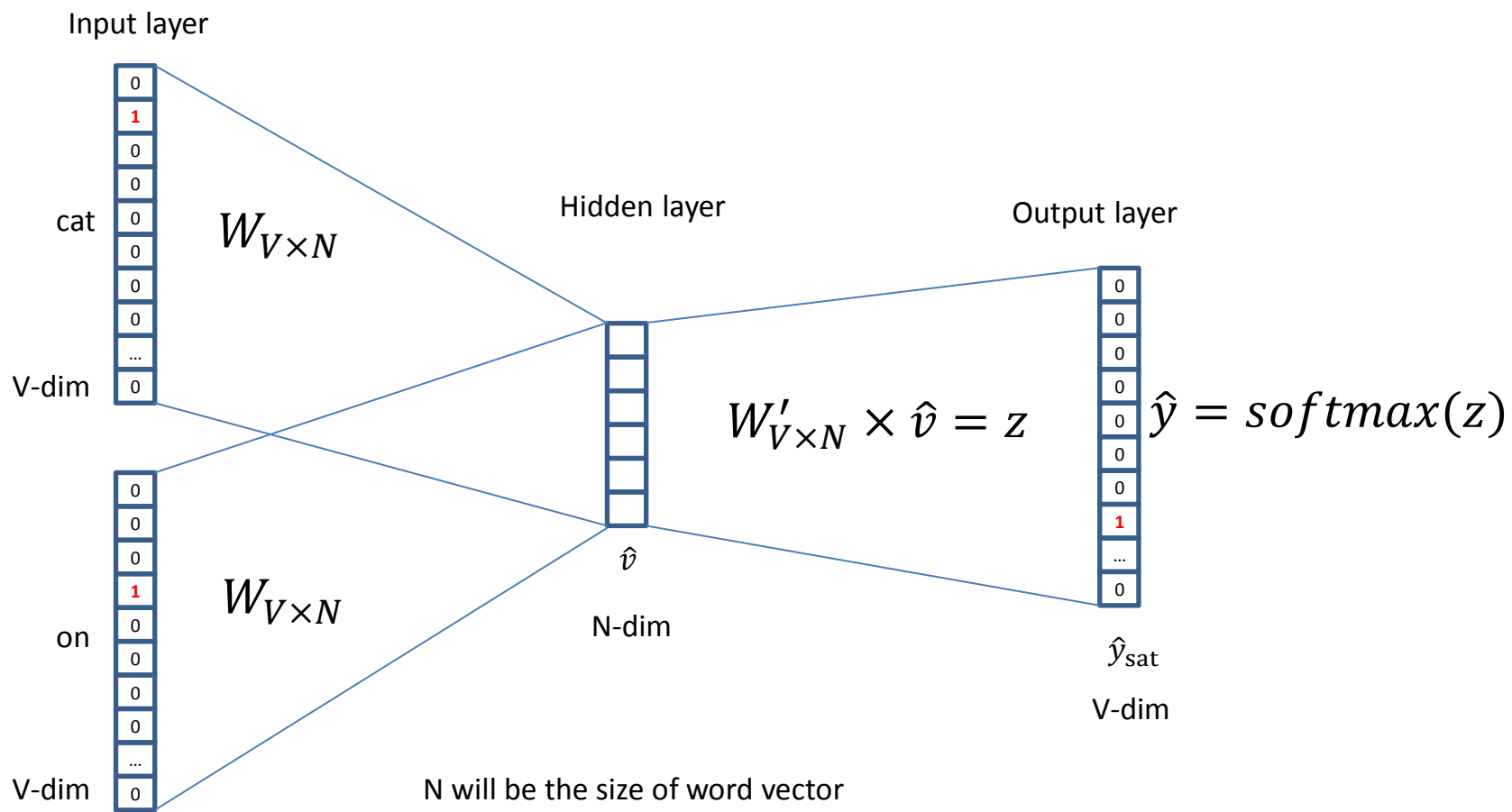
N is the number of dimensions we choose to represent our word in.
Also, N is the number of neurons in the hidden layer.
 V can be the Vocabulary

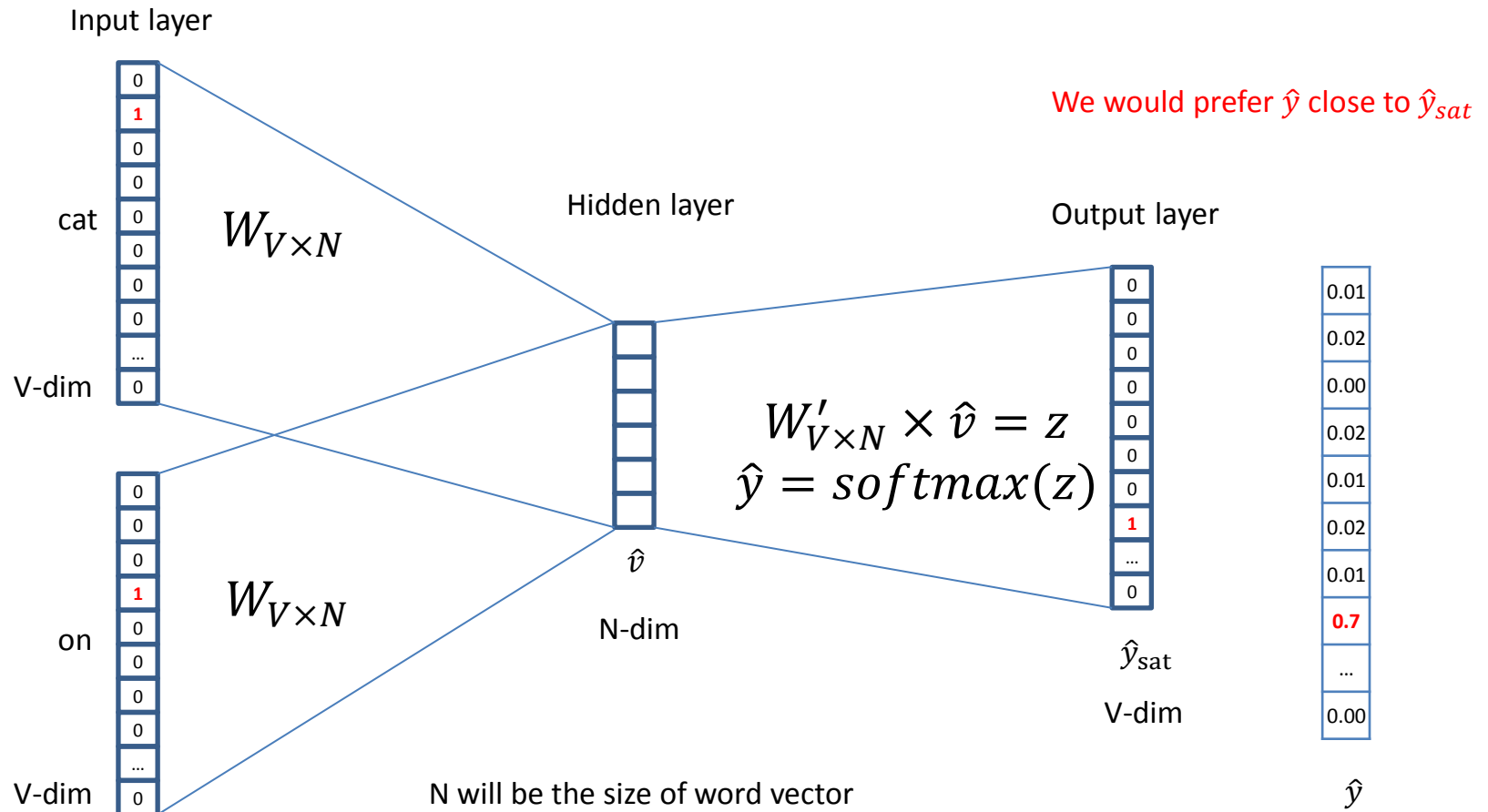


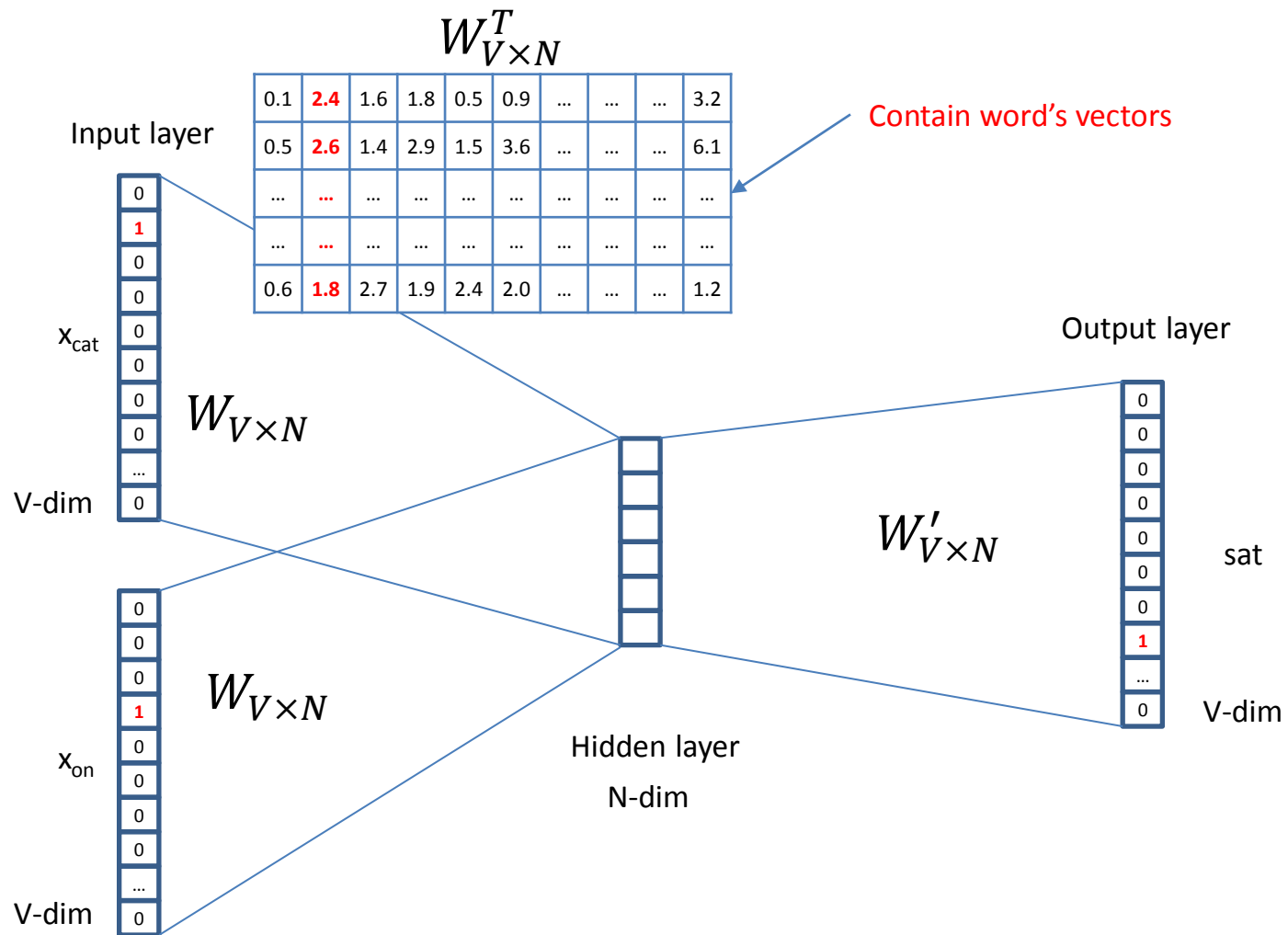
Input-Hidden layer matrix size = $[V \times N]$,
 hidden-Output layer matrix size = $[N \times V]$









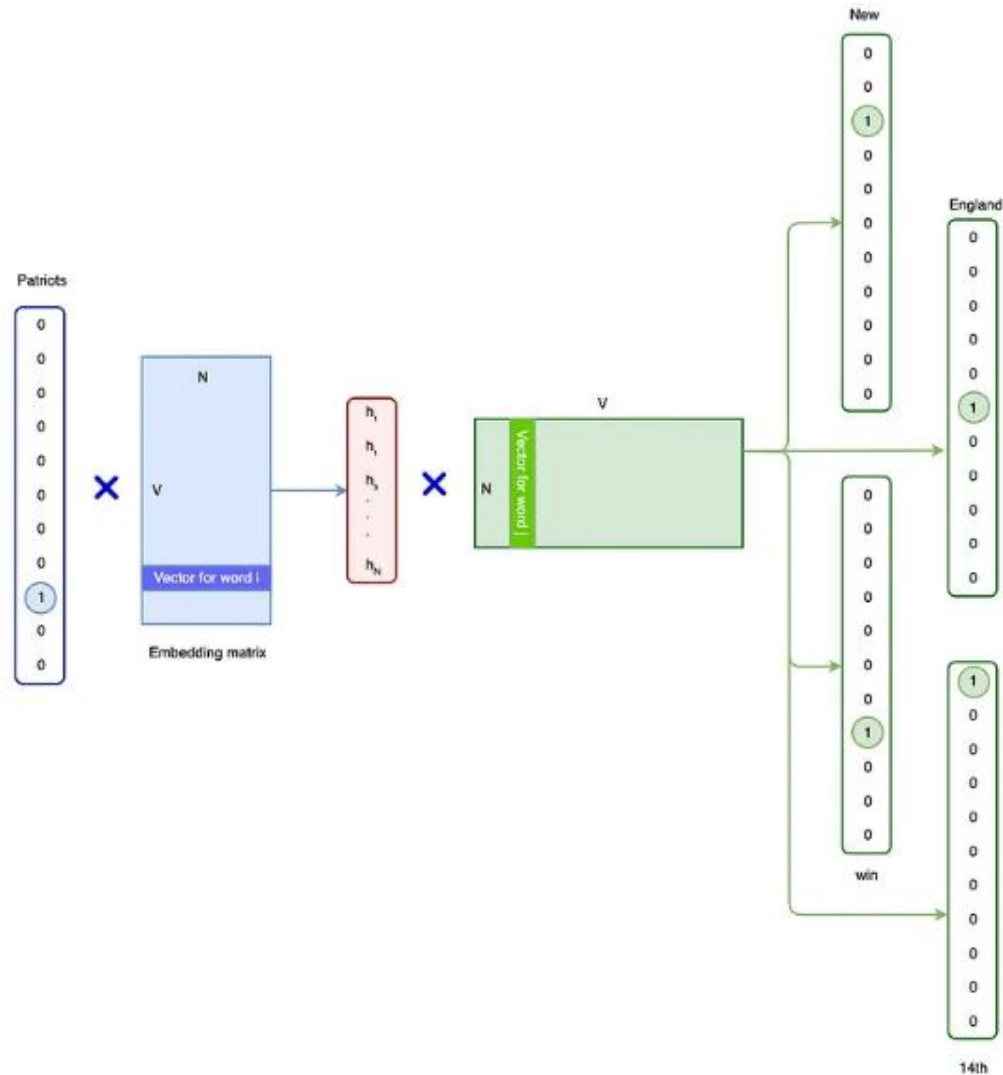


We can consider either W or W' as the word's representation. Or even take the average.

Skip gram Model

- In Skip-gram model, we take a **centre** word and a window of **context (neighbor)** words
- Here, we try to **predict context words** out to **some window size for each centre word**.
- So, our model is going to define a **probability distribution** i.e. **probability of a word appearing in context given a centre word** and we are going to choose our vector representations to **maximize the probability**.

- **Example:** *New England Patriots win 14th straight regular-season game at home.*



- The log-likelihood for the predicted words given the target word t (“Patriots”) will be:

A sequence of training words (w_1, w_2, w_3, \dots)

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t)$$

Surrounding words for word t

New England **Patriots** win 14th straight regular-season game at home.

Patriots → New
 Patriots → England
 Patriots → win
 Patriots → 14th

GloVe

- GloVe stands for global vectors for word representation.
- It is an unsupervised learning algorithm developed by Stanford for generating word embeddings by aggregating global word-word co-occurrence matrix from a corpus.

- In GloVe, a word co-occurrence matrix is generated, where rows represent the words and columns represent the context.
- When co-occurrence frequencies are extracted at a sentence level, every word is said to be in the context of another word in the same sentence, and the corpus can be represented in the following matrix form.

	<i>I</i>	<i>love</i>	<i>Chemistry</i>	<i>Maths</i>	<i>tolerate</i>	<i>Biology</i>
1. <i>I love chemistry.</i>	0	2	1	1	1	1
2. <i>I love maths.</i>	2	0	1	1	0	0
3. <i>I tolerate biology.</i>	1	1	0	0	0	0
	1	1	0	0	0	0
	1	0	0	0	0	1
	1	0	0	0	1	0

- X represents the matrix of the co-occurrence frequencies of words in the corpus. Each value in X is interpreted as how frequently a word co-occurs with its context. Factorization of the co-occurrence matrix results in a low-dimensional matrix, where rows represent words and columns represent features.

- Let the matrix of word-word co-occurrence counts be denoted by X
- whose entries X_{ij} tabulate the number of times word j occurs in the context of word i .
- Let $X_i = \sum_k X_{ik}$ be the number of times any word appears in the context of word i .
- $P_{ij} = P(j|i) = X_{ij}/X_i$ probability that word j appear in the context of word i .

Table 1: Co-occurrence probabilities for target words *ice* and *steam* with selected context words from a 6 billion token corpus. Only in the ratio does noise from non-discriminative words like *water* and *fashion* cancel out, so that large values (much greater than 1) correlate well with properties specific to ice, and small values (much less than 1) correlate well with properties specific of steam.

Probability and Ratio	$k = \text{solid}$	$k = \text{gas}$	$k = \text{water}$	$k = \text{fashion}$
$P(k \text{ice})$	1.9×10^{-4}	6.6×10^{-5}	3.0×10^{-3}	1.7×10^{-5}
$P(k \text{steam})$	2.2×10^{-5}	7.8×10^{-4}	2.2×10^{-3}	1.8×10^{-5}
$P(k \text{ice})/P(k \text{steam})$	8.9	8.5×10^{-2}	1.36	0.96

- Let $i = \text{ice}$ and $j = \text{steam}$. The relationship of these words can be examined by studying the ratio of their co-occurrence probabilities with various probe words, k .
- For words k related to **ice** but not steam, say $k = \text{solid}$, we expect the ratio P_{ik}/P_{jk} will be **large**.
- For words k related to **steam** but not ice, say $k = \text{gas}$, the ratio should be **small**.
- For words k like water or fashion, that are either related to both ice and steam, or to neither, the ratio should be close to one.

$w \in \mathbb{R}^d$ are word vectors

probe word

$$F(w_i, w_j, \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}$$

co-relations between the word w_i and w_j

co-occurrence probabilities for the word w_j and w_k

- The objective function (weighted least square) in the GloVe model

$$J = \sum_{i,j=1}^V f(X_{ij})(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$$

- Here, V refers to the vocabularies in the corpora and f is a weighting function, which assigns lower weights to rare co- occurrences .
- X_{ij} is the co-occurrence matrix for a target word (i) and its context word (j), and w_i , w_j , b_i and b_j are a set of trainable parameters for i and j ,
- where w_i and w_j are the embeddings,
- and b_i and b_j are their corresponding biases.
- GloVe learns neural embeddings by minimizing the reconstruction error between co-occurrence statistics predicted by the model and global co-occurrence statistics observed in the training corpus

Topic Analysis

- Topic Analysis
 - A technique that allows us to automatically **extract meaning from texts by identifying recurrent themes or topics.**
 - It is technique that **organizes and understands large collections of text data**, by **assigning tags or categories** according to each individual text's topic or theme
 - **Two Approaches to it:**
 - **Topic Modeling**
 - unsupervised machine learning technique.
 - It can **infer patterns and cluster similar expressions without** needing to define **topic tags or train data** beforehand
 - **Topic Classification**
 - We must know topics of a text before starting the analysis.
 - Tag data in order to train a topic classifier

Topic Analysis

- **Topic analysis** can be applied at different levels of scope:
 - **Document-level**: the topic model obtains the different topics within a complete text. **For example, the topics of an email or an article.**
 - **Sentence-level**: the topic model obtains the topic of a single sentence. For example, **the topic of a news article headline.**
 - **Sub-sentence level**: the topic model obtains the topic of sub-expressions within a sentence. **For example, different topics within a single sentence of a product review.**

Topic Models

Gordon Roque @gordonroque May 27
Today, I am #thankful for my car. I drive a red 2008 #Toyota #Scion
xD. Its compact size allows it... [instagram.com/p/3NcvlgQUsk/](https://www.instagram.com/p/3NcvlgQUsk/)

8:14 PM - 27 May 2015 - Details



Topics:

- Vehicle
- Compact Car

Strong interest in car.

Topic Models

Topics

```
gene      0.04
dna       0.02
genetic   0.01
...
```

```
life      0.02
evolve    0.01
organism  0.01
...
```

```
brain      0.04
neuron     0.02
nerve      0.01
...
```

```
data      0.02
number    0.02
computer  0.01
...
```

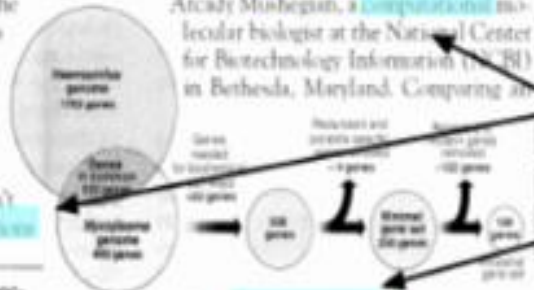
Documents

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here, two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 252 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson at the University in Sweden. "I just arrived at the 800 number, but coming up with a concrete answer may be more than just a numbers game, particularly as more and more genomes are completely sequenced and sequenced." "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing all



* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments



Introduction

- Topic models are a method of extracting information from bodies of text to see what “**topics**” occur across all the documents.
- It is a method of **uncovering hidden (latent) structure** in a collection of texts.
- The “**topics**” produced by topic modeling techniques are **clusters of similar words**
- **Topic models** are also referred to as **probabilistic topic models**, which refers to statistical algorithms for **discovering the latent semantic structures of an extensive text body**.

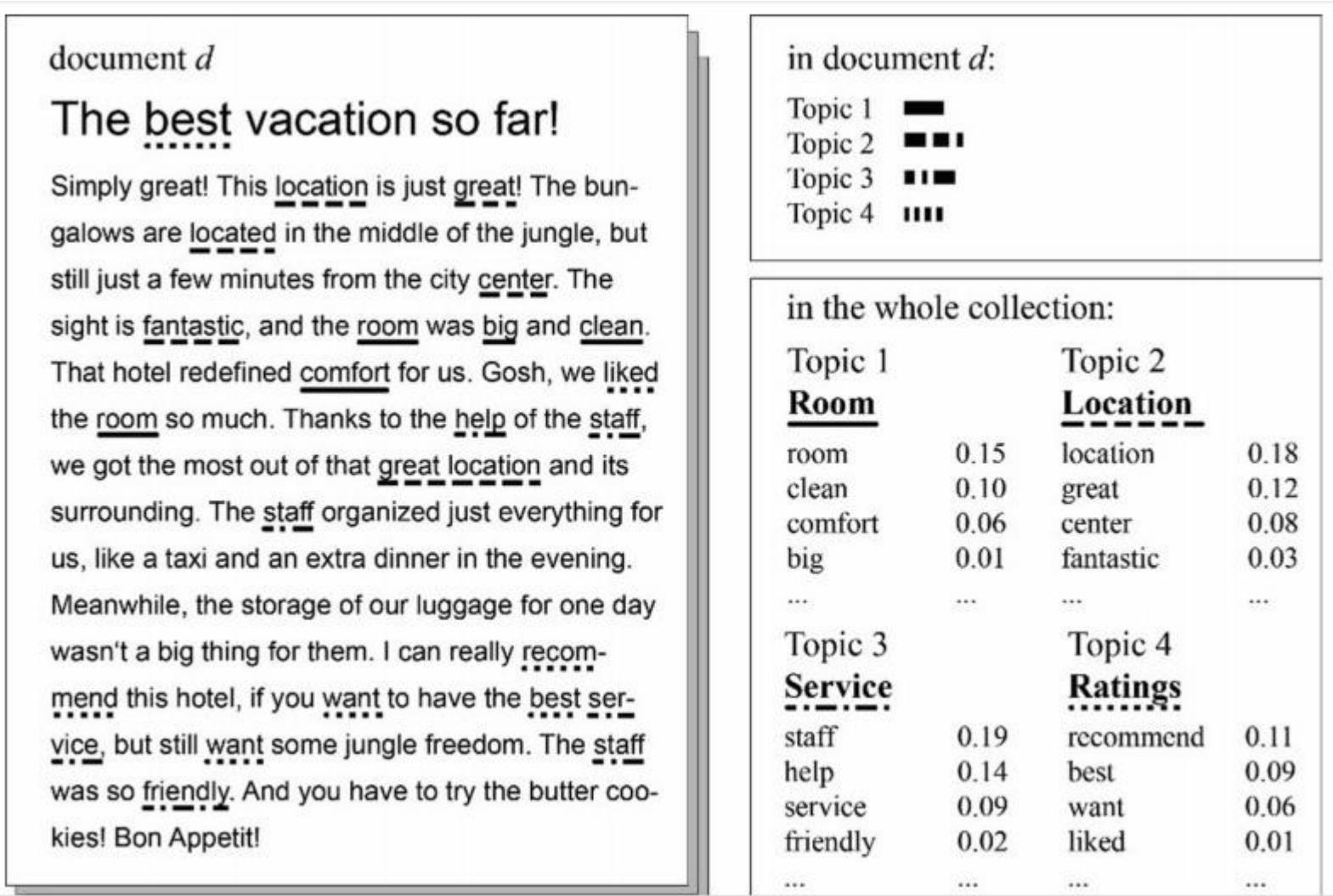
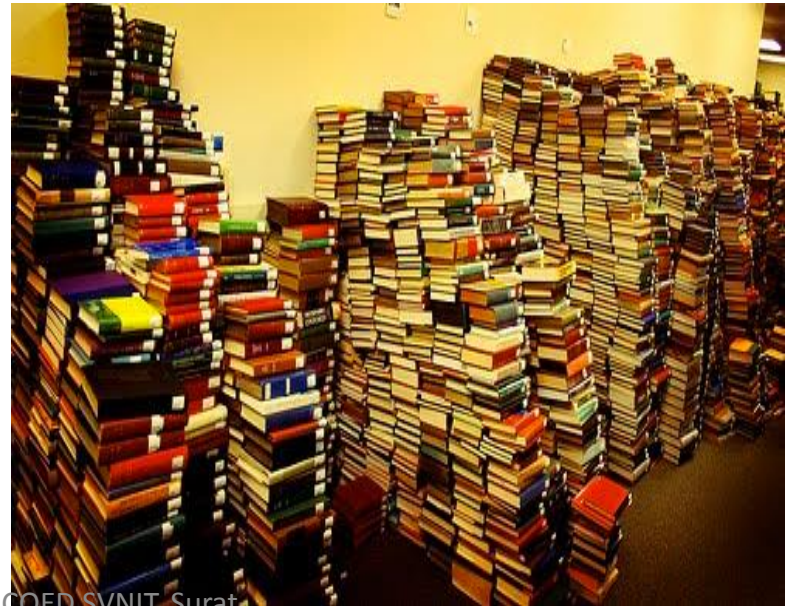


Fig. Ref: Paper: Martin Reisenbichler, Thomas Reutterer, "Topic modeling in marketing: recent advances and research opportunities", Journal of Business Economics, Springer, 2019.

Need of Topic Models

- In the age of information, **the amount of the written material we encounter each day is simply beyond our processing capacity.**
- Topic models can help to **organize and offer insights** for us to understand **large collections of unstructured text bodies.**



Topic Models-Applications

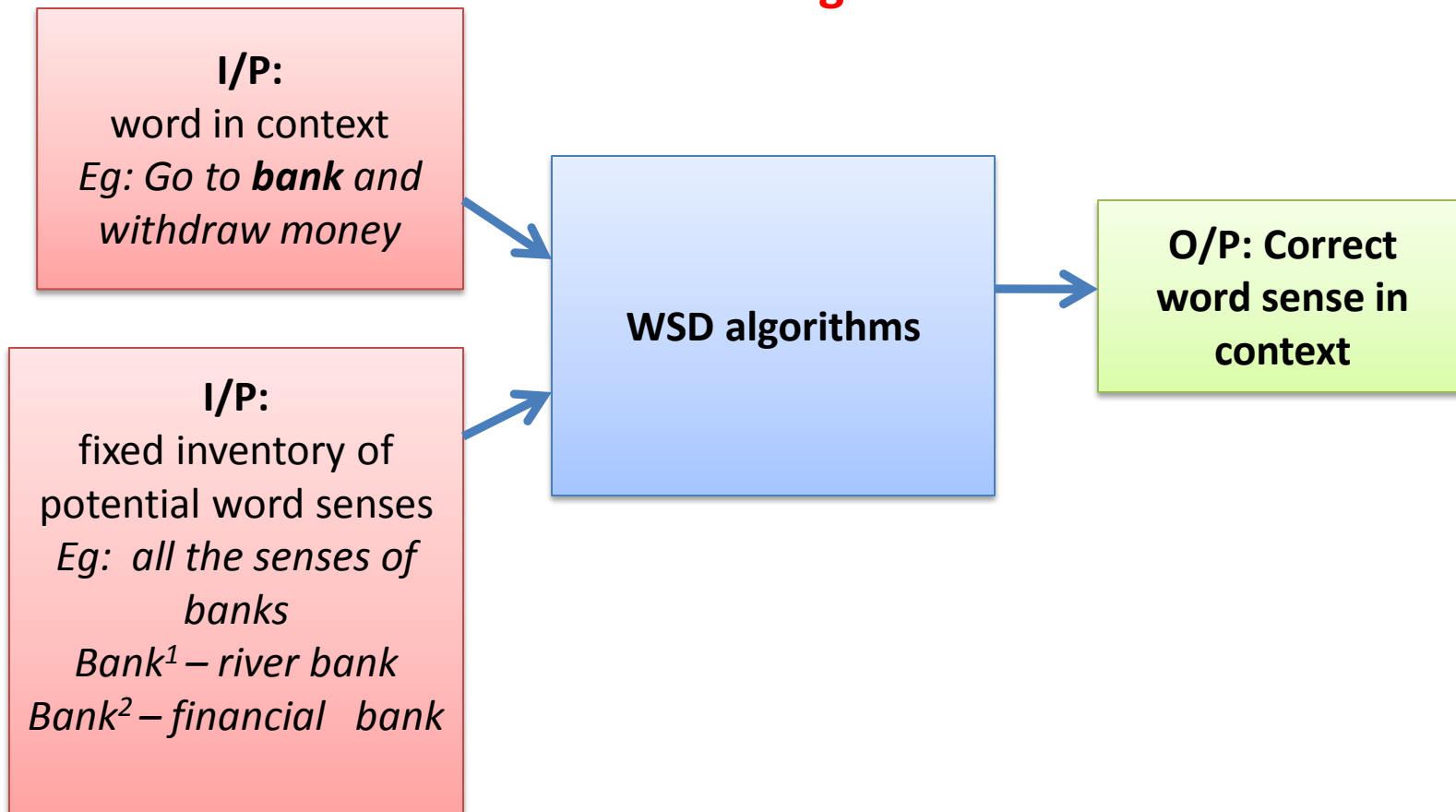
- **Organize** the documents into **thematic** categories
- **Describe** the evolution of those categories **over time**
- Enable a **domain expert** to **analyze and understand** the content
- Information **Retrieval**
- **Classify** the documents
- Find **relationships** between the categories
- Understand how **authorship** influences the content
- It also helps in:
 - Discovering **hidden topical patterns** that are present across the collection
 - **Annotating** documents according to the topics
 - Using **these annotations to organize, search and summarize texts**
 - Find a **short description of the data**.
 - Preserve the **essential statistical relationships**

Topic Models-Applications

Word Sense Disambiguation

➤ **Task : Select the correct sense for a word**

➤ **Given : Context and meaning of a word**



Flowchart

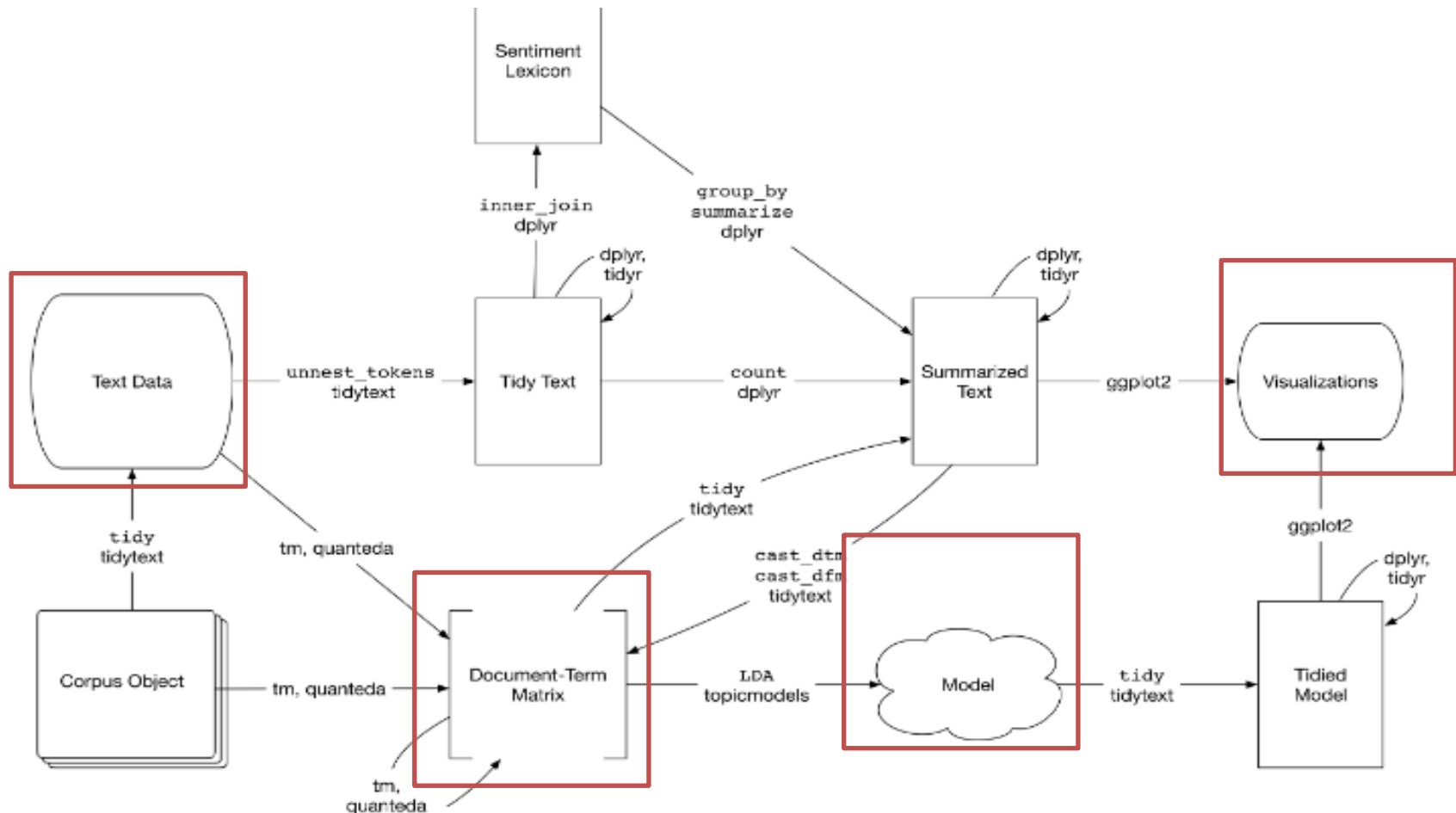


Figure 6.1: A flowchart of a text analysis that incorporates topic modeling. The `topicmodels` package takes a Document-Term Matrix as input and produces a model that can be tidied by `tidytext`, such that it can be manipulated and visualized with `dplyr` and `ggplot2`.

Thank you!