# Word Embedding

# Text Analysis



Text analytics
...into **high-quality information** or...
...**human effort** (on consuming text...
...owledge for **optimal decision ma**...
...**t retrieval**, which is an essential c...
...g system
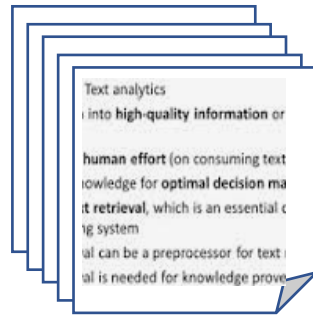...al can be a preprocessor for text...
...al is needed for knowledge prove...

# Text Analysis



**Various Applications**
- Text Recommendation
    - Information Retrieval
- Text Classification
    - Sentiment analysis
    - Hate/Communal speech detection
    - Fake news detection
- NLP
    - POS
    - NER
    - Machine transliteration
    - Machine Translation
- Image/Video Tagging

# Text Analysis



**We need effective representation of :**
- **Words**
- **Sentences**
- **Documents**

**Various Applications**
- Text Recommendation
    - Information Retrieval
- Text Classification
    - Sentiment analysis
    - Hate/Communal speech detection
    - Fake news detection
- NLP
    - POS
    - NER
    - Machine transliteration
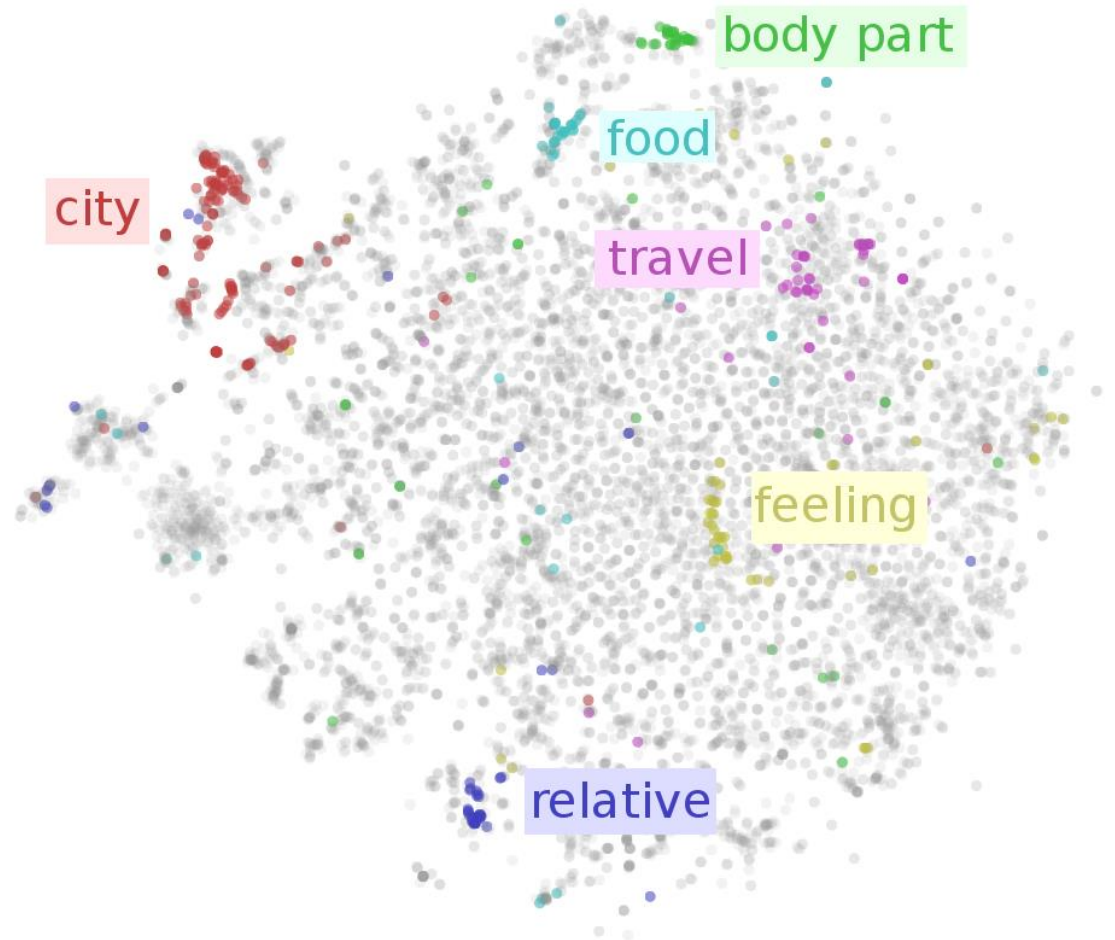    - Machine Translation
- Image/Video Tagging

# Text Analysis



We need effective representation of :
- Words
- Sentences
- Document

# Traditional Form of Text Representation

$$t_1 \ t_2 \ t_3 \qquad\qquad t_m$$

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| $d_1$ | 1 | 0 | 1 | 1 | …… . …. | 0 | 0 0 1 |
| $d_2$ | 1 | 1 | 0 | 0 | …… . …. | 0 | 1 1 0 |
| $d_3$ | | | | | | | |
| $d_4$ | | | | | | | |
| $d_n$ | 0 | 0 | 1 | 0 | …… . …. | 0 | 1 0 0 |

$n_x m$

**Document and Term Matrix**

**Possible Terms are**

- Unigram
- Bigram
- Tri-gram

# Traditional Form of Text Representation

$$
\begin{array}{c c c c c c c c c c c}
 & t_1 & t_2 & t_3 & & & & & & & t_1 \\
d_1 & 1 & 0 & 1 & 1 & \ldots\ldots & . & \ldots & 0 & 0 & 0 & 1 \\
d_2 & 1 & 1 & 0 & 0 & \ldots\ldots & . & \ldots & 0 & 1 & 1 & 0 \\
d_3 & & & & & & & & & & & \\
d_4 & & & & & & & & & & & \\
 & & & & & & & & & & & \\
d_n & 0 & 0 & 1 & 0 & \ldots\ldots & . & \ldots & 0 & 1 & 0 & 0 \\
\end{array}
$$

**A document is represented by term features**

**Document and Term Matrix**

# Traditional Form of Text Representation

|       | $t_1$ | $t_2$ | $t_3$ |        |       | $t_1$ |   |   |   |
|-------|-------|-------|-------|--------|-------|-------|---|---|---|
| $d_1$ | 1     | 0     | 1     | 1      | …… . …. | 0   | 0 | 0 | 1 |
| $d_2$ | 1     | 1     | 0     | 0      | …… . …. | 0   | 1 | 1 | 0 |
| $d_3$ |       |       |       |        |       |       |   |   |   |
| $d_4$ |       |       |       |        |       |       |   |   |   |
| $d_n$ | 0     | 0     | 1     | 0      | …… . …. | 0   | 1 | 0 | 0 |

→ **A term is represented by document features**

**Document and Term Matrix**

# Traditional Form of Text Representation

$$\begin{array}{cccccccccc} & t_1 & t_2 & t_3 & & & & & & t_1 \\ d_1 & \boxed{1 & 0 & 1 & 1 & \ldots\ldots.\ldots & 0 & 0 & 0 & 1} \\ d_2 & 1 & 1 & 0 & 0 & \ldots\ldots.\ldots & 0 & 1 & 1 & 0 \\ d_3 & & & & & & & & & \\ d_4 & & & & & & & & & \\ & & & & & & & & & \\ d_n & 0 & 0 & 1 & 0 & \ldots\ldots.\ldots & 0 & 1 & 0 & 0 \end{array}$$

**Document and Term Matrix**

**Find the most similar documents**

- Estimate similarity between the documents

$$Cos\theta(d_i, d_j) = \frac{\bar{d}_i . \bar{d}_j}{\|\bar{d}_i\| . \|\bar{d}_j\|}$$

**Find the most similar words**

- Estimate similarity between the documents

# Traditional Form of Text Representation

|     | $t_1$ | $t_2$ | $t_3$ |     |     |     |     | $t_1$ |
|-----|-------|-------|-------|-----|-----|-----|-----|-------|
| $d_1$ | 1 0 1 1 …… . …. 0 0 0 1 |
| $d_2$ | 1 1 0 0 …… . …. 0 1 1 0 |
| $d_3$ |
| $d_4$ |
|     |
| $d_n$ | 0 0 1 0 …… . …. 0 1 0 0 |

**Document and Term Matrix**

**Issues with document-term representations**

- Large feature set
- Many of the features may not be useful
- Sparse
- Curse of dimensionality

# What is word embedding?

**Word embedding** is a dense representation of words in the form of numeric vectors in lower dimensional space.
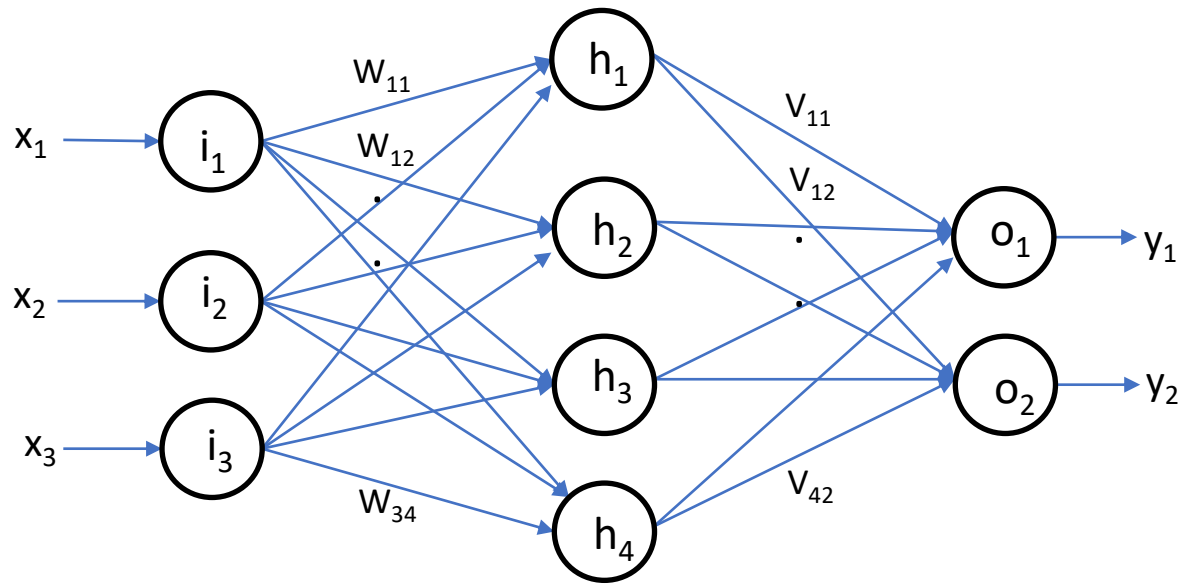
# What is word embedding?

**Word embedding** is a dense representation of words in the form of numeric vectors in lower dimensional space.

**Earlier Methods**
- Principal component Analysis (PCA)
- Singular Value Decomposition (SVD) – Latent Semantic Indexing
- Latent Dirichlet Allocation

**Recent Methods**
- Word2Vec (Skipgram/CBOW)
- Glove
- FastText
- BERT

This model can be used for various tasks
- Classification
- Regression
- Probability estimate
- Representation
- ….

This model can be used for various tasks
- Classification
- Regression
- **Probability estimate**
- Representation
- ….

Example Sentence : The man sat on the floor

Training samples : (The, man), (man, sat), (sat, on), (on, the),……..

Pr( man | The) = ?

This model can be used for various tasks
- Classification
- Regression
- **Probability estimate**
- Representation
- ....

Example Sentence : The man sat on the floor

Training samples : (The, man), (man, sat), (sat, on), (on, the),……..
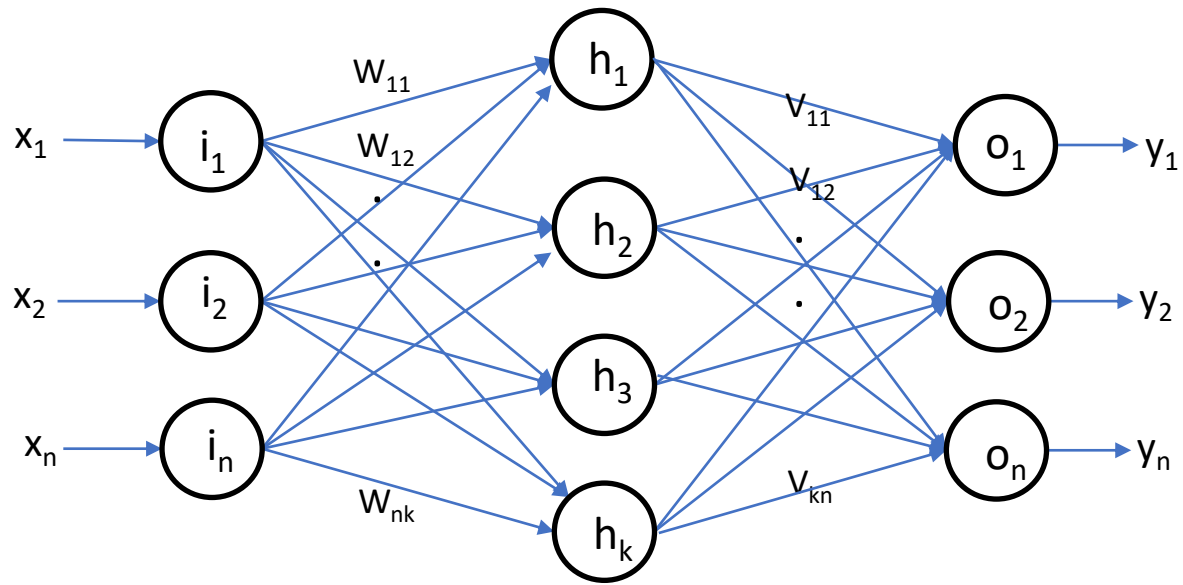
Pr( man | The) = ?

This model can be used for various tasks
- Classification
- Regression
- **Probability estimate**
- Representation
- ….

Example Sentence : The man sat on the floor

Training samples : (The, man), (man, sat), (sat, on), (on, the),……..

Pr( man | The) = ?

# Word2vec – Word Embedding

- Objective : Predict the target word given the words in the context

- Example Sentence : The man sat on the floor

- Target Word: man
- Window Size: 2
- Context Words: The, sat, on

- Training samples : (man, the), (man, sat), (man, on)

# Word2vec – Word Embedding

The man sat on the floor

Training samples : (the, man), (the, sat)

The man sat on the floor

Training samples : (man, the), (man, sat)
                                   (man, on)

The man sat on the floor

Training samples : (sat, the), (sat, man)
                                   (sat, on), (sat, the)

The man sat on the floor

Training samples : (on, man), (on, sat)
                                   (on, the), (on, floor)

# word2vec

**2 basic neural network models:**

- **Continuous Bag of Word (CBOW)**: Pr(target|context)
- **Skip-gram (SG):** Pr(context|target).

# Word2vec – Continuous Bag of Word

- E.g. "The man sat on floor"
  - Window size = 2

Input layer

Index of man in vocabulary

man

Hidden layer

Output layer

one-hot vector

on

sat   one-hot vector

We must learn W and V

Input layer

Hidden layer

Output layer

cat

$W_{n \times d}$

| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| … |
| 0 |

n-dim

$V_{d \times n}$

sat

| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |
| … |
| 0 |

n-dim

on

$W_{n \times d}$

| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| … |
| 0 |

n-dim

d-dim

Linear Activation function

d will be the size of word vector

22

$$W_{n \times d}^{T} \qquad \times \qquad x_{cat} = h_{cat}$$

Input layer

| 0.1 | **2.4** | 1.6 | 1.8 | 0.5 | 0.9 | ... | ... | ... | 3.2 |
|-----|---------|-----|-----|-----|-----|-----|-----|-----|-----|
| 0.5 | **2.6** | 1.4 | 2.9 | 1.5 | 3.6 | ... | ... | ... | 6.1 |
| ... | **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| ... | **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| 0.6 | **1.8** | 2.7 | 1.9 | 2.4 | 2.0 | ... | ... | ... | 1.2 |

$W_{n \times d}^{T} \times x_{cat} = h_{cat}$

$W_{n \times d}^{T} \times x_{on} = h_{on}$

$+$

$\hat{h} = \dfrac{h_{cat} + h_{on}}{2}$

Output layer

$x_{cat}$

n-dim

$x_{on}$

n-dim

Hidden layer

d-dim

sat

n-dim

23

$$W_{n \times d}^T \quad \times \quad x_{on} = h_{on}$$

| 0.1 | 2.4 | 1.6 | **1.8** | 0.5 | 0.9 | … | … | … | 3.2 |
|-----|-----|-----|---------|-----|-----|---|---|---|-----|
| 0.5 | 2.6 | 1.4 | **2.9** | 1.5 | 3.6 | … | … | … | 6.1 |
| … | … | … | **…** | … | … | … | … | … | … |
| … | … | … | **…** | … | … | … | … | … | … |
| 0.6 | 1.8 | 2.7 | **1.9** | 2.4 | 2.0 | … | … | … | 1.2 |

Input layer

$x_{cat}$

n-dim

$W_{n \times d}^T \times x_{cat} = h_{cat}$

$W_{n \times d}^T \times x_{on} = h_{on}$

$x_{on}$

n-dim

$+$

$\hat{h} = \dfrac{h_{cat} + h_{on}}{2}$

Hidden layer
d-dim

Output layer

sat

n-dim

24

$V_{d \times n}^T \qquad \times \qquad \hat{h} \qquad = \qquad z$

| 0.01 | 1.2 | ... | ... | 1.2 |
| 0.51 | 1.6 | ... | ... | 5.1 |
| 1.4 | 2.1 | ... | ... | 6.4 |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 0.26 | 0.8 | ... | ... | 1.7 |

$n \times d$

$\times$

| 2.1 |
| 2.75 |
| ... |
| ... |
| 1.85 |

$d \times 1$

$=$

| 1.2 |
| 7.5 |
| 1.6 |
| ... |
| ... |
| 1.5 |

$n \times 1$

$\hat{y} = softmax(z)$

$\hat{y} = \dfrac{e^z}{\sum e^z}$

Input layer

cat

n-dim

$W_{n \times d}$

on

n-dim

$W_{n \times d}$

Hidden layer

$V_{d \times n}^T \times \hat{h} = z$

$\hat{h}$

d-dim

d will be the size of word vector

Output layer

| 0.0001 |
| 0.002 |
| 0.002 |
| 0.001 |
| 0.7 |
| ... |
| 0.00 |

$\hat{y}_{sat}$

V-dim

$$V^T_{d \times n} \quad \times \quad \hat{h} \quad = \quad z$$

| 0.01 | 1.2 | ... | ... | 1.2 |
|------|-----|-----|-----|-----|
| 0.51 | 1.6 | ... | ... | 5.1 |
| 1.4 | 2.1 | ... | ... | 6.4 |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 0.26 | 0.8 | ... | ... | 1.7 |

$n \times d$

| 2.1 |
|------|
| 2.75 |
| ... |
| ... |
| 1.85 |

$d \times 1$

$\times$

$=$

| 1.2 |
|------|
| 7.5 |
| 1.6 |
| ... |
| ... |
| 1.5 |

$n \times 1$

$$\hat{y} = softmax(z)$$
$$\hat{y} = \frac{e^z}{\sum e^z}$$

Input layer

| 0 |
|---|
| **1** |
| 0 |
| 0 |

cat

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| ... |

n-dim

| 0 |

$W_{n \times d}$

Hidden layer

Output layer

| 0 |
|---|
| 0 |
| 0 |
| **1** |

on

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| ... |

n-dim

| 0 |

$W_{n \times d}$

$\hat{h}$

d-dim

$$V^T_{d \times n} \times \hat{h} = z$$

| 0.0001 |
|--------|
| 0.002 |
| 0.002 |
| 0.001 |
| **0.7** |
| ... |
| 0.00 |

$\hat{y}_{sat}$

V-dim

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| **1** |
| ... |
| 0 |

$y_{sat}$

d will be the size of word vector

$$V_{d \times n}^T \quad \times \quad \hat{h} \quad = \quad z$$

| 0.01 | 1.2 | ... | ... | 1.2 |
|------|-----|-----|-----|-----|
| 0.51 | 1.6 | ... | ... | 5.1 |
| 1.4  | 2.1 | ... | ... | 6.4 |
| ...  | ... | ... | ... | ... |
| ...  | ... | ... | ... | ... |
| 0.26 | 0.8 | ... | ... | 1.7 |

**n× $d$**

| 2.1 |
|------|
| 2.75 |
| ... |
| ... |
| 1.85 |

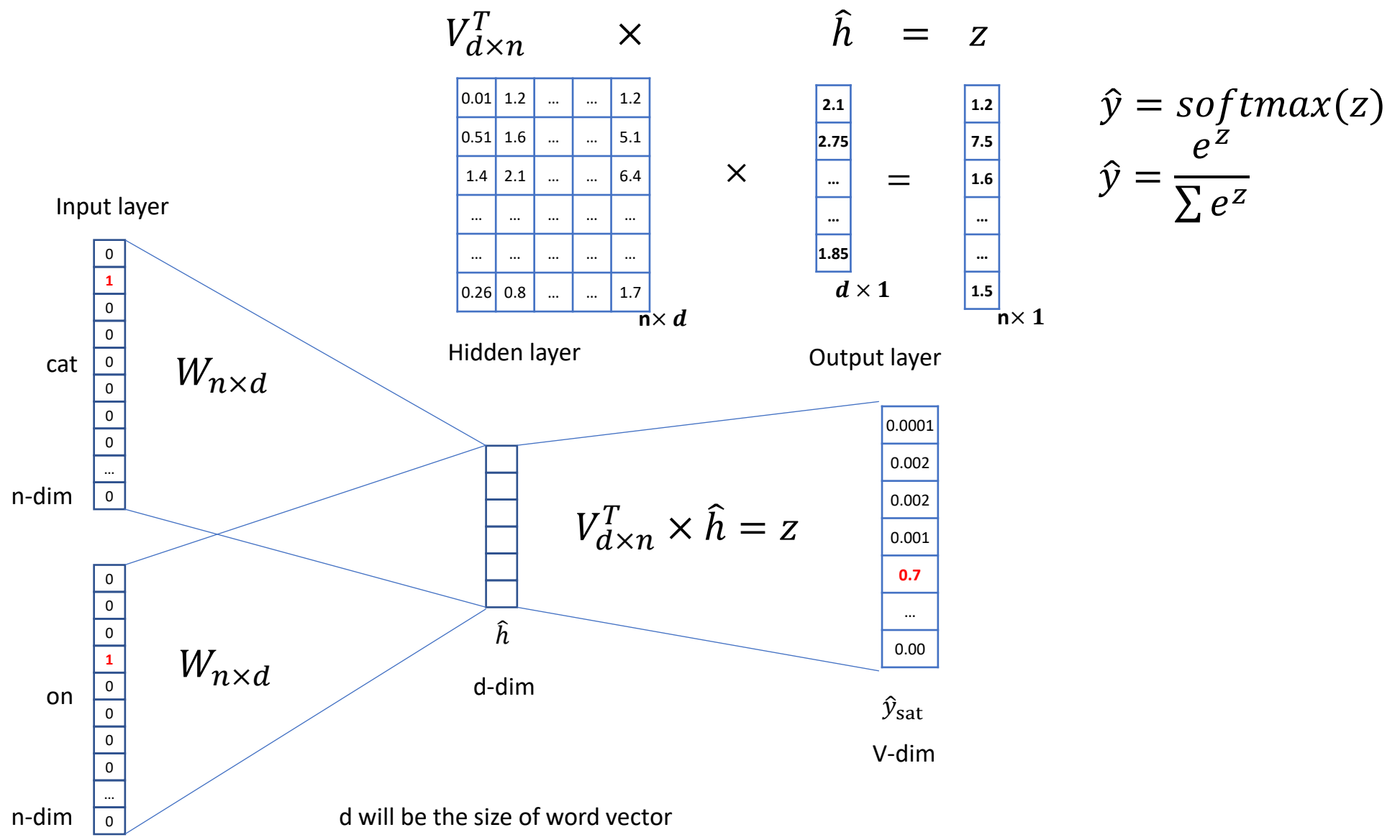**$d \times 1$**

×

=

| 1.2 |
|------|
| 7.5 |
| 1.6 |
| ... |
| ... |
| 1.5 |

**n× 1**

$$\hat{y} = softmax(z)$$

$$\hat{y} = \frac{e^z}{\sum e^z}$$

Input layer

| 0 |
|---|
| **1** |
| 0 |
| 0 |

cat

| 0 |
|---|
| 0 |
| 0 |
| 0 |

n-dim

| ... |
|---|
| 0 |

$W_{n \times d}$

Hidden layer

Output layer

$$V_{d \times n}^T \times \hat{h} = z$$

| 0.0001 |
|---|
| 0.002 |
| 0.002 |
| 0.001 |
| **0.7** |
| ... |
| 0.00 |

$\hat{y}_{sat}$

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| **1** |
| ... |
| 0 |

$y_{sat}$

$$E(\hat{y}, y) = - \sum_{i=1..n} y_i \log \hat{y}_i$$

Loss Function

| 0 |
|---|
| 0 |
| 0 |
| **1** |
| 0 |

on

| 0 |
|---|
| 0 |
| 0 |
| ... |
| 0 |

n-dim

$W_{n \times d}$

$\hat{h}$

d-dim

V-dim

$$W_{n \times d}$$

| 0.1 | **2.4** | ... | ... | 3.2 |
|-----|---------|-----|-----|-----|
| **2.1** | **1.4** | ... | ... | **4.1** |
| 0.5 | **2.6** | ... | ... | 6.1 |
| ... | **...** | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 0.6 | **1.8** | ... | ... | 1.2 |

Embedding of "cat"

Input layer

| 0 |
|---|
| **1** |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| ... |
| 0 |

$x_{cat}$

n-dim

$$W_{n \times d}$$

Output layer

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| **1** |
| ... |
| 0 |

sat

n-dim

$$V_{d \times n}$$

Hidden layer

d-dim

| 0 |
|---|
| 0 |
| 0 |
| **1** |
| 0 |
| 0 |
| 0 |
| 0 |
| ... |
| 0 |

$x_{on}$

n-dim

$$W_{n \times d}$$

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

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| **0.1** | 2.4 | 1.6 | **1.8** | 0.5 | 0.9 | … | … | … | 3.2 |
| **0.5** | 2.6 | 1.4 | **2.9** | 1.5 | 3.6 | … | … | … | 6.1 |
| **…** | … | … | **…** | … | … | … | … | … | … |
| **…** | … | … | **…** | … | … | … | … | … | … |
| **0.6** | 1.8 | 2.7 | **1.9** | 2.4 | 2.0 | … | … | … | 1.2 |

**kxn**

$W^T$

.

| 1 |
|---|
| 0 |
| 0 |
| **0** |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| … |
| 0 |

**x**

=

| **0.1** |
|---|
| **0.5** |
| **…** |
| **…** |
| **0.6** |

**h**

Embedding of $X_1$



$1$ $x_1$ → $i_1$  $W_{11}$ → $h_1$  $V_{11}$ → $o_1$ → $y_1$ $0$

$W_{12}$ → $h_2$  $V_{12}$ → $o_2$ → $y_2$ $1$

$0$ $x_2$ → $i_2$  $h_3$

$0$

$0$ $x_n$ → $i_n$  $W_{nk}$ → $h_k$  $V_{kn}$ → $o_n$ → $y_n$ $0$

$0$

| 0.1 |
| 0.5 |
| ... |
| ... |
| 0.6 |

| 0.01 | 0.51 | 1.4 | ... | ... | ... | 3.2 |
|------|------|-----|-----|-----|-----|-----|
| 1.2  | 1.6  | 2.1 | ... | ... | ... | 6.1 |
| ...  | ...  | ... | ... | ... | ... | ... |
| ...  | ...  | ... | ... | ... | ... | ... |
| 1.2  | 5.1  | 6.4 | ... | ... | ... | 1.2 |

V

30