

HW4 : Updates on website & ed stem

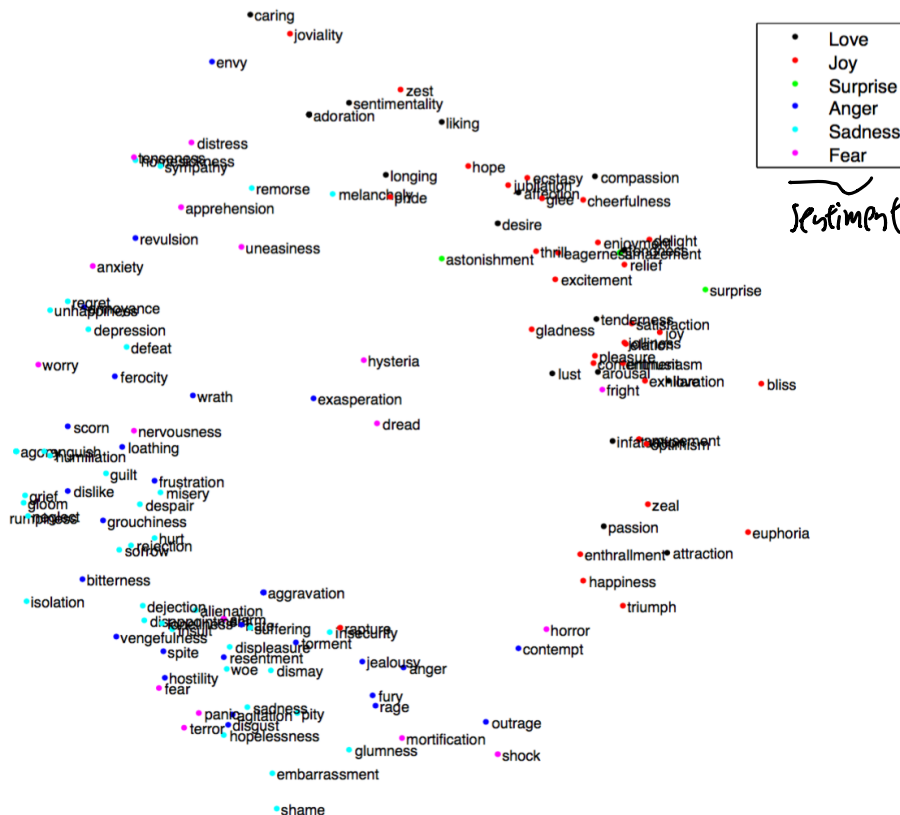
# Feature extraction given Text data

---

W

~~~~~ Verkn

**word2vec** is a popular unsupervised learning approach that just uses a text corpus (e.g. [nytimes.com](https://www.nytimes.com))



# Word embeddings, word2vec

| Source Text                                                                                                             |       | Training Samples |       |                                          |                                                |                                                                  |                                                             |
|-------------------------------------------------------------------------------------------------------------------------|-------|------------------|-------|------------------------------------------|------------------------------------------------|------------------------------------------------------------------|-------------------------------------------------------------|
| <table><tr><td>The</td><td>quick</td><td>brown</td></tr></table> fox jumps over the lazy dog. →                         | The   | quick            | brown | $(x, y)$<br>(the, quick)<br>(the, brown) |                                                |                                                                  |                                                             |
| The                                                                                                                     | quick | brown            |       |                                          |                                                |                                                                  |                                                             |
| <table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td></tr></table> jumps over the lazy dog. →                 | The   | quick            | brown | fox                                      | (quick, the)<br>(quick, brown)<br>(quick, fox) |                                                                  |                                                             |
| The                                                                                                                     | quick | brown            | fox   |                                          |                                                |                                                                  |                                                             |
| <table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td><td>jumps</td></tr></table> over the lazy dog. →         | The   | quick            | brown | fox                                      | jumps                                          | (brown, the)<br>(brown, quick)<br>(brown, fox)<br>(brown, jumps) |                                                             |
| The                                                                                                                     | quick | brown            | fox   | jumps                                    |                                                |                                                                  |                                                             |
| <table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td><td>jumps</td><td>over</td></tr></table> the lazy dog. → | The   | quick            | brown | fox                                      | jumps                                          | over                                                             | (fox, quick)<br>(fox, brown)<br>(fox, jumps)<br>(fox, over) |
| The                                                                                                                     | quick | brown            | fox   | jumps                                    | over                                           |                                                                  |                                                             |

data:

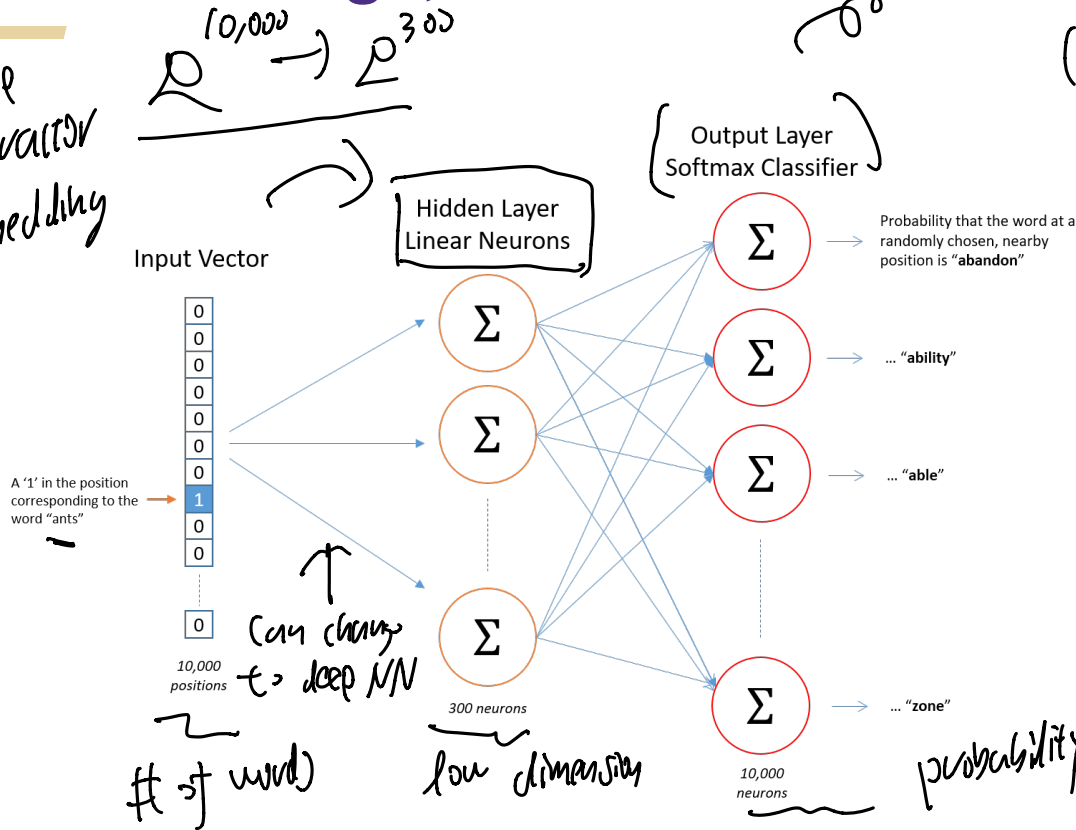
sentence

from wiki,  
usnews, ...

self-supervised learning  
pre-training

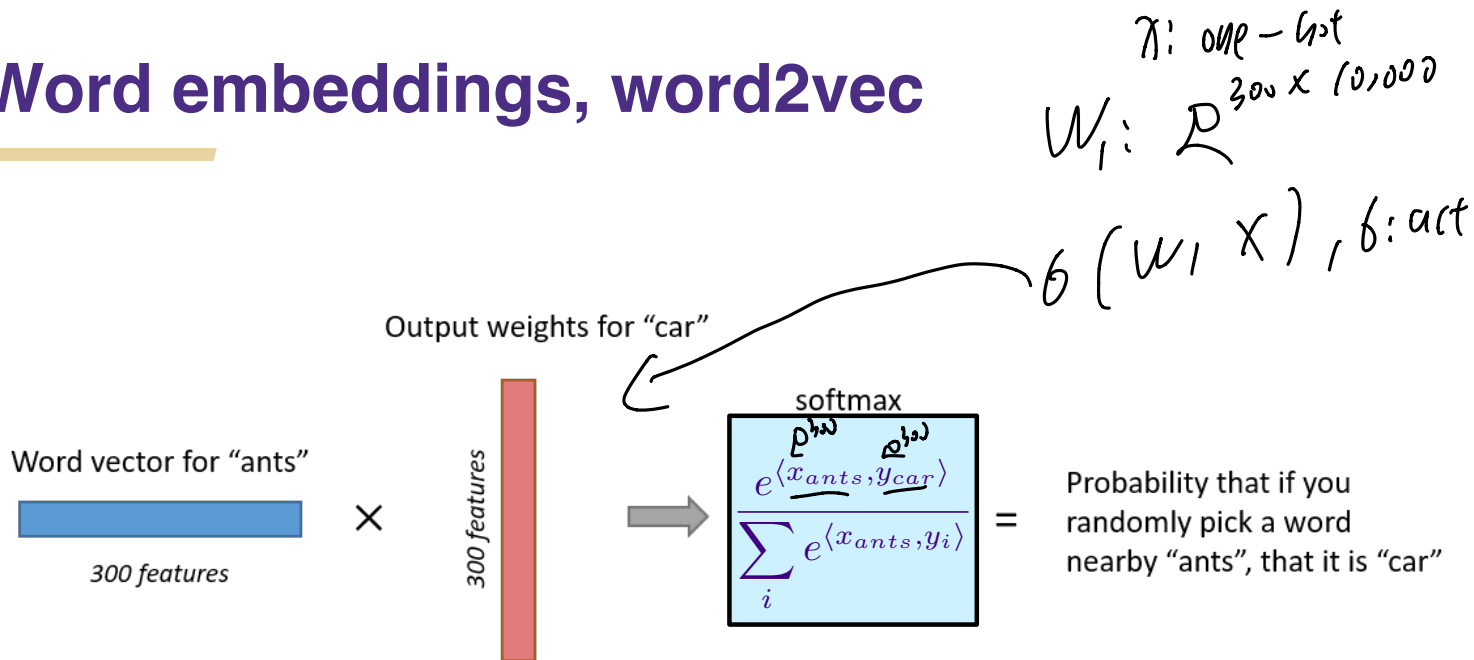
# Word embeddings, word2vec

feature  
extractor  
/ embedding



Training neural network to predict co-occurring words. Use first layer weights as embedding, throw out output layer

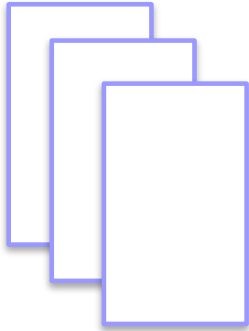
# Word embeddings, word2vec



use: predict the missing word in a sentence

Training neural network to predict co-occurring words. Use first layer weights as embedding, throw out output layer

# Bag of Words



n documents/articles with lots of text

predict sentiment  
(+, -)

Questions:

- How to get a feature representation of each article?
- How to cluster documents into topics?

Bag of words model:

$D$ : size of vocabulary

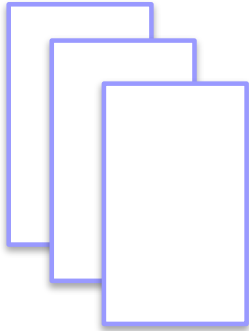
$i$ th document:  $x_i \in \mathbb{R}^D$

$x_{i,j}$  = proportion of times  $j$ th word occurred in  $i$ th document

$$x_{i,j} = \frac{\text{\# of word } j \text{ occurred in document } i}{\text{total \# of words in document } i}$$

# Bag of Words

---



n documents/articles with lots of text

- **Can we embed each document into a feature space?**

Bag of words model:

*i*th document:  $x_i \in \mathbb{R}^D$

$x_{i,j}$  = proportion of times *j*th word occurred in *i*th document

Given vectors, run k-means or Gaussian mixture model to find k clusters/topics

# Nonnegative matrix factorization (NMF)

$m$ : # of articles

topic space:  
"player", "basketball",

$A \in \mathbb{R}^{m \times D}$   $A_{i,j}$  = frequency of  $j$ th word in document  $i$   
 $d \ll m, D$

**Nonnegative  
Matrix factorization:**

$$\min_{W \in \mathbb{R}_+^{m \times d}, H \in \mathbb{R}_+^{d \times D}}$$

$$\|A - WH^T\|_F^2$$



$d$  is number of topics

→ frequencies of words in each topic

Each column of  $H$  represents a cluster of a topic,  
Each row  $W$  is some weights a combination of topics

article

topic space, words

Also see latent Dirichlet factorization (LDA)





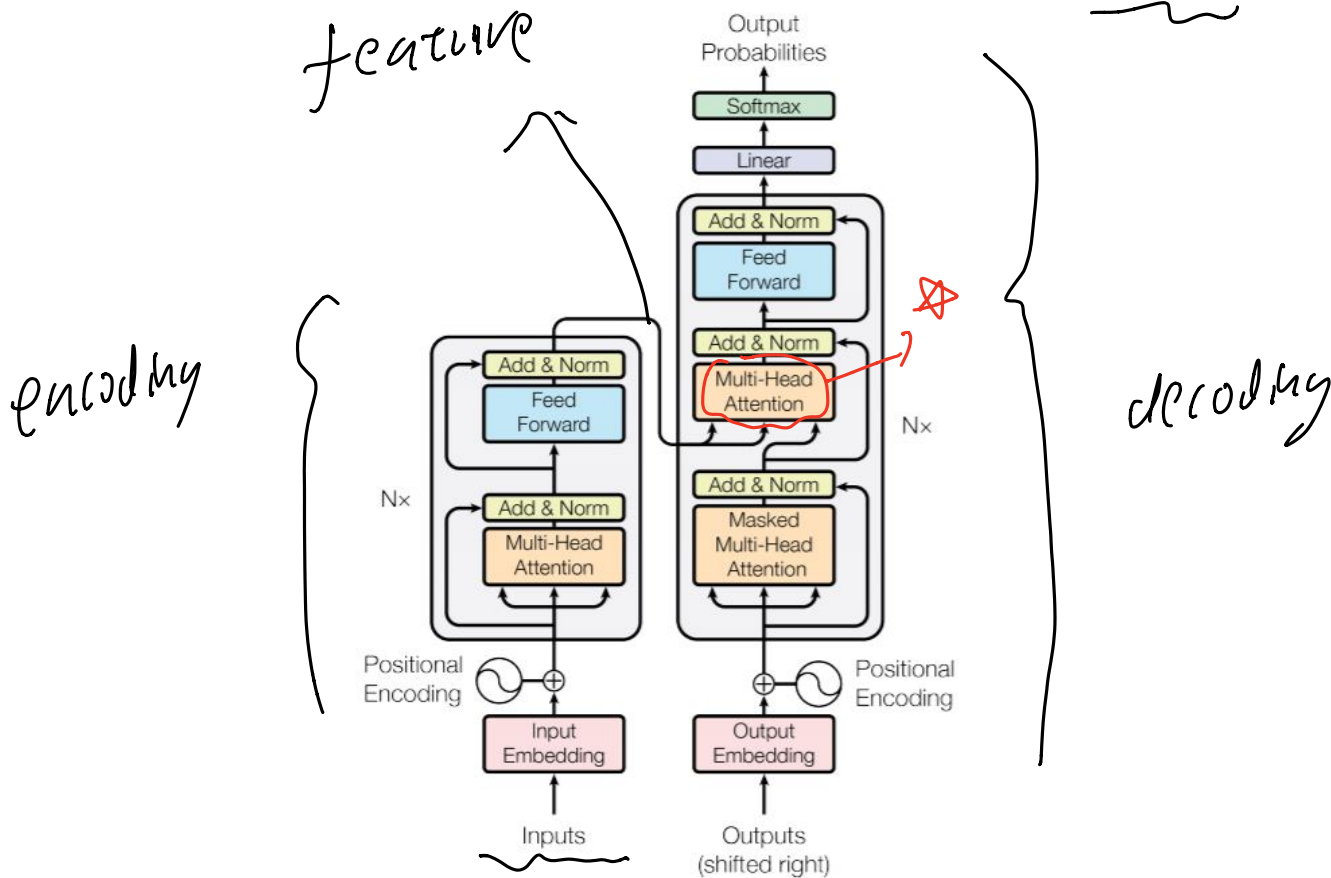
# BERT

2019

data: w/te

masked word modeling

— —   x     x   — —



# Feature extraction given sequential data

---



# Time-dependent data



$x_t \in \mathbb{R}$  : AAPL stock price at time  $t$

Prediction model:  $p(x_{t+1} | \underbrace{x_t, x_{t-1}, x_{t-2}, \dots})$

# Time-dependent data



$x_t \in \mathbb{R}$  : AAPL stock price at time  $t$

$h_t \in \mathbb{R}^d$  : hidden latent state of AAPL

Prediction model:  $p(\underbrace{x_{t+1}}_{\text{predicted}} | \underbrace{x_t, x_{t-1}, x_{t-2}, \dots}_{\text{history}})$   
 $\approx p(\underbrace{x_{t+1} | x_t, h_{t+1}}_{\text{prediction}})$

Markov property  
related reinforcement learning

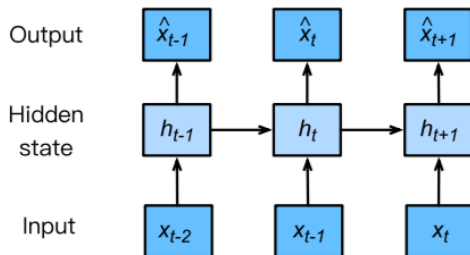
# Time-dependent data



$x_t \in \mathbb{R}$  : AAPL stock price at time  $t$

$h_t \in \mathbb{R}^d$ : hidden latent state of AAPL

Prediction model:  $p(x_{t+1} | x_t, x_{t-1}, x_{t-2}, \dots)$   
 $\approx p(x_{t+1} | x_t, h_{t+1})$



$$\underbrace{h_{t+1}} = \underbrace{g(h_t, x_t)}_{\text{learnable}}, \text{ } h_0 \text{ is known}$$

Hidden state and  $g$  never observed, but learned!

# Time-dependent data

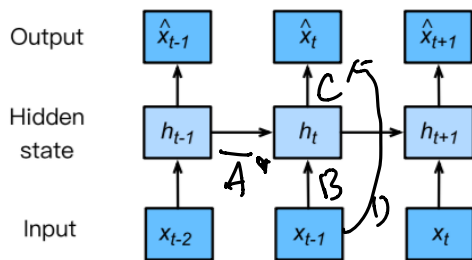
$x_t \in \mathbb{R}$  : AAPL stock price at time  $t$

$h_t \in \mathbb{R}^d$ : hidden latent state of AAPL

$h_t \in \mathbb{R}^d$   
 $A: d \times d, B: d \times 1$   
 $C: 1 \times d, D: 1 \times 1$

Prediction model:  $p(x_{t+1} | x_t, x_{t-1}, x_{t-2}, \dots)$   $h_t = \text{given}$

$$\approx p(x_{t+1} | x_t, h_{t+1}) \quad \{x_1, \dots, x_t\}$$



$$h_{t+1} = g(h_t, x_t) \quad \begin{matrix} \text{predict } x_{t+1} \\ \text{given } x_t, h_t \end{matrix}$$

Hidden state and  $g$  never observed, but learned!

Explicit:

activation function,  $\hat{x}_{t+1} = \mathbb{E}_{p(\cdot | x_t, h_t)} [x_{t+1}]$

hidden dynamics

$$h_{t+1} = \sigma(Ah_t + Bx_t)$$

emission

$$\hat{x}_{t+1} = Ch_{t+1} + Dx_t$$

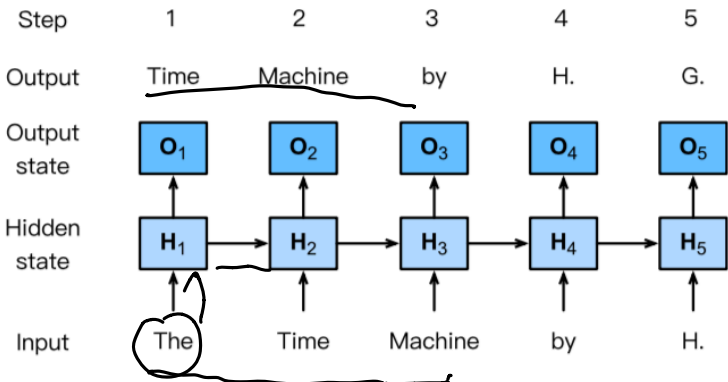
$$\min_{A, B, C, D} \sum_t (x_t - \hat{x}_t)^2$$

Prediction model:  $p(x_{t+1} | x_t, x_{t-1}, x_{t-2}, \dots)$

$$h_{t+1} = g(h_t, x_t)$$

Hidden state and  $g$  never observed, but learned!

## Model also works with text!



# Time-dependent data

Prediction model:  $p(x_{t+1} | x_t, x_{t-1}, x_{t-2}, \dots)$   
 $\approx p(x_{t+1} | x_t, h_{t+1})$

$$h_{t+1} = g(h_t, x_t)$$

Hidden state and  $g$  never observed, but learned!

Recurrent Neural Network

