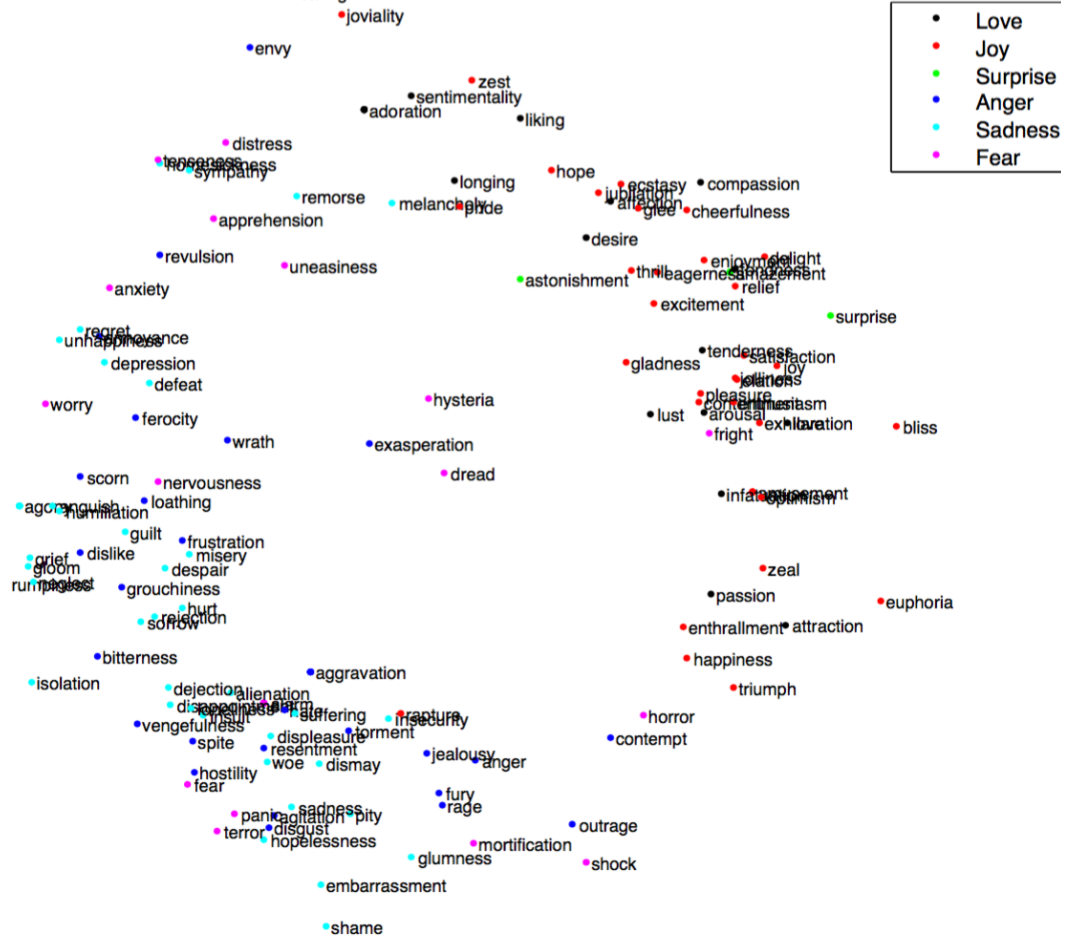


# Feature extraction given Text data

---

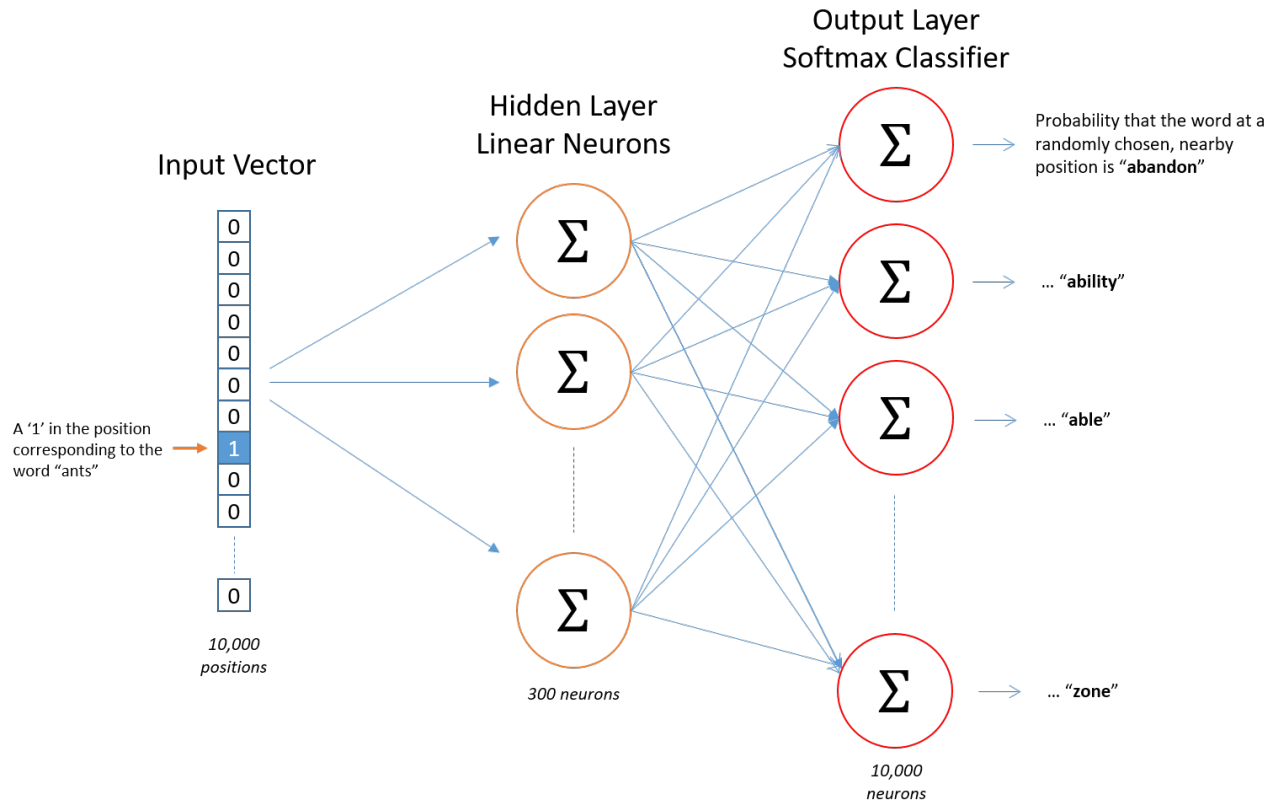




# Word embeddings, word2vec

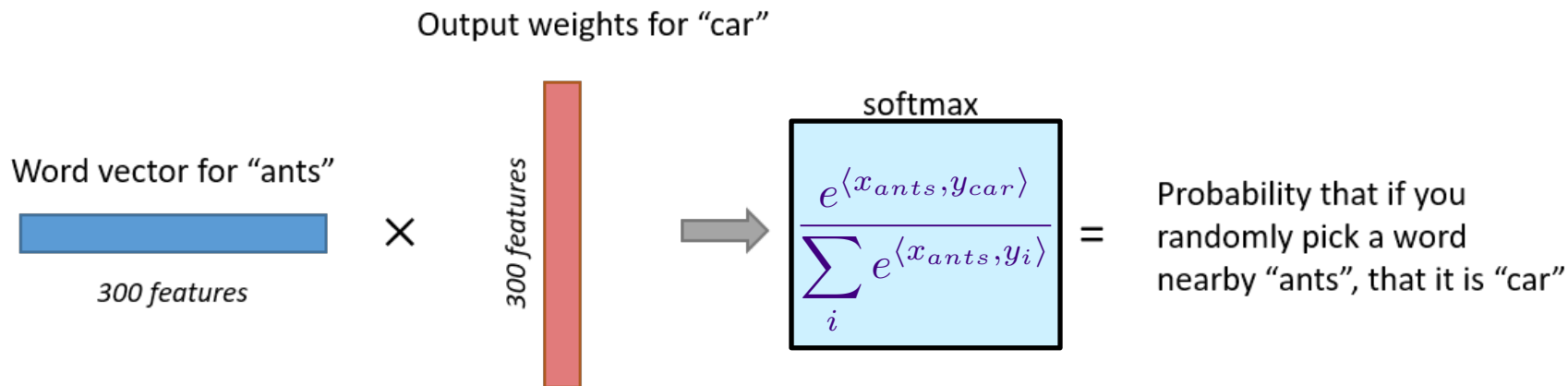
Source Text	Training Samples					
<table><tr><td>The</td><td>quick</td><td>brown</td></tr></table> fox jumps over the lazy dog. ➡	The	quick	brown	(the, quick) (the, brown)		
The	quick	brown				
The <table><tr><td>quick</td><td>brown</td><td>fox</td></tr></table> jumps over the lazy dog. ➡	quick	brown	fox	(quick, the) (quick, brown) (quick, fox)		
quick	brown	fox				
The quick <table><tr><td>brown</td><td>fox</td><td>jumps</td></tr></table> over the lazy dog. ➡	brown	fox	jumps	(brown, the) (brown, quick) (brown, fox) (brown, jumps)		
brown	fox	jumps				
The <table><tr><td>quick</td><td>brown</td><td>fox</td><td>jumps</td><td>over</td></tr></table> the lazy dog. ➡	quick	brown	fox	jumps	over	(fox, quick) (fox, brown) (fox, jumps) (fox, over)
quick	brown	fox	jumps	over		

# Word embeddings, word2vec



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

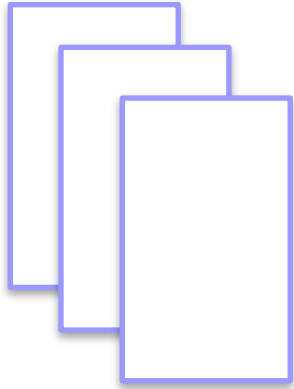
# Word embeddings, word2vec



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

Questions:

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

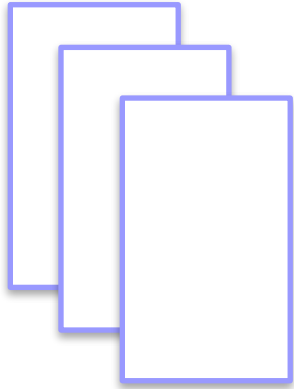
Bag of words model:

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

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

# 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)

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

**Nonnegative  
Matrix factorization:**

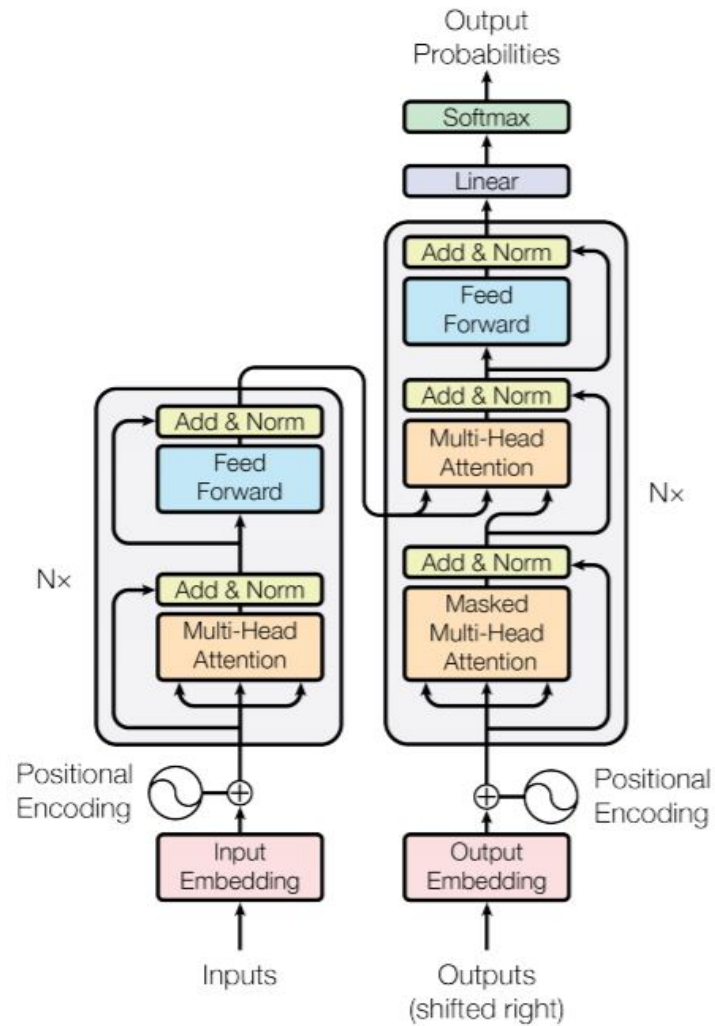
$$\min_{W \in \mathbb{R}_+^{m \times d}, H \in \mathbb{R}_+^{n \times d}} \|A - WH^T\|_F^2$$

$d$  is number of topics

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

Also see latent Dirichlet factorization (LDA)

# BERT



# Feature extraction given sequential data

---



# Time-dependent data



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

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

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

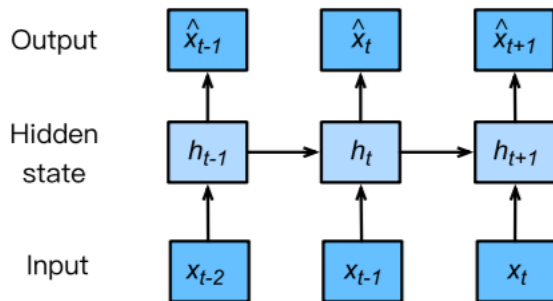
# Time-dependent data



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

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

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



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

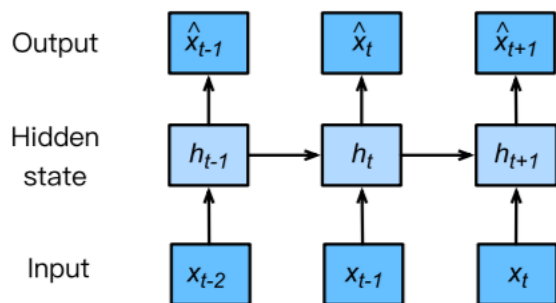
Hidden state and  $g$  never observed, but learned!

# 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

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!

Explicit:

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

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

$$\sum_t (x_t - \hat{x}_t)^2$$

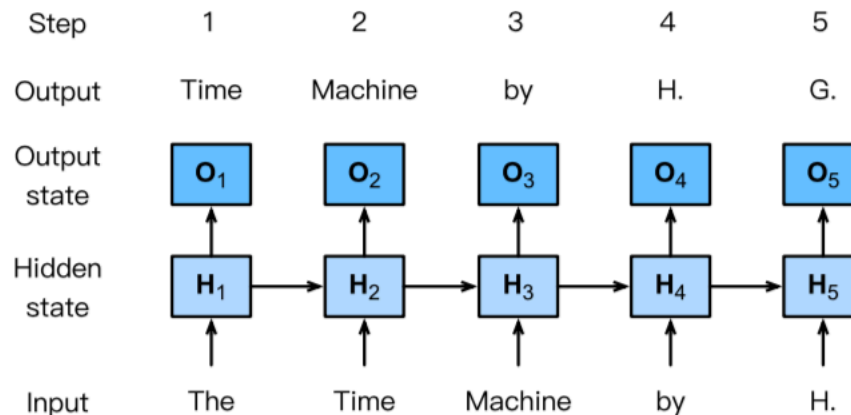
# Time-dependent data

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!

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

