




# DS 4440 Embeddings

## Today

- Our first foray into representation learning
- How to encode TEXT
- Word2Vec (and beyond)

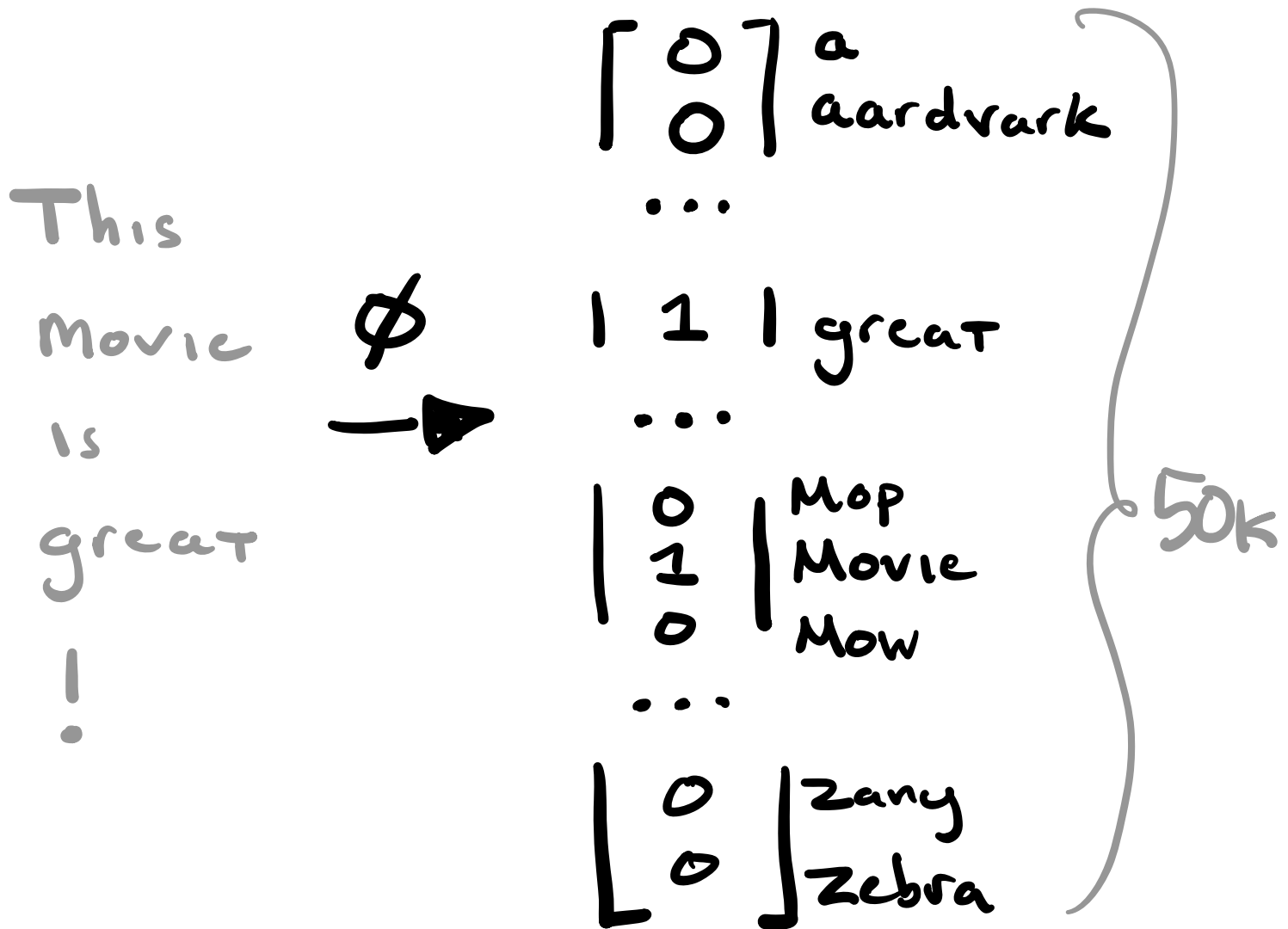
[  ]  
"Cat"

[  ]  
"Dog"

[  ]  
"Pancake"

Wayyy back ~ 2013

## Bag-of-Words (BoW)



Word  $w$  → Position  $w_{idx}$

$$x_{w_{idx}} = \begin{cases} 1 & \text{if } w \text{ in input} \\ 0 & \text{otherwise} \end{cases}$$

Long, sparse X vectors.

## Problems with BoW?

We often want to measure  
Text Similarity

$$\text{Cos}(a, b) = \frac{\phi(a)^T \phi(b)}{\|\phi(a)\| \|\phi(b)\|}$$

For cat, dog

$$\begin{array}{c} \text{"cat"} \\ [0 \dots 1 \dots 0] \cdot \begin{array}{c} [0] \\ \dots \\ 1 \text{ "dog"} \\ \dots \\ 0 \end{array} = 0. \end{array}$$

So

$$\text{Sim}(\text{cat}, \text{dog}) =$$

$$\text{Sim}(\text{cat}, \text{pancake}) = 0.$$

## Distributional Semantics

"you shall know a word  
by the company it keeps"



Words to "low dimensional"  
embeddings

# Similar words ~ Nearby

Word2Vec [Mikolov 2013]

Target word and Context

Two variants / objectives

## Skip-gram

The man his son



loves

$$P(\text{The} | \text{loves}) \cdot$$

$$P(\text{Man} | \text{loves}) \cdot$$

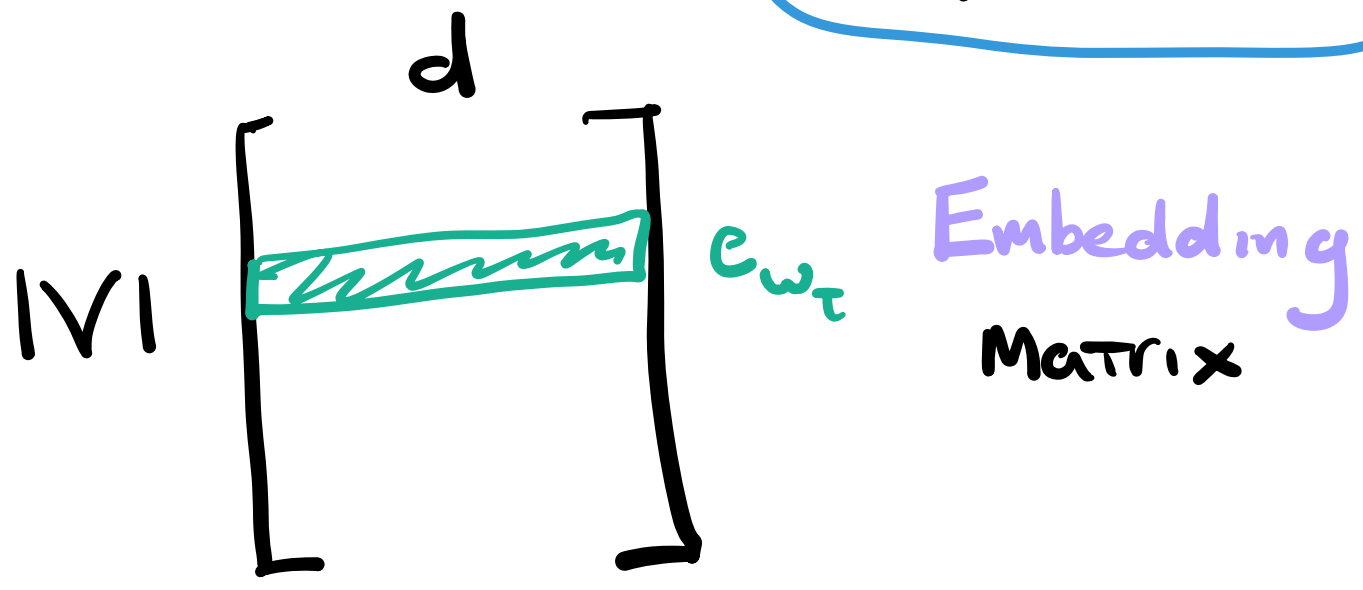
$$P(\text{his} | \text{loves}) \cdot$$

$$P(\text{son} | \text{loves})$$

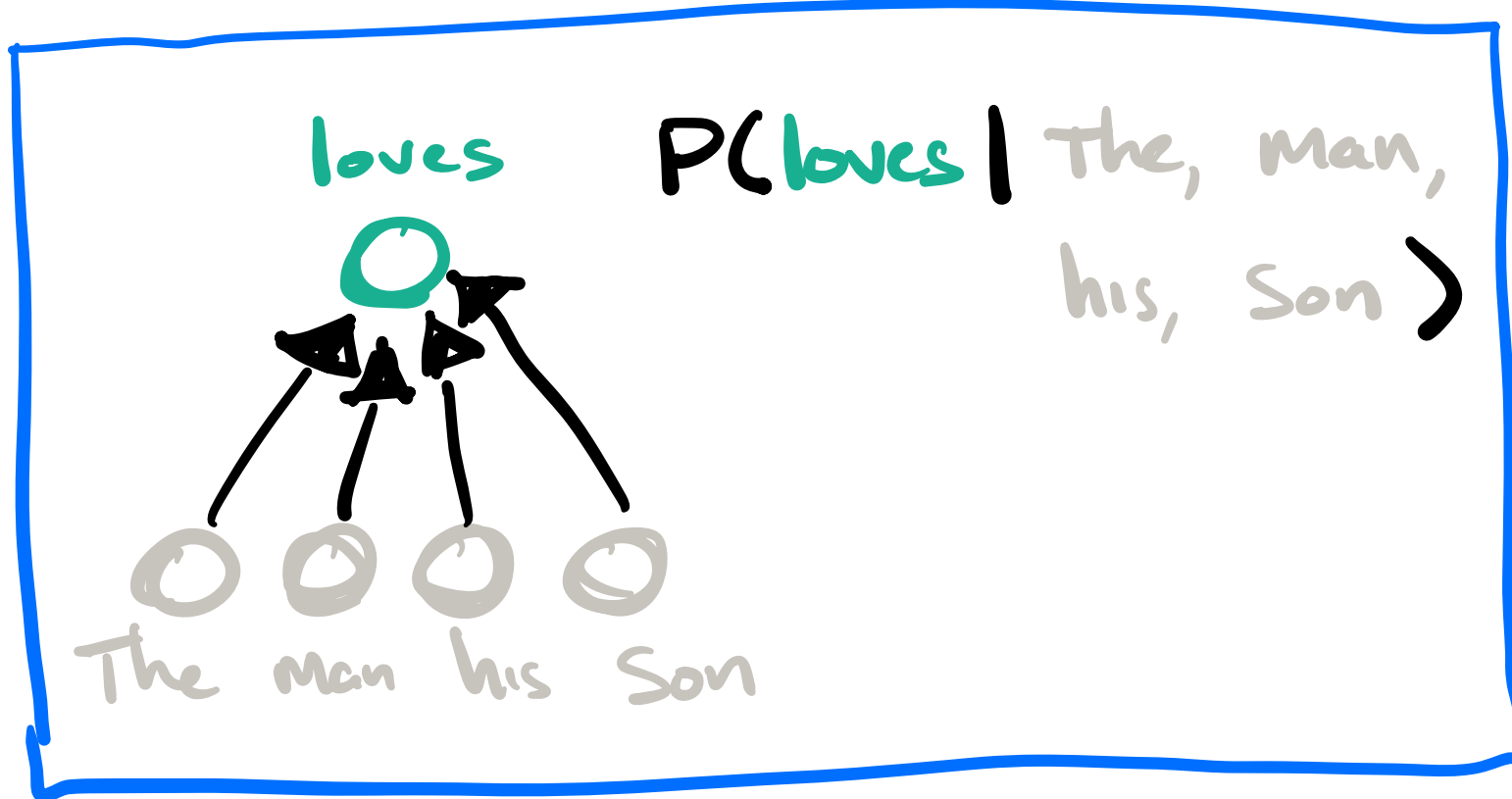
Loss: 
$$-\prod_{w_c} \prod_{w_c \text{ in context}} P(w_c | w_t)$$

$$P(w_c | w_\tau) \equiv \frac{\exp\{e_{w_\tau} \cdot e_{w_c}\}}{\sum_{\tilde{w} \in V} \exp\{e_{w_\tau} \cdot e_{\tilde{w}}\}}$$

Softmax!



# CBOW



$$\text{Loss: } -P(\omega_\tau | \omega_{c_1} \dots \omega_{c_k})$$

$$\equiv \frac{\exp\{\epsilon_{\omega_\tau} \cdot \bar{\omega}_c\}}{\sum_{\tilde{\omega} \in V} \exp\{\epsilon_{\tilde{\omega}} \cdot \bar{\omega}_c\}}$$

Where:

$$\bar{\omega}_c \equiv \frac{1}{k} \sum_{j=1}^k \epsilon_{\omega_{c_j}}$$

OK, BUT  $V$  is BIG!

→ Very slow

Idea Negative Sampling

First Sample

$(\omega_\tau, \omega_{c_1} \dots \omega_{c_k})$

# Skip-gram

↑ Maximize Sim b/w  $w_\tau, w_{cl}$   $S_1$   
(for all  $l$ )

↓ Minimize Sim b/w  $w_\tau$  and  
Words from other (random)  
CONTEXTS  $S_2$

$$S_1 \uparrow \sigma(e_{w_\tau} \cdot e_{w_{cl}})$$

$$S_2 \downarrow \sigma(e_{w_\tau} \cdot e_{\tilde{w}})$$

$$\text{Loss} = -S_1 + S_2$$



# CBow

↑ Maximize Sim b/w  $w_t, \bar{w}_c$

↓ Minimize Sim b/w  $\tilde{w}, \bar{w}_c$   
random words

$$S_1 \uparrow \sigma(e_{w_t} \cdot e_{\bar{w}_c})$$

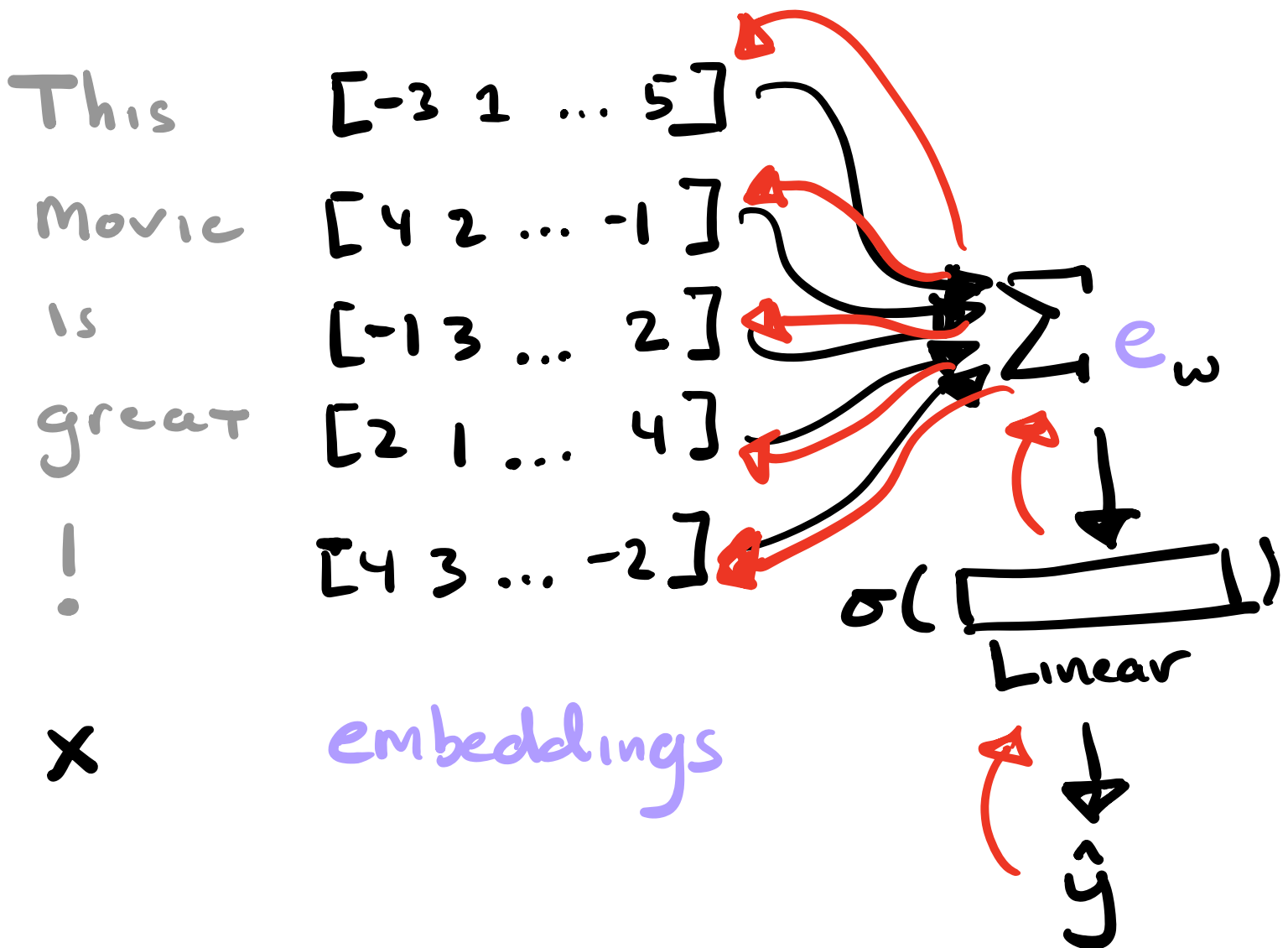
$$S_2 \downarrow \sigma(e_{\tilde{w}} \cdot e_{\bar{w}_c})$$

$$\text{Loss} = -S_1 + S_2$$

Let's see in Colab...

# How to use embeddings (for TEXT Classification)?

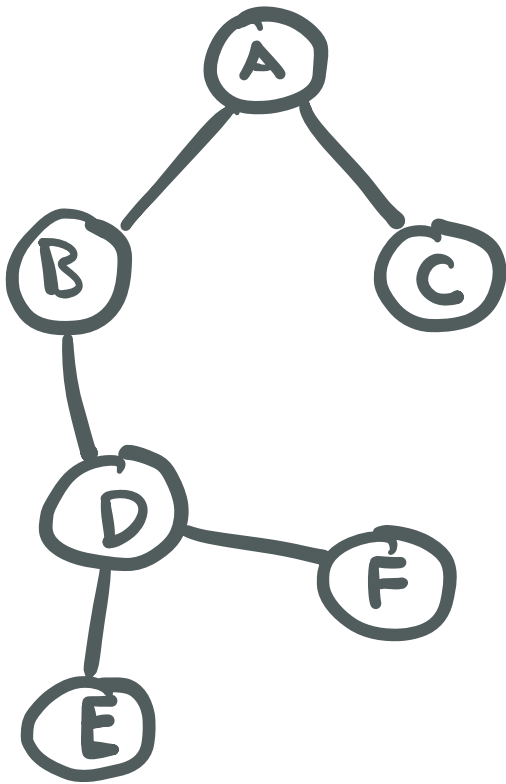
One way: Adopt CBOW



# Beyond "Words"

Embeddings not just for NLP!

Example DeepWalk embeds graph nodes (Perozzi et al.)



- Sample walks randomly
- Treat as "sentences"

A B D E

A C

...

B D F

Back to NLP, Can Treat  
Paragraphs or documents as  
Special "Words"

