# NLP and Word Embeddings

deeplearning.ai

## Learning word embeddings

# Neural language model

I  want  a  glass  of  orange  juice .

4343 9665  1  3852  6163  6257   apple juice.

$e_{4343} = E \, o_{4343}$

| | | | |
|---|---|---|---|
| I | $o_{4343}$ | $\rightarrow$ $E$ $\rightarrow$ | $e_{4343}$ |
| want | $o_{9665}$ | $\rightarrow$ $E$ $\rightarrow$ | $e_{9665}$ |
| a | $o_{1}$ | $\rightarrow$ $E$ $\rightarrow$ | $e_{1}$ |
| glass | $o_{3852}$ | $\rightarrow$ $E$ $\rightarrow$ | $e_{3852}$ |
| of | $o_{6163}$ | $\rightarrow$ $E$ $\rightarrow$ | $e_{6163}$ |
| orange | $o_{6257}$ | $\rightarrow$ $E$ $\rightarrow$ | $e_{6257}$ |

Softmax

10,000

$W^{[1]}, b^{[1]}$

$W^{[2]}, b^{[2]}$

1800 → 1200

[Bengio et. al., 2003, A neural probabilistic language model]

Andrew Ng

# Other context/target pairs

I want a [glass] of [orange] juice to go along with my cereal.

context    target

Context: Last 4 words.

- 4 words on left & right
- Last 1 word
- Nearby 1 word

a glass of orange ___? to go alg with

orange ___?

glass ___?

skip gram

Andrew Ng

# NLP and Word Embeddings

## Word2Vec

# Skip-grams

I want a glass of orange juice to go along with my cereal.

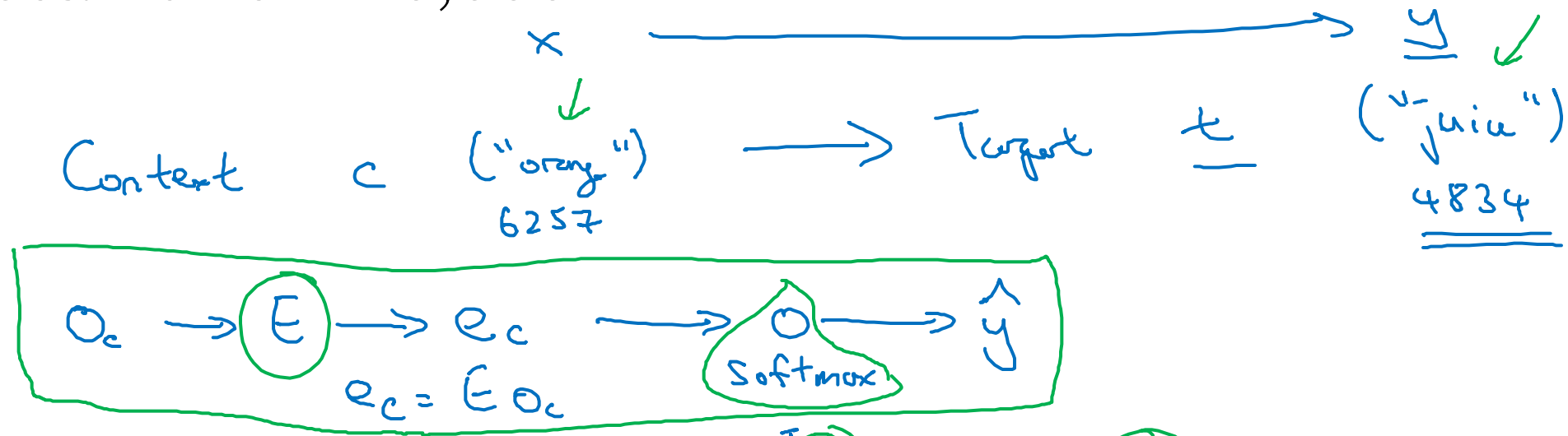| Context | Target |
|---------|--------|
| orange  | juice  |
| orange  | glass  |
| orange  | my     |

[Mikolov et. al., 2013. Efficient estimation of word representations in vector space.]

Andrew Ng

# Model

Vocab size = 10,000k

$x$

Context $c$ ("orang") $\longrightarrow$ Target $t$ $\longrightarrow y$ ("juice")

6257

4834

$$O_c \rightarrow \boxed{E} \rightarrow e_c \rightarrow \bigcirc \rightarrow \hat{y}$$

$e_c = E O_c$

Softmax

Softmax: $p(t|c) = \dfrac{e^{\theta_t^T e_c}}{\sum\limits_{j=1}^{10,000} e^{\theta_j^T e_c}}$

$\Theta_t$ = parameter associated with output $t$

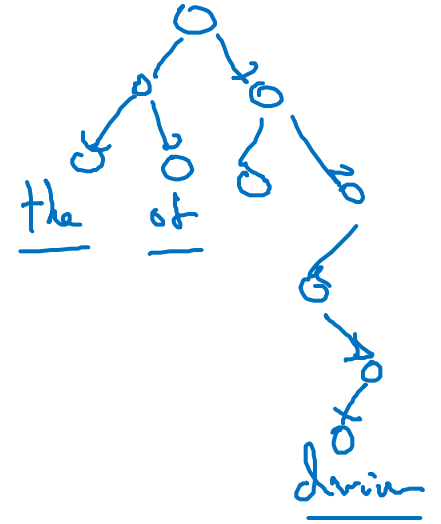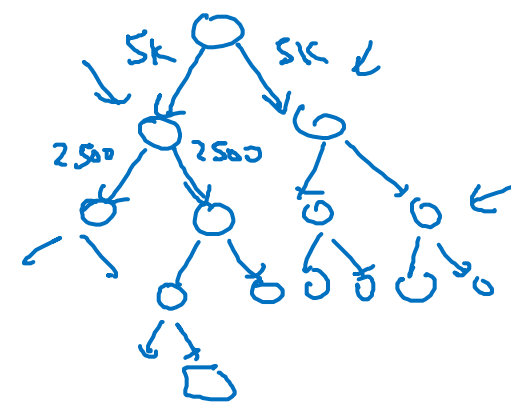$\rightarrow \mathcal{L}(\hat{y}, y) = -\sum\limits_{i=1}^{10,000} y_i \log \hat{y}_i$

$$y = \begin{bmatrix} 0 \\ \vdots \\ i \\ \vdots \\ 0 \end{bmatrix} \leftarrow 4834$$

Andrew Ng

# Problems with softmax classification

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

Hierachial    softmax.

$\log |V|$

5k    5k

2500    2500

the    of

durian

How to sample the context $c$?

$\rightarrow$ the, of, a, and, to, ...

$\rightarrow$ orange, apple, durian

durian

$t$

$c \rightarrow t$

$P(c)$

# NLP and Word Embeddings

## Negative sampling

deeplearning.ai

# Defining a new learning problem

I want a glass of orange juice to go along with my cereal.



$x \longrightarrow y$

| Context | word | target? |
|---------|------|---------|
| orange | juice | 1 |
| orange | king | |
| orange | book | |
| orange | the | |
| orange | of | |

$k = 5 - 20$  smaller. datasets

$k = 2 - 5$  larger dataset

[Mikolov et. al., 2013. Distributed representation of words and phrases and their compositionality]
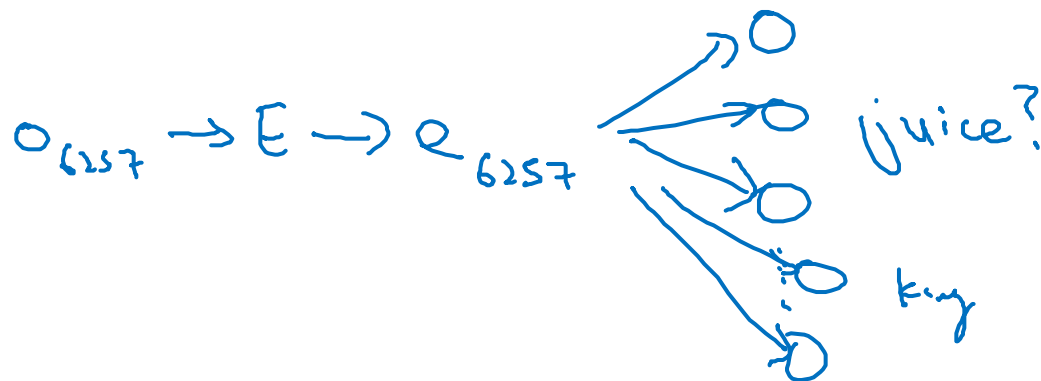
Andrew Ng

# Model

Softmax: $$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

10,000-way softmax

$$P(y=1 \mid c, t) = \sigma\left(\theta_t^T e_c\right)$$

Orange
6257

$O_{6257} \rightarrow E \rightarrow e_{6257}$

juice?

king

10,000

| | $x$ | | $y$ |
| context | word | target? |
| --- | --- | --- |
| orange | juice | 1 |
| orange | king | 0 |
| orange | book | 0 |
| orange | the | 0 |
| orange | of | 0 |

$c$     $t$     $y$

10,000 binary classification problem

$k+1$

# Selecting negative examples

$\overbrace{\phantom{xxxxxxxxxxxxxxxxxxxxx}}^{c}$ $\overbrace{\phantom{xxxxxxx}}^{y}$

| context | word | target? |
|---------|------|---------|
| orange | juice | 1 |
| orange | king | 0 |
| orange | book | 0 |
| orange | the | 0 |
| orange | of | 0 |

$\uparrow$
$t$

the , of, and, ...

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}}$$

$$\frac{1}{|V|}$$
$\uparrow$

Andrew Ng

# GloVe (global vectors for word representation)

I want a glass of orange juice to go along with my cereal.

$c, t$

$X_{ij}$ = # times $i$ appears in context of $j$.

$$X_{ij} = X_{ji}$$

[Pennington et. al., 2014. GloVe: Global vectors for word representation]

Andrew Ng

# Model

$$\text{Minimize} \quad \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) \left( \Theta_i^T e_j + b_i + b_j' - \log X_{ij} \right)^2 \quad \overset{>0?}{\longleftarrow}$$

$t \downarrow \quad c \downarrow$ (above $\Theta_i^T e_j$)

$t \downarrow \quad c \downarrow$ (above $b_i + b_j'$)

$> 0?$ (above right, pointing to $\log X_{ij}$)

$$\overset{t \quad c}{\text{``} \Theta_t^T e_c \text{''}}$$

weighting term

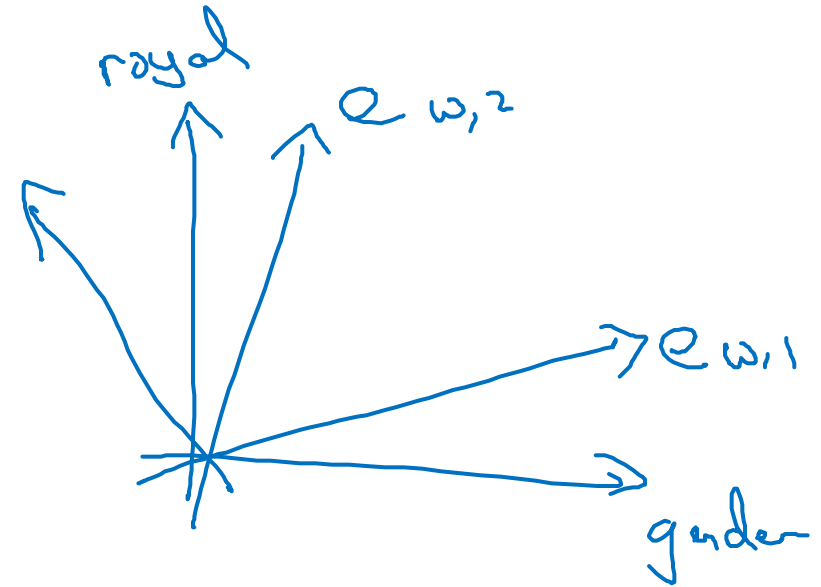$f(X_{ij}) = 0$ at $X_{ij} = 0$.      "$0 \log 0$" $= 0$

this, is, of, a, ....

durian

$\Theta_i$, $e_j$ are symmetric

$$e_w^{(final)} = \frac{e_w + \Theta_w}{2}$$

Andrew Ng

# A note on the featurization view of word embeddings

|          | Man (5391) | Woman (9853) | King (4914) | Queen (7157) |
|----------|------------|--------------|-------------|--------------|
| Gender   | $-1$       | 1            | -0.95       | 0.97         |
| Royal    | 0.01       | 0.02         | 0.93        | 0.95         |
| Age      | 0.03       | 0.02         | 0.70        | 0.69         |
| Food     | 0.09       | 0.01         | 0.02        | 0.01         |

$$\text{minimize} \ \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij})\left(\theta_i^T e_j + b_i - b_j' - \log X_{ij}\right)^2$$

$(A\Theta_i)^T (A^{-T} e_j) = \theta_i^T A^T A^{-T} e_j$

royal

$e_{w,2}$

$e_{w,1}$

gender

Andrew Ng