# Word Embeddings

- 1-hot representation

- word embedding

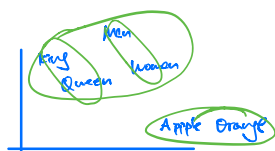e.g.

| | 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 | 0.7 | 0.69 |
| ⋮ | | | | |

↑ 300 ↓

$e_{5391}$   $e_{9853}$

300D → 2D Visualization



---

NER example



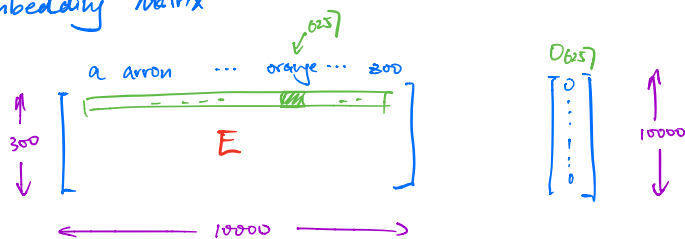| Sally | Johnson | is | an | orange | farmer |
| Robert | Lin | is | a | durian | cultivator |

⇒ transfer learning

1. learn word embeddings from large text corpus   (1 - 100B words)
   (or download pre-trained embedding)

2. transfer embedding to new task w/ smaller training set
   (say 100 k words)

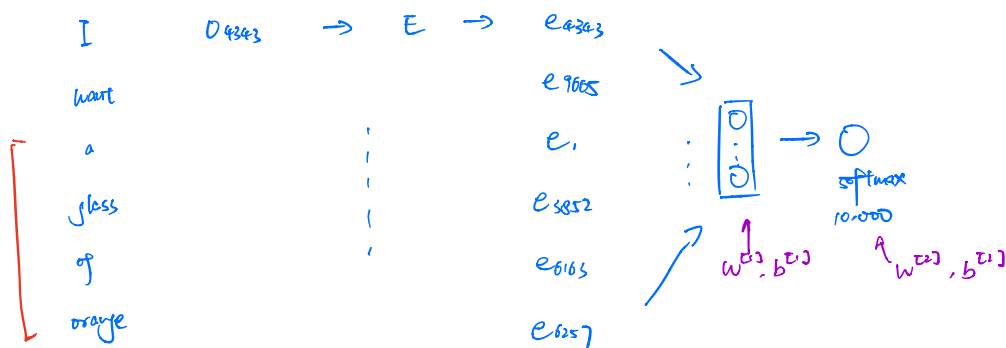3 Optional. Fine-tune word embeddings w/ new data

# Embedding Matrix

$$\underset{(300,10k)}{E} \cdot \underset{(10k,1)}{O_{6257}} = \begin{bmatrix} \boxed{12} \\ \end{bmatrix}_{(300,1)} = e_{6157}$$

---

# Learning word embeddings

e.j.    I    want    a    glass    of    orange    ——— .

4343   9665   1   3852   6163   6257

I        $O_{4343}$ → E → $e_{4343}$

want                        $e_{9665}$

a                           $e_{1}$          → softmax
                                               10,000
glass                       $e_{3852}$

of                          $e_{6163}$        $W^{[1]}, b^{[1]}$    $W^{[2]}, b^{[2]}$

orange                      $e_{6257}$

                            $300 \times 6 = 1800 \rightarrow 1800$ window

backprop to learn E

—  Context :    last   4   words

                4 words   on   left & right

                last   1   word

                nearby   1   word

## Word 2 Vec

### Skip - Gram

e.g.   I   want   a   glass   of   <u>orange</u>   juice   to   go   along   my   cereal.

|  |  |
|---|---|
| <u>Context</u> | <u>target</u> | ← sampled from nearby window |
| orange | juice |
| orange | glass |
| orange | my |

Model:

Vocab   size   =   10,000   K

Context   c ("orange")   →   Target   t   ("juice")
             6257                              4834

$$O_c \rightarrow \boxed{E} \rightarrow e_c \rightarrow \text{softmax} \rightarrow \hat{y} \leftarrow \text{one-hot}$$
$$e_c = E O_c$$

Softmax:   $P(t|c) = \dfrac{e^{\theta_t^T e_c}}{\sum\limits_{j=1}^{10000} e^{\theta_j^T e_c}}$        $\theta_t$ = param  associated w/ output t

$$\mathcal{L}(\hat{y}, y) = - \sum_{i=1}^{10000} y_i \log \hat{y}_i$$

Problem :

Calc   denominator   of   softmax   is   slow

Solution :

hierachical   softmax
negative   sampling

---

### Negative Sampling

e.g.   I   want   a   glass   of   <u>orange</u>   juice   to   go   along   my   cereal.

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

$\left\{\begin{array}{l} \uparrow \\ k \\ \downarrow \end{array}\right.$  (x brace over Context/word, y brace over target)

$k = 5 - 20$   smaller dataset

$k = 2 - 5$   larger dataset

Softmax :  $P(t|c) = \dfrac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$   $\Big\}$  10000-way softmax

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

much faster

$O_{6257} \rightarrow E \rightarrow e_{6257} \Rightarrow \begin{matrix} 0 \\ 0 \\ \vdots \\ 0 \end{matrix}$   juice ?

10,000 logistic regression

Select  negative sampling

$P(w_i) = \dfrac{f(w_i)^{3/4}}{\sum_{j=1}^{10000} f(w_j)^{3/4}}$   $\leftarrow$  between uniform distribution & word distribution

## GloVe word vectors

e.g.  I want a glass of <u>orange</u> juice to go along my cereal.

$c, t$     $\underset{c}{t}$      $\underset{i}{c}$

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

$x_{ij} = x_{ji}$

Symmetric $\Rightarrow e_w^{(final)} = \dfrac{e_w + \theta_w}{2}$

Model :

minimize  $\displaystyle\sum_{i=1}^{10000} \sum_{j=1}^{10000} f(x_{ij}) \,(\,\theta_i^T e_j + b_i + b_j' - \log x_{ij}\,)^2$

weighting term  $f(x_{ij}) = 0$ if $x_{ij} = 0$