

Deep learning

Dr. Aissa Boulmerka
a.boulmerka@centre-univ-mila.dz

2023-2024

CHAPTER 9

NLP AND WORD EMBEDDINGS

Word representation

$V = [a, aaron, \dots, zulu, \langle \text{UNK} \rangle]$

$|V| = 10000$

1-hot representation

Man (5391) Woman (9853) King (4914) Queen (7157) Apple (456) Orange (6257)

$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$
---	---	---	---	---	---

O_{5391}

O_{9853}

O_{4914}

O_{7157}

O_{456}

O_{6257}

I want a glass of orange juice.

I want a glass of apple ?.

Featurized representation: word embedding

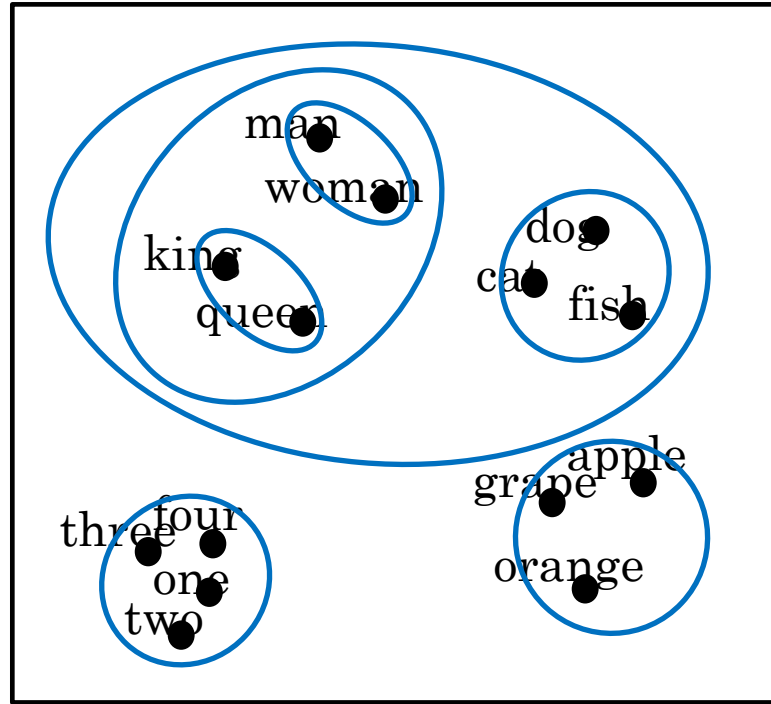
	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97
⋮						
Size						
Cost						

I want a glass of orange juice.

I want a glass of apple juice.

300

Visualizing word embeddings

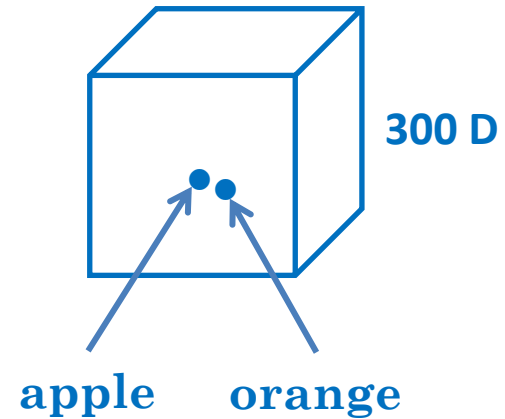


t-SNE

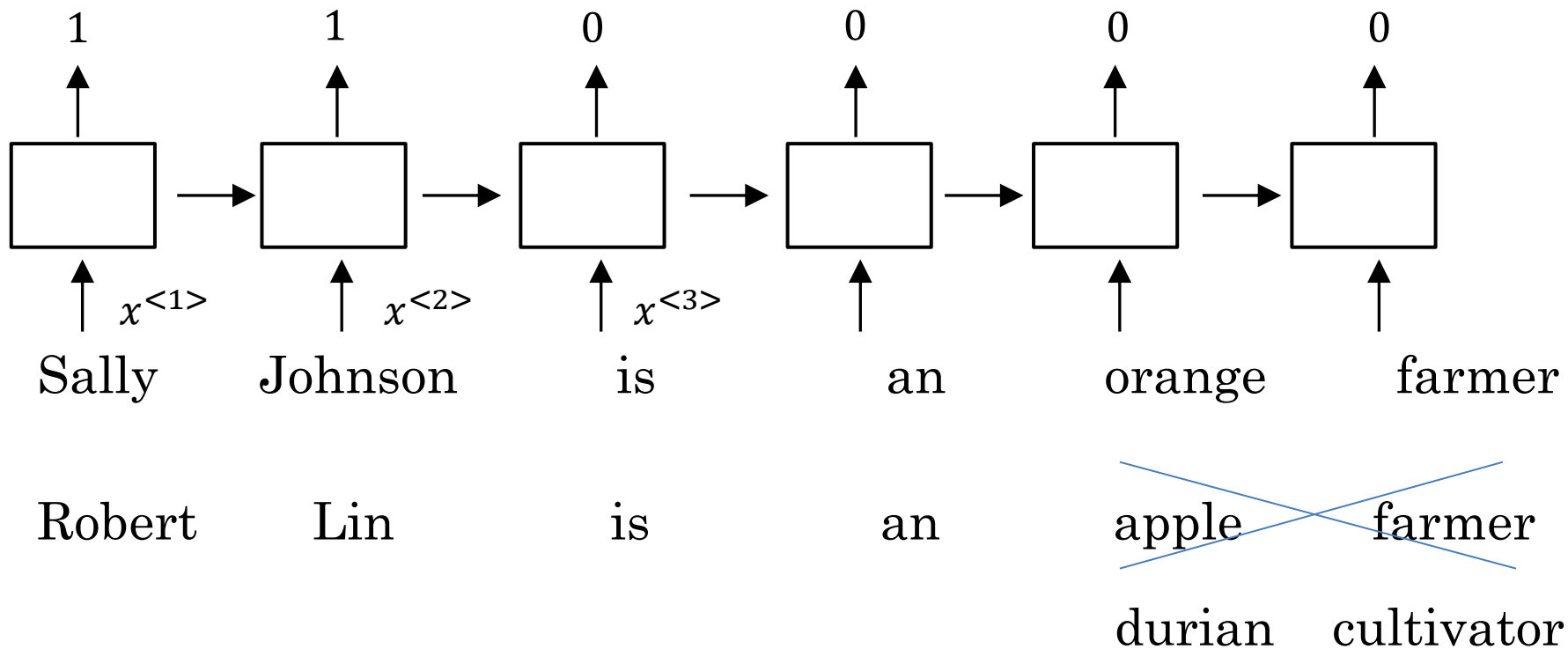
300 D



2 D



Named entity recognition example



BRNN

Transfer learning and word embeddings

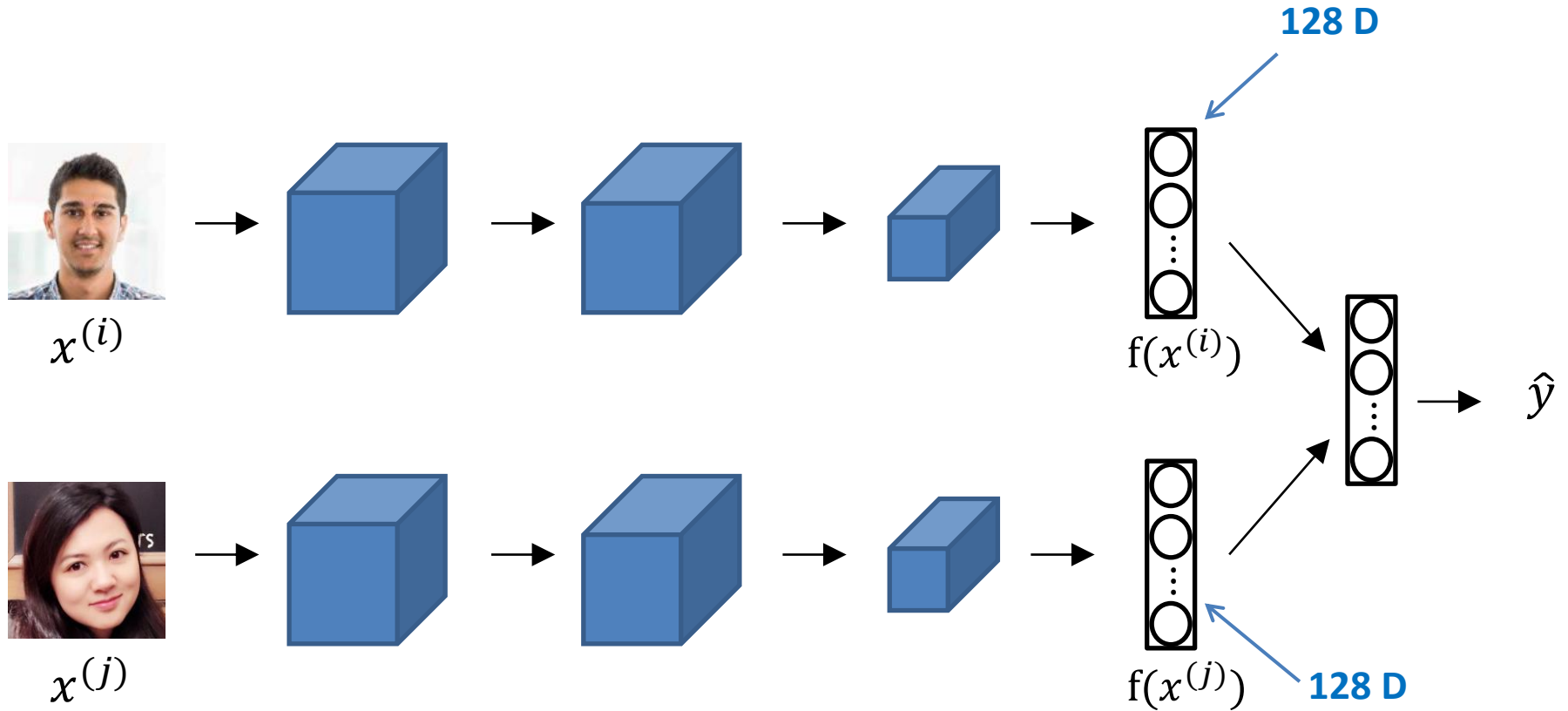
1. Learn word embeddings from large text corpus. (1-100B words)

(Or download pre-trained embedding online.)

2. Transfer embedding to new task with smaller training set.
(say, 100k words)

3. Optional: Continue to finetune the word embeddings with new data.

Relation to face embedding



Analogies

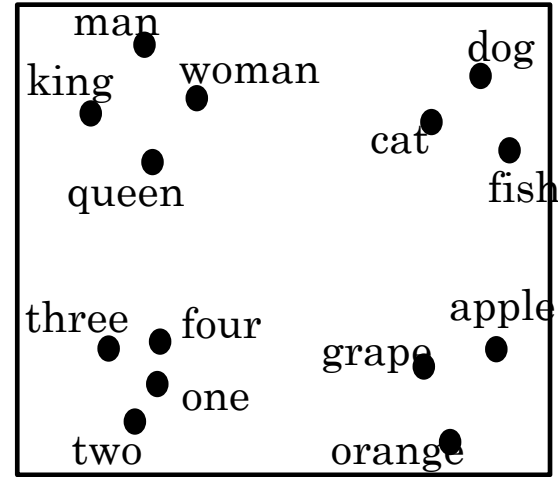
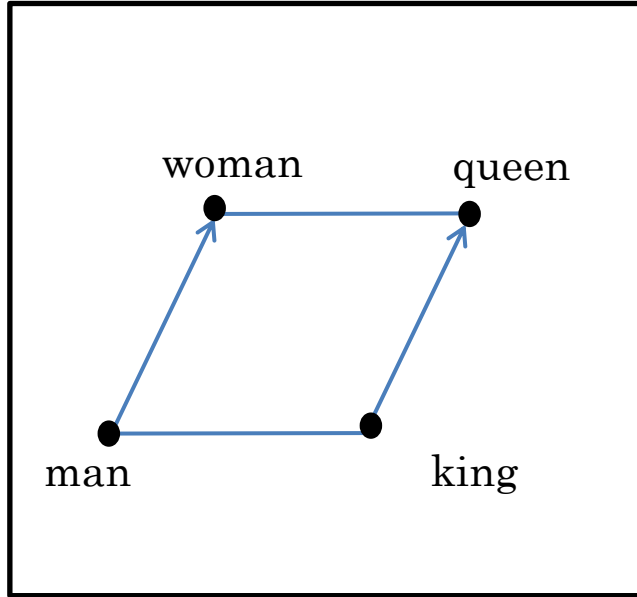
	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

$$e_{man} - e_{woman} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$e_{king} - e_{queen} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$e_{man} - e_{woman} \approx e_{king} - e_{queen}$$

Analogies using word vectors



t-SNE

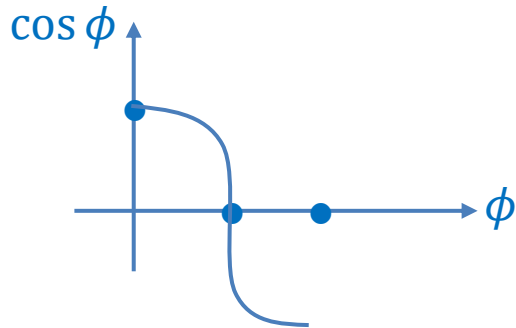
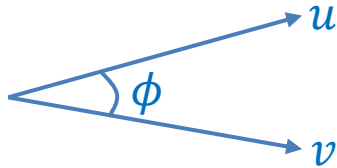
$$e_{man} - e_{woman} \approx e_{king} - e_w \quad ?$$

Find word w : $\underset{w}{\operatorname{argmax}} \operatorname{sim}(e_w, e_{king} - e_{man} + e_{woman})$

Cosine similarity

$$\text{sim}(e_w, e_{king} - e_{man} + e_{woman})$$

$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$



Man:Woman as Boy:Girl

Ottawa:Canada as Nairobi:Kenya

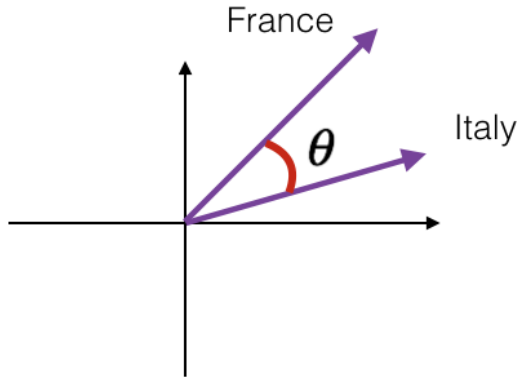
Big:Bigger as Tall:Taller

Yen:Japan as Ruble: Russia

Cosine similarity

$$\text{sim}(e_w, e_{king} - e_{man} + e_{woman})$$

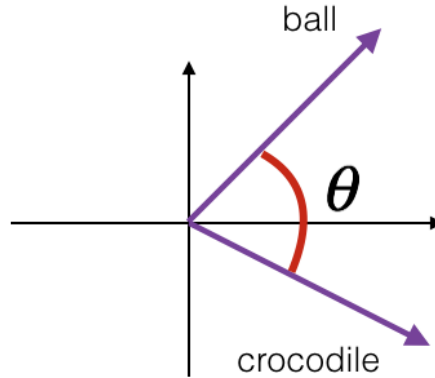
$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$



France and Italy are quite similar

θ is close to 0°

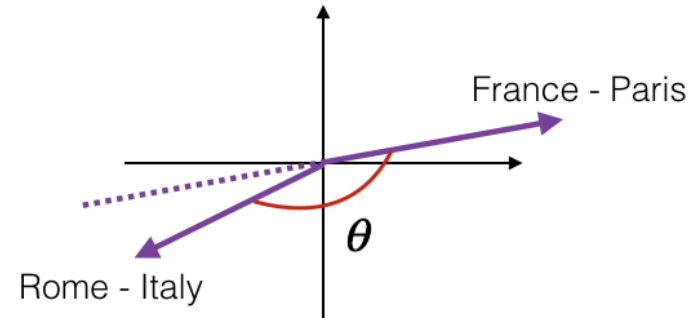
$\cos(\theta) \approx 1$



ball and crocodile are not similar

θ is close to 90°

$\cos(\theta) \approx 0$

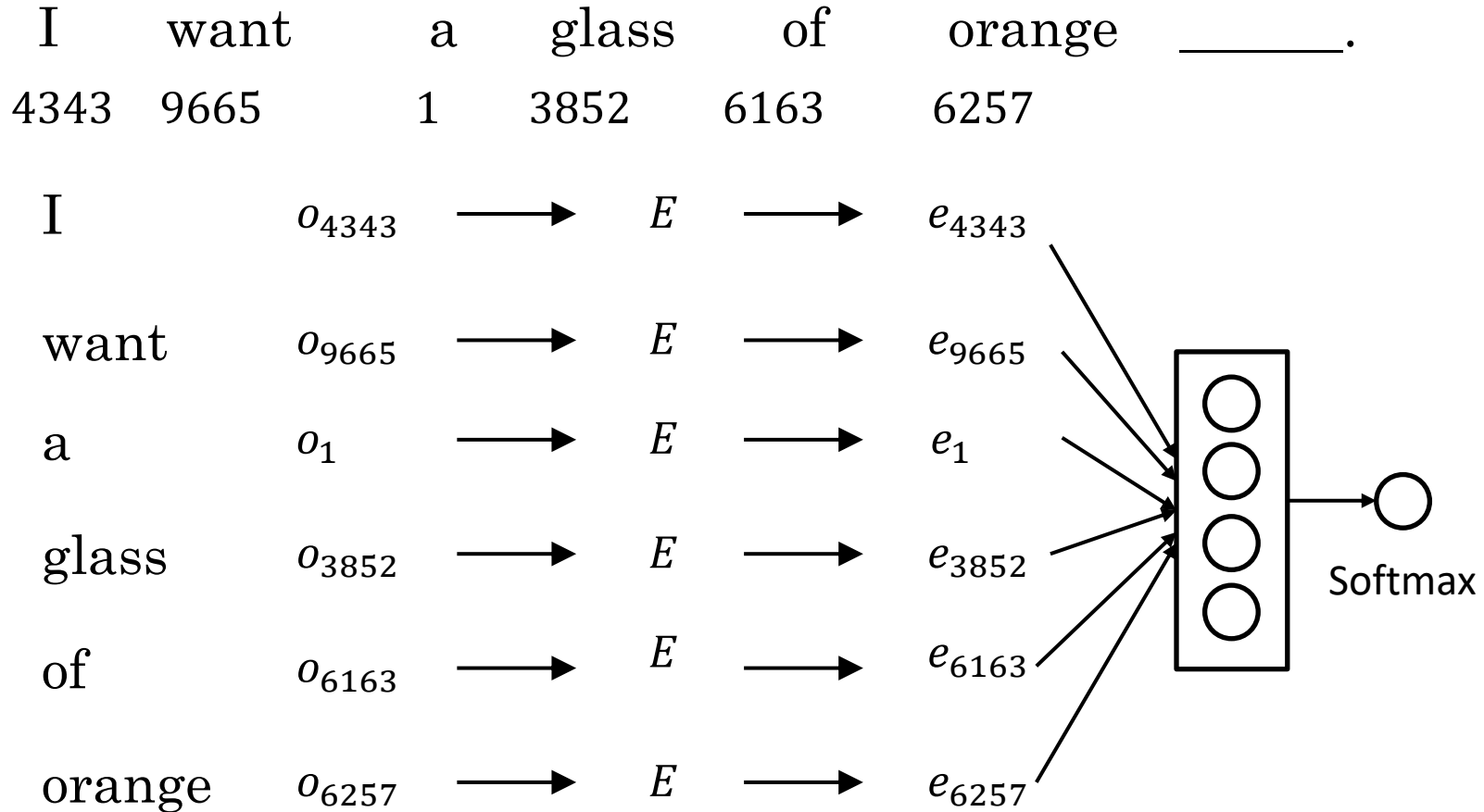


the two vectors are similar but opposite
the first one encodes (city - country)
while the second one encodes (country - city)

θ is close to 180°

$\cos(\theta) \approx -1$

Neural language model



Other context/target pairs

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

Context: Last 4 words.

4 words on left & right

Last 1 word

Nearby 1 word

WORD2VEC

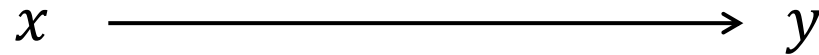
Skip-grams

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

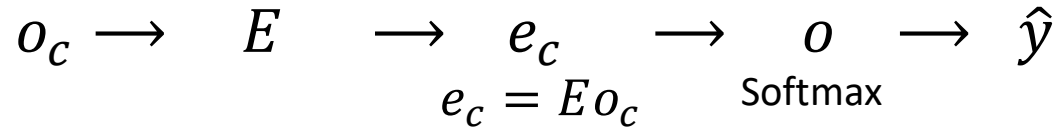
<u>Context</u>	<u>Target</u>
orange	juice
orange	glass
orange	my

Model

Vocab size = 10k



Context c (« orange ») Target t (« juice »)
6257 4834



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

θ_t : parameter associated with the output t

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

$$y = \begin{bmatrix} \vdots \\ 1 \\ \vdots \end{bmatrix} \longleftarrow 4834$$

Problems with softmax classification

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

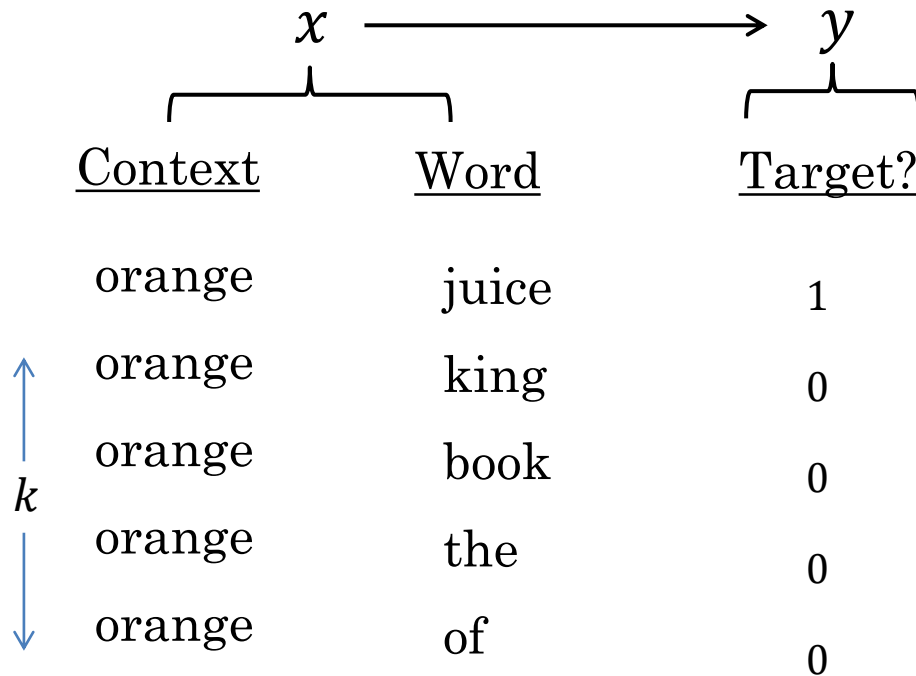
Very slow to compute !!!!

How to sample the context c ?

NEGATIVE SAMPLING

Defining a new learning problem

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



$k = 5 - 20$: smaller dataset

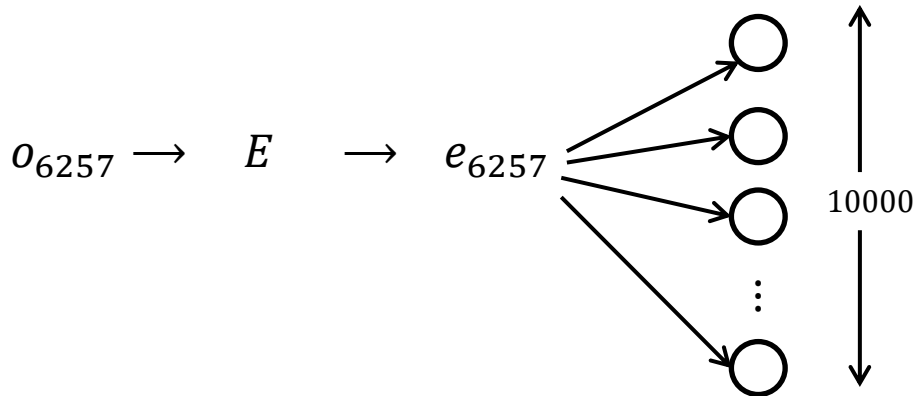
$k = 2 - 5$: larger dataset

Model

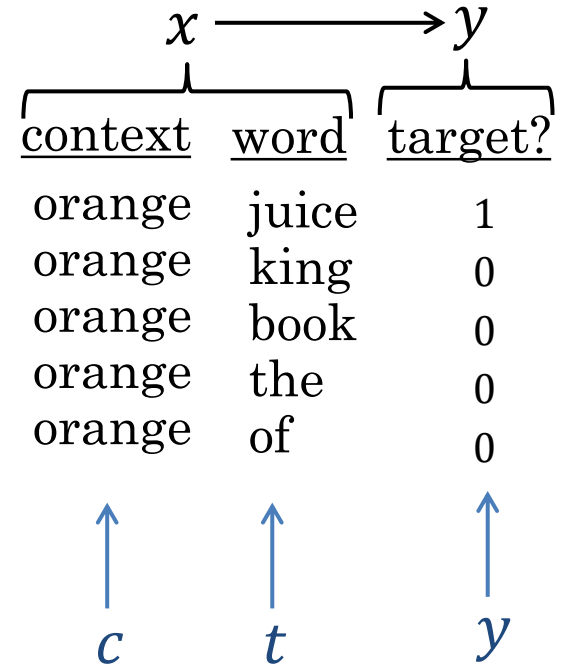
Softmax:
$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

Softmax: 10000 classes

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



10000 binary classification problem



Selecting negative examples

<u>context</u>	<u>word</u>	<u>target?</u>
orange	juice	1
orange	king	0
orange	book	0
orange	the	0
orange	of	0

$$p(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10000} f(w_j)^{3/4}}$$

GLOVE WORD VECTORS

GloVe (global vectors for word representation)

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

c, t

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

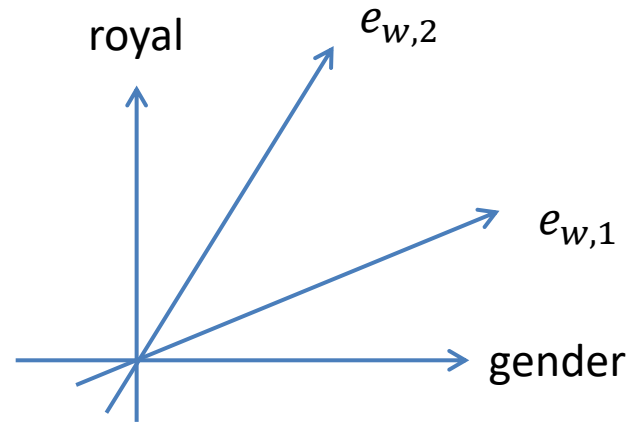
$$X_{ij} = X_{ji}$$

Model

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

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



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

SENTIMENT CLASSIFICATION

Sentiment classification problem

x

The dessert is excellent.

Service was quite slow.

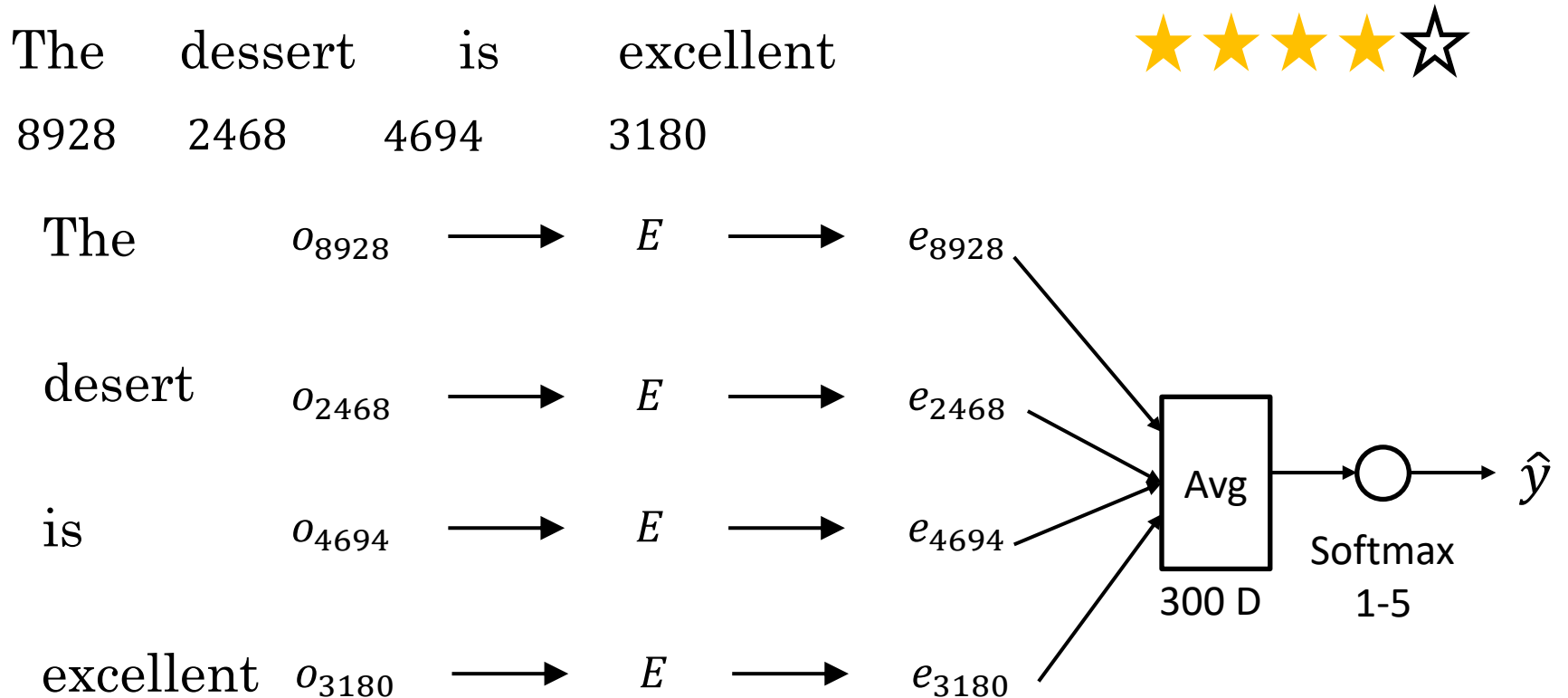
Good for a quick meal, but nothing special.

Completely lacking in good taste, good service, and good ambience.

y

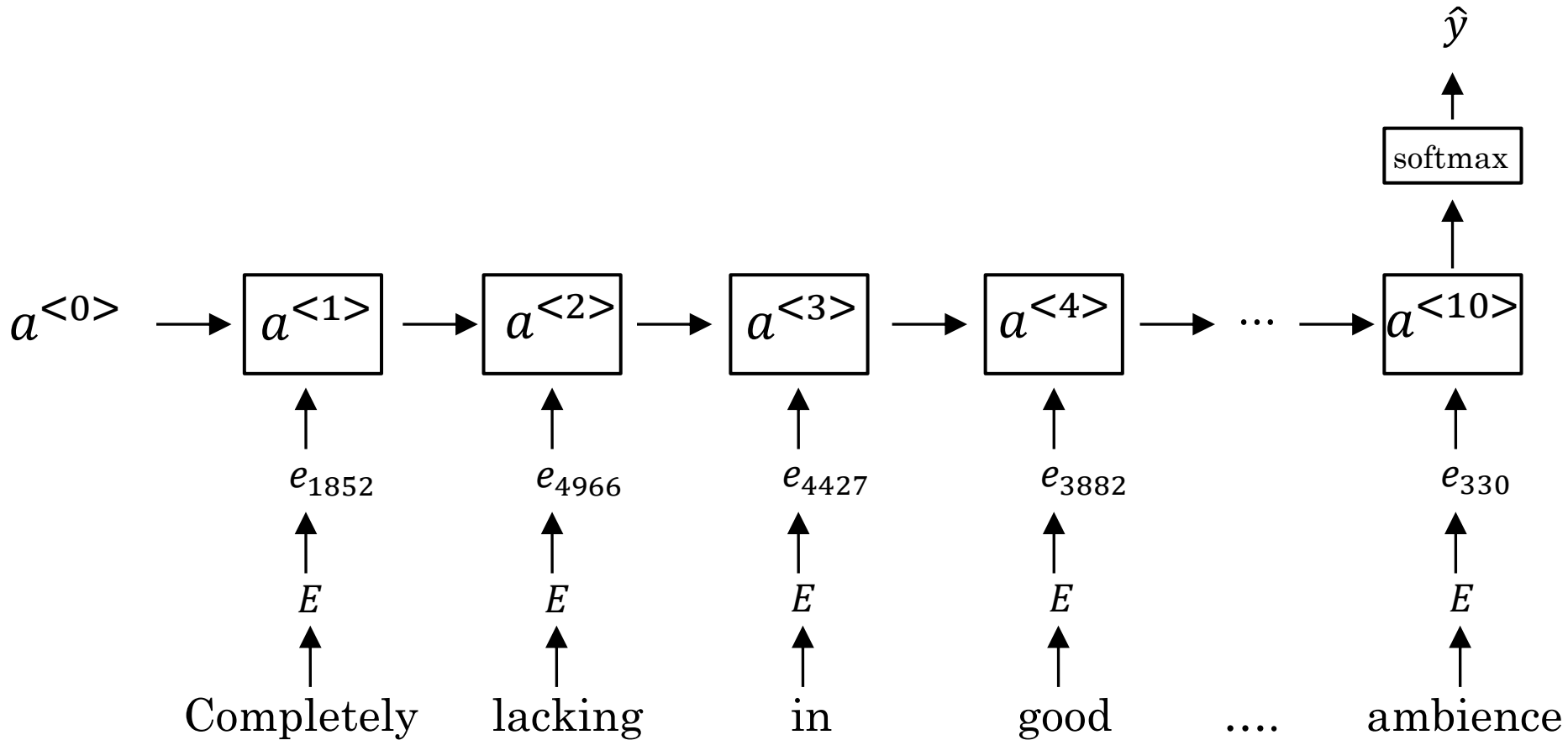


Simple sentiment classification model



“Completely lacking in
good taste, **good** service,
and **good** ambience.”

RNN for sentiment classification



Many-to-one

DEBIASING WORD EMBEDDINGS

The problem of bias in word embeddings

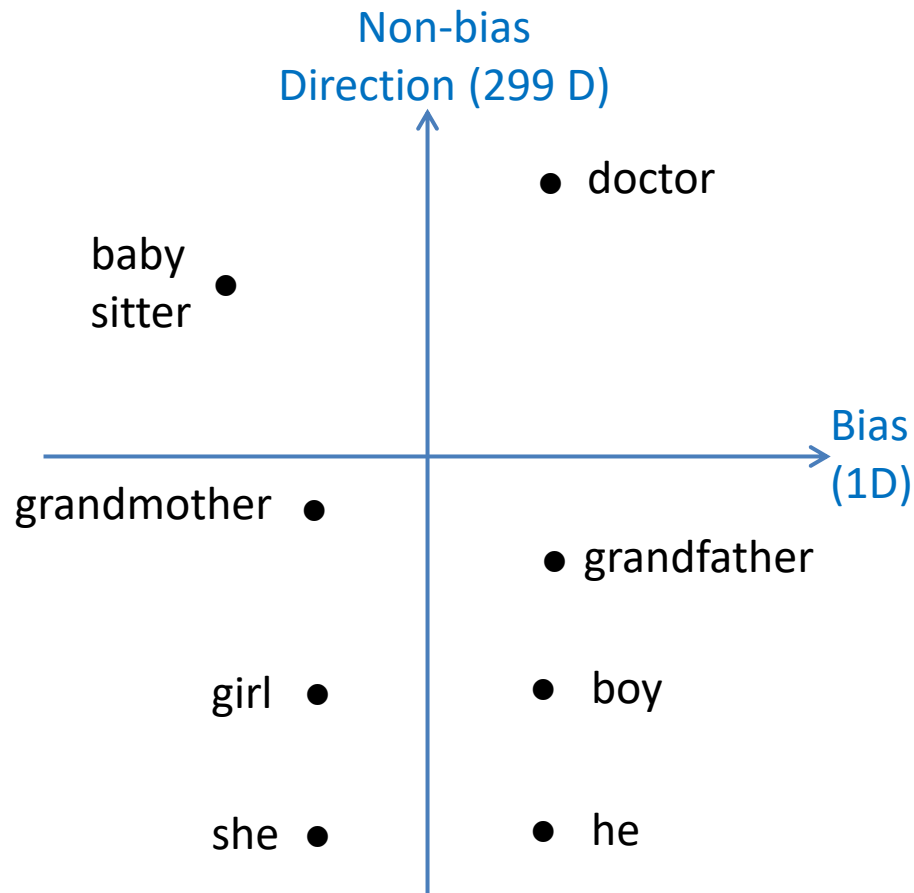
Man:Woman as King:Queen

Man:Computer_Programmer as Woman:Homemaker

Father:Doctor as Mother: Nurse

Word embeddings can reflect gender, ethnicity, age, and other biases of the text used to train the model.

Addressing bias in word embeddings



1. Identify bias direction.

$$\left[\begin{array}{l} e_{he} - e_{she} \\ e_{male} - e_{female} \\ \vdots \end{array} \right]$$

Average

2. Neutralize: For every word that is not definitional, project to get rid of bias.

3. Equalize pairs.

References

- Andrew Ng. Deep learning. Coursera.
- Geoffrey Hinton. Neural Networks for Machine Learning.
- Kevin P. Murphy. Probabilistic Machine Learning An Introduction. MIT Press, 2022.
- MIT Deep Learning 6.S191 (<http://introtodeeplearning.com/>)