



deeplearning.ai

NLP and Word Embeddings

Sentiment classification

Sentiment classification problem



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.



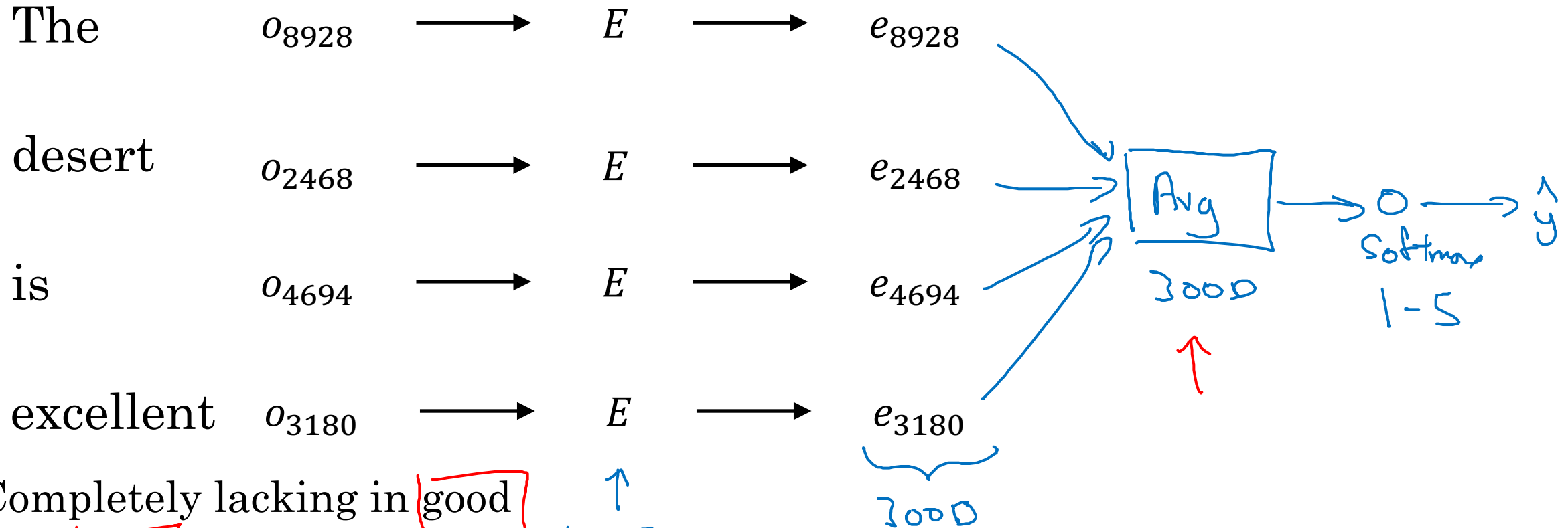
10,000  100,000 words

Simple sentiment classification model

The dessert is excellent



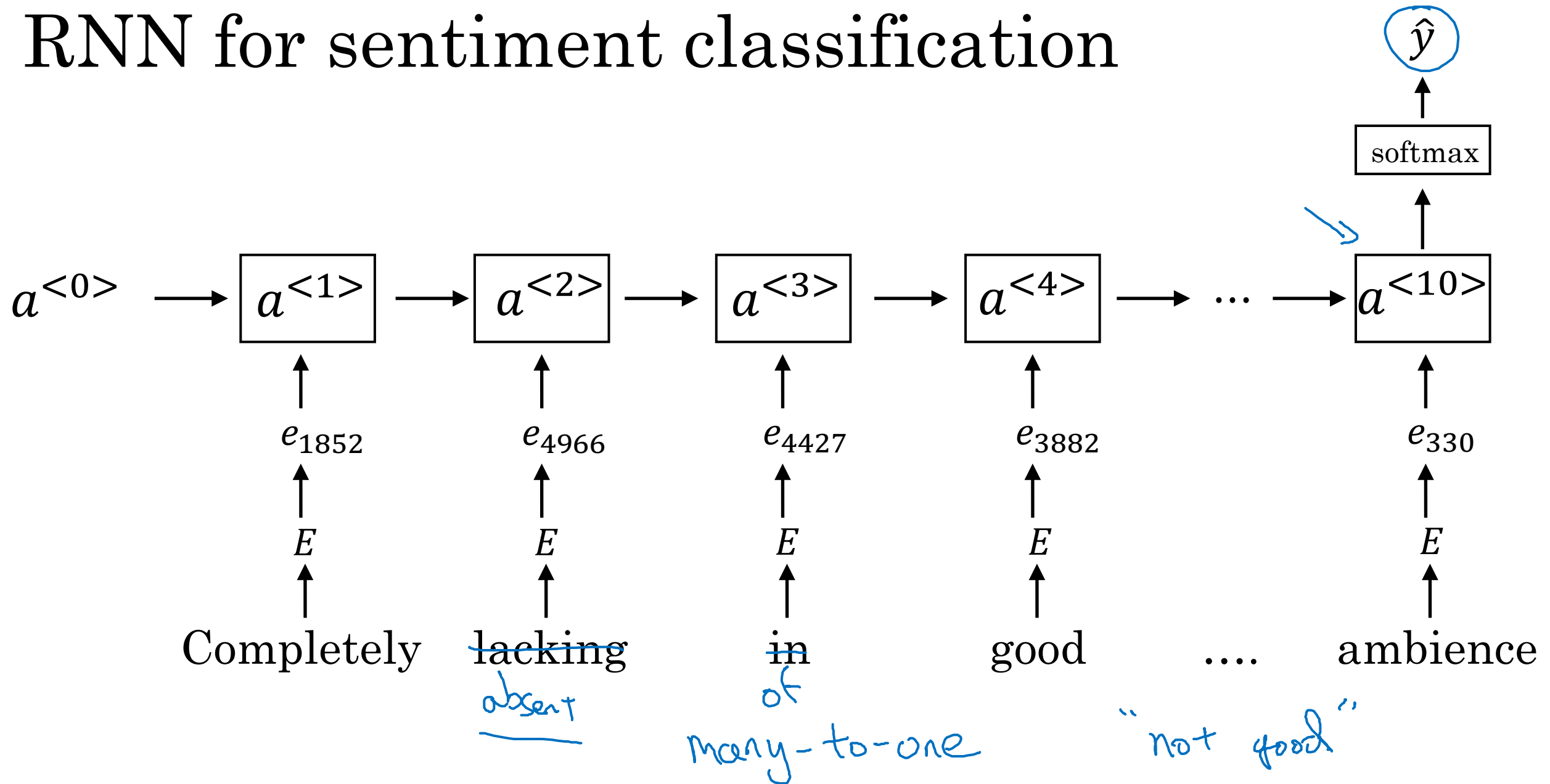
8928 2468 4694 3180



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

↑
1000 words

RNN for sentiment classification





deeplearning.ai

NLP and Word Embeddings

Debiasing word embeddings

The problem of bias in word embeddings

Man:Woman as King:Queen

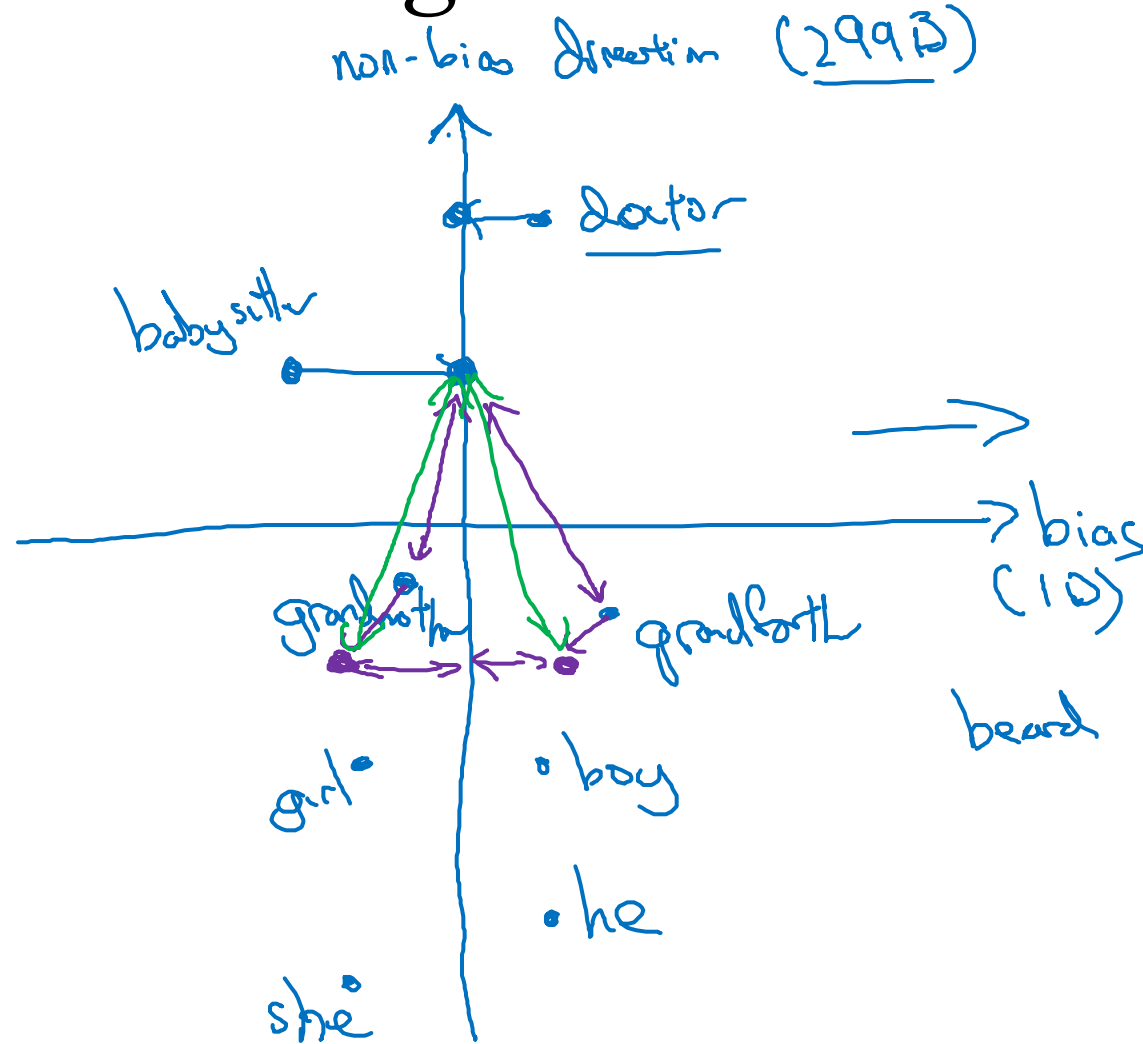
Man:Computer_Programmer as Woman:Homemaker X

Father:Doctor as Mother:Nurse X

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



Addressing bias in word embeddings



1. Identify bias direction.

$\{ \begin{aligned} &e_{he} - e_{she} \\ &e_{male} - e_{female} \\ &\vdots \end{aligned} \}$
→ average

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

3. Equalize pairs.

→ $\left. \begin{array}{cc} \text{grandmother} & \text{grandfather} \\ \text{girl} & \text{boy} \end{array} \right\}$