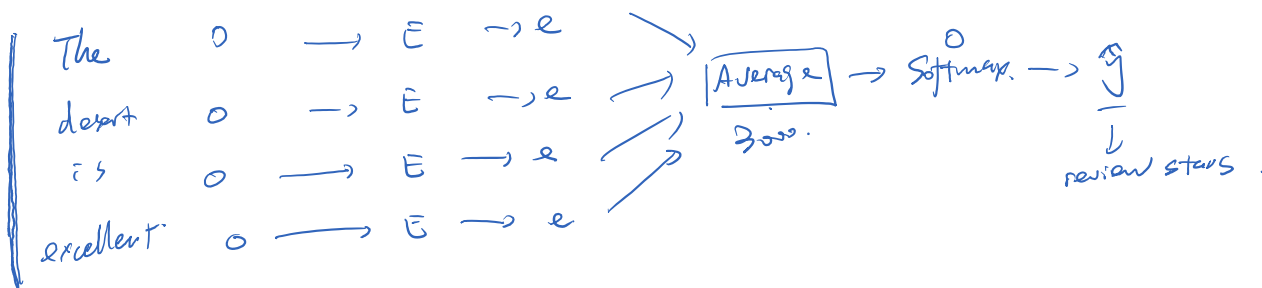


Sentiment classification

Challenge: Don't have a huge labeled dataset.

x \longrightarrow y
 The dissent is excellent $\star\star\star\star$
 slow service $\star\star$

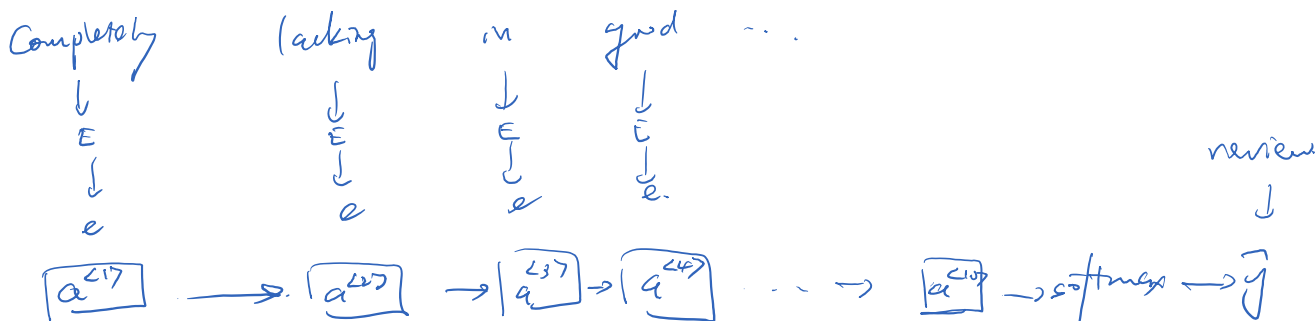
Challenge: you might not have a huge dataset.



problem: ignore order of word.

completely lacking of good food good ~~was~~ ~~not~~ not good

Solution: RNN



Because word embedding can be trained from a larger dataset.
this will still do a good job.

$\frac{1}{2}$ bias \rightarrow $\frac{1}{2}$ bias

Debiasing word embedding

Diminishing bias. [gender, ethnicity] -- bias.

Teacher as doctor | mother as nurse.

Reminding ...

Father as doctor | mother as nurse.
Man programmer | woman homemaker

Word embedding can reflect gender ethnicity ... bias in the text
used to train the model.

procedure

non-directional (good).
doctor
librarian
bias (1D)

① Identify bias direction

$e_{he} - e_{she}$
 $e_{male} - e_{female}$

get direction gender direction

② Neutralization: for each word that is not directional,
project to get rid of bias.

how to decide which
word to neutralize?

train a classifier to tell which word is directional

gender is a
part by def.

e.g. a directional word
grandmother / grandfather

③ Equalize pairs.

grandmother
girl

grandfather
boy.

The paper

Bolukbasi et al. 2016

Man is to computer program ...