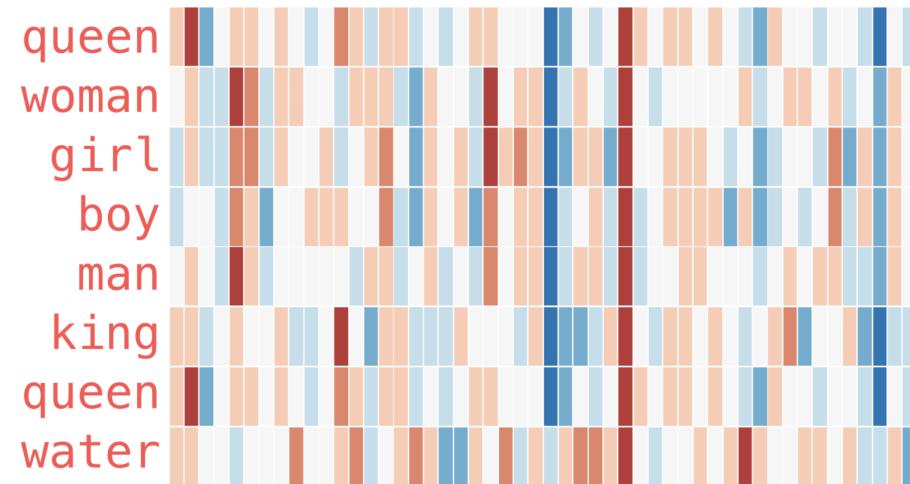


Natural Language Processing (NLP)



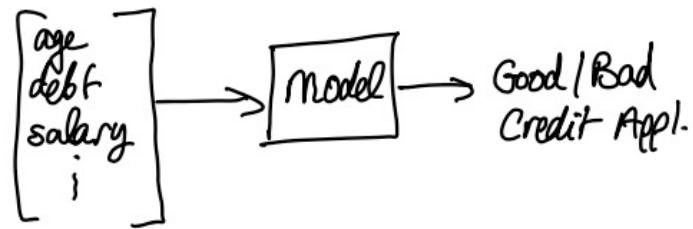
<https://jalammar.github.io/illustrated-word2vec/>

Slide credits:

- <https://www.coursera.org/learn/nlp-sequence-models/home/week/2>
- <https://www.youtube.com/watch?v=yvMgcLKuvVg>

Many good figs from this great blog:

- <https://jalammar.github.io/illustrated-word2vec/>



Numeric input
Fixed length

Numeric features - Fixed Length Input.

What happens when the input is text? \Rightarrow NATURAL LANGUAGE UNDERSTANDING

"The movie was great" \rightarrow Model \rightarrow ★★★★☆
(5-stars)

PROCESSING
(NLP)

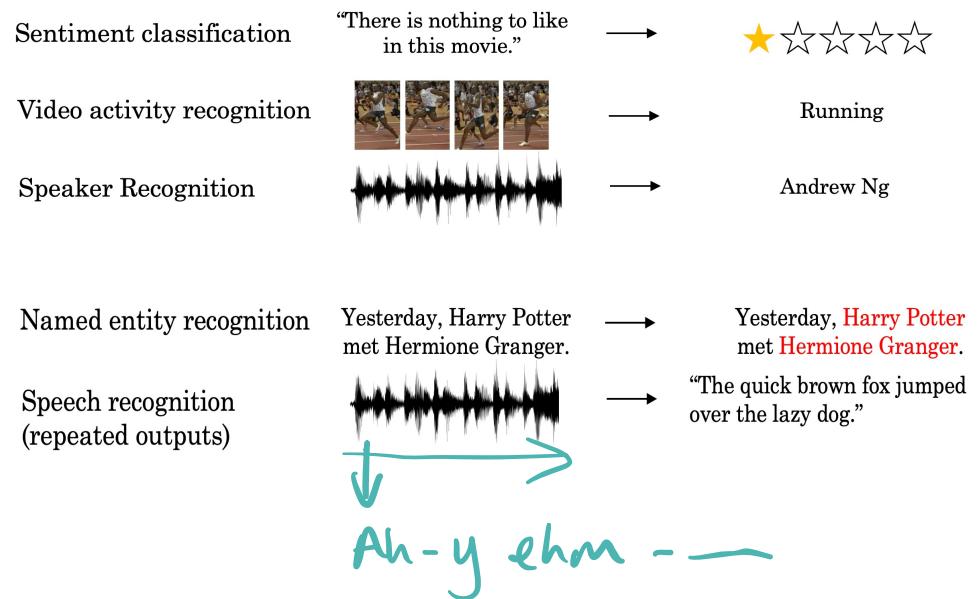
"The movie was great, " \rightarrow Model \rightarrow ★★☆
but the acting was so so

- sentiment analysis
↓ - text summariz.
- Q&A

Text input
Variable length

Various NLP & Seq2Seq Tasks

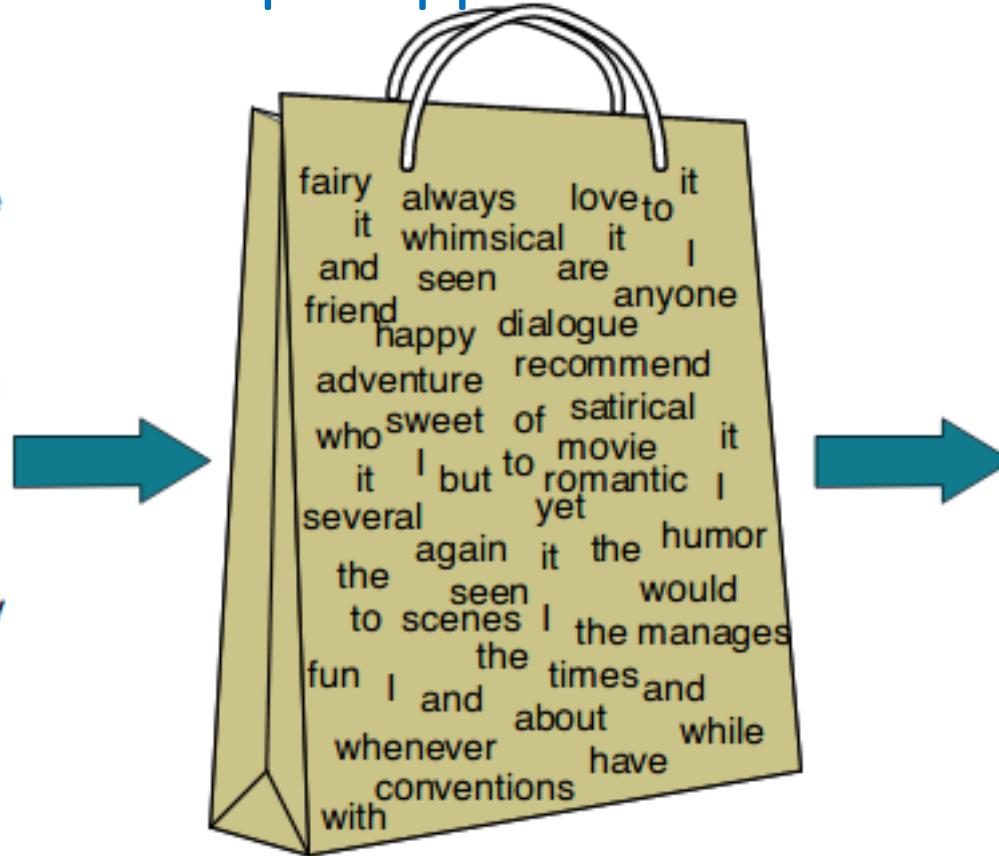
- Sentiment Analysis (**many-to-one**)
- Video activity recognition (**many-to-one**)
- Speaker recognition
- Named Entity Recognition (**many-to-many, aligned**)
- Speech recognition (**many-to-many, aligned**)
- Machine translation (**many-to-many, unaligned**)



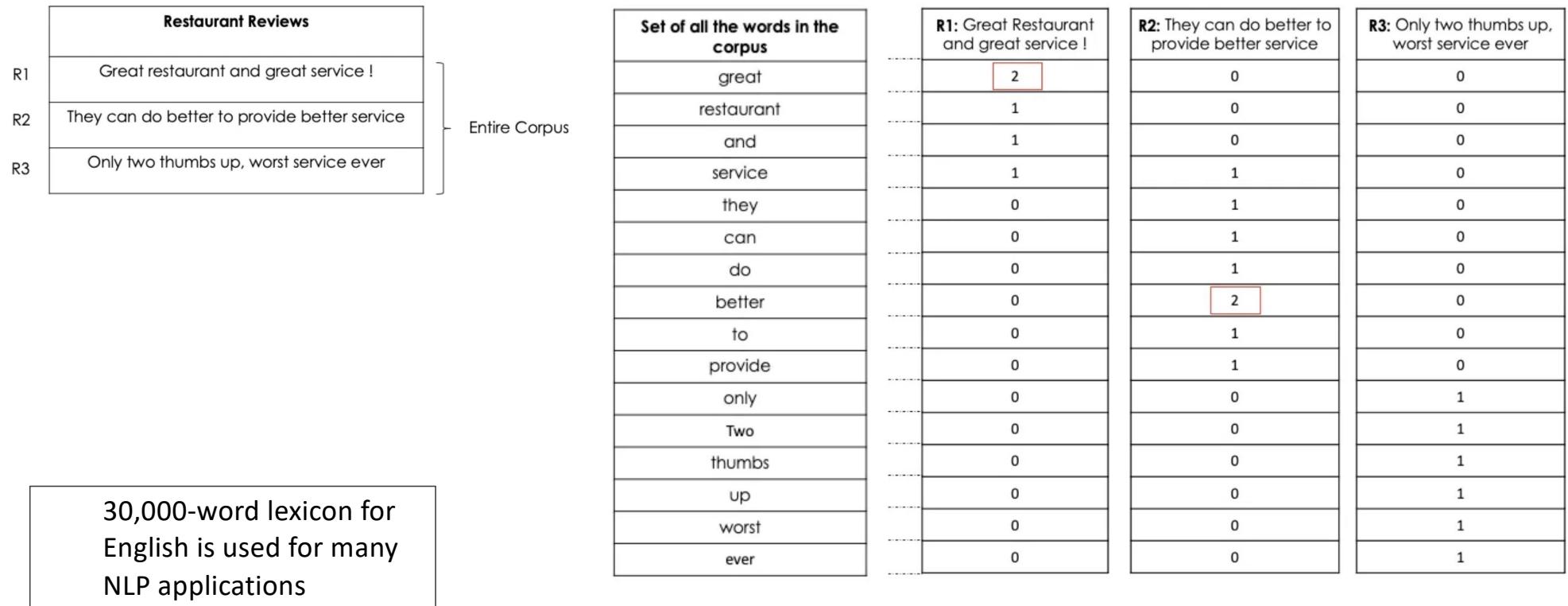
- How do you represent words in NLP applications?
- How do we deal with the variable length sequence?

Sentiment Analysis – A Simple Approach

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



Sentiment Analysis – A Simple Approach



Images from: <https://towardsdatascience.com/text-data-representation-with-one-hot-encoding-tf-idf-count-vectors-co-occurrence-vectors-and-f1bccbd98bef>

Sentiment Analysis with tf-idf vectors

$$TF(t, d) = \frac{\text{number of times } t \text{ appears in } d}{\text{total number of terms in } d}$$

$$IDF(t) = \log \frac{N}{1 + df}$$

$$TF-IDF(t, d) = TF(t, d) * IDF(t)$$

TF term ($\text{TF}(t,d)$) indicates the frequency of term t in document d

IDF indicates the inverse document frequency of the term t in the whole collection D

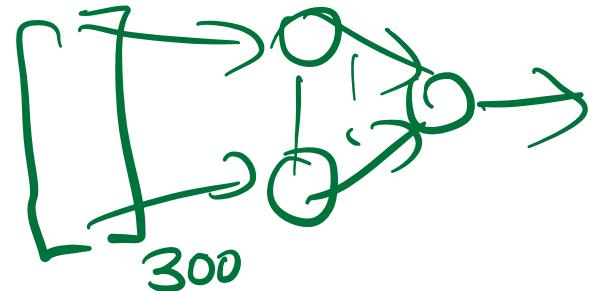
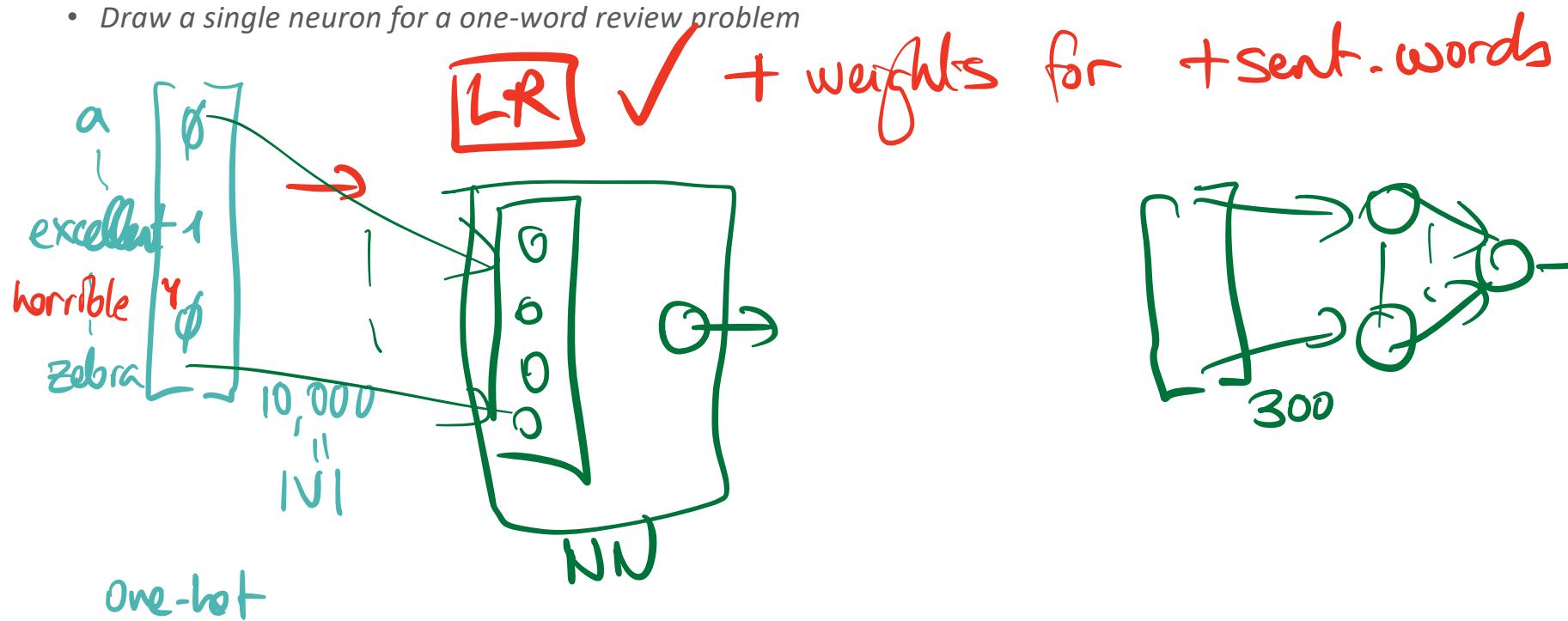
$df(t) = \text{number of documents term } t \text{ appears} / |D|$

R1: Great Restaurant and great service !	R2: They can do better to provide better service	R3: Only two thumbs up, worst service ever
1	0	0
1	0	0
1	0	0
1	1	0
0	1	0
0	1	0
0	1	0
0	1	0
0	1	0
0	1	0
0	1	0
0	0	1
0	0	1
0	0	1
0	0	1
0	0	1

Term frequencies as features

- Train a model with this feature representation

- *Draw a single neuron for a one-word review problem*



One-Hot Encoding

“Hotel was great, a bit small but very comfortable”

a ... great ... hotel motel zone

motel = [0 0 0 0 0 0 0 0 1 0 0 0]

hotel = [0 0 0 0 0 0 1 0 0 0 0 0]

great =

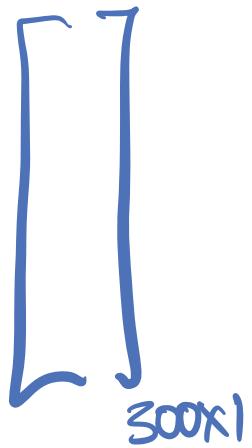
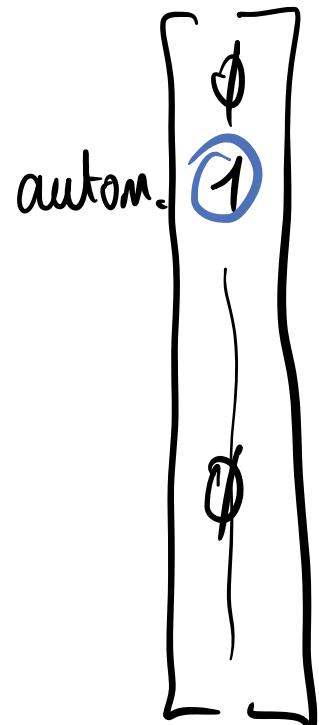
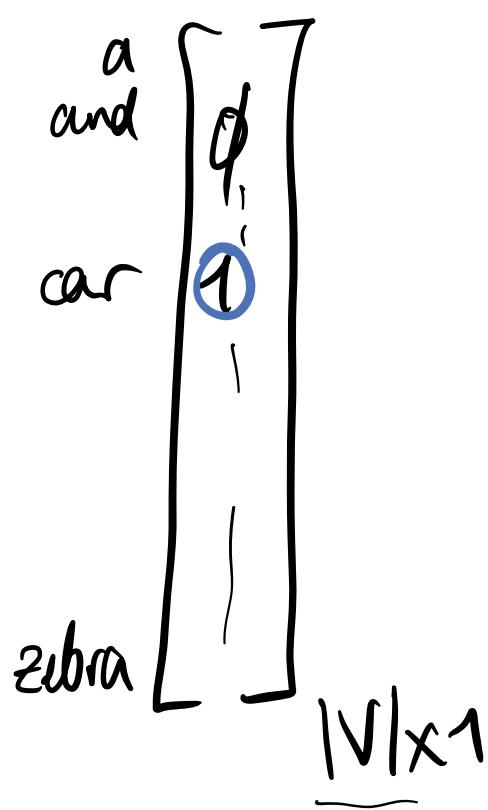
Major Shortcomings:

- Any two vectors are orthogonal
- No notion of similarity between similar words

One-Hot Encoding

$$V = \{a, \dots, zebra\}$$

Distributed



Challenges

Variable-length: A text has variable length which makes it infeasible to use ML models that expect a fixed-length input



- Bag-of-words is a fixed-length vector

Large lexicon: Turkish has agglutinative morphology

gelebiliyor
gelebiliyorum
gelebiliyorsun
...



Word similarities are not captured



- “Great” is close to “Excellent”
- “Motel” is close to “Hotel”
-



Distributed Representation

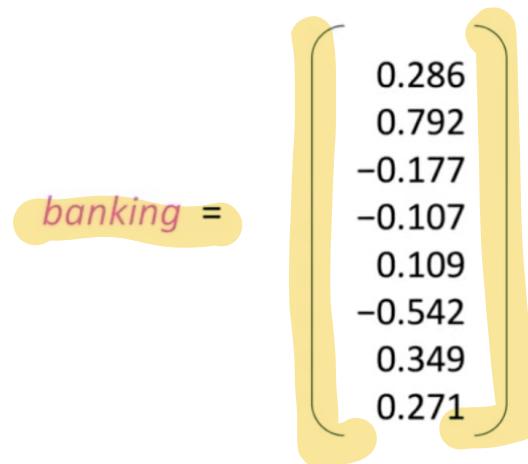
- Can we map the meaning of a **word** or even a **sentence** to a vector?
- *Show example with Excellent, Great, So so, Horrible or King, Queen, Man, Woman...*
- Find **vector representations** that represent the **meaning** in text or images.
- Then we can take these meaning representations (embeddings) and use them as features in any ML model.



Distributed Representation

- Can we map the meaning in a word or even a sentence to a vector?

We will build a dense vector for each word, chosen so that it is similar to vectors of words that appear in similar contexts

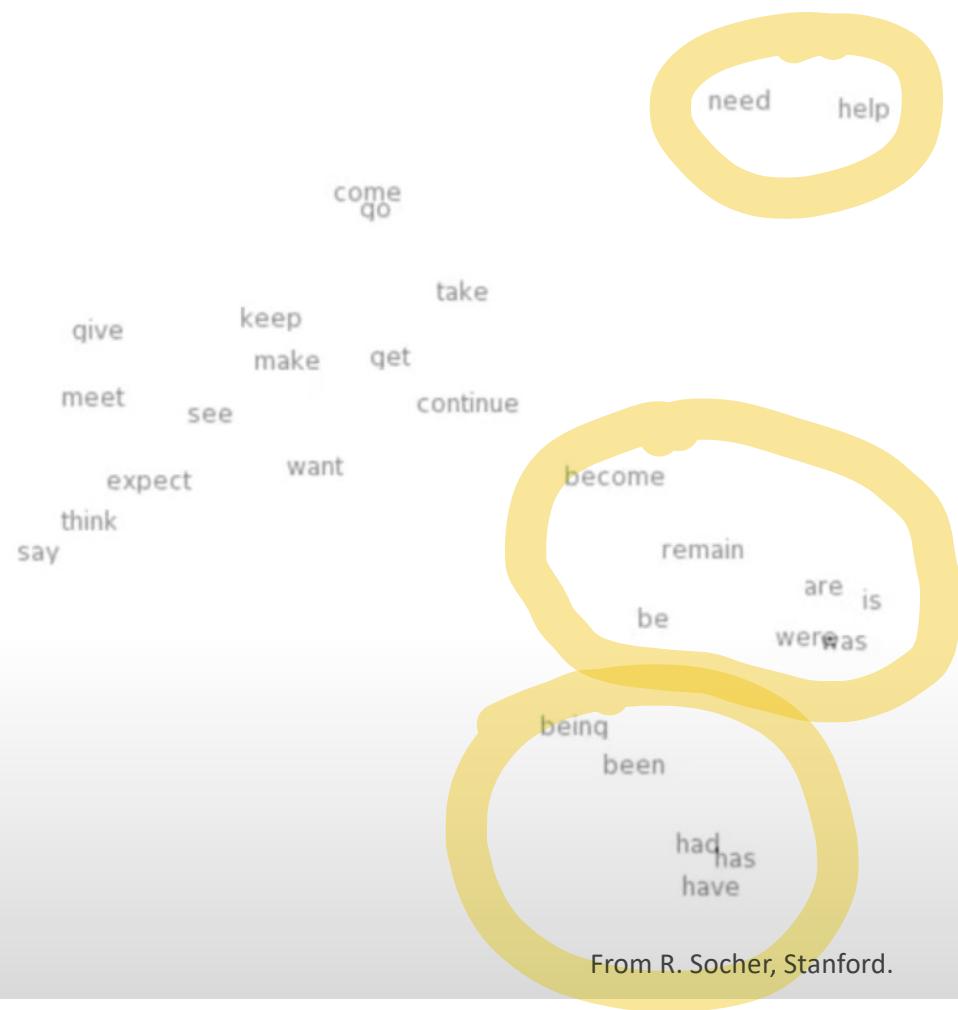


Note: word vectors are also called word embeddings or (neural) word representations
They are a distributed representation

From Manning, Stanford.

Word meaning as a neural word vector – visualization

$$\text{expect} = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \\ 0.487 \end{pmatrix}$$



Distributed Representations

But...

- The dimensions are **not interpretable**
- A given aspect of meaning may be distributed over a combination of many dimensions
- A given dimension may contribute to capturing several aspects of meaning

One-hot vs Distributional Representations

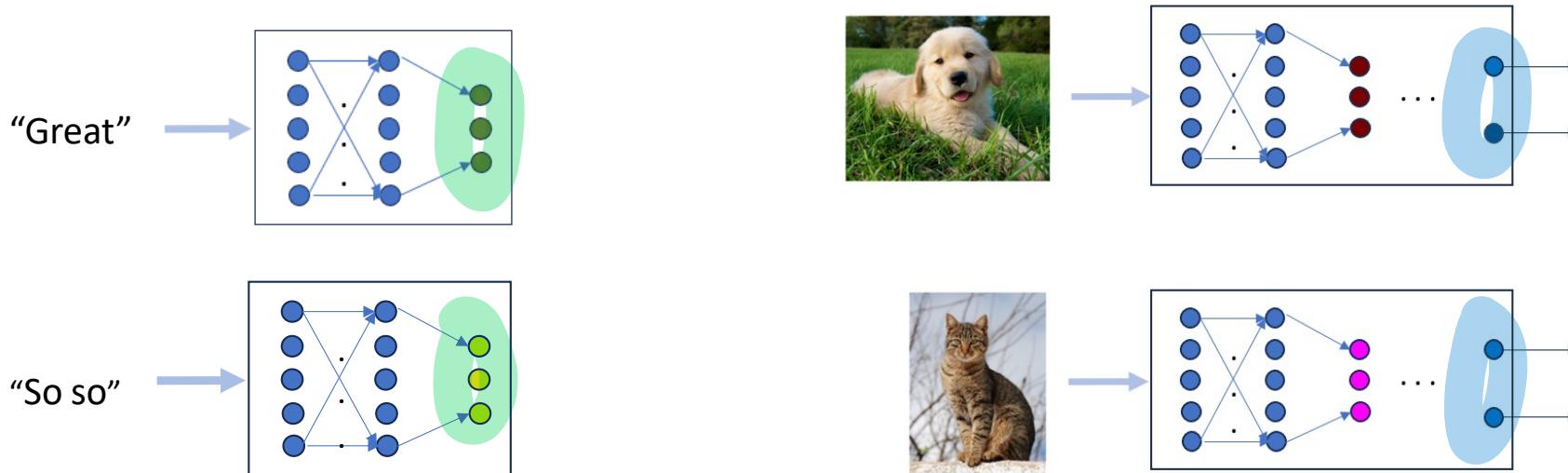
- Local Representation (One-hot Representation)
 - Entities as discrete symbols
 - Interaction between entities as discrete relations
- Distributed Representations [Hinton et al., 1987]
 - Entity as a **vector of values**
 - **Meaning and relation to other entities are captured by** these vectors
 - Words mapped to **shared low dimensional space**

★★★ Embeddings

Find vector representations that represent the meaning in text or images.

Can we map the meaning in a word or even a sentence to a vector?

- We will use deep learning approaches



One-hot vs Distributional Representations

- Local Representation (**One-hot** Representation)
 - Entities as discrete symbols
 - Interaction between entities as discrete relations
 - **Distributed** Representations [Hinton et al., 1987]
 - Entity as a **vector of values**
 - **Meaning and relation to other entities are captured** by these vectors
 - Words mapped to **shared low dimensional space**
 - The dimensions are not interpretable
 - A given aspect of meaning may be distributed over a combination of many dimensions
 - A given dimension may contribute to capturing several aspects of meaning
- Word2Vec ... and now **transformer** based word embeddings

Examples of Learned Relationships

Table 8: *Examples of the word pair relationships, using the best word vectors from Table 4 (Skip-gram model trained on 783M words with 300 dimensionality).*

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

What can you do with word embeddings?

- *Show the toy NN with word embeddings*

How to learn Word Embeddings

Distributional semantics

First observation:

- Distributional semantics: A word's meaning is given by words that appear close by.

“A red sports was speeding through the intersection.”

Corollary (implication):

- Words in similar contexts should have similar embeddings (e.g. car and automobile)

$w_1 \ w_2 \ - \ - \ - \ w_t \ w_{t+1}^?$

$w_1 \ w_2 \ - \ - \ w_j \ w_{j+1} \ - \ - \ - \ w_t$

Language Modelling !

- unigrams
- bigrams
- $P(u|q^{t+1})$
- $P(\underline{to} | "I am" \underline{going})$

Language Model

- A language model is a probability distribution over sequences of words.

$$P(w_1, w_2, \dots, w_T)$$

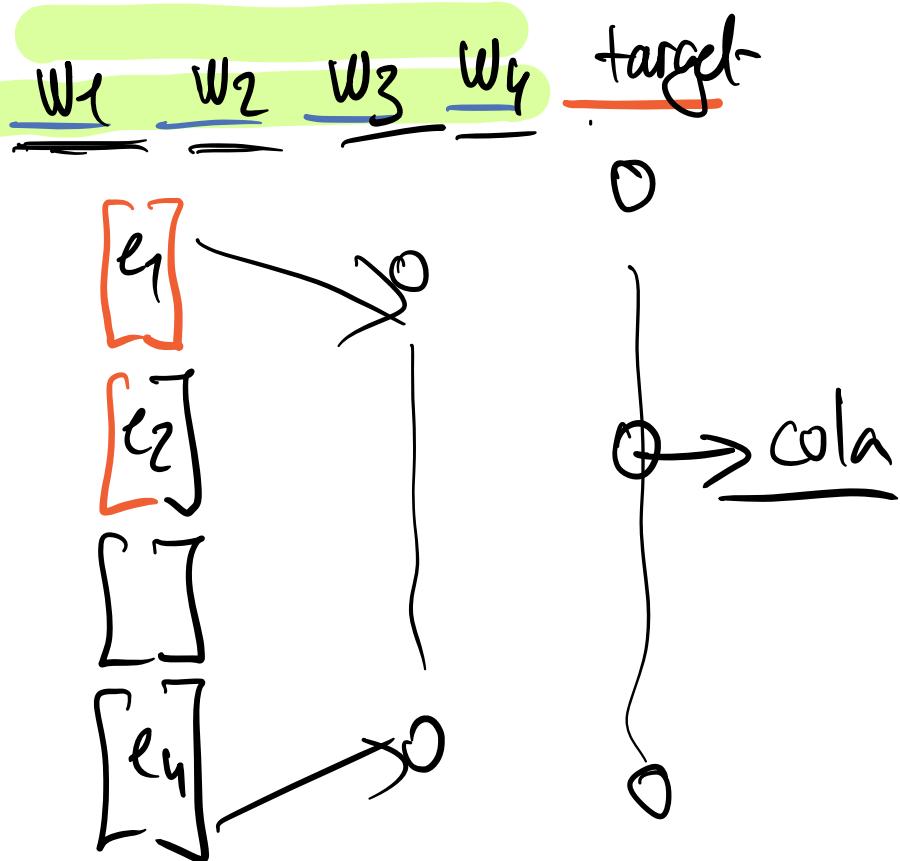
- Given any sequence of words of length T, a language model assigns a probability P to the whole sequence.
- Equivalently, find $P(w_T | w_1, \dots, w_{T-1})$ – **What is the next word?**

Thou shalt not

➤ *kill / steal / commit ...*

Bengio et al. 2003 - A neural probabilistic language model. Journal of Machine Learning Research, 3:1137-1155, 2003.

- Learn word embeddings while learning a **language model**
- First algorithm to learn word embeddings



$$P(w_t | \underline{w_{t-1}, w_{t-2}, \dots})$$

- Unigram
- bigram
- trigram

Mikolov et al. 2013

Efficient Estimation of Word Representations in Vector Space

Tomas Mikolov

Google Inc., Mountain View, CA

tmikolov@google.com

Kai Chen

Google Inc., Mountain View, CA

kaichen@google.com

Greg Corrado

Google Inc., Mountain View, CA

gcorrado@google.com

Jeffrey Dean

Google Inc., Mountain View, CA

jeff@google.com

- Word2Vec -

Large corpus

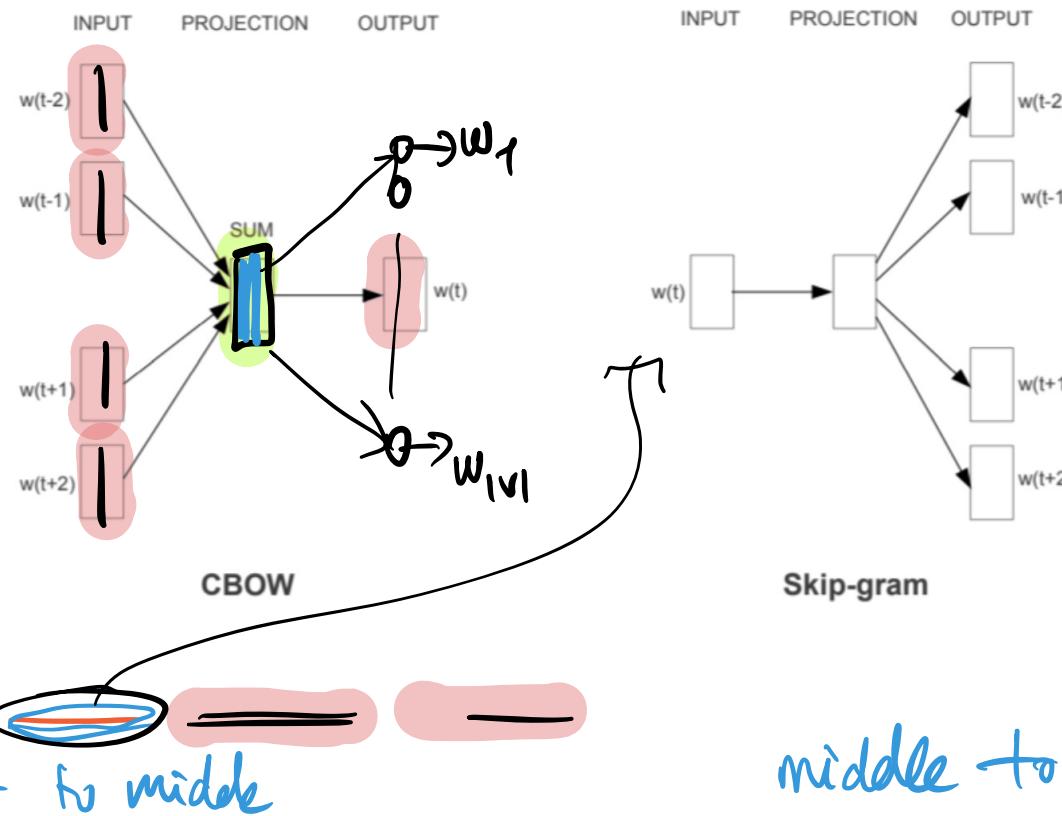
The figure displays three separate Wikipedia article pages arranged vertically. Each page includes the standard Wikipedia header with links for 'Article' and 'Talk', and buttons for 'Read', 'Edit', 'View history', and 'Search Wikipedia'.

1. **Hyperion Cantos**: A series of science fiction novels by Dan Simmons. The title *The Fall of Hyperion* was later used for the overall storyline, including *Edymion*, *The Rise of Endymion*, and a number of short stories.
2. **Dune (novel)**: A 1965 science fiction novel by Frank Herbert, originally published as two serials in *Analog* magazine. It tied with Roger Zelazny's *This Immortal* for the Hugo Award in 1966, and won the Nebula Award for Best Novel.
3. **The Matrix**: A 1999 science fiction action film written and directed by The Wachowskis. It depicts a dystopian future where humans are controlled by machines. The film features Keanu Reeves, Laurence Fishburne, and Carrie-Anne Moss.

The diagram illustrates the concept of a large corpus by showing how multiple Wikipedia articles are interconnected. A red arrow points from the right side of the **Hyperion Cantos** page to the right margin of the **Dune (novel)** page. Another red arrow points from the right side of the **Dune (novel)** page to the right margin of the **The Matrix** page. This visual representation emphasizes the interconnected nature of the information contained within these articles, which is a key characteristic of a large corpus.

CBOW vs Skip-Gram

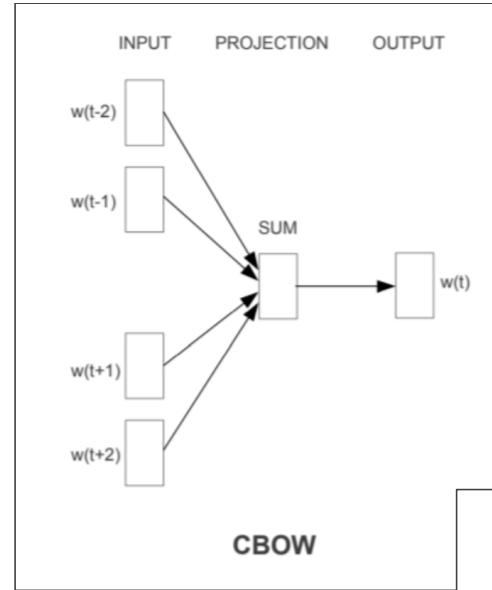
- The distributed representations (word embeddings) of **context words** are combined to **predict the word in the middle (target word)**.
- The distributed representation of the input word is used to **predict the context**.



CBOW vs Skip-Gram

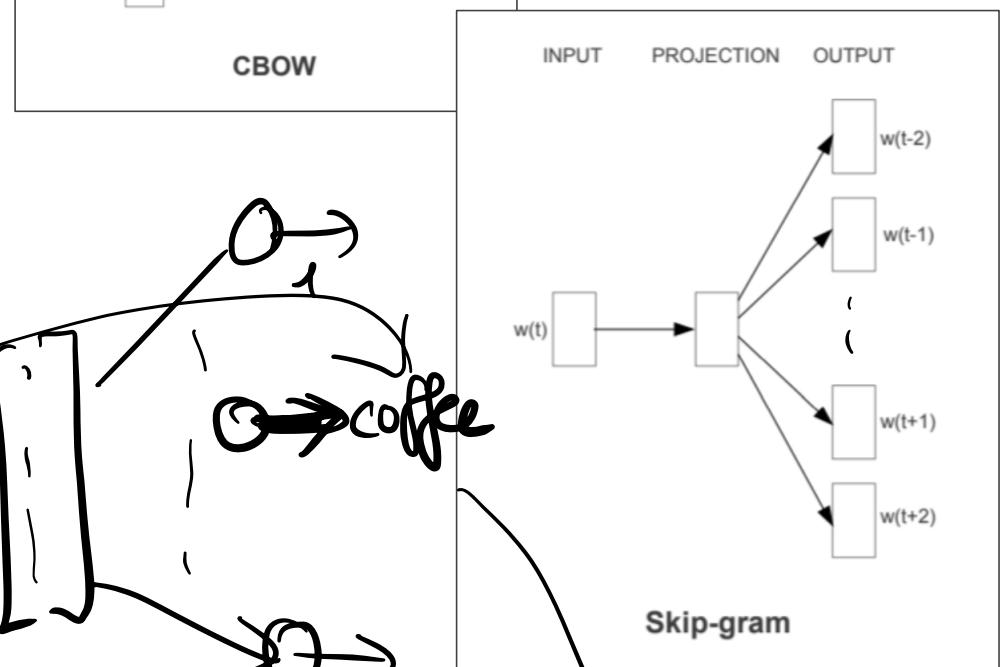
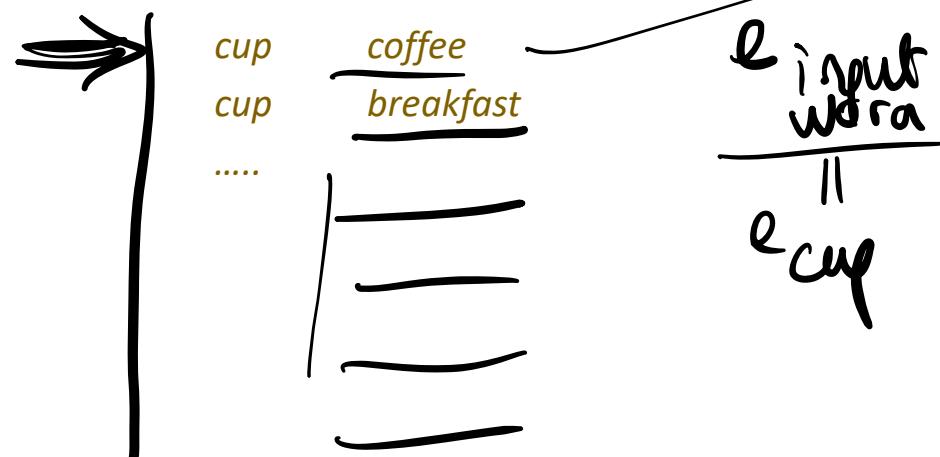
- **CBOW:** The distributed representations of **context words** are given to **predict the target word** (coffee).

"I like to have a cup of coffee at breakfast"



- **Skip-Gram:** The distributed representation of the **input word** (cup) is used to **predict the context words**.

"I like to have a cup of coffee at breakfast"

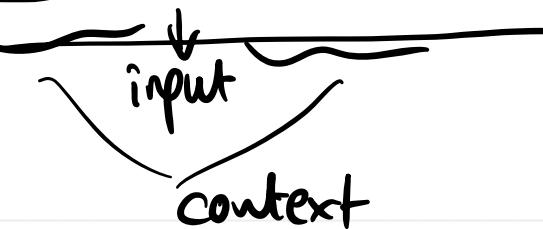


We cross-link
to learn loss :-

Skipgram Training set

Thou shalt not make a machine in the likeness of a human mind

thou	shalt	not	make	a	machine	in	the	...
------	-------	-----	------	---	---------	----	-----	-----



input word	target word
not	thou
not	shalt
not	make
not	a

Thou shalt not make a machine in the likeness of a human mind

thou	shalt	not	make	a	machine	in	the	...
------	-------	-----	------	---	---------	----	-----	-----

thou	shalt	not	make	a	machine	in	the	...
------	-------	-----	------	---	---------	----	-----	-----

input word	target word
not	thou
not	shalt
not	make
not	a

input word	target word
not	thou
not	shalt
not	make
not	a
make	not
make	a
make	make
make	machine
a	not
a	make
a	machine
a	in
machine	make
machine	a
machine	in
machine	the
in	a
in	machine
in	the
in	likeness

Cbow – Continuous Bag of Words Training set

Thou shalt not make a machine in the likeness of a human mind

Sliding window across running text

thou	shalt	not	make	a	machine	in	the	...
------	-------	-----	------	---	---------	----	-----	-----

Dataset

input 1	input 2	output
thou	shalt	not

Thou shalt not make a machine in the likeness of a human mind

Sliding window across running text

thou	shalt	not	make	a	machine	in	the	...
thou	shalt	not	make	a	machine	in	the	

Dataset

input 1	input 2	output
thou	shalt	not
shalt	not	make

Dataset

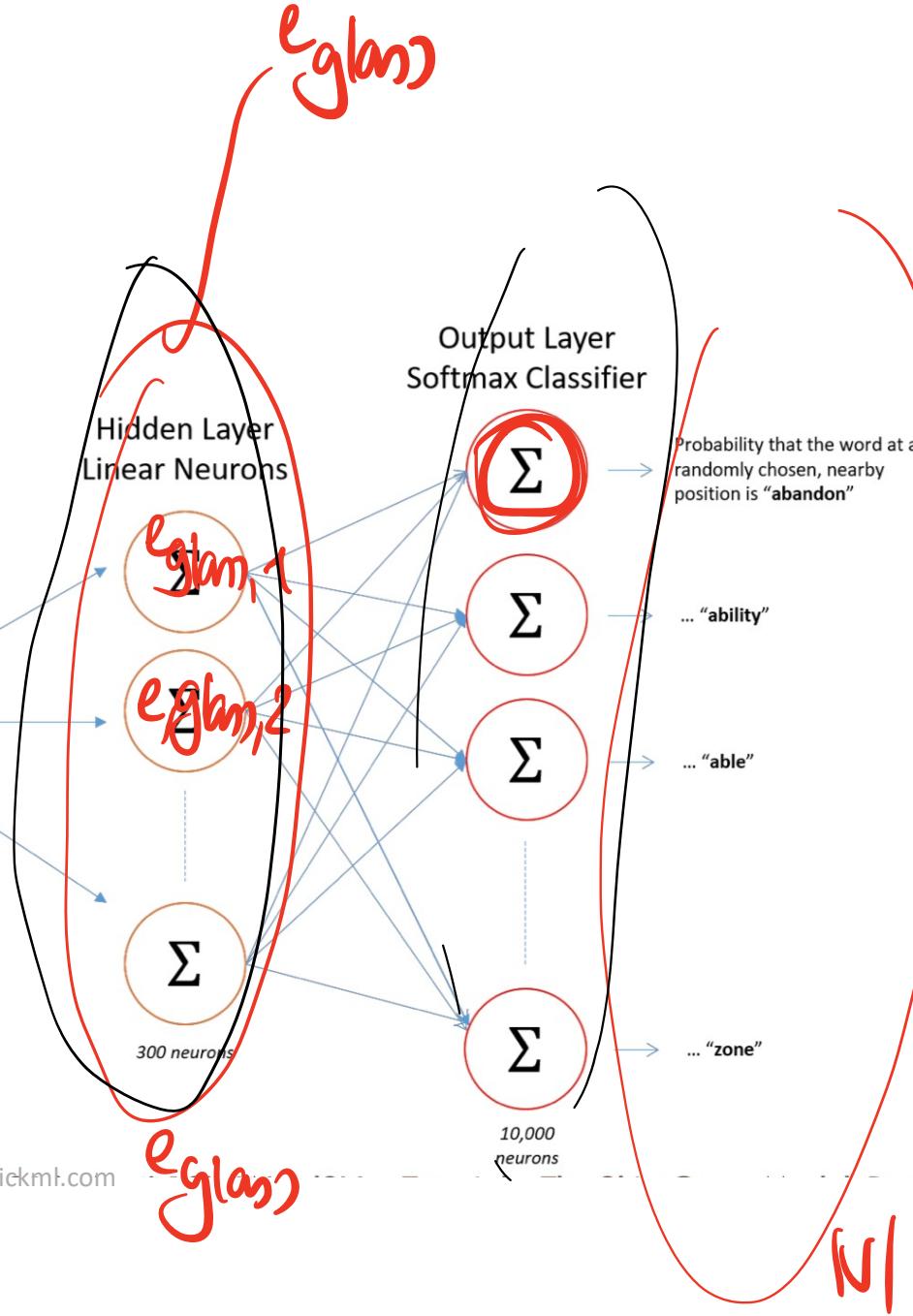
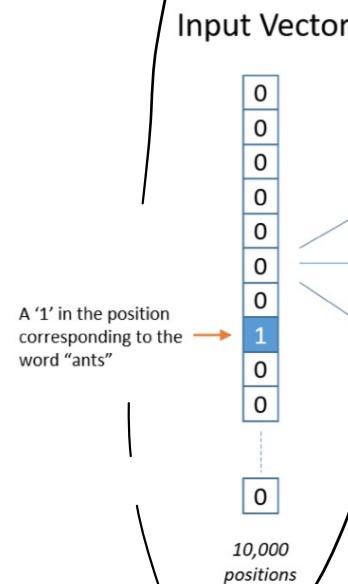
input 1	input 2	output
thou	shalt	not
shalt	not	make
not	make	a
make	a	machine
a	machine	in

SkipGram - Learning Task

- If two different words (*e.g. car and automobile*) have similar “contexts”, then our model needs to **output similar results** (posterior probabilities) **for these two words**.
- One way for the network to output similar predictions for these two words is **if the word vectors are similar**.
- We will predict the posterior probabilities of **context words** (called **outside** words in Manning slides) based on the **similarity (dot product)** of its **word embedding with that of the center word**.

Skip-gram

glass



e_{glam}

- θ_{abandon}

e_{glam}, e_{glass}, θ_{juice}

e_{glam}, θ_{apple}

e_{glam}, θ_{zebra}



↓

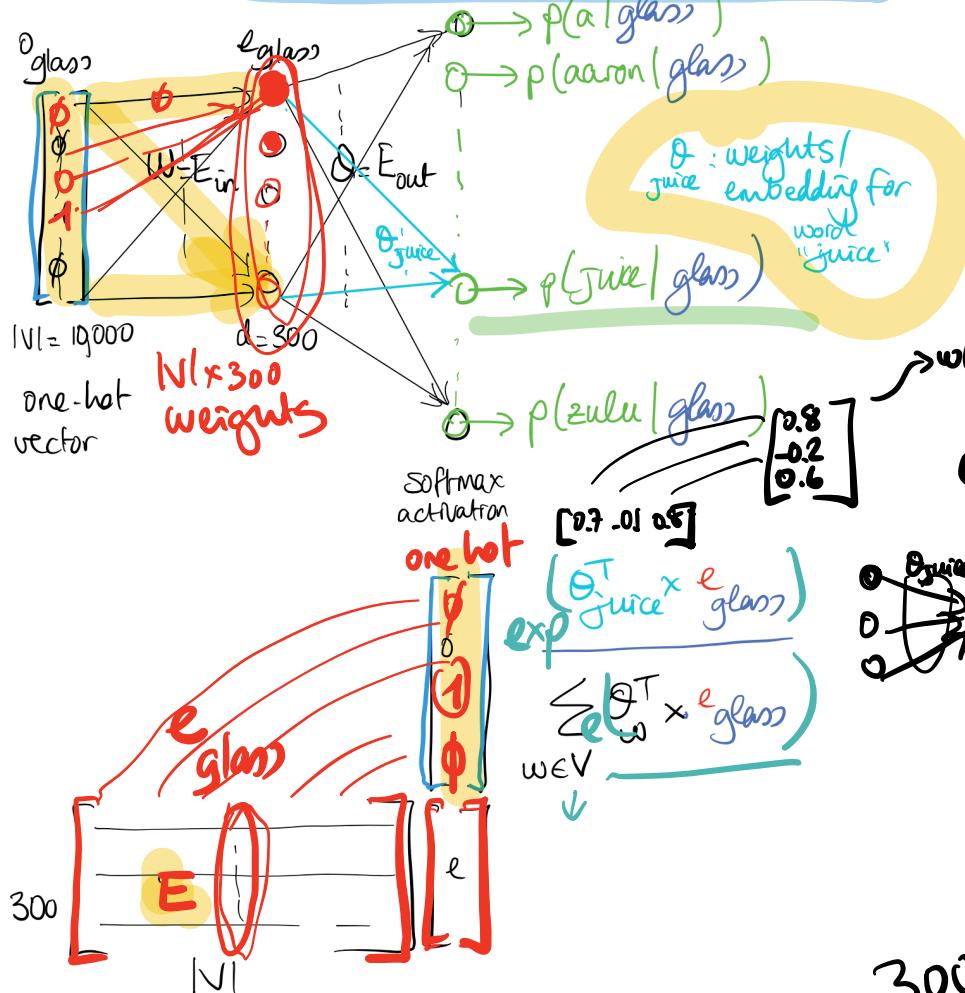
Apply softmax!

Image credit: <http://www.mccormickml.com>

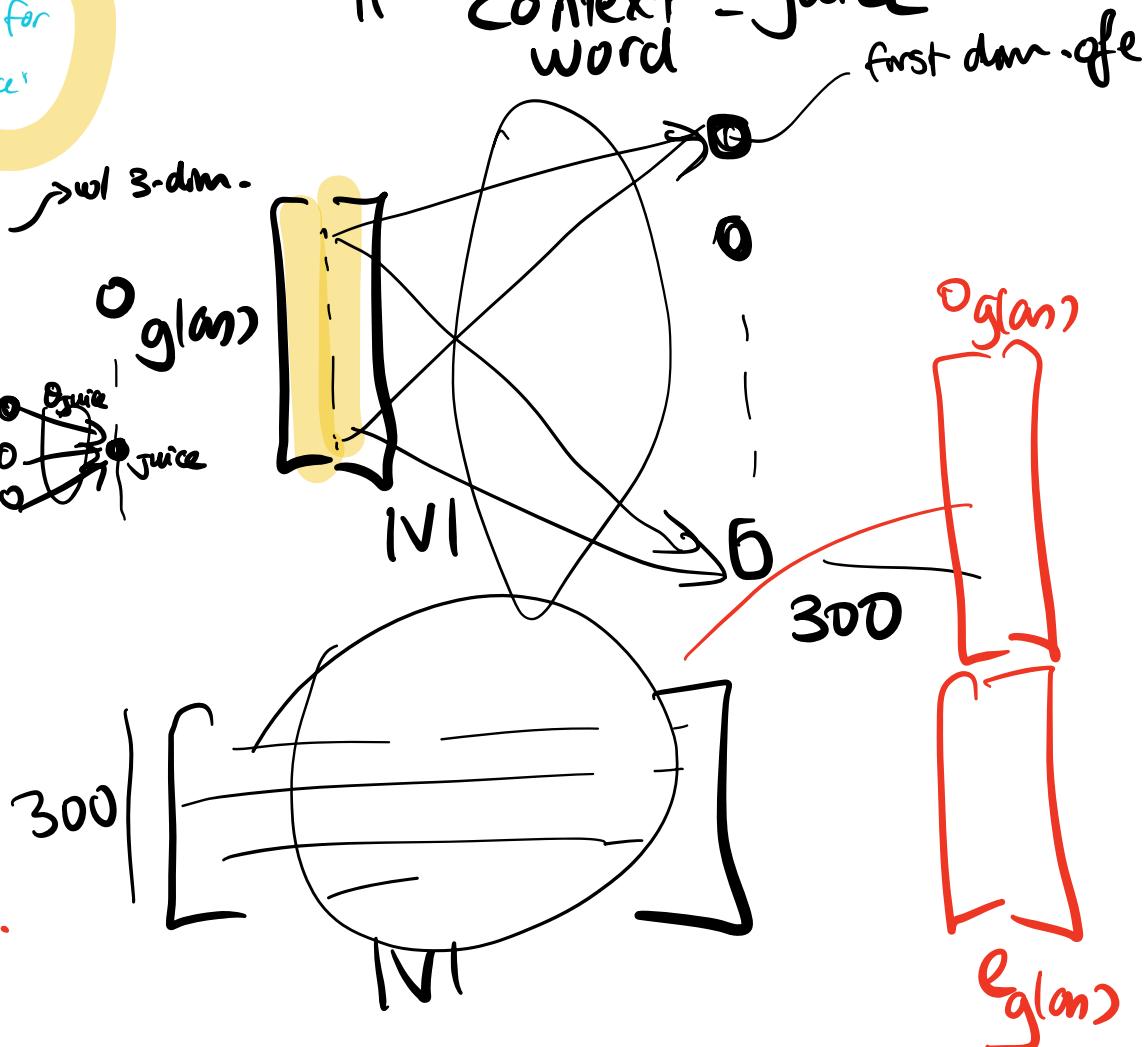
$$e_1 = \phi \cdot w_1 + \phi \cdot w_2 + \dots + \phi \cdot w_{\text{glass}} + \phi \dots$$

Corpus: I want to drink --

I want to drink a glass of orange juice.



$$300 \times N \times N \times 1 \Rightarrow 300 \text{ dim. emb.}$$

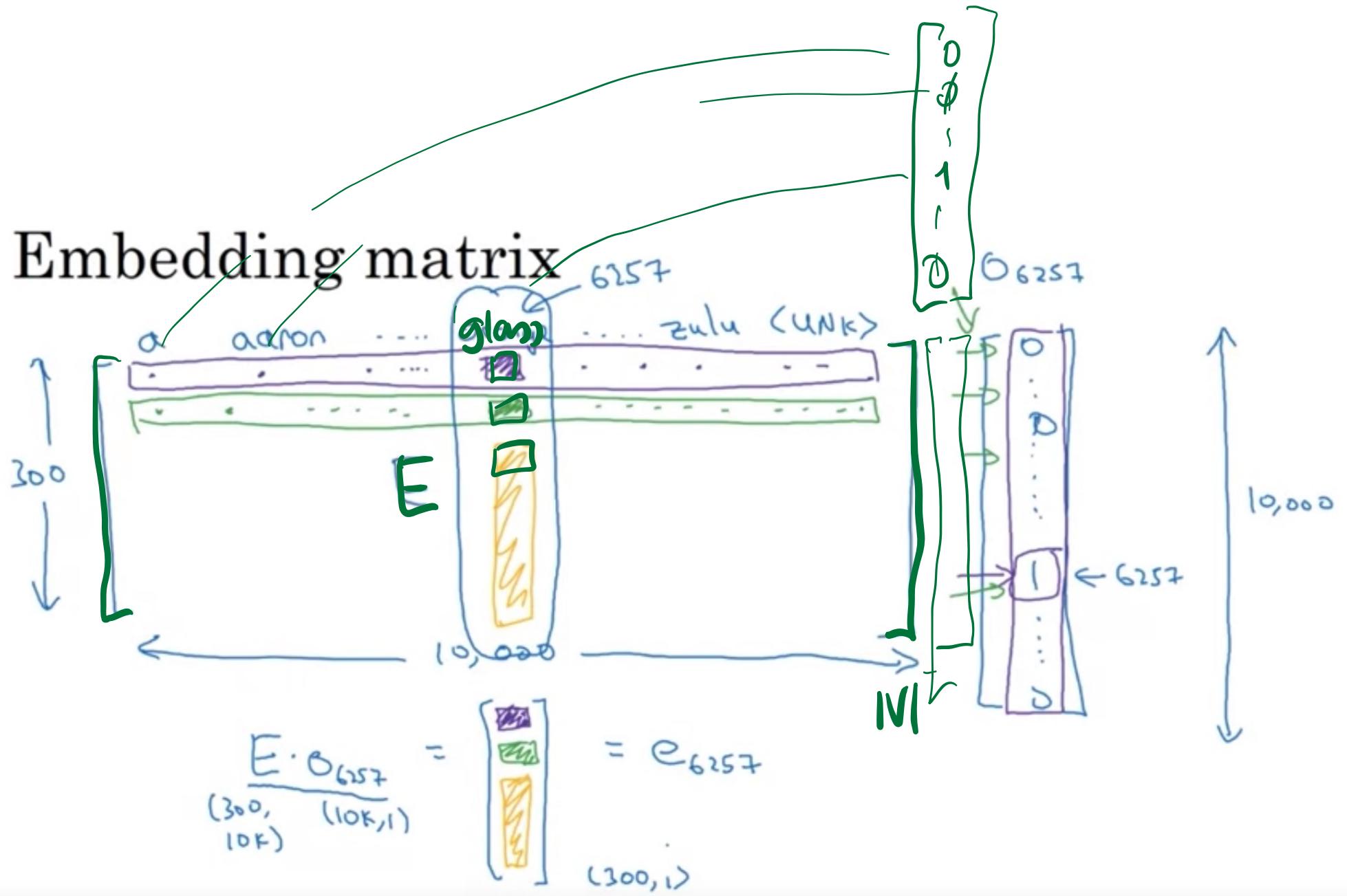


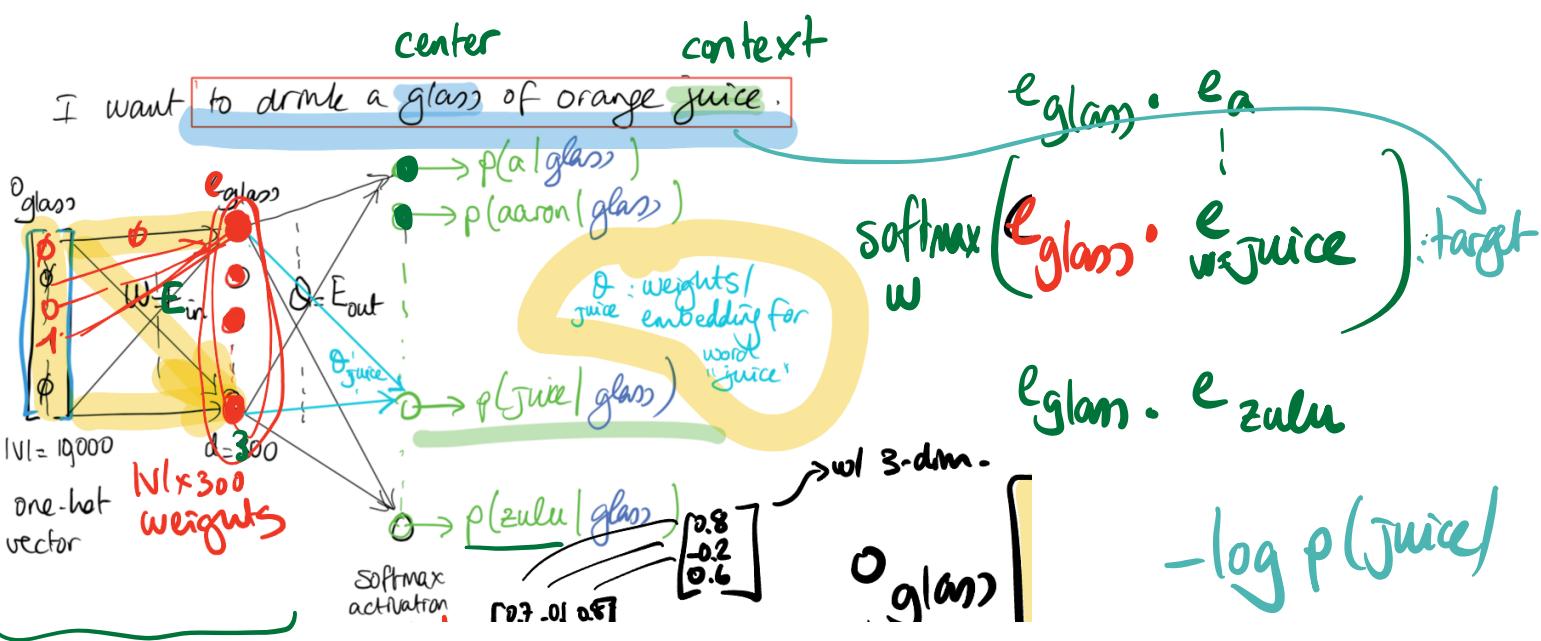
Chosen input = glass

II context = juice
word

first dim. off

Embedding matrix





$${}^0_{\text{glass}} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

e_{glass} extracted by

$$\overline{E} \cdot {}^0_{\text{glass}}$$

weight tensor of the
hidden layer -

hidden
activations

Word2vec: prediction function

- ② Exponentiation makes anything positive

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

① Dot product compares similarity of o and c .
 $u^T v = u \cdot v = \sum_{i=1}^n u_i v_i$
Larger dot product = larger probability

③ Normalize over entire vocabulary
to give probability distribution

- This is an example of the **softmax function** $\mathbb{R}^n \rightarrow (0,1)^n$

$$\text{softmax}(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)} = p_i$$

- The softmax function maps arbitrary values x_i to a probability distribution p_i

- “max” because amplifies probability of largest x_i
- “soft” because still assigns some probability to smaller x_i
- Frequently used in Deep Learning

Open region

But sort of a weird name
because it returns a distribution!

Shortcomings

- Huge number of weights in the embedding matrix and the output layer weights
 - $|V| \times D$ where $|V|$ is the vocabulary size and D is the embedding dimensions.
 - 3M for $|V|=10,000$ $D=3000$
- Training with a large corpus, this is prohibitive.
 - Proposed suggestions in the next paper:
 - Hierarchical softmax
 - Negative sampling
 - Fixed embeddings for each word – that does not depend on the context

Distributed Representations of Words and Phrases and their Compositionality

Tomas Mikolov

Google Inc.

Mountain View

mikolov@google.com

Ilya Sutskever

Google Inc.

Mountain View

ilyasu@google.com

Kai Chen

Google Inc.

Mountain View

kai@google.com

Greg Corrado

Google Inc.

Mountain View

gcorrado@google.com

Jeffrey Dean

Google Inc.

Mountain View

jeff@google.com

Negative Sampling

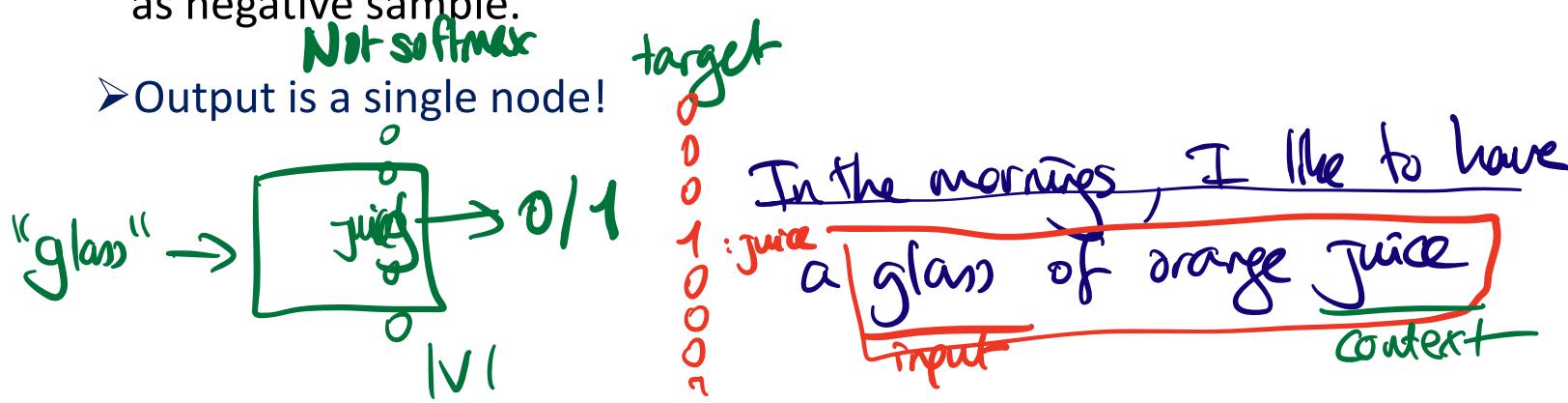
given predict
Task 1: input word \rightarrow context word
"glam" "juice"
out of all $w \in V$

Negative Sampling Idea

Task 2:

Instead of a $|V|$ -way softmax which is costly, change the task to predicting “whether two words are in context of one another”

- So, for the context words around the center word, target will be 1
- Then pick a negative sample randomly from the dictionary (V) and set its target to 0 as it is not in context (at least not now).
- But instead of using all words in the dictionary (except for context) as negative samples, it randomly selects a handful ([2-20]) of words depending on the training size, as negative sample.



dataset-

glass	juice	1
glass	car	0
glass	sleep	0
glass	orange	1