

Intro to Word Embeddings/Language Model

Tf-idf, One hot encodings are sparse representations

tf-idf vectors are

- **long** (length $|V| = 20,000$ to $50,000$)
- **sparse** (most elements are zero)

Sparse versus dense vectors

Why dense vectors?

- Short vectors may be easier to use as **features** in machine learning (less weights to tune)
- Dense vectors may **generalize** better than storing explicit counts
- They may do better at capturing synonymy:
 - *car* and *automobile* are synonyms; but are distinct dimensions
 - a word with *car* as a neighbor and a word with *automobile* as a neighbor should be similar, but aren't
- **In practice, they work better**

Dense embeddings you can download!

Word2vec (Mikolov et al.)

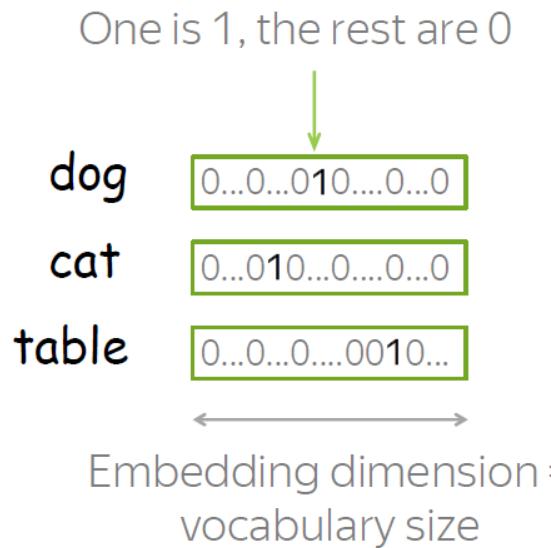
<https://code.google.com/archive/p/word2vec/>

Fasttext <http://www.fasttext.cc/>

Glove (Pennington, Socher, Manning)

<http://nlp.stanford.edu/projects/glove/>

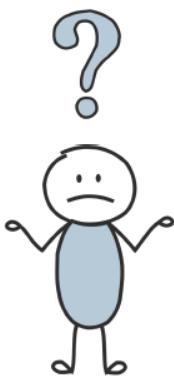
One-Hot Vectors: Represent Words as Discrete Symbols



Problems:

- Vector size is too large
- Vectors know nothing about **meaning**

e.g., **cat** is as close to **dog** as it is to **table!**



What is meaning?

What is meaning?

Do you know what the word **tezgüino** means ?

(We hope you do not)



What is meaning?

Now look how this word is used in different contexts:

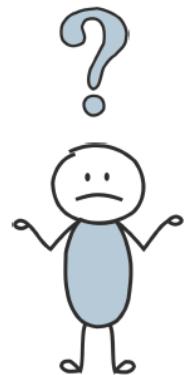
A bottle of **tezgüino** is on the table.

Everyone likes **tezgüino**.

Tezgüino makes you drunk.

We make **tezgüino** out of corn.

Can you understand what **tezgüino** means ?



What is meaning?

Now look how this word is used in different contexts:

A bottle of **tezgüino** is on the table.

Everyone likes **tezgüino**.

Tezgüino makes you drunk.

We make **tezgüino** out of corn.



Tezgüino is a kind of alcoholic beverage made from corn.



With context, you can understand the meaning!

What is meaning?

- (1) A bottle of _____ is on the table.
- (2) Everyone likes _____ .
- (3) _____ makes you drunk.
- (4) We make _____ out of corn.

What other words fit into these contexts ?

	(1)	(2)	(3)	(4)	...	← contexts
tezgüino	1	1	1	1		
loud	0	0	0	0		← rows show contextual
motor oil	1	0	0	1		properties: 1 if a word can
tortillas	0	1	0	1		appear in the context, 0 if not
wine	1	1	1	0		



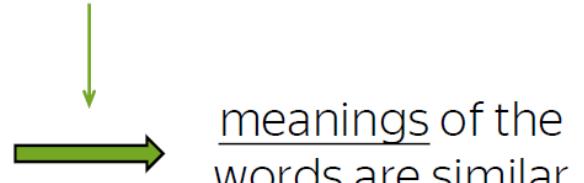
What is meaning?

- (1) A bottle of _____ is on the table.
- (2) Everyone likes _____ .
- (3) _____ makes you drunk.
- (4) We make _____ out of corn.

	(1)	(2)	(3)	(4)	...
tezgüino	1	1	1	1	
loud	0	0	0	0	
motor oil	1	0	0	1	
tortillas	0	1	0	1	
wine	1	1	1	0	

This is the distributional hypothesis

rows are
similar



Distributional Hypothesis

Words which frequently appear in **similar contexts** have **similar meaning**.

(Harris 1954, Firth 1957)

This can be used in practice to build word vectors!

Word2vec

Popular embedding method

Very fast to train

Code available on the web

Idea: **predict** rather than **count**

Word2vec

- Instead of **counting** how often each word w occurs near "*apricot*"
- Train a classifier on a **binary prediction** task:
 - Is w likely to show up near "*apricot*"?
- We don't actually care about this task
 - But we'll take the learned classifier weights as the word embeddings

Brilliant insight: Use running text as implicitly supervised training data!

- A word s near *apricot*
 - Acts as gold ‘correct answer’ to the question
 - “Is word w likely to show up near *apricot*? ”
- No need for hand-labeled supervision

Word2Vec: Skip-Gram Task

Word2vec provides a variety of options. Let's do

- "skip-gram with negative sampling" (SGNS)

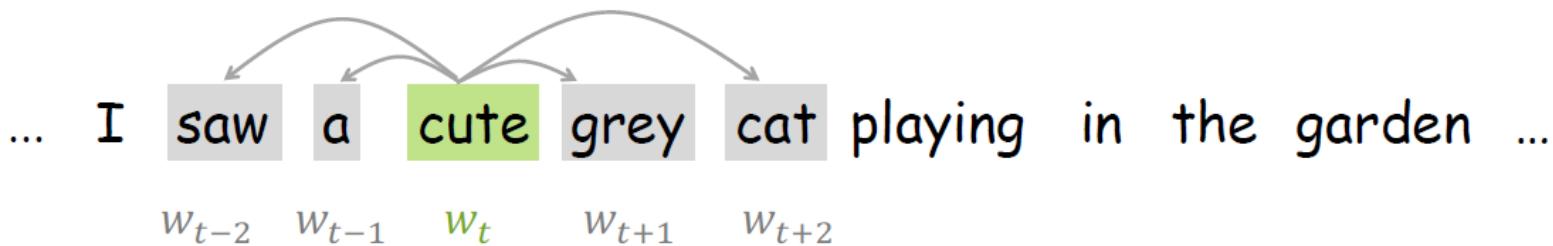
Skip-gram algorithm

1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the weights as the embeddings

Skip-Gram Training Data



context words central word context words



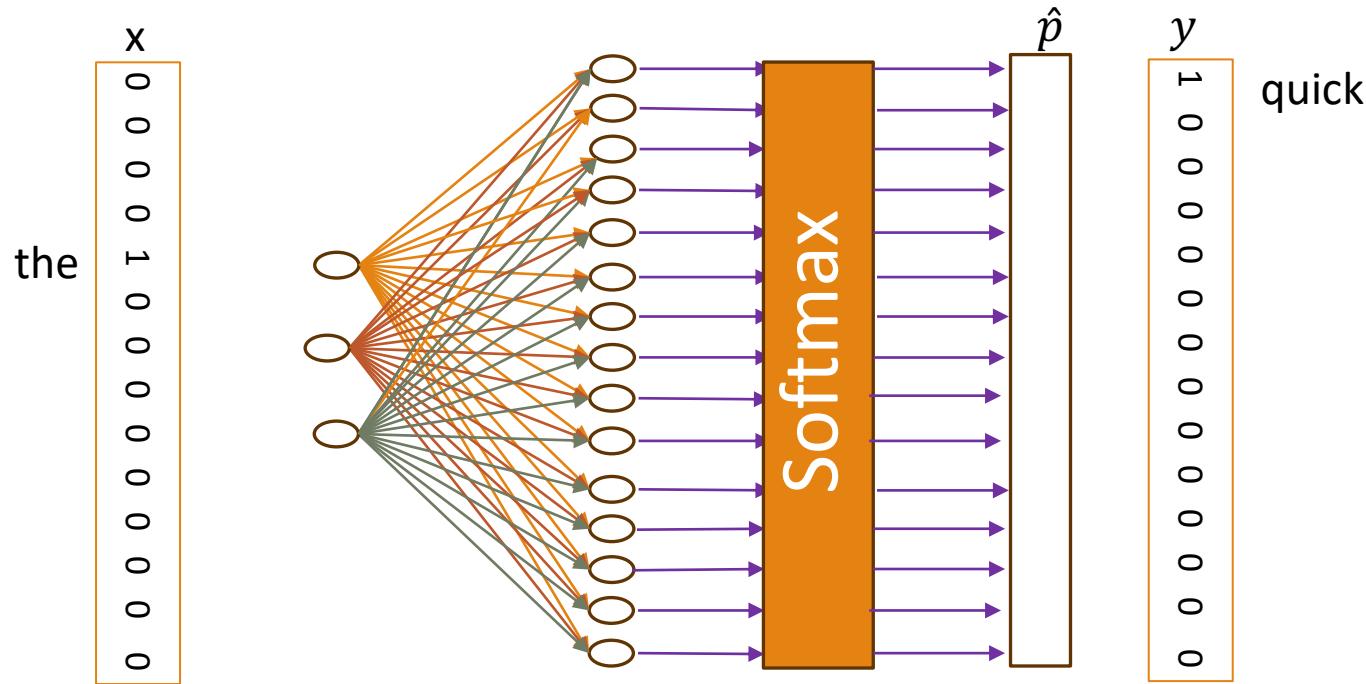
context words central word context words

Given a central word (X), predict context words (y).

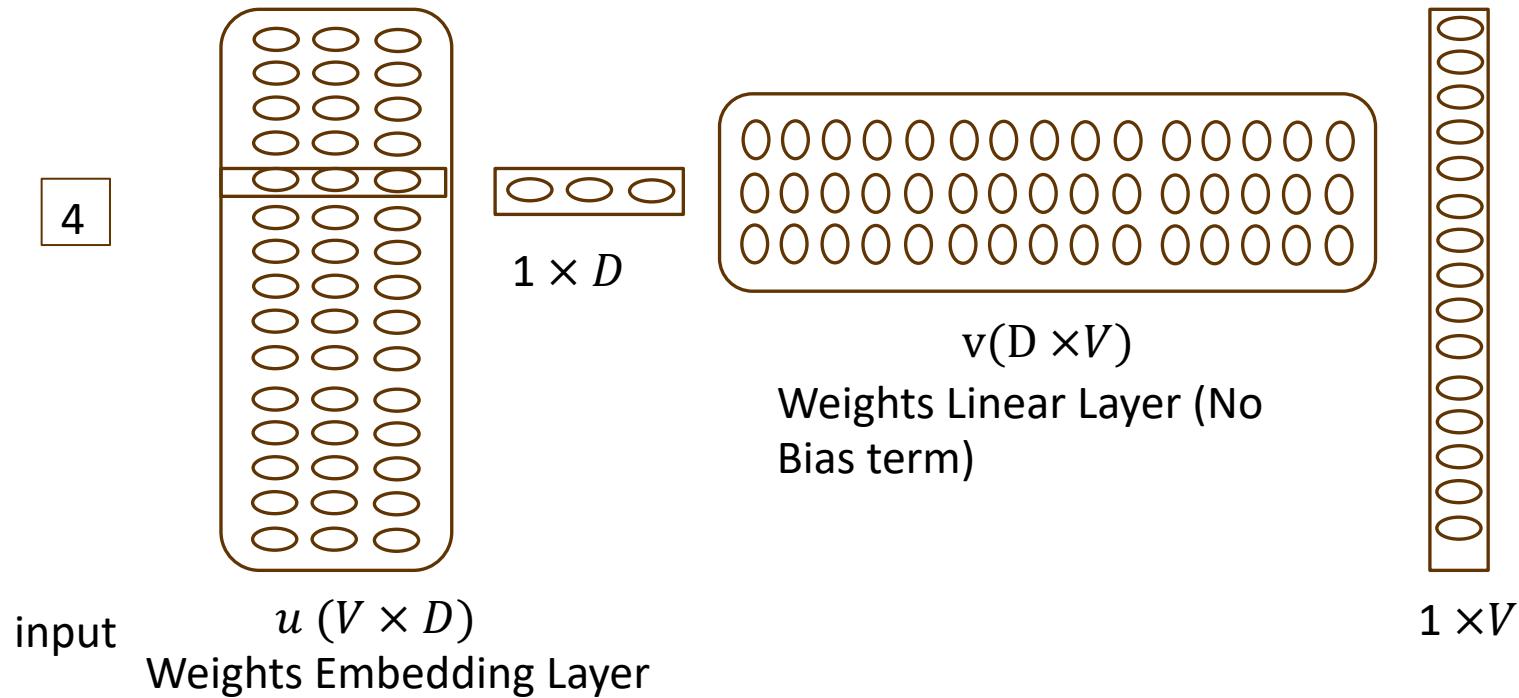
Skip-Gram Training Data

Source Text

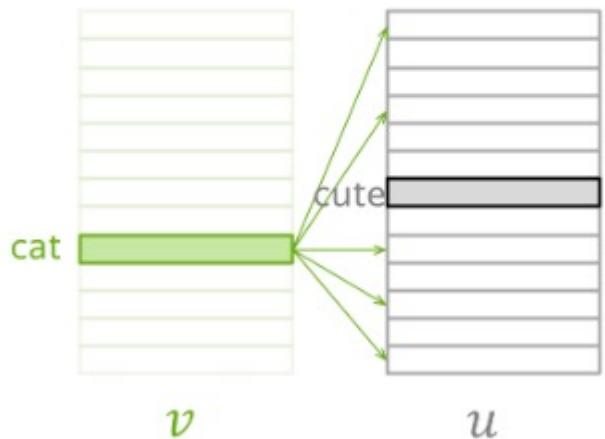
(X, y)	Training Samples
The quick brown fox jumps over the lazy dog. →	(the, quick) (the, brown)
The quick brown fox jumps over the lazy dog. →	(quick, the) (quick, brown) (quick, fox)
The quick brown fox jumps over the lazy dog. →	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox jumps over the lazy dog. →	(fox, quick) (fox, brown) (fox, jumps) (fox, over)



4



1. Take dot product of v_{cat} with all u



2. exp

The diagram shows the second step. It lists several dot products: $u_{w1}^T v_{cat}$, $u_{w3}^T v_{cat}$, followed by a vertical ellipsis, then $u_{cute}^T v_{cat}$ (with a box around it), another vertical ellipsis, and finally $u_{wn}^T v_{cat}$. Each dot product is followed by an arrow pointing to its exponential value: $\exp(u_{w1}^T v_{cat})$, $\exp(u_{w3}^T v_{cat})$, and so on. These individual exponentials are then summed up, as indicated by a large green bracket and the label "2" in a circle.

3. sum all

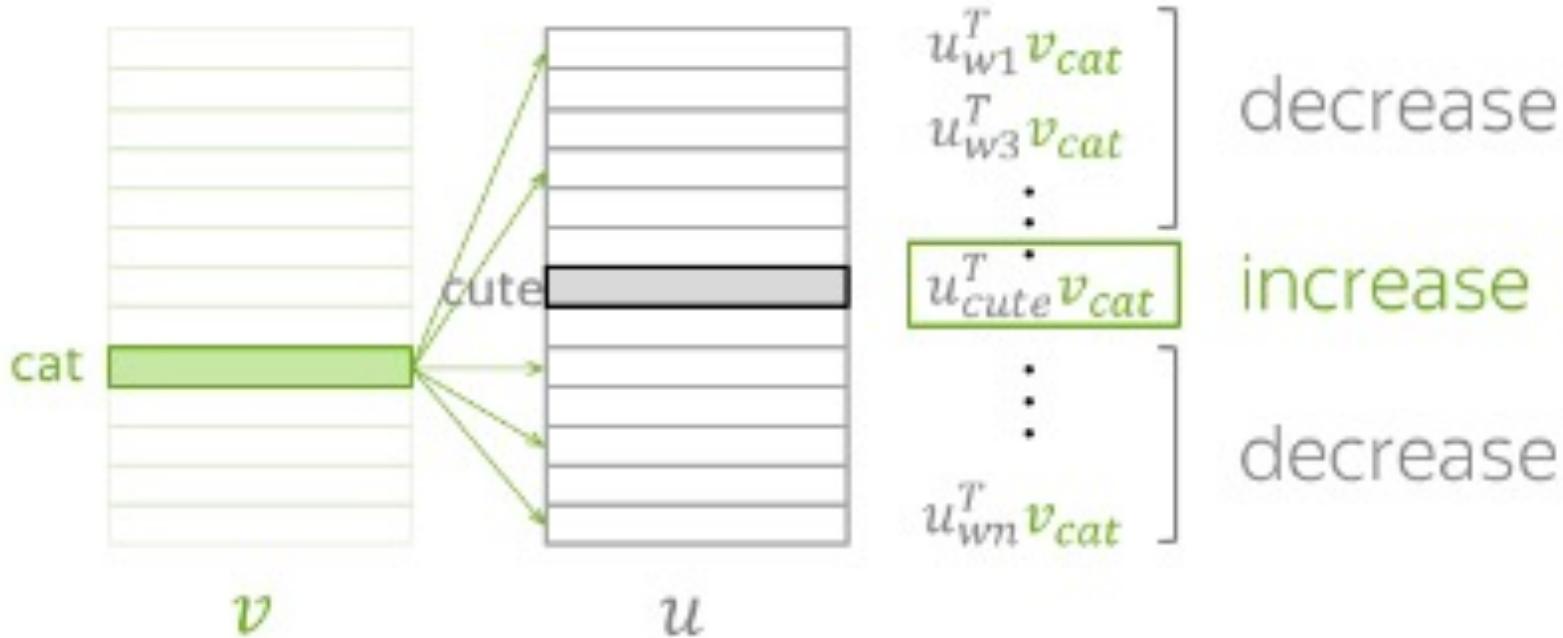
Softmax: $-(\log(\text{num}/\text{den})) = -(\log(\text{num}) - \log(\text{den})) = -\log(\text{num}) + \log(\text{den})$

4. get loss (for this one step)

$$J_{t,j}(\theta) = \underbrace{-u_{cute}^T v_{cat}}_1 + \log \underbrace{\sum_{w \in V} \exp(u_w^T v_{cat})}_2$$

5. evaluate the gradient,
make an update

$$v_{cat} := v_{cat} - \alpha \frac{\partial J_{t,j}(\theta)}{\partial v_{cat}}$$
$$u_w := u_w - \alpha \frac{\partial J_{t,j}(\theta)}{\partial u_w} \quad \forall w \in V$$



- The doc product increases as the vectors become similar
- The words which appear in a similar context, the embeddings become similar

Negative Sampling

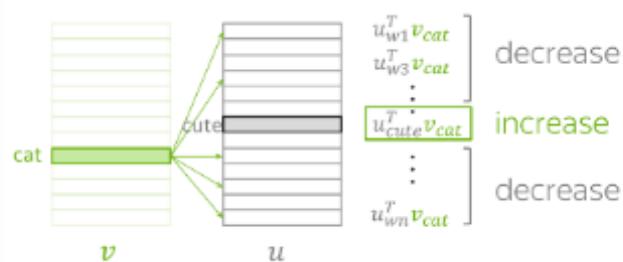
Dot product of v_{cat} :

- with u_{cute} - increase,
- with all other u - decrease

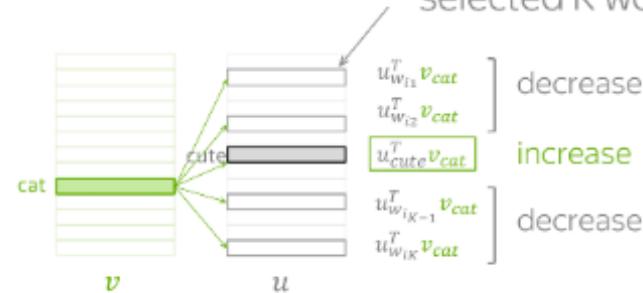


Dot product of v_{cat} :

- with u_{cute} - increase,
- with a subset of other u - decrease



Negative samples: randomly selected K words



Parameters to be updated:

- v_{cat}
 - u_w for all w in the vocabulary
- $|V| + 1$ vectors

Parameters to be updated:

- v_{cat}
 - u_{cute} and u_w for w in K negative examples
- K + 2 vectors

Selecting Negative Samples

Could pick w according to their unigram frequency $P(w)$

More common to chosen then according to $p_\alpha(w)$

$$P_\alpha(w) = \frac{\text{count}(w)^\alpha}{\sum_w \text{count}(w)^\alpha}$$

$\alpha = \frac{3}{4}$ works well because it gives rare noise words slightly higher probability

To show this, imagine two events $p(a) = .99$ and $p(b) = .01$:

$$P_\alpha(a) = \frac{.99^{.75}}{.99^{.75} + .01^{.75}} = .97$$

$$P_\alpha(b) = \frac{.01^{.75}}{.99^{.75} + .01^{.75}} = .03$$

Sample

Source Text

The quick brown fox jumps over the lazy dog. →

The quick brown fox jumps over the lazy dog. →

The quick brown fox jumps over the lazy dog. →

The quick brown fox jumps over the lazy dog. →

Training Samples

(the, quick)
(the, brown)

(quick, the)
(quick, brown)
(quick, fox)

(brown, the)
(brown, quick)
(brown, fox)
(brown, jumps)

(fox, quick)
(fox, brown)
(fox, jumps)
(fox, over)

Problem – frequent Words

There are two "problems" with common words like "the":

- When looking at word pairs, ("fox", "and the") do not tell us much about the meaning of "fox".
"the" appears in the context of pretty much every word.
- We will have many more samples of ("the", ...) than we need to learn a good vector for "the".

Solution – Sub Sampling

If we have a window size of 10, and we remove a specific instance of "the" from our text:

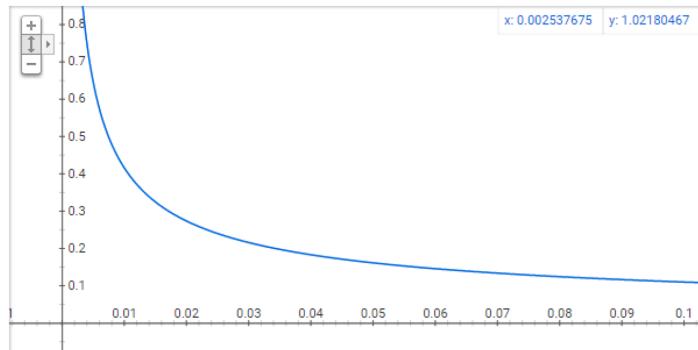
- As we train on the remaining words, "the" will not appear in any of their context windows.
- We'll have 10 fewer training samples where "the" is the input word.

Sampling rate

probability of *keeping* the word:

$$P(w_i) = \left(\sqrt{\frac{z(w_i)}{0.001}} + 1 \right) \cdot \frac{0.001}{z(w_i)}$$

Graph for [\(\$\text{sqrt}\(x/0.001\)+1\$ \)* \$0.001/x\$](#)



the fraction of the total words in the corpus
that are that word.

- $P(w_i) = 1.0$ if $z(w_i) < -0.0026$
- $P(w_i) = 0.5$ if $z(w_i)$ when $z(w_i) = 0.00746$
- $P(w_i) = 0.033$ if $z(w_i)$ when $z(w_i) = 1$
- For Wikipedia only top 27 most frequent words would have any subsampling applied.

Rank	Word	Portion of Corpus	Probability to Keep
0	the	7.80%	12.60%
1	of	3.71%	19.10%
2	and	3.19%	20.84%

Context Position Weighting

Randomly shrink the size of context window

For example if window size is 5 – This means we can pick any size window with equal probability.

This means that words closest to Target word are more likely to be chosen as positive samples.

Random Window Size = 4

"Since the more distant words are usually less related to the current word than those close to it, we give less weight to the distant words by sampling less from those words in our training examples."

Random Window Size = 2

"Since the more distant words are usually less related to the current word than those close to it, we give less weight to the distant words by sampling less from those words in our training examples."

Random Window Size = 3

"Since the more distant words are usually less related to the current word than those close to it, we give less weight to the distant words by sampling less from those words in our training examples."

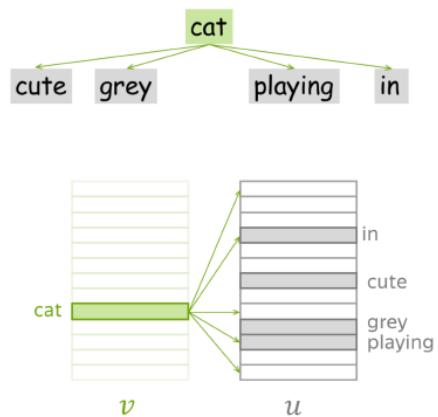
Random Window Size = 1

"Since the more distant words are usually less related to the current word than those close to it, we give less weight to the distant words by sampling less from those words in our training examples."

Word2Vec Variants: Skip-Gram and CBOW

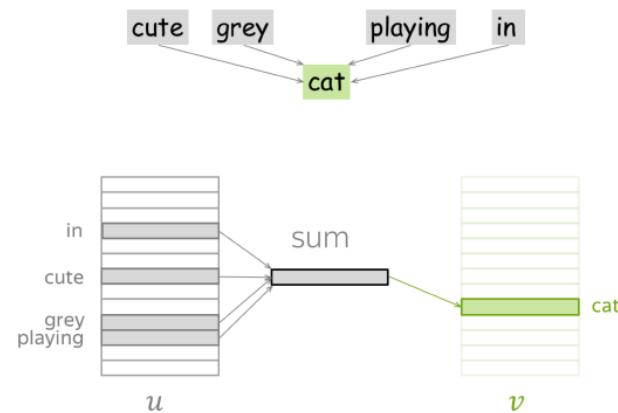
... I saw a cute grey cat playing in the garden ...

Skip-Gram: from **central** predict context
(one at a time)



(this is what we did so far)

CBOW: from sum of context predict **central**



(Continuous Bag of Words)

CBOW

The quick brown fox jumps over the lazy dog.

Positive Training Samples

Sentence

The quick brown fox jumps over the lazy dog. →

Skip-Gram
(the, quick)
(the, brown)

The quick brown fox jumps over the lazy dog. →

Skip-Gram
(quick, the)
(quick, brown)
(quick, fox)

The quick brown fox jumps over the lazy dog. →

Skip-Gram
(brown, the)
(brown, quick)
(brown, fox)
(brown, jumps)

The quick brown fox jumps over the lazy dog. →

Skip-Gram
(fox, quick)
(fox, brown)
(fox, jumps)
(fox, over)

CBOW

((quick, brown), the)

((the, brown, fox), quick)

((the, quick, fox, jumps), brown)

((quick, brown, jumps, over), fox)

Short Comings of Word2Vec

- Out of Vocabulary Words



- Morphology

Shared radical

eat eats eaten eater eating

FastText – Sub-word Generation

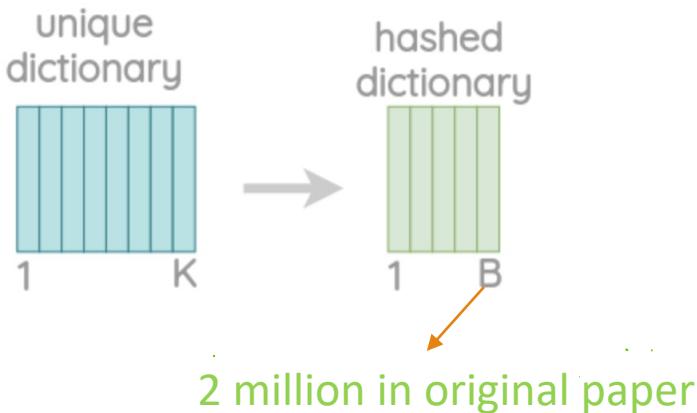
Steps

- (1) Take a word and add angular brackets to denote the beginning and end of a word
eating → <eating>
- (2) Generate character n-grams of length n

<eating>
3-grams <ea eat ati tin ing ng>

Word	Length(n)	Character n-grams
eating	3	<ea, eat, ati, tin, ing, ng>
eating	4	<eat, eati, atin, ting, ing>
eating	5	<eati, eatin, ating, ting>
eating	6	<eatin, eating, ating>

- 3a. Instead of learning an embedding for each unique n-gram, we learn total B embeddings where B denotes the bucket size



Source :

<https://amitness.com/2020/06/fasttext-embeddings/>

FastText – Embedding

Steps

- 3b. Each character n-gram is hashed to an integer between 1 to B

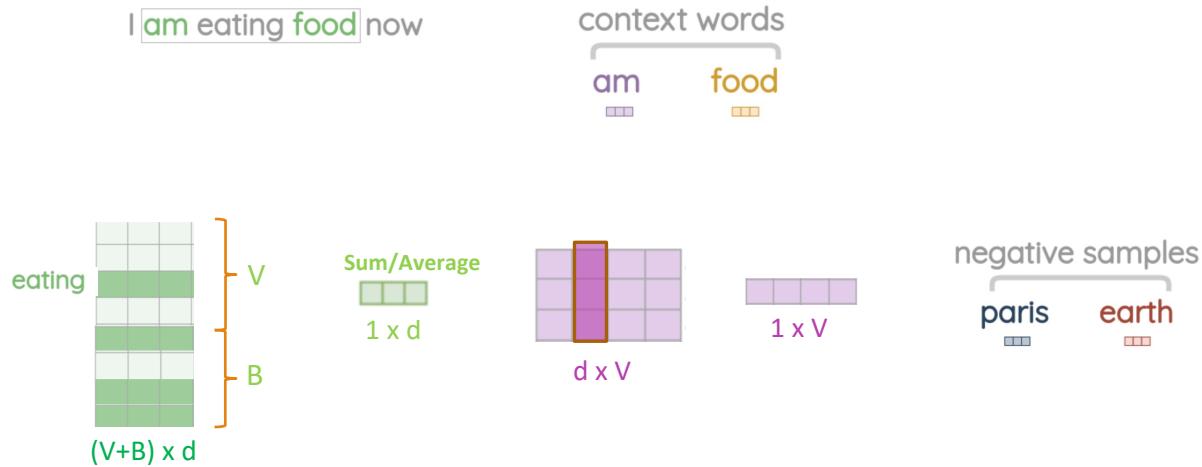


- Many to One Mapping – There could be collisions
- embedding for the center word is calculated by taking a sum of vectors for the character n-grams and the whole word itself

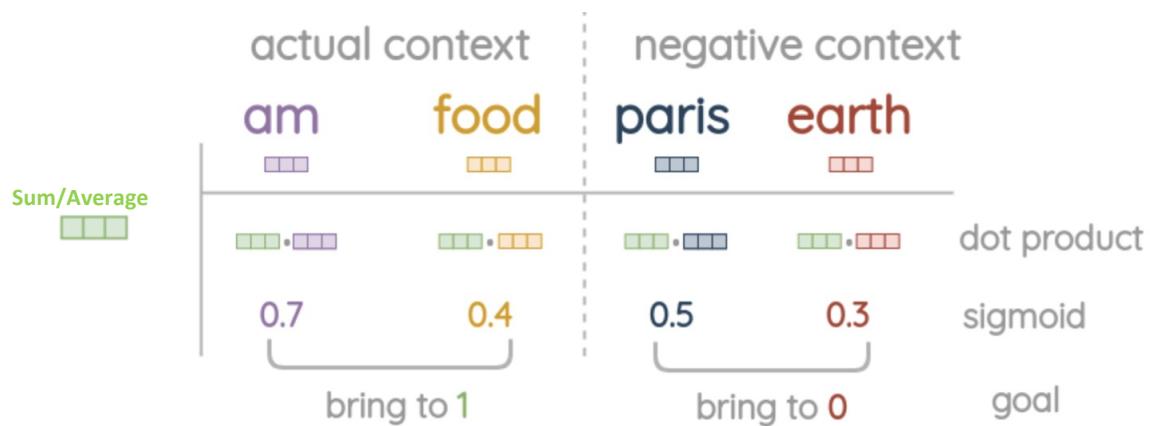


Source : <https://amitness.com/2020/06/fasttext-embeddings/>

FastText – Training



V = Vocab size, B = Bucket Size, d = embedding size



FastText vs Word2Vec

- FastText – more useful to learn embeddings for rare words and can deal with OOV (Out of off vocabulary) words
- Sub words do not improve representations of common words and may sometimes lead to slightly worse representation
- FastText improves performance on syntactic word analogy tasks significantly for morphologically rich language like Czech and German
- FastText has degraded performance on semantic analogy tasks compared to Word2Vec.
- FastText can be slower to train than regular skipgram

Summary:

- Fasttext is an extension of Word2Vec and has pros and cons.

Capture co-occurrence counts directly?

With a co-occurrence matrix X

- 2 options: windows vs. full document
- Window: Similar to word2vec, use window around each word → captures both syntactic (POS) and semantic information
- Word-document co-occurrence matrix will give general topics (all sports terms will have similar entries) leading to “Latent Semantic Analysis”

Window based co-occurrence matrix

- Example corpus:
 - I like deep learning.
 - I like NLP.
 - I enjoy flying

counts	I	like	enjoy	deep	learning	NLP	flying	.
I	0	2	1	0	0	0	0	0
like	2	0	0	1	0	1	0	0
enjoy	1	0	0	0	0	0	1	0
deep	0	1	0	0	1	0	0	0
learning	0	0	0	1	0	0	0	1
NLP	0	1	0	0	0	0	0	1
flying	0	0	1	0	0	0	0	1
.	0	0	0	0	1	1	1	0

Source: <http://web.stanford.edu/class/cs224n/slides/cs224n-2020-lecture02-wordvecs2.pdf>

Method: Dimensionality Reduction on X

Singular Value Decomposition of co-occurrence matrix X

Factorizes X into $U\Sigma V^T$, where U and V are orthonormal

$$\underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_{X^k} = \underbrace{\begin{bmatrix} * & * \\ * & * \\ * & * \end{bmatrix}}_U \underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_{\Sigma} \underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_{V^T}$$

Retain only k singular values, in order to generalize.

Expensive to compute for large matrices.

Towards GloVe: Count based vs. direct prediction

- LSA, HAL (Lund & Burgess),
- COALS, Hellinger-PCA (Rohde et al, Lebret & Collobert)

- Fast training
- Efficient usage of statistics
- Primarily used to capture word similarity
- Disproportionate importance given to large counts

- Skip-gram/CBOW (Mikolov et al)
- NNLM, HLBL, RNN (Bengio et al; Collobert & Weston; Huang et al; Mnih & Hinton)

- Scales with corpus size
- Inefficient usage of statistics
- Generate improved performance on other tasks
- Can capture complex patterns beyond word similarity

Source: <http://web.stanford.edu/class/cs224n/slides/cs224n-2020-lecture02-wordvecs2.pdf>

GloVe (Global Vectors) model

- Combines the advantages of the two major model families in the literature:
 - global matrix factorization and
 - local context window methods
- It leverages statistical information by training only on the nonzero elements in a word-word co-occurrence matrix, rather than on the entire sparse matrix or on individual context windows in a large corpus.

Original Glove Paper :

<https://nlp.stanford.edu/pubs/glove.pdf>

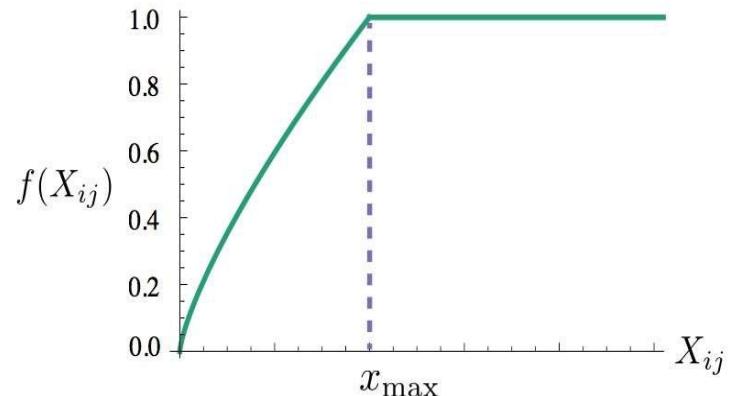
Combining the best of both worlds GloVe

[Pennington et al., EMNLP 2014]

$$J = \sum_{i,j=1}^V f(X_{ij}) \left(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij} \right)^2$$

$$f(x) = \begin{cases} (x/x_{\max})^{\alpha} & \text{if } x < x_{\max} \\ 1 & \text{otherwise} \end{cases}$$

- Fast training
- Scalable to huge corpora
- Good performance even with small corpus and small vectors



Original Glove Paper :
<https://nlp.stanford.edu/pubs/glove.pdf>

Figure 1: Weighting function f with $\alpha = 3/4$.

GloVe Results

Nearest words to **frog**:

1. frogs
2. toad
3. litoria
4. leptodactylidae
5. rana
6. lizard
7. eleutherodactylus



litoria

leptodactylidae



rana



eleutherodactylus

Source: <http://web.stanford.edu/class/cs224n/slides/cs224n-2020-lecture02-wordvecs2.pdf>

Comparison

Feature/Model	Word2Vec	FastText	GloVe
Core Idea	Predictive model using local context windows	Predictive model using local context windows and subword information	Count-based model using aggregated global word co-occurrence statistics from the entire corpus
Context Window	Fixed size	Fixed size, but includes character n-grams	Uses a fixed-size context window to build the co-occurrence matrix but aggregates statistics across the entire corpus
Handling of OOV Words	Cannot handle OOV words	Can handle OOV words through subword information	Cannot handle OOV words without additional methods
Morphological Richness	Poor (does not consider subword information)	Excellent (captures morphological information)	Poor (does not consider subword information)
Semantic Relationships	Good (especially syntactic analogies)	Good (captures morphological semantics)	Good (captures both global co-occurrence and local context semantics)
Use Cases	General NLP tasks, syntactic relationships	Morphologically rich languages, text classification	Tasks requiring understanding of global word usage patterns
Strengths	Efficient, well-researched, good at capturing syntactic relationships	Handles OOV and morphologically complex words, fast predictions	Captures global statistics efficiently, good for word analogy and similarity
Weaknesses	Struggles with OOV words, does not use global statistics	Larger model size, potential noise from Subwords	Requires significant memory for co-occurrence matrix, less effective with simple morphologies

Evaluating embeddings -Intrinsic

Compare to human scores on word similarity-type tasks:

- WordSim-353 (Finkelstein et al., 2002)
- SimLex-999 (Hill et al., 2015)
- Stanford Contextual Word Similarity (SCWS) dataset (Huang et al., 2012)
- TOEFL dataset: *Levied is closest in meaning to: imposed, believed, requested, correlated*

Evaluating embeddings -Extrinsic

Use embeddings in Downstream tasks like
Classification

Properties of embeddings

Similarity depends on window size C

C = ± 2 The nearest words to *Hogwarts*:

- *Sunnydale*
- *Evernight*

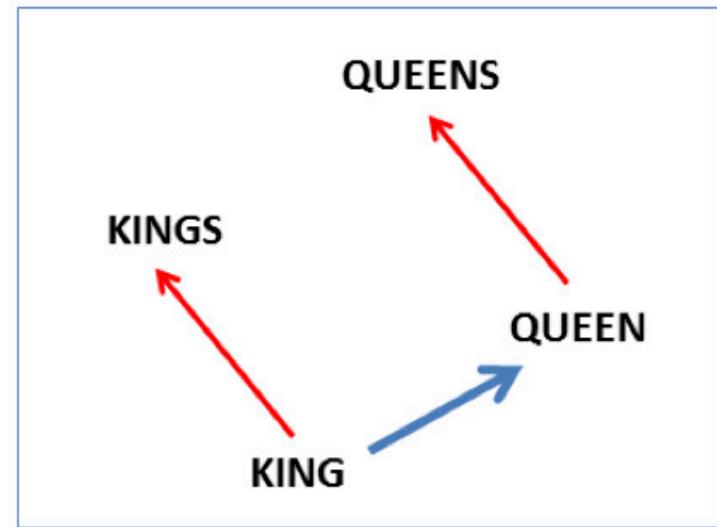
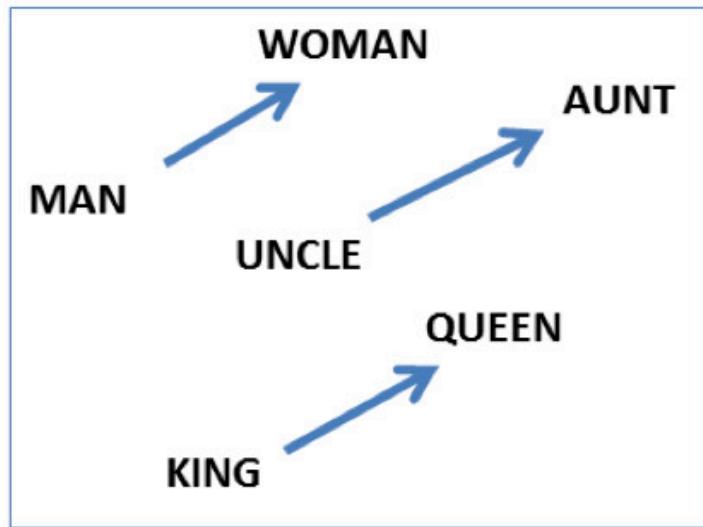
C = ± 5 The nearest words to *Hogwarts*:

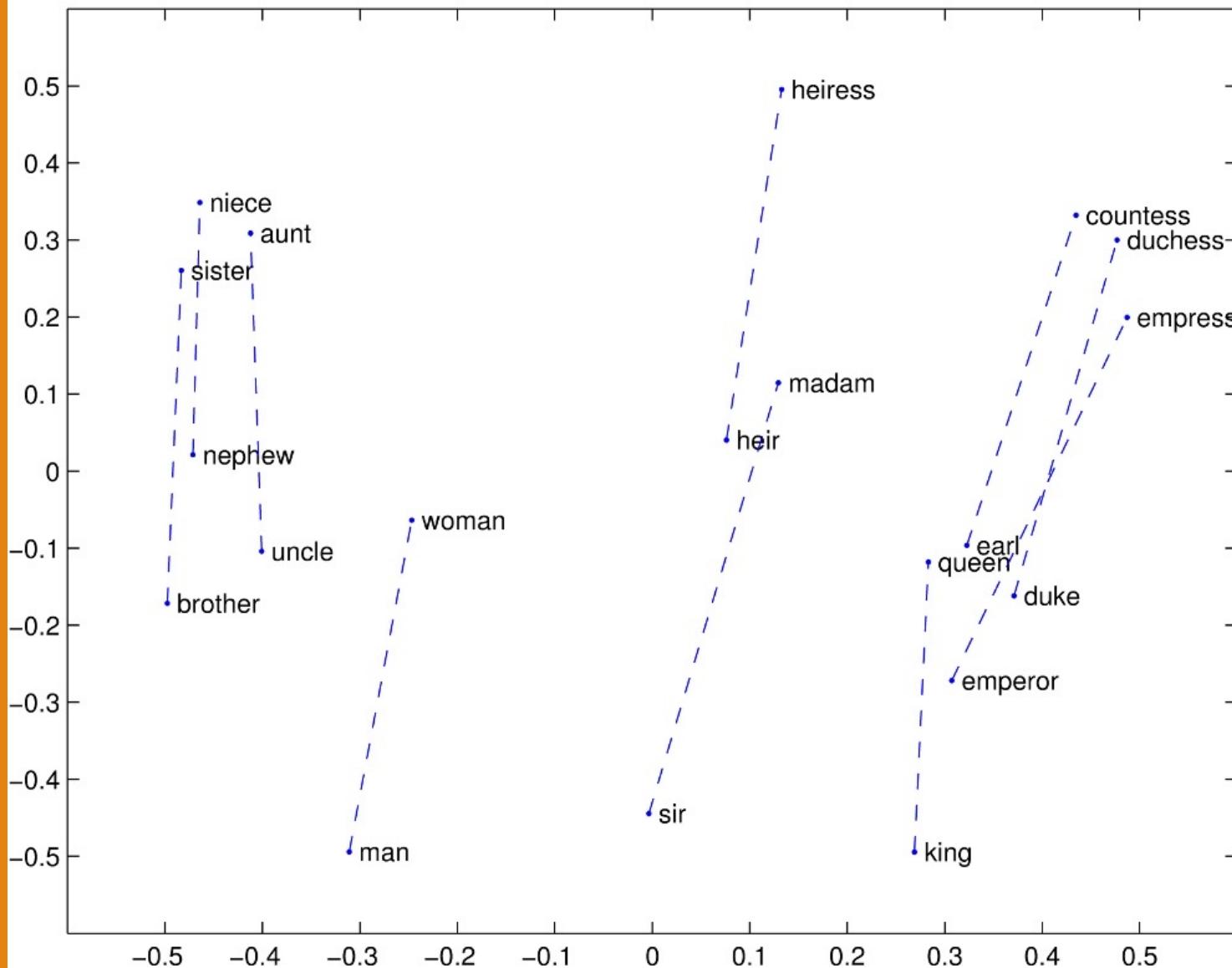
- *Dumbledore*
- *Malfoy*
- *halfblood*

Analogy: Embeddings capture relational meaning!

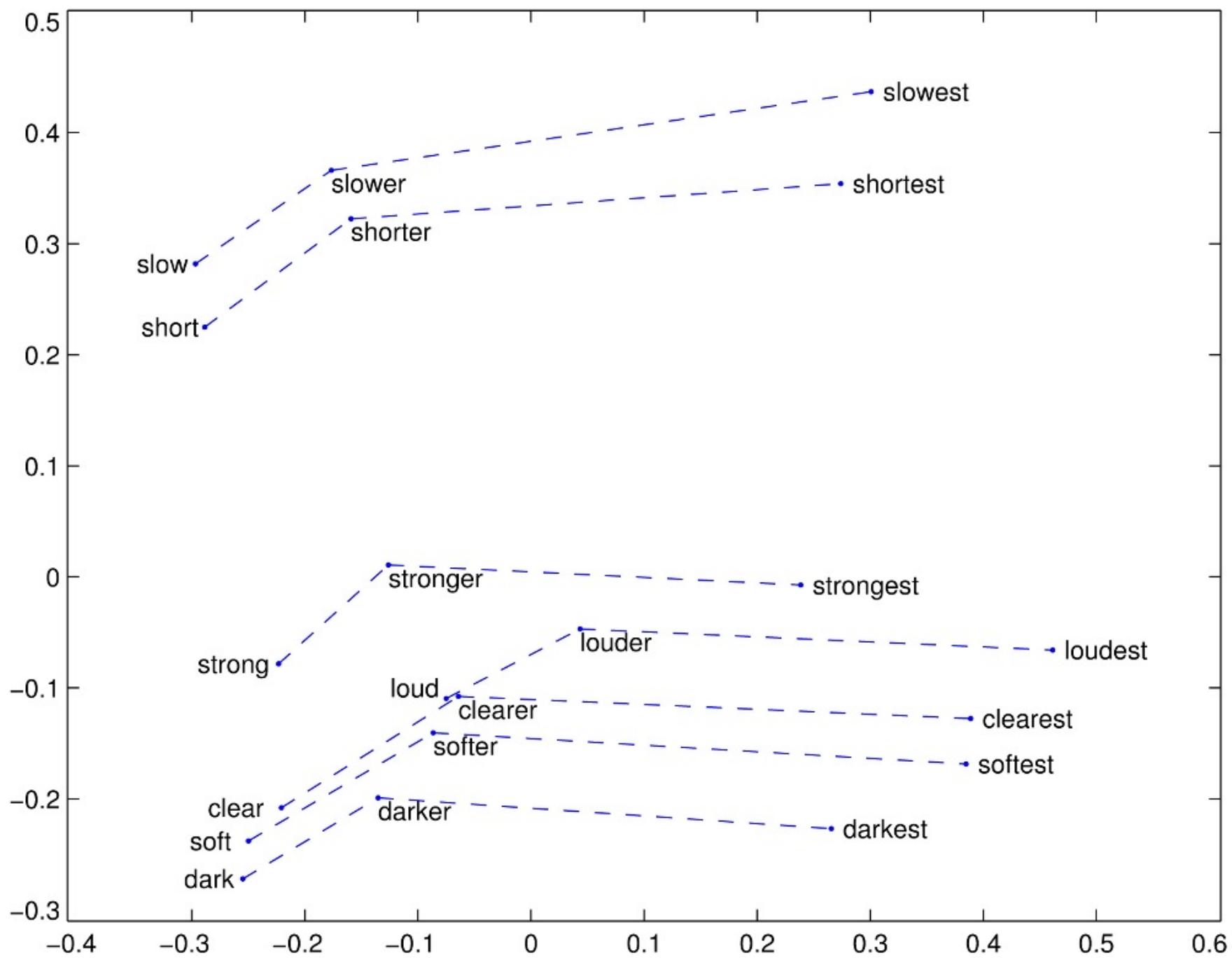
$\text{vector('king')} - \text{vector('man')} + \text{vector('woman')} \approx \text{vector('queen')}$

$\text{vector('Paris')} - \text{vector('France')} + \text{vector('Italy')} \approx \text{vector('Rome')}$





Source : Chapter 6 of Book: <https://web.stanford.edu/~jurafsky/slp3/>



Embeddings and bias

Embeddings reflect cultural bias

Bolukbasi, Tolga, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam T. Kalai. "Man is to computer programmer as woman is to homemaker? debiasing word embeddings." In *Advances in Neural Information Processing Systems*, pp. 4349-4357. 2016.

Ask “Paris : France :: Tokyo : x”

- x = Japan

Ask “father : doctor :: mother : x”

- x = nurse

Ask “man : computer programmer :: woman : x”

- x = homemaker

Embeddings reflect cultural bias

Caliskan, Aylin, Joanna J. Bruson and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science* 356:6334, 183-186.

Implicit Association test (Greenwald et al 1998): How associated are

- concepts (*flowers, insects*) & attributes (*pleasantness, unpleasantness*)?
- Studied by measuring timing latencies for categorization.

Psychological findings on US participants:

- African-American names are associated with unpleasant words (more than European-American names)
- Male names associated more with math, female names with arts
- Old people's names with unpleasant words, young people with pleasant words.

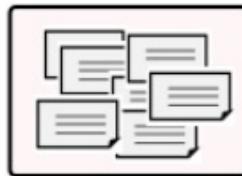
Caliskan et al. replication with embeddings:

- African-American names (*Leroy, Shaniqua*) had a higher GloVe cosine with unpleasant words (*abuse, stink, ugly*)
- European American names (*Brad, Greg, Courtney*) had a higher cosine with pleasant words (*love, peace, miracle*)

Embeddings reflect and replicate all sorts of pernicious biases.

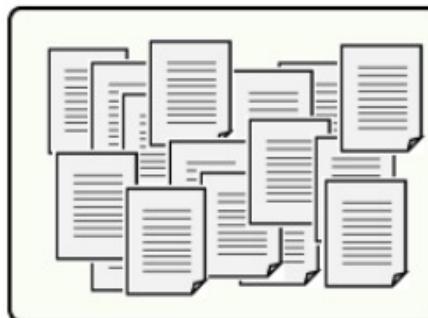
Source : Chapter 6 of Book: <https://web.stanford.edu/~jurafsky/slp3/>

Using Dense Embeddings



Training data for text classification (labeled)

- Not huge, or not diverse, or both
- Domain: task-specific



Training data for word embeddings (unlabeled)

- Huge diverse corpus (e.g., Wikipedia)
- Domain: general

Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

Pre-Trained Embeddings

- Train from scratch

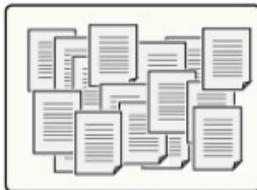
What they will know:



May be not enough to learn relationships between words

- Take pretrained (Word2Vec, GloVe)

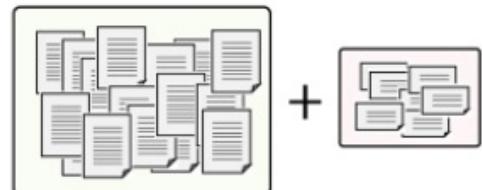
What they will know:



Know relationships between words, but are not specific to the task

- Initialize with pretrained, then fine-tune

What they will know:

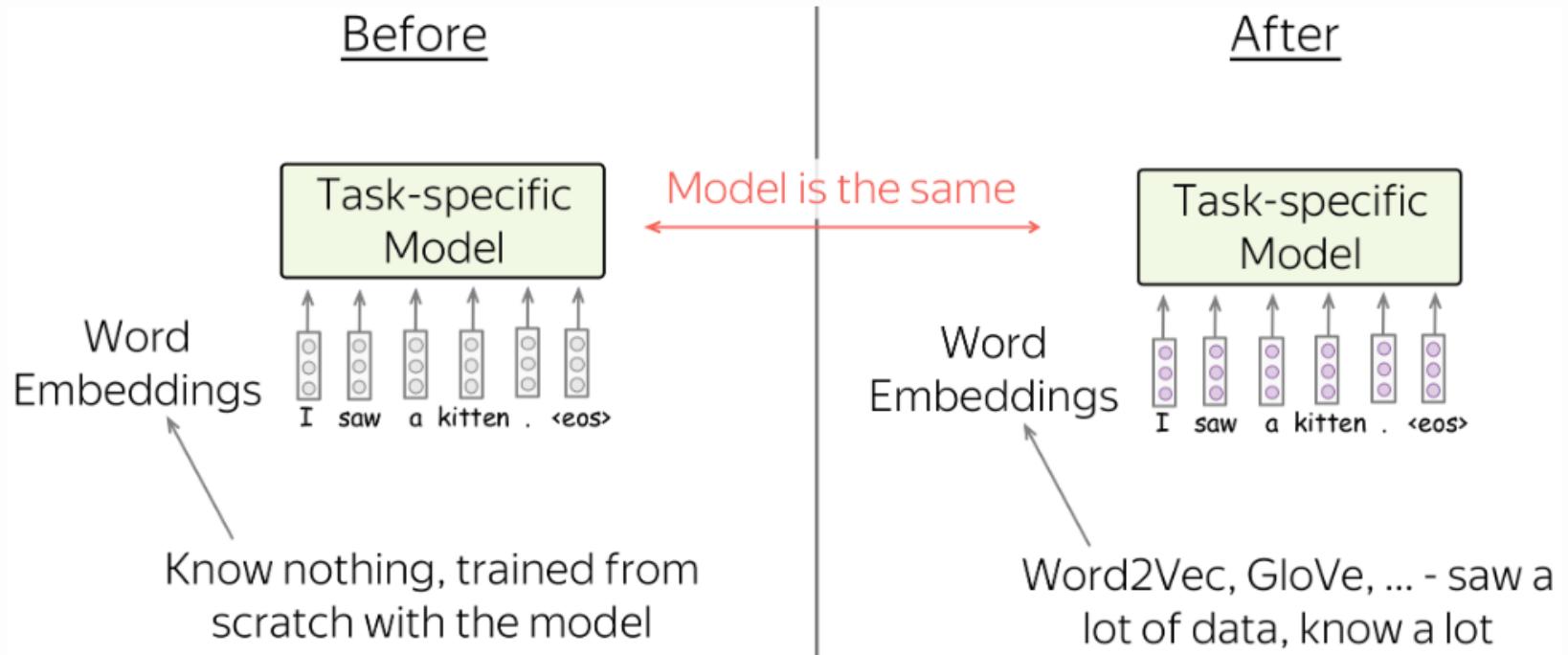


Know relationships between words and adapted for the task

"Transfer" knowledge from a huge unlabeled corpus to your task-specific model

Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

Pre-Trained Embeddings



Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

Trained embeddings

- Word2vec
<https://code.google.com/archive/p/word2vec/>
- Glove
<http://nlp.stanford.edu/projects/glove/>
- Levy/Goldberg dependency embeddings
<https://levyomer.wordpress.com/2014/04/25/dependency-based-word-embeddings/>

Source : <https://people.ischool.berkeley.edu/~dbamman/info256.html>

Improvement over tfidf

- **Semantic Meaning:** Word2Vec and GloVe embeddings capture deep semantic meaning by considering the context in which words appear. This allows words with similar meanings to have similar embeddings, which isn't the case with TF-IDF.
- **Word Relationships:** These models can capture complex relationships between words, such as analogies ("king" is to "queen" as "man" is to "woman"), which TF-IDF cannot do.
- **Dimensionality:** Word2Vec and GloVe embeddings are typically of a much lower dimension than TF-IDF, making them more computationally efficient while retaining more useful information.
- **Handling Synonyms:** By considering word contexts, Word2Vec and GloVe are better at associating synonyms with similar vectors, while TF-IDF treats all words as independent

Limitation still remaining

Word2Vec and GloVe provided significant advancements over TF-IDF embeddings, but there were still issues and limitations that remained even with these more sophisticated models:

- **OOV Words:** Word2Vec and GloVe can't create embeddings for words not in their training data, a challenge for the ever-changing vocabulary of human languages.
- **Polysemy:** Both assign one vector per word, unable to differentiate meanings for words like "bank," which can signify a river's edge or a financial institution.
- **Static Representations:** Word2Vec and GloVe produce fixed embeddings, not adapting to varying word meanings in different contexts.
- **Word Order Sensitivity:** These models miss the word order, so phrases like "free drug" and "drug free" might get similar representations despite contrasting meanings.

Two kinds of “training” data

- The labeled data for a specific task (e.g., labeled sentiment for movie reviews): ~ 2K labels/reviews, ~1.5M words → used to train a supervised model → Supervised training
- General text (Wikipedia, the web, books, etc.), ~ trillions of words → used to train word distributed representations → Self-supervised training (We can train the model on generic text – no manually labelled data required)

Source : <https://people.ischool.berkeley.edu/~dbamman/info256.html>

Self Supervised training

Self-supervised learning is a type of machine learning that does not require explicit external labeling of data. Instead, it generates its own supervisory signal based on the input data.

Can we think of another such task?

Source : <https://people.ischool.berkeley.edu/~dbamman/info256.html>

Self Supervised training

Self-supervised learning is a type of machine learning that does not require explicit external labeling of data. Instead, it generates its own supervisory signal based on the input data.

Can we think of other such task?

Language Model

Source : <https://people.ischool.berkeley.edu/~dbamman/info256.html>



Language Modeling & RNN

Probabilistic Language Models

Goal: assign a probability to a sentence

- Machine Translation:
 - $P(\text{high winds tonite}) > P(\text{large winds tonite})$
- Spell Correction
 - The office is about fifteen **minuets** from my house
 - $P(\text{about fifteen minutes from}) > P(\text{about fifteen minuets from})$
- Speech Recognition
 - $P(\text{I saw a van}) \gg P(\text{eyes awe of an})$
 - + Summarization, question-answering, etc., etc.!!

Why?

Probabilistic Language Modeling

Goal: compute the probability of a sentence or sequence of words:

$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$

Related task: probability of an upcoming word:

$$P(w_5 | w_1, w_2, w_3, w_4)$$

A model that computes either of these:

$P(W)$ or $P(w_n | w_1, w_2 \dots w_{n-1})$ is called a **language model**.

Better: **the grammar** But **language model** or **LM** is standard

Why should we care about Language Modeling?

- Language Modeling is a **benchmark task** that helps us **measure our progress** on understanding language
- Language Modeling is a **subcomponent** of many NLP tasks, especially those involving **generating text** or **estimating the probability of text**:
 - Predictive typing
 - Speech recognition
 - Handwriting recognition
 - Spelling/grammar correction
 - Authorship identification
 - Machine translation
 - Summarization
 - Dialogue
 - etc.

How to compute $P(W)$

How to compute this joint probability:

- $P(\text{its, water, is, so, transparent, that})$

Intuition: let's rely on the Chain Rule of Probability

Reminder: The Chain Rule

Recall the definition of conditional probabilities

$$p(B|A) = P(A,B)/P(A) \quad \text{Rewriting: } P(A,B) = P(A)P(B|A)$$

More variables:

$$P(A,B,C,D) = P(A)P(B|A)P(C|A,B)P(D|A,B,C)$$

The Chain Rule in General

$$P(x_1, x_2, x_3, \dots, x_n) = P(x_1)P(x_2|x_1)P(x_3|x_1, x_2)\dots P(x_n|x_1, \dots, x_{n-1})$$

The Chain Rule applied to compute joint probability of words in sentence

$$P(w_1 w_2 \dots w_n) = \prod_i P(w_i | w_1 w_2 \dots w_{i-1})$$

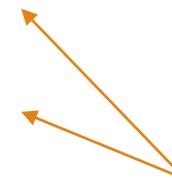
$P(\text{"its water is so transparent"}) =$

$P(\text{its}) \times P(\text{water}|\text{its}) \times P(\text{is}|\text{its water})$

$\times P(\text{so}|\text{its water is}) \times P(\text{transparent}|\text{its water is so})$

How to estimate these probabilities

$$P(w_n | w_1, w_2, \dots, w_{n-1}) = \frac{N(w_1, w_2, \dots, w_n)}{N(w_1, w_2, \dots, w_{n-1})}$$



where $N(w_1, w_2, \dots, w_n)$ is the number of times
 (w_1, w_2, \dots, w_n) occurs in the corpus

These are too long, and
many of the counts will
be zero

$P(\text{the lits water is so transparent that}) =$

$\frac{\text{Count}(\text{its water is so transparent that the})}{\text{Count}(\text{its water is so transparent that})}$

$\text{Count}(\text{its water is so transparent that})$

Markov Assumption

The probability of a word only depends on a **fixed** number of previous words.

For an n-gram model,

$$P(y_t | \underbrace{y_1, y_2, \dots, y_{t-1}}_{\text{all previous tokens}}) = P(y_t | \underbrace{y_{t-n+1}, \dots, y_{t-1}}_{n-1 \text{ previous token}})$$

For example,

- n=3 (trigram model): $P(y_t | y_1, y_2, \dots, y_{t-1}) = P(y_t | y_{t-2}, y_{t-1})$
- n=2 (bigram model): $P(y_t | y_1, y_2, \dots, y_{t-1}) = P(y_t | y_{t-1})$
- n=1 (unigram model): $P(y_t | y_1, y_2, \dots, y_{t-1}) = P(y_t)$

Markov Assumption: Example

$P(I \text{ saw a cat on a mat}) = ?$

Before

$P(I \text{ saw a cat on a mat}) =$

$P(I)$

- $P(\text{saw} | I)$
- $P(a | I \text{ saw})$
- $P(\text{cat} | I \text{ saw a})$
- $P(\text{on} | I \text{ saw a cat})$
- $P(a | I \text{ saw a cat on})$
- $P(\text{mat} | I \text{ saw a cat on a})$

After (3-gram)

$P(I \text{ saw a cat on a mat}) =$



$P(I)$

- $P(\text{saw} | I)$ → $P(I)$
 - $P(a | I \text{ saw})$ → $\cdot P(\text{saw} | I)$
 - $P(\text{cat} | I \text{ saw a})$ → $\cdot P(a | I \text{ saw})$
 - $P(\text{on} | I \text{ saw a cat})$ → $\cdot P(\text{cat} | \text{saw a})$
 - $P(a | I \text{ saw a cat on})$ → $\cdot P(\text{on} | \text{a cat})$
 - $P(\text{mat} | I \text{ saw a cat on a})$ → $\cdot P(\text{a} | \text{cat on})$
- ignore use

Estimating Probabilities using bigrams

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

<s> I am Sam </s>
<s> Sam I am </s>
<s> I do not like green eggs and ham </s>

$$P(\text{I} | \text{<s>}) = \frac{2}{3} = .67 \quad P(\text{Sam} | \text{<s>}) = \frac{1}{3} = .33 \quad P(\text{am} | \text{I}) = \frac{2}{3} = .67$$
$$P(\text{</ s>} | \text{Sam}) = \frac{1}{2} = 0.5 \quad P(\text{Sam} | \text{am}) = \frac{1}{2} = .5 \quad P(\text{do} | \text{I}) = \frac{1}{3} = .33$$

Raw bigram counts

Out of 9222 sentences

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

Raw bigram probabilities

$$P\left(\frac{\text{food}}{\text{chinese}}\right) = \frac{\text{count(chinese, food)}}{\text{count(chinese)}} = \frac{82}{158} = 0.52$$

$$P\left(\frac{w}{\text{chinese}}\right) = \frac{\text{count(chinese, } w\text{)}}{\text{count(chinese)}}$$

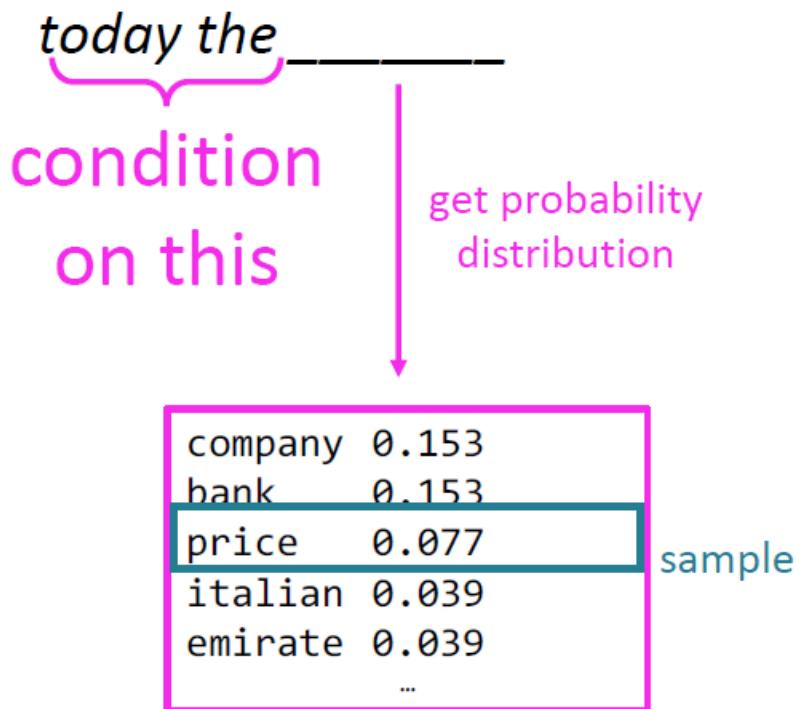
Normalize by unigrams:

i	want	to	eat	chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

Result:

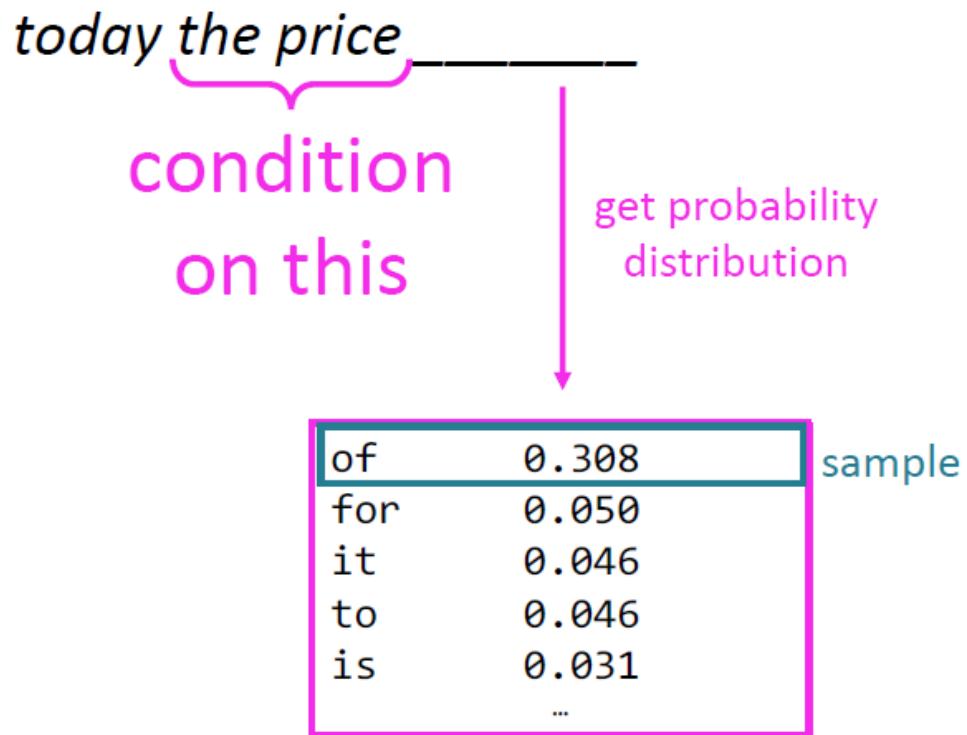
	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

Generating text with n-gram Language Model



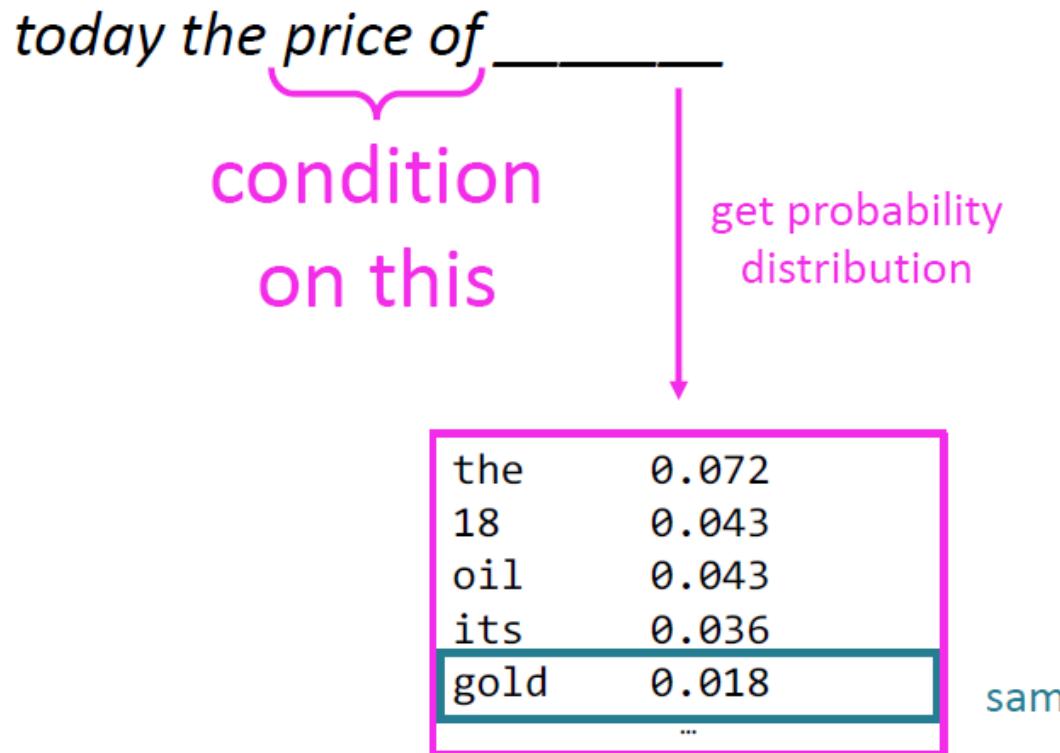
Slide credit: Christopher Manning , Abigail See and John Hewitt

Generating text with n-gram Language Model



Slide credit: Christopher Manning , Abigail See and John Hewitt

Generating text with n-gram Language Model



Slide credit: Abigail See and John Hewitt

Generating text with n-gram Language Model

Try Yourself:

<https://nlpforhackers.io/language-models/>

Issues n-gram Language Models

Suppose we are learning a **4-gram** Language Model.

~~as the proctor started the clock, the students opened their~~ _____
discard _____
condition on this _____

$$P(\mathbf{w}|\text{students opened their}) = \frac{\text{count(students opened their } \mathbf{w}\text{)}}{\text{count(students opened their)}}$$

For example, suppose that in the corpus:

- “students opened their” occurred 1000 times
 - “students opened their books” occurred 400 times
 - $\rightarrow P(\text{books} \mid \text{students opened their}) = 0.4$
 - “students opened their exams” occurred 100 times
 - $\rightarrow P(\text{exams} \mid \text{students opened their}) = 0.1$

Should we have discarded
the “proctor” context?

Can we increase n to improve model ?

Sparsity problem- n-gram Language Models

$$P(\text{mat} \mid \text{I saw a cat on a}) = P(\text{mat} \mid \text{cat on a}) = \frac{N(\text{cat on a mat})}{N(\text{cat on a})}$$

$$P(\text{mat} \mid \text{cat on a}) = \frac{N(\text{cat on a mat})}{N(\text{cat on a})} = ?$$

zero

not good: can not compute the probability

$$P(\text{mat} \mid \text{cat on a}) = \frac{\text{zero}}{N(\text{cat on a})} = ?$$

not good: zeros out probability of the sentence

Sparsity Problem 1

Sparsity Problem 2

Increasing n makes sparsity problem worse

SLIDE CREDIT: LENA VOITA (GITHUB.IO/NLP_COURSE.HTML#PREVIEW_LANG_MODELS)

Storage problem- n-gram Language Models

Storage: Need to store count for all n -grams you saw in the corpus.

$$P(\mathbf{w}|\text{students opened their}) = \frac{\text{count(students opened their } \mathbf{w}\text{)}}{\text{count(students opened their)}}$$

Increasing n or increasing corpus increases model size!

Fixed Window Neural Language Model

A fixed-window neural Language Model

as the proctor started the clock
discard

the students opened their _____
fixed window

Fixed Window Neural Language Model

A fixed-window neural Language Model

output distribution

$$\hat{y} = \text{softmax}(\mathbf{U}\mathbf{h} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden layer

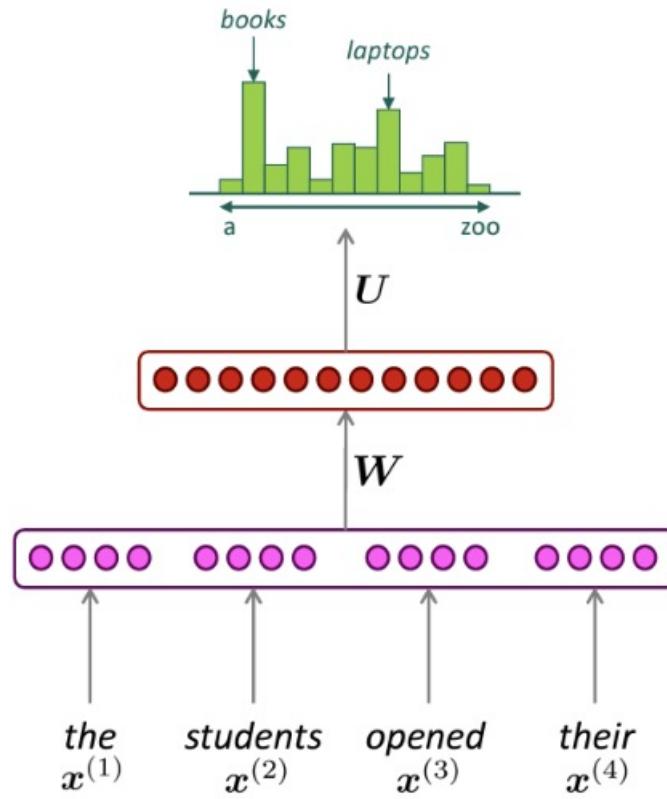
$$\mathbf{h} = f(\mathbf{W}\mathbf{e} + \mathbf{b}_1)$$

concatenated word embeddings

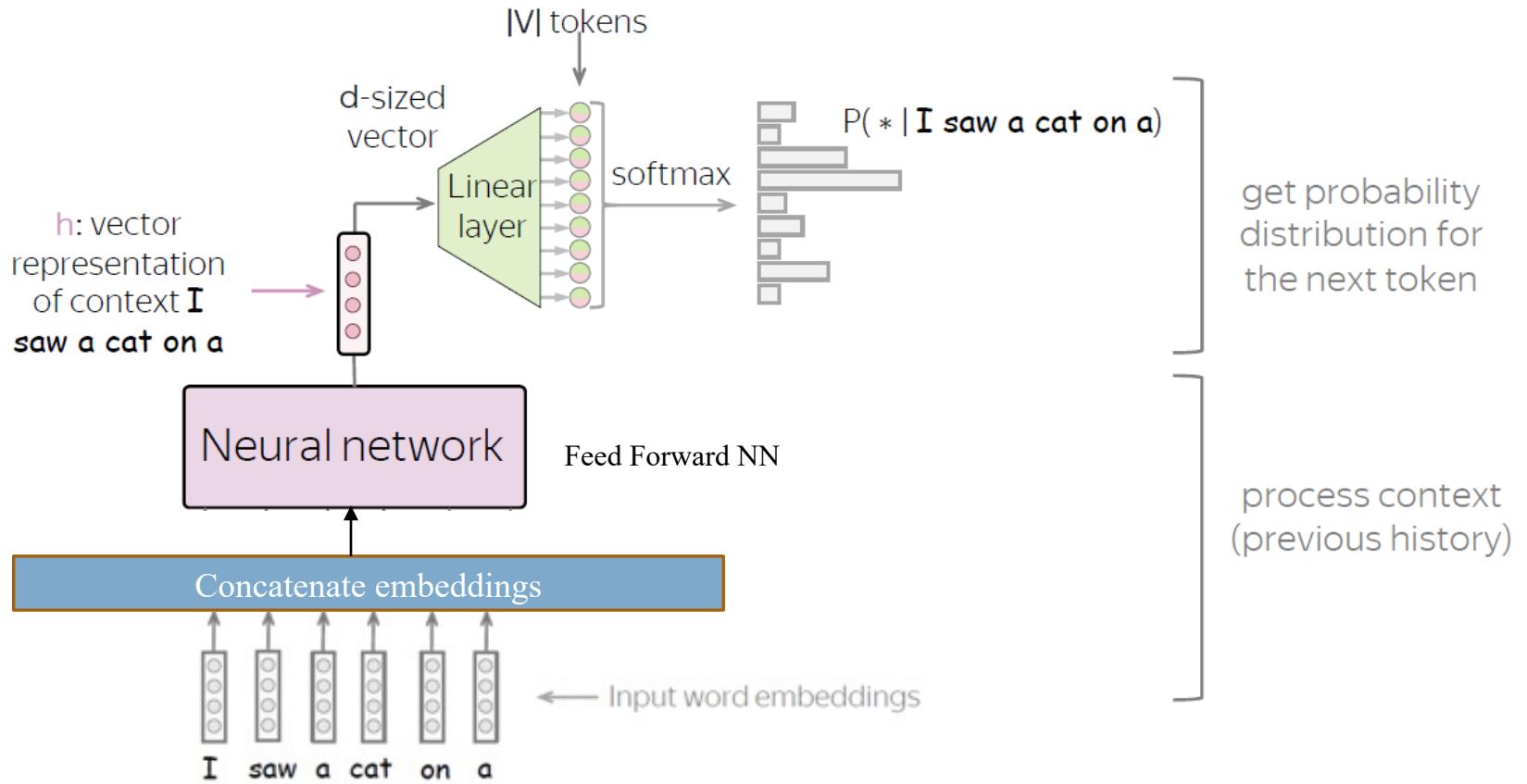
$$\mathbf{e} = [\mathbf{e}^{(1)}; \mathbf{e}^{(2)}; \mathbf{e}^{(3)}; \mathbf{e}^{(4)}]$$

words / one-hot vectors

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \mathbf{x}^{(4)}$$



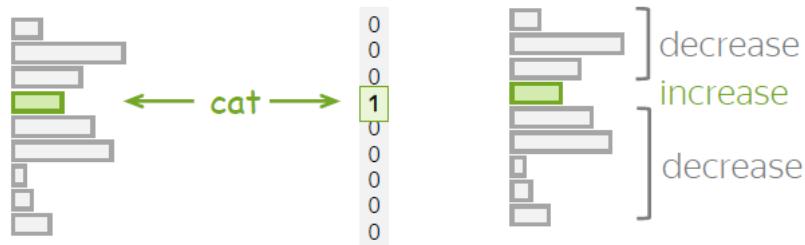
Fixed Window Neural Language Model



Training: Fixed Window Neural Language Model

Target at this step
↓
Training example: I saw a **cat** on a mat <eos>

Model prediction: Target: Loss = $-\log(p(\text{cat})) \rightarrow \min$
 $p(* | \text{I saw a})$ p^*



Cross-entropy loss:

$$-\sum_{i=1}^{|V|} p_i^* \cdot \log P(y_t = i | x) \rightarrow \min (p_k^* = 1, p_i^* = 0, i \neq k)$$

For one-hot targets, this is equivalent to

$$-\log P(y_t = \text{cat} | x) \rightarrow \min$$

Fixed Window Neural Language Model

Improvements: The model does not depend on the counts

- No Sparsity Problem
- No Storage problem

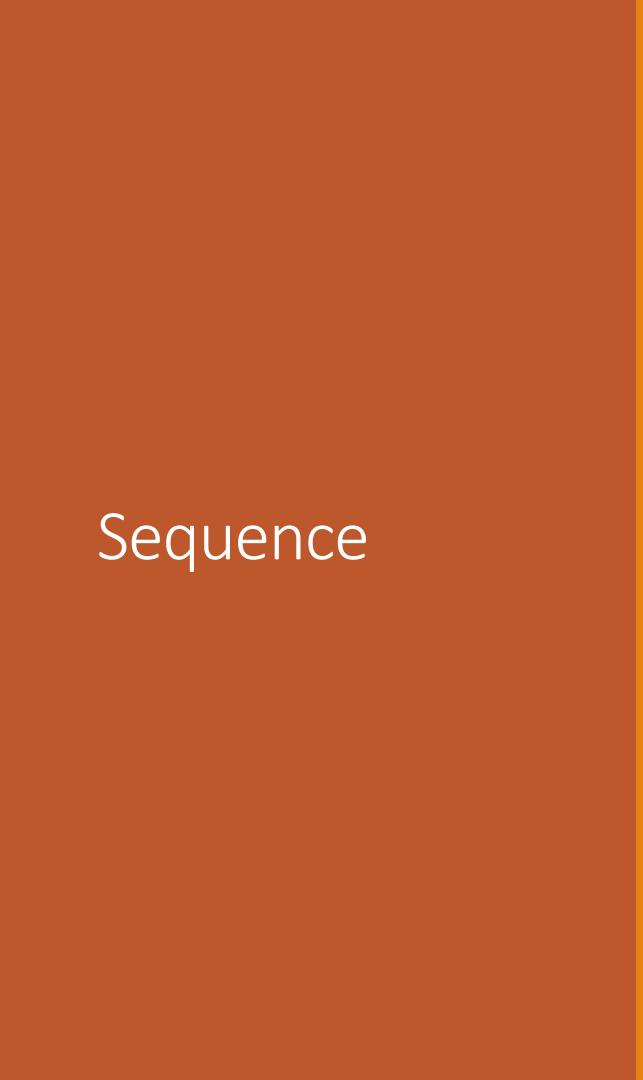
Issues in using Feed Forward Neural Network:

- Fixed window – Window is still small
- Increasing Window size will increase number of parameters to be estimated

We need a network architecture that can process a **sequence** of **any length**.



Recurrent Neural Networks

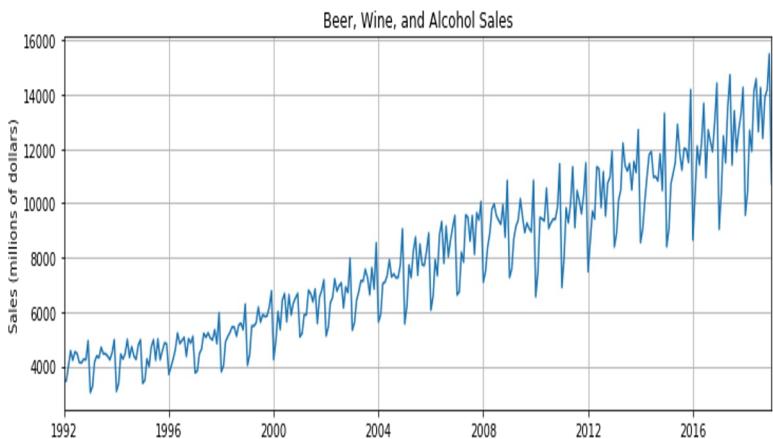


Sequence

We've used Neural Networks to solve Classification and Regression problems, but we still haven't seen how Neural Networks can deal with sequence information.

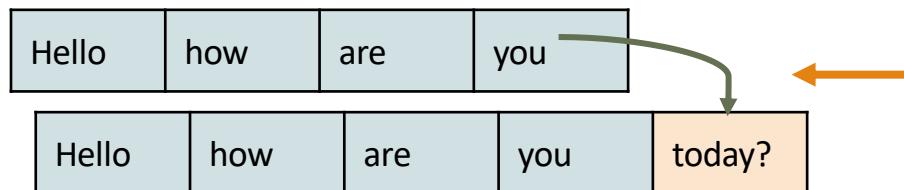
Sequences

Time Series



Examples of Sequences

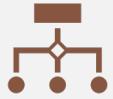
- Time Series Data (Sales)
- Sentences
- Audio
- Car Trajectories
- Music
- Patient Records



Language Model



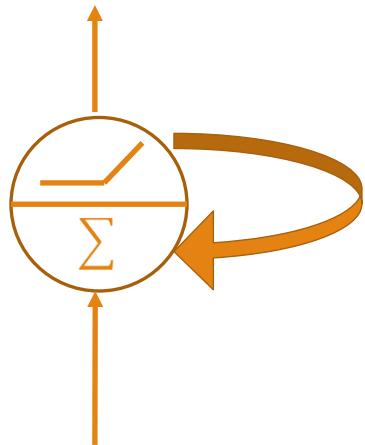
To process sequence, we need to somehow let the neuron “know” about its previous history of outputs.



One easy way to do this is to simply feed its output back into itself as an input!

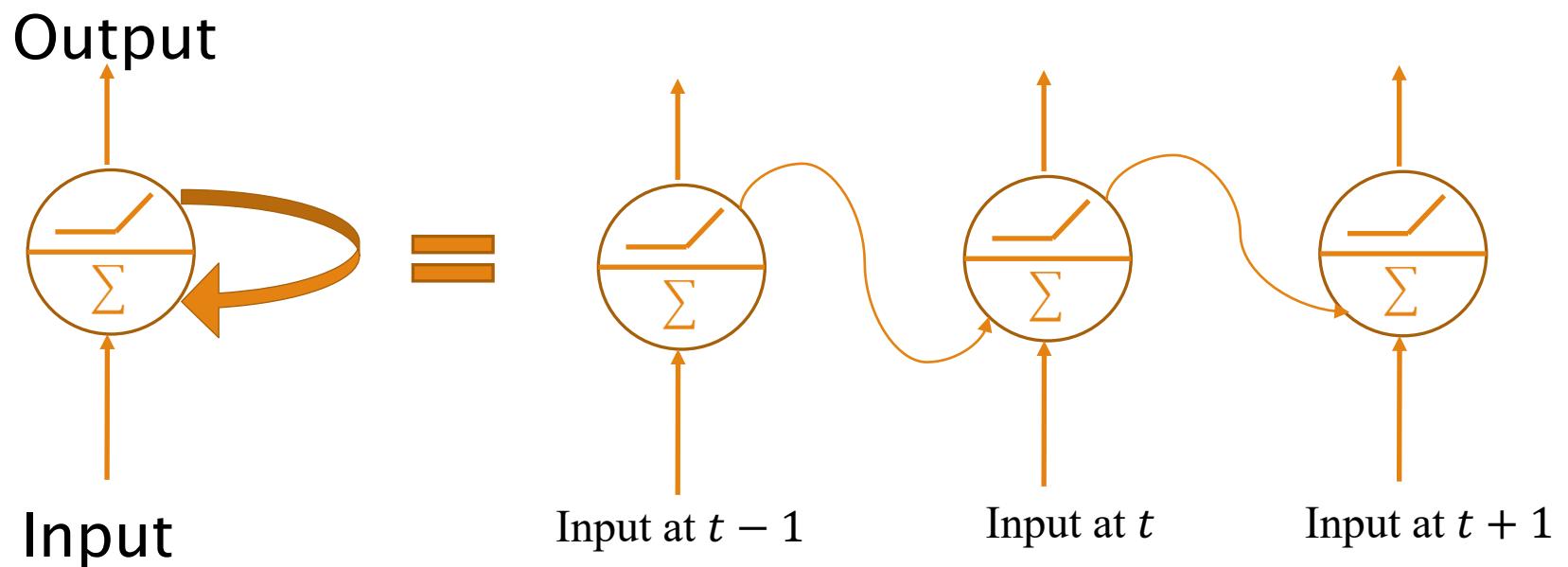
RNN - Intuition

Recurrent Neuron

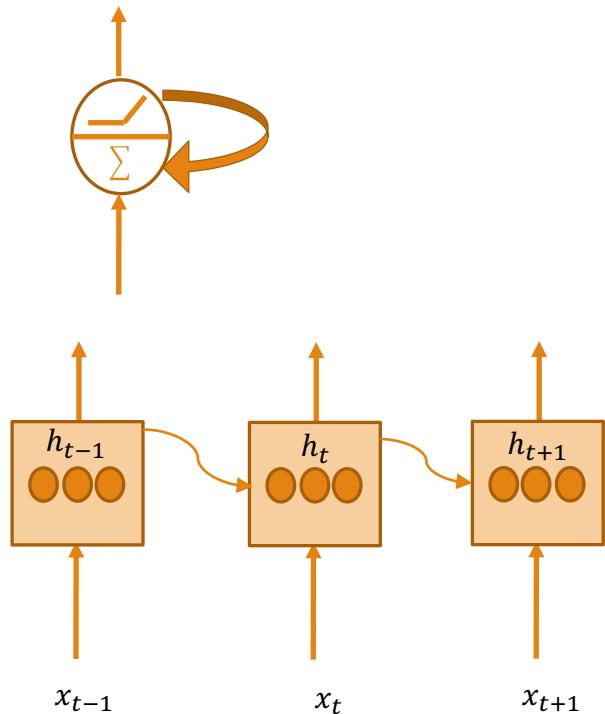


- RNNs have an '*hidden state*' that is updated as the sequence is processed

RNN Intuition



RNN Intuition



We can process a sequence of vectors \mathbf{x} by applying a **recurrence formula** at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

new state / old state input vector at
some function some time step
with parameters W

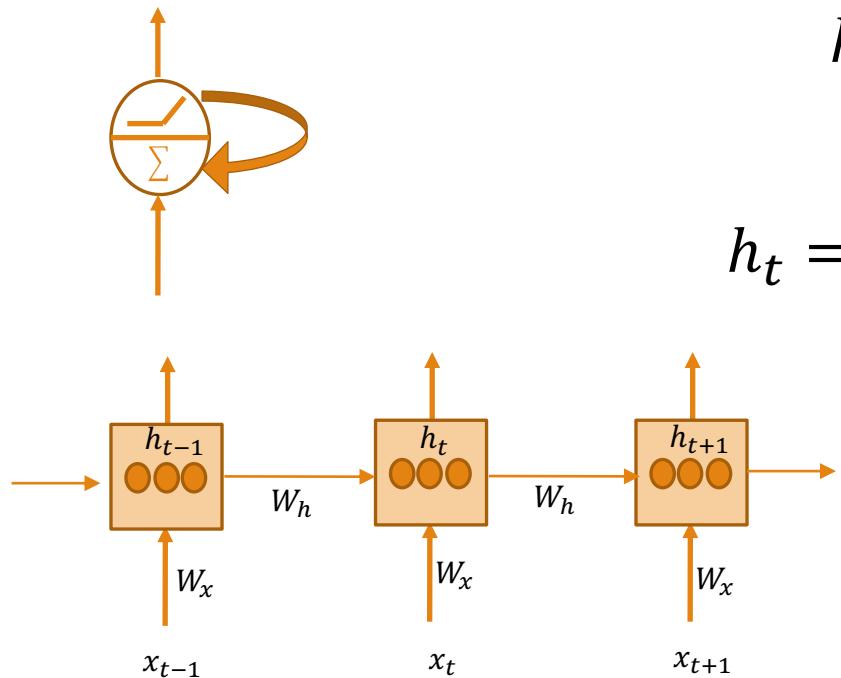
SLIDE CREDIT: JUSTIN JOHNSON (UNIVERSITY OF MICHIGAN)

RNN Intuition

The state consists of a “hidden” vector h :

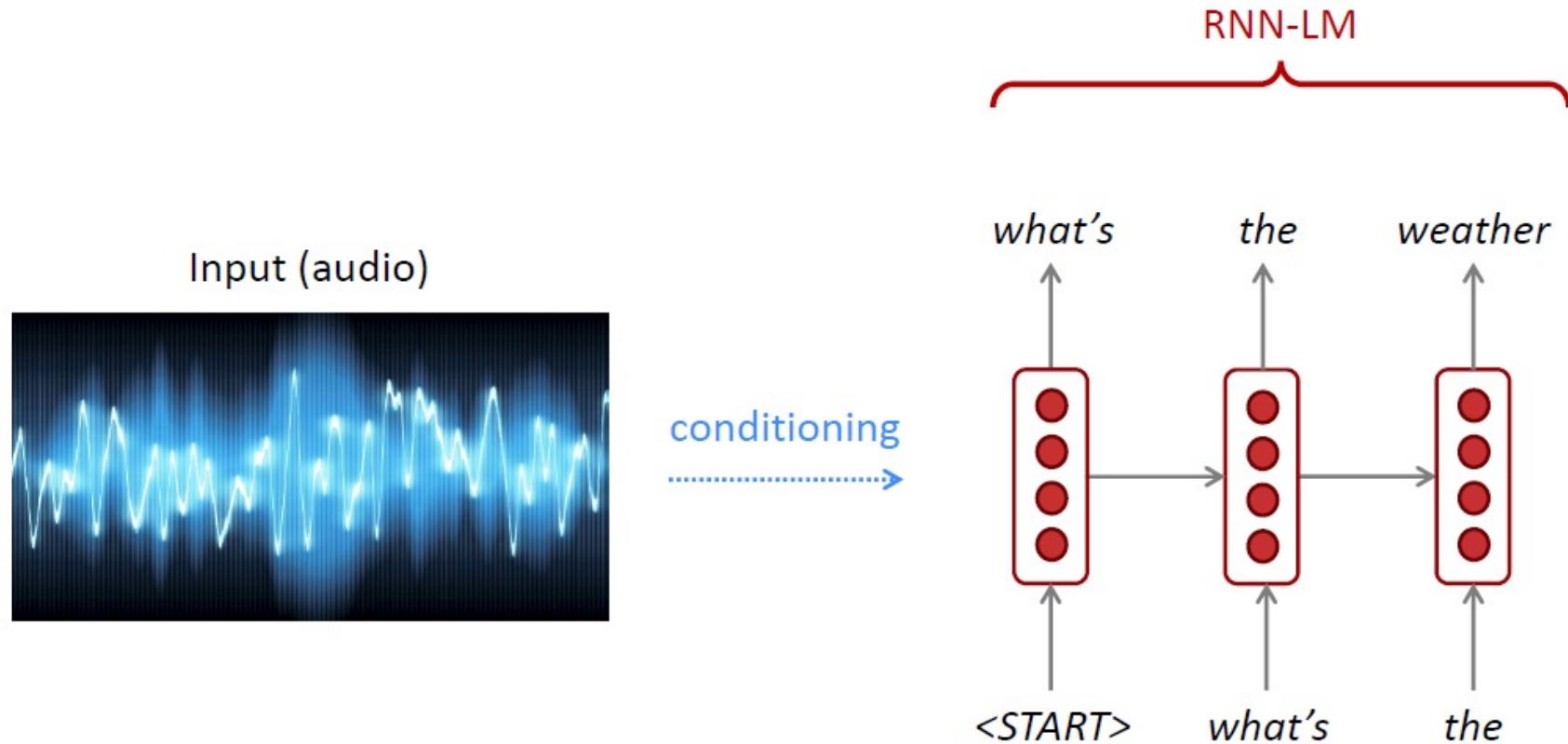
$$h_t = f_W(h_{t-1}, x_t)$$

$$h_t = \tanh(W_h h_{t-1} + b_h + W_x x_t + b_x)$$



RNN-LMs can be used to generate text

e.g. speech recognition, machine translation, summarization





Language Modelling with RNN

Language Model with RNN

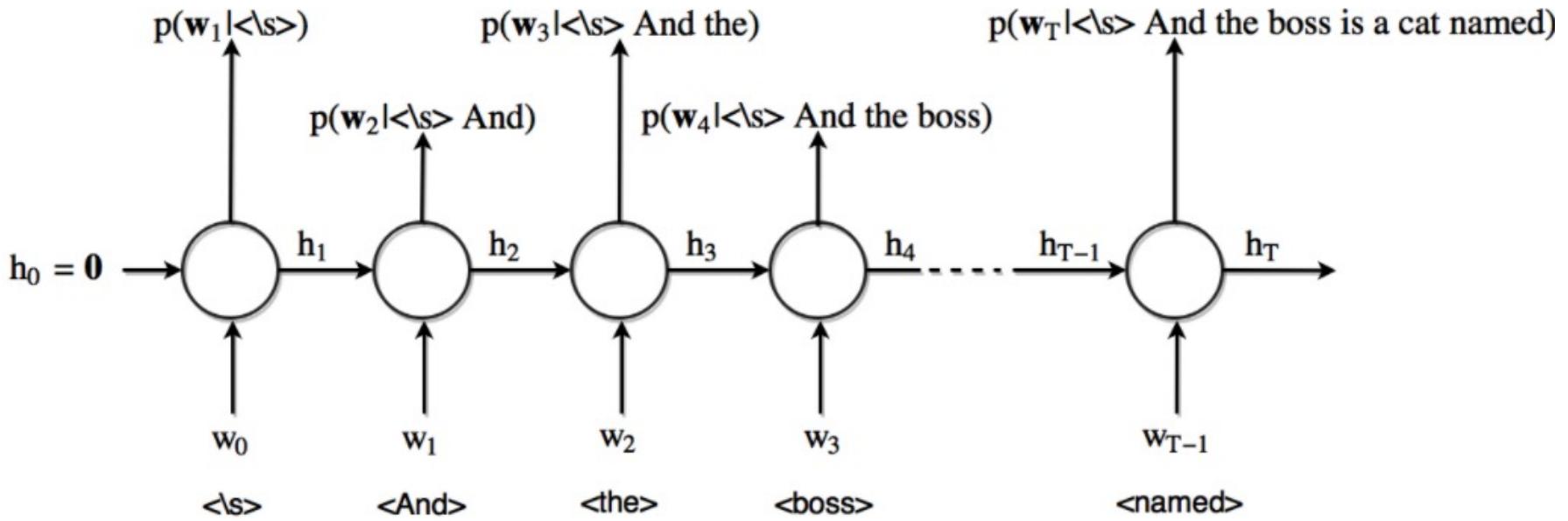


Figure Source <https://blog.paperspace.com/recurrent-neural-networks-part-1-2/>

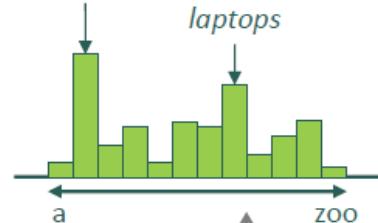
A Simple RNN Language Model

output distribution

$$\hat{y}^{(t)} = \text{softmax}(\mathbf{U}\mathbf{h}^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

$$\hat{y}^{(4)} = P(x^{(5)} | \text{the students opened their})$$

books



hidden states

$$\mathbf{h}^{(t)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1)$$

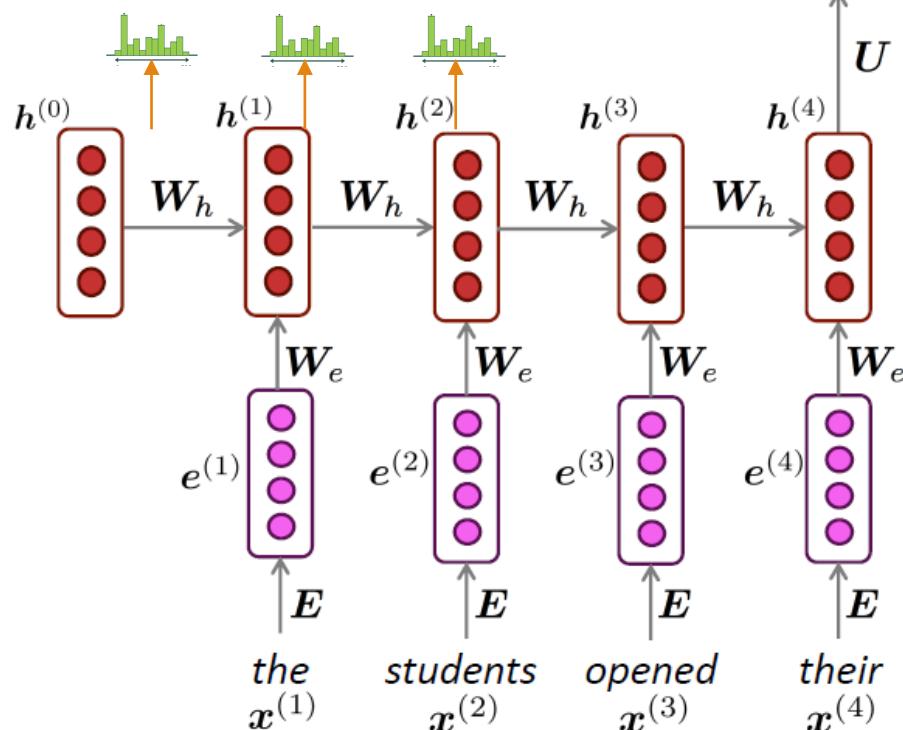
$\mathbf{h}^{(0)}$ is the initial hidden state

word embeddings

$$\mathbf{e}^{(t)} = \mathbf{E} \mathbf{x}^{(t)}$$

words / one-hot vectors

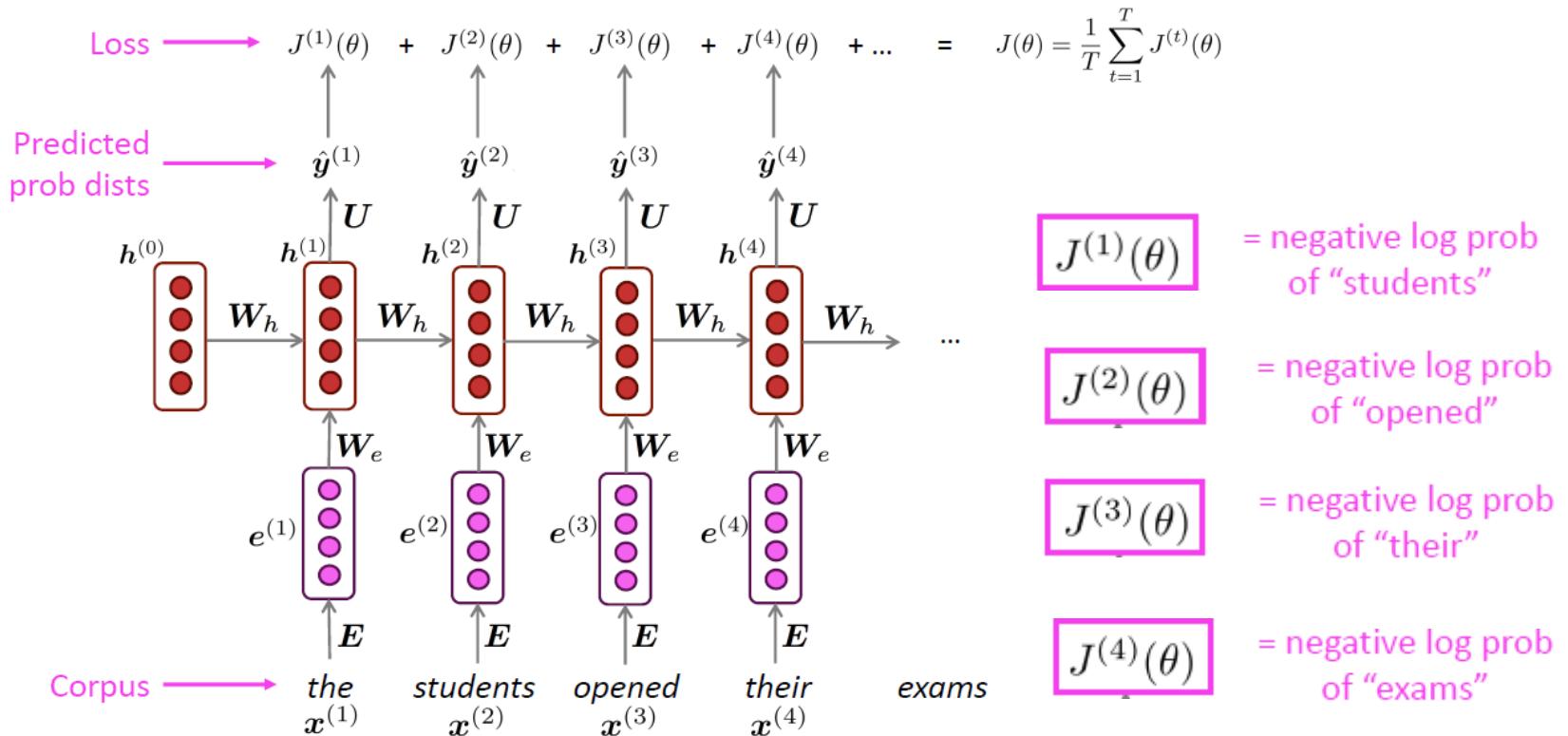
$$\mathbf{x}^{(t)} \in \mathbb{R}^{|V|}$$



Note: this input sequence could be much longer now!

Training an RNN Language Model

“Teacher forcing”



Training a RNN Language Model

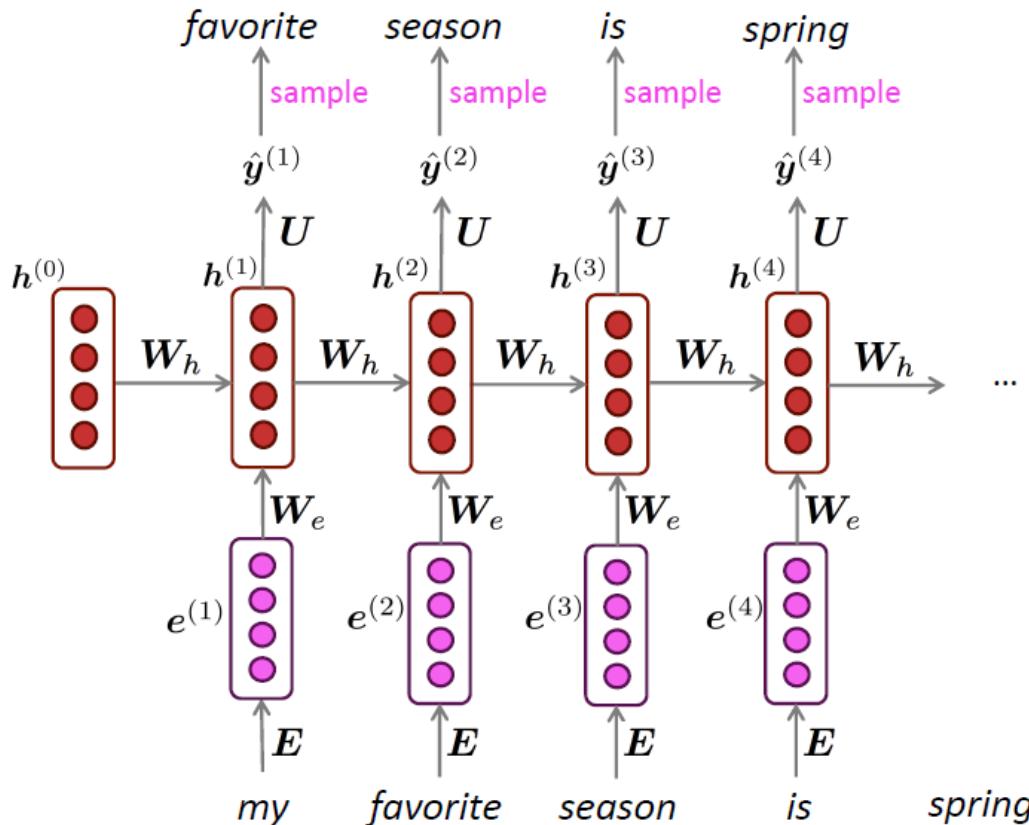
- However: Computing loss and gradients across entire corpus $x^{(1)}, \dots, x^{(T)}$ is too expensive!

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta)$$

- In practice, consider $x^{(1)}, \dots, x^{(T)}$ as a sentence (or a document)
- Recall: Stochastic Gradient Descent allows us to compute loss and gradients for small chunk of data, and update.
- Compute loss $J(\theta)$ for a sentence (actually, a batch of sentences), compute gradients and update weights. Repeat.

Generating text with a RNN Language Model

Just like a n-gram Language Model, you can use a RNN Language Model to generate text by **repeated sampling**. Sampled output becomes next step's input.



RNN Language Models

RNN Advantages:

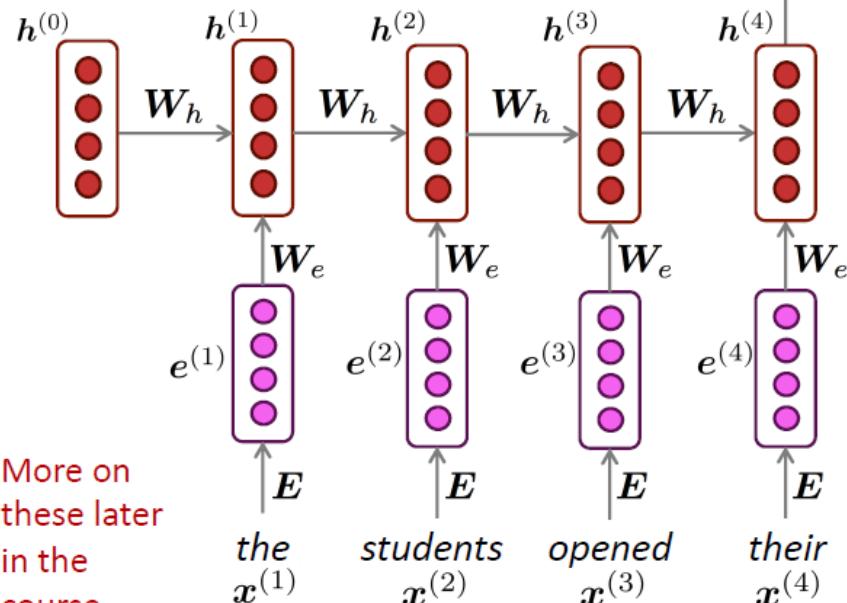
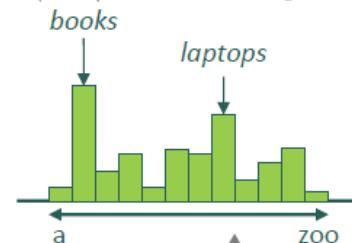
- Can process **any length** input
- Computation for step t can (in theory) use information from **many steps back**
- **Model size doesn't increase** for longer input context
- Same weights applied on every timestep, so there is **symmetry** in how inputs are processed.

RNN Disadvantages:

- Recurrent computation is **slow**
- In practice, difficult to access information from **many steps back**

More on
these later
in the
course

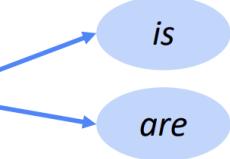
$$\hat{y}^{(4)} = P(\mathbf{x}^{(5)} | \text{the students opened their})$$



Slide credit: Abigail See and John Hewitt

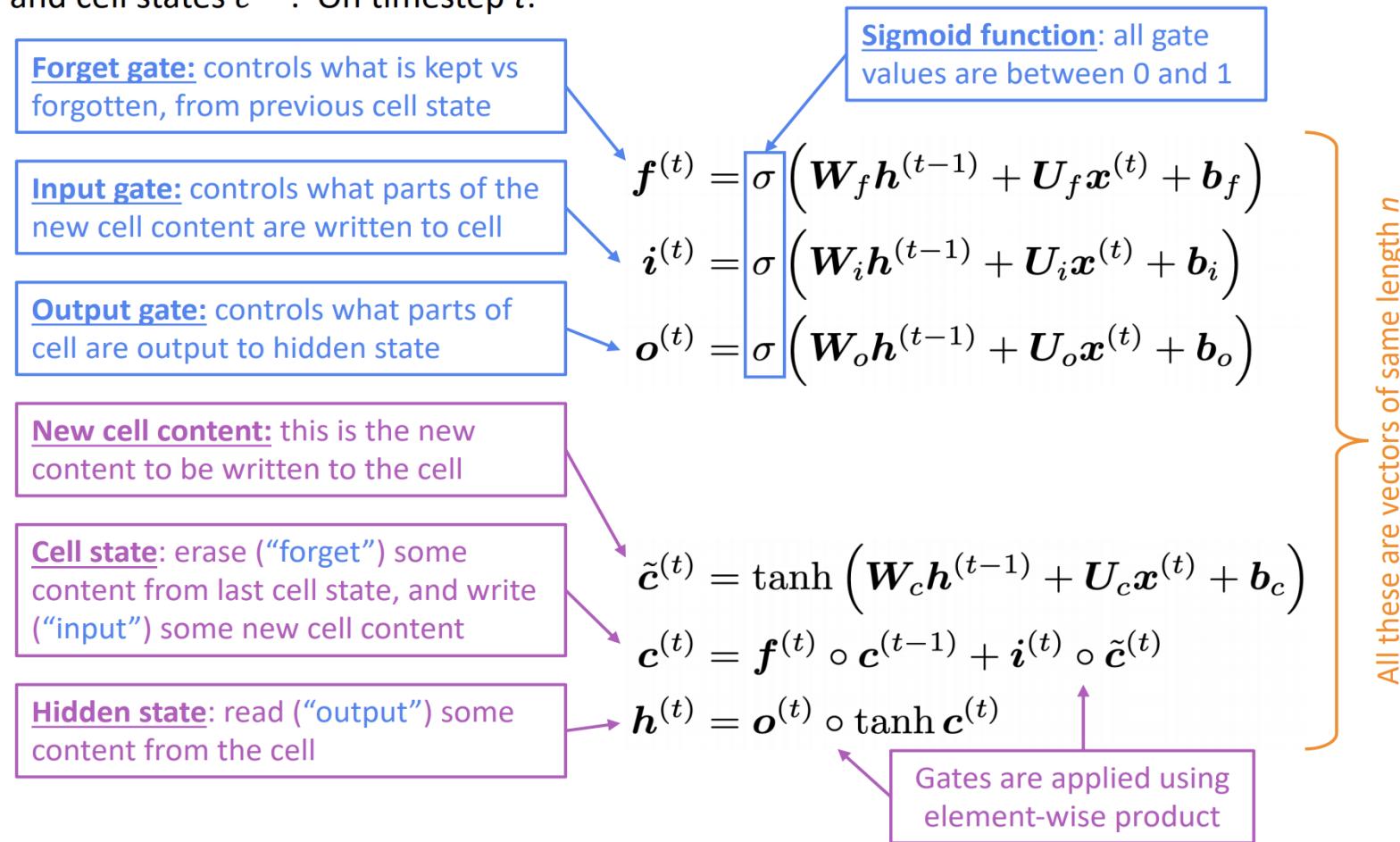
Effect of Vanishing Gradients

Effect of vanishing gradient on RNN-LM

- LM task: *The writer of the books __*
- Correct answer: *The writer of the books is planning a sequel*
- Syntactic recency: *The writer of the books is* (correct)
- Sequential recency: *The writer of the books are* (incorrect)
- Vanishing gradient problems may bias RNN-LMs towards learning from sequential recency, so they make this type of error more often than we'd like. [Linzen et al 2016]

Long Short-Term Memory (LSTM)

We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states $c^{(t)}$. On timestep t :



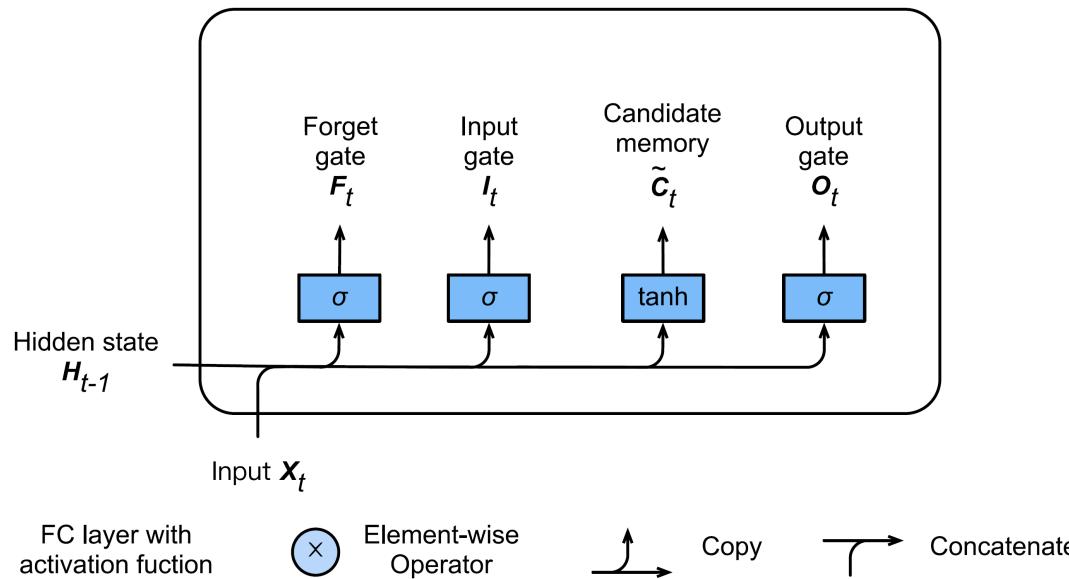
LSTM

Gates

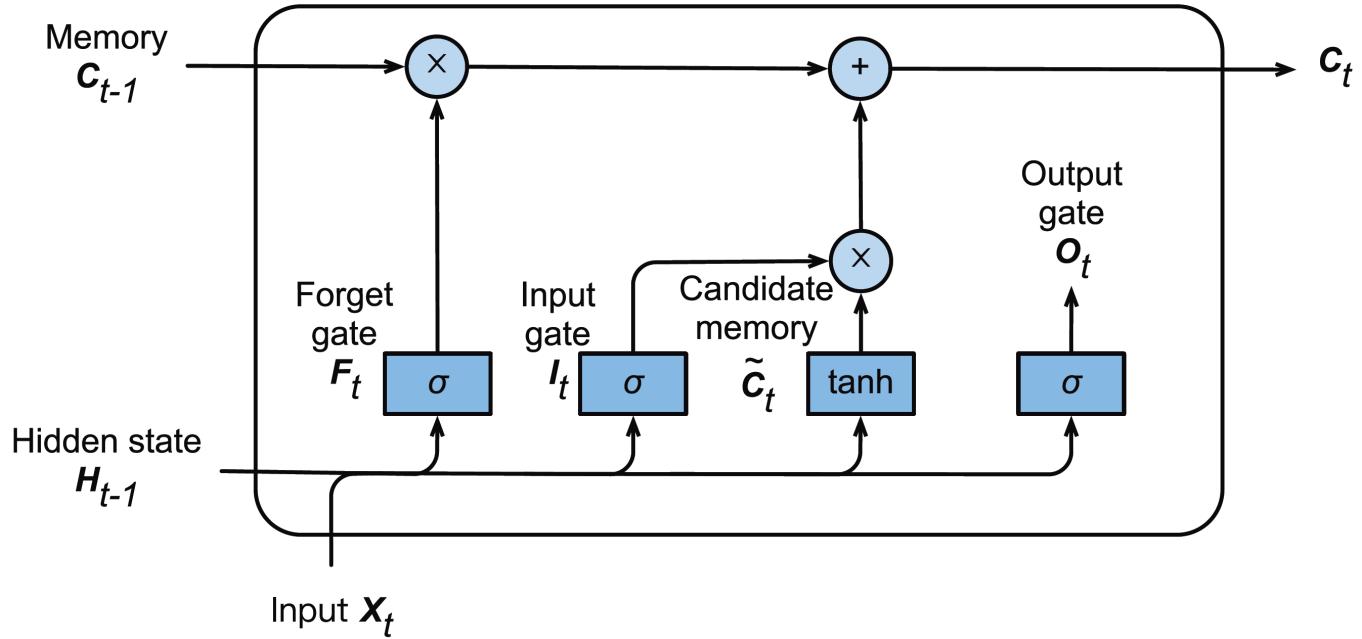
$$I_t = \sigma(X_t W_{xi} + H_{t-1} W_{hi} + b_i)$$

$$F_t = \sigma(X_t W_{xf} + H_{t-1} W_{hf} + b_f)$$

$$O_t = \sigma(X_t W_{xo} + H_{t-1} W_{ho} + b_o)$$

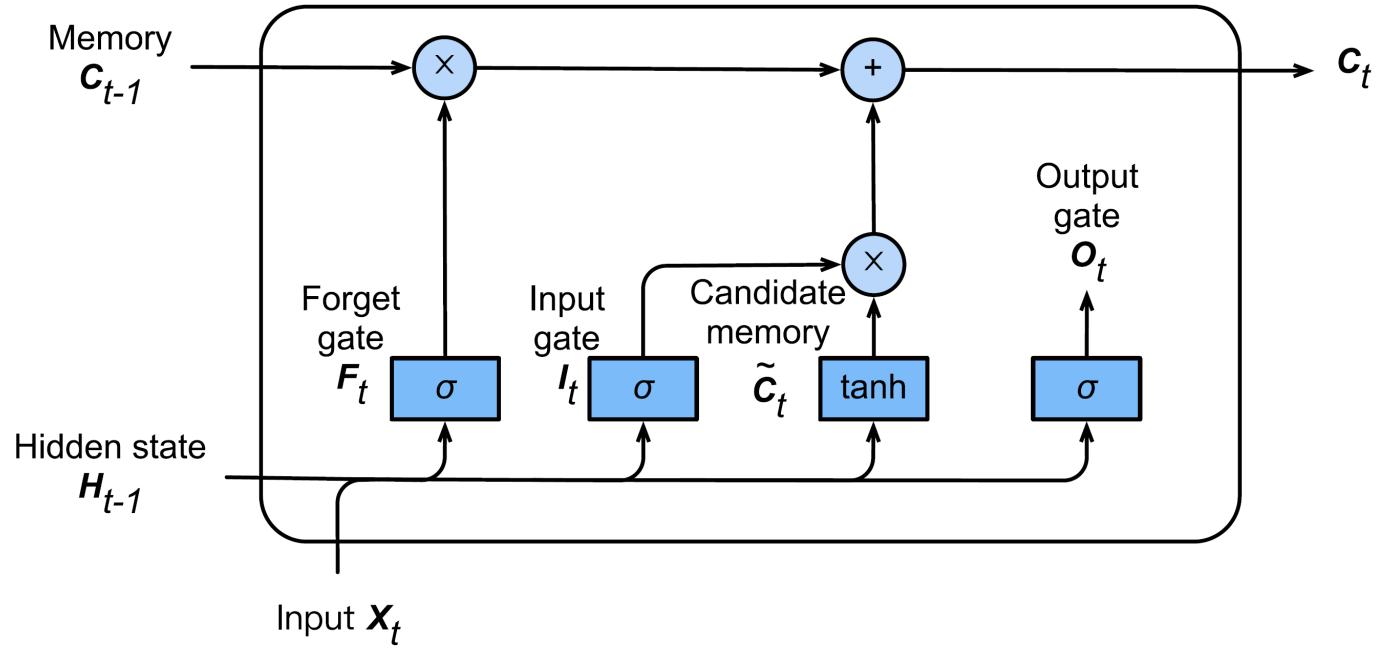


$$\tilde{C}_t = \tanh(X_t W_{xc} + H_{t-1} W_{hc} + b_c)$$



Slide credit: Alex Smola and Mu Li (Berkeley)

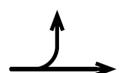
$$C_t = F_t \odot C_{t-1} + I_t \odot \tilde{C}_t$$



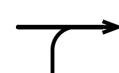
FC layer with
activation fuction



Element-wise
Operator



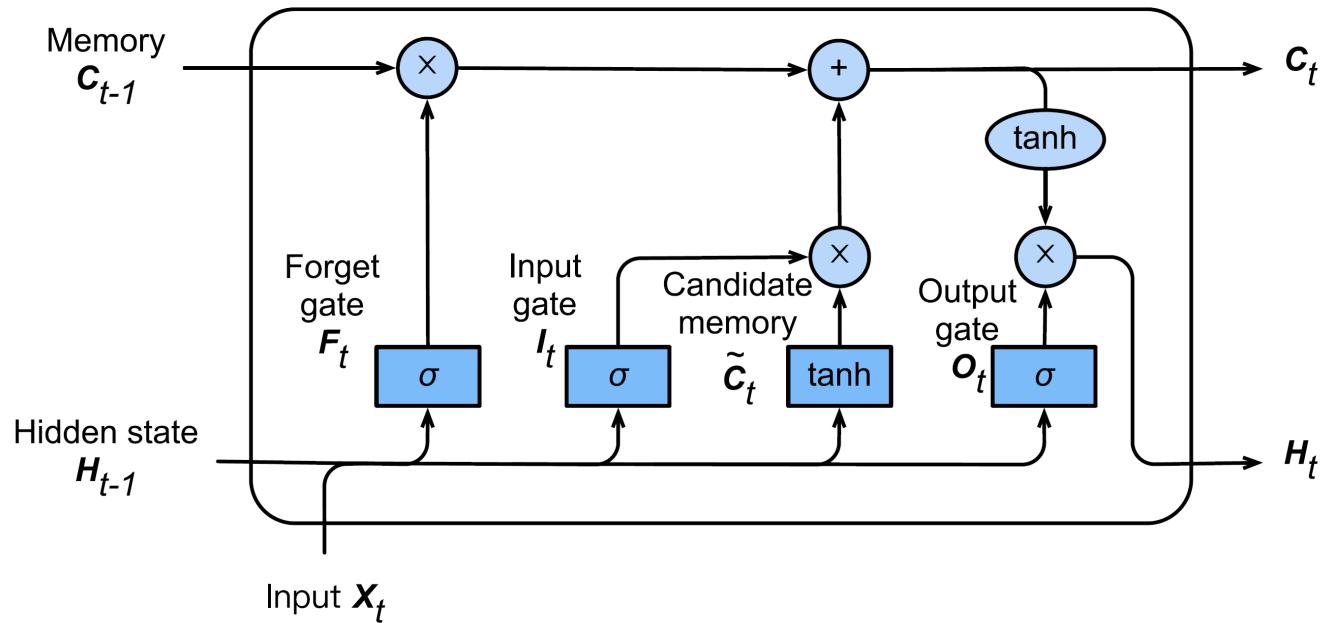
Copy



Concatenate

Hidden State / Output

$$H_t = O_t \odot \tanh(C_t)$$



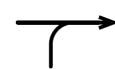
FC layer with
activation fuction



Element-wise
Operator

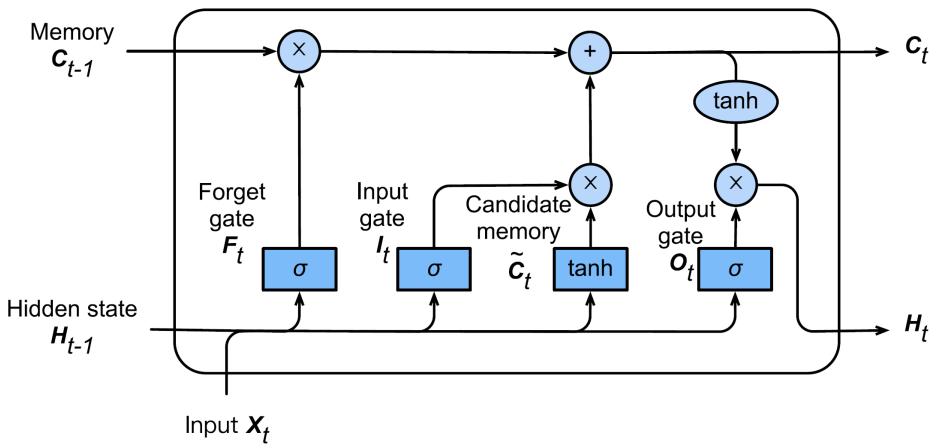


Copy



Concatenate

Hidden State / Output



$$I_t = \sigma(X_t W_{xi} + H_{t-1} W_{hi} + b_i)$$

$$F_t = \sigma(X_t W_{xf} + H_{t-1} W_{hf} + b_f)$$

$$O_t = \sigma(X_t W_{xo} + H_{t-1} W_{ho} + b_o)$$

$$\tilde{C}_t = \tanh(X_t W_{xc} + H_{t-1} W_{hc} + b_c)$$

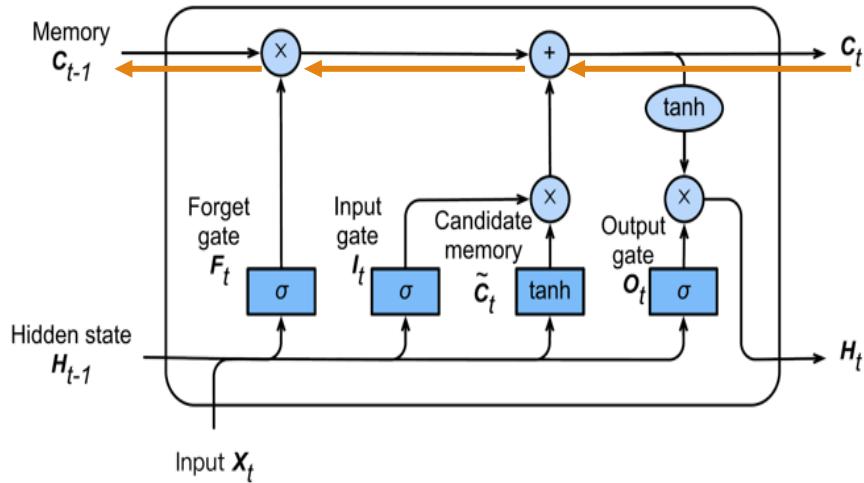
$$C_t = F_t \odot C_{t-1} + I_t \odot \tilde{C}_t$$

$$H_t = O_t \odot \tanh(C_t)$$

How does LSTM solve vanishing gradients?

- The LSTM architecture makes it easier for the RNN to preserve information over many timesteps
 - e.g. if the forget gate is set to 1 for a cell dimension and the input gate set to 0, then the information of that cell is preserved indefinitely.
 - By contrast, it's harder for vanilla RNN to learn a recurrent weight matrix W_h that preserves info in hidden state
- LSTM doesn't *guarantee* that there is no vanishing/exploding gradient, but it does provide an easier way for the model to learn long-distance dependencies

How does LSTM solve vanishing gradients?

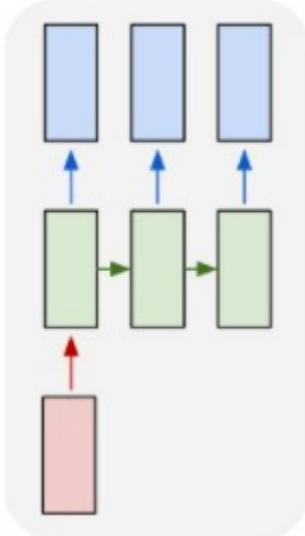


Back propagation from c_t to c_{t-1} : only elementwise multiplication with f, no matrix multiplication by W

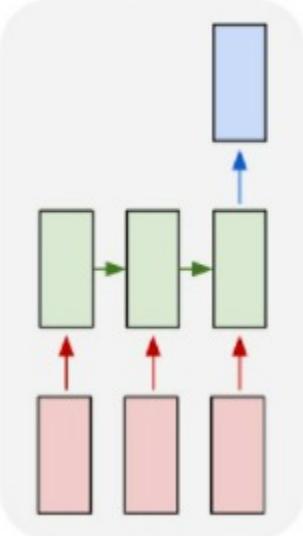
$$C_t = F_t \odot C_{t-1} + I_t \odot \tilde{C}_t$$

RNN/LSTM Applications

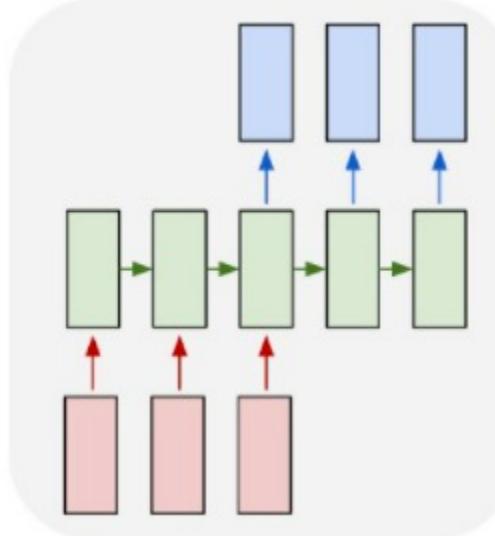
one to many



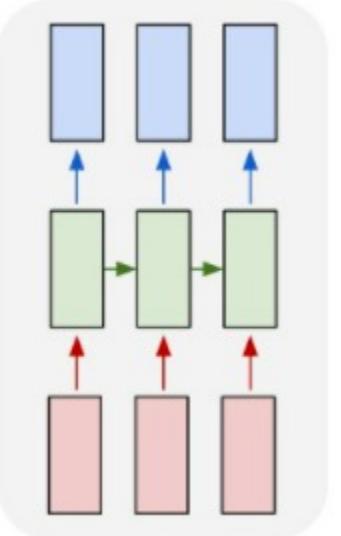
many to one



many to many



many to many



LSTMs: real-world success

- In 2013–2015, LSTMs started achieving state-of-the-art results
 - Successful tasks include handwriting recognition, speech recognition, machine translation, parsing, and image captioning, as well as language models
 - LSTMs became the dominant approach for most NLP tasks
- Now (2021), other approaches (e.g., Transformers) have become dominant for many tasks
 - For example, in WMT (a Machine Translation conference + competition):
 - In WMT 2016, the summary report contains “RNN” 44 times
 - In WMT 2019: “RNN” 7 times, “Transformer” 105 times

+ .
o

What we Learnt

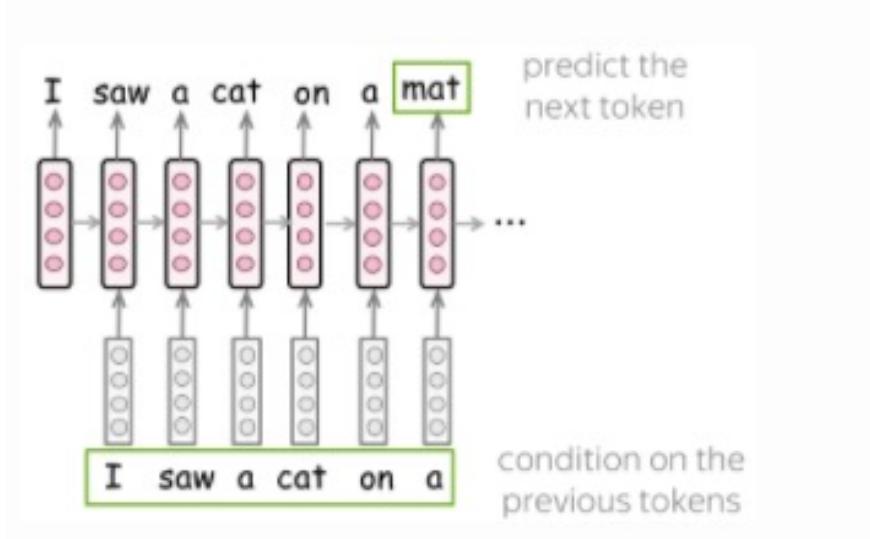
Context is key



- Previous models (GloVe, Word2Vec, etc.)
 - + only have one representation per word
 - They can't capture these ambiguities.
- We just have one representation for a word, but words have different **aspects**, including semantics, syntactic behavior, and register/connotations.
- When you only have one representation, all levels of meaning are combined.
- Solution: have multiple levels of understanding.
 - Embeddings from *Language Model* representations

Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

Language Model



- Those RNN layers are trained to predict the next word
- But those language models are producing context-specific word representations at each position

Image source https://lena-voita.github.io/nlp_course/seq2seq_and_attention.htm

Language Model – Context Specific Embeddings

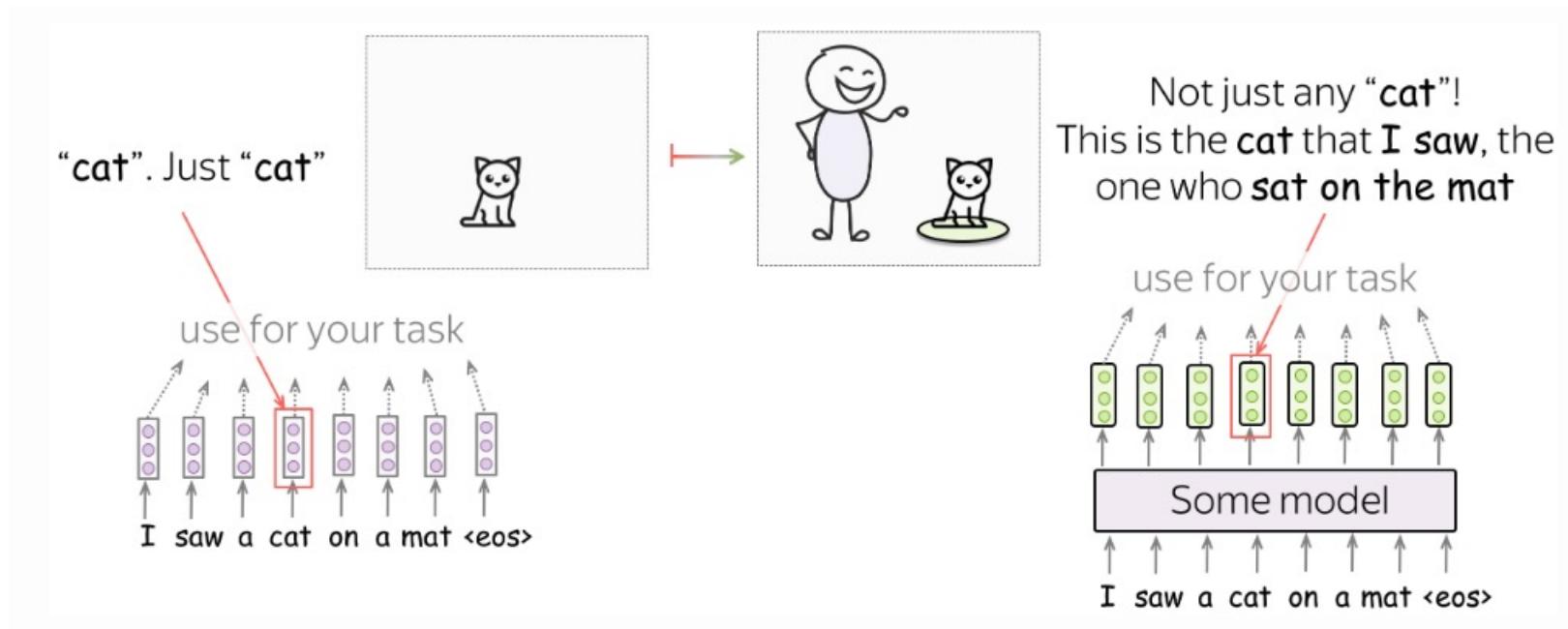
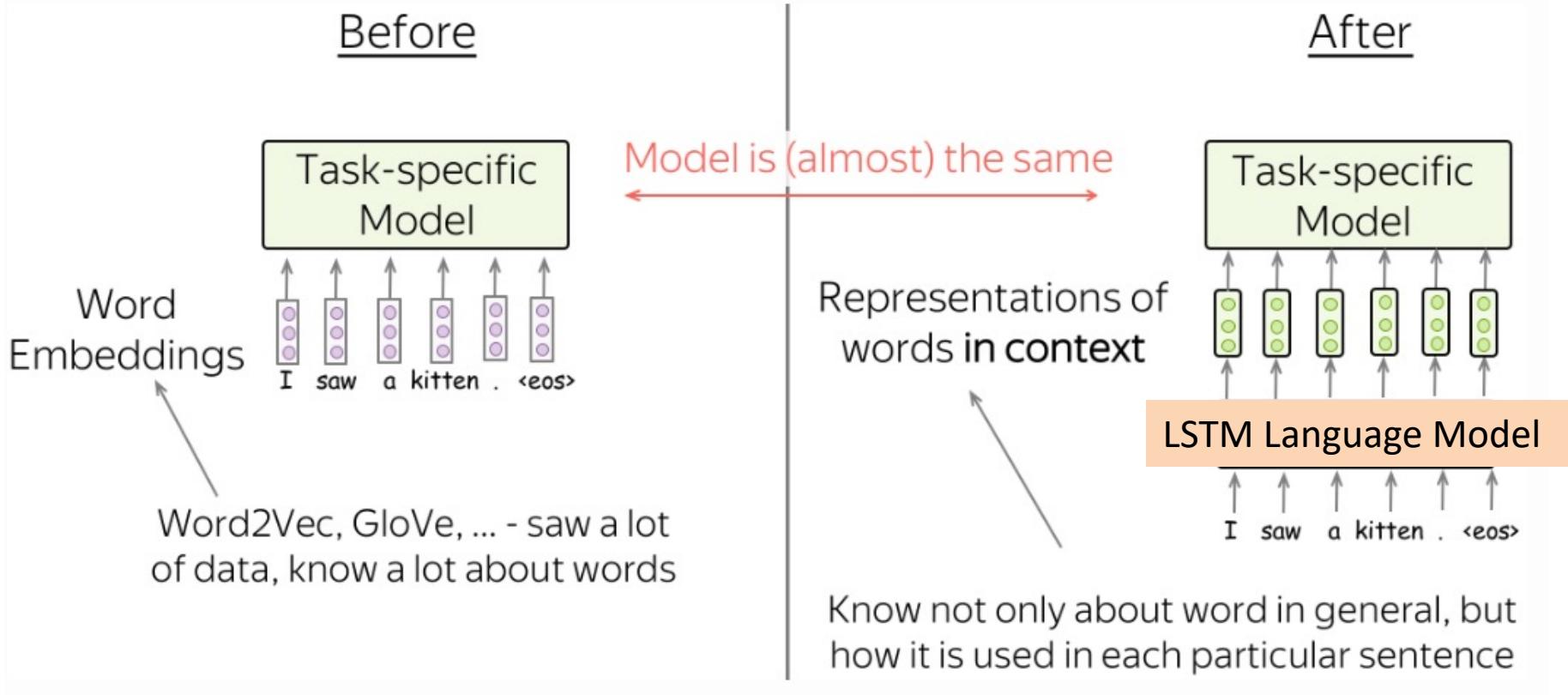


Image source https://lena-voita.github.io/nlp_course/seq2seq_and_attention.htm

Language Model – Context Specific Embeddings



Embeddings from Language/Translation Models

GloVe, Word2Vec → CoVe, ELMo

Great idea 1

What is encoded

words → words in context

How it is used for downstream tasks

Input for task-specific models

Input for task-specific models

Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

RNN/LSTM solves following issues

Polysemy: Both assign one vector per word, unable to differentiate meanings for words like "bank," which can signify a river's edge or a financial institution.

Static Representations: Word2Vec and GloVe produce fixed embeddings, not adapting to varying word meanings in different contexts.

Word Order Sensitivity: These models miss the word order, so phrases like "free drug" and "drug free" might get similar representations despite contrasting meanings.

Issues With RNN/LSTM

Recurrent computation is slow:

RNNs process data sequentially, which means each step must wait for the previous step to complete. This sequential dependency slows down the computation, particularly compared to models that can process data in parallel.

In practice, difficult to access information from many steps back:

RNNs are theoretically capable of handling long-range dependencies, but in practice, they often struggle to learn and retain information from many steps back in the input sequence. This is due to issues like vanishing gradients, where the contribution of information decays geometrically over time, making it hard for the RNN to maintain long-term dependencies.

Still needs Task Specific Models

ULMfit: A New Approach to Transfer Learning

Fine-tune the Language model for downstream tasks instead of training Task specific models

ULMfit

- **Pretraining:** Train LM on big general domain corpus
- **Domain Adaptation:** Tune LM on target task data (Task is same as a pre-training task)
- Fine-tune as classifier on target task

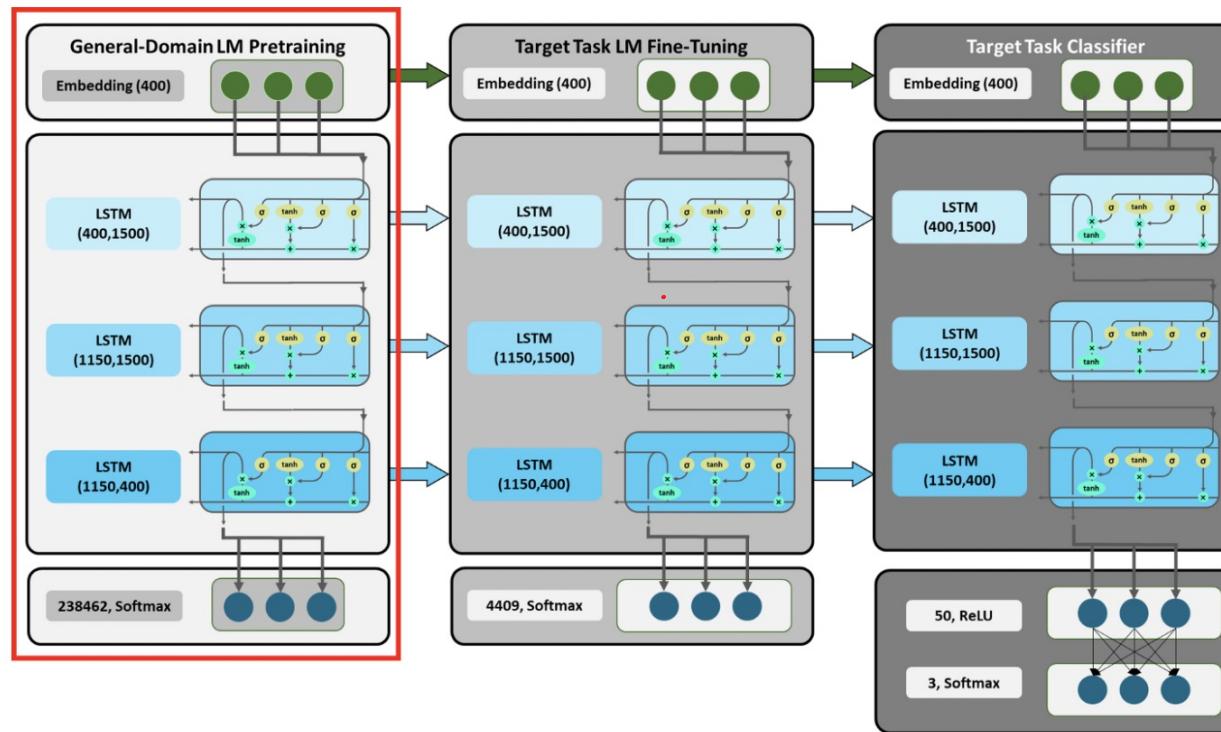
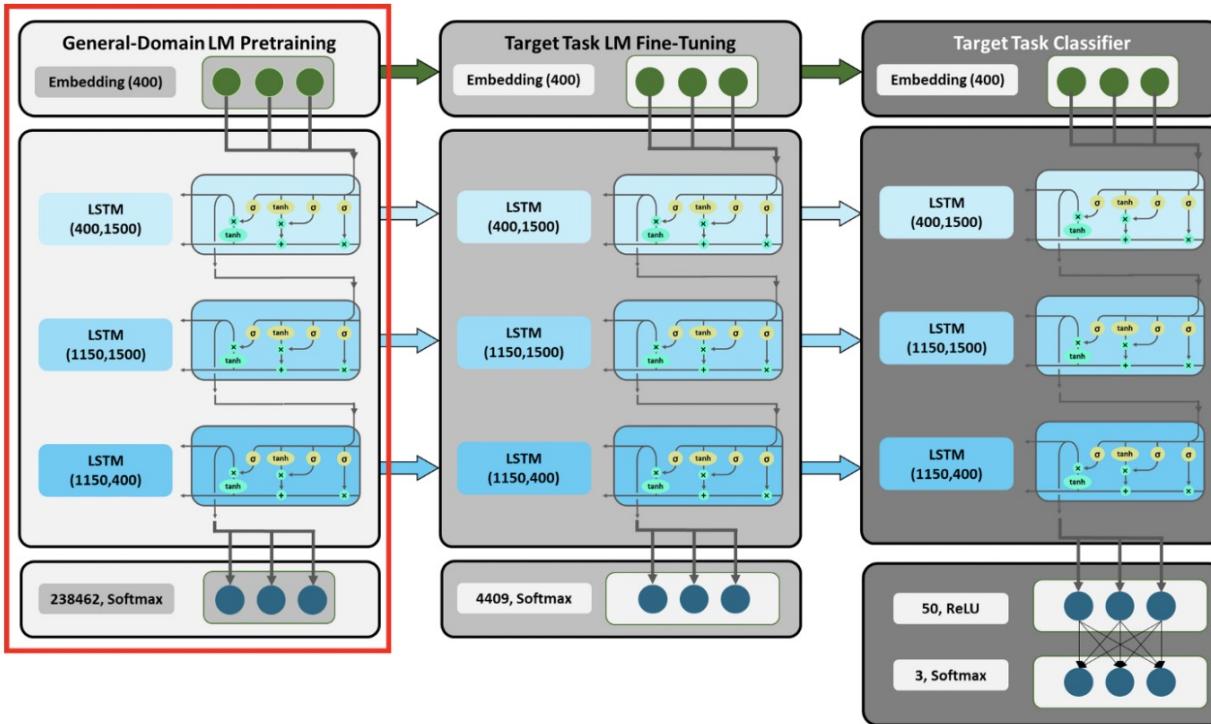


Image source: <https://towardsdatascience.com/understanding-ulmfit-and-elmo-the-shift-towards-transfer-learning-in-nlp-b5d8e2e3f664>

ULMfit emphases

- A lot of care in LM fine-tuning
 - Different per-layer learning rates
 - Slanted triangular learning rate (STLR) schedule
- Gradual layer unfreezing and STLR when learning classifier



ULMfit performance

Model	Test	Model	Test
CoVe (McCann et al., 2017)	8.2	TREC-6	CoVe (McCann et al., 2017) 4.2
IMDb oh-LSTM (Johnson and Zhang, 2016)	5.9		TBCNN (Mou et al., 2015) 4.0
Virtual (Miyato et al., 2016)	5.9		LSTM-CNN (Zhou et al., 2016) 3.9
ULMFiT (ours)	4.6		ULMFiT (ours) 3.6

ULMfit transfer learning

