

Natural Language Processing 1

Lecture 7: Word embeddings and sentence representations

Katia Shutova

ILLC
University of Amsterdam

19 November 2018

Distributional semantic models

1. Count-based models:

- ▶ Explicit vectors: dimensions are elements in the context
- ▶ **long sparse** vectors with **interpretable** dimensions

2. Prediction-based models:

- ▶ Train a model to predict plausible contexts for a word
- ▶ learn word representations in the process
- ▶ **short dense** vectors with **latent** dimensions

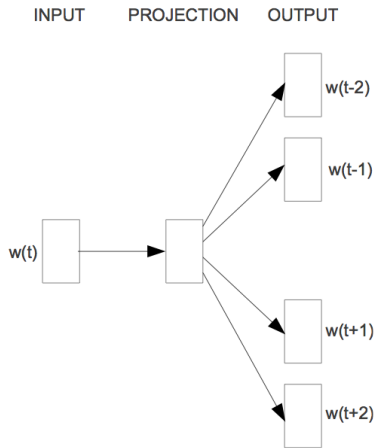
Prediction-based distributional models

Mikolov et. al. 2013. *Efficient Estimation of Word Representations in Vector Space*.

word2vec: **Skip-gram** model

- ▶ inspired by work on neural language models
- ▶ train a neural network to predict neighboring words
- ▶ learn dense embeddings for the words in the training corpus in the process

Skip-gram



Slide credit: Tomas Mikolov

Skip-gram

Intuition: words with similar meanings often occur near each other in texts

Given a word w_t :

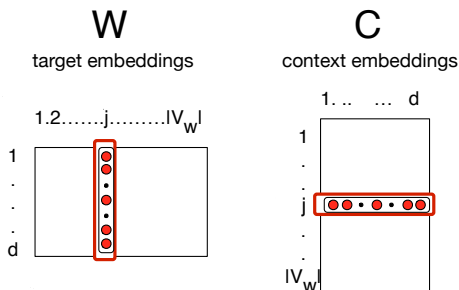
- ▶ Predict each neighbouring word
 - ▶ in a context window of $2L$ words
 - ▶ from the current word.
- ▶ For $L = 2$, we predict its 4 neighbouring words:

$$[w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2}]$$

Skip-gram: Parameter matrices

Learn 2 embeddings for each word $w_j \in V_w$:

- ▶ **word embedding** v , in word matrix W
- ▶ **context embedding** c , in context matrix C



Skip-gram: Setup

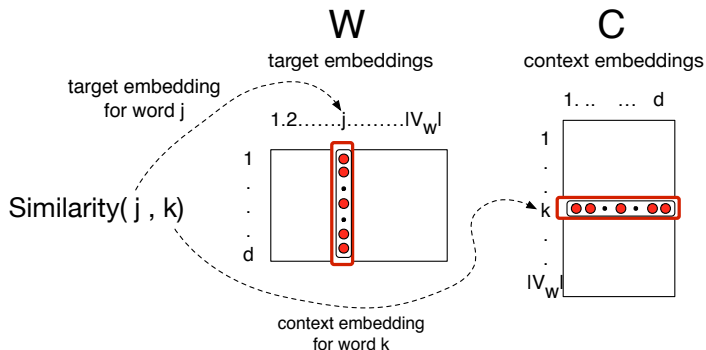
- ▶ Walk through the corpus pointing at word $w(t)$, whose index in the vocabulary is j — we will call it w_j
- ▶ our goal is to predict $w(t + 1)$, whose index in the vocabulary is k — we will call it w_k
- ▶ to do this, we need to compute

$$p(w_k | w_j)$$

- ▶ **Intuition** behind skip-gram: to compute this probability we need to compute similarity between w_j and w_k

Skip-gram: Computing similarity

Similarity as dot-product between the target vector and context vector



Slide credit: Dan Jurafsky

Skip-gram: Similarity as dot product

- ▶ Remember cosine similarity?

$$\cos(v_1, v_2) = \frac{\sum v_{1k} * v_{2k}}{\sqrt{\sum v_{1k}^2} * \sqrt{\sum v_{2k}^2}} = \frac{v_1 \cdot v_2}{\|v_1\| \|v_2\|}$$

It's just a normalised dot product.

- ▶ Skip-gram: Similar vectors have a high dot product

$$\text{Similarity}(c_k, v_j) \propto c_k \cdot v_j$$

Skip-gram: Compute probabilities

- ▶ Compute similarity as a dot product

$$\textit{Similarity}(c_k, v_j) \propto c_k \cdot v_j$$

- ▶ Normalise to turn this into a probability
- ▶ by passing through a softmax function:

$$p(w_k | w_j) = \frac{e^{c_k \cdot v_j}}{\sum_{i \in V} e^{c_i \cdot v_j}}$$

Skip-gram: Learning

- ▶ Start with some initial embeddings (usually random)
- ▶ At training time, walk through the corpus
- ▶ iteratively make the embeddings for each word
 - ▶ more like the embeddings of its neighbors
 - ▶ less like the embeddings of other words.

Skip-gram: Objective

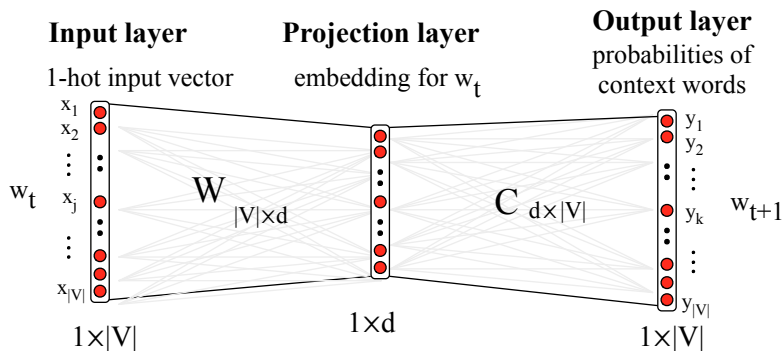
Learn parameters C and W that maximize the overall corpus probability:

$$\arg \max \prod_{(w_j, w_k) \in D} p(w_k | w_j)$$

$$p(w_k | w_j) = \frac{e^{c_k \cdot v_j}}{\sum_{i \in V} e^{c_i \cdot v_j}}$$

$$\arg \max \sum_{(w_j, w_k) \in D} \log p(w_k | w_j) = \sum_{(w_j, w_k) \in D} (\log e^{c_k \cdot v_j} - \log \sum_{c_i \in V} e^{c_i \cdot v_j})$$

Visualising skip-gram as a network



Slide credit: Dan Jurafsky

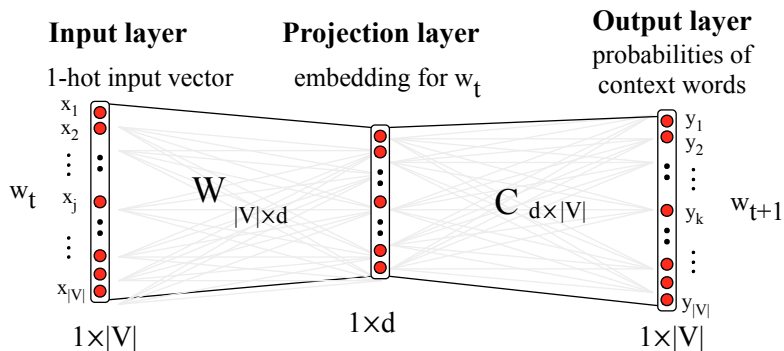
One hot vectors

- ▶ A vector of length $|V|$
- ▶ 1 for the target word and 0 for other words
- ▶ So if “bear” is vocabulary word 5
- ▶ The one-hot vector is $[0,0,0,0,1,0,0,0,0,\dots,0]$

w_0 w_1 w_j $w_{|V|}$

0 0 0 0 0 ... 0 0 0 0 1 0 0 0 0 0 ... 0 0 0 0

Visualising skip-gram as a network



Slide credit: Dan Jurafsky

Skip-gram with negative sampling

Problem with softmax: expensive to compute the denominator for the whole vocabulary

$$p(w_k | w_j) = \frac{e^{c_k \cdot v_j}}{\sum_{i \in V} e^{c_i \cdot v_j}}$$

Approximate the denominator: **negative sampling**

- ▶ At training time, walk through the corpus
- ▶ for each target word and positive context
- ▶ sample k noise samples or negative samples, i.e. other words

Skip-gram with negative sampling

- ▶ Objective in training:

- ▶ Make the word like the context words

lemon, a [tablespoon of apricot preserves or] jam.

c_1 c_2 w c_3 c_4

- ▶ And not like the k negative examples

[cement idle dear coaxial apricot attendant whence forever puddle]

n_1 n_2 n_3 n_4 w n_5 n_6 n_7 n_8

Skip-gram with negative sampling: Training examples

Convert the dataset into word pairs:

▶ **Positive (+)**

(apricot, tablespoon)

(apricot, of)

(apricot, jam)

(apricot, or)

▶ **Negative (-)**

(apricot, cement)

(apricot, idle)

(apricot, attendant)

(apricot, dear)

...

Skip-gram with negative sampling

- ▶ instead of treating it as a **multi-class problem** (and returning a probability distribution over the whole vocabulary)
- ▶ **return a probability** that word w_k is a valid context for word w_j

$$P(+|w_j, w_k)$$

$$P(-|w_j, w_k) = 1 - P(+|w_j, w_k)$$

Skip-gram with negative sampling

- ▶ model similarity as dot product

$$\text{Similarity}(c_k, v_j) \propto c_k \cdot v_j$$

- ▶ turn this into a probability using the **sigmoid function**:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$P(+|w_j, w_k) = \frac{1}{1 + e^{-c_k \cdot v_j}}$$

$$P(-|w_j, w_k) = 1 - P(+|w_j, w_k) = 1 - \frac{1}{1 + e^{-c_k \cdot v_j}} = \frac{1}{1 + e^{c_k \cdot v_j}}$$

Skip-gram with negative sampling

- ▶ model similarity as dot product

$$\textit{Similarity}(c_k, v_j) \propto c_k \cdot v_j$$

- ▶ turn this into a probability using the **sigmoid function**:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$P(+|w_j, w_k) = \frac{1}{1 + e^{-c_k \cdot v_j}}$$

$$P(-|w_j, w_k) = 1 - P(+|w_j, w_k) = 1 - \frac{1}{1 + e^{-c_k \cdot v_j}} = \frac{1}{1 + e^{c_k \cdot v_j}}$$

Skip-gram with negative sampling: Objective

- ▶ make the word like the context words
- ▶ and not like the negative examples

$$\arg \max \prod_{(w_j, w_k) \in D_+} p(+|w_k, w_j) \prod_{(w_j, w_k) \in D_-} p(-|w_k, w_j)$$

$$\arg \max \sum_{(w_j, w_k) \in D_+} \log p(+|w_k, w_j) + \sum_{(w_j, w_k) \in D_-} \log p(-|w_k, w_j) =$$

$$\arg \max \sum_{(w_j, w_k) \in D_+} \log \frac{1}{1 + e^{-c_k \cdot v_j}} + \sum_{(w_j, w_k) \in D_-} \log \frac{1}{1 + e^{c_k \cdot v_j}}$$

Skip-gram with negative sampling: Objective

- ▶ make the word like the context words
- ▶ and not like the negative examples

$$\arg \max \prod_{(w_j, w_k) \in D_+} p(+|w_k, w_j) \prod_{(w_j, w_k) \in D_-} p(-|w_k, w_j)$$

$$\arg \max \sum_{(w_j, w_k) \in D_+} \log p(+|w_k, w_j) + \sum_{(w_j, w_k) \in D_-} \log p(-|w_k, w_j) =$$

$$\arg \max \sum_{(w_j, w_k) \in D_+} \log \frac{1}{1 + e^{-c_k \cdot v_j}} + \sum_{(w_j, w_k) \in D_-} \log \frac{1}{1 + e^{c_k \cdot v_j}}$$

Skip-gram with negative sampling: Objective

- ▶ make the word like the context words
- ▶ and not like the negative examples

$$\arg \max \prod_{(w_j, w_k) \in D_+} p(+|w_k, w_j) \prod_{(w_j, w_k) \in D_-} p(-|w_k, w_j)$$

$$\arg \max \sum_{(w_j, w_k) \in D_+} \log p(+|w_k, w_j) + \sum_{(w_j, w_k) \in D_-} \log p(-|w_k, w_j) =$$

$$\arg \max \sum_{(w_j, w_k) \in D_+} \log \frac{1}{1 + e^{-c_k \cdot v_j}} + \sum_{(w_j, w_k) \in D_-} \log \frac{1}{1 + e^{c_k \cdot v_j}}$$

Properties of embeddings

They capture similarity

FRANCE	JESUS	XBOX	REDDISH	SCRATCHED	MEGABITS
454	1973	6909	11724	29869	87025
AUSTRIA	GOD	AMIGA	GREENISH	NAILED	OCTETS
BELGIUM	SATI	PLAYSTATION	BLUISH	SMASHED	MB/S
GERMANY	CHRIST	MSX	PINKISH	PUNCHED	BIT/S
ITALY	SATAN	IPOD	PURPLISH	POPPED	BAUD
GREECE	KALI	SEGA	BROWNISH	CRIMPED	CARATS
SWEDEN	INDRA	PSNUMBER	GREYISH	SCRAPED	KBIT/S
NORWAY	VISHNU	HD	GRAYISH	SCREWED	MEGAHERTZ
EUROPE	ANANDA	DREAMCAST	WHITISH	SECTIONED	MEGAPIXELS
HUNGARY	PARVATI	GEFORCE	SILVERY	SLASHED	GBIT/S
SWITZERLAND	GRACE	CAPCOM	YELLOWISH	RIPPED	AMPERES

Slide credit: Ronan Collobert

Properties of embeddings

They capture **analogy**

Analogy task: ***a** is to **b** as **c** is to **d***

The system is given words *a*, *b*, *c*, and it needs to find *d*.

“apple” is to “apples” as “car” is to ?

“man” is to “woman” as “king” is to ?

Solution: capture analogy via vector offsets

$$a - b \approx c - d$$

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

$$d_w = \operatorname{argmax}_{d'_w \in V} \cos(a - b, c - d')$$

Properties of embeddings

They capture **analogy**

Analogy task: ***a** is to **b** as **c** is to **d***

The system is given words *a*, *b*, *c*, and it needs to find *d*.

“apple” is to “apples” as “car” is to ?

“man” is to “woman” as “king” is to ?

Solution: capture analogy via vector offsets

$$a - b \approx c - d$$

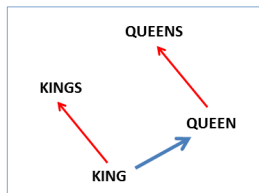
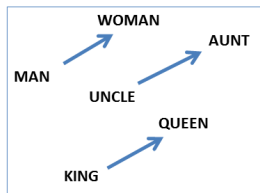
$$man - woman \approx king - queen$$

$$d_w = \operatorname{argmax}_{d'_w \in V} \cos(a - b, c - d')$$

Properties of embeddings

Capture analogy via vector offsets

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



Mikolov et al. 2013. *Linguistic Regularities in Continuous Space Word Representations*

Properties of embeddings

They capture a range of semantic relations

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

Mikolov et al. 2013. *Efficient Estimation of Word Representations in Vector Space*

Word embeddings in practice

Word2vec is often used for pretraining in other tasks.

- ▶ It will help your models start from an **informed** position
- ▶ Requires only **plain text** - which we have a lot of
- ▶ Is very **fast** and easy to use
- ▶ Already **pretrained** vectors also available (trained on 100B words)

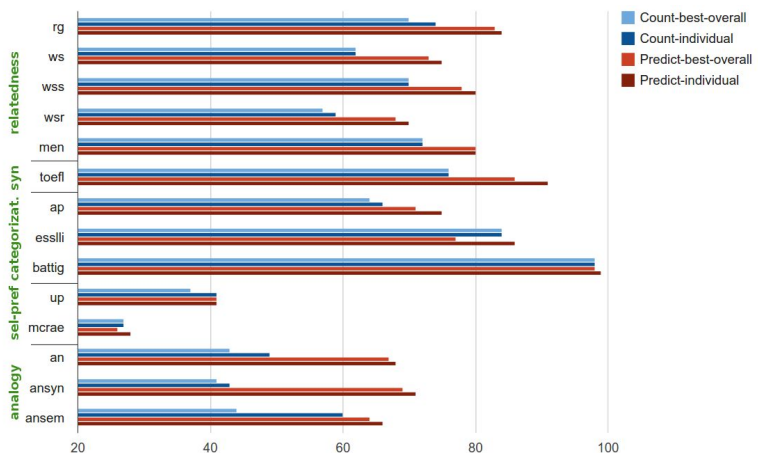
However, for best performance it is important to continue training, fine-tuning the embeddings for a specific task.

Count-based models vs. skip-gram word embeddings

Baroni et. al. 2014. *Don't count, predict! A systematic comparison of context-counting vs. context-predicting semantic vectors.*

- ▶ Comparison of count-based and neural word vectors on 5 types of tasks and 14 different datasets:
 1. Semantic relatedness
 2. Synonym detection
 3. Concept categorization
 4. Selectional preferences
 5. Analogy recovery

Count-based models vs. skip-gram word embeddings



Some of these findings were later disputed by Levy et. al. 2015. *Improving Distributional Similarity with Lessons Learned from Word Embeddings*

Acknowledgement

Some slides were adapted from Dan Jurafsky



Encoding Sentences with **Recurrent** and **Tree Recursive** Neural Networks

Joost Bastings

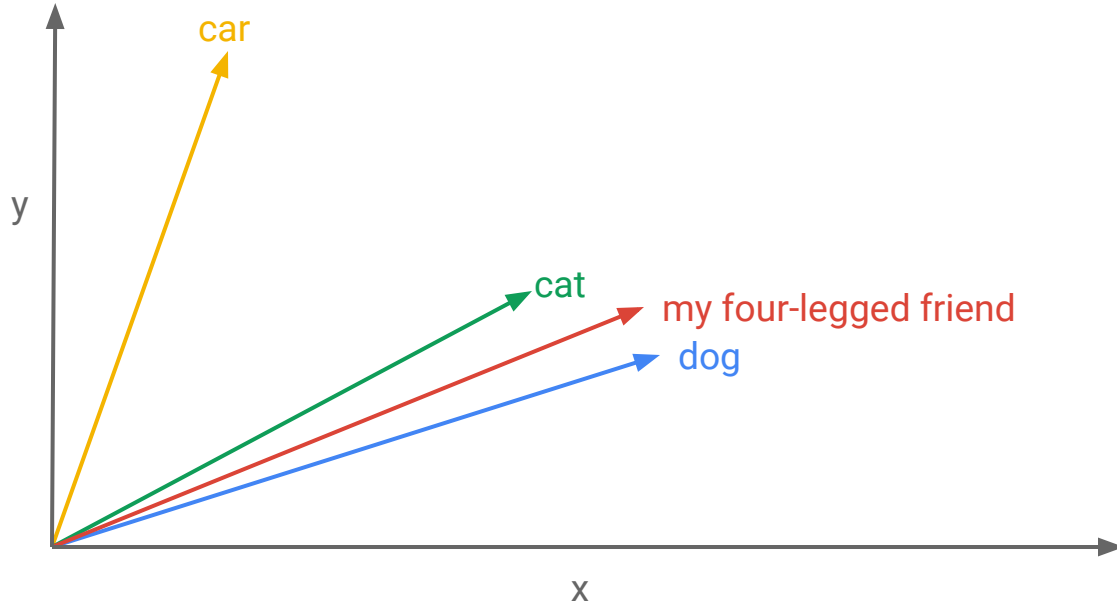
bastings.github.io

Today

How do we learn a **representation** of a **sentence** with a **neural network**?

How do we make a **prediction** from that representation, e.g. **sentiment**?

A vector space of words and sentences



Turning words into numbers

We want to **feed words** to a neural network
How to turn **words** into **numbers**?

Bad idea: number sequence

```
cat 1  
tree 2  
chair 3  
dog 4  
mat 5
```

cat is closer to **tree**
than to **dog**?!

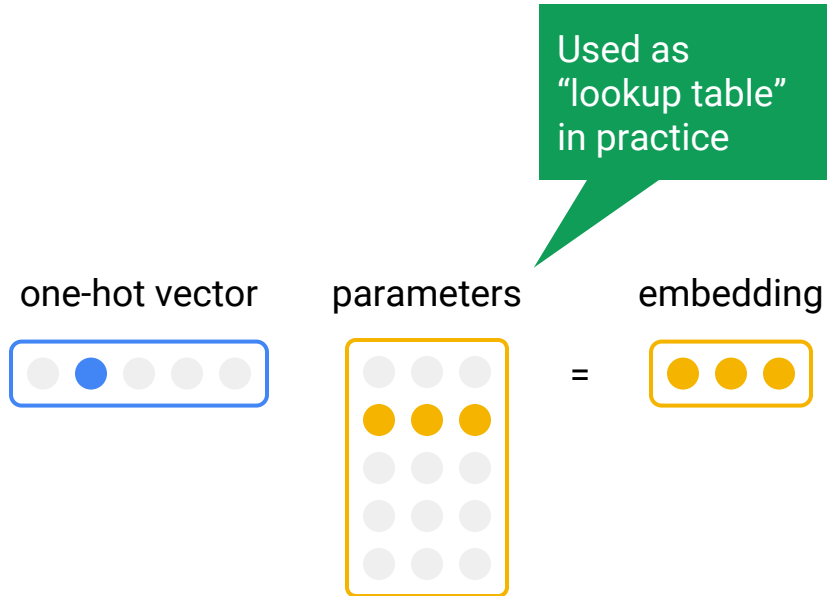


Good idea: one-hot vectors

```
cat [0, 0, 0, 0, 1]  
tree [0, 0, 0, 1, 0]  
chair [0, 0, 1, 0, 0]  
dog [0, 1, 0, 0, 0]  
mat [1, 0, 0, 0, 0]
```



One-hot vectors select word embeddings



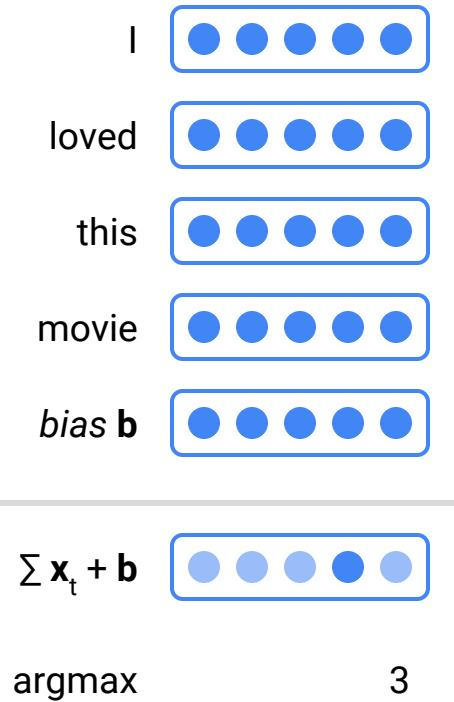
Bag of Words

Bag of Words at CMU



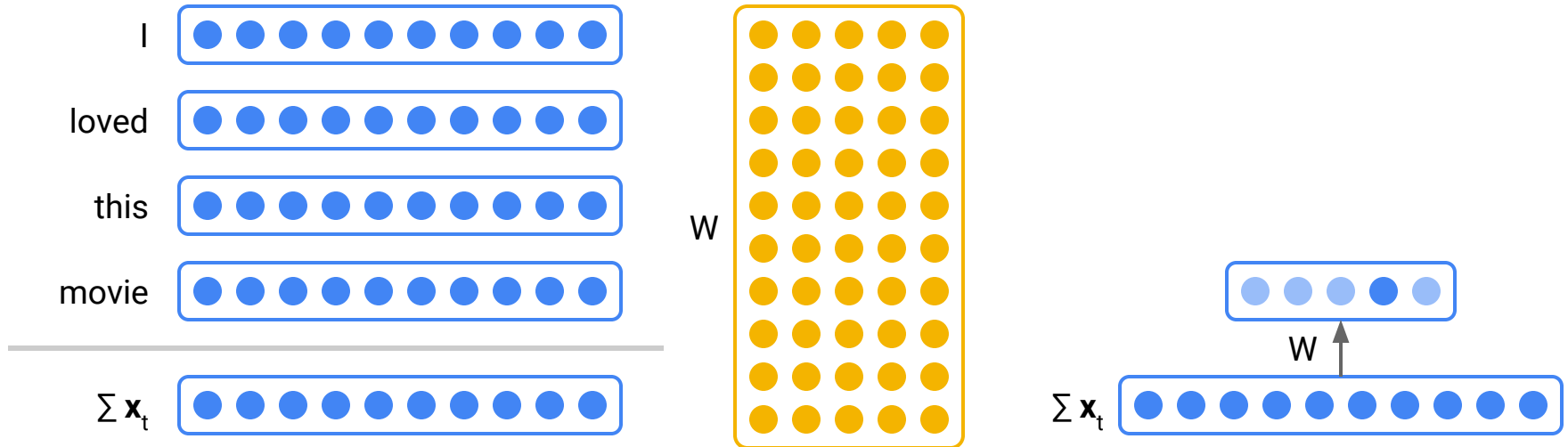
Bag of Words

Sum word embeddings, add bias

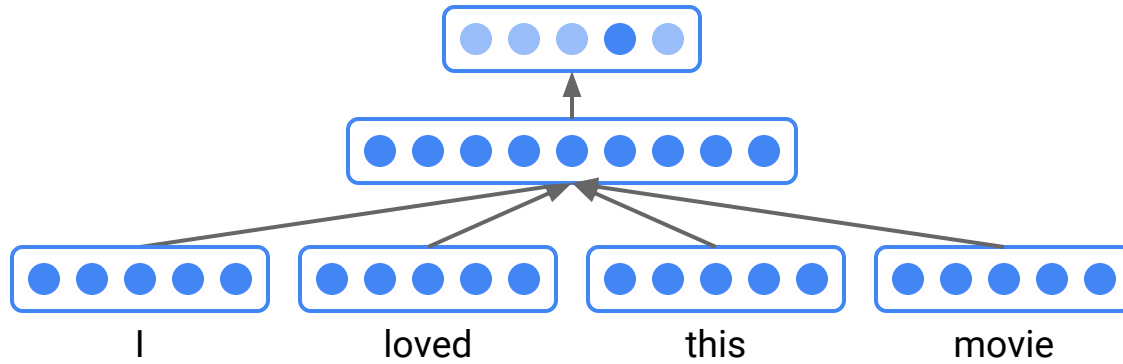


Continuous Bag of Words (CBOW)

Sum word embeddings, project to 5D using W , add bias: $W (\sum \mathbf{x}_t) + \mathbf{b}$

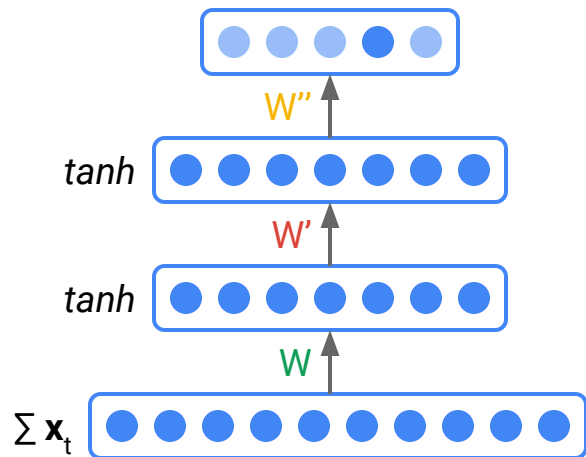
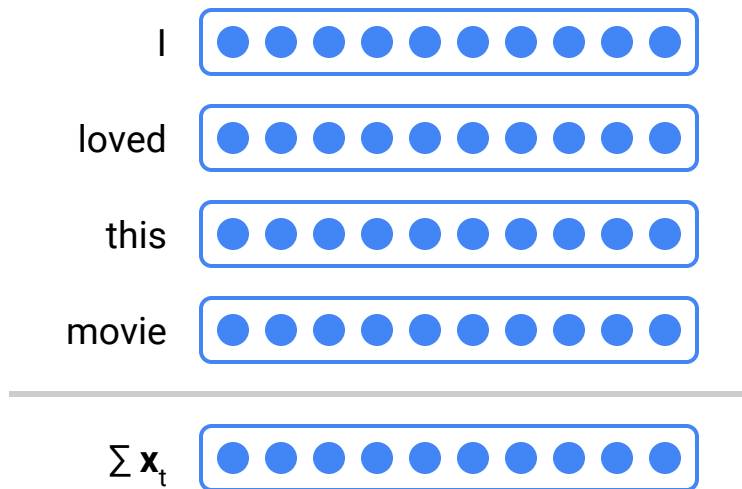


Why not this?



Deep CBOW

$$W'' \tanh(W' \tanh(W (\sum \mathbf{x}_t) + \mathbf{b}) + \mathbf{b}') + \mathbf{b}''$$



Softmax

We don't need a softmax for **prediction**, there we simply take the **argmax**

$$\mathbf{o} = [-0.1, 0.1, 0.1, \mathbf{2.4}, 0.2]$$

$$\text{softmax}(o_i) = \exp(o_i) / \sum_j \exp(o_j)$$

This makes \mathbf{o} sum to 1.0:

$$\text{softmax}(\mathbf{o}) = [0.0589, 0.0720, 0.0720, \mathbf{0.7177}, 0.0795]$$

Training a neural network

We train our network with Stochastic Gradient Descent (SGD):

1. Sample a training example
2. Forward pass
 - a. Compute network activations, output vector
3. Compute loss
 - a. Compare output vector with true label using a **loss function**
4. Backward pass (backpropagation)
 - a. Compute gradient of loss w.r.t. parameters
5. Take a small step in the opposite direction of the gradient

Cross Entropy Loss

Given:

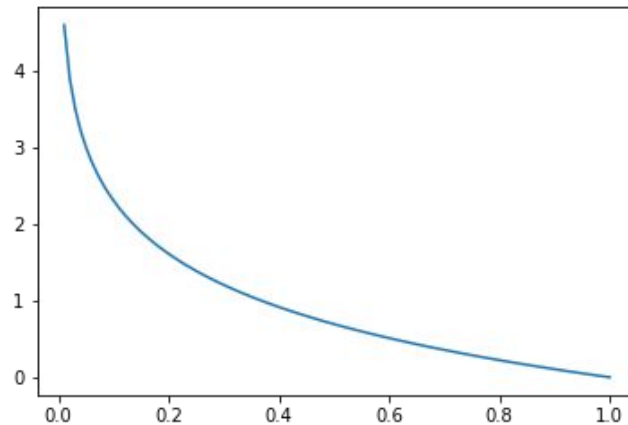
$$\hat{\mathbf{y}} = [0.1, 0.1, 0.1, 0.5, 0.2] \quad \text{output vector (after softmax) from forward pass}$$
$$\mathbf{y} = [0, 0, 0, 1, 0] \quad \text{target / label (} y_3 = 1 \text{)}$$

When our output is **categorical** (i.e. a number of classes), we can use a **Cross Entropy** loss:

$$\text{CE}(\mathbf{y}, \hat{\mathbf{y}}) = -\sum y_i \log \hat{y}_i$$

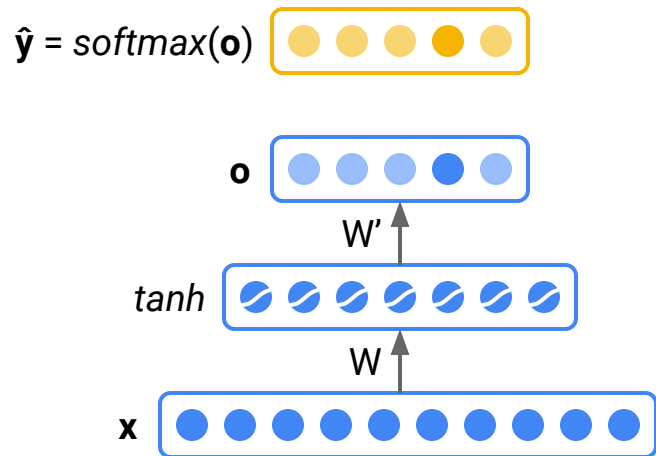
$$\text{SparseCE}(y = 3, \hat{\mathbf{y}}) = -\log \hat{y}_y$$

`torch.nn.CrossEntropyLoss`
works like this and does the
softmax on `o` for you!



Backpropagation example

the **chain rule** is your friend!
 $L = f(g(x))$
 $\delta L / \delta x = \delta f(g(x)) / \delta g(x) \cdot \delta g(x) / \delta x$



$$\hat{\mathbf{y}} = [0.1, 0.1, 0.1, 0.5, 0.2]$$
$$\mathbf{y} = [0, 0, 0, 1, 0]$$

$$\text{loss } L = \text{CE}(\hat{\mathbf{y}}, \mathbf{y}) = -\log(\hat{y}_3) = -\log(0.5)$$

compute gradients, e.g. for W' :

$$\begin{aligned} \delta L / \delta W' &= \delta L / \delta \mathbf{o} \cdot \delta \mathbf{o} / \delta W' \\ \delta L / \delta \mathbf{o} &= \delta L / \delta \hat{\mathbf{y}} \cdot \delta \hat{\mathbf{y}} / \delta \mathbf{o} \\ &= -1 / \hat{y}_3 \cdot \delta \text{softmax}(\mathbf{o}) / \delta \mathbf{o} \end{aligned}$$

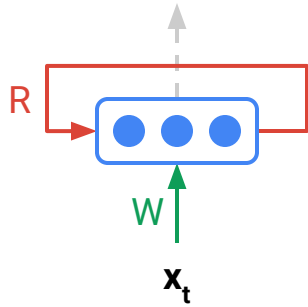
update weights:

$$W' = W' - \text{eta} * \delta L / \delta W'$$

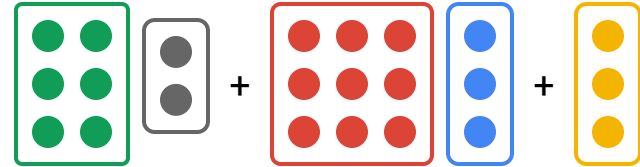
Sequences

Recurrent Neural Network (RNN)

RNNs model **sequential data** - one input \mathbf{x}_t per time step t



$$\begin{aligned} \mathbf{h}_t &= f(\mathbf{x}_t, \mathbf{h}_{t-1}) \\ &= \sigma(\mathbf{W}\mathbf{x}_t + \mathbf{R}\mathbf{h}_{t-1} + \mathbf{b}) \end{aligned}$$



Recurrent Neural Network (RNN)

Example:

the cat sat on the mat

\mathbf{x}_1 \mathbf{x}_2 \mathbf{x}_3 \mathbf{x}_4 \mathbf{x}_5 \mathbf{x}_6

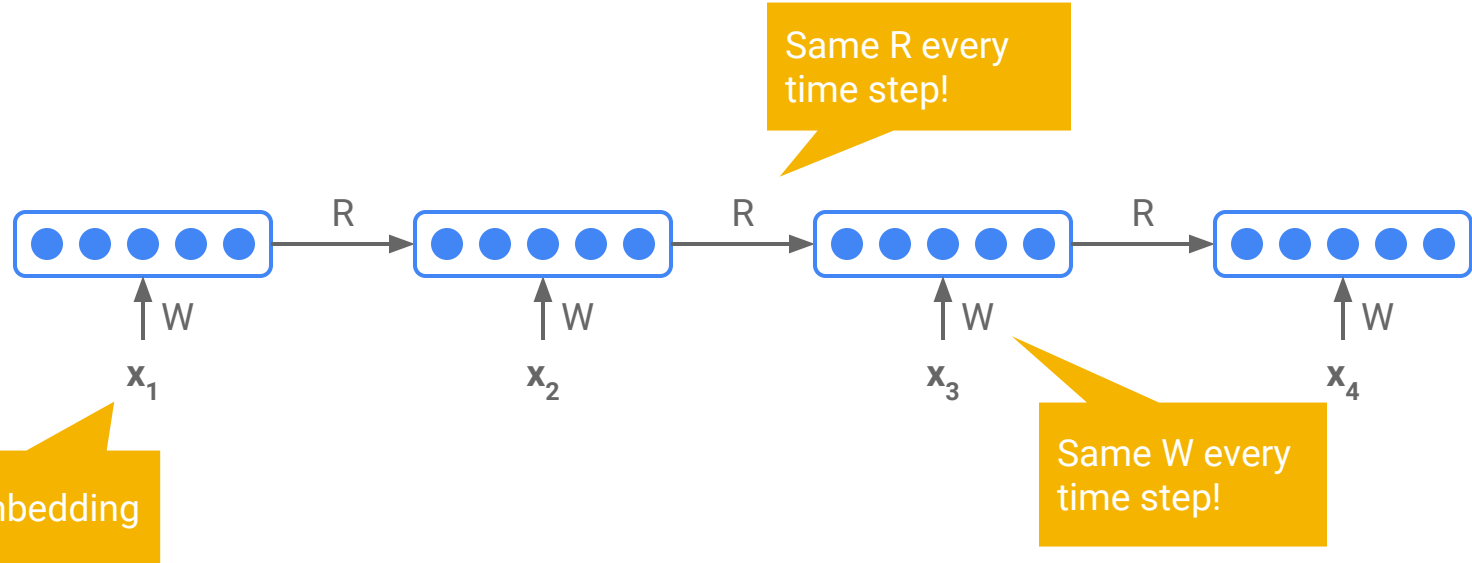
Let's compute the RNN state after reading in this sentence.

Remember:

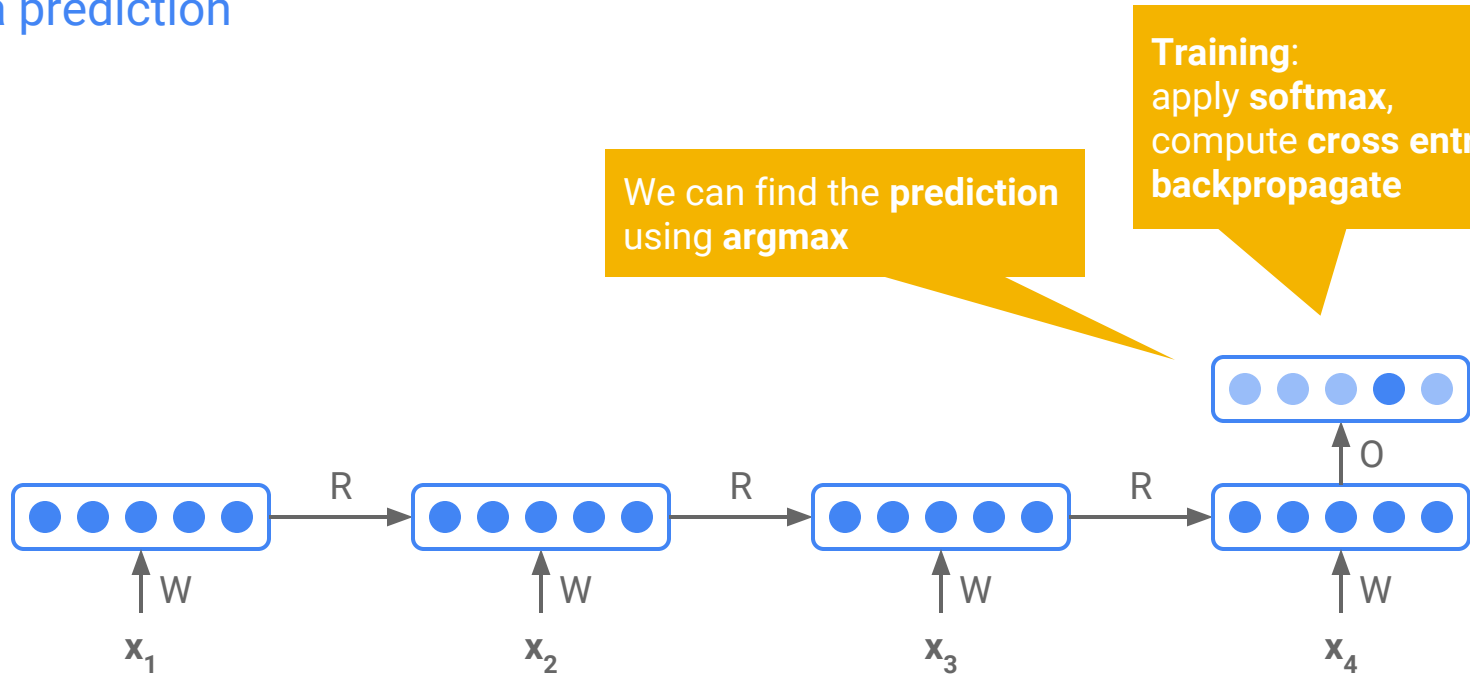
$$\mathbf{h}_t = f(\mathbf{x}_t, \mathbf{h}_{t-1})$$

$$\begin{aligned} \mathbf{h}_1 &= f(\mathbf{x}_1, \mathbf{h}_\theta) \\ \mathbf{h}_2 &= f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_\theta)) \\ \mathbf{h}_3 &= f(\mathbf{x}_3, f(\mathbf{x}_2, f(\mathbf{x}_1, \mathbf{h}_\theta))) \\ &\dots \\ \mathbf{h}_6 &= f(\mathbf{x}_6, f(\mathbf{x}_5, f(\mathbf{x}_4, \dots))) \end{aligned}$$

Unfolding the RNN



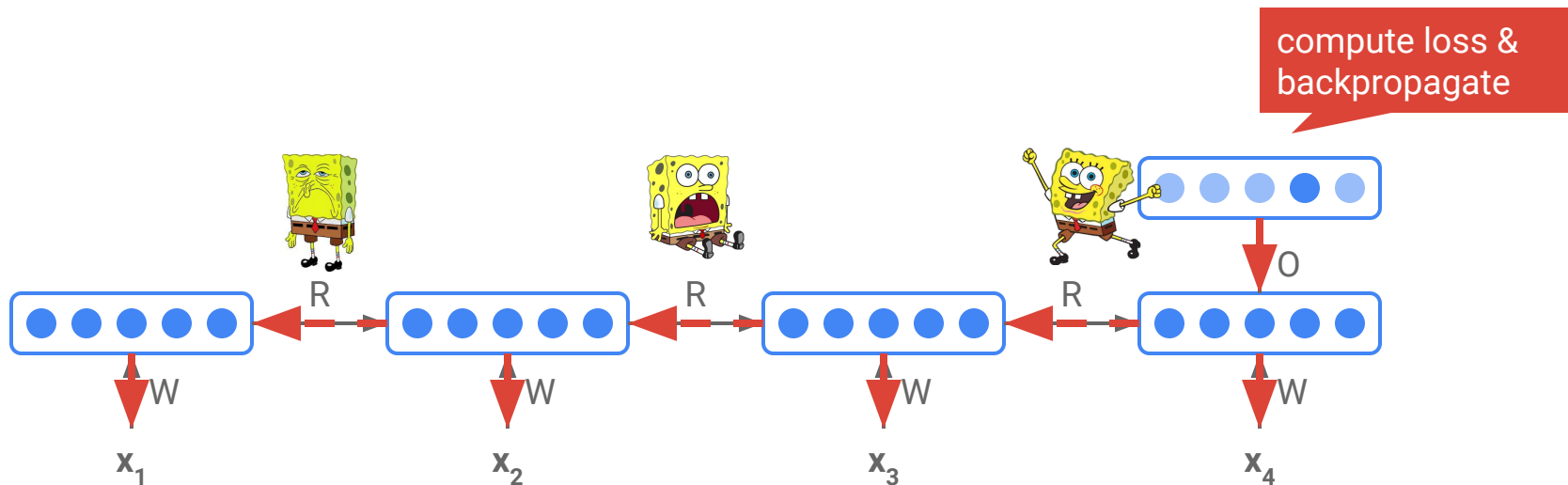
Making a prediction



The vanishing gradient problem

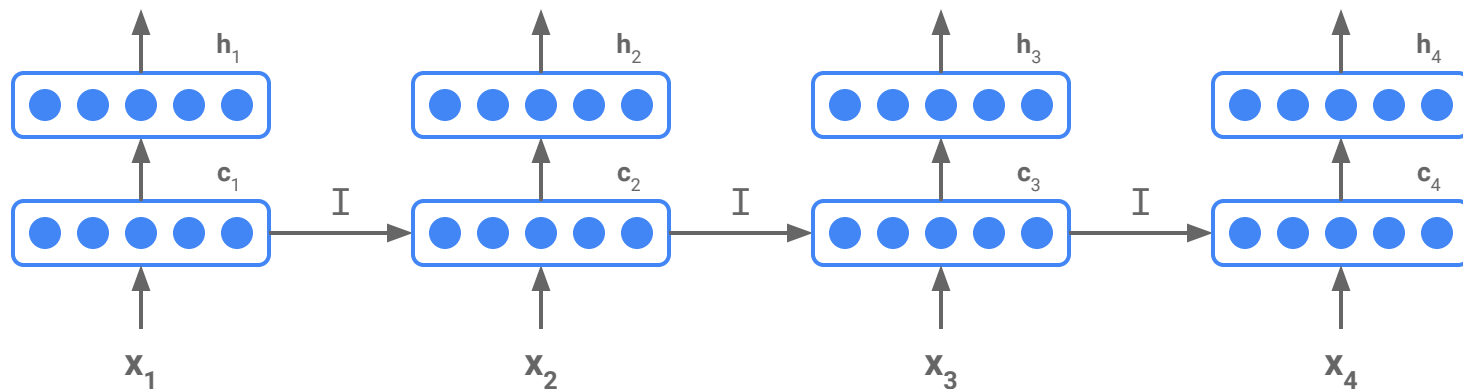
Simple RNNs are hard to train because of the **vanishing gradient** problem.

During backpropagation, **gradients** can quickly become **small**, as they **repeatedly** go through multiplications (R) & non-linear functions (e.g. sigmoid or tanh)



Intuition to solving the vanishing gradient

Let's use an extra vector, cell state \mathbf{c}



$$\mathbf{c}_t = \mathbf{c}_{t-1} + f(\mathbf{x}_t)$$

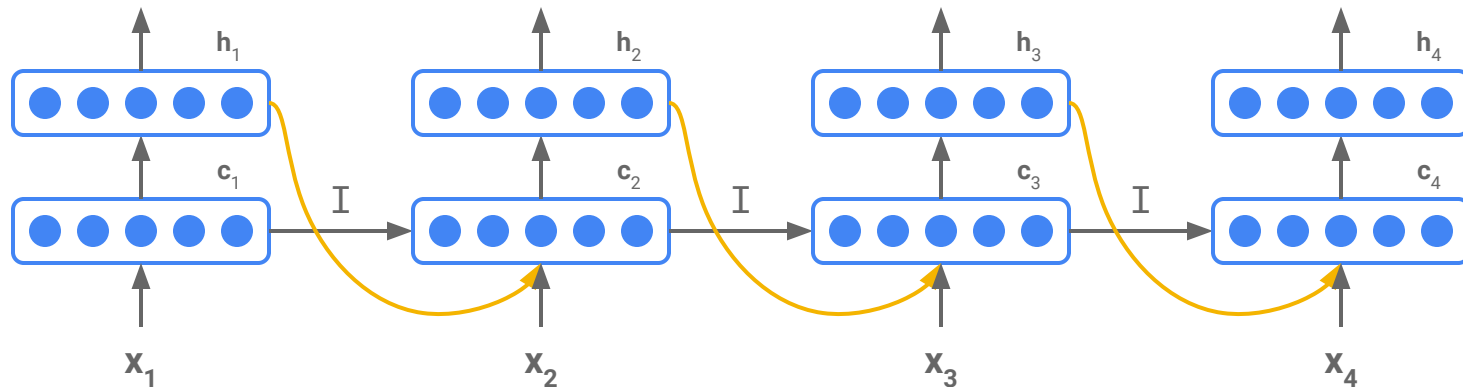
$$\mathbf{h}_t = \tanh(\mathbf{c}_t)$$

$$\frac{\delta \mathbf{c}_t}{\delta \mathbf{c}_{t-1}} = I$$

A small improvement



Better gradient propagation is possible when you use **additive** rather than multiplicative/highly non-linear recurrent dynamics



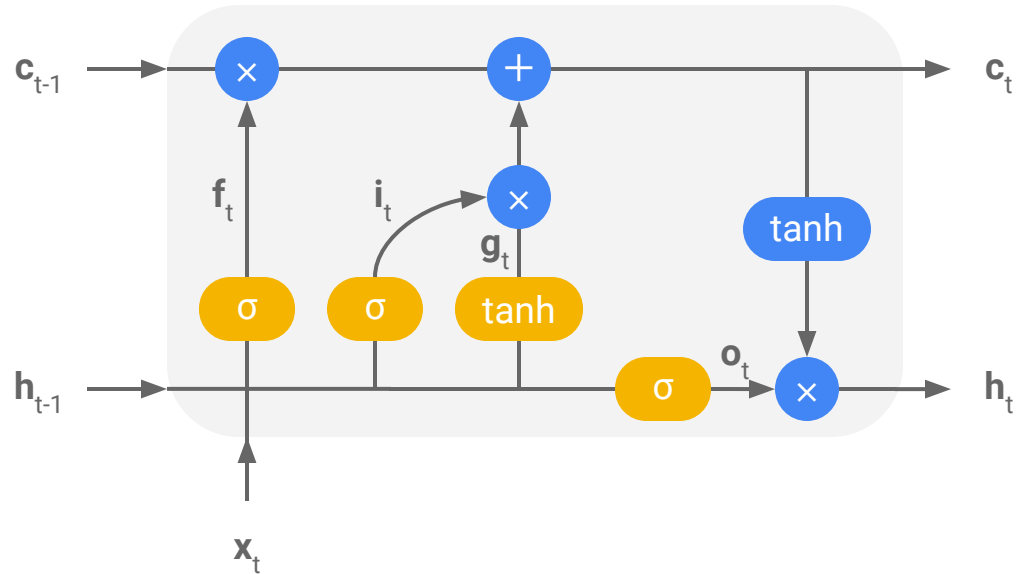
$$\mathbf{c}_t = \mathbf{c}_{t-1} + f(\mathbf{x}_t, \mathbf{h}_{t-1})$$

$$\mathbf{h}_t = \tanh(\mathbf{c}_t)$$

$$\frac{\delta \mathbf{c}_t}{\delta \mathbf{c}_{t-1}} = I + \epsilon$$

Long Short-Term Memory (LSTM)

LSTMs are a special kind of RNN that can deal with **long-term dependencies** in the data



Long Short-Term Memory (LSTM)

hidden state

cell state

previous hidden state and cell state

$$\mathbf{h}_t, \mathbf{c}_t = \text{lstm}(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{c}_{t-1})$$

input gate $\mathbf{i}_t = \sigma(W_i \mathbf{x}_t + R_i \mathbf{h}_{t-1} + \mathbf{b}_i)$

forget gate $\mathbf{f}_t = \sigma(W_f \mathbf{x}_t + R_f \mathbf{h}_{t-1} + \mathbf{b}_f)$

candidate $\mathbf{g}_t = \tanh(W_g \mathbf{x}_t + R_g \mathbf{h}_{t-1} + \mathbf{b}_g)$

output gate $\mathbf{o}_t = \sigma(W_o \mathbf{x}_t + R_o \mathbf{h}_{t-1} + \mathbf{b}_o)$

cell state $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{g}_t$

hidden state $\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$

Trees

Exploiting tree structure

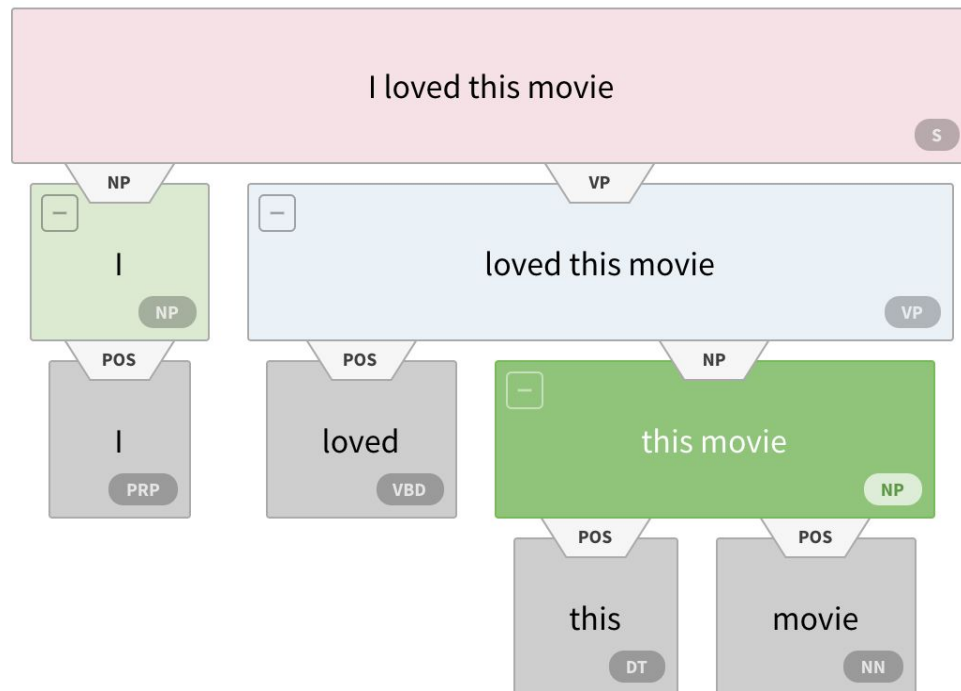
Instead of treating our input as a **sequence**, we can take an alternative approach: assume a **tree structure** and use the principle of **compositionality**.

The meaning (vector) of a sentence is determined by:

1. the meanings of its **words** and
2. the **rules** that combine them

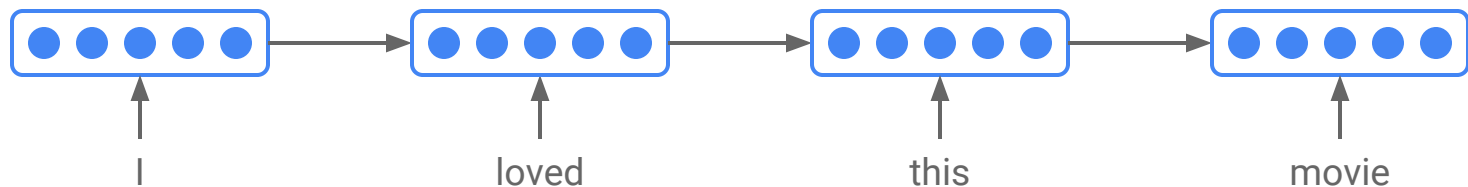
Constituency Parse

Can we obtain a sentence vector using the tree structure given by a parse?

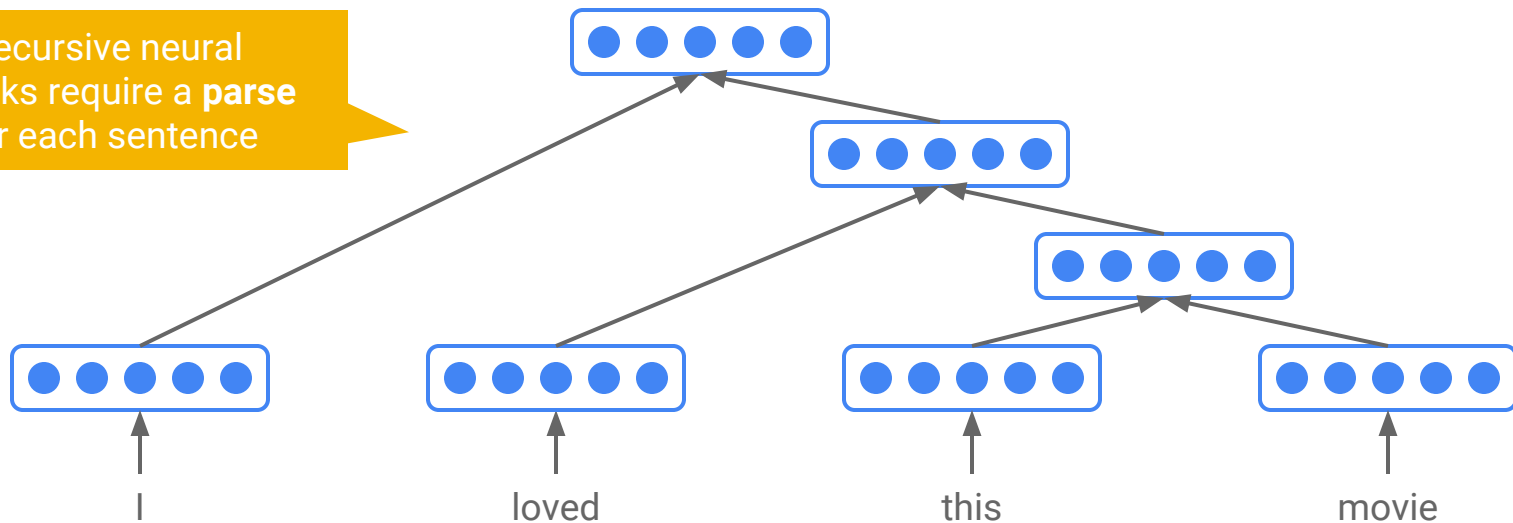


Recurrent vs Tree Recursive NN

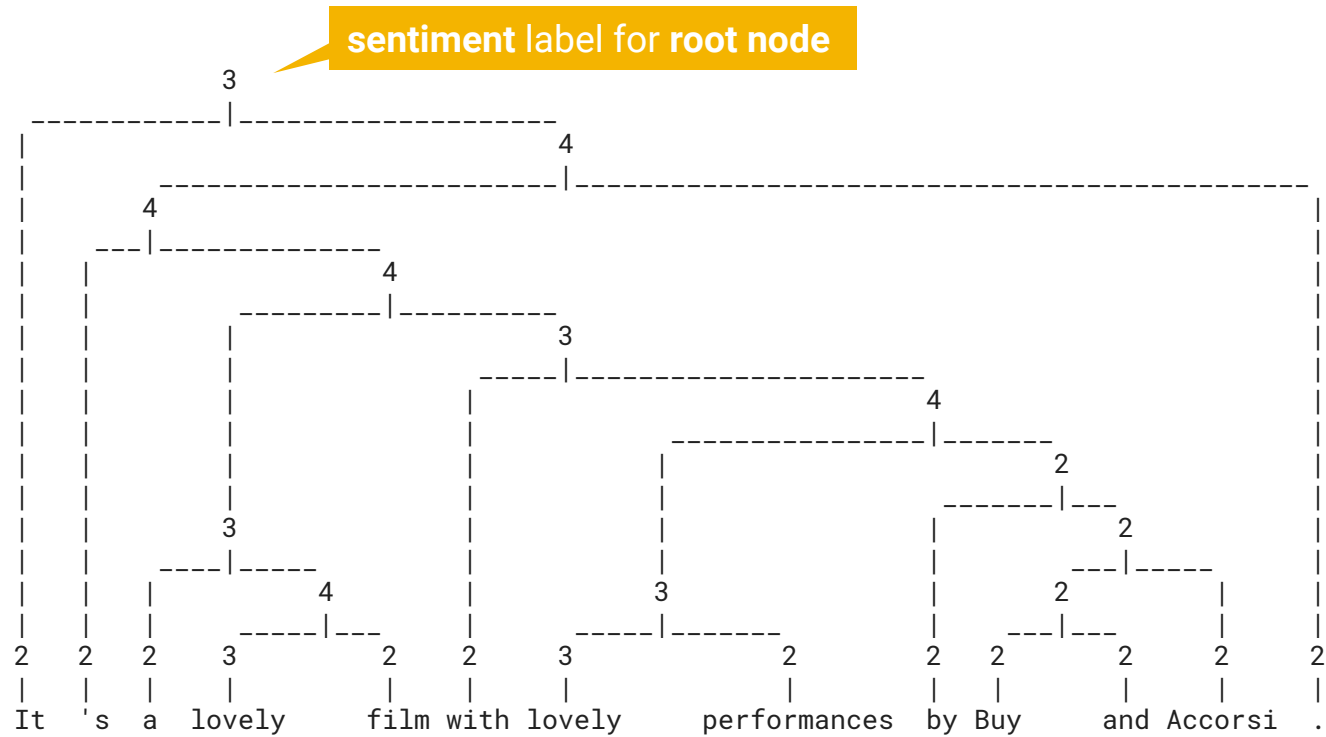
RNNs cannot capture phrases **without prefix context** and often capture too much of **last words** in final vector



Tree Recursive neural networks require a **parse tree** for each sentence



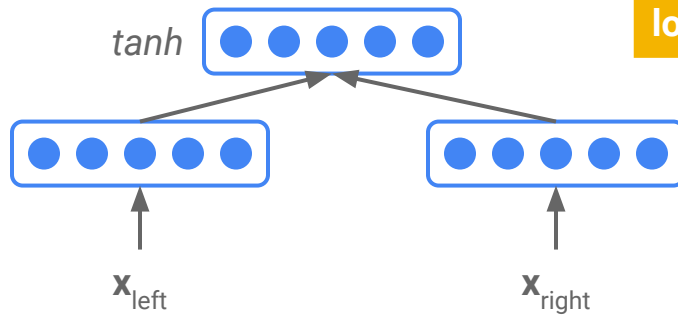
Practical II data set: Stanford Sentiment Treebank (SST)



A naive recursive NN

Combine every two children (left and right) into a parent node \mathbf{p} :

$$\mathbf{p} = \tanh(W_{\text{left}} \mathbf{x}_{\text{left}} + W_{\text{right}} \mathbf{x}_{\text{right}} + \mathbf{b})$$



a bit **simplistic** and
does not work well for
longer sentences

Tree LSTM next time!

Extra

Recap: Matrix Multiplication

Rows multiply with **columns**

1	2	3
4	5	6

2x3

×

1	2
1	2
1	2

3x2

=

1×1 + 2×1 + 3×1	1×2 + 2×2 + 3×2
4×1 + 5×1 + 6×1	4×2 + 5×2 + 6×2

2x2

Recap: Activation functions

