

THE UNIVERSITY OF BRITISH COLUMBIA

Topics in AI (CPSC 532L): **Multimodal Learning with Vision, Language and Sound**

Lecture 7: Word2Vec, Language Models and RNNs



Course Logistics

- Assignment 1 grades (available on **Connect** ???)
 - Solutions will be posted over the weekend
- Assignment 2 was due Yesterday
- Assignment 3 will be out Friday, January 26th
 - The due deadline will be extended

- Paper choices will be due **next week** (google form) - **Projects** groups and short description (google form)

Representing a Word: One Hot Encoding

Vocabulary

dog

cat

person

holding

tree

computer

using

Representing a Word: One Hot Encoding

Vocabulary

- dog
- 2 cat
- 3 person
- holding 4
- 5 tree
- computer 6
- using 7

Representing a Word: One Hot Encoding

Vocabulary

- dog
- 2 cat
- 3 person
- holding 4
- 5 tree
- computer 6
- using

one-hot encodings

[1, 0, 0, 0, 0, 0, 0, 0, 0, 0] [0, 1, 0, 0, 0, 0, 0, 0, 0, 0][0, 0, 1, 0, 0, 0, 0, 0, 0][0, 0, 0, **1**, 0, 0, 0, 0, 0, 0] [0, 0, 0, 0, 1, 0, 0, 0, 0][0, 0, 0, 0, 0, 1, 0, 0, 0][0,0,0,0,0,0,1,0,0]

bag-of-words representation

Vocabulary

dog	1
cat	2
person	3
holding	4
tree	5
computer	6
using	7

dog cat person holding tree tree using



person holding dog $\{3, 4, 1\}$ [1, 0, 1, 1, 0, 0, 0, 0, 0, 0]

bag-of-words representation

Vocabulary

dog	1
cat	2
person	3
holding	4
tree	5
computer	6
using	7

dog cat person holding tree tree tree using



person holding dog

person holding cat

bag-of-words representation **{3, 4, 1} [1, 0, 1, 1, 0, 0, 0, 0, 0, 0]**

 $\{3, 4, 2\}$ [1, 1, 0, 1, 0, 0, 0, 0, 0]

dog cat person holding tree tree tree using

Vocabulary

dog	1
cat	2
person	3
holding	4
tree	5
computer	6
using	7



person holding dog $\{3, 4, 1\}$ [1, 0, 1, 1, 0, 0, 0, 0, 0, 0]

- person holding cat
- person using computer $\{3, 7, 6\}$ [0, 0, 0, 1, 0, 1, 1, 0, 0, 0]

bag-of-words representation

 $\{3, 4, 2\}$ [1, 1, 0, 1, 0, 0, 0, 0, 0]

- dog cat person holding tree tree using

Vocabulary

dog	1
cat	2
person	3
holding	4
tree	5
computer	6
using	7



person holding dog

person holding cat

person using computer $\{3, 7, 6\}$ [0, 0, 0, 1, 0, 1, 1, 0, 0, 0]

person using computer person holding cat

bag-of-words representation **{3, 4, 1} [1, 0, 1, 1, 0, 0, 0, 0, 0, 0]**

- $\{3, 4, 2\}$ [1, 1, 0, 1, 0, 0, 0, 0, 0]

 - dog cat person holding tree computer using

$\{3, 3, 7, 6, 2\}$ [0, 1, 2, 1, 0, 1, 1, 0, 0, 0]

Vocabulary

dog	1
cat	2
person	3
holding	4
tree	5
computer	6
using	7



Distributional Hypothesis

- At least certain aspects of the meaning of lexical expressions depend on their distributional properties in the linguistic contexts

- The degree of semantic similarity between two linguistic expressions is a function of the similarity of the two linguistic contexts in which they can appear



* Adopted from slides by Louis-Philippe Morency

What is the meaning of "bardiwac"?

- He handed her glass of **bardiwac**.
- Beef dishes are made to complement the bardiwacs.
- Nigel staggered to his feet, face flushed from too much bardiwac.
- Malbec, one of the lesser-known **bardiwac** grapes, responds well to Australia's sunshine.
- I dined off bread and cheese and this excellent bardiwac.
- -The drinks were delicious: blood-red **bardiwac** as well as light, sweet Rhenish.

* Adopted from slides by Louis-Philippe Morency



What is the meaning of "bardiwac"?

- He handed her glass of **bardiwac**.
- Beef dishes are made to complement the **bardiwacs**.
- Nigel staggered to his feet, face flushed from too much bardiwac.
- Malbec, one of the lesser-known **bardiwac** grapes, responds well to Australia's sunshine.
- I dined off bread and cheese and this excellent bardiwac.
- -The drinks were delicious: blood-red **bardiwac** as well as light, sweet Rhenish.

bardic is an alcoholic beverage made from grapes

* Adopted from slides by Louis-Philippe Morency



Geometric Interpretation: Co-occurrence as feature

 Row vector describes usage of word in a corpus of text

 Can be seen as coordinates o the point in an n-dimensional Euclidian space

	get	see	use	hear	eat	kil
knife	51	20	84	0	3	0
cat	52	58	4	4	6	26
dog	115	83	10	42	33	17
boat	59	39	23	4	0	0
cup	98	14	6	2	1	0
pig	12	17	3	2	9	27
banana	11	2	2	0	18	0

Co-occurrence Matrix



Geometric Interpretation: Co-occurrence as feature

 Row vector describes usage of word in a corpus of text

Can be seen as coordinates o the point in an n-dimensional Euclidian space

	get	see	use	hear	eat	kil
knife	51	20	84	0	3	0
cat	52	58	4	4	6	26
dog	115	83	10	42	33	17
boat	59	39	23	4	0	0
cup	98	14	6	2	1	0
pig	12	17	3	2	9	27
banana	11	2	2	0	18	0

Co-occurrence Matrix



Distance and Similarity

Illustrated in two dimensions

 Similarity = spatial proximity (Euclidian distance)

 Location depends on frequency of **NOUN** (dog is 27 times as frequent as ca)



Angle and Similarity

direction is more important than location

normalize length of vectors

- or use angle as a distance measure



Geometric Interpretation: Co-occurrence as feature

 Row vector describes usage of word in a corpus of text

Can be seen as coordinates of the point in an n-dimensional Euclidian space

	get	see	use	hear	eat	kil
knife	51	20	84	0	3	0
cat	52	58	4	4	6	26
dog	115	83	10	42	33	17
boat	59	39	23	4	0	0
cup	98	14	6	2	1	0
pig	12	17	3	2	9	27
banana	11	2	2	0	18	0

Co-occurrence Matrix



Geometric Interpretation: Co-occurrence as feature

 Row vector describes usage of word in a corpus of text

Can be seen as coordinates of the point in an n-dimensional Euclidian space

Way too high dimensional!

	get	see	use	hear	eat	kil
knife	51	20	84	0	3	0
cat	52	58	4	4	6	26
dog	115	83	10	42	33	17
boat	59	39	23	4	0	0
cup	98	14	6	2	1	0
pig	12	17	3	2	9	27
banana	11	2	2	0	18	0

Co-occurrence Matrix



SVD for Dimensionality Reduction



Learned Word Vector Visualization

We can also use other methods, like LLE here:



Nonlinear dimensionality reduction by locally linear embedding. Sam Roweis & Lawrence Saul. Science, v.290,2000

[Roweis and Saul, 2000]



Issues with SVD

Computational cost for a $d \times n$ matrix is $\mathcal{O}(dn^2)$, where d < n

It is hard to incorporate out of sample (**new**) words or documents

Makes it not possible for large number of word vocabularies or documents

word2vec: Representing the Meaning of Words [Mikolov et al., 2013]

Key idea: Predict surrounding words of every word

Benefits: Faster and easier to incorporate new document, words, etc.



word2vec: Representing the Meaning of Words [Mikolov et al., 2013]

Key idea: Predict surrounding words of every word

Benefits: Faster and easier to incorporate new document, words, etc.

middle word

Skip-gram: use the middle word to predict surrounding ones in a window



CBOW

Skip-gram

Continuous Bag of Words (**CBOW**): use context words in a window to predict





Example: "The cat sat on floor" (window size 2)



[Mikolov et al., 2013]







[Mikolov et al., 2013]

sat (one-hot vector)







[Mikolov et al., 2013]









[Mikolov et al., 2013]

Parameters to be learned









[Mikolov et al., 2013]

Parameters to be learned

Size of the word vector (e.g., 300)







[Mikolov et al., 2013]







[Mikolov et al., 2013]

$\mathbf{W}_{ V imes N }^{T}$	\times	\mathbf{x}_{cat}	=	\mathbf{v}_{co}
$ V \times N $	\mathbf{A}	\mathbf{x}_{cat}		▼ C(

2.4		0		3.2	 	 0.9	0.5	1.8	1.6	2.4).1
2.6		1 0		6.1	 	 3.6	1.5	2.9	1.4	2.6).5
•••	=	0	×		 	 					
•••		0			 	 					
1.8		0		1.2	 	 2.0	2.4	1.9	2.7	1.8).6
		0									
		0									
		0									



.7	ţ			





[Mikolov et al., 2013]

$\mathbf{W}_{ V imes N }^{T}$	×	\mathbf{x}_{on}	=	$\mathbf{V}_{O'}$

1.8		0		3.2	 	 0.9	0.5	1.8	1.6	2.4).1
2.9		0		6.1	 	 3.6	1.5	2.9	1.4	2.6).5
	=	1	×		 	 					
		0			 	 					
1.9		0		1.2	 	 2.0	2.4	1.9	2.7	1.8).6
		0									
		0									
		0									









[Mikolov et al., 2013]







[Mikolov et al., 2013]







[Mikolov et al., 2013]











[Mikolov et al., 2013]


CBOW: Interesting Observation





[Mikolov et al., 2013]

*slide from Vagelis Hristidis



Skip-Gram Model



[Mikolov et al., 2013]



Comparison

Model	Vector	Training	Ac		
	Dimensionality	words			
			Semantic	Syntactic	Total
Collobert-Weston NNLM	50	660M	9.3	12.3	11.0
Turian NNLM	50	37M	1.4	2.6	2.1
Turian NNLM	200	37M	1.4	2.2	1.8
Mnih NNLM	50	37M	1.8	9.1	5.8
Mnih NNLM	100	37M	3.3	13.2	8.8
Mikolov RNNLM	80	320M	4.9	18.4	12.7
Mikolov RNNLM	640	320M	8.6	36.5	24.6
Huang NNLM	50	990M	13.3	11.6	12.3
Our NNLM	20	6B	12.9	26.4	20.3
Our NNLM	50	6B	27.9	55.8	43.2
Our NNLM	100	6B	34.2	64.5	50.8
CBOW	300	783M	15.5	53.1	36.1
Skip-gram	300	783M	50.0	55.9	53.3

[Mikolov et al., 2013]

- CBOW is not great for rare words and typically needs less data to train - Skip-gram better for rate words and needs more data to train the model



Interesting Results: Word Analogies

Test for linear relationships, examined by Mikolov et al. (2014)

a:b :: c:?

man:woman :: king:?

- + king [0.300.70]
- man [0.200.20]
- + woman [0.60 0.30]

queen [0.70 0.80]





Interesting Results: Word Analogies



[Mikolov et al., 2013]



Model the **probability of a sentence**; ideally be able to sample plausible sentences

Model the **probability of a sentence**; ideally be able to sample plausible sentences

Why is this useful?

Model the **probability of a sentence**; ideally be able to sample plausible sentences

Why is this useful?

arg max P(wordsequence | acoustics) = wordsequence



arg max *P*(*acoustics* | *wordsequence*) × *P*(*wordsequence*)

wordsequence

$P(acoustics | wordsequence) \times P(wordsequence)$ P(acoustics)

Model the **probability of a sentence**; ideally be able to sample plausible sentences

Why is this useful?

wordsequence



arg max *P*(*acoustics* | *wordsequence*) × *P*(*wordsequence*)

wordsequence

arg max P(wordsequence | acoustics) =

$P(acoustics | wordsequence) \times P(wordsequence)$ P(acoustics)

Simple Language Models: N-Grams

Given a word sequence: $w_{1:n} = [w_1, w_2, ..., w_n]$

We want to estimate $p(w_{1:n})$

Simple Language Models: N-Grams

Given a word sequence: $w_{1:n} = [w_1, w_2, ..., w_n]$

We want to estimate $p(w_{1:n})$

Using **Chain Rule** of probabilities:

 $p(w_{1:n}) = p(w_1)p(w_2|w_1)p(w_3|w_1, w_2) \cdots p(w_n|w_{1:n-1})$

Simple Language Models: N-Grams

Given a word sequence: $w_{1:n} = [w_1, w_2, ..., w_n]$

We want to estimate $p(w_{1:n})$

Using **Chain Rule** of probabilities:

$$p(w_{1:n}) = p(w_1)p(w_2|w_1)p(w_3|w_1, w_2) \cdots p(w_n|w_{1:n-1})$$

Bi-gram Approximation:

$$p(w_{1:n}) = \prod_{k=1}^{n} p(w_k | w_{k-1})$$

N-gram Approximation:

$$p(w_{1:n}) = \prod_{k=1}^{n} p(w_k | w_{k-N+1:k-1})$$

Estimating **Probabilities**

N-gram conditional probabilities can counts in the observed sequences

Bi-gram:

 $p(w_n|w_{n-1}) =$

N-gram:

 $p(w_n | w_{n-N-1:n-1}) =$

N-gram conditional probabilities can be estimated based on raw concurrence

$$\frac{C(w_{n-1}w_n)}{C(w_{n-1})}$$

$$\frac{C(w_{n-N-1:n-1}w_n)}{C(w_{n-N-1:n-1})}$$

Neural-based Unigram Language Mode







Neural-based Unigram Language Mode





Problem: Does not model sequential information (too local)

Neural-based Unigram Language Mode



We need sequence modeling!



Problem: Does not model sequential information (too local)

Why Model Sequences?







Image Credit: Alex Graves and Kevin Gimpel

* slide from Dhruv Batra

Multi-modal tasks



[Vinyals *et al.*, 2015]



Sequences where you don't expect them ...

Classify images by taking a series of "glimpses"

[Gregor et al., ICML 2015] [Mnih et al., ICLR 2015]

2	10	8	2	9	1	ł	1	ļ	8
3	3	3	8	6	9	6	5	1	3
8	8	1	8	2	6	9	¥	3	4
F	0	2	1	6	\mathcal{O}	9	ŀ	4	5
7	/	4	4	4	A	4	ų	7	9
3	1	8	9	3	4	2	4	7	3
6	6	1	6	З	- An	3	3	-	0
b	1	۵	Б	3	5	1	8	3	4
9	9	ł	1	3	0	5	9	5	4
1	1	8	4	9	8	20	2		R

Sequences where you don't expect them ...

Classify images by taking a series of "glimpses"

[Gregor et al., ICML 2015] [Mnih et al., ICLR 2015]

2	10	8	2	9	1	ł	1	ļ	8
3	3	3	8	6	9	6	5	1	3
8	8	1	8	2	6	9	¥	3	4
F	0	2	1	6	\mathcal{O}	9	ŀ	4	5
7	/	4	4	4	A	4	ų	7	9
3	1	8	9	3	4	2	4	7	3
6	6	1	6	З	- An	3	3	-	0
b	1	۵	Б	3	5	1	8	3	4
9	9	ł	1	3	0	5	9	5	4
1	1	8	4	9	8	20	2		R

Sequences in Inputs or Outputs?

one to one



Input: No sequence Output: No seq.

Example:

"standard" classification / regression problems

Sequences in Inputs or Outputs?



Input: No sequence Output: No seq.

Example:

"standard" classification / regression problems

Input: No sequence **Output:** Sequence **Example:** Im₂Caption

Input: Sequence Output: No seq. **Example:** sentence classification, multiple-choice question answering

- **Input:** Sequence **Output:** Sequence
- **Example:** machine translation, video captioning, open-ended question answering, video question answering
- * slide from Fei-Dei Li, Justin Johnson, Serena Yeung, cs231n Stanford



Key Conceptual Ideas

Parameter Sharing

- in computational graphs = adding gradients

"Unrolling"

in computational graphs with parameter sharing

Parameter Sharing + "Unrolling"

- Allows modeling arbitrary length sequences!
- Keeps number of parameters in check

* slide from Dhruv Batra



y RNN

X

usually want to predict a vector at some time steps

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



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

$h_{t} = f_{W}(h_{t-1}, x_{t})$

Note: the same function and the same set of parameters are used at every time step



(Vanilla) Recurrent Neural Network

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



(Vanilla) **Recurrent** Neural Network

$h_t = f_W(h_{t-1}, x_t)$ $h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$

* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, cs231n Stanford

V

RNN

X

(Vanilla) **Recurrent** Neural Network

$y_t = W_{hy}h_t + b_y$

$h_t = f_W(h_{t-1}, x_t)$ $h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$









Re-use the same weight matrix at every time-step



RNN Computational Graph: Many to Many



RNN Computational Graph: Many to Many


RNN Computational Graph: Many to Many



RNN Computational Graph: Many to One



RNN Computational Graph: One to Many



Sequence to Sequence: Many to One + One to Many

Many to one: Encode input sequence in a single vector





Sequence to Sequence: Many to One + One to Many

Many to one: Encode input sequence in a single vector



One to many: Produce output sequence from single input vector



Example: Character-level Language Model

Vocabulary: ['h', 'e', 'l', 'o']

Example training sequence: "hello"

Example: Character-level Language Model

Vocabulary: ['h', 'e', 'l', 'o']

Example training sequence: "hello"

$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$

Example: Character-level Language Model

Vocabulary: ['h', 'e', 'l', 'o']

Example training sequence: "hello"

Vocabulary: ['h', 'e', 'l', 'o']

At test time sample one character at a time and feed back to the model

Vocabulary: ['h', 'e', 'l', 'o']

At test time sample one character at a time and feed back to the model

Vocabulary: ['h', 'e', 'l', 'o']

At test time sample one character at a time and feed back to the model

Vocabulary: ['h', 'e', 'l', 'o']

At test time sample one character at a time and feed back to the model

BackProp Through Time

sequence to compute gradient

Forward through entire sequence to compute loss, then backward through entire

Truncated BackProp Through Time

instead of the whole sequence

Run backwards and forwards through (fixed length) chunks of the sequence,

Truncated BackProp Through Time

Run backwards and forwards through (fixed length) chunks of the sequence, instead of the whole sequence

Carry hidden states forward, but only BackProp through some smaller number of steps

Truncated BackProp Through Time

instead of the whole sequence

Run backwards and forwards through (fixed length) chunks of the sequence,

Implementation: Relatively Easy

... you will have a chance to experience this in the Assignment 3

Learning to Write Like Shakespeare

THE SONNETS

by William Shakespeare

From fairest creatures we desire increase, That thereby beauty's rose might never die, But as the riper should by time decease, His tender heir might bear his memory: But thou, contracted to thine own bright eyes, Feed'st thy light's flame with self-substantial fuel, Making a famine where abundance lies, Thyself thy foe, to thy sweet self too cruel: Thou that art now the world's fresh ornament, And only herald to the gaudy spring, Within thine own bud buriest thy content, And tender churl mak'st waste in niggarding: Dity the world, or else this glutten be

Pity the world, or else this glutton be, To eat the world's due, by the grave and thee.

When forty winters shall besiege thy brow, And dig deep trenches in thy beauty's field, Thy youth's proud livery so gazed on now, Will be a tatter'd weed of small worth held: Then being asked, where all thy beauty lies, Where all the treasure of thy lusty days; To say, within thine own deep sunken eyes, Were an all-eating shame, and thriftless praise. How much more praise deserv'd thy beauty's use, If thou couldst answer 'This fair child of mine Shall sum my count, and make my old excuse,' Proving his beauty by succession thine!

This were to be new made when thou art old, And see thy blood warm when thou feel'st it cold.

Learning to Write Like Shakespeare ... after training a bit

at first:

tyntd-iafhatawiaoihrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e plia tklrgd t o idoe ns,smtt h ne etie h,hregtrs nigtike,aoaenns lng

"Tmont thithey" fomesscerliund Keushey. Thom here sheulke, anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

Aftair fall unsuch that the hall for Prince Velzonski's that me of her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort how, and Gogition is so overelical and ofter.

"Why do what that day," replied Natasha, and wishing to himself the fact the princess, Princess Mary was easier, fed in had oftened him. Pierre aking his soul came to the packs and drove up his father-in-law women.

train more

train more

train more

Learning to Write Like Shakespeare ... after training

PANDARUS:

Alas, I think he shall be come approached and the day When little srain would be attain'd into being never fed, And who is but a chain and subjects of his death, I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul, Breaking and strongly should be buried, when I perish The earth and thoughts of many states.

DUKE VINCENTIO: Well, your wit is in the care of side and that.

```
Second Lord:
They would be ruled after this chamber, and
my fair nues begun out of the fact, to be conveyed,
Whose noble souls I'll have the heart of the wars.
```

Clown: Come, sir, I will make did behold your worship.

VIOLA: I'll drink it.

VIOLA:

Why, Salisbury must find his flesh and thought That which I am not aps, not a man and in fire, To show the reining of the raven and the wars To grace my hand reproach within, and not a fair are hand, That Caesar and my goodly father's world; When I was heaven of presence and our fleets, We spare with hours, but cut thy council I am great, Murdered and by thy master's ready there My power to give thee but so much as hell: Some service in the noble bondman here, Would show him to her wine.

KING LEAR:

O, if you were a feeble sight, the courtesy of your law, Your sight and several breath, will wear the gods With his heads, and my hands are wonder'd at the deeds, So drop upon your lordship's head, and your opinion Shall be against your honour.

Learning Code

```
static void do_command(struct seq_file *m, void *v)
  int column = 32 << (cmd[2] & 0x80);</pre>
  if (state)
    cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
  else
    seq = 1;
  for (i = 0; i < 16; i++) {</pre>
    if (k & (1 << 1))
      pipe = (in_use & UMXTHREAD_UNCCA) +
        ((count & 0x0000000fffffff8) & 0x000000f) << 8;
    if (count == 0)
      sub(pid, ppc_md.kexec_handle, 0x2000000);
    pipe_set_bytes(i, 0);
  /* Free our user pages pointer to place camera if all dash */
  subsystem info = &of_changes[PAGE_SIZE];
  rek_controls(offset, idx, &soffset);
  /* Now we want to deliberately put it to device */
  control_check_polarity(&context, val, 0);
  for (i = 0; i < COUNTER; i++)</pre>
    seq_puts(s, "policy ");
}
```

Trained on entire source code of Linux kernel

DopeLearning: Computational Approach to Rap Lyrics

Everybody got one And all the pretty mommies want some And what i told you all was But you need to stay such do not touch They really do not want you to vote what do you condone Music make you lose control What you need is right here all oh This is for you and me I had to dedicate this song to you Mami Now I see how you can be I see u smiling i kno u hattig Best I Eva Had x4 That I had to pay for Do I have the right to take yours Trying to stay warm

- (2 Chainz Extremely Blessed)
- (Mos Def Undeniable)
- (Lil Wayne Welcome Back)
- (Common Heidi Hoe)
- (KRS One The Mind)
- (Cam'ron Bubble Music)
- (Missy Elliot Lose Control)
- (Wiz Khalifa Right Here)
- (Missy Elliot Hit Em Wit Da Hee)
- (Fat Joe Bendicion Mami)
- (Lil Wayne How To Hate)
- (Wiz Khalifa Damn Thing)
- (Nicki Minaj Best I Ever Had)
- (Ice Cube X Bitches)
- (Common Retrospect For Life)
- (Everlast 2 Pieces Of Drama)

[Malmi et al., KDD 2016]

Sunspring: First movie generated by Al

Sunspring, a short science fiction movie written entirely by AI, debuts exclusively on Ars today.

Sunspring | A Sci-Fi Short Film Starring Thomas Middleditch

Multilayer RNNs

$$\begin{aligned} h^l_t &= \tanh W^l \begin{pmatrix} h^{l-1}_t \\ h^l_{t-1} \end{pmatrix} \\ h \in \mathbb{R}^n, \qquad W^l \ [n \times 2n] \end{aligned}$$

[Bengio et al., 1994] [Pascanu et al., ICML 2013]

$$h_{t} = \tanh(W_{hh}h_{t-1} + W_{xh}x_{t})$$
$$= \tanh\left(\left(W_{hh} \quad W_{hx}\right) \begin{pmatrix}h_{t-1}\\x_{t}\end{pmatrix}\right)$$
$$= \tanh\left(W\begin{pmatrix}h_{t-1}\\x_{t}\end{pmatrix}\right)$$

Backpropagation from h_t to h_{t-1} multiplies by W (actually W_{hh}^{T})

[Bengio et al., 1994] [Pascanu et al., ICML 2013]

$$h_{t} = \tanh(W_{hh}h_{t-1} + W_{xh}x_{t})$$
$$= \tanh\left(\left(W_{hh} \quad W_{hx}\right) \begin{pmatrix}h_{t-1}\\x_{t}\end{pmatrix}\right)$$
$$= \tanh\left(W\begin{pmatrix}h_{t-1}\\x_{t}\end{pmatrix}\right)$$

Computing gradient of h₀ involves many factors of W (and repeated tanh)

[Bengio et al., 1994] [Pascanu et al., ICML 2013]

Computing gradient of h₀ involves many factors of W (and repeated tanh)

Exploding gradients

Vanishing gradients

[Bengio et al., 1994] [Pascanu et al., ICML 2013]

Largest singular value > 1:

Largest singular value < 1:

Computing gradient of h₀ involves many factors of W (and repeated tanh)

Largest singular value > 1: **Exploding gradients**

Largest singular value < 1: Vanishing gradients

[Bengio et al., 1994] [Pascanu et al., ICML 2013]

Gradient clipping: Scale gradient if its norm is too big

> grad_norm = np.sum(grad * grad) if grad_norm > threshold: grad *= (threshold / grad_norm)

Computing gradient of h₀ involves many factors of W (and repeated tanh)

Exploding gradients

Vanishing gradients

[Bengio et al., 1994] [Pascanu et al., ICML 2013]

Largest singular value > 1:

Largest singular value < 1: Change RNN architecture

Long-Short Term Memory (LSTM)

Vanilla RNN

$$h_t = \tanh\left(W\begin{pmatrix}h_{t-1}\\x_t\end{pmatrix}\right)$$

LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

[Hochreiter and Schmidhuber, NC **1977**]

Long-Short Term Memory (LSTM)

Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

* slide from Dhruv Batra

Long-Short Term Memory (LSTM)

Cell state / **memory**

Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

* slide from Dhruv Batra

LSTM Intuition: Forget Gate

Should we continue to **remember** this "bit" of information or not?

Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$f_t = \sigma \left(W_f \cdot [h_{t-1}, x_t] + b_f \right)$

* slide from Dhruv Batra

, Dotr

LSTM Intuition: Forget Gate

Should we continue to **remember** this "bit" of information or not?

Intuition: memory and forget gate output multiply, output of forget gate can be though of as binary (0 or 1) anything x 1 = anything (remember) anything x 0 = 0 (forget)

Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$f_t = \sigma \left(W_f \cdot [h_{t-1}, x_t] + b_f \right)$

* slide from Dhruv Batra

LSTM Intuition: Input Gate

Should we **update** this "bit" of information or not? If yes, then what should we **remember**?

Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$$i_t = \sigma \left(W_i \cdot [h_{t-1}, x_t] + b_i \right)$$
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

* slide from Dhruv Batra

, Dotr
LSTM Intuition: Memory Update



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

Forget what needs to be forgotten + memorize what needs to be remembered

$C_t = f_t * C_{t-1} + i_t * C_t$



LSTM Intuition: Output Gate

Should we output this bit of information (e.g., to "deeper" LSTM layers)?



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$o_t = \sigma \left(W_o \left[h_{t-1}, x_t \right] + b_o \right)$ $h_t = o_t * \tanh(C_t)$

* slide from Dhruv Batra

, Dotr

LSTM Intuition: Additive Updates

Backpropagation from c_t to c_{t-1} only elementwise multiplication by f, no matrix multiply by W



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

* slide from Dhruv Batra

, Dotr

LSTM Intuition: Additive Updates



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

Uninterrupted gradient flow!



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

Uninterrupted gradient flow!



LSTM Variants: with Peephole Connections

Lets gates see the cell state / memory



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$$f_t = \sigma \left(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f \right)$$

$$i_t = \sigma \left(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i \right)$$

$$o_t = \sigma \left(W_o \cdot [C_t, h_{t-1}, x_t] + b_o \right)$$

* slide from Dhruv Batra

, Dotr

LSTM Variants: with Coupled Gates

Only memorize new information when you're forgetting old



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$C_t = f_t * C_{t-1} + (1 - f_t) * \tilde{C}_t$

Gated Recurrent Unit (GRU)

No explicit memory; memory = hidden output



z = memorize new and forget old

Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$$z_t = \sigma \left(W_z \cdot [h_{t-1}, x_t] \right)$$
$$r_t = \sigma \left(W_r \cdot [h_{t-1}, x_t] \right)$$
$$\tilde{h}_t = \tanh \left(W \cdot [r_t * h_{t-1}, x_t] \right)$$
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Phased LSTM

Gates are controlled by **phased** (periodic) **oscillations**



[Neil et al., 2016]



Skip-thought Vectors



word2vec but for sentences, where each sentence is processed by an LSTM

[Kiros et al., 2015]