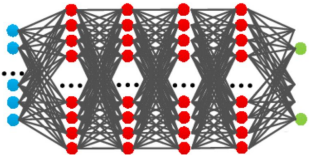# Deep Learning

Week 02: Word Embeddings

---

## Logistics

- HW1 is due on Thursday
- Submit on gradescope
  - If you worked in a group, create a group and then submit
- Clarifications are on Ed
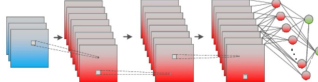- Come to office hours if you have questions

---

MLPs:

- Fully connected layers
- Require more parameters and computational resources
- Flexible and can handle various input types

CNNs:

- Convolutional layers with filters
- Designed specifically for structured input like images
- Inherently translation invariant due to shared weights
- Requires fewer parameters



https://www.researchgate.net/figure/The-architecture-of-MLP-and-CNN-MLP-is-consisted-of-fullyconnected-FC-layers-and-CNN_fig2_334489445

---

*I don't know how to parallel park.*
*I'm taking my dog for a walk at the park.* — Homonyms

*We ate outside and swam in the lake all week.*
*We ate outside and in the lake all week.* — Typos

*Biden speaks to the media in Illinois.*
*The president greets the press in Chicago.* — Paraphrases/ Synonyms

*Although interchangeable, the body pieces on the 2 cars are not similar.*
*Although similar, the body pieces are not interchangeable on the 2 cars.* — Word order

---

## How to handle text data?

---

## Language Modeling: predict the next word

**Assign probabilities to text.**

Given a sequence $(\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T)$, we want to **maximize** $P(\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T)$.

$$P(\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T) = P(\mathbf{x}_1)P(\mathbf{x}_2|\mathbf{x}_1)P(\mathbf{x}_3|\mathbf{x}_1, \mathbf{x}_2)P(\mathbf{x}_4|\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)\ldots P(\mathbf{x}_T|\mathbf{x}_1, \mathbf{x}_2, \ldots \mathbf{x}_{T-1})$$

P(I like cats because they look cute) = P(I) P(like | I) P(cats | I like) P(as | I like cats)  P(they | I like cats because)

P(look | I like cats because they)  P(cute | I like cats because they look)

Predict the next word given current text!

---

## *n*-Gram Language Model

*n*-Gram: chunk of *n* consecutive words

**Count** the frequency of each *n*-grams and predict next word!

**Assume** each word only depends on previous *n - 1* words.

**Uni-gram:** "I" "like" "cats" "as" "they" "look" "cute"

**Bi-gram:** "I like" "like cats" "cats as" "as they" …

**Tri-gram:** "I like cats" "like cats as" "cats as they" …

$$P(\mathbf{x}_t|\mathbf{x}_1, \ldots, \mathbf{x}_{t-1}) = P(\mathbf{x}_t|\mathbf{x}_{t-n+1}, \ldots, \mathbf{x}_{t-1})$$
$$= \frac{\text{count}(\mathbf{x}_{t-n+1}, \ldots, \mathbf{x}_{t-1}, \mathbf{x}_t)}{\text{count}(\mathbf{x}_{t-n+1}, \ldots, \mathbf{x}_{t-1})}$$

In *bi-gram* LM

P(I like cats as they look cute) = P(I) P(like | I) P(cats | like) P(as | cats) P(they | because) P(look | they) P(cute | look)

Discuss: Do you want to have a large n or a small n in a
n-gram model?

---

## *n*-Gram Language Model: issue

*n*-Gram: chunk of *n* consecutive words

**Count** the frequency of each *n*-grams and
predict next word!

**Uni-gram:** "I" "like" "cats" "as" "they" "look" "cute"

**Bi-gram:** "I like" "like cats" "cats as" "as they" …

**Tri-gram:** "I like cats" "like cats as" "cats as they" …

**Assume** each word only depends on
previous *n* - 1 words.

$$P(\mathbf{x}_t|\mathbf{x}_1,\ldots,\mathbf{x}_{t-1}) = P(\mathbf{x}_t|\mathbf{x}_{t-n+1},\ldots,\mathbf{x}_{t-1})$$
$$= \frac{\text{count}(\mathbf{x}_{t-n+1},\ldots,\mathbf{x}_{t-1},\mathbf{x}_t)}{\text{count}(\mathbf{x}_{t-n+1},\ldots,\mathbf{x}_{t-1})}$$

**Increase *n* provides contextual information, but exponentially
increase the size of the counting table!**

---

## Bag of Words



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

| | |
|---|---|
| it | 6 |
| I | 5 |
| the | 4 |
| to | 3 |
| and | 3 |
| seen | 2 |
| yet | 1 |
| would | 1 |
| whimsical | 1 |
| times | 1 |
| sweet | 1 |
| satirical | 1 |
| adventure | 1 |
| genre | 1 |
| fairy | 1 |
| humor | 1 |
| have | 1 |
| great | 1 |
| … | … |

https://koushik1102.medium.com/nlp-bag-of-words-and-tf-idf-explained-fd1f49dce7c4
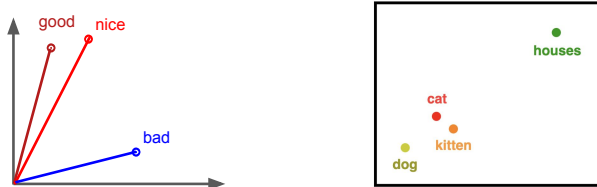
---

## What are word embeddings

- What are Word Embeddings?
  - vector representations of words that capture semantic relationships

## Semantic similarity

- Motivation
  - Put words into vectors so we can measure the similarity between words
  - Use cosine similarity

## Why Do We Need Word Embeddings?

- Why Do We Need Word Embeddings?
  - Numerical Input
  - Shows Similarity and Distance



| | living being | feline | human | gender | royalty | verb | plural |
|---|---|---|---|---|---|---|---|
| cat | 0.6 | 0.9 | 0.1 | 0.4 | -0.7 | -0.3 | -0.2 |
| kitten | 0.5 | 0.8 | -0.1 | 0.2 | -0.6 | -0.5 | -0.1 |
| dog | 0.7 | -0.1 | 0.4 | 0.3 | -0.4 | -0.1 | -0.3 |
| houses | -0.8 | -0.4 | -0.5 | 0.1 | -0.9 | 0.3 | 0.8 |
| man | 0.6 | -0.2 | 0.8 | 0.9 | -0.1 | -0.9 | -0.7 |
| woman | 0.7 | 0.3 | 0.9 | -0.7 | 0.1 | -0.5 | -0.4 |
| king | 0.5 | -0.4 | 0.7 | 0.8 | 0.9 | -0.7 | -0.6 |
| queen | 0.8 | -0.1 | 0.8 | -0.9 | 0.8 | -0.5 | -0.9 |

**embedding using features of words**

## What does ong choi mean?

Suppose you see these sentences:

- Ong choi is delicious sautéed with garlic.
- Ong choi is superb over rice
- Ong choi leaves with salty sauces

And you've also seen these:

- ...spinach sautéed with garlic over rice
- Chard stems and leaves are delicious
- Collard greens and other salty leafy greens

https://web.stanford.edu/~jurafsky/slp3/

Generative AI is experimental. Learn more

Ong choy is **a leafy green vegetable with long, hollow stems and slender leaves**. It's also known as Chinese water spinach, Chinese water spinach, or hollow stem spinach.

The Seasoned Wok
Ong Choy (Water Spinach) Recipe with Fermented Bean...
Nov 3, 2022 — What is Ong Choy, Rau Muong or Water Spinach? Ong choy i...

The Woks of Life
Ong Choy with XO sauce - The Woks of Life
May 2, 2017 — Ong Choy is a popular Chinese leafy green vegetable that's...

Onolicious Hawai'i
Garlic and Fish Sau... - Onolicious Hawai...
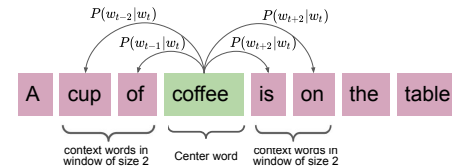Jan 7, 2021 — What is O... Hawaii everyone know...

## Word2Vec

- We want vectors for words so that the context of a word can suggest the vector of this word, and vice versa
- Idea: **Similar words appear in similar contexts**

> A cup of **coffee** is on the table.
> **Coffee** helps me focus.
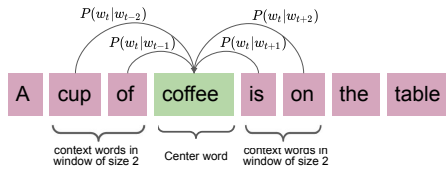> Espresso is my favorite type of **coffee**.
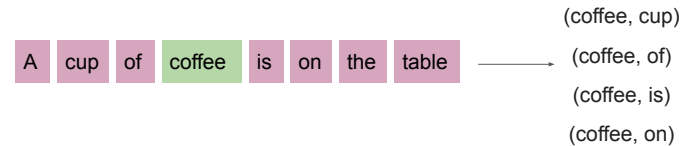
---

## Word2Vec - Training

SkipGram - Predict context from target



$P(w_{t-2}|w_t)$    $P(w_{t+2}|w_t)$
$P(w_{t-1}|w_t)$    $P(w_{t+1}|w_t)$

| A | cup | of | coffee | is | on | the | table |

context words in window of size 2 — Center word — context words in window of size 2

---

## Word2Vec - Training

Continuous Bag of Words (CBOW) - predict target from context



$P(w_t|w_{t-2})$    $P(w_t|w_{t+2})$
$P(w_t|w_{t-1})$    $P(w_t|w_{t+1})$

| A | cup | of | coffee | is | on | the | table |

context words in window of size 2 — Center word — context words in window of size 2

---

## SkipGram - Training samples

| A | cup | of | coffee | is | on | the | table |

(coffee, cup)
(coffee, of)
(coffee, is)
(coffee, on)

## Word2Vec Architecture - SkipGram

Predict every target word from each context word!

coffee → cup

## Discuss: Word2Vec Architecture - CBOW

d = embedding size

What is the output of multiplying
the one-hot vector
[0,1,0,0,0,0,0] with **W**?

|V| = vocab size

| 0.1 | 0.5 | 0.3 | -0.9 | 0.4 |
| 0.8 | -0.4 | 0.7 | 0.3 | -0.1 |
| 0.4 | 0.3 | -0.9 | -0.2 | 0.7 |
| -0.5 | -0.1 | 0.7 | 0.8 | 0.6 |
| -0.9 | 0.6 | -0.5 | 0.6 | -0.2 |
| 0.2 | 0.8 | 0.6 | 0.3 | 0.6 |
| 0.2 | 0.2 | -0.9 | -0.5 | 0.3 |

**Matrix W (learnt from training)**

## Word2Vec Architecture - SkipGram

Predict every target word from each context word!

cup → coffee

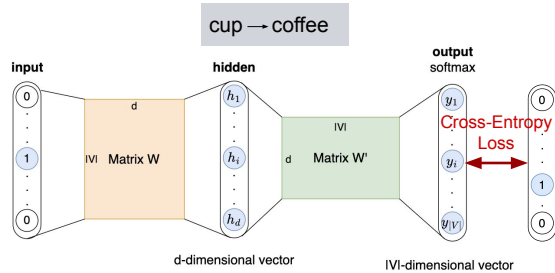

softmax of |V| dimensional output to get probabilities

## Looking closer…

● We observe that every row of the W matrix corresponds to a target word and every column of the W' matrix corresponds to a context word.

● We compute the probability of a target-context pair as:

$$p(w_c|w_t) = \frac{exp(W_t {W'}_c^T)}{\sum_{i=1}^{|V|} exp(W_t {W'}_i^T)}$$

## Word2Vec Architecture - SkipGram

Predict every target word from each context word!

cup → coffee



input    hidden    output softmax

d-dimensional vector    |V|-dimensional vector
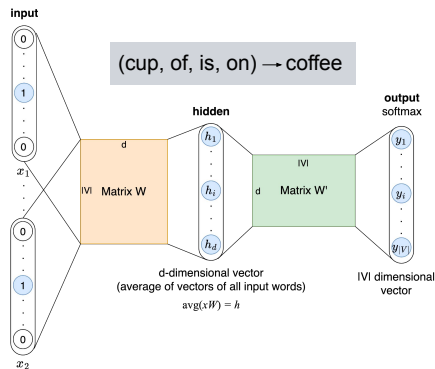
Cross-Entropy Loss

## Cross Entropy

- Cross Entropy: lower cross entropy indicates high similarity between two distributions

$$\mathcal{L}_\theta = -\sum_{i=1}^{|V|} y_i \log p(w_i|w_t) = -log\, p(w_c|w_t)$$
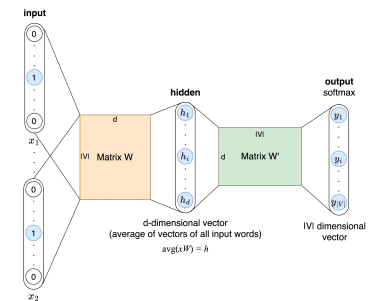
- So the loss function is:

$$\mathcal{L}_\theta = -log\frac{exp(W_t W'^T_c)}{\sum_{i=1}^{|V|} exp(W_t W'^T_i)}$$

## Word2Vec Architecture - CBOW (continuous bag of words)

input

(cup, of, is, on) → coffee



hidden    output softmax

d-dimensional vector
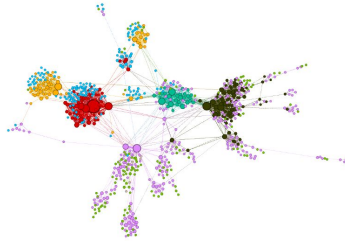(average of vectors of all input words)
avg(xW) = h

|V| dimensional vector

## Word2Vec Architecture - CBOW (continuous bag of words)

Where do we get the word embeddings from in this version of Word2Vec (CBOW)?



input    hidden    output softmax

d-dimensional vector
(average of vectors of all input words)
avg(xW) = h

|V| dimensional vector

## Slide 1

### X 2 vec

- Generate vector representations (embeddings) for various data types
- Examples:
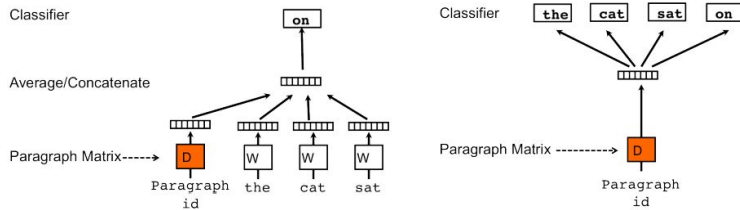  - Word2Vec
  - Doc2Vec
  - Node2Vec
  - Item2Vec
  - Sent2Vec

## Slide 2

### Demo

Visualize: https://projector.tensorflow.org/

Explore: http://epsilon-it.utu.fi/wv_demo/

## Slide 3

### Doc2Vec

- A vector to represent a paragraph, regardless of length
  - embeddings for paragraph and words
  - Applications: Document classification, sentiment analysis, recommendation systems, and information retrieval

Classifier — on

Average/Concatenate

Paragraph Matrix ------> D    W    W    W

Paragraph id    the    cat    sat

Classifier — the   cat   sat   on

Paragraph Matrix ---------> D

Paragraph id

## Slide 4

### In vector space…

## Word embeddings are time-dependent (why?)
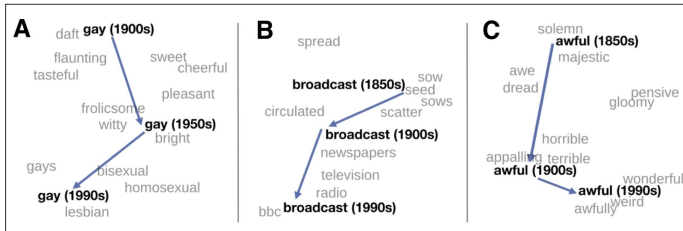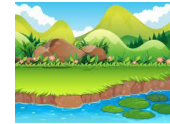
- Semantic similarity of words depends on *time.*



A
daft **gay (1900s)**
flaunting sweet cheerful
tasteful
pleasant
frolicsome
witty **gay (1950s)**
bright
gays
bisexual
homosexual
**gay (1990s)**
lesbian

B
spread
**broadcast (1850s)** sow
seed
circulated sows
scatter
**broadcast (1900s)**
newspapers
television
radio
bbc **broadcast (1990s)**

C
solemn
**awful (1850s)**
majestic
awe
dread
pensive
gloomy
horrible
appalling terrible
**awful (1900s)**
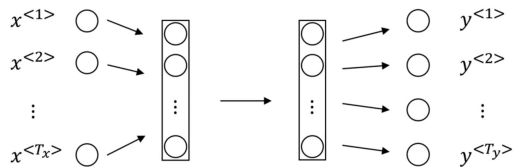wonderful
**awful (1990s)** weird
awfully

---

## Problems with word2vec

- Words with multiple meanings only have one representation
  - eg. **bank** of river or **bank** of money
  - Need contextual information
- Limited Context
  - only trained on words within the context window



---

## How to use word vectors with neural networks?

$x^{<1>}$
$x^{<2>}$
⋮
$x^{<T_x>}$

$y^{<1>}$
$y^{<2>}$
⋮
$y^{<T_y>}$
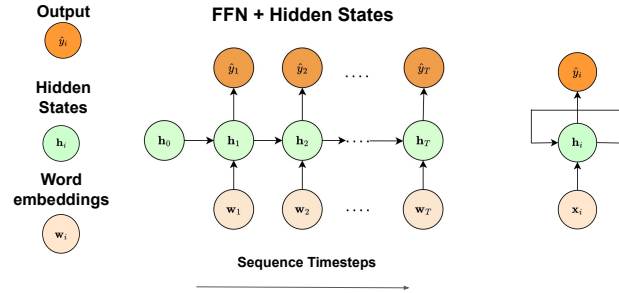
- Inputs and outputs don't have fixed lengths
- Features are not shared

---

## Let's simplify!

What if we have a single word and a single output?

$x_i^0$
⋮
$x_i^n$
$\hat{y}_i$

$\hat{y}_i$
$\mathbf{h}_i$
$\mathbf{x}_i$

## Recurrent neural network (RNN)

**Output**
$\hat{y}_i$

**Hidden States**
$\mathbf{h}_i$

**Word embeddings**
$\mathbf{w}_i$

**FFN + Hidden States**

$\hat{y}_1$ $\hat{y}_2$ .... $\hat{y}_T$ $\hat{y}_i$

$\mathbf{h}_0 \to \mathbf{h}_1 \to \mathbf{h}_2 \to \cdots \to \mathbf{h}_T$ $\mathbf{h}_i$

$\mathbf{w}_1$ $\mathbf{w}_2$ .... $\mathbf{w}_T$ $\mathbf{x}_i$

**Sequence Timesteps**

---

## Parameterize RNN

$\hat{y}_1$ $\hat{y}_2$ .... $\hat{y}_T$

$\mathbf{w}^1_{h,\hat{y}}$ $\mathbf{w}^2_{h,\hat{y}}$ $\mathbf{w}^T_{h,\hat{y}}$

$\mathbf{h}_0 \xrightarrow{\mathbf{w}^0_{h,h}} \mathbf{h}_1 \xrightarrow{\mathbf{w}^1_{h,h}} \mathbf{h}_2 \xrightarrow{\mathbf{w}^2_{h,h}} \cdots \xrightarrow{\mathbf{w}^T_{h,h}} \mathbf{h}_T$

$\mathbf{w}^1_{x,h}$ $\mathbf{w}^2_{x,h}$ $\mathbf{w}^T_{x,h}$
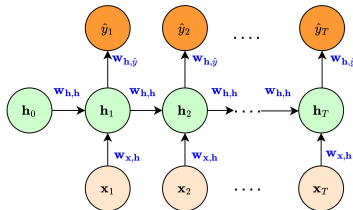
$\mathbf{x}_1$ $\mathbf{x}_2$ .... $\mathbf{x}_T$

- Too many parameters if we have a long sequence!
- Longer sequence parameters will not receive many updates

---

## RNN w/ parameter-sharing

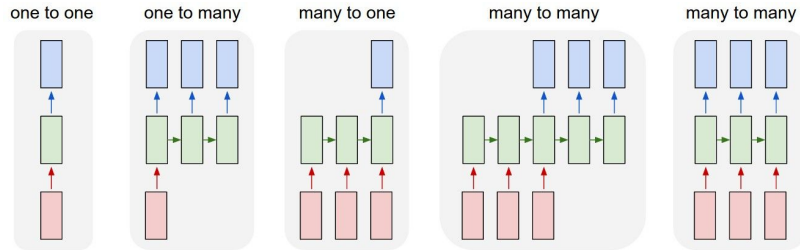Simple fix: use **the same parameters** across different timesteps.

$\hat{y}_1$ $\hat{y}_2$ .... $\hat{y}_T$

$\mathbf{w}_{h,\hat{y}}$ $\mathbf{w}_{h,\hat{y}}$ $\mathbf{w}_{h,\hat{y}}$

$\mathbf{h}_0 \xrightarrow{\mathbf{w}_{h,h}} \mathbf{h}_1 \xrightarrow{\mathbf{w}_{h,h}} \mathbf{h}_2 \xrightarrow{\mathbf{w}_{h,h}} \cdots \xrightarrow{\mathbf{w}_{h,h}} \mathbf{h}_T$

$\mathbf{w}_{x,h}$ $\mathbf{w}_{x,h}$ $\mathbf{w}_{x,h}$

$\mathbf{x}_1$ $\mathbf{x}_2$ .... $\mathbf{x}_T$

---

## Discuss: RNN w/ parameter-sharing

Simple fix: use **the same parameters** across different timesteps.

$\hat{y}_1$ $\hat{y}_2$ .... $\hat{y}_T$

$\mathbf{w}_{h,\hat{y}}$ $\mathbf{w}_{h,\hat{y}}$ $\mathbf{w}_{h,\hat{y}}$

$\mathbf{h}_0 \xrightarrow{\mathbf{w}_{h,h}} \mathbf{h}_1 \xrightarrow{\mathbf{w}_{h,h}} \mathbf{h}_2 \xrightarrow{\mathbf{w}_{h,h}} \cdots \xrightarrow{\mathbf{w}_{h,h}} \mathbf{h}_T$

$\mathbf{w}_{x,h}$ $\mathbf{w}_{x,h}$ $\mathbf{w}_{x,h}$

$\mathbf{x}_1$ $\mathbf{x}_2$ .... $\mathbf{x}_T$

A non-linearity is applied to the output of the recurrent unit before it is passed to the next time step or to the output layer of the network.

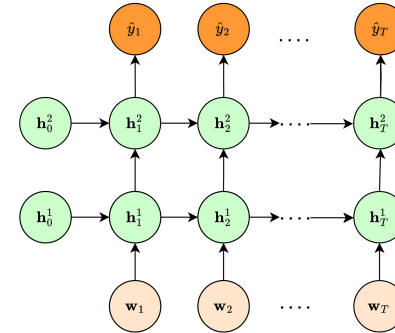Write a closed-form solution for $\mathbf{h}_i$ and $\hat{y}_i$

## Types of RNNs



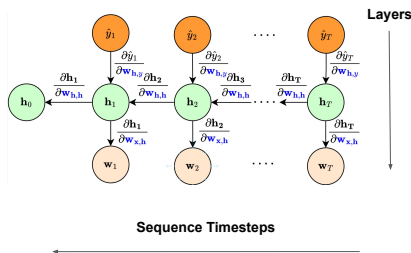| one to one | one to many | many to one | many to many | many to many |

https://www.analyticsvidhya.com/blog/2021/06/time-series-analysis-recurrence-neural-network-in-python/
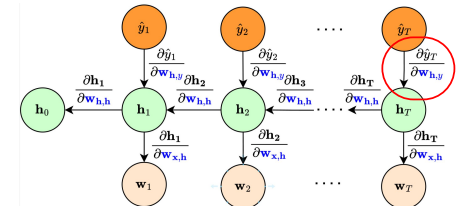
## Stacking RNN Layers

## Backpropagation through the Time (BPTT)



**Layers**

- Unfold a recurrent neural network in time
- Gradients are accumulated across all time steps by applying the chain rule
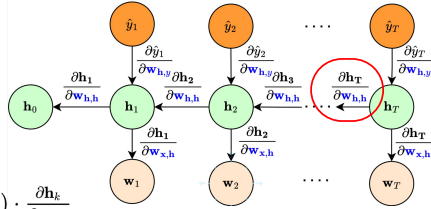- Propagate gradients backwards through time steps

**Sequence Timesteps**

## Backpropagation through the Time (BPTT)

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_{\mathbf{h},\hat{y}}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_T} \cdot \frac{\partial \hat{y}_T}{\partial \mathbf{w}_{\mathbf{h},\hat{y}}}$$

## Backpropagation through the Time (BPTT)

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w_{h,h}}} = \frac{\partial \mathcal{L}}{\partial \hat{y}_T} \cdot \frac{\partial \hat{y}_T}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \cdot \frac{\partial \mathbf{h}_{T-1}}{\partial \mathbf{w}_{h,h}}$$

$$= \sum_{k=1}^{T} \frac{\partial \mathcal{L}}{\partial \hat{y}_T} \cdot \frac{\partial \hat{y}_T}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_k} \cdot \frac{\partial \mathbf{h}_k}{\partial \mathbf{w}_{h,h}}$$

$$= \sum_{k=1}^{T} \frac{\partial \mathcal{L}}{\partial \hat{y}_T} \cdot \frac{\partial \hat{y}_T}{\partial \mathbf{h}_T} \cdot \left( \prod_{j=k}^{T-1} \frac{\partial \mathbf{h}_{j+1}}{\partial \mathbf{h}_j} \right) \cdot \frac{\partial \mathbf{h}_k}{\partial \mathbf{w}_{h,h}}$$

## Recap

- **N-gram models**
- **Bag-of-words representations**
- **Word2Vec**
  - CBOW: use context to predict target word
  - SkipGram: use target word to predict context
- **RNN**
  - Has an internal state (memory)
  - Can handle arbitrary sequences of inputs
  - Trained with back propagation through time

## Image credits:

https://web.stanford.edu/~jurafsky/slp3/6.pdf

https://lilianweng.github.io/posts/2017-10-15-word-embedding/