# Midterm Review

CS224N/Ling284: Natural Language Processing with Deep Learning
Feb 8, 2018

# Midterm

- Feb 13, 4:30-5:50, Memorial Auditorium
- Alternate exam: Feb 9(tomorrow), 4:00 - 5:20 pm, 200-303 (Lane History Corner)
- One cheatsheet allowed (letter sized, double-sided)
- Bring a pencil and eraser to the midterm
- Covers all the lectures so far
- Approximate questions breakdown:
  - multiple choice and true false
  - short answers
  - more involved questions
- SCPD: Either show up or have an exam monitor pre-registered with SCPD!!

# Review Outline

- Word Vector Representations
- Neural Networks
- Backpropagation / Gradient Calculation
- Dependency Parsing
- RNNs

# Word Vector Representations
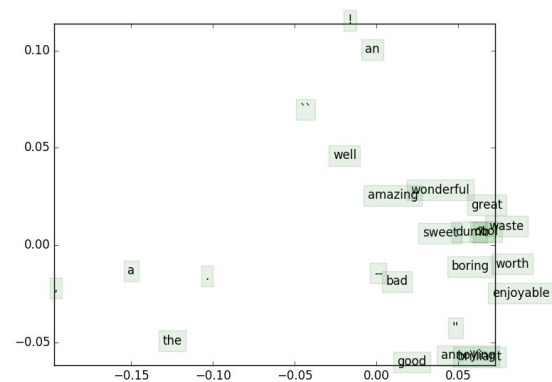
CS224N Midterm Review

Michelle Mei

# Word Vectors

Definition: A vector that captures the meaning of a word.

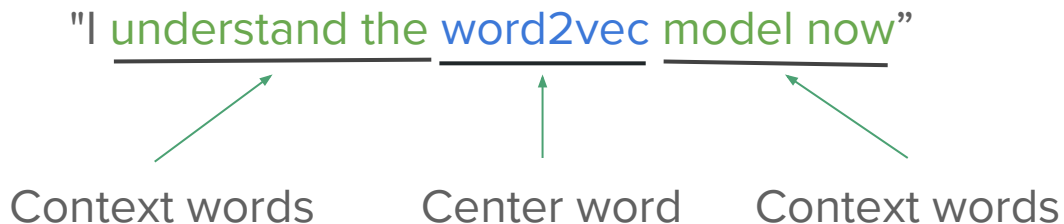Sometimes can also be called as word embeddings or word representations.

We will be reviewing: Word2Vec and GloVe.

# Word2Vec

**Task:** Learn word vectors to encode the probability of a word given its context.

Consider the following example with context window size = 2:

"I understand the word2vec model now"

Context words        Center word        Context words

# Word2Vec

**Task:** Learn word vectors to encode the probability of a word given its context.

For each word, we want to learn 2 vectors:

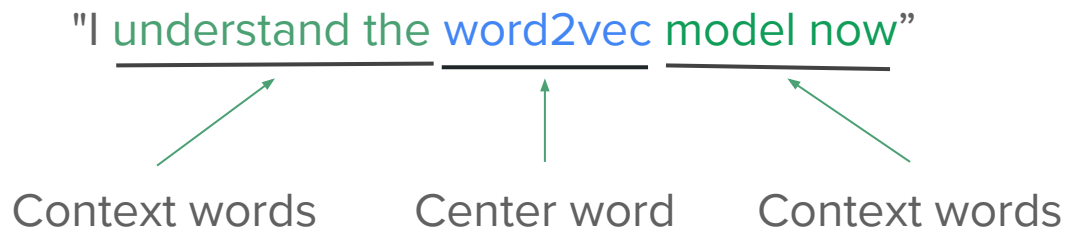$v$ : **input vector**          $u$ : **output vector**

**Two algorithms:**

- **Skipgram:** predicts the probability of context words from a center word.
- **Continuous Bag-of-Words (CBOW):** predicts a center word from the surrounding context in terms of word vectors.
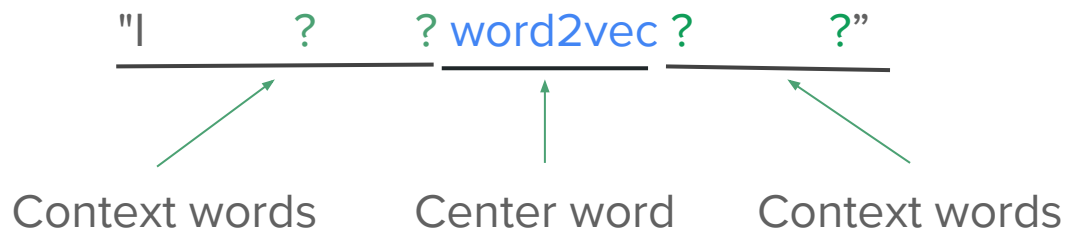
# Word2Vec - Skipgram

- Predicts the probability of context words from a center word.

- Let's look at the previous example again:

"I understand the word2vec model now"

Context words    Center word    Context words

# Word2Vec - Skipgram

- Predicts the probability of context words from a center word.

- Let's look at the previous example again:

"I        ?        ? word2vec ?        ?"

Context words        Center word        Context words

# Word2Vec - Skipgram

"I _____ ? _____ ? word2vec ? _____ ?"

- Generate a one-hot vector, $\mathbf{w}_c$ of the center word, "word2vec". It is a |VocabSize|-dim vector with a 1 at the word index and 0 elsewhere.
- Look up the input vector, $\mathbf{v}_c$ in $\mathbf{V}$ using $\mathbf{w}_c$. $\mathbf{V}$ is the input vector matrix.
- Generate a score vector, $\mathbf{z} = \mathbf{U}\mathbf{v}_c$ where $\mathbf{U}$ is the output vector matrix.
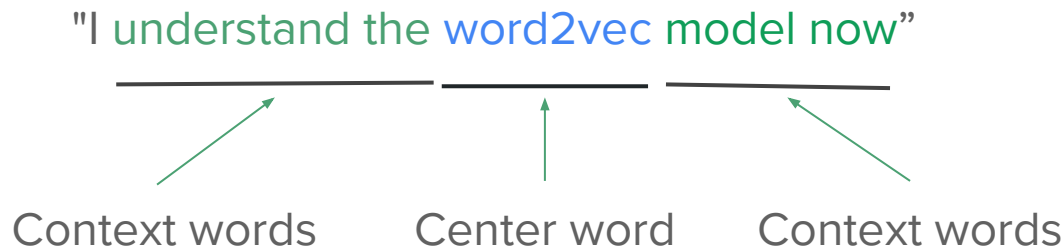
# Word2Vec - Skipgram

"I _____ ? _ ? word2vec ? _____ ?"

- Turn the score vector into probabilities, $\hat{\mathbf{y}}$ = **softmax(z)**.

- $[\hat{\mathbf{y}}_{c-m}, \dots, \hat{\mathbf{y}}_{c-1}, \hat{\mathbf{y}}_{c+1}, \dots, \hat{\mathbf{y}}_{c+m}]$ : probabilities of observing each context word (**m** is the context window size)

- Minimize cost given by: (**F** can be neg-sample or softmax-CE)

$$J_{\text{skip-gram}}(\text{word}_{c-m\dots c+m}) = \sum_{-m \leq j \leq m, j \neq 0} F(\boldsymbol{w}_{c+j}, \boldsymbol{v}_c)$$

# Word2Vec - Continuous Bag-Of-Words (CBOW)

● Predicts a center word from the surrounding context
● Let's look at the previous example again:

"I understand the word2vec model now"
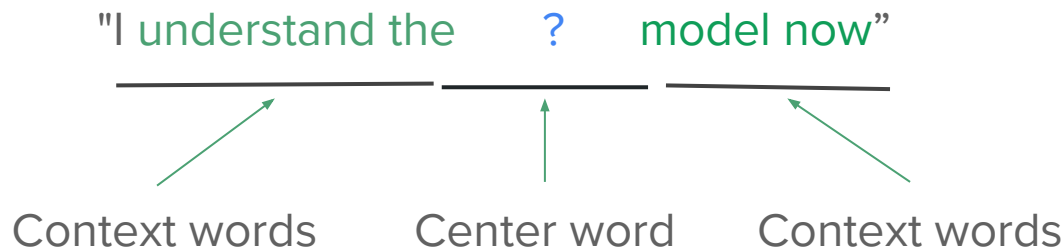
Context words　　　Center word　　　Context words

# Word2Vec - Continuous Bag-Of-Words (CBOW)

- Predicts a center word from the surrounding context
- Let's look at the previous example again:

"I understand the      ?      model now"

Context words          Center word          Context words

# Word2Vec - Continuous Bag-Of-Words (CBOW)

"I understand the    ?    model now"

_____  _____  _____

- Generate one-hot vectors, $w_{c-m}, \ldots, w_{c-1}, w_{c+1}, \ldots, w_{c+m}$ for the context words.

- Look up the input vectors, $v_{c-m}, \ldots, v_{c-1}, v_{c+1}, \ldots, v_{c+m}$ in **V** using the one-hot vectors. **V** is the input vector matrix.

- Average these vectors to get $v_{avg} = (v_{c-m} + \ldots + v_{c-1} + v_{c+1} + \ldots + v_{c+m})/2m$

# Word2Vec - Continuous Bag-Of-Words (CBOW)

"I understand the     ?     model now"

‾‾‾‾‾‾‾‾‾‾‾ ‾‾‾‾‾‾‾ ‾‾‾‾‾‾‾

- Generate a score vector, $\mathbf{z = Uv_{avg}}$ where $\mathbf{U}$ is the output vector matrix.
- Turn the score vector into probabilities, $\mathbf{\hat{y} = softmax(z)}$.
- $\mathbf{\hat{y}}$: probability of the center word.
- Minimize cost given by: ($\mathbf{F}$ can be neg-sample or softmax-CE)

$$J_{CBOW}(word_{c-m...c+m}) = F(w_c, v_{avg})$$

# GloVe

- Like Word2Vec, GloVe is a set of vectors that capture the semantic information (i.e. meaning) about words.
- Unlike Word2Vec, Glove makes use of global co-occurrence statistics.
- Fast Training
- Scalable to huge corpora
- Good Performance even with small corpus and small vectors

"GloVe consists of a weighted least squares model that trains on global word-word co-occurrence counts."

# GloVe

Co-occurrence Matrix (window-based):

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 |

# GloVe

- Let **X** be the word-word co-occurrence counts matrix.
  - $X_i$ is the number of times any word **k** appears in the context of word **i**.
  - $X_{ij}$ is the number of times word **j** occurs in the context of word **i**.
- Like the case in Word2Vec, each word has 2 vectors, **input (v)** and **output (u)**.
- The cost function:
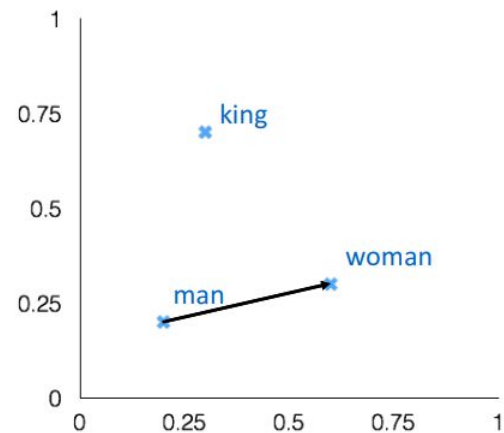
$$\hat{J} = \sum_{i=1}^{W} \sum_{j=1}^{W} X_i (\vec{u}_j^T \vec{v}_i - \log X_{ij})^2$$

# GloVe

- In the end, we have **V** and **U** from all the input and output vectors, **v** and **u.**

- Both capture similar co-occurrence information, and so the word vector for a word can be simply obtained by summing **u** and **v** up!

# Evaluate Word Vectors

- Intrinsic Method

  - Word Vector Analogies: Evaluate word vectors by how well their cosine distance after addition captures intuitive semantic and syntactic analogy questions.

- Extrinsic Method

  - Entity recognition

# Neural Networks

CS224N Midterm Review

Amani V. Peddada

# Loss Functions

- Prediction of category or label (classification)

  - *Softmax + Cross-Entropy Loss*: optimize correct class probability

  $$\text{softmax}(\theta)_i = \hat{y}_i = \frac{e^{\theta_i}}{\sum_{j=1}^{C} e^{\theta_j}} \quad \text{CE}(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{i=1}^{C} y_i \log(\hat{y}_i)$$

  - *Max-Margin Loss:* optimize margin between correct class score and incorrect class scores.

  $$J = \max(0, 1 - s + s_c)$$

- Prediction of real values or continuous outputs (regression)

  - L2 Loss: $L_2(y, \theta) = ||\theta - y||_2^2$

  - Others: L1, etc.

# Network Structure

- Recall: forward pass of a neural network. Hidden layers computed as follows:
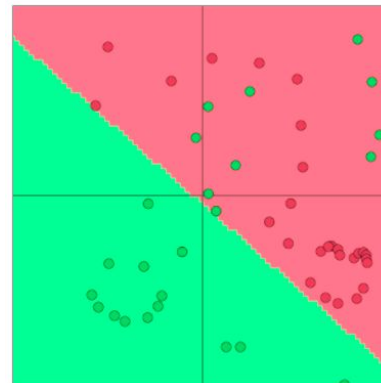
$$\mathbf{h_1} = f(\mathbf{W_1 x + b_1})$$
$$\mathbf{h_2} = f(\mathbf{W_2 h_1 + b_2})$$
$$\cdots$$
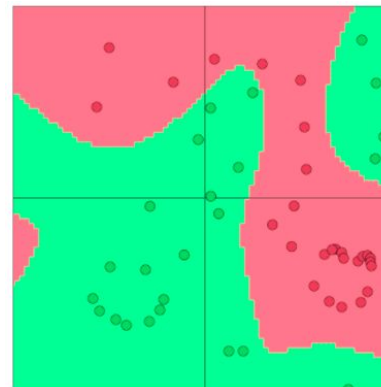$$\mathbf{h_i} = f(\mathbf{W_i h_{i-1} + b_i})$$

- Number of hidden layers/size of each hidden layer affects representational power. More parameters => more expressive model.

- Initialization is important
  - Small random numbers (e.g. Xavier/Glorot) for weight matrices.
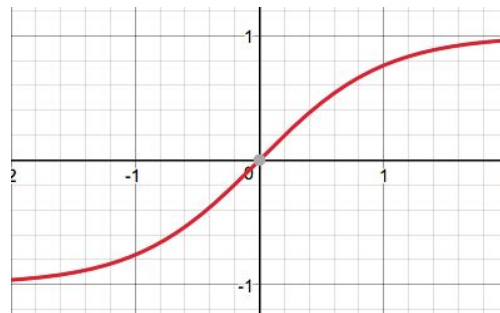
Linear Model
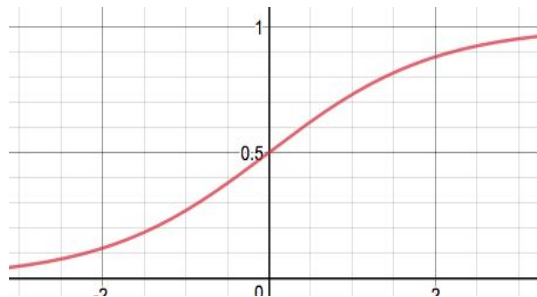


Multilayer Neural Network

# Non-Linearities

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

$$\tanh(x) = 2\sigma(2x) - 1$$

$$\mathrm{ReLU}(x) = \max(0, x)$$

$$\mathrm{LeakyReLU}(x) = \max(\alpha x, x)$$

- Responsible for network's expressiveness -- otherwise just a linear model
- Beware of saturation and "dead" neurons
- Other variants: PreLu, Maxout, Hard Tanh

# Gradient Check

- Used to verify correctness of the **analytic gradient**

- Compute **numerical gradient** using the *central difference formula*:

$$\frac{\partial f}{\partial x} \approx \frac{f(x+h) - f(x-h)}{2h}$$

- Vary one dimension of parameters at a time, observe change in output function (loss)

- Potentially very expensive to compute over large numbers of parameters; can sanity check by checking only a few dimensions a time

# Optimization

- Optimize loss function with Gradient Descent to compute parameter updates:

$$\theta^{new} = \theta^{old} - \alpha \nabla_\theta J(\theta)$$

- Taking gradient over entire training set is expensive, so use mini-batches (Stochastic Gradient Descent)

- In addition to SGD, there are more complicated updates: *Adam* (see PA2), *AdaGrad*, *RMSProp*, *Nesterov Momentum,* etc.

- Sanity check: If network is working properly, should be able to get close to 0 loss on small subset of training data.

- May be helpful to randomize order of examples

# Monitoring Learning Curves

- Plot training loss as a function of iteration/time.

- Adjusting learning rate
  - Training loss increases => learning rate too high
  - Training loss plateaus at high value => learning rate too high
  - Linear decrease in training loss => learning rate too low
  - May be helpful to *anneal* learning rate over time

# Monitoring Learning Curves

- Also should compare training and validation loss/accuracies

- Large gap => **Overfitting**: Model does not generalize well to unseen data

- Bad training performance=> **Underfitting**: Model is not powerful enough to learn the training data, resulting in bad performance on both training and validation datasets.

- Note: do not compare to test set, which is reserved for final evaluation.

# Handling Overfitting

- Add Dropout
  - Constrain each neuron to learn more meaningful information.
  - Can also be interpreted as an "ensemble" of smaller networks.
  - Need to scale activations to maintain expected value (see PA2)
- L2 Regularization
  - Add $+\lambda||\theta||_2^2$ with tunable lambda for non-bias parameters
  - Encourages weights to be more spread out, place less emphasis on any one input dimension
- Reduce Network depth/size
- Reduce input feature dimensionality
- Early Stopping
- Others: Max-Norm, L1 penalty, DropConnect etc.



(a) Standard Neural Net      (b) After applying dropout.

# Handling Underfitting

- Increase model complexity/size
- Decreasing regularization effects
- Reducing L2 penalty weight
- Reducing Dropout probability
- Usually opposite of overfitting solutions

# Other Helpful Techniques

- Ensembling
  - Combine separately trained models for more robust predictions
- Data Preprocessing
  - Mean-centering data
- Batch Normalization
  - Encourage outputs after hidden layer to have zero mean, unit variance
- Curriculum Learning
  - During training, present examples in a certain order to speed up optimization
- Data Augmentation
  - Can augment training set with additional examples by applying transformations to input

# Backpropagation / Gradient Calculation

CS224N Midterm Review

Huseyin A. Inan

# Matrix Calculus Primer

Scalar-by-Vector

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \dfrac{\partial y}{\partial x_1} & \dfrac{\partial y}{\partial x_2} & \cdots & \dfrac{\partial y}{\partial x_n} \end{bmatrix}$$

(We can transpose it to convert it to column shape)

Vector-by-Vector

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \dfrac{\partial y_1}{\partial x_1} & \dfrac{\partial y_1}{\partial x_2} & \cdots & \dfrac{\partial y_1}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \dfrac{\partial y_m}{\partial x_1} & \dfrac{\partial y_m}{\partial x_2} & \cdots & \dfrac{\partial y_m}{\partial x_n} \end{bmatrix}$$

Scalar-by-Matrix

$$\frac{\partial y}{\partial A} = \begin{bmatrix} \dfrac{\partial y}{\partial A_{11}} & \dfrac{\partial y}{\partial A_{12}} & \cdots & \dfrac{\partial y}{\partial A_{1n}} \\ \vdots & \vdots & \ddots & \vdots \\ \dfrac{\partial y}{\partial A_{m1}} & \dfrac{\partial y}{\partial A_{m2}} & \cdots & \dfrac{\partial y}{\partial A_{mn}} \end{bmatrix}$$

# Backpropagation Shape Rule and Dimension Balancing

The gradient at each intermediate step has **shape of denominator**

$$X \in \mathbb{R}^{m \times n} \iff \delta_X = \frac{\delta Scalar}{\delta X} \in \mathbb{R}^{m \times n}$$

- **Dimension balancing** is the "cheap" but efficient way to calculate gradients in most practical settings
- Read **gradient computation notes** to understand how to derive matrix expressions for gradients from first principles
- Dimension balancing approach should be used with a good understanding of what is happening behind it

$$z = Wx$$

$$\frac{\partial z}{\partial x} = W$$

$$z = xW$$

$$\frac{\partial z}{\partial x} = W^T$$

$$z = Wx \qquad \frac{\partial J}{\partial z} = \delta$$

$$\frac{\partial J}{\partial W} = \delta x^T$$

$$z = xW \qquad \frac{\partial J}{\partial z} = \delta$$

$$\frac{\partial J}{\partial W} = x^T \delta$$

# Activation Function

$\boldsymbol{h} = f(\boldsymbol{z})$, what is $\dfrac{\partial \boldsymbol{h}}{\partial \boldsymbol{z}}$? $\qquad\qquad \boldsymbol{h}, \boldsymbol{z} \in \mathbb{R}^n$

$h_i = f(z_i)$

$$\frac{\partial \boldsymbol{h}}{\partial \boldsymbol{z}} = \begin{pmatrix} f'(z_1) & & 0 \\ & \ddots & \\ 0 & & f'(z_n) \end{pmatrix} = \operatorname{diag}(\boldsymbol{f}'(\boldsymbol{z})) \qquad \boldsymbol{f}'(\boldsymbol{z}) = [f'(z_1), f'(z_2), \ldots, f'(z_n)]$$

$$\frac{\partial \mathcal{J}}{\partial \boldsymbol{h}} = \delta_h \qquad\Longrightarrow\qquad \frac{\partial \mathcal{J}}{\partial \boldsymbol{z}} = \frac{\partial \mathcal{J}}{\partial \boldsymbol{h}} \frac{\partial \boldsymbol{h}}{\partial \boldsymbol{z}} = \delta_h \operatorname{diag}(\boldsymbol{f}'(\boldsymbol{z})) = \delta_h \circ \boldsymbol{f}'(\boldsymbol{z})$$

# Backpropagation

$$h_1 = \sigma(xW_1 + b_1)$$

$$\hat{y} = softmax(h_1 W_2 + b_2)$$

$$J = CE(\hat{y}, y)$$

$$x \in \mathbb{R}^{D_x}$$

$$W_1 \in \mathbb{R}^{D_x \times D_z}$$

$$b_1 \in \mathbb{R}^{D_z}$$

$$h_1 \in \mathbb{R}^{D_z}$$

$$W_2 \in \mathbb{R}^{D_z \times D_y}$$

$$b_2 \in \mathbb{R}^{D_y}$$

1. Identify intermediate functions (forward prop)

2. Compute local gradients

3. Combine with downstream error signal to get full gradient

# Backpropagation

Loss Function:

$$h_1 = \sigma(xW_1 + b_1)$$

$$\hat{y} = softmax(h_1 W_2 + b_2)$$

$$J = CE(\hat{y}, y)$$

Intermediate Variables:
(forward propagation)

$$z_1 = xW_1 + b_1$$

$$h_1 = \sigma(z_1)$$

$$\theta = h_1 W_2 + b_2$$

$$\hat{y} = softmax(\theta)$$

$$J = CE(\hat{y}, y)$$

Intermediate Variables:
(forward propagation)

$$z_1 = xW_1 + b_1$$
$$h_1 = \sigma(z_1)$$
$$\theta = h_1 W_2 + b_2$$
$$\hat{y} = softmax(\theta)$$
$$J = CE(\hat{y}, y)$$

$$x \in \mathbb{R}^{D_x}$$
$$W_1 \in \mathbb{R}^{D_x \times D_z}$$
$$b_1 \in \mathbb{R}^{D_z}$$
$$h_1 \in \mathbb{R}^{D_z}$$
$$W_2 \in \mathbb{R}^{D_z \times D_y}$$
$$b_2 \in \mathbb{R}^{D_y}$$

Let's do it for x first:

$$\frac{\partial J}{\partial x} = \frac{\partial J}{\partial z_1}\frac{\partial z_1}{\partial x} = \delta_2 W_1^T \in \mathbb{R}^{D_x}$$
$$\delta_2 = \delta_1 \circ h_1 \circ (1 - h_1) \in \mathbb{R}^{D_z}$$
$$\frac{\partial J}{\partial z_1} = \frac{\partial J}{\partial h_1}\frac{\partial h_1}{\partial z_1} = \delta_1 \circ \sigma'(z_1)$$
$$\delta_1 = (\hat{y} - y)W_2^T \in \mathbb{R}^{D_z}$$
$$\frac{\partial J}{\partial h_1} = (\hat{y} - y)W_2^T$$
$$\frac{\partial J}{\partial \theta} = \hat{y} - y$$

Intermediate Variables:
(forward propagation)

$$z_1 = xW_1 + b_1$$
$$h_1 = \sigma(z_1)$$
$$\theta = h_1W_2 + b_2$$
$$\hat{y} = softmax(\theta)$$
$$J = CE(\hat{y}, y)$$

$$x \in \mathbb{R}^{D_x}$$
$$W_1 \in \mathbb{R}^{D_x \times D_z}$$
$$b_1 \in \mathbb{R}^{D_z}$$
$$h_1 \in \mathbb{R}^{D_z}$$
$$W_2 \in \mathbb{R}^{D_z \times D_y}$$
$$b_2 \in \mathbb{R}^{D_y}$$

Let's continue with: $W_2, b_2, W_1, b_1$

$$\frac{\partial J}{\partial b_1} = \frac{\partial J}{\partial z_1}\frac{\partial z_1}{\partial b_1} = \delta_2$$
$$\frac{\partial J}{\partial W_1} = x^T\delta_2 \in \mathbb{R}^{D_x \times D_z}$$
$$\frac{\partial J}{\partial z_1} = \delta_2$$
$$\frac{\partial J}{\partial b_2} = \frac{\partial J}{\partial \theta}\frac{\partial \theta}{\partial b_2} = (\hat{y} - y)$$
$$\frac{\partial J}{\partial W_2} = h_1^T(\hat{y} - y) \in \mathbb{R}^{D_z \times D_y}$$
$$\frac{\partial J}{\partial \theta} = \hat{y} - y$$

Here is one such model to evaluate how similar two input words are using Euclidean distance. There are two input word vectors $x_1$, $x_2 \in \mathbb{R}^n$, shared parameters $W \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$, and a single hidden layer associated with *each* input:

$$h_1 = \sigma(W x_1 + b)$$

$$h_2 = \sigma(W x_2 + b)$$

We evaluate the distance between the two activations $h_1$, $h_2$ using Euclidean distance as our similarity metric. The model objective $J$ is

$$J = \frac{1}{2}\|h_1 - h_2\|_F^2 + \frac{\lambda}{2}\|W\|_F^2$$

where $\lambda$ is a given regularization parameter. (The Frobenius norm $\|.\|_F$ is a matrix norm defined by $\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} |A_{ij}|^2}$)

Calculate $\nabla_W J$ and $\nabla_b J$.

Our model is:

$$h_1 = \sigma(W_1 x + b_1)$$

$$h_2 = \text{relu}(W_2 x + b_2)$$

$$\hat{y} = \text{softmax}(W_3(h_1 + h_2) + b_3)$$

where $x \in \mathbb{R}^n$, $W_1, W_2 \in \mathbb{R}^{m \times n}$, $W_3 \in \mathbb{R}^{k \times m}$, $b_1, b_2 \in \mathbb{R}^m$, and $b_3 \in \mathbb{R}^k$. We evaluate this model for $N$ examples and $k$ classes with cross entropy loss

$$J = -\frac{1}{N} \sum_{j=1}^{N} \sum_{i=1}^{k} y_j^i log(\hat{y}_j^i)$$

where $y_j$ is the one-hot vector for example $j$ with all probability mass on the correct class and $\hat{y}_j$ are the softmax scores for example $j$.

Find $\nabla_{h_1} J$, $\nabla_{h_2} J$, and $\nabla_x J$.

# Summary

- Identify intermediate functions (forward prop)

- Compute local gradients from top to bottom

- Use Dimension Balancing to double check ( or use it to achieve the final result in "hacky" way :) )

# Dependency Parsing

CS224N Midterm Review

Rex Ying

# Two views of Linguistic Structure



Constituency Structure uses phrase structure grammar to organize words into nested constituents.



Dependency Structure uses dependency grammar to identify which words depend on which other words (and how).

# Dependency Parsing

- Asymmetric relations between words (head of the dependency to the dependent).
- Typed with the name of the grammatical relation.
- Usually forms a connected, single-head tree.
- Ambiguities exist

# Greedy deterministic transition based parsing

- Bottom up actions analogous to shift-reduce parser
- States defined as a triple of words in buffer, words in stack and set of parsed dependencies.
- Discriminative classification
- Evaluation metrics: UAS, LAS
- MaltParser



Left Arc

[root] | I | ate → [root] | ate     A += nsubj(ate → I)

Shift

[root] | ate | fish → [root] | ate | fish

Right Arc

[root] | ate | fish → [root] | ate     A += obj(ate → fish)

# Projectivity

Projective arcs have no crossing arcs when the words are laid in linear order.

However, some sentences have non-projective dependency structure

# Handling non-projectivity

- Declare defeat
- Use post-processor to identify and resolve these non-projective dependencies
- Add extra transitions
- Use a parsing mechanism that doesn't have projectivity constraint.

# Neural Dependency Parsing

- Instead of sparse, one-hot vector representations used in the previous methods, we use embedded vector representations for each feature.
- Features used:
  - Vector representation of first few words in buffer and stack and their dependents
  - POS tags for those words
  - Arc labels for dependents

**Softmax layer:**
$$p = \text{softmax}(W_2 h)$$
**Hidden layer:**
$$h = (W_1^w x^w + W_1^t x^t + W_1^l x^l + b_1)^3$$
**Input layer:** $[x^w, x^t, x^l]$

words      POS tags      arc labels

# RNNs

CS224N Midterm Review

Sam Kim

# Overview

- Language models

- Applications of RNNs

- Backpropagation of RNNs

- Vanishing gradient problem

- GRUs and LSTMs

# A fixed-window neural Language Model

output distribution
$$\hat{y} = \mathrm{softmax}(\boldsymbol{U}\boldsymbol{h} + \boldsymbol{b}_2) \in \mathbb{R}^{|V|}$$

hidden layer
$$\boldsymbol{h} = f(\boldsymbol{W}\boldsymbol{e} + b_1)$$

concatenated word embeddings
$$\boldsymbol{e} = [\boldsymbol{e}^{(1)}; \boldsymbol{e}^{(2)}; \boldsymbol{e}^{(3)}; \boldsymbol{e}^{(4)}]$$

words / one-hots
$$\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \boldsymbol{x}^{(3)}, \boldsymbol{x}^{(4)}$$

books

laptops

a                    zoo

$\boldsymbol{U}$

$\boldsymbol{W}$

the          students      opened        their
$\boldsymbol{x}^{(1)}$     $\boldsymbol{x}^{(2)}$       $\boldsymbol{x}^{(3)}$       $\boldsymbol{x}^{(4)}$

# A fixed-window neural Language Model

**Improvements** over *n*-gram LM:
- No sparsity problem
- Model size is O(*n*) not O(exp(*n*))

Remaining **problems**:
- Fixed window is too small
- Enlarging window enlarges
- Window can never be large enough!
- Each $\boldsymbol{x}^{(i)}$ uses different rows of $\boldsymbol{W}$. We don't share weights across the window.

# Recurrent Neural Networks (RNN)

Core idea: Apply the same weights *repeatedly*



outputs (optional)

$\hat{\boldsymbol{y}}^{(1)}$   $\hat{\boldsymbol{y}}^{(2)}$   $\hat{\boldsymbol{y}}^{(3)}$   $\hat{\boldsymbol{y}}^{(4)}$   …

$\boldsymbol{h}^{(1)}$   $\boldsymbol{h}^{(2)}$   $\boldsymbol{h}^{(3)}$   $\boldsymbol{h}^{(4)}$

hidden states

$\boldsymbol{W}$   $\boldsymbol{W}$   $\boldsymbol{W}$   $\boldsymbol{W}$   …

input sequence (any length)

$\boldsymbol{x}^{(1)}$   $\boldsymbol{x}^{(2)}$   $\boldsymbol{x}^{(3)}$   $\boldsymbol{x}^{(4)}$   …

# RNN Language Model

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

**output distribution**

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

**hidden states**

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

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

**word embeddings**

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

**words / one-hot vectors**

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

# RNN Language Model

RNN **Advantages**:
- Can process any length input
- Model size doesn't increase for longer input
- Computation for step *t* can (in theory) use information from many steps back
- Weights are shared across timesteps → representations are shared

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

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



$\boldsymbol{h}^{(0)}$ $\boldsymbol{h}^{(1)}$ $\boldsymbol{h}^{(2)}$ $\boldsymbol{h}^{(3)}$ $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$ $\boldsymbol{W}_h$

$\boldsymbol{U}$

$\boldsymbol{W}_e$ $\boldsymbol{W}_e$ $\boldsymbol{W}_e$ $\boldsymbol{W}_e$

$\boldsymbol{e}^{(1)}$ $\boldsymbol{e}^{(2)}$ $\boldsymbol{e}^{(3)}$ $\boldsymbol{e}^{(4)}$

$\boldsymbol{E}$ $\boldsymbol{E}$ $\boldsymbol{E}$ $\boldsymbol{E}$

*the* *students* *opened* *their*

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

*books* *laptops*

a      zoo

# RNNs can be used for tagging
e.g. part-of-speech tagging, named entity recognition

# RNNs can be used for classification

e.g. sentiment analysis

**positive**

Sentence
encoding

How to compute
sentence encoding?

Basic way:
Use final hidden state

*equals*

*overall*  *I*  *enjoyed*  *the*  *movie*  *a*  *lot*

# RNNs can be used for classification

e.g. sentiment analysis



**positive**

Sentence encoding

How to compute sentence encoding?

Usually better:
Take element-wise max or mean of all hidden states

overall    I    enjoyed    the    movie    a    lot

# RNNs can be used to generate text
e.g. speech recognition, machine translation, summarization



Can use a RNN Language Model to generate text by repeated sampling. Sampled output is next step's input, typically by taking the *argmax* of each probability distribution.

# Multivariable Chain Rule

- Given a multivariable function $f(x, y)$, and two single variable functions $x(t)$ and $y(t)$, here's what the multivariable chain rule says:

$$\underbrace{\frac{d}{dt} f(x(t), y(t))}_{\text{Derivative of composition function}} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

One final output  $f(x(t), y(t))$

Two intermediate outputs  $x(t)$     $y(t)$

One input  $t$

# Backpropagation for RNNs

- Given a multivariable function $f(x, y)$, and two single variable functions $x(t)$ and $y(t)$, here's what the multivariable chain rule says:

$$\underbrace{\frac{d}{dt} f(x(t), y(t))}_{\text{Derivative of composition function}} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

In our example:

$$J^{(t)}(\theta)$$

$$W_h|_{(1)} \quad W_h|_{(2)} \quad \ldots \quad W_h|_{(t)}$$

equals    equals    equals

$$W_h$$

Apply the multivariable chain rule:

$$\frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^{t} \frac{\partial J^{(t)}}{\partial W_h}\bigg|_{(i)} \boxed{\frac{\partial W_h|_{(i)}}{\partial W_h}} = 1$$

$$= \sum_{i=1}^{t} \frac{\partial J^{(t)}}{\partial W_h}\bigg|_{(i)}$$

# Backpropagation for RNNs

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

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

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

$$\theta^{(t)} = \boldsymbol{U}h^{(t)} + \boldsymbol{b}_2$$



$J^{(t)}(\boldsymbol{\theta})$

$\boldsymbol{h}^{(t-1)}$  $\boldsymbol{h}^{(t)}$

$\boldsymbol{W}_h$  $\boldsymbol{W}_h$

…  …

Recall $\boldsymbol{W}_h$ appears at every time step. Calculate the sum of gradients w.r.t each time it appears

**Question:** Consider only the last two time steps, t and t-1.
What's the derivative $\dfrac{\partial J^{(t)}}{\partial \boldsymbol{W}_h}$ ? Leave as a chain rule

# Backpropagation for RNNs

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

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

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

$$\theta^{(t)} = \boldsymbol{U}h^{(t)} + \boldsymbol{b}_2$$



$J^{(t)}(\theta)$

$\boldsymbol{h}^{(t-1)}$  $\boldsymbol{h}^{(t)}$

$\boldsymbol{W}_h$  $\boldsymbol{W}_h$  $\boldsymbol{W}_h$

…  …

Recall $\boldsymbol{W}_h$ appears at every time step. Calculate the sum of gradients w.r.t each time it appears

**Question:** Consider only the last two time steps, t and t-1.
What's the derivative $\dfrac{\partial J^{(t)}}{\partial \boldsymbol{W}_h}$ ? Leave as a chain rule

**Answer:** $\dfrac{\partial J^{(t)}}{\partial \boldsymbol{W}_h} = \displaystyle\sum_{i=t-1}^{t} \left.\dfrac{\partial J^{(t)}}{\partial \boldsymbol{W}_h}\right|_{(i)} = \dfrac{\partial J^{(t)}}{\partial \boldsymbol{\theta}^{(t)}} \dfrac{\partial \boldsymbol{\theta}^{(t)}}{\partial \boldsymbol{h}^{(t)}} \left( \dfrac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{z}^{(t)}} \dfrac{\partial \boldsymbol{z}^{(t)}}{\partial \boldsymbol{W}_h} + \dfrac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} \dfrac{\partial \boldsymbol{h}^{(t-1)}}{\partial \boldsymbol{z}^{(t-1)}} \dfrac{\partial \boldsymbol{z}^{(t-1)}}{\partial \boldsymbol{W}_h} \right)$

# Backpropagation for RNNs

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

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

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

$$\theta^{(t)} = \boldsymbol{U}h^{(t)} + \boldsymbol{b}_2$$



$J^{(t)}(\boldsymbol{\theta})$

$\boldsymbol{h}^{(t-1)}$    $\boldsymbol{h}^{(t)}$

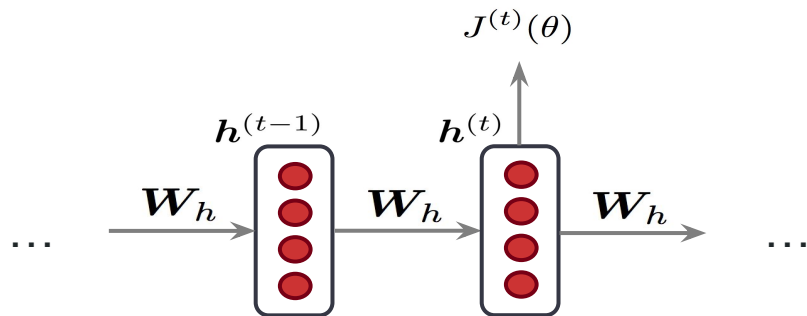… $\boldsymbol{W}_h$   $\boldsymbol{W}_h$   $\boldsymbol{W}_h$ …

Recall $\boldsymbol{W}_h$ appears at every time step. Calculate the sum of gradients w.r.t each time it appears

**Question:** Consider only the last two time steps, t and t-1.
What's the derivative $\dfrac{\partial J^{(t)}}{\partial \boldsymbol{W}_h}$ ? Leave as a chain rule

**Answer:** $\dfrac{\partial J^{(t)}}{\partial \boldsymbol{W}_h} = \sum_{i=t-1}^{t} \left.\dfrac{\partial J^{(t)}}{\partial \boldsymbol{W}_h}\right|_{(i)} = \dfrac{\partial J^{(t)}}{\partial \boldsymbol{\theta}^{(t)}} \dfrac{\partial \boldsymbol{\theta}^{(t)}}{\partial \boldsymbol{h}^{(t)}} \left( \dfrac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{z}^{(t)}} \dfrac{\partial \boldsymbol{z}^{(t)}}{\partial \boldsymbol{W}_h} + \dfrac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} \dfrac{\partial \boldsymbol{h}^{(t-1)}}{\partial \boldsymbol{z}^{(t-1)}} \dfrac{\partial \boldsymbol{z}^{(t-1)}}{\partial \boldsymbol{W}_h} \right)$

Looks scary!

# Backpropagation for RNNs

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

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

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

$$\theta^{(t)} = \boldsymbol{U}h^{(t)} + \boldsymbol{b}_2$$

$J^{(t)}(\theta)$

$\boldsymbol{h}^{(t-1)}$   $\boldsymbol{h}^{(t)}$

$\boldsymbol{W}_h$   $\boldsymbol{W}_h$   $\boldsymbol{W}_h$

...   ...

$$\frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h} = \sum_{i=t-1}^{t} \left.\frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h}\right|_{(i)} = \frac{\partial J^{(t)}}{\partial \boldsymbol{\theta}^{(t)}}\frac{\partial \boldsymbol{\theta}^{(t)}}{\partial \boldsymbol{h}^{(t)}}\left(\frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{z}^{(t)}}\frac{\partial \boldsymbol{z}^{(t)}}{\partial \boldsymbol{W}_h} + \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}}\frac{\partial \boldsymbol{h}^{(t-1)}}{\partial \boldsymbol{z}^{(t-1)}}\frac{\partial \boldsymbol{z}^{(t-1)}}{\partial \boldsymbol{W}_h}\right)$$

However, similarly in HW2, let $\gamma^{(t)} = \dfrac{\partial J^{(t)}}{\partial h^{(t)}}$   and   $\gamma^{(t-1)} = \dfrac{\partial J^{(t)}}{\partial h^{(t-1)}}$

This simplifies to,

$$\frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h} = \sum_{i=t-1}^{t} \left.\frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h}\right|_{(i)} = \gamma^{(t)}\frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{z}^{(t)}}\frac{\partial \boldsymbol{z}^{(t)}}{\partial \boldsymbol{W}_h} + \gamma^{(t-1)}\frac{\partial \boldsymbol{h}^{(t-1)}}{\partial \boldsymbol{z}^{(t-1)}}\frac{\partial \boldsymbol{z}^{(t-1)}}{\partial \boldsymbol{W}_h}$$

# Backpropagation for RNNs

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

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

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

$$\theta^{(t)} = \boldsymbol{U}h^{(t)} + \boldsymbol{b}_2$$

$$J^{(t)}(\theta)$$

$$\boldsymbol{h}^{(t-1)} \qquad \boldsymbol{h}^{(t)}$$

$$\dots \quad \boldsymbol{W}_h \qquad \boldsymbol{W}_h \qquad \boldsymbol{W}_h \qquad \dots$$

$$\frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h} = \sum_{i=t-1}^{t} \frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h}\bigg|_{(i)} = \frac{\partial J^{(t)}}{\partial \boldsymbol{\theta}^{(t)}}\frac{\partial \boldsymbol{\theta}^{(t)}}{\partial \boldsymbol{h}^{(t)}}\left(\frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{z}^{(t)}}\frac{\partial \boldsymbol{z}^{(t)}}{\partial \boldsymbol{W}_h} + \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}}\frac{\partial \boldsymbol{h}^{(t-1)}}{\partial \boldsymbol{z}^{(t-1)}}\frac{\partial \boldsymbol{z}^{(t-1)}}{\partial \boldsymbol{W}_h}\right)$$
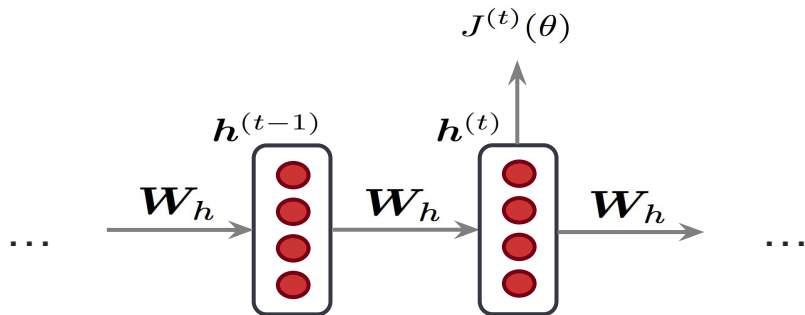
However, similarly in HW2, let $\gamma^{(t)} = \dfrac{\partial J^{(t)}}{\partial h^{(t)}}$ and $\gamma^{(t-1)} = \dfrac{\partial J^{(t)}}{\partial h^{(t-1)}}$

Can see how the gradient can be unrolled for T time steps

This simplifies to,

$$\frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h} = \sum_{i=t-1}^{t} \frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h}\bigg|_{(i)} = \gamma^{(t)}\frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{z}^{(t)}}\frac{\partial \boldsymbol{z}^{(t)}}{\partial \boldsymbol{W}_h} + \gamma^{(t-1)}\frac{\partial \boldsymbol{h}^{(t-1)}}{\partial \boldsymbol{z}^{(t-1)}}\frac{\partial \boldsymbol{z}^{(t-1)}}{\partial \boldsymbol{W}_h}$$

# Backpropagation for RNNs

$$\hat{\boldsymbol{y}}^{(t)} = \mathrm{softmax}\left(\boldsymbol{U}\boldsymbol{h}^{(t)} + \boldsymbol{b}_2\right) \in \mathbb{R}^{|V|}$$

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

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

$$\theta^{(t)} = \boldsymbol{U}h^{(t)} + \boldsymbol{b}_2$$

$J^{(t)}(\theta)$

$\boldsymbol{h}^{(t-1)}$  $\boldsymbol{h}^{(t)}$

$\cdots$  $\boldsymbol{W}_h$  $\boldsymbol{W}_h$  $\boldsymbol{W}_h$  $\cdots$

$$\frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h} = \sum_{i=t-1}^{t} \frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h}\bigg|_{(i)} = \frac{\partial J^{(t)}}{\partial \boldsymbol{\theta}^{(t)}} \frac{\partial \boldsymbol{\theta}^{(t)}}{\partial \boldsymbol{h}^{(t)}} \left( \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{z}^{(t)}} \frac{\partial \boldsymbol{z}^{(t)}}{\partial \boldsymbol{W}_h} + \boxed{\frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}}} \frac{\partial \boldsymbol{h}^{(t-1)}}{\partial \boldsymbol{z}^{(t-1)}} \frac{\partial \boldsymbol{z}^{(t-1)}}{\partial \boldsymbol{W}_h} \right)$$

can lead to vanishing or exploding gradients!

# Gradient Problems

- Backprop in RNNs have a recursive gradient call for hidden layer
- Magnitude of gradients of typical activation functions (sigmoid, tanh) lie between 0 and 1. Also depends on repeated multiplications of W matrix.
- If gradient magnitude is small/big, increasing timesteps decreases/increases the final magnitude.
- RNNs fail to learn long term dependencies.

**How to solve:**

*Exploding Gradients*
- gradient clipping

*Vanishing Gradients*
- use GRUs or LSTMs

# Vanishing Gradients

**Question**: (True/False) Adding L2-regularization will help with vanishing gradients

# Vanishing Gradients

**Question**: (True/False) Adding L2-regularization will help with vanishing gradients

**Answer:** False. This will pull the weights toward 0, which can make vanishing gradients worse

# Vanishing Gradients

**Question**: (True/False) Adding L2-regularization will help with vanishing gradients

# Vanishing Gradients

**Question**: (True/False) Adding more hidden layers will solve the vanishing gradient problem for a 2 layer neural network

**Answer:** False. Making the network deeper by adding hidden layers will increase the chance of vanishing gradient problems

# Gated Recurrent Units (GRU)

- Reset gate, $r_t$
- Update gate, $z_t$
- Intuition:
  - High $r_t$ => Short-term dependencies
- High $z_t$ => Long-term dependencies (solves vanishing gradients problem)

$$z_t = \sigma \left( W^{(z)} x_t + U^{(z)} h_{t-1} \right)$$

$$r_t = \sigma \left( W^{(r)} x_t + U^{(r)} h_{t-1} \right)$$

$$\widetilde{h}_t = \tanh \left( W x_t + r_t \circ U h_{t-1} \right)$$

$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \widetilde{h}_t$$

# Gated Recurrent Units (GRU)

**Question**: (True/False) If the update gate $z_t$ is close to 1, the net does not update its current state significantly

# Gated Recurrent Units (GRU)

**Question**: (True/False) If the update gate $z_t$ is close to 1, the net does not update its current state significantly

**Answer:** True. In this case, $h_t \approx h_{t-1}$

# Gated Recurrent Units (GRU)

**Question**: (True/False) If the update gate $z_t$ is close to 0 and the reset gate $r_t$ is close to 0, the net remembers the past state very well.

# Gated Recurrent Units (GRU)

**Question**: (True/False) If the update gate $z_t$ is close to 1 and the reset gate $r_t$ is close to 0, the net remembers the past state very well.

**Answer:** False. In this case, $h_t$ depends strongly on $x_t$ and not on $h_{t-1}$

# Long-Short-Term-Memories (LSTM)

- $i_t$: Input gate - How much does current input matter
- $f_t$: Forget gate - How much does past matter
- $o_t$: Output gate - How much should current cell be exposed
- $c_t$: New memory - Memory from current cell

$$i_t = \sigma\left(W^{(i)}x_t + U^{(i)}h_{t-1}\right)$$

$$f_t = \sigma\left(W^{(f)}x_t + U^{(f)}h_{t-1}\right)$$

$$o_t = \sigma\left(W^{(o)}x_t + U^{(o)}h_{t-1}\right)$$

$$\widetilde{c_t} = \tanh\left(W^{(c)}x_t + U^{(c)}h_{t-1}\right)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \widetilde{c_t}$$

$$h_t = o_t \circ \tanh\left(c_t\right)$$

# Long-Short-Term-Memories (LSTM)

$$i_t = \sigma\left(W^{(i)}x_t + U^{(i)}h_{t-1}\right)$$
$$f_t = \sigma\left(W^{(f)}x_t + U^{(f)}h_{t-1}\right)$$
$$o_t = \sigma\left(W^{(o)}x_t + U^{(o)}h_{t-1}\right)$$
$$\widetilde{c}_t = \tanh\left(W^{(c)}x_t + U^{(c)}h_{t-1}\right)$$
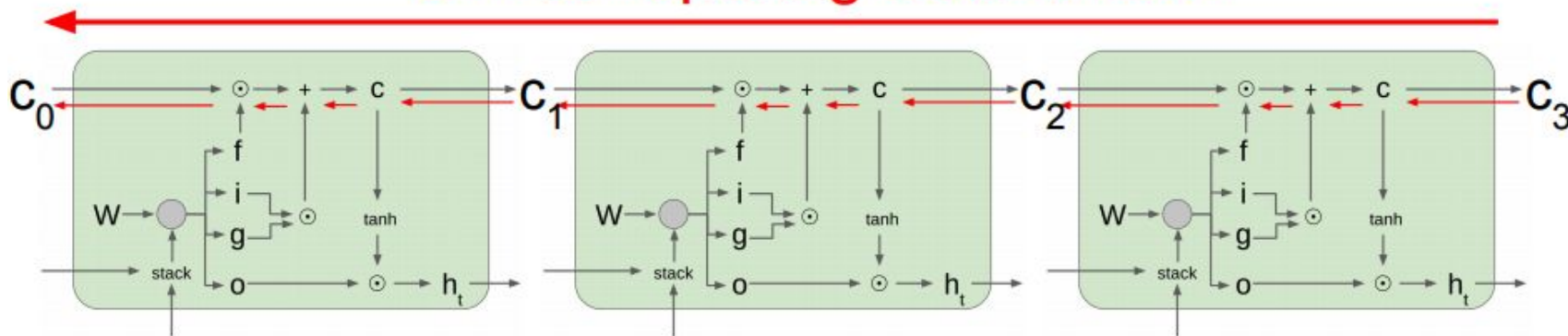$$c_t = f_t \circ c_{t-1} + i_t \circ \widetilde{c}_t$$
$$h_t = o_t \circ \tanh\left(c_t\right)$$

Backpropagation from $c_t$ to $c_{t-1}$ only elementwise multiplication by $f_t$. No longer only depends on $dh_t/dh_{t-1}$



Uninterrupted gradient flow!

# Long-Short-Term-Memories (LSTM)

**Question**: (True/False) If $f_t$ is very small or zero, then error will not be back-propagated to earlier time steps

# Long-Short-Term-Memories (LSTM)

**Question**: (True/False) If $f_t$ is very small or zero, then error will not be back-propagated to earlier time steps

**Answer:** False. $i_t$ and $\tilde{c}_t$ also still depend on $h_{t-1}$

# Long-Short-Term-Memories (LSTM)

**Question**: (True/False) The entries of $f_t$, $i_t$ and $o_t$ are non-negative.

# Long-Short-Term-Memories (LSTM)

**Question**: (True/False) The entries of $f_t$, $i_t$ and $o_t$ are non-negative.

**Answer:** True. The range of sigmoid is (0,1)

# Long-Short-Term-Memories (LSTM)

**Question**: (True/False) $f_t$, $i_t$ and $o_t$ can be viewed as probability distributions (entries sum to 1 and each entry is between 0 and 1)

# Long-Short-Term-Memories (LSTM)

**Question**: (True/False)  $f_t$, $i_t$ and $o_t$ can be viewed as probability distributions (entries sum to 1 and each entry is between 0 and 1)

**Answer:** False. Sigmoid is applied independently element-wise. The sum need not be 1.