

Natural Language Processing

Deep Learning for Sequence Processing

Prof. Iacopo Masi and Prof. Stefano Faralli

```

import matplotlib.pyplot as plt
import scipy
import numpy as np
import os
pd.set_option('display.coheader_justify', 'center')

import matplotlib
import matplotlib.pyplot as plt
import numpy as np
matplotlib inline
plt.style.use('seaborn-whitegrid')
font = {'family': 'Times',
        'weight': 'bold',
        'size': 12}
matplotlib.rcParams['font'] = font

# Aux functions
def plot_grid(Xs, Ys, axis=None):
    """Aux function to plot a grid"""
    if axis is None:
        axis = plt.gca()
    axis.plot(0, 0, marker='x', color='r', linestyle='none') # plot origin
    axis.scatter(Xs,Ys, c=ct, cmap='jet', marker='o') # scatter x vs y
    axis.axis('scaled') # axis scaled
    else:
        plt.plot(0, 0, marker='x', color='r', linestyle='none') # plot origin
        plt.scatter(Xs,Ys, c=ct, marker='o') # scatter x vs y
        plt.axis('scaled') # axis scaled

def linear_map(A, Xs, Ys):
    """Map src points with A to target points along 3-rd axis, like adding another layer
    src = np.stack([Xs,Ys], axis=Xs.ndim)
    axis = np.arange(0,src.shape[0])
    # (NNU)2
    src_r = src.reshape(-1,src.shape[-1]) # mask reshape to keep last dimension and adjust the rest
    axis = np.arange(0,src_r.shape[0])
    dst = A @ src_r.T # ZMM
    dst = dst.T # transpose as NNN2
    dst = (dst.T).reshape(src.shape)
    # access 3 and -3
    return dst[...,:], dst[...,-1]
    return dst[...,:], dst[...,-1]

def plot_points(ax, Xs, Ys, col='red', unit=None, linestyle='solid'):
    ax.set_aspect('equal')
    ax.grid(True, which='both')
    ax.axhline(0, color='gray', linestyle='--')
    ax.axvline(0, color='gray', linestyle='--')
    if unit is None:
        plotvectors(ax, [0,1],[1,0], ['gray']+2, alpha=1, linestyle=linestyle)
    else:
        plotvectors(ax, unit, [col]+2, alpha=1, linestyle=linestyle)

def plotVector(ax, vecs, cols, alpha=1, linestyle='solid'):
    for i in range(len(vecs)):
        x = np.concatenate(([0,0], vecs[i]))
        ax.quiver(*x[[0,1]], *x[[1,2]], *x[[2,3]], angles='xy', scale_units='xy', scale=1, color=cols[i], alpha=alpha, linestyle=linestyle)

def plot_pointrc(ax, Xs, Ys, col='red', unit=None, linestyle='solid'):
    ax.set_aspect('equal')
    ax.grid(True, which='both')
    ax.axhline(0, color='gray', linestyle='--')
    ax.axvline(0, color='gray', linestyle='--')
    if unit is None:
        plotvectorsrc(ax, [0,1],[1,0], ['gray']+2, alpha=1, linestyle=linestyle)
    else:
        plotvectorsrc(ax, unit, [col]+2, alpha=1, linestyle=linestyle)

def plotvectorsrc(ax, unit, [col]+2, alpha=1, linestyle='solid'):
    for i in range(len(unit)):
        x = np.concatenate(([0,0], unit[i]))
        ax.quiver(*x[[0,1]], *x[[1,2]], *x[[2,3]], angles='xy', scale_units='xy', scale=1, color=col[i], alpha=alpha, linestyle=linestyle)
    
```

Language Model: the order matters!

the cat is small
 $\begin{matrix} 0 & 1 & 2 & 3 \end{matrix}$
the cat is small
 $p(\text{the}, \text{cat}, \text{is}, \text{small}) = 0.035$

small the cat
 $\begin{matrix} 0 & 1 & 2 & 3 \end{matrix}$
small the cat is
 $p(\text{small}, \text{the}, \text{cat}, \text{is}) = 0.000001$

Language Model (LM)

More formally, given a sequence of words w_1, \dots, w_t , compute the probabilities distributions over the text words w_{t+1} where the support of the probability is a vocabulary V .

$$p(w_{t+1}|w_1, \dots, w_t)$$

where w_{t+1} can be any word in $V = \{w_1, \dots, w_T\}$

A system that does this is called a **Language Model**

Language Model (LM)

You can also think a LM as a **machinery** that assign a probability to a piece of text of T words:

$$\begin{aligned} p(w_1, \dots, w_t) &= \underbrace{p(w_1|w_0)}_{\text{baseline}} \cdot \underbrace{p(w_2|w_1)}_{\text{recursion}} \cdot \dots \cdot p(w_t|w_{t-1}) \\ &= p(w_1, \dots, w_t|w_0) = p(w_1) \cdot p(w_2|w_1) \cdot p(w_3|w_2) \cdot \dots \cdot p(w_t|w_{t-1}) \\ &= p(w_1, \dots, w_t|w_1, w_2) \cdot p(w_2|w_1) \cdot \prod_{i=3}^t p(w_i|w_{i-1}, \dots, w_1) \cdot p(w_t) \end{aligned}$$

LM can be interpreted as an Autoregressive Model

Such models that regress the value of a signal on the previous values of that same signal are naturally called **autoregressive models**.

$$p(w_{t+1}|w_1, \dots, w_t)$$

What LM can do

Language models prove useful for all sorts of reasons.

- Evaluate the likelihood of sentences
 - For example, we might wish to compare the naturalness of two candidate outputs generated by a machine translation system or by a speech recognition system.
 - Ability to sample sequences, and even to optimize for the most likely sequences. Act as **Generators**.
- LM are not **ONLY generators (samplers)** yet are statistical tools that models joint density of words where word order matters.

My own latex definitions

Today's lecture

- Recap on Language Models
- From Feedforward Nets to Recurrent Neural Nets (RNN)
- RNN Application: Parts of Speech (POS) as Sequence Labeling
- RNN Application: Text Generation as Autoregressive model

This lecture material is taken from

- Chapter 9 Jurafsky Book
- Chapter 6 Eisenstein Book
- Stanford Slide RNN
- Stanford Lecture RNN
- Stanford Notes on Word2Vec
- Andrej Karpathy Lecture on RNN
- Andrej Karpathy Slides on RNN

Another resource with code is [\[G2T\] RNN](#)

Deep Learning for Sequence Processing

Language Model (LM)

Language Modeling (LM) is the task of predicting what words comes next:

the students opened their $___$
books, laptops, exams, minds?

books? laptops? exams? minds?

Language Model: the order matters!

the cat is small
 $\begin{matrix} 0 & 1 & 2 & 3 \end{matrix}$
the cat is small
 $p(\text{the}, \text{cat}, \text{is}, \text{small}) = 0.035$

small the cat
 $\begin{matrix} 0 & 1 & 2 & 3 \end{matrix}$
small the cat is
 $p(\text{small}, \text{the}, \text{cat}, \text{is}) = 0.000001$

A fixed-window Neural LM (Markovian Assumption)

$$p(w_{t+1}|w_1, \dots, w_t) \approx p(w_{t+1}|w_t, \dots, w_{t-N})$$

The length of the context N impacts the LM!

The Length of the Context matters! (Markovian Assumption)

Language Modeling (LM) is the task of predicting what words comes next:

the students opened their $___$

What is your prediction here?

As the proctor started the clock, the students opened their $___$

What is your prediction here?

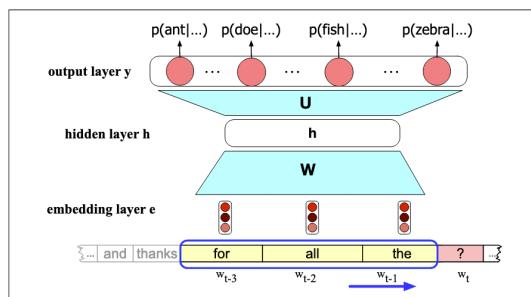


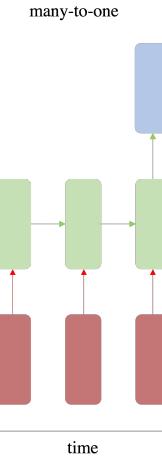
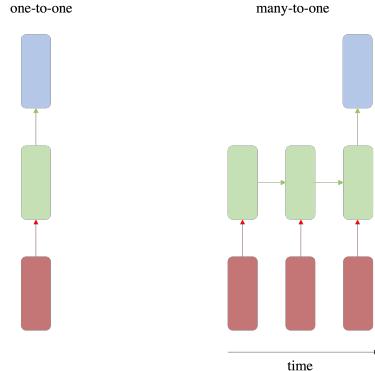
Figure 9.1 Simplified sketch of a feedforward neural language model moving through a text. At each time step t the network converts N context words, each to a d -dimensional embedding, and concatenates the N embeddings together to get the $Nd \times 1$ unit input vector \mathbf{x} for the network. The output of the network is a probability distribution over the vocabulary representing the model's belief with respect to each word being the next possible word.

Improvements over n-gram LM

- No sparsity problem
- No need to store all observed n-grams!

Remaining problems:

- Fixed window is too small
- Enlarging window enlarges \mathbf{W}
- Window can never be large enough
- No symmetry in how the inputs are processed.



RNN: basic principle

- We still use Markovian assumption to consider only N previous steps in time
- Yet we do not model explicitly $p(w_t | w_{t-1}, \dots, w_{t-N+1})$. Instead we use a **latent variable model**:

$$p(w_t | w_{t-1}, \dots, w_1) \approx p(w_t | h_{t-1}),$$

where h_{t-1} is a hidden state that "summarizes" the sequence information up to time step $t-1$.

$$p(w_t | w_{t-1}, \dots, w_1) \approx p(w_t | h_{t-1}),$$

RNN: hidden state

$p(w_t | w_{t-1}, \dots, w_1) \approx p(w_t | h_{t-1})$,

where h_{t-1} is a hidden state that stores the sequence information up to time step $t-1$.

In general, the hidden state at any time step t could be computed based on both the current input w_t and the previous hidden state h_{t-1} :

$$h_t = f(w_t, h_{t-1})$$

Informal: think h_t as a sort of memory of what happened so far up to $t-1$ plus the current input w_t .

Hidden layer vs Hidden states

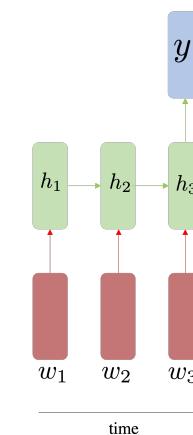
Recall that we have discussed hidden layers with hidden units in [window-based Neural LM](#).

It is noteworthy that hidden layers and hidden states refer to two very different concepts.

- Hidden layers are, as explained, layers that are hidden from view on the path from input to output.
- Hidden states are technically speaking inputs to whatever we do at a given step, and they can only be computed by looking at data at previous time steps.

Sequence Modeling

many-to-one



How do we implement $f(w_t, h_{t-1})$?

Assume e_t is the word embedding of word token at time t and w_t corresponds to the one-hot encoding of the word, as $e_t = \theta w_t$. Then:

$$f(w_t, h_{t-1}) \doteq \sigma(\mathbf{W}_h h_{t-1} + \mathbf{W}_w e_t + \mathbf{b})$$

This formulation is called [Elman Networks](#) [Elman, 1990] or simple recurrent neural networks.

How do we implement $f(w_t, h_{t-1})$?

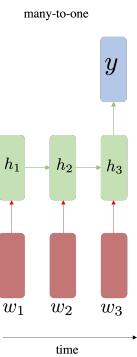
Assume e_t is the word embedding of word token at time t and w_t corresponds to the one-hot encoding of the word, as $e_t = \theta w_t$. Then:

$$f(w_t, h_{t-1}) \doteq \sigma\left(\frac{\text{bias prev. step}}{\text{current input}} + \mathbf{W}_h h_{t-1} + \mathbf{W}_w e_t + \mathbf{b}\right)$$

where $\sigma(\cdot)$ is a non-linear activation function such as hyperbolic tan or sigmoid.

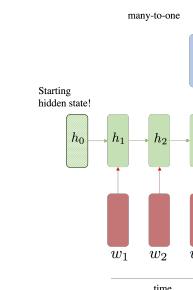
Sequence Modeling: something is missing here

Can you spot what?



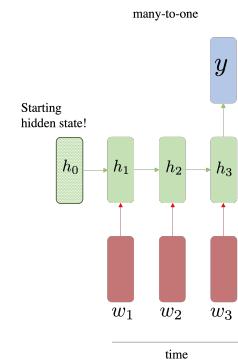
Sequence Modeling: something is missing here

Can you spot what? The base case: Starting hidden state!



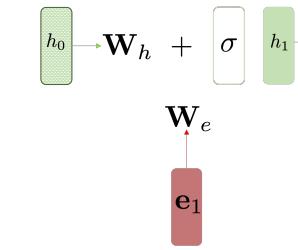
Sequence Modeling: Something is missing here, can you spot what?

You need the "base case": the starting hidden state, usually \mathbf{h}_0 is set to all zeros.



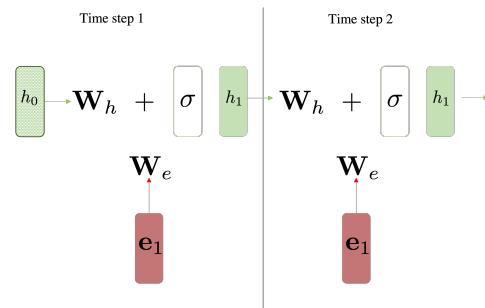
What happens at each time step

Note: Bias omitted for simplicity



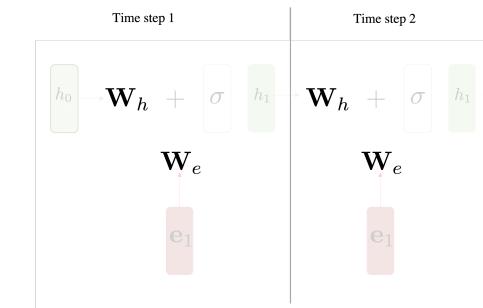
What happens at each time step

Note: Bias omitted for simplicity



The weights are always the same!

Note: Bias omitted for simplicity



\mathbf{W}_e and \mathbf{W}_h do NOT depend on t !

$$h_t = \sigma(\underbrace{\mathbf{W}_h h_{t-1}}_{\text{from prev. step}} + \underbrace{\mathbf{W}_e e_t}_{\text{current input}} + b)$$

Good Scalability: for space complexity.

The model scales w.r.t to the time dimension, so we can allow our context to be long!

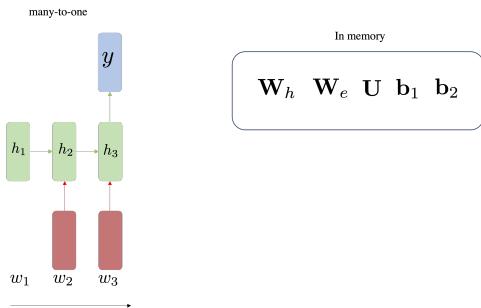
From the hidden state to an output

$$h_t = \sigma(\underbrace{\mathbf{W}_h h_{t-1}}_{\text{from prev. step}} + \underbrace{\mathbf{W}_o o_t}_{\text{next output}} + b)$$

$$y_t = \text{softmax}(\mathbf{U}h_t + b_t)$$

Dynamics vs what we really store in memory

Therefore, the parameterization cost of an RNN does not grow as the number of time steps increases.



Implementation detail

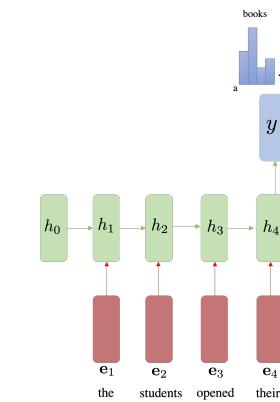
$$h_t = \sigma(\underbrace{W_h h_{t-1}}_{\text{from prev. step}} + \underbrace{W_e e_t}_{\text{current input}} + b)$$

Can be implemented as **fully-connected layer** matrix multiplication using concatenation:

$$h_t = \sigma(W_h; W_e)[h_{t-1}; e_t] + b$$

RNN Forward pass in pseudo code

```
# Input x (the sequence of tokens)
h[0] = zeros # h is an array of all zeros!
for i in range(1, len(x)): # unroll the RNN in time (this is the slow part)
    h[i] = sigmoid(W_h * h[i-1] + W_e * e[i] + b_1) # W_h and W_e shared across time
    y[i] = softmax(h[i] * b_2) # h[i], y[i] are time dependent
```



RNN as a Neural Language Model

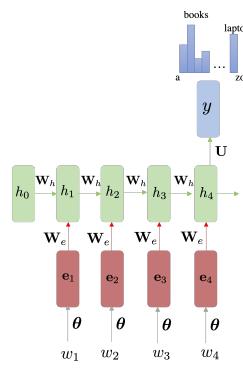
The RNN models $p(w_t | \text{the, students, opened, their})$.

Output distribution
 $y = \text{softmax}(U h_t + b_2)$

Hidden states
 $h_t = \sigma(W_h h_{t-1} + W_e e_t + b)$

Word embeddings
 $e_t = \theta w_t$

Word token as one-hot encoding
 $w_t \in \{0, 1\}^{|V|}$

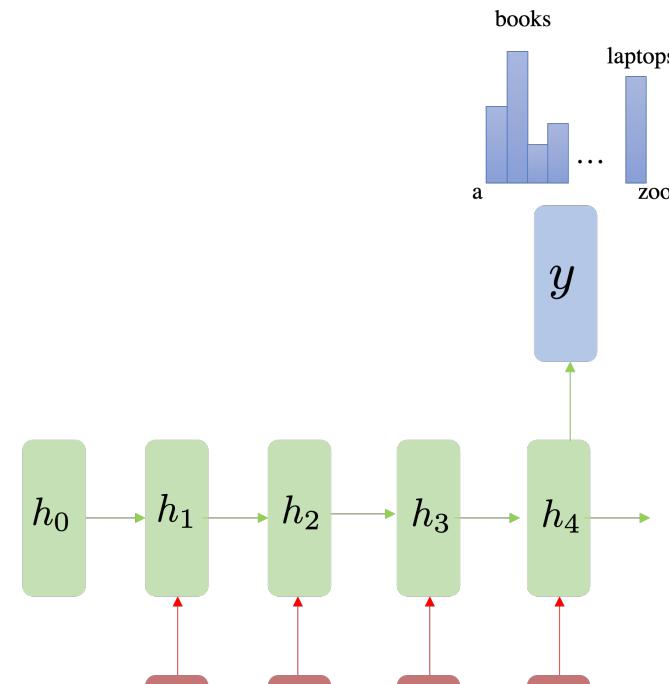


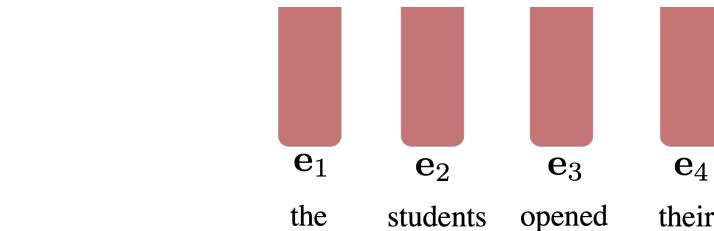
RNN Language Model

Improvements over window-based LM

- Can process any length input
- Computation for step t can (in theory) use information from many steps back
- Bounded model size relative to context size
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

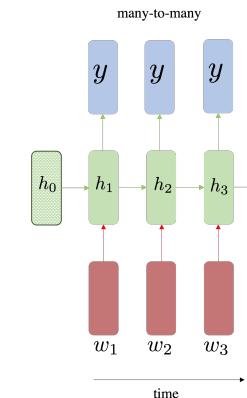
- Remaining problems:
 - Recurrent computation is slow (unroll in time)
 - In practice, difficult to access information from many steps back





Training a RNN Language Model

RNN LM is many-to-many
Sequence to Sequence modeling

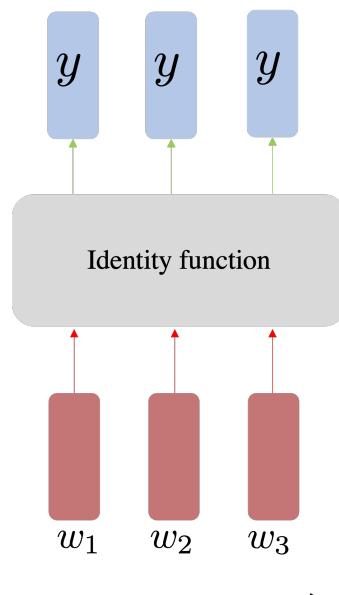


RNN LM is many-to-many

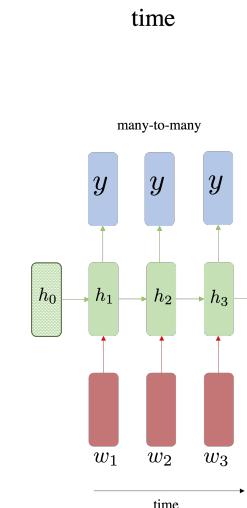
- The most "stupid" sequence to sequence modeling is the **identity function** of the data
- The simplest yet useless way to generate the data is propagate the data itself
- You can think this also as a form of **memorization** (**overfitting**) compared to learning.

LM : NLP = Generative Models : Vision

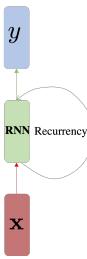
many-to-many



RNN LM is many-to-many
Sequence to Sequence modeling

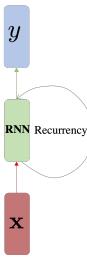


Recurrency (Bad Visualization)



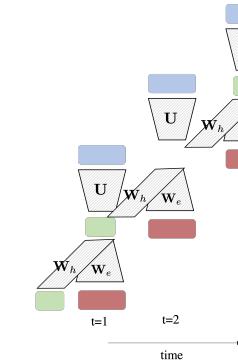
Recurrency (Bad Visualization)

In theory, recurrency with cycles is a problem for training, but we can show that we can unfold the RNN over time and still create a DAG.



Unfolding the recurrency over time

Explicitly unrolling a recurrent network into a feedforward computational graph eliminates any explicit recurrences, allowing the network weights to be trained directly.
The step of unfolding is an hyper-param that you have to decide.



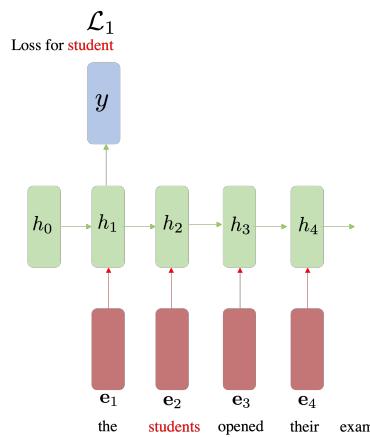
Training a RNN Language Model

- Get a big corpus of text which is a sequence of words $\{w_1, \dots, w_N\}$
- Feed a window till time t to the RNN-LM, compute probability distribution over V at time $t+1$
 - Predict a probability distribution \hat{y} over words, given words so far
- Loss function is same as word2vec: we compare two discrete distributions with cross-entropy

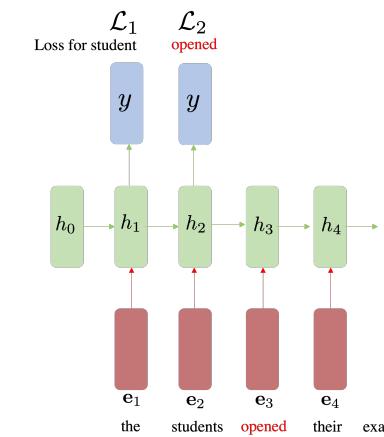
$$\mathcal{L}(\{\mathbf{W}\}) = CE(\hat{y}, y) = - \sum_{v \in V} y \log(\hat{y}) = - \log(\hat{y})[t+1]$$

- Note: $[t+1]$ selects the index in V of the next word, at time $t+1$
- We use again self supervision

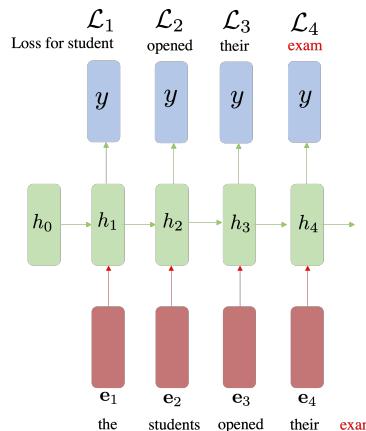
Training a RNN Language Model



Training a RNN Language Model

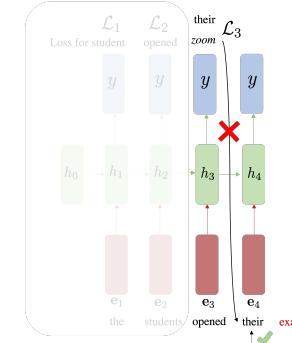


Training a RNN Language Model



Teacher Forcing at training time

When training, **the autoregressive capability is turned OFF**. This is called **teacher forcing**. At time $t = 3$ the model predicts `zoom` but the word should have been `their`. As input to time $t = 3$ we use the value from the corpus: `their`. NOT the prediction `zoom`.



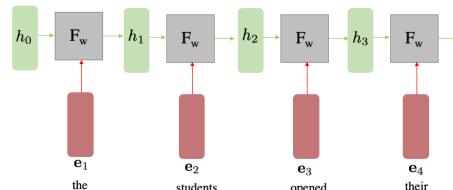
How do we get the gradients on the weights?

We use Back Propagation Through Time (BPTT)

Back Propagation Through Time (BPTT)

We can approach BPTT: 1) using the computational graph (engineering point of view/intuition) 2) mathematically (more formal, gives a motivation)

Forward pass in RNN



Training a RNN Language Model

- Computing the loss and gradients across the entire corpus $\{w_1, \dots, w_N\}$ is too expensive
 - We would do a simple GD update after seeing the entire corpus!
 - So we use an approximation using SGD.
 - We consider $\{w_1, \dots, w_N\}$ as a sentence or a document but not the entire corpus.

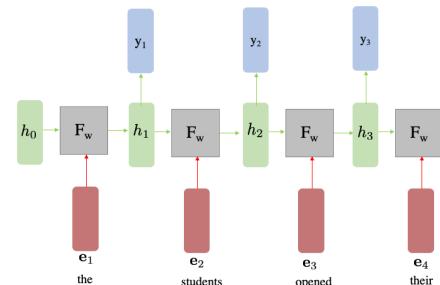
$$\frac{1}{T} \sum_{t=1}^T \mathcal{L}_t(\theta)$$

- SGD allows us to compute loss and gradients for a small chunk of text, and then update the model.

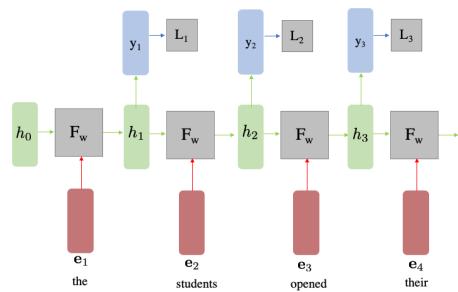
Training a RNN Language Model

1. Compute the loss for a sentence (Note: in the implementation is a batch of sentences)
2. Compute the gradients on weights (Most difficult part to understand! [ai](#))
3. Update the weights
4. Get another batch of sentences and repeat.

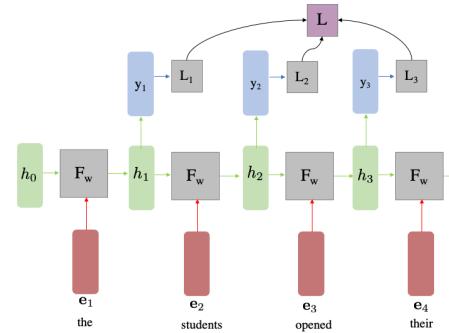
Forward pass in RNN



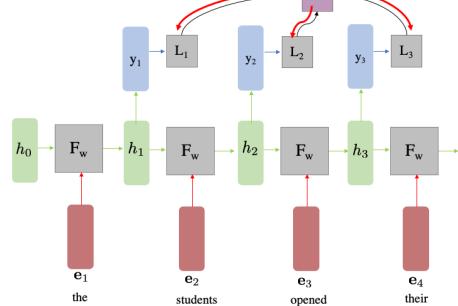
Forward pass in RNN



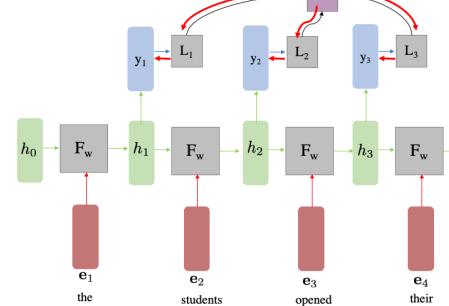
Forward pass in RNN



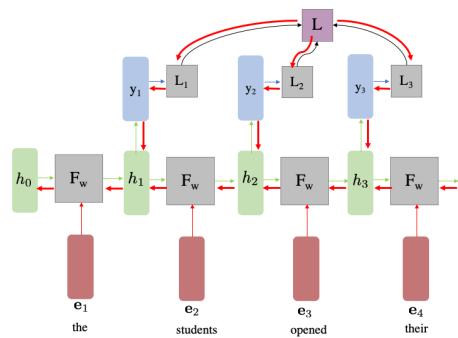
Back Propagation Through Time (BPTT) - Backward



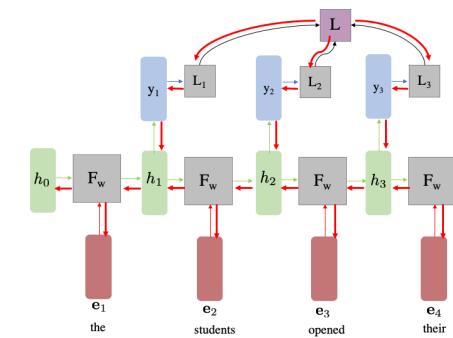
Back Propagation Through Time (BPTT) - Backward



Back Propagation Through Time (BPTT) - Backward



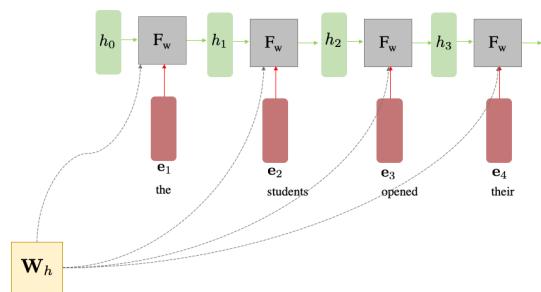
Back Propagation Through Time (BPTT) - Backward



Back Propagation Through Time (BPTT) - Backward

Eventually we will arrive to receive gradients on weights.

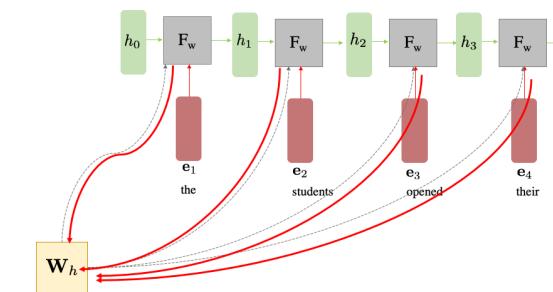
Remember at each time step F_w received different inputs yet the parameters W_h are the same across time.



Back Propagation Through Time (BPTT) - Backward

Eventually we will arrive to receive gradients on weights.

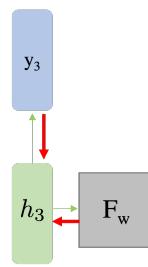
Remember at each time step F_w received different inputs yet the parameters W_h are the same across time.



Gradients sum at outward branches

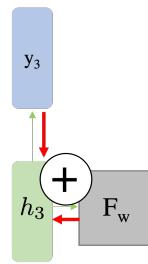
h_3 is input to two functions "in the future", one produce y_3 , the other will produce h_4 .

So h_3 will receive two gradients that will be summed together



Gradients sum at outward branches

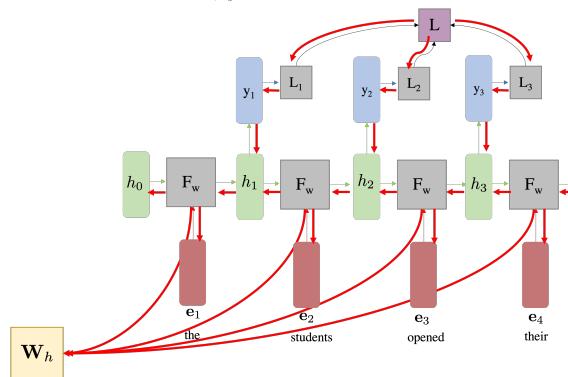
$$\frac{\partial f(x(t), y(t))}{\partial t} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial t}$$



Back Propagation Through Time (BPTT)

Let's first agree on what we have to compute:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_h} = ? \quad \text{we use chain rule to compute it}$$



Back Propagation Through Time (BPTT)

Mathematical: apply chain rule. Note: I simplified by NOT showing:

- the biases
- the weights that transform the input \mathbf{W}_h and the weights that do classification at each output \mathbf{U}_i
- Keep in mind that we have to receive gradients over them because we need to update them as well

I do show the weights \mathbf{W}_h because this is the hardest part of BPTT.

Let's first agree on what we have to compute:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_h} = ? \quad \text{we use chain rule to compute it}$$

Note this is a function that maps a matrix \mathbf{W}_h to a scalar value, the loss \mathcal{L} .

This makes sense to compute because the loss \mathcal{L} depends on \mathbf{W}_h .

We want to know what is the "infinitesimal" influence on the loss by perturbing a bit of \mathbf{W}_h (informal definition of gradient).

Back Propagation Through Time (BPTT)

Let's first agree on what we have to compute:

$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_h} = ?$ we use chain rule to compute it

Note this is a function that maps a matrix \mathbf{W}_h to a scalar value, the loss \mathcal{L} .

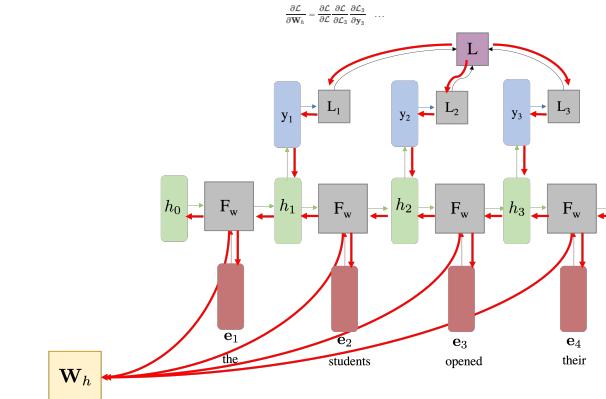
This makes sense to compute because the loss \mathcal{L} depends on \mathbf{W}_h .

We want to know what is the "infinitesimal" influence on the loss by perturbing a bit of \mathbf{W}_h (informal definition of gradient).

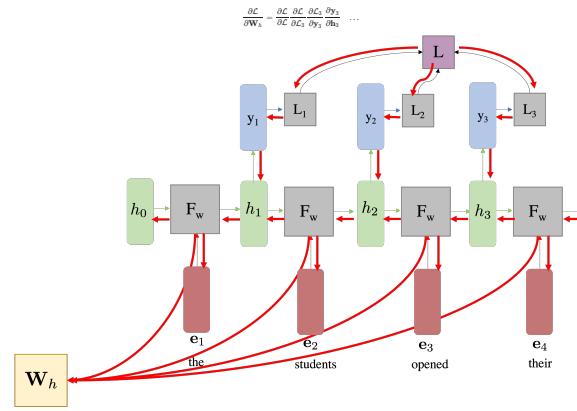
Back Propagation Through Time (BPTT)

Let's first agree on what we have to compute:

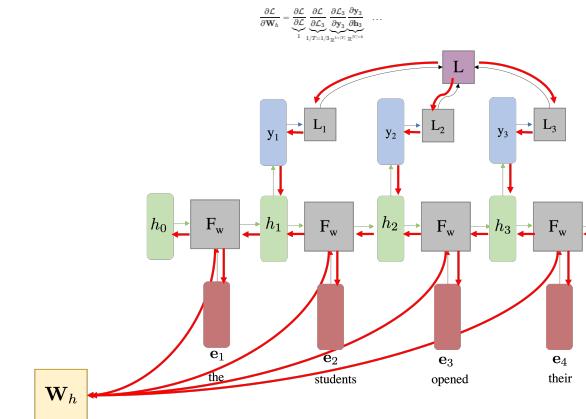
$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_h} = \frac{\partial \mathcal{L}}{\partial \mathbf{C}} \frac{\partial \mathbf{C}}{\partial \mathbf{y}_3} \frac{\partial \mathbf{y}_3}{\partial \mathbf{y}_2} \dots$$



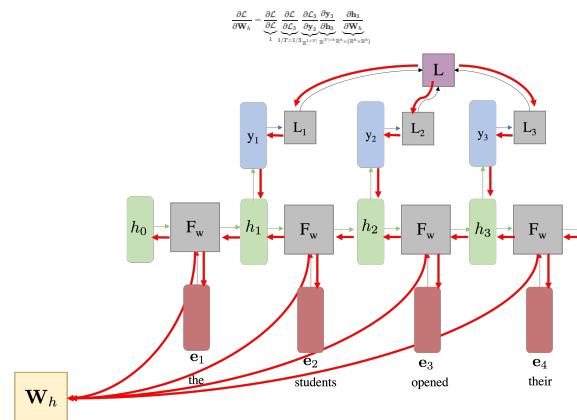
Back Propagation Through Time (BPTT)



Back Propagation Through Time (BPTT)



Back Propagation Through Time (BPTT)



BPTT

The tricky part is $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_h}$. Remember that $h_1 = f(h_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{x}_1)$
 $h_2 \text{ depends on } h_1$

BPTT

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_h} = \frac{\partial \mathcal{L}}{\partial \mathcal{L}_1} \frac{\partial \mathcal{L}_1}{\partial y_1} \frac{\partial y_1}{\partial \mathbf{y}_2} \frac{\partial \mathbf{y}_2}{\partial h_1} \frac{\partial h_1}{\partial \mathbf{W}_h}$$

$$h_1 = f(h_0, \mathbf{W}_1, \dots)$$

BPTT and Vanishing Gradient problem

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_h} = \frac{\partial \mathcal{L}}{\partial \mathbf{C}} \frac{\partial \mathcal{L}}{\partial \mathbf{C}_t} \frac{\partial \mathcal{L}_t}{\partial \mathbf{y}_t} \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h}$$

$$\mathbf{h}_t = f(\underbrace{(\mathbf{h}_{t-1}, \mathbf{W}_{h, \dots})}_{\mathbf{h}_0}, \mathbf{W}_h, \dots)$$

This is the essence of BPTT

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_h} = \frac{\partial \mathcal{L}}{\partial \mathbf{C}} \frac{\partial \mathcal{L}}{\partial \mathbf{C}_t} \frac{\partial \mathcal{L}_t}{\partial \mathbf{y}_t} \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h}$$

$$\mathbf{h}_t = f(\underbrace{(\mathbf{h}_{t-1}, \mathbf{W}_{h, \dots})}_{\mathbf{h}_0}, \mathbf{W}_h, \dots)$$

If you want to go deeper, click on BPTT (https://ufal.mff.cuni.cz/~archer/recurrent-neural-networks/bptt.html)

This is the essence of BPTT

Memory cost:

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} = ?$$

So far:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_h} = \frac{\partial \mathcal{L}}{\partial \mathbf{C}} \frac{\partial \mathcal{L}}{\partial \mathbf{C}_t} \frac{\partial \mathcal{L}_t}{\partial \mathbf{y}_t} \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h}$$

$$\mathbf{h}_t = f(\underbrace{(\mathbf{h}_{t-1}, \mathbf{W}_{h, \dots})}_{\mathbf{h}_0}, \mathbf{W}_h, \dots)$$

If you want to go deeper, click on BPTT (https://ufal.mff.cuni.cz/~archer/recurrent-neural-networks/bptt.html)

This is the essence of BPTT

Memory cost:

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} \text{ at time step } t = \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} \quad \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_2} \quad \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_2}$$

note: another rec.

So far:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_h} = \frac{\partial \mathcal{L}}{\partial \mathbf{C}} \frac{\partial \mathcal{L}}{\partial \mathbf{C}_t} \frac{\partial \mathcal{L}_t}{\partial \mathbf{y}_t} \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h}$$

$$\mathbf{h}_t = f(\underbrace{(\mathbf{h}_{t-1}, \mathbf{W}_{h, \dots})}_{\mathbf{h}_0}, \mathbf{W}_h, \dots)$$

If you want to go deeper, click on BPTT (https://ufal.mff.cuni.cz/~archer/recurrent-neural-networks/bptt.html)

BPTT

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} \text{ at time step } t = \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} / \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_2}$$

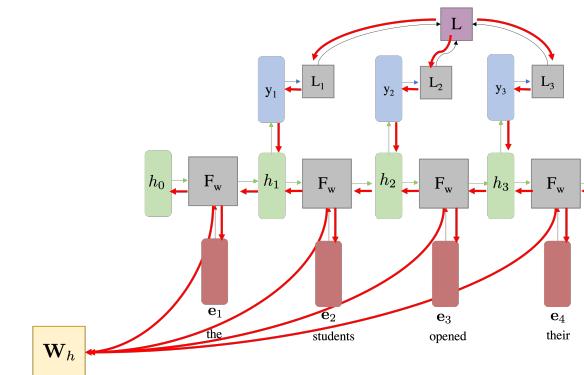
$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} \text{ at time step } t = \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} / \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}_h} \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}_h}$$

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} \text{ at time step } t = \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} / \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}_h} \left[\frac{\partial \mathbf{h}_2}{\partial \mathbf{h}_1} + \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}_h} \right]$$

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} \text{ at time step } t = \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} + \frac{\partial \mathbf{h}_1}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}_h} + \frac{\partial \mathbf{h}_2}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}_h}$$

Back Propagation Through Time (BPTT)

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_h} = \frac{\partial \mathcal{L}}{\partial \mathbf{C}} \frac{\partial \mathcal{L}}{\partial \mathbf{C}_t} \frac{\partial \mathcal{L}_t}{\partial \mathbf{y}_t} \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h}$$



BPTT and Vanishing Gradient Problem

If the recursive chain is too long, the product of matrices $\frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} \dots$ can make the gradient to vanish especially if using tanh activation function.

That is multiplying $a \cdot b \cdot b \cdot b \cdot b \dots$ where $0 < b < 1$ (not for matrices).

At the end the norm of the gradient will be so small that will get to zero numerically!

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} \text{ at time step } t = \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} + \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_2}{\partial \mathbf{W}_h} + \frac{\partial \mathbf{h}_2}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}_h}$$

Vanishing Gradient and RNN

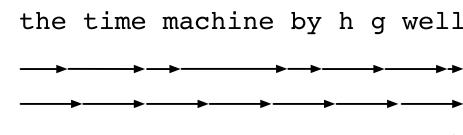
- For this reason RNN cannot capture long context dependency in the text and they are bounded to model moderately small sequences of text.
- Still better than window-based approach.

Vanishing Gradient and NLP

- Example: When she tried to print her tickets, she found that the printer was out of toner. She went to the stationery store to buy more toner. It was very overpriced. After installing the toner into the printer, she finally printed her.
- To learn this training example, the RNN LM needs to model the dependency between "tickets" on the 7th step and the target word "tickets" at the end.
- But if the gradient is small, the model cannot learn this dependency.
- So, the model is unable to predict similar long-distance dependencies at test time.

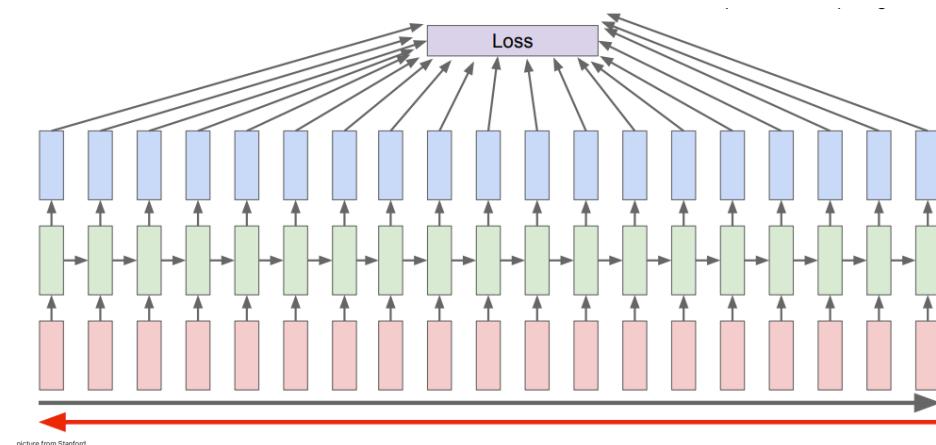
Vanishing Gradient Workaround: Truncated BPTT

- We can SGD compute the gradients over a small window (e.g. 25 steps) and then update the parameters.
- We avoid computing this over longer sequences.

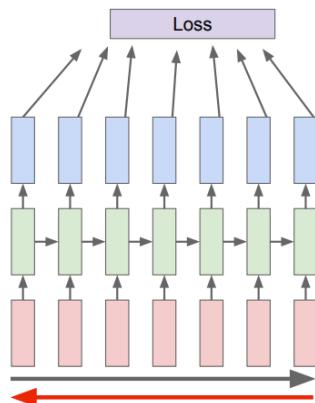


Computing strategies for computing gradients in RNNs. From top to bottom: randomized truncation, regular truncation, and full computation.

BPTT

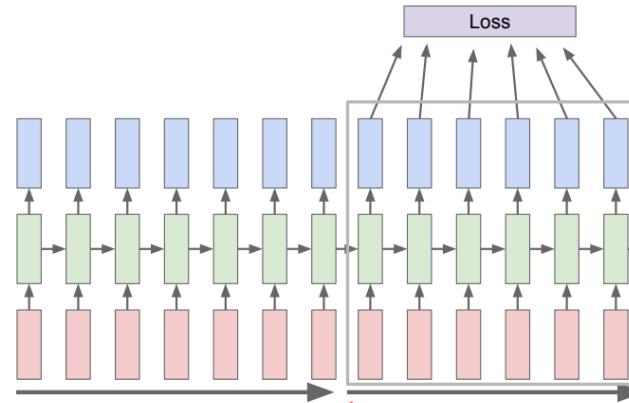


Truncated BPTT



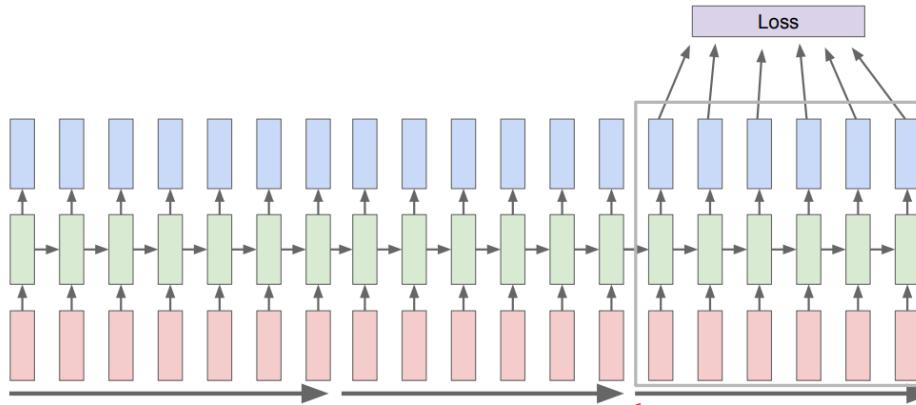
picture from Stanford

Truncated BPTT



Carry hidden states forward in time forever, but only backpropagate for some smaller number of steps

Truncated BPTT



picture from Stanford

Opposite problem: Exploding gradient

- The norm of the gradient becomes so big that it's force that is applied too strong and arbitrary.
- Machine generation cannot regenerate if "too strong". You quit it if it in the loss you encounter. Null.
- Faster problem than vanishing gradient

Easy work around is to clip the gradient

```
if grad.norm() higher than a maximum threshold then
    grad = grad / max_norm * max_norm
```

```
grad = grad.clip(-1, 1) # clip to mitigate exploding gradients
```

Gradient Clipping for Exploding Gradient

Algorithm 1 Pseudo-code for norm clipping

```
ĝ ← ∂E / ∂θ
if ||ĝ|| ≥ threshold then
    ĝ ← threshold * ĝ / ||ĝ||
end if
```

OK, so now that we know how to train it how do we generate text?

Remember LM is a Autoregressive Generative model

Generation from a LM - "Generating roll outs"

Equivalent to 1) generate samples from a generative model 2) use the autoregressive capability of the LM

Look up: [short](#)

1. Input a starting token
2. Input a hidden state. You have two choice: deterministic or probabilistic
 - Deterministic: take argmax over the softmax probabilities
 - Probabilistic: we use inverse transform sampling to sample a word from the softmax probabilities
3. Once you sampled at time t , feed it as input the sampled word at time $t+1$ and continue
4. The LM ends when reading a [\[END\]](#) special token

Inverse Transform Sampling

```
Given a list of cumulative probabilities V
1. Draw a random number U from uniform(0, 1)
2. Find the first index i for which u > v[i]
3. Sample uniformly in v ~ U(0, 1) give u to index i using the CDF
4. According to the interval you are in, now select the word
```

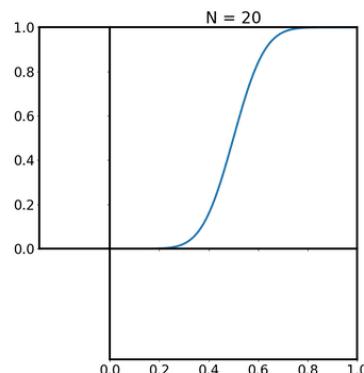
Inverse Transform Sampling

```
If n = random.rand(1000) #random sample uniformly from [0,1]
#n is a float between 0 and 1
#Find the first bit for which n > v[i]
#sampled_idx = np.argmax(n > v, axis=1)
```

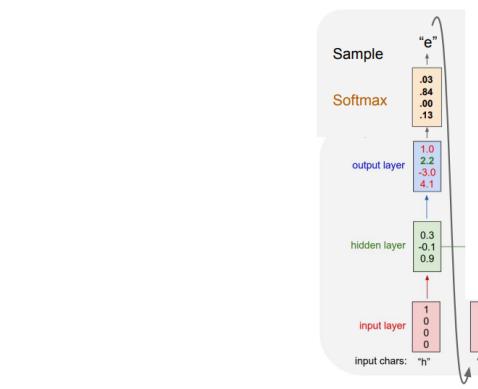
Inverse Transform Sampling even simpler with just Python

```
n = random.choices(v, weights=wprobabilities)
```

Inverse Transform Sampling in Action



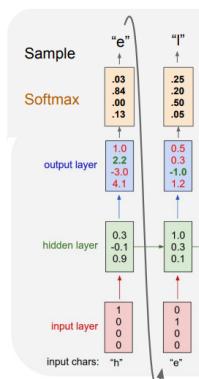
Generation from a LM - Character-level Language Model



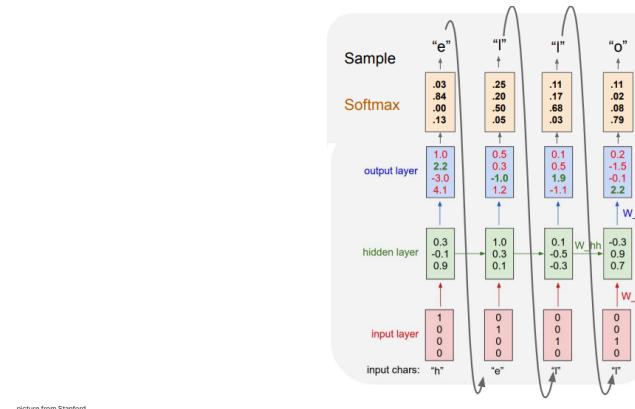
Generation from a LM - Character-level Language Model

- The vocabulary $V = \{h, e, l, o\}$ we have learned how to say hello
- We can start from the special token $\langle \text{start} \rangle$ (read a char as input)
- Long story short: Predict + sample (if probabilistic) + autoregress

Generation from a LM - Character-level Language Model

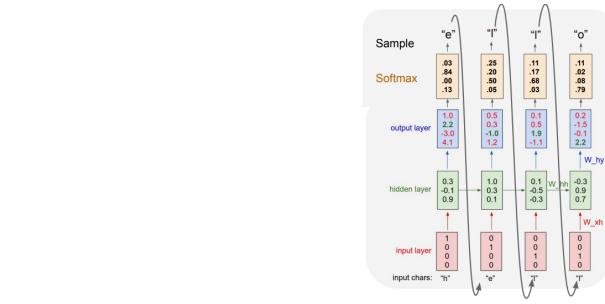


Generation from a LM - Character-level Language Model



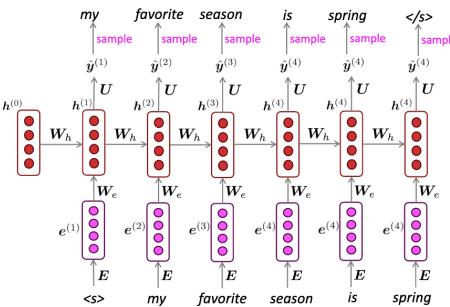
Generation from a LM - Character-level Language Model

Note that sampling makes the generation probabilistic. Even if we start from the same char, we may end up with completely different words (this gives variability in the output), yet the words should match the "distribution" of the text on which was trained on.



picture from Stanford

Generation from a LM - It works with words too



picture from Stanford

LM can be seen as an Encoder Decoder

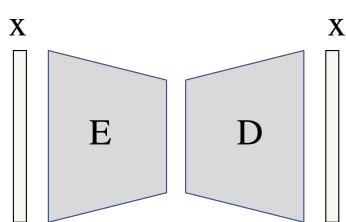
Important idea for when we will cover *Transformers* and *BERT*

$h = \text{encode}(x)$ we go from x to another more useful representation h
 $x' = \text{decode}(h)$ from h we can go back to an approximation of x

$x = \text{decode}(\{\text{encode}(x)\})$

Encoder Decoder for Vision

Here x is an image. The encoder maps the image to a low dimensional space using convolution.



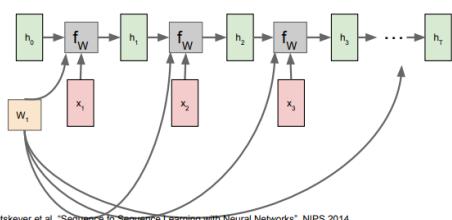
LM is an Encoder Decoder for NLP

Text is Sequence to Sequence but we can decompose it as:

Key Idea: We decompose Sequence-to-Sequence = Many-to-one + one-to-many

Encoder for NLP

Many-to-one: Encode input sequence in a single vector $f_w(x)$, the last hidden state

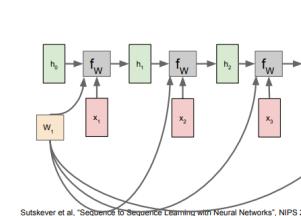


Sutskever et al., "Sequence to Sequence Learning with Neural Networks", NIPS 2014

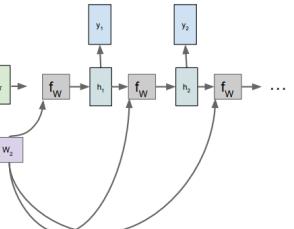
Encoder for NLP

Many to one: Encode input sequence in a single vector (e.g. the last hidden state)

Many to one: Encode input sequence in a single vector



One to many: Produce output sequence from single input vector



Now the fun 🎉

Minimal RNN implementation in numpy in ~100 lines

Code from Karpathy's "RNNs in 100 lines"

Same but updated to py3

Converges but Code is slower than Keras (MIT License)

We will just study this for practicing, in real scenario please use [pytorch](#)!

[min-char-rnn.py](#) gist: 112 lines of Python

```

/*
  Recurrent character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
  CC-BY-NC-SA
  import numpy as np

  # Data I/O
  data = open('input.txt', 'r').read() # should be simple plain text file
  chars = list(data)
  n_chars = len(chars)
  print("data has %d characters, %d unique." % (n_chars, len(set(chars))))
  char_to_ix = {ch: i for i, ch in enumerate(chars)}
  ix_to_char = dict(ix_to_char.items())
  print("chars: %s" % str(chars[:10]))
  print("char_to_ix: %s" % str(char_to_ix))
  print("ix_to_char: %s" % str(ix_to_char))

  # Hyperparameters
  hidden_size = 300 # size of hidden layer of neurons
  seq_length = 20 # number of steps to unroll the RNN for
  learning_rate = 1.0

  # Model parameters
  np.random.seed(1234)

  wh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
  bh = np.random.randn(hidden_size, 1) # bias for hidden state
  hb = np.zeros(hidden_size, 1) # bias for hidden state to output
  h0 = np.zeros(hidden_size, 1) # initial hidden state
  by = np.zeros(vocab_size, 1) # initial bias

  def loss(inputs, targets, hprev):
    """Compute total loss for a batch of inputs & targets.

    Inputs, targets are both list of integers.
    hprev is initial hidden state
    returns the loss, gradients on model parameters, and last hidden state
    """
    loss = 0
    next_h = h0
    for i in range(len(inputs)):
      x_t = inputs[i]
      y_t = targets[i]

      if i > 0:
        next_h = hprev

      hprev, loss = forward(next_h, x_t, wh, bh, by, hprev)

      if y_t is not None:
        loss += softmax_loss(next_h, y_t)

      next_h = np.tanh(next_h) * softmax(next_h)

    return loss, next_h

  def forward(inputs, targets, hprev):
    """Compute total loss for a batch of inputs & targets.

    Inputs, targets are both list of integers.
    hprev is initial hidden state
    returns the loss, gradients on model parameters, and last hidden state
    """
    loss = 0
    next_h = hprev
    for i in range(len(inputs)):
      x_t = inputs[i]
      y_t = targets[i]

      if i > 0:
        next_h = hprev

      hprev, loss = forward(next_h, x_t, wh, bh, by, hprev)

      if y_t is not None:
        loss += softmax_loss(next_h, y_t)

      next_h = np.tanh(next_h) * softmax(next_h)

    return loss, next_h

  def backward(inputs, targets, hprev):
    """Compute gradients of the loss function w.r.t. each parameter.

    Inputs, targets are both list of integers.
    hprev is initial hidden state
    returns a list of gradients for w,b,h
    """
    # Unpack parameters
    w, b, h0, y0, p, o, i, h, hprev, loss, next_h = hprev

    # Compute gradients
    grad_w = np.zeros_like(w)
    grad_b = np.zeros_like(b)
    grad_h0 = np.zeros_like(h0)
    grad_y0 = np.zeros_like(y0)
    grad_p = np.zeros_like(p)
    grad_o = np.zeros_like(o)
    grad_i = np.zeros_like(i)
    grad_h = np.zeros_like(h)
    grad_hprev = np.zeros_like(hprev)

    for t in range(len(inputs)-1, -1, -1):
      grad_y0 = np.zeros_like(y0)
      grad_p = np.zeros_like(p)
      grad_o = np.zeros_like(o)
      grad_i = np.zeros_like(i)
      grad_h = np.zeros_like(h)
      grad_hprev = np.zeros_like(hprev)

      if t < len(inputs)-1:
        grad_y0 = targets[t+1]
        grad_p = softmax(next_h) * np.exp(next_h)
        grad_o = np.tanh(next_h) * grad_p
        grad_i = np.tanh(next_h) * np.sinh(next_h) * grad_o
        grad_h = np.tanh(next_h) * np.cosh(next_h) * grad_i
        grad_hprev = np.tanh(next_h) * np.cosh(next_h) * grad_h

      grad_w += np.outer(grad_h, inputs[t])
      grad_b += grad_h
      grad_h0 += grad_h
      grad_y0 += grad_h
      grad_p += grad_h
      grad_o += grad_h
      grad_i += grad_h
      grad_h += grad_h
      grad_hprev += grad_h

    grad_w /= len(inputs)
    grad_b /= len(inputs)
    grad_h0 /= len(inputs)
    grad_y0 /= len(inputs)
    grad_p /= len(inputs)
    grad_o /= len(inputs)
    grad_i /= len(inputs)
    grad_h /= len(inputs)
    grad_hprev /= len(inputs)

    return grad_w, grad_b, grad_h0, grad_y0, grad_p, grad_o, grad_i, grad_h, grad_hprev, loss

```

(<https://gist.github.com/karpathy/d4dee566867f8291f086>)

2. Init of the model parameters and hyper-parameters

```
# hyperparameters
hidden_size = 100 # size of hidden layer of neurons
seq_length = 25 # number of steps to unroll the RNN for
learning_rate = 1e-1

# model parameters
Wih = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
bh = np.zeros((hidden_size, 1)) # hidden bias
by = np.zeros((vocab_size, 1)) # output bias
```

3. Forward pass

```

def lossFn(inputs, targets, hsPrev):
    """Compute the loss between target and predicted hidden states.

    inputs: a list of input vectors
    targets: a list of target vectors
    hsPrev: the initial hidden state
    Returns the loss, gradients on model parameters, and last hidden state
    """
    hs, vs, ps = ([], [], [])
    loss = 0
    hs.append(hsPrev)
    for i in range(len(inputs)):
        x = np.array(inputs[i])
        vs.append(np.copy(x))
        hs.append(np.copy(hs[-1]))
        loss += np.sum((hs[-1] - np.array(targets[i])) ** 2)
        hs[-1] = np.tanh(np.dot(hs[-1], w1) + np.dot(vs[-1], b1)) # hidden state
        vs[-1] = np.dot(w1.T, hs[-1]) + b1 # unnormalized log probability for next char
        ps.append(np.exp(vs[-1]) / np.sum(np.exp(vs[-1]))) # softmax layer
        hs[-1] = np.dot(w2, ps[-1]) + b2 # final hidden state
    return loss, hs, ps

```

4. Back pass

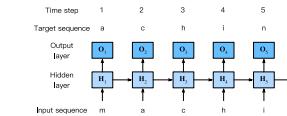
In practical environment: Pytorch or Pytorch Lightning



RNN in 100 lines of code

I trains a RNN at the character level with unfolding period of 25 chars, hidden layer size is 100.
We will model the input just with one-hot encoding; we will not use word embedding to keep it simple.

Requirements: numpy only



Six pages

1. I/O: corpus parsing and vocabulary construction
 2. Init of the model parameters and hyper-parameters
 3. Forward pass (input to loss computation)
 4. Backward pass (most complex part, Truncated BPTT)
 5. Main Training Loop
 6. Generation part (sampling) fun part! 🎉

1. I/O: corpus parsing and vocabulary construction

```
import numpy as
```

```
# data = open('input.txt', 'r').read() # should be simple plain text file
chars = list(set(data))
data_size, vocab_size = len(data), len(chars)
print(f'Data has {len(data)} characters. {len(unique)} unique. ({data_size}, {vocab_size})')
# map from character to index
char_to_ix = {ch: i for i, ch in enumerate(chars)}
# map from index to character
ix_to_char = [i for i, ch in enumerate(chars)]
```

5a. Main loop

5b. Main loop

88-11

```

# sample from the model now and then
sample = np.random.randint(0, 100, 100)
for i in range(100): print("for %d is sample %d" % (i, sample[i]))
```



```

def sample(b_size, n):
    """sample a sequence of integers from the model
    b_size many times, and n is next letter for first time
    """
    x = np.zeros((b_size, 1))
    i = 0
    for t in range(n):
        a = np.random.randint(0, 100, 1)
        b = np.random.randint(0, 100, 1) + a
        c = np.exp(a) / np.sum(np.exp(a))
        b = np.random.choice(np.arange(b_size), 1, p=c).ravel()
        x[i] = b
        i += 1
    return x

```

Details on the implementation of the backward pass

[Details on the backward pass: mktfl.github.io/2019/07/08/minimalist-RNN](#)

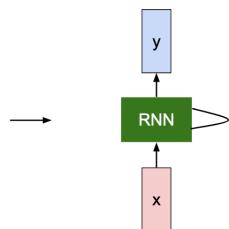
RNN learns to generate Sonnets

THE SONNETS

by William Shakespeare

Fare fairer creatures we desire no more,
That beauty bounte we might never see,
But as the riper should by time deserve,
So let me have some right to be desired.
But then, constrained to these own bright eyes,
That look on beauty, though they themselves
Holding a feature where abundance lies,
No feature like them where want lies,
Then being asked, where all thy beauty lies,
To see, within there own deep native eyes,
How beautie like her selfe doth beautie lie,
How much more proue deserved the beaute's eye,
Than eye deserved to beaute's selfe,
Shall my eye count, and make my old receve,
This were to be new made when thon art old,
And set by thond wane when thon forth it cold.

When forty winters shall besiege these bones,
Till you should be as green as grass we see,
The youth's proud liver so gaud to see,
Then being asked, where all thy beauty lies,
To see, within there own deep native eyes,
How beautie like her selfe doth beautie lie,
How much more proue deserved the beaute's eye,
Than eye deserved to beaute's selfe,
Shall my eye count, and make my old receve,
This were to be new made when thon art old,
And set by thond wane when thon forth it cold.



picture from Stanford

RNN learns to generate Sonnets

at first:

tyntd-lafhatwlaohrdemot lytdws e ,fti, astai f ogoh eaesse rrarbyne 'nhthnee e
pia tklrgd t o idoe ns,smth h ne etie h,hrgrts nigtike,aoeanns lng

| train more

"Tmont thithay" fomescerliund
Keushey. Thon here
sheulke, ammerenith ol sivh I laltherthend Bleipile shuyv fil on aseterlome
coainogennc Phe lism thond hon at. MeDimorotion in ther thize."

| train more

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

| train more

"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.

picture from Stanford

RNN learns to generate Sonnets

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

Second Senator:
They are awry this miserles, 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 nose begin 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.

picture from Stanford

...or generate Latex

The Stacks Project: open source algebraic geometry textbook

The Stacks Project	
home	about
tags explained	tag lookup
browse	search
bibliography	recent comments
blog	add slogan
Browse chapters	
Part	Chapter
Preliminaries	online Tex source view pdf
1. Introduction	online tex pdf
2. Conventions	online tex pdf
3. Set Theory	online tex pdf
4. Categories	online tex pdf
5. Topology	online tex pdf
6. Sheaves on Spaces	online tex pdf
7. Sites and Sheaves	online tex pdf
8. Stacks	online tex pdf
9. Fields	online tex pdf
10. Commutative Algebra	online tex pdf

<http://stacks.math.columbia.edu/>
The stacks project is licensed under the [GNU Free Documentation License](#)

picture from Stanford

Latex source

[...or generate Latex](#)

<p>For $\bigoplus_{i=1}^n U_i = 0$, where $C_{\alpha\beta} = 0$, hence we can find a closed subset N in \mathbb{N} and any sets J on X', U is a closed immersion of S, then $U + T$ is a separated algebraic space.</p> <p><i>Proof.</i> Proof of (1). It also start we get $S = \text{Spec}(R) = U \times_X U \times_X U$ and the compatibility in the fiber product commutes we have to prove the lemma generates by $\coprod_{i=1}^n U_i = U$ in the fiber category of S in \mathcal{C}. Moreover we can take any U affine, see Morphism Lemma 77. Hence we obtain a scheme S and any open subset $W \subset U$ in $\text{Sh}(G)$ and $W \rightarrow S$ is smooth or an open subspace $Z \subset X$ of X' where Z in X' is proper (some defining as a closed immersion, some defining as a closed subspace). Then Z is a closed subspace of X.</p> <p>(1) f is locally of finite type. Since $S = \text{Spec}(R)$ and $Y = \text{Spec}(R)$.</p> <p>which has a nonzero morphism we may assume that f_i is of finite presentation over S. We claim that $\mathcal{O}_{X,i}$ is a scheme where $x, x', x'' \in S$ such that $\mathcal{O}_{X,i'} \rightarrow \mathcal{O}_{X,i''}$ is injective. By Algebra, Lemma 10.12 we can endow a map of complexes $\text{GL}_{\mathcal{O}}(r/S)^{\wedge i}$ and we will prove it.</p> <p>$\tilde{M} = \mathbb{Z}^n \otimes_{\mathcal{O}_{X,i}} \mathcal{O}_{X,i} \rightarrow i_2^* \mathcal{F}$</p> <p>is a unique morphism of algebraic stacks. Note that $\text{Arrows} = (\mathbf{Sht}/S)^{\text{opp}} = (\mathbf{Sch}/S)_{/\text{perf}}$ and</p> <p>$V = \Gamma(S, \mathcal{O}) \rightarrow \Gamma(\text{Spec}(A))$</p> <p>is an open subset of X. Thus U is affine. This is a continuous map of X is the inverse, the groupoid scheme S.</p> <p><i>Proof.</i> See discussion of sheaves of sets.</p> <p>The result follows from open covering follows from the less of Example 77. If U is regular, then by (1) which gives an open subspace of X and T equal to S_{reg} we see Descent, Lemma 77. Namely, by Lemma 77 we see that R is geometrically regular over S.</p>	<p>Lemma 0.1. Assume (2) and (3) by the construction in the description. Suppose $X = \lim(X)$ is the formal open covering X and a single map $\text{Proj}_X(A) = \text{Spec}(B)$ over U compatible with the covering.</p> <p>Since $\mathcal{F} = \varprojlim \mathcal{F}_i$ in the second conditions of (1), and (3). This finishes the proof. By Definition 77 (further element in the fiber category of S in \mathcal{C}) T is surjective we can take any U to U in $\text{Sh}(G)$ with residue field of S. Moreover we can take any U affine, see Morphism Lemma 77. Hence we obtain a scheme S and any open subset $W \subset U$ in $\text{Sh}(G)$ and $W \rightarrow S$ is smooth or an open subspace $Z \subset X$ of X' where Z in X' is proper (some defining as a closed immersion, some defining as a closed subspace). Then Z is a closed subspace of X.</p> <p>(1) f is locally of finite type. Since $S = \text{Spec}(R)$ and $Y = \text{Spec}(R)$.</p> <p>This is form all sheaves of sheaves on X. But given a scheme U and a respective site-morphism $U \rightarrow X$. Let $U/U = \coprod_{i=1}^n U_i$ be the scheme X over S at the schemes $X_i \rightarrow X$ and $U = \lim(X_i)$.</p> <p>The following lemma recursive retrocomposes this implies that $\mathcal{F}_{x_0} = \mathcal{F}_{x_n}$ for $i > 0$ and \mathcal{F}_i exists and let \mathcal{F}_i be a prelocal of \mathcal{O}_X-modules on \mathcal{C} as \mathcal{F}-module. In particular $\mathcal{F} = \cup \mathcal{F}_i$.</p> <p>Lemma 0.2. Let X be a locally Noetherian scheme over S, $E = \mathcal{F}_{X/S}$. Set $\mathcal{I} = \mathcal{J}_i \subset \mathcal{I}_n$. Since $\mathcal{I} \subset \mathcal{I}'$ are nonzero over $i \leq p$ is a subset of $\mathcal{J}_0 \cup \mathcal{J}_1$ works.</p> <p>Lemma 0.3. In Situation 77. Hence we may assume $q = 0$.</p> <p><i>Proof.</i> We start the proof by we see that j is the next functor (?)?. On the other hand, by Lemma 77 we see that</p> <p>$D(\mathcal{O}_V) = \mathcal{O}_X(D)$</p> <p>where K is an F-algebra where δ_{n+1} is a scheme over S.</p>
---	--

ecture from Stanford

[...or C GNU/Linux Kernel code](#)

```
static void do_command(struct seq_file *m, void *v)
{
    int column = 32 << (cmd[2] & 0x80);
    if (!state)
        cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & 2) : 1);
    else
        seq = 1;
    for (i = 0; i < 16; i++) {
        if (K & (1 << i))
            pipe = (in_use & UNTHREAD_UNCA) +
                (count & 0xfffffffffffffe) & 0x00000000;
        if (count == 0)
            subpid, pcp_md.xeccc_handle, 0x20000000);
        pipe_set_byter(i, 0);
    }
    /* Free our user page pointer to place camera if all dash */
    subsystem_info = acf_change(PAGE_SIZE);
    subsystem_info |= 0x0000000000000000;
    /* Now we want to deliberately put it to device */
    control_check_polarity(&context, val, 0);
    for (i = 0; i < COUNTER; i++)
        seq_putz(s, "policy");
}
```

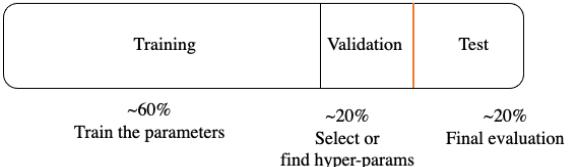
ecture from Stanford

[How to evaluate a LM](#)

[LM Evaluation \(you already seen it in Part I\)](#)

[Train on train split](#)
[Select hyper-params on validation or dev set](#)
[Test on the test split](#)

Labeled data



[Set size and partitioning: not a clear definition](#)

RNNs greatly improved perplexity over prior art

Informally how well you predict next words on a text corpus W you never trained on. **PP** the lower the better.

Model	Test perplexity
Interpolated Kneser-Ney 5-gram (Chelba et al., 2013)	67.6
Feedforward NN + D-Softmax (Chen et al., 2015)	91.2
4-layer IRNN-512 (Le et al., 2015)	69.4
RNN-2048 + BlackOut sampling (Ji et al., 2015)	68.3
Sparse Non-negative Matrix Language Model (Shazeer et al., 2015)	52.9
RNN-1024 + MaxEnt 9-gram (Chelba et al., 2013)	51.3
LSTM-2048-512 (Jozefowicz et al., 2016)	43.7
2-layer LSTM-8192-1024 + CNN inputs (Jozefowicz et al., 2016)	30.0
Ours (LSTM-2048)	43.9
Ours (2-layer LSTM-2048)	39.8

Table 2. One Billion Word benchmark. Perplexity on the test set for single models. Our result is obtained after 5 epochs.

Efficient softmax computation for GPUs (<https://arxiv.org/pdf/0908.0459.pdf>)

Applications of Sequence Modeling: Stacked Bi-directional RNN

Today's lecture

-Recap on RNN

-We will see a few application of RNN in NLP and vision

-Bidirectional RNN and Stacked RNN

This lecture material is taken from

[Chapter 9 Jurafsky Book](#)

[Chapter 4.3 Eisner Book](#)

- [Stanford Slides RNN](#)
- [Stanford Lecture RNN](#)
- [Stanford Notes on Word2Vec](#)
- [Andrea Esuli's notes on RNN](#)
- [Another Slides about RNN](#)

Another resource with code is [\[GitHub\]](#)

Generated
C code

Recap on RNN and LM

Language Model: A system that predicts the next word

Recurrent Neural Network: A family of neural networks that:

- Take an initial input or state
- Apply the same weights on each step
- Can optionally produce output on each step
- Recurrent Neural Networks = Language Model

We have shown that RNNs are a great way to build a LM (despite some problems).

RNNs are also useful for much more! We see it today.

Why should we care about Language Modeling

Language Modeling is a benchmark task that helps us measure our progress on predicting language use.

Language Modeling is a component of many NLP tasks, especially those involving generating text or estimating the probability of text.

- Predictive text
- Speech recognition
- Handwriting recognition
- Document classification
- Entity identification
- Machine translation
- Summarization
- Dialogue

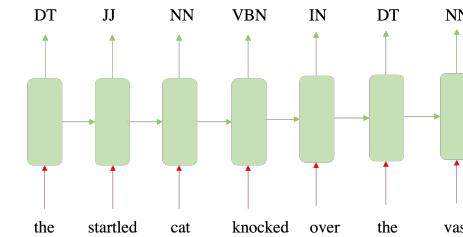
Everything else in NLP has now been built upon Language Modeling. GPT-3 is an LM! Though remember that NLP ≠ LMs.

RNN can be used for...

RNN can be used for Sequence Tagging

Solve problems such as part-of-speech (POS) tagging, named entity recognition

You can solve it with RNN as a many-to-many method with supervision at the word level



pos is the process of assigning a part-of-speech to each word in part-of-speech tagging a text. Tagging is a disambiguation task: words are ambiguous.

Good rule for Deep Learning Researchers

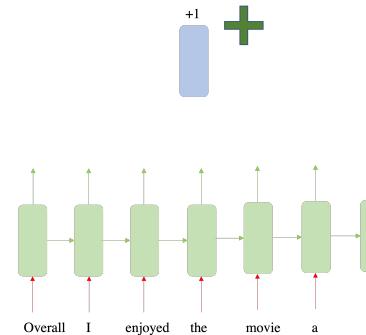
Most Frequent Class Baseline:

||| Most Frequent Class Baseline: Always compare a classifier against a baseline at least as good as the most frequent class baseline, *i.e.* every token in the class is covered in most often in the training set.

RNN can be used for Sentence Classification

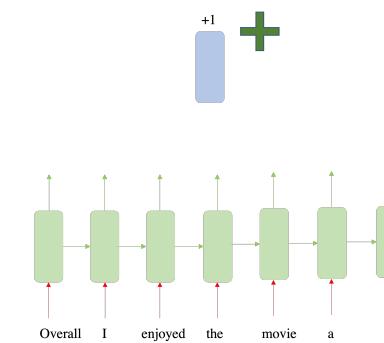
Solve problems such as sentiment classification or sentiment analysis

You can solve it with RNN as a many-to-one method with supervision at the word level



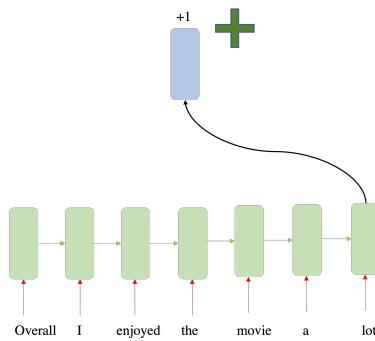
RNN can be used for Sentence Classification

How to compute sentence encoding?



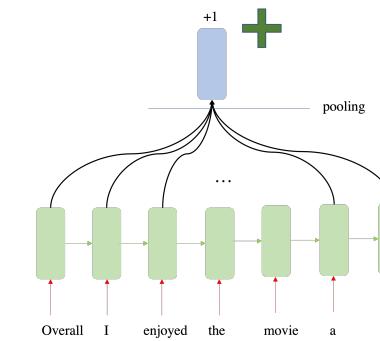
RNN can be used for Sentence Classification

[How to compute sentence encoding?](#) Simple: take last hidden state vector



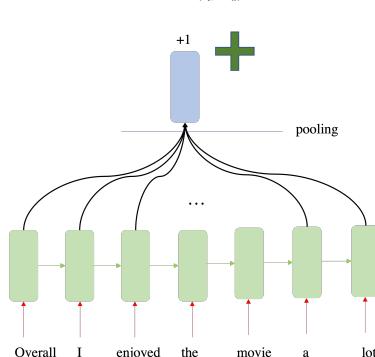
RNN can be used for Sentence Classification

[How to compute sentence encoding?](#) Simple: do pooling (i.e. average of all intermediate hidden states)



RNN can be used for Encoding for other tasks

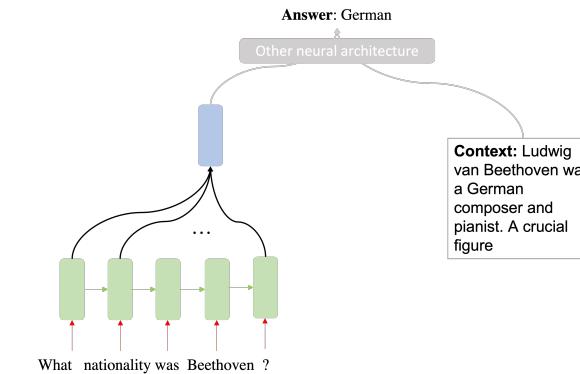
$x = \text{encoder}(w_1, \dots, w_N)$



RNN can be used for Encoding for other tasks

[Useful for: question answering, machine translation](#)

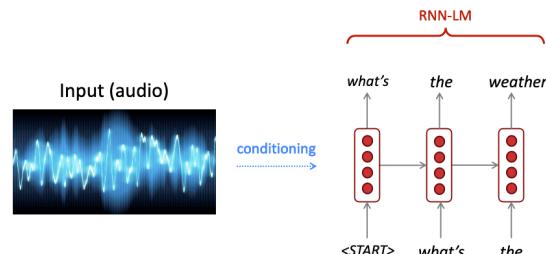
$x = \text{encoder}(w_1, \dots, w_N)$



RNN can be used for Decoding for other tasks

x = decoder(encoder(input), < start >)

Useful for speech recognition, machine translation, summarization, image captioning
This is an example of a conditional language model. The LM is conditioned on the speech representation.



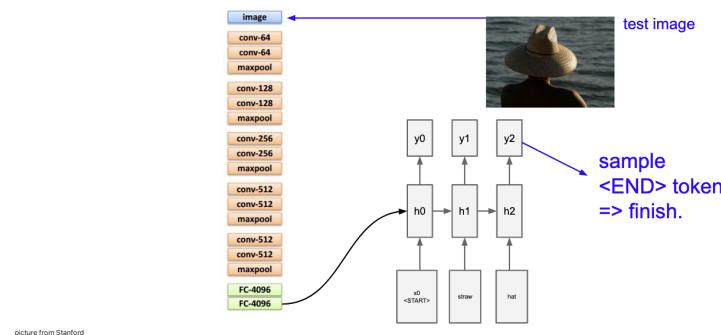
picture from Stanford

RNN as decoder: Image Captioning

Given an image x return a sentence that describe the image. We need:

1. Something that gets a $f(x)$ representation of the image, e.g. $f(x)$ where f can be a ConvNet trained on ImageNet or with Self-Supervision.
2. Conditioner of $f(x)$, namely $\{s_1, s_2, \dots, s_n\}$ we train an RNN to match the ground-truth sentence.
3. $s_i \rightarrow \{w_1, \dots, w_n\}$

RNN as decoder: Image Captioning



picture from Stanford

RNN as decoder: Image Captioning

Image Captioning: Example Results



picture from Stanford

RNN as decoder: Image Captioning

Image Captioning: Failure Cases



picture from Stanford

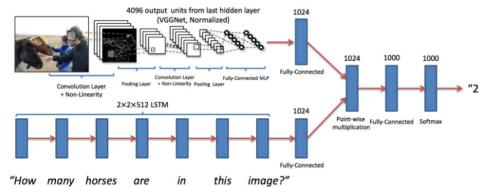
RNN as Encoder: Visual Question Answering



Agrawal et al., "VQA: Visual Question Answering", ICCV 2015
Zhu et al., "Visual 7W: Grounded Question Answering in Images", CVPR 2016
Figure from Zhu et al., copyright IEEE 2016. Reproduced for educational purposes.

Figure from Stanford

RNN as Encoder: Visual Question Answering

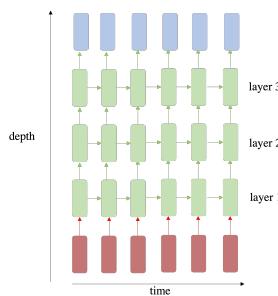


Agrawal et al., "Visual 7W: Grounded Question Answering in Images", CVPR 2015
Figure from Agrawal et al., copyright IEEE 2015. Reproduced for educational purposes.

Figure from Stanford

Stacked RNN (or Multi-layer RNN)

The hidden states from RNN layer i are the inputs to RNN layer $i + 1$.

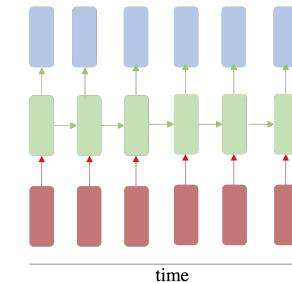


Stacked RNN (or Multi-layer RNN)

- Multi-layer or stacked RNNs allow a network to compute more complex representations; they work better than just having one layer of high-dimensional encodings.
- Hand-designed RNNs are usually multi-layer, BUT NOT as deep as convolutional or feed-forward networks
- In a 2017 paper, Bratt et al. find that for Neural Machine Translation, 2 to 4 layers is best for the encoder RNN, and 4 layers is best for the decoder RNN.
 - Often 2 layers is a lot better than 1, and 3 might be a little better than 2
 - Highly skip connections (dense connections) are needed to train deeper RNNs (e.g. 8 layers).

Stacked Bidirectional RNN

RNN so far



Stacked RNN (or Multi-layer RNN)

- RNNs are already “deep” on one dimension, the time dimension – they unroll over many timesteps.
- We can also make them “deep” in another dimension (the representation dimension)! We can do so by applying multiple RNNs – this is a multi-layer RNN (stacked RNN).
- This allows the network to compute more complex representations.
- The lower RNNs should compute lower-level features and the higher RNNs should compute higher-level features.