

CS6101 - Deep Learning for NLP

Week 5: Recurrent Neural Networks & Language Models

Kee Yuan Chuan Zhou Yizhuo Wu Jiacheng Liu Juncheng Jin Zhe



Today we will:

- Introduce a new NLP task
 - Language Modeling

motivates

- Introduce a new family of neural networks
 - Recurrent Neural Networks (RNNs)

THE most important idea for the rest of the class!

Language Modeling

Language Modeling is the task of predicting what word comes next

the students opened their



• More formally: given a sequence of words $x^{(1)}, x^{(2)}, \ldots, x^{(t)}$, compute the probability distribution of the next word $x^{(t+1)}$:

$$P(\boldsymbol{x}^{(t+1)} = \boldsymbol{w}_j \mid \boldsymbol{x}^{(t)}, \dots, \boldsymbol{x}^{(1)})$$

where $oldsymbol{w}_j$ is a word in the vocabulary $V = \{oldsymbol{w}_1,...,oldsymbol{w}_{|V|}\}$

• A system that does this is called a Language Model.

You use Language Models every day!



You use Language Models every day!



what is the			Ļ
what is the weather what is the meanin what is the dark we what is the xfl what is the doomse what is the doomse what is the weather what is the keto die what is the speed o what is the bill of r	r g of life eb day clock r today et an dream of light ights		
	Google Search	I'm Feeling Lucky	

n-gram Language Models

the students opened their

- **<u>Question</u>**: How to learn a Language Model?
- <u>Answer</u> (pre- Deep Learning): learn a *n*-gram Language Model!
- <u>Definition</u>: A *n*-gram is a chunk of *n* consecutive words.
 - unigrams: "the", "students", "opened", "their"
 - bigrams: "the students", "students opened", "opened their"
 - trigrams: "the students opened", "students opened their"
 - 4-grams: "the students opened their"
- <u>Idea:</u> Collect statistics about how frequent different n-grams are, and use these to predict next word.

n-gram Language Models

First we make a simplifying assumption: x^(t+1) depends only on the preceding (n-1) words

λ

$$P(\boldsymbol{x}^{(t+1)}|\boldsymbol{x}^{(t)},\ldots,\boldsymbol{x}^{(1)}) = P(\boldsymbol{x}^{(t+1)}|\boldsymbol{x}^{(t)},\ldots,\boldsymbol{x}^{(t-n+2)})$$

(assumption)

prob of a n-gram
$$= P(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})$$
 (definition of prob of a (n-1)-gram
$$= P(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})$$
 conditional prob)

- <u>Question</u>: How do we get these *n*-gram and (*n*-1)-gram probabilities?
- **Answer:** By counting them in some large corpus of text!

$$\approx \frac{\operatorname{count}(\boldsymbol{x}^{(t+1)}, \boldsymbol{x}^{(t)}, \dots, \boldsymbol{x}^{(t-n+2)})}{\operatorname{count}(\boldsymbol{x}^{(t)}, \dots, \boldsymbol{x}^{(t-n+2)})}$$
 (statistical approximation)

Why is it an approximation?

 $P(x^{(t+1)},x^{(t)}...,x^{(t-n+2)})$ is the joint probability of the words sequence

"Out of sequence of length n, how many of them is the sequence of interest?"

Recall conditional probability of 2 random variables, x and a known y is

$$P(X|Y = y_i) = \frac{P(X, Y = y_i)}{P(Y = y_i)}$$

Rearranging

$$P(X,Y) = P(X|Y=y_i)P(Y=y_i)\;\;$$
 aka chain rule

Generalising

lising
$$P(X_1...X_n) = \prod_{k=1}^n P(X_k | X_1^{k-1})$$

P("book", "their", "opened", "students") = P("book"|"their", "opened", "students")P("their"|"opened", "students")P("opened"|"students")P("students

https://lagunita.stanford.edu/c4x/Engineering/CS-224N/asset/slp4.pdf

n-gram Language Models: Example



 $P(w_j | \text{students opened their}) = \frac{\text{count}(\text{students opened their})}{\text{count}(\text{students opened their})}$

In the corpus:

- "students opened their" occurred 1000 times
- "students opened their books" occurred 400 times
 - \rightarrow P(books | students opened their) = 0.4
- "students opened their exams" occurred 100 times
 - \rightarrow P(exams | students opened their) = 0.1

Should we have discarded the "proctor" context?

Problems with n-gram Language Models

Sparsity Problem 1

Problem: What if *"students* opened their w_j " never occurred in data? Then w_i has probability 0!

(Partial) Solution: Add small δ to count for every $w_j \in V$. This is called *smoothing*.

count(students opened their \boldsymbol{w}_i) $P(\boldsymbol{w}_i|\text{students opened their}) =$

count(students opened their)

Sparsity Problem 2

Problem: What if *"students* opened their" never occurred in data? Then we can't calculate probability for any w_i !

(Partial) Solution: Just condition on "opened their" instead. This is called *backoff*.

Note: Increasing *n* makes sparsity problems *worse*. Typically we can't have *n* bigger than 5.

Problems with n-gram Language Models



Increasing *n* makes model size huge!

n-gram Language Models in practice

 You can build a simple trigram Language Model over a 1.7 million word corpus (Reuters) in a few seconds on your laptop*



Otherwise, seems reasonable!

https://nlpforhackers.io/language-models/

• You can also use a Language Model to generate text.



• You can also use a Language Model to generate text.



• You can also use a Language Model to generate text.



• You can also use a Language Model to generate text.

today the price of gold

• You can also use a Language Model to generate text.

today the price of gold per ton , while production of shoe lasts and shoe industry , the bank intervened just after it considered and rejected an imf demand to rebuild depleted european stocks , sept 30 end primary 76 cts a share .

Incoherent! We need to consider more than 3 words at a time if we want to generate good text.

But increasing *n* worsens sparsity problem, and exponentially increases model size...

How to build a *neural* Language Model?

- Recall the Language Modeling task:
 - Input: sequence of words $oldsymbol{x}^{(1)},oldsymbol{x}^{(2)},\ldots,oldsymbol{x}^{(t)}$
 - Output: prob dist of the next word $P(\boldsymbol{x}^{(t+1)} = \boldsymbol{w}_j \mid \boldsymbol{x}^{(t)}, \dots, \boldsymbol{x}^{(1)})$
- How about a window-based neural model?
 - We saw this applied to Named Entity Recognition in Lecture 4



the students opened their ______

- A.K.A normalized exponential function
- is a generalization of the *logistic function* AKA *Sigmoid function*
- Used to provide a probabilistic interpretation of classification prediction
- More appropriate for mutually-exclusive classes because during training, values of correct classes pushed towards +ive inf, values of wrong classes pushed towards -ive inf ⇒ Only one correct class, all other classes are wrong!

Say we have 4 classes, so final layer fan-out is 4



$$\sigma(z)_j = \frac{e^{z_j}}{\sum_{k=1}^k e^{z_k}} \text{ for } j = 1, \dots, k$$

Question: Why do we use output values as exponent to the natural log base?

Hint: Think about the value range of layer output values

In practice...

$$p_i = rac{e^{a_i}}{\sum_{k=1}^N e^{a_k}} \ = rac{Ce^{a_i}}{C\sum_{k=1}^N e^{a_k}} \ = rac{e^{a_i+\log(C)}}{\sum_{k=1}^N e^{a_k+\log(C)}}$$

Where *log(C)* is taken to be *-max(a)*. Why?

Cross Entropy

 A measure of distance between what the model believes the output distribution should be ("unnatural" distribution), and the original distribution ("natural" distribution)

For discrete outcomes:



Cross Entropy
$$H(p,q) = -\sum_x p(x) \, \log q(x).$$

Say for a given input, the correct class is 3 out of classes {1,2,3,4}

$$\left(\begin{array}{c}
0\\
0\\
1\\
0
\end{array}\right) \bullet \left(\begin{array}{c}
0.23\\
0.01\\
0.67\\
0.67\\
0.09
\end{array}\right) = -1 \times 0.67 = -0.67$$

Cross Entropy $H(p,q) = -\sum_{x} p(x) \log q(x).$ $= -\log q(x)$

output distribution

$$\hat{oldsymbol{y}} = ext{softmax}(oldsymbol{U}oldsymbol{h} + oldsymbol{b}_2) \in \mathbb{R}^{|V|}$$

hidden layer $h = f(We + b_1)$

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

words / one-hot vectors $oldsymbol{x}^{(1)},oldsymbol{x}^{(2)},oldsymbol{x}^{(3)},oldsymbol{x}^{(4)}$



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 $oldsymbol{W}$
- Window can never be large enough!
- Each x⁽ⁱ⁾uses different rows cols of W. We don't share weights across the window.

We need a neural architecture that can process any length input



How does our weight matrix look like?

 $\times e =$ $W_{1,1} \dots W_{1,d} \dots \dots W_{1,nd}$ X W_{fan_out,1} fan out.nd fan $out \times nd$ fan out $\times 1$ *n* is window size $nd \times 1$

d dimensions in each word embedding

Question: Why do we want to share weights across the window?

Why do we want to share weights across windows?

- Else the number of weights would grow linearly with the number of time steps (our window size) like a feed-forward network
- We want to capture shared representations across sequences of text

Recurrent Neural Networks (RNN)

A family of neural architectures



Recurrent Neural Networks (RNN) A family of neural architectures

Core idea: Apply the same weights *W* repeatedly

visualized as feedforward networks "unrolled across time":



Recurrent Neural Networks (RNN)

How does our weight matrix $W_{\rm h}$ look like?

 $\times h =$ h $\boldsymbol{W}_{1,1} \, \boldsymbol{W}_{1,2} \, \dots \, \dots$ h_1 1,fan out X h
fan_ou fan_out.1 ··· ··· W fan_out.fan_out fan out $\times 1$ fan out $\times 1$ fan out × fan out

Recurrent Neural Networks (RNN)

How does our weight matrix $W_{\rm h}$ look like?

$W_{a} \times e =$		
$\left(\begin{array}{c} \boldsymbol{e} \\ \boldsymbol{W}_{1,1} \boldsymbol{W}_{1,2} \dots \boldsymbol{W}_{1,d} \end{array}\right)$	$\begin{bmatrix} e_1 \end{bmatrix}$	<pre></pre>
	×	=
fan_out,1 ··· ·· fan_out,d	$\left(\begin{array}{c} e_{d} \end{array} \right)$	
$fan_out \times d$	$d \times 1$	fan_out × 1




Question: So what is *recurrent* about it?



A 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

More on these next week

 $oldsymbol{h}^{(0)}$

 \bigcirc

 \bigcirc



- Get a big corpus of text which is a sequence of words $m{x}^{(1)},\ldots,m{x}^{(T)}$
- Feed into RNN-LM; compute output distribution ŷ^(t) for every step t.
 i.e. predict probability dist of every word, given words so far
- Loss function on step t is usual cross-entropy (CE) between our predicted probability distribution $\hat{y}^{(t)}$, and the true next word : $y^{(t)} = x^{(t+1)}$

$$J^{(t)}(\theta) = CE(\boldsymbol{y}^{(t)}, \hat{\boldsymbol{y}}^{(t)}) = -\sum_{j=1}^{|V|} y_j^{(t)} \log \hat{y}_j^{(t)}$$

Average this to get overall loss for entire training set:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta)$$





•••







- However: Computing loss and gradients across entire corpus is too expensive!
- <u>Recall:</u> Stochastic Gradient Descent allows us to compute loss and gradients for small chunk of data, and update.
- \rightarrow In practice, consider $x^{(1)}, \ldots, x^{(T)}$ as a sentence

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta)$$

• Compute loss $J(\theta)$ for a sentence (actually usually a batch of sentences), compute gradients and update weights. Repeat.

Example: Character-level Language Model

Vocabulary: [h,e,l,o]

Example training sequence: "hello"



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 36 University 017



Imagine training the entire wikipedia corpus to just get 1 gradient update...

Backpropagation for RNNs



<u>Question</u>: What's the derivative of $J^{(t)}(\theta)$ w.r.t. the repeated weight matrix W_h ?

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

"The gradient w.r.t. a repeated weight is the sum of the gradient w.r.t. each time it appears"

Why?

Backpropagation Review

http://cs231n.github.io/optimization-2/#staged

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:

$$\underbrace{\frac{d}{dt}f(\boldsymbol{x}(t),\boldsymbol{y}(t))}_{dt} = \frac{\partial f}{\partial \boldsymbol{x}}\frac{d\boldsymbol{x}}{dt} + \frac{\partial f}{\partial \boldsymbol{y}}\frac{d\boldsymbol{y}}{dt}$$

Derivative of composition function



Source:

https://www.khanacademy.org/math/multivariable-calculus/multivariable-derivatives/differentiating-vector-valued-functions/a/multivariable-chain-r ule-simple-version

Backpropagation for RNNs: Proof sketch

- 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:

$$rac{d}{dt} f(x(t), y(t)) = rac{\partial f}{\partial x} rac{dx}{dt} + rac{\partial f}{\partial y} rac{dy}{dt}$$

Derivative of composition function



Source:

https://www.khanacademy.org/math/multivariable-calculus/multivariable-derivatives/differentiating-vector-valued-functions/a/multivariable-chain-r ule-simple-version

Backpropagation for RNNs



we

$$\frac{\partial J^{(t)}}{\partial W_{h}} = \left[\sum_{i=1}^{t} \frac{\partial J^{(t)}}{\partial W_{h}} \right|_{(i)}$$
Question: How do calculate this?

Answer: Backpropagate over timesteps *i*=*t*,...,0, summing gradients as you go. This algorithm is called **"backpropagation through time"**

37

BPTT(backpropagation through time)

$$egin{aligned} s_t &= tanh(Ux_t + Ws_{t-1}) \ \hat{y}_t &= softmax(Vs_t) \end{aligned}$$

$$egin{aligned} E_t(y_t, \hat{y}_t) &= -y_t \log \hat{y}_t \ E(y, \hat{y}) &= \sum_t E_t(y_t, \hat{y}_t) \ &= -\sum_t y_t \log \hat{y}_t \end{aligned}$$

BPTT(backpropagation through time)



 $egin{aligned} rac{\partial E_3}{\partial V} &= rac{\partial E_3}{\partial {\hat y}_3} rac{\partial {\hat y}_3}{\partial V} \ &= rac{\partial E_3}{\partial {\hat y}_3} rac{\partial {\hat y}_3}{\partial z_3} rac{\partial z_3}{\partial V} \ &= ({\hat y}_3 - y_3) \otimes s_3 \end{aligned}$

 $rac{\partial E_3}{\partial W} = rac{\partial E_3}{\partial {\hat y}_3} rac{\partial {\hat y}_3}{\partial s_3} rac{\partial s_3}{\partial W}$

 $rac{\partial E_3}{\partial W} = \sum_{k=0}^3 rac{\partial E_3}{\partial {\hat y}_3} rac{\partial {\hat y}_3}{\partial s_3} rac{\partial s_3}{\partial s_k} rac{\partial s_k}{\partial W}$

BPTT(backpropagation through time)



```
def bptt(self, x, y):
    T = len(y)
    # Perform forward propagation
    o, s = self.forward_propagation(x)
    # We accumulate the gradients in these variables
    dLdU = np.zeros(self.U.shape)
    dLdV = np.zeros(self.V.shape)
    dLdW = np.zeros(self.W.shape)
    delta o = o
    delta_0[np.arange(len(y)), y] = 1.
    # For each output backwards...
    for t in np.arange(T)[::-1]:
        dLdV += np.outer(delta_o[t], s[t].T)
        # Initial delta calculation: dL/dz
        delta_t = self.V.T.dot(delta_o[t]) * (1 - (s[t] ** 2))
        # Backpropagation through time (for at most self.bptt_truncate steps)
        for bptt_step in np.arange(max(0, t-self.bptt_truncate), t+1)[::-1]:
            # print "Backpropagation step t=%d bptt step=%d " % (t, bptt_step)
            # Add to gradients at each previous step
            dLdW += np.outer(delta_t, s[bptt_step-1])
            dLdU[:,x[bptt_step]] += delta_t
            # Update delta for next step dL/dz at t-1
            delta_t = self.W.T.dot(delta_t) * (1 - s[bptt_step-1] ** 2)
    return [dLdU, dLdV, dLdW]
```

Vanishing/exploding gradient problem

- Multiply the same W at each time step during backprop.
- The gradient is a product of Jacobian matrices, each associated with a step in the forward computation. This can become very small or very large quickly, and the locality assumption of gradient descent breaks down. → Vanishing or exploding gradient.
- Gradients can be seen as a measure of influence of the past on the future.
- The vanishing gradient problem can cause problems: When predicting the next word, information from many time steps in the past is not taken into consideration.

Vanishing/exploding gradient problem

- Gradient clipping method
- Initialization and ReLus
- Gated Recurrent Units (GRU) introduced by [Cho et al. 2014] and LSTMs [Hochreiter & Schmidhuber, 1999]

Truncated BPTT

Truncated BPTT processes the sequence one timestep at a time, and every k1 timesteps, it runs BPTT for k2 timesteps, so a parameter update can be cheap if k2 is small.

- k1: The number of forward-pass timesteps between updates. Generally, this influences how slow or fast training will be, given how often weight updates are performed.
- k2: The number of timesteps to which to apply BPTT. Generally, it should be large enough to capture the temporal structure in the problem for the network to learn. Too large a value results in vanishing gradients.

TBPTT (k1, k2)

- **TBPTT(n,n)**: Updates are performed at the end of the sequence across all timesteps in the sequence (e.g. classical BPTT).
- **TBPTT(1,n)**: timesteps are processed one at a time followed by an update that covers all timesteps seen so far (e.g. classical TBPTT by Williams and Peng).
- **TBPTT(k1,1)**: The network likely does not have enough temporal context to learn, relying heavily on internal state and inputs.
- **TBPTT(k1,k2), where k1<k2<n**: Multiple updates are performed per sequence which can accelerate training.
- **TBPTT(k1,k2), where k1=k2**: A common configuration where a fixed number of timesteps are used for both forward and backward-pass timesteps (e.g. 10s to 100s).

Prepare sequence data

The way that you break up your sequence data will define the number of timesteps used in the forward and backward passes of BPTT.

- Use data as-is
- Naive data split
- Domain-specific data split
 - In natural language processing problem, the input sequence could be divided by sentence and then padded to a fixed length, or split according to the average sentence length in the domain.
- Systematic data split (e.g. grid search)
 - perform a grid search over each sub-sequence length and adopt the configuration that results in the best performing model on average.
- Lean heavily on internal states with TBPTT(1, 1)

Pseudo code truncated version of BPTT

Back_Propagation_Through_Time(a, y) // a[t] is the input at time t. y[t] is the output

Unfold the network to contain k instances of f

do until stopping criteria is met:

x = the zero-magnitude vector;// x is the current context

for t from 0 to n - k // t is time. n is the length of the training sequence

Set the network inputs to x, a[t], a[t+1], ..., a[t+k-1]

p = forward-propagate the inputs over the whole unfolded network<math>e = y[t+k] - p; // error = target - prediction Back-propagate the error, e, back across the whole unfolded network Sum the weight changes in the k instances of f together. Update all the weights in f and g.

x = f(x, a[t]); // compute the context for the next time-step

RNN Applications

- Language Model
- Sentiment Classification
- Machine Translation
- Question Answering
- Speech recognition
- Time series prediction

Just like a n-gram Language Model, you can use a RNN Language Model to generate text by repeated sampling. Sampled output is next step's input.



- Let's have some fun!
- You can train a RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on Obama speeches:



The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done.

- Let's have some fun!
- You can train a RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on *Harry Potter*:



"Sorry," Harry shouted, panicking—"I'll leave those brooms in London, are they?"

"No idea," said Nearly Headless Nick, casting low close by Cedric, carrying the last bit of treacle Charms, from Harry's shoulder, and to answer him the common room perched upon it, four arms held a shining knob from when the spider hadn't felt it seemed. He reached the teams too.

Source:

https://medium.com/deep-writing/harry-potter-written-by-artificial-intelligence-8a9431803d a6

- Let's have some fun!
- You can train a RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on Seinfeld scripts:



He slams his hand on the door. KRAMER enters dancing with garbage.

KRAMER Hey hey hey, great idea for a big sponge: Make it so large you think it's got a fat clock in the middle.

JERRY (takes off his bones) Kramer, do you have a fun flashback to do?

Source:

https://www.avclub.com/a-bunch-of-comedy-writers-teamed-up-with-a-computer-to-18186332 42

- Let's have some fun!
- You can train a RNN-LM on any kind of text, then generate text in that style.
- (character-level) RNN-LM trained on paint colors:

Ghasty Pink 231 137 165
Power Gray 151 124 112
Navel Tan 199 173 140
Bock Coe White 221 215 236
Horble Gray 178 181 196
Homestar Brown 133 104 85
Snader Brown 144 106 74
Golder Craam 237 217 177
Hurky White 232 223 215
Burf Pink 223 173 179
Rose Hork 230 215 198

Sand Dan 201 172 143
Grade Bat 48 94 83
Light Of Blast 175 150 147
Grass Bat 176 99 108
Sindis Poop 204 205 194
Dope 219 209 179
Testing 156 101 106
Stoner Blue 152 165 159
Burble Simp 226 181 132
Stanky Bean 197 162 171
Turdly 190 164 116

Source: http://aiweirdness.com/post/160776374467/new-paint-colors-invented-by-neural-netwo rk

Evaluating Language Models

• The traditional evaluation metric for Language Models is perplexity.



Inverse probability of dataset

- Lower is better!
- minimizing perplexity and minimizing the loss function are equivalent.
- $log(PP) = J(\theta)$

RNNs have greatly improved perplexity

	Model	Perplexity
<i>n</i> -gram model—	→ Interpolated Kneser-Ney 5-gram (Chelba et al., 2013)	67.6
	RNN-1024 + MaxEnt 9-gram (Chelba et al., 2013)	51.3
crossingly	RNN-2048 + BlackOut sampling (Ji et al., 2015)	68.3
omplex RNNs	Sparse Non-negative Matrix factorization (Shazeer et al., 2015)	52.9
	LSTM-2048 (Jozefowicz et al., 2016)	43.7
	2-layer LSTM-8192 (Jozefowicz et al., 2016)	30
	Ours small (LSTM-2048)	43.9
	Ours large (2-layer LSTM-2048)	39.8

Perplexity improves (lower is better)

Source:https://research.fb.com/building-an-efficient-neural-language-model-over-a-billion-w ords/
Why should we care about Language Modeling?

- Language Modeling is a subcomponent of other NLP systems:
 - Speech recognition
 - Use a LM to generate transcription, conditioned on audio
 - Machine Translation
 - Use a LM to generate translation, conditioned on original text
 - Summarization
 - Use a LM to generate summary, conditioned on original text
- Language Modeling is a benchmark task that helps us measure our progress on understanding language

These systems are called conditional Language Models

Recap

- Language Model: A system that predicts the next word
- <u>Recurrent Neural Network</u>: A family of neural networks that:
 - Take sequential input of any length
 - Apply the same weights on each step
 - Can optionally produce output on each step
- Recurrent Neural Network ≠ Language Model
- We've shown that RNNs are a great way to build a LM.
- But RNNs are useful for much more!

RNNs can be used for tagging

e.g. part-of-speech tagging, named entity recognition



RNNs can be used for sentence classification

e.g. sentiment classification



RNNs can be used for sentence classification

e.g. sentiment classification



RNNs can be used for sentence



RNNs can be used to generate text

e.g. speech recognition, machine translation, summarization



Remember: these are called "conditional language models". We'll see Machine Translation in much more detail later.

RNNs can be used as an encoder module

e.g. question answering, machine translation



Thank you