CS 1678: Intro to Deep Learning Recurrent Neural Networks

Prof. Adriana Kovashka University of Pittsburgh March 16, 2021

Plan for this lecture

- Recurrent neural networks
 - Basics
 - Training (backprop through time, vanishing gradient)
 - Recurrent networks with gates (GRU, LSTM)
- Applications in NLP and vision
 - Neural machine translation (beam search, attention)
 - Image/video captioning

Recurrent neural networks

Some pre-RNN captioning results



This is a picture of one sky, one road and one sheep. The gray sky is over the gray road. The gray sheep is by the gray road.



This is a picture of two dogs. The first dog is near the second furry dog.



Here we see one road, one sky and one bicycle. The road is near the blue sky, and near the colorful bicycle. The colorful bicycle is within the blue sky.

Results with Recurrent Neural Networks



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



vanilla neural networks

Andrej Karpathy



e.g. image captioning image -> sequence of words

Andrej Karpathy



e.g. **sentiment classification** sequence of words -> sentiment







many to many



e.g. **machine translation** seq of words -> seq of words



e.g. video classification on frame level



Andrej Karpathy





We can process a sequence of vectors \mathbf{x} by applying a recurrence formula at every time step:

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

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



(Vanilla) Recurrent Neural Network

The state consists of a single "hidden" vector h:



Character-level language model example

Vocabulary: [h,e,l,o]

Example training sequence: **"hello"**



Character-level language model example

Vocabulary: [h,e,l,o]

Example training sequence: **"hello"**



Character-level language model example

Example training sequence: **"hello"**

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



Character-level language model example

Vocabulary: [h,e,l,o]

Example training sequence: **"hello"**



What do we still need to specify, for this to work? What kind of loss can we formulate?

Andrej Karpathy

Training a Recurrent Neural Network

- Get a big corpus of text which is a sequence of words $x^{(1)}, \ldots, x^{(T)}$
- Feed into RNN; compute output distribution $\hat{y}^{(t)}$ for every step t.
 - i.e. predict probability distribution of *every word*, given words so far
- Loss function on step t is cross-entropy between predicted probability distribution $\hat{y}^{(t)}$, and true next word $y^{(t)}$ (one-hot); V is vocabulary

$$J^{(t)}(\theta) = CE(\boldsymbol{y}^{(t)}, \hat{\boldsymbol{y}}^{(t)}) = -\sum_{w \in V} \boldsymbol{y}_{w}^{(t)} \log \hat{\boldsymbol{y}}_{w}^{(t)} = -\log \hat{\boldsymbol{y}}_{\boldsymbol{x}_{t+1}}^{(t)}$$

• Average this to get overall loss for entire training set:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^{T} -\log \hat{y}_{\boldsymbol{x}_{t+1}}^{(t)}$$

The vanishing/exploding gradient problem

- The error at a time step ideally can tell a previous time step from many steps away to change during backprop
- Multiply the same matrix at each time step during backprop



The vanishing gradient problem

• Total error is the sum of each error at time steps t

$$\frac{\partial E}{\partial W} = \sum_{t=1}^{T} \frac{\partial E_t}{\partial W}$$

• Chain rule:

$$\frac{\partial E_t}{\partial W} = \sum_{k=1}^t \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W}$$

• More chain rule:

$$\frac{\partial h_t}{\partial h_k} = \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}}$$

 Derivative of vector wrt vector is a Jacobian matrix of partial derivatives; norm of this matrix can become very small or very large quickly [Bengio et al 1994, Pascanu et al. 2013], leading to vanishing/exploding gradient What uses of language models from everyday life can you think of?

Now in more detail...

Language Modeling

• Language Modeling is the task of predicting what word comes next.



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

$$P(\boldsymbol{x}^{(t+1)} | \boldsymbol{x}^{(t)}, \dots, \boldsymbol{x}^{(1)})$$

where $m{x}^{(t+1)}$ can be any word in the vocabulary $V = \{m{w}_1,...,m{w}_{|V|}\}$

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

n-gram Language Models

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

n-1 words

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

(assumption)

prob of a n-gram
$$\begin{array}{c} & & \\ & = & P(\boldsymbol{x}^{(t+1)}, \boldsymbol{x}^{(t)}, \dots, \boldsymbol{x}^{(t-n+2)}) \\ & & P(\boldsymbol{x}^{(t)}, \dots, \boldsymbol{x}^{(t-n+2)}) \end{array}$$
 (definition of conditional prob)

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

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

Sparsity Problems with n-gram Language Models

Sparsity Problem 1

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

<u>(Partial)</u> Solution: Add small δ to the count for every $w \in V$. This is called *smoothing*.

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

count(students opened their w)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*!

(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.

A fixed-window neural Language Model

books

laptops output distribution $\hat{y} = \operatorname{softmax}(Uh + b_2) \in \mathbb{R}^{|V|}$ ิล ZO0 IJ hidden layer 000000000000 $\boldsymbol{h} = f(\boldsymbol{W}\boldsymbol{e} + \boldsymbol{b}_1)$ W concatenated word embeddings 0000 0000 0000 0000 $e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$ words / one-hot vectors the students their opened $x^{(1)}$ $m{x}^{(1)},m{x}^{(2)},m{x}^{(3)},m{x}^{(4)}$ $x^{(4)}$ $x^{(2)}$ $x^{(3)}$

A fixed-window neural Language Model

Improvements over *n*-gram LM:

- No sparsity problem
- Don't need to store all observed *n*-grams

Remaining **problems**:

- Fixed window is too small
- Enlarging window enlarges W
- Window can never be large enough!
- x⁽¹⁾ and x⁽²⁾ are multiplied by completely different weights in W. No symmetry in how the inputs are processed.

We need a neural architecture that can process *any length input*



Recurrent Neural Networks (RNN)

A family of neural architectures

Core idea: Apply the same weights *W repeatedly*



A RNN Language Model

output distribution

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

hidden states $h^{(t)} = \sigma \left(W_h h^{(t-1)} + W_e e^{(t)} + b_1 \right)$ $h^{(0)}$ is the initial hidden state

word embeddings

$$oldsymbol{e}^{(t)} = oldsymbol{E}oldsymbol{x}^{(t)}$$

words / one-hot vectors $oldsymbol{x}^{(t)} \in \mathbb{R}^{|V|}$



A RNN Language Model

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

 \bigcirc

 \bigcirc

 \bigcirc

RNN Advantages:

- Can process any length input
- Computation for step t can (in theory) use information from many steps back
- Model size doesn't increase for longer input
- Same weights applied on every timestep, so there is symmetry in how inputs are processed

RNN **Disadvantages**:

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





Recall: Training a RNN Language Model

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

$$J^{(t)}(\theta) = CE(\boldsymbol{y}^{(t)}, \hat{\boldsymbol{y}}^{(t)}) = -\sum_{w \in V} \boldsymbol{y}_{w}^{(t)} \log \hat{\boldsymbol{y}}_{w}^{(t)} = -\log \hat{\boldsymbol{y}}_{\boldsymbol{x}_{t+1}}^{(t)}$$

• Average this to get overall loss for entire training set:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^{T} -\log \hat{y}_{\boldsymbol{x}_{t+1}}^{(t)}$$

Training a RNN Language Model



...

Training a RNN Language Model



...

Training a RNN Language Model



...
Training a RNN Language Model



...

Training a RNN Language Model



Backpropagation for RNNs



Question: 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"

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:

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

Derivative of composition function



Source:

https://www.khanacademy.org/math/multivariable-calculus/multivariable-derivatives/differentiating-vector-valued-functions/a/multivariable-chain-rule-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:

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

Derivative of composition function



Source:

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

Backpropagation for RNNs



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

$$\uparrow$$
Question: How do we calculate this?

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













Vanishing gradient proof sketch

• Recall:
$$\boldsymbol{h}^{(t)} = \sigma \left(\boldsymbol{W}_h \boldsymbol{h}^{(t-1)} + \boldsymbol{W}_x \boldsymbol{x}^{(t)} + \boldsymbol{b}_1 \right)$$

- Therefore: $\frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} = \operatorname{diag}\left(\sigma'\left(\boldsymbol{W}_{h}\boldsymbol{h}^{(t-1)} + \boldsymbol{W}_{x}\boldsymbol{x}^{(t)} + \boldsymbol{b}_{1}\right)\right)\boldsymbol{W}_{h}$ (chain rule)
- Consider the gradient of the loss $J^{(i)}(\theta)$ on step *i*, with respect to the hidden state $h^{(j)}$ on some previous step *j*.

$$\begin{aligned} \frac{\partial J^{(i)}(\theta)}{\partial h^{(j)}} &= \frac{\partial J^{(i)}(\theta)}{\partial h^{(i)}} \prod_{j < t \leq i} \frac{\partial h^{(t)}}{\partial h^{(t-1)}} & \text{(chain rule)} \\ &= \frac{\partial J^{(i)}(\theta)}{\partial h^{(i)}} \overline{W_h^{(i-j)}} \prod_{j < t \leq i} \text{diag} \left(\sigma' \left(W_h h^{(t-1)} + W_x x^{(t)} + b_1 \right) \right) & \text{(value of } \frac{\partial h^{(t)}}{\partial h^{(t-1)}} \right) \\ & \text{If } W_h \text{ is small, then this term gets} \\ & \text{vanishingly small as } i \text{ and } j \text{ get further apart} \end{aligned}$$

Source: "On the difficulty of training recurrent neural networks", Pascanu et al, 2013. <u>http://proceedings.mlr.press/v28/pascanu13.pdf</u>

Vanishing gradient proof sketch

• Consider matrix L2 norms:

$$\frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(j)}} \left\| \leq \left\| \frac{\partial J^{(i)}(\theta)}{\partial \boldsymbol{h}^{(i)}} \right\| \left\| \boldsymbol{W}_h \right\|^{(i-j)} \prod_{j < t \leq i} \left\| \operatorname{diag} \left(\sigma' \left(\boldsymbol{W}_h \boldsymbol{h}^{(t-1)} + \boldsymbol{W}_x \boldsymbol{x}^{(t)} + \boldsymbol{b}_1 \right) \right) \right\|$$

- Pascanu et al showed that that if the largest eigenvalue of W_h is less than 1, then the gradient $\left\|\frac{\partial J^{(i)}(\theta)}{\partial h^{(j)}}\right\|$ will shrink exponentially
- There's a similar proof relating a largest eigenvalue >1 to exploding gradients

<u>Source</u>: "On the difficulty of training recurrent neural networks", Pascanu et al, 2013. <u>http://proceedings.mlr.press/v28/pascanu13.pdf</u>

Why is vanishing gradient a problem?



Effect of vanishing gradient on RNN-LM

- LM task: 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 from 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 gradient is small, the model can't learn this dependency
 - So the model is unable to predict similar long-distance dependencies at test time

Effect of vanishing gradient on RNN-LM

- LM task: The writer of the books _____
- **Correct answer**: The writer of the books <u>is</u> planning a sequel
- Syntactic recency: The <u>writer</u> of the books is
- Sequential recency: The writer of the books are

(incorrect)

(correct)

 Due to vanishing gradient, RNN-LMs are better at learning from sequential recency than syntactic recency, so they make this type of error more often than we'd like [Linzen et al 2016]

İS

are

Why is <u>exploding</u> gradient a problem?

• If the gradient becomes too big, then the SGD update step becomes too big:

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

loorning rate

- This can cause bad updates: we take too large a step and reach a bad parameter configuration (with large loss)
- In the worst case, this will result in Inf or NaN in your network (then you have to restart training from an earlier checkpoint)

Gradient clipping: solution for exploding gradient

• <u>Gradient clipping</u>: if the norm of the gradient is greater than some threshold, scale it down before applying SGD update

 Algorithm 1 Pseudo-code for norm clipping

 $\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$

 if $\|\hat{\mathbf{g}}\| \ge threshold$ then

 $\hat{\mathbf{g}} \leftarrow \frac{threshold}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$

 end if

• <u>Intuition</u>: take a step in the same direction, but a smaller step

RNNs with Gates

How to fix vanishing gradient problem?

- The main problem is that *it's too difficult for the RNN to learn to preserve information over many timesteps.*
- In a vanilla RNN, the hidden state is constantly being rewritten

$$oldsymbol{h}^{(t)} = \sigma \left(oldsymbol{W}_h oldsymbol{h}^{(t-1)} + oldsymbol{W}_x oldsymbol{x}^{(t)}
ight)$$

• How about a RNN with separate memory?

- More complex hidden unit computation in recurrence!
- Introduced by Cho et al. 2014
- Main ideas:
 - keep around memories to capture long distance dependencies
 - allow error messages to flow at different strengths depending on the inputs

- Standard RNN computes hidden layer at next time step directly: $h_t = f\left(W^{(hh)}h_{t-1} + W^{(hx)}x_t\right)$
- GRU first computes an update gate (another layer) based on current input word vector and hidden state

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

Compute reset gate similarly but with different weights

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

- Update gate $z_t = \sigma \left(W^{(z)} x_t + U^{(z)} h_{t-1} \right)$
- Reset gate $r_t = \sigma \left(W^{(r)} x_t + U^{(r)} h_{t-1} \right)$
- New memory content: $\tilde{h}_t = \tanh(Wx_t + r_t \circ Uh_{t-1})$ If reset gate unit is ~0, then this ignores previous memory and only stores the new word information
- Final memory at time step combines current and previous time steps: $h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$



$$z_t = \sigma \left(W^{(z)} x_t + U^{(z)} h_{t-1} \right)$$
$$r_t = \sigma \left(W^{(r)} x_t + U^{(r)} h_{t-1} \right)$$
$$\tilde{h}_t = \tanh \left(W x_t + r_t \circ U h_{t-1} \right)$$
$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

Richard Socher

 If reset *r* is close to 0, ignore previous hidden state: Allows model to drop information that is irrelevant in the future

$$z_t = \sigma \left(W^{(z)} x_t + U^{(z)} h_{t-1} \right)$$
$$r_t = \sigma \left(W^{(r)} x_t + U^{(r)} h_{t-1} \right)$$
$$\tilde{h}_t = \tanh \left(W x_t + r_t \circ U h_{t-1} \right)$$
$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

- If update z is close to 1, can copy information through many time steps, i.e. copy-paste state: Less vanishing gradient!
- Units with short-term dependencies often have reset gates (r) very active; ones with long-term dependencies have active update gates (z)

Long-short-term-memories (LSTMs)

- Proposed by Hochreiter and Schmidhuber in 1997
- We can make the units even more complex
- Allow each time step to modify
 - Input gate (current cell matters)
 - Forget (gate 0, forget past)
 - Output (how much cell is exposed) $o_t = \sigma \left(W^{(o)} x_t + U^{(o)} h_{t-1} \right)$
 - New memory cell
- Final memory cell:
- Final hidden state:

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

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

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

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

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

Long-short-term-memories (LSTMs)



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

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

Intuition: memory cells can keep information intact, unless inputs makes them forget it or overwrite it with new input

Cell can decide to output this information or just store it

Richard Socher, figure from wildml.com

Review on your own: Gated Recurrent Units (GRU)

- Proposed by Cho et al. in 2014 as a simpler alternative to the LSTM.
- On each timestep t we have input $x^{(t)}$ and hidden state $h^{(t)}$ (no cell state).

Update gate: controls what parts of hidden state are updated vs preserved $\boldsymbol{\succ} \boldsymbol{u}^{(t)} = \sigma \left(\boldsymbol{W}_{u} \boldsymbol{h}^{(t-1)} + \boldsymbol{U}_{u} \boldsymbol{x}^{(t)} + \boldsymbol{b}_{u}
ight)$ **Reset gate:** controls what parts of $\rightarrow \boldsymbol{r}^{(t)} = \sigma \left(\boldsymbol{W}_r \boldsymbol{h}^{(t-1)} + \boldsymbol{U}_r \boldsymbol{x}^{(t)} + \boldsymbol{b}_r \right)$ previous hidden state are used to compute new content New hidden state content: reset gate $\boldsymbol{J} = anh\left(oldsymbol{W}_h(oldsymbol{r}^{(t)} \circ oldsymbol{h}^{(t-1)}) + oldsymbol{U}_holdsymbol{x}^{(t)} + oldsymbol{b}_h
ight)$ selects useful parts of prev hidden state. Use this and current input to $h^{(t)} = (1 - u^{(t)}) \circ h^{(t-1)} + u^{(t)} \circ \tilde{h}^{(t)}$ compute new hidden content. Hidden state: update gate How does this solve vanishing gradient? simultaneously controls what is kept GRU makes it easier to retain info long-term

from previous hidden state, and what is updated to new hidden state content

"Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation", Cho et al. 2014, https://arxiv.org/pdf/1406.1078v3.pdf

(e.g. by setting update gate to 0)

Review on your own: Long Short-Term Memory (LSTM)

We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states $c^{(t)}$. On timestep t:



All these are vectors of same length *n*

Review on your own: Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:



LSTM vs GRU

- Researchers have proposed many gated RNN variants, but LSTM and GRU are the most widely-used
- The biggest difference is that GRU is quicker to compute and has fewer parameters
- There is no conclusive evidence that one consistently performs better than the other
- LSTM is a good default choice (especially if your data has particularly long dependencies, or you have lots of training data)
- <u>Rule of thumb</u>: start with LSTM, but switch to GRU if you want something more efficient

LSTMs: real-world success

- In 2013-2015, LSTMs started achieving state-of-the-art results for sequence modeling
 - Successful tasks include: handwriting recognition, speech recognition, machine translation, parsing, image captioning
 - LSTM became the dominant approach
- Starting in 2019, other approaches (e.g. Transformers) became more dominant for certain NLP tasks (will discuss next lecture)
 - For example in **WMT** (machine translation competition):
 - In WMT 2016, the summary report contains "RNN" 44 times
 - In WMT 2018, the report contains "RNN" 9 times and "Transformer" 63 times

Source: "Findings of the 2016 Conference on Machine Translation (WMT16)", Bojar et al. 2016, http://www.statmt.org/wmt16/pdf/W16-2301.pdf Source: "Findings of the 2018 Conference on Machine Translation (WMT18)", Bojar et al. 2018, http://www.statmt.org/wmt16/pdf/W16-2301.pdf

Is vanishing/exploding gradient just a RNN problem?

- No! It can be a problem for all neural architectures (including feed-forward and convolutional), especially deep ones.
 - Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
 - Thus lower layers are learnt very slowly (hard to train)
 - Solution: lots of new deep feedforward/convolutional architectures that add more direct connections (thus allowing the gradient to flow)

For example:

- Residual connections aka "ResNet"
- Also known as skip-connections
- The identity connection preserves information by default
- This makes deep networks much easier to train



Is vanishing/exploding gradient just a RNN problem?

- No! It can be a problem for all neural architectures (including feed-forward and convolutional), especially deep ones.
 - Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
 - Thus lower layers are learnt very slowly (hard to train)
 - Solution: lots of new deep feedforward/convolutional architectures that add more direct connections (thus allowing the gradient to flow)

For example:

- Dense connections aka "DenseNet"
- Directly connect everything to everything!



Bidirectional RNNs: motivation

Task: Sentiment Classification


Bidirectional RNNs

This contextual representation of "terribly" has both left and right context!



Abigail See

Bidirectional RNNs



Multi-layer RNNs

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

 \bigcirc

 \bigcirc

 \bigcirc

 \bigcirc

 \bigcirc

 \bigcirc

 \bigcirc

0

 \bigcirc

 \bigcirc

 \bigcirc \bigcirc 0 0 0 \bigcirc **RNN** layer 3 0 \bigcirc \bigcirc 0 \bigcirc \bigcirc 0 \bigcirc \bigcirc 0 0 \bigcirc \bigcirc **RNN** layer 2 0 \bigcirc () \bigcirc () \bigcirc \bigcirc \bigcirc \bigcirc 0 0 \bigcirc \bigcirc \bigcirc **RNN** layer 1 0 \bigcirc \bigcirc ()0 \bigcirc \bigcirc the terribly exciting movie was

Abigail See

Evaluating Language Models

• The standard evaluation metric for Language Models is perplexity.

$$perplexity = \prod_{t=1}^{T} \left(\frac{1}{P_{\text{LM}}(\boldsymbol{x}^{(t+1)} \mid \boldsymbol{x}^{(t)}, \dots, \boldsymbol{x}^{(1)})} \right)^{1/T} \underbrace{\qquad \text{Normalized by}}_{\text{number of words}}$$

• This is equal to the exponential of the cross-entropy loss $J(\theta)$:

$$= \prod_{t=1}^{T} \left(\frac{1}{\hat{y}_{x_{t+1}}^{(t)}} \right)^{1/T} = \exp\left(\frac{1}{T} \sum_{t=1}^{T} -\log \hat{y}_{x_{t+1}}^{(t)} \right) = \exp(J(\theta))$$

Lower perplexity is better!

Recap thus far

- 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
- Vanishing gradient problem: what it is, why it happens, and why it's bad for RNNs
- LSTMs and GRUs: more complicated RNNs that use gates to control information flow; more resilient to vanishing gradients

Plan for this lecture

- Recurrent neural networks
 - Basics
 - Training (backprop through time, vanishing gradient)
 - Recurrent networks with gates (GRU, LSTM)
- Applications in NLP and vision
 - Neural machine translation (beam search, attention)
 - Image/video captioning

Applications

Why should we care about Language Modeling?

- Language Modeling is a benchmark task that helps us measure our progress on understanding language
- Language Modeling is a subcomponent of many NLP tasks, especially those involving generating text or estimating the probability of text:
 - Predictive typing
 - Speech recognition
 - Handwriting recognition
 - Spelling/grammar correction
 - Authorship identification
 - Machine translation
 - Summarization
 - Dialogue
 - etc.

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.

Source: https://medium.com/@samim/obama-rnn-machine-generated-political-speeches-c8abd18a2ea0

- 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-8a9431803da6

- 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 paint color names:

 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 |

This is an example of a character-level RNN-LM (predicts what character comes next)

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

Generating poetry with RNNs



Generating poetry with RNNs

| at first: | tyntd-iafhatawiaoihrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e plia tklrgd t o idoe ns,smtt – h ne etie h,hregtrs nigtike,aoaenns lng |
|-----------|--|
| | train more |
| | "Tmont thithey" fomesscerliund Keushey. Thom here sheulke, anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize." |
| | train more |
| | Aftair fall unsuch that the hall for Prince Velzonski's that me of her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort how, and Gogition is so overelical and ofter. |
| | 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. |

More info: <u>http://karpathy.github.io/2015/05/21/rnn-effectiveness/</u>

Generating poetry with RNNs

PANDARUS:

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

Second Senator:

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

DUKE VINCENTIO:

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

Second Lord:

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

Clown:

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

VIOLA:

I'll drink it.

VIOLA:

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

KING LEAR:

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

Generating textbooks with RNNs

For $\bigoplus_{n=1,...,m}$ where $\mathcal{L}_{m_{\bullet}} = 0$, hence we can find a closed subset \mathcal{H} in \mathcal{H} and any sets \mathcal{F} on X, U is a closed immersion of S, then $U \to T$ is a separated algebraic space.

Proof. Proof of (1). It also start we get

 $S = \operatorname{Spec}(R) = U \times_X U \times_X U$

and the comparicoly in the fibre product covering we have to prove the lemma generated by $\coprod Z \times_U U \to V$. Consider the maps M along the set of points Sch_{fppf} and $U \to U$ is the fibre category of S in U in Section, ?? and the fact that any U affine, see Morphisms, Lemma ??. Hence we obtain a scheme S and any open subset $W \subset U$ in Sh(G) such that $Spec(R') \to S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that f_i is of finite presentation over S. We claim that $\mathcal{O}_{X,x}$ is a scheme where $x, x', s'' \in S'$ such that $\mathcal{O}_{X,x'} \to \mathcal{O}'_{X',x'}$ is separated. By Algebra, Lemma ?? we can define a map of complexes $\operatorname{GL}_{S'}(x'/S'')$ and we win.

To prove study we see that $\mathcal{F}|_U$ is a covering of \mathcal{X}' , and \mathcal{T}_i is an object of $\mathcal{F}_{X/S}$ for i > 0 and \mathcal{F}_p exists and let \mathcal{F}_i be a presheaf of \mathcal{O}_X -modules on \mathcal{C} as a \mathcal{F} -module. In particular $\mathcal{F} = U/\mathcal{F}$ we have to show that

$$\widetilde{M}^{\bullet} = \mathcal{I}^{\bullet} \otimes_{\operatorname{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1} \mathcal{F})$$

is a unique morphism of algebraic stacks. Note that

Arrows = $(Sch/S)_{fppf}^{opp}, (Sch/S)_{fppf}$

and

 $V = \Gamma(S, \mathcal{O}) \longmapsto (U, \operatorname{Spec}(A))$

is an open subset of X. Thus U is affine. This is a continuous map of X is the inverse, the groupoid scheme S.

Proof. See discussion of sheaves of sets.

The result for prove any open covering follows from the less of Example ??. It may replace S by $X_{spaces, \acute{e}tale}$ which gives an open subspace of X and T equal to S_{Zar} , see Descent, Lemma ??. Namely, by Lemma ?? we see that R is geometrically regular over S.

Lemma 0.1. Assume (3) and (3) by the construction in the description.

Suppose $X = \lim |X|$ (by the formal open covering X and a single map $\underline{Proj}_X(\mathcal{A}) = \operatorname{Spec}(B)$ over U compatible with the complex

 $Set(\mathcal{A}) = \Gamma(X, \mathcal{O}_{X, \mathcal{O}_X}).$

When in this case of to show that $\mathcal{Q} \to C_{Z/X}$ is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition ?? (without element is when the closed subschemes are catenary. If T is surjective we may assume that T is connected with residue fields of S. Moreover there exists a closed subspace $Z \subset X$ of X where U in X' is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem

(1) f is locally of finite type. Since S = Spec(R) and Y = Spec(R).

Proof. This is form all sheaves of sheaves on X. But given a scheme U and a surjective étale morphism $U \to X$. Let $U \cap U = \coprod_{i=1,...,n} U_i$ be the scheme X over S at the schemes $X_i \to X$ and $U = \lim_i X_i$.

The following lemma surjective restrocomposes of this implies that $\mathcal{F}_{x_0} = \mathcal{F}_{x_0} = \mathcal{F}_{\mathcal{X},\dots,0}$.

Lemma 0.2. Let X be a locally Noetherian scheme over $S, E = \mathcal{F}_{X/S}$. Set $\mathcal{I} = \mathcal{J}_1 \subset \mathcal{I}'_n$. Since $\mathcal{I}^n \subset \mathcal{I}^n$ are nonzero over $i_0 \leq \mathfrak{p}$ is a subset of $\mathcal{J}_{n,0} \circ \overline{A}_2$ works.

Lemma 0.3. In Situation ??. Hence we may assume q' = 0.

Proof. We will use the property we see that \mathfrak{p} is the mext functor (??). On the other hand, by Lemma ?? we see that

$$D(\mathcal{O}_{X'}) = \mathcal{O}_X(D)$$

where K is an F-algebra where δ_{n+1} is a scheme over S.

Andrej Karpathy

Generating code with RNNs

```
static void do_command(struct seq file *m, void *v)
{
  int column = 32 << (cmd[2] & 0x80);</pre>
  if (state)
    cmd = (int)(int state ^ (in 8(&ch->ch flags) & Cmd) ? 2 : 1);
  else
    seq = 1;
  for (i = 0; i < 16; i++) {
    if (k & (1 << 1))
     pipe = (in use & UMXTHREAD UNCCA) +
        ((count & 0x0000000fffffff8) & 0x000000f) << 8;
    if (count == 0)
      sub(pid, ppc_md.kexec_handle, 0x2000000);
    pipe set bytes(i, 0);
  }
  /* Free our user pages pointer to place camera if all dash */
  subsystem info = &of changes[PAGE SIZE];
  rek controls(offset, idx, &soffset);
  /* Now we want to deliberately put it to device */
  control check polarity(&context, val, 0);
  for (i = 0; i < COUNTER; i++)</pre>
    seq puts(s, "policy ");
}
```

Generated C code

Neural Machine Translation

- Neural Machine Translation (NMT) is a way to do Machine Translation with a *single neural network*
- The neural network architecture is called sequence-to-sequence (aka seq2seq) and it involves *two* RNNs.

Neural Machine Translation (NMT)



Greedy decoding

• We saw how to generate (or "decode") the target sentence by taking argmax on each step of the decoder



- This is greedy decoding (take most probable word on each step)
- Problems with this method?

Problems with greedy decoding

- Greedy decoding has no way to undo decisions!
 - Input: *il a m'entarté* (he hit me with a pie)
 - \rightarrow he _____
 - \rightarrow he hit ____
 - \rightarrow he hit a ____

(whoops! no going back now...)

• How to fix this?

Exhaustive search decoding

• Ideally we want to find a (length 7) translation y that maximizes

$$P(y|x) = P(y_1|x) P(y_2|y_1, x) P(y_3|y_1, y_2, x) \dots, P(y_T|y_1, \dots, y_{T-1}, x)$$
$$= \prod_{t=1}^T P(y_t|y_1, \dots, y_{t-1}, x)$$

- We could try computing all possible sequences *y*
 - This means that on each step t of the decoder, we're tracking V^T possible partial translations, where V is vocabulary size
 - This O(V^T) complexity is far too expensive!

Beam search decoding

- <u>Core idea:</u> On each step of decoder, keep track of the k most probable partial translations (*hypotheses*)
 - k is the beam size (in practice around 5 to 10)
- A hypothesis y_1, \ldots, y_t has a score which is its log probability: $\operatorname{score}(y_1, \ldots, y_t) = \log P_{\mathrm{LM}}(y_1, \ldots, y_t | x) = \sum_{i=1}^t \log P_{\mathrm{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$
 - Scores are all negative, and higher score is better
 - We search for high-scoring hypotheses, tracking top *k* on each step
- Beam search is not guaranteed to find optimal solution
- But much more efficient than exhaustive search!

Beam size = k = 2. Blue numbers = $score(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{LM}(y_i|y_1, \dots, y_{i-1}, x)$



Calculate prob dist of next word

Abigail See

Beam size = k = 2. Blue numbers = $score(y_1, \dots, y_t) = \sum_{i=1}^{r} \log P_{LM}(y_i|y_1, \dots, y_{i-1}, x)$



Beam size = k = 2. Blue numbers = $score(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{LM}(y_i|y_1, \dots, y_{i-1}, x)$



Beam size = k = 2. Blue numbers = $score(y_1, \dots, y_t) = \sum_{i=1}^{r} \log P_{LM}(y_i|y_1, \dots, y_{i-1}, x)$



Abigail See

Beam size = k = 2. Blue numbers = $score(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{LM}(y_i|y_1, \dots, y_{i-1}, x)$



Beam size = k = 2. Blue numbers = $score(y_1, \dots, y_t) = \sum_{i=1}^{t} \log P_{LM}(y_i|y_1, \dots, y_{i-1}, x)$



Beam size = k = 2. Blue numbers = $score(y_1, \dots, y_t) = \sum_{i=1}^{r} \log P_{LM}(y_i|y_1, \dots, y_{i-1}, x)$



For each of the *k* hypotheses, find top *k* next words and calculate scores

Beam size = k = 2. Blue numbers = $score(y_1, \dots, y_t) = \sum_{i=1}^{r} \log P_{LM}(y_i|y_1, \dots, y_{i-1}, x)$



Of these *k*² hypotheses, just keep *k* with highest scores

Beam size = k = 2. Blue numbers = $score(y_1, \dots, y_t) = \sum_{i=1}^{r} \log P_{LM}(y_i|y_1, \dots, y_{i-1}, x)$



For each of the *k* hypotheses, find top *k* next words and calculate scores

Beam size = k = 2. Blue numbers = $score(y_1, \dots, y_t) = \sum_{i=1}^{r} \log P_{LM}(y_i|y_1, \dots, y_{i-1}, x)$



Of these k^2 hypotheses, just keep k with highest scores

Beam size = k = 2. Blue numbers = $score(y_1, \dots, y_t) = \sum_{i=1}^{r} \log P_{LM}(y_i|y_1, \dots, y_{i-1}, x)$



For each of the *k* hypotheses, find top *k* next words and calculate scores

Beam size = k = 2. Blue numbers = $score(y_1, \dots, y_t) = \sum_{i=1}^{r} \log P_{LM}(y_i|y_1, \dots, y_{i-1}, x)$



This is the top-scoring hypothesis!

Beam size = k = 2. Blue numbers = $score(y_1, \dots, y_t) = \sum_{i=1}^{r} \log P_{LM}(y_i|y_1, \dots, y_{i-1}, x)$



Backtrack to obtain the full hypothesis
Beam search decoding: finishing up

- We have our list of completed hypotheses.
- How to select top one with highest score?
- Each hypothesis y_1, \ldots, y_t on our list has a score $\operatorname{score}(y_1, \ldots, y_t) = \log P_{\operatorname{LM}}(y_1, \ldots, y_t | x) = \sum_{i=1}^t \log P_{\operatorname{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$
- <u>Problem with this:</u> longer hypotheses have lower scores
- <u>Fix:</u> Normalize by length. Use this to select top one instead:

$$\frac{1}{t}\sum_{i=1}^t \log P_{\mathrm{LM}}(y_i|y_1,\ldots,y_{i-1},x)$$

How do we evaluate Machine Translation?

BLEU (Bilingual Evaluation Understudy)

- BLEU compares the <u>machine-written translation</u> to one or several <u>human-written translation(s)</u>, and computes a <u>similarity score</u> based on:
 - *n*-gram precision (usually for 1, 2, 3 and 4-grams)
 - Plus a penalty for too-short system translations
- BLEU is useful but imperfect
 - There are many valid ways to translate a sentence
 - So a good translation can get a poor BLEU score because it has low *n*-gram overlap with the human translation ⁽²⁾

MT progress over time

[Edinburgh En-De WMT newstest2013 Cased BLEU; NMT 2015 from U. Montréal]



Source: http://www.meta-net.eu/events/meta-forum-2016/slides/09_sennrich.pdf

NMT: the biggest success story of NLP Deep Learning

Neural Machine Translation went from a fringe research activity in **2014** to the leading standard method in **2016**

- 2014: First seq2seq paper published
- 2016: Google Translate switches from SMT to NMT
- This is amazing!
 - SMT systems, built by hundreds of engineers over many years, outperformed by NMT systems trained by a handful of engineers in a few months

So is Machine Translation solved?

- Nope!
- Many difficulties remain:
 - Out-of-vocabulary words
 - Domain mismatch between train and test data
 - Maintaining context over longer text
 - Low-resource language pairs

So is Machine Translation solved?

- Nope!
- Using common sense is still hard •



Open in Google Translate





Abigail See

So is Machine Translation solved?

- Nope!
- NMT picks up biases in training data



Didn't specify gender

Source: https://hackernoon.com/bias-sexist-or-this-is-the-way-it-should-be-ce1f7c8c683c

Abigail See

NMT research continues

NMT is the **flagship task** for NLP Deep Learning

- NMT research has pioneered many of the recent innovations of NLP Deep Learning
- In **2019**: NMT research continues to thrive
 - Researchers have found *many, many* improvements to the "vanilla" seq2seq NMT system
 - But one improvement is so integral that it is the new vanilla...

ATTENTION

Sequence-to-sequence: the bottleneck problem



Attention

- Attention provides a solution to the bottleneck problem.
- <u>Core idea</u>: on each step of the decoder, use *direct* connection to the encoder to focus on a particular part of the source sequence



 First we will show via diagram (no equations), then we will show with equations



Decoder RNN





Abigail See







Decoder RNN





Abigail See



Use the attention distribution to take a **weighted sum** of the encoder hidden states.

The attention output mostly contains information from the hidden states that received high attention.

Decoder RNN



Abigail See



Abigail See

Sometimes we take the attention output from the previous step, and also feed it into the decoder (along with the usual decoder input).



Attention: in equations

- We have encoder hidden states $h_1, \ldots, h_N \in \mathbb{R}^h$
- On timestep *t*, we have decoder hidden state $s_t \in \mathbb{R}^h$
- We get the attention score e^t for this step:

$$\boldsymbol{e}^t = [\boldsymbol{s}_t^T \boldsymbol{h}_1, \dots, \boldsymbol{s}_t^T \boldsymbol{h}_N] \in \mathbb{R}^N$$

 We take softmax to get the attention distribution α^t for this step (this is a probability distribution and sums to 1)

$$\alpha^t = \operatorname{softmax}(\boldsymbol{e}^t) \in \mathbb{R}^N$$

• We use α^t to take a weighted sum of the encoder hidden states to get the attention output a_t

$$\boldsymbol{a}_t = \sum_{i=1}^{n} \alpha_i^t \boldsymbol{h}_i \in \mathbb{R}^h$$

• Finally we concatenate the attention output a_t with the decoder hidden state s_t and proceed as in the non-attention seq2seq model

$$[oldsymbol{a}_t;oldsymbol{s}_t]\in\mathbb{R}^{2h}$$

Abigail See

Attention is great

- Attention significantly improves NMT performance
 - It's very useful to allow decoder to focus on certain parts of the source
- Attention solves the bottleneck problem
 - Attention allows decoder to look directly at source; bypass bottleneck
- Attention helps with vanishing gradient problem
 - Provides shortcut to faraway states
- Attention provides some interpretability
 - By inspecting attention distribution, we can see what the decoder was focusing on
 - We get (soft) alignment for free!
 - This is cool because we never explicitly trained
 an alignment system
 - The network just learned alignment by itself



Attention is a general Deep Learning technique

- We've seen that attention is a great way to improve the sequence-to-sequence model for Machine Translation.
- <u>However</u>: You can use attention in many architectures (not just seq2seq) and many tasks (not just MT)
- More general definition of attention:
 - Given a set of vector values, and a vector query, attention is a technique to compute a weighted sum of the values, dependent on the query.
- We sometimes say that the query attends to the values.
- For example, in seq2seq + attention model, each decoder hidden state (query) attends to all encoder hidden states (values).

Abigail See



CVPR 2015:

Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei Show and Tell: A Neural Image Caption Generator, Vinyals et al.

Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al. Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

Adapted from Andrej Karpathy

Recurrent Neural Network



Convolutional Neural Network

Andrej Karpathy



<START>

Andrej Karpathy





Andrej Karpathy





Andrej Karpathy



Andrej Karpathy





"man in black shirt is playing guitar."



"a young boy is holding a baseball bat."



"construction worker in orange safety vest is working on road."



"a cat is sitting on a couch with a remote control."



"two young girls are playing with lego toy."





"a horse is standing in the middle of a road."



"a woman holding a teddy bear in front of a mirror."



Key Insight:

Generate feature representation of the video and "decode" it to a sentence

Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015



Venugopalan et al., "Translating Videos to Natural Language using Deep Recurrent Neural Networks", NAACL-HTL 2015



FGM: A person is dancing with the person on the stage. YT: A group of men are riding the forest.

- I+V: A group of people are dancing.
- GT: Many men and women are dancing in the street.



FGM: A person is walking with a person in the forest. YT: A monkey is walking.

I+V: A bear is eating a tree.

GT: Two bear cubs are digging into dirt and plant matter at the base of a tree.



FGM: A person is cutting a potato in the kitchen.YT: A man is slicing a tomato.I+V: A man is slicing a carrot.GT: A man is slicing carrots.



FGM: A person is riding a horse on the stage. YT: A group of playing are playing in the ball. I+V: A basketball player is playing.

GT: Dwayne wade does a fancy layup in an allstar game.



Venugopalan et al., "Sequence to Sequence - Video to Text", ICCV 2015



Venugopalan et al., "Sequence to Sequence - Video to Text", ICCV 2015