

Lecture 10 Recap

• Digit recognition: 10 classes

60k parameters



- Conv -> Pool -> Conv -> Pool -> Conv -> FC
- As we go deeper: Width, Height ↓ Number of Filters ↓

AlexNet



• Softmax for 1000 classes

[Krizhevsky et al., ANIPS'12] AlexNet

I2DL: Prof. Dai

VGGNet

- Striving for **simplicity**
 - Conv -> Pool -> Conv -> Pool -> Conv -> FC
 - Conv=3x3, s=1, same; Maxpool=2x2, s=2
- As we go deeper: Width, Height + Number of Filters +
- Called VGG-16: 16 layers that have weights

138M parametersLarge but simplicity makes it appealing

Residual Block



I2DL: Prof. Dai

Inception Layer



(a) Inception module, naïve version

(b) Inception module with dimensionality reduction

[Szegedy et al., CVPR'15] GoogleNet



Lecture 11



Transfer Learning

Transfer Learning

• Training your own model can be difficult with limited data and other resources

e.g.,

- It is a laborious task to manually annotate your own training dataset
- Why not reuse already pre-trained models?



Transfer Learning for Images



[Zeiler al., ECCV'14] Visualizing and Understanding Convolutional Networks

Trained on	Transfer Learning
ImageNet	$\mathbf{\circ}$
FC-1000 FC-4096 FC-4096	
Conv-512 Conv-512	
Conv-512 Conv-512	Footuro
MaxPool Conv-256 Conv-256	extraction
MaxPool Conv-128 Conv-128	
MaxPool Conv-64 Conv-64	
Image	$[D_{D_{T}}}}}}}}}$





Transfer Learning

FRC

If the dataset is big enough train more layers with a low learning rate



When Transfer Learning Makes Sense

- When task T1 and T2 have the same input (e.g. an RGB image)
- When you have more data for task T1 than for task T2

• When the low-level features for T1 could be useful to learn T2

Now you are:

• Ready to perform image classification on any dataset

• Ready to design your own architecture

Ready to deal with other problems such as semantic segmentation (Fully Convolutional Network)



Representation Learning

Learning Good Features

- Good features are essential for successful machine
 learning
- (Supervised) deep learning depends on training data used: input/target labels
- Change in inputs (noise, irregularities, etc) can result in drastically different results

Representation Learning

Allows for discovery of representations required for various tasks

Deep representation learning: model maps input X to output Y

Deep Representation Learning

• Intuitively, deep networks learn multiple levels of abstraction



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

How to Learn Good Features?

• Determine desired feature invariances

• Teach machines to distinguish between similar and dissimilar things

Match the correct animal







https://amitness.com/2020/03/illustrated-simclr/ ²³

How to Learn Good Features?



[Chen et al., ICML'20] SimCLR, https://amitness.com/2020/03/illustrated-simclr/ 24



Transfer & Representation Learning

- Transfer learning can be done via representation learning
- Effectiveness of representation learning often demonstrated by transfer learning performance (but also other factors, e.g., smoothness of the manifold)



Recurrent Neural Networks

Processing Sequences

• Recurrent neural networks process sequence data

• Input/output can be sequences



Classical neural networks for image classification



Image captioning

I2DL: Prof. Dai



Language recognition



Machine translation



Event classification



Event classification

Basic Structure of an RNN

• Multi-layer RNN



Basic Structure of an RNN

• Multi-layer RNN

The hidden state will have its own internal dynamics More expressive model!



Basic Structure of an RNN

• We want to have notion of "time" or "sequence"


• We want to have notion of "time" or "sequence"



 $A_t = \theta_c A_{t-1} + \theta_x x_t$ Parameters to be learned

• We want to have notion of "time" or "sequence"



$$\boldsymbol{A}_t = \boldsymbol{\theta}_c \boldsymbol{A}_{t-1} + \boldsymbol{\theta}_x \boldsymbol{x}_t$$

 $\boldsymbol{h}_t = \boldsymbol{\theta}_{\boldsymbol{h}} \boldsymbol{A}_t$

Note: non-linearities ignored for now

• We want to have notion of "time" or "sequence"



$$A_t = \theta_c A_{t-1} + \theta_x x_t$$

$$h_t = \theta_h A_t$$

Same parameters for each time step = generalization!

• Unrolling RNNs

Same function for the hidden layers





• Unrolling RNNs



• Unrolling RNNs as feedforward nets



Backprop through an RNN

• Unrolling RNNs as feedforward nets



Add the derivatives at different times for each weight



[Olah, https://colah.github.io '15] Understanding LSTMs 45

• Simple recurrence $A_t = \theta_c A_{t-1} + \theta_x x_t$

• Let us forget the input $A_t = \theta_c^t A_0$

Same weights are multiplied over and over again t

Xt.

• Simple recurrence $A_t = \theta_c^{\ t} A_0$

What happens to small weights? Vanishing gradient

What happens to large weights? Exploding gradient

 $\boldsymbol{\theta} = \boldsymbol{Q} \boldsymbol{\Lambda} \boldsymbol{Q}^T$

• Simple recurrence $A_t = \theta_c^{\ t} A_0$

• If $\boldsymbol{\theta}$ admits eigendecomposition

Diagonal of this matrix are the eigenvalues t

Xt

• Simple recurrence $A_t = \theta^t A_0$

• If $\boldsymbol{\theta}$ admits eigendecomposition

 $\boldsymbol{\theta} = \boldsymbol{Q} \boldsymbol{\Lambda} \boldsymbol{Q}^T$



- Orthogonal $\boldsymbol{\theta}$ allows us to simplify the recurrence

$$\boldsymbol{A}_t = \boldsymbol{Q}\boldsymbol{\Lambda}^t \boldsymbol{Q}^T \boldsymbol{A}_0$$

• Simple recurrence $A_t = Q \Lambda^t Q^T A_0$

What happens to eigenvalues with magnitude less than one? Vanishing gradient

What happens to eigenvalues with magnitude larger than one?

Exploding gradient 🔍

Gradient clipping

• Simple recurrence $A_t = \theta_c^{\ t} A_0$

Let us just make a matrix with eigenvalues = 1

Allow the **cell** to maintain its "*state*"

Vanishing Gradient

• 1. From the weights $A_t = \theta_c^t A_0$

• 2. From the activation functions (*tanh*)



tanh.x 1.0 H

0.5

-1.0

Vanishing Gradient

- 1. From the weights $A_t = \mathbf{A}^t A_0$
- 2. From the activation functions (*tanh*)





Long Short Term Memory

[Hochreiter et al., Neural Computation'97] Long Short-Term Memory

Long-Short Term Memory Units

• Simple RNN has **tanh** as non-linearity





Long-Short Term Memory Units

- Key ingredients
- Cell = transports the information through the unit



Long-Short Term Memory Units

- Key ingredients
- Cell = transports the information through the unit
- Gate = remove or add information to the cell state



• Forget gate $f_t = sigm(\theta_{xf}x_t + \theta_{hf}h_{t-1} + b_f)$



Decides when to erase the cell state

Sigmoid = output between **0** (forget) and **1** (keep)

• Input gate $i_t = sigm(\theta_{xi}x_t + \theta_{hi}h_{t-1} + b_i)$



Decides which values will be updated

New cell state, output from a tanh (-1,1)

• Element-wise operations



• Output gate $h_t = o_t \odot \tanh(C_t)$



Decides which values will be outputted

Output from a tanh (-1, 1)

- Forget gate $f_t = sigm(\theta_{xf}x_t + \theta_{hf}h_{t-1} + b_f)$
- Input gate
- Output gate
- Cell update

- $i_{t} = sigm(\theta_{xi}x_{t} + \theta_{hi}h_{t-1} + b_{i})$ $o_{t} = sigm(\theta_{xo}x_{t} + \theta_{ho}h_{t-1} + b_{o})$ $g_{t} = tanh(\theta_{xg}x_{t} + \theta_{hg}h_{t-1} + b_{g})$
- Cell $C_t = f_t \odot C_{t-1} + i_t \odot g_t$
- Output $h_t = o_t \odot \tanh(C_t)$

- Forget gate
- Input gate
- Output gate
- Cell update
- Cell
- Output

 $\boldsymbol{f}_t = sign(\boldsymbol{\theta}_{x_f} \boldsymbol{x}_t + \boldsymbol{\theta}_{h_f} \boldsymbol{h}_{t-1} + \boldsymbol{b}_f)$ $\mathbf{i}_t = sigm(\boldsymbol{\theta}_{xi}\mathbf{x}_t + \boldsymbol{\theta}_{hi}\mathbf{h}_{t-1} + \mathbf{b}_i)$ $\boldsymbol{o}_t = sigm(\boldsymbol{\theta}_{xo}\boldsymbol{x}_t + \boldsymbol{\theta}_{ho}\boldsymbol{h}_{t-1} + \boldsymbol{b}_o)$ $\boldsymbol{g}_t = tan\boldsymbol{h}(\boldsymbol{\theta}_{xg}\boldsymbol{x}_t + \boldsymbol{\theta}_{hg}\boldsymbol{h}_{t-1} + \boldsymbol{b}_g)$ $\boldsymbol{C}_t = \boldsymbol{f}_t \odot \boldsymbol{C}_{t-1} + \boldsymbol{i}_t \odot \boldsymbol{g}_t$

 $\boldsymbol{h}_t = \boldsymbol{o}_t \odot \tanh(\boldsymbol{C}_t)$

Learned through backpropagation

LSTM

• Highway for the gradient to flow



LSTM: Dimensions

• Cell update $\begin{array}{ccc} 128 & 128 \\ \boldsymbol{g}_t = tanh(\boldsymbol{\theta}_{xg}\boldsymbol{x}_t + \boldsymbol{\theta}_{hg}\boldsymbol{h}_{t-1} + \boldsymbol{b}_g) \end{array}$



[Olah, https://colah.github.io '15] Understanding LSTMs 67

LSTM in code

def lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b): def lstm step backward(dnext h. dnext c. cache): Forward pass for a single timestep of an LSTM. Backward pass for a single timestep of an LSTM. The input data has dimension D, the hidden state has dimension H, and we use Inputs: a minibatch size of N. - dnext h: Gradients of next hidden state, of shape (N, H) - dnext c: Gradients of next cell state, of shape (N, H) Inputs: - cache: Values from the forward pass - x: Input data, of shape (N, D) - prev h: Previous hidden state, of shape (N, H) Returns a tuple of: - prev c: previous cell state, of shape (N, H) - dx: Gradient of input data, of shape (N, D) - Wx: Input-to-hidden weights, of shape (D, 4H) - dprev h: Gradient of previous hidden state, of shape (N, H) - Wh: Hidden-to-hidden weights, of shape (H, 4H) - dprev c: Gradient of previous cell state, of shape (N, H) - b: Biases, of shape (4H,) - dWx: Gradient of input-to-hidden weights, of shape (D, 4H) - dWh: Gradient of hidden-to-hidden weights, of shape (H. 4H) Returns a tuple of: - db: Gradient of biases, of shape (4H,) - next h: Next hidden state, of shape (N, H) next c: Next cell state, of shape (N, H) dx. dh. dc. dWx. dWh. db = None, None, None, None, None, None - cache: Tuple of values needed for backward pass. i, f, o, g, a, ai, af, ao, ag, Wx, Wh, b, prev h, prev c, x, next c, next h = cache next h, next c, cache = None, None, None # backprop into step 5 N, H = prev h.shape do = np.tanh(next c) * dnext h# 1 a = np.dot(x, Wx) + np.dot(prev h, Wh) + bdnext $c \neq 0 * (1 - np.tanh(next c) ** 2) * dnext h$ # backprop into 4 # 2 ai = a[:, :H] df = prev c * dnext c af = a[:, H:2*H]dprev c = f * dnext cao = a[:, 2*H:3*H]di = q * dnext c aq = a[:, 3*H:]dq = i * dnext c# 3 # backprop into 3 i = sigmoid(ai) dai = sigmoid(ai) * (1 - sigmoid(ai)) * di f = sigmoid(af) daf = sigmoid(af) * (1 - sigmoid(af)) * df o = sigmoid(ao)dao = sigmoid(ao) * (1 - sigmoid(ao)) * do g = np.tanh(ag) dag = (1 - np.tanh(ag) ** 2) * dg# 4 # backprop into 2 next c = f * prev c + i * q da = np.hstack((dai, daf, dao, dag)) # 5 # backprop into 1 next h = o * np.tanh(next c) db = np.sum(da, axis = 0)dprev h = np.dot(Wh, da.T).Tcache = i, f, o, g, a, ai, af, ao, ag, Wx, Wh, b, prev h, prev c, x, next c, next h dWh = np.dot(prev h.T, da)dx = np.dot(da, Wx.T)return next h, next c, cache dWx = np.dot(x.T, da)

return dx, dprev_h, dprev_c, dWx, dWh, db



Attention

Attention is all you need

Attention Is All You Need

Ashish Vaswani* Google Brain avaswani@google.com Noam Shazeer* Google Brain noam@google.com n:

Niki Parmar* Google Research nikip@google.com Jakob Uszkoreit* Google Research usz@google.com

Llion Jones* Google Research llion@google.com Aidan N. Gomez^{*†} University of Toronto aidan@cs.toronto.edu Łukasz Kaiser* Google Brain lukaszkaiser@google.com

Illia Polosukhin* [‡] illia.polosukhin@gmail.com

Attention is all you need

Attention Is All You Need

~62,000 citations in 5 years!

Ashish Vaswani* Google Brain avaswani@google.com Noam Shazeer* Google Brain noam@google.com

Niki Parmar* Google Research nikip@google.com

Jakob Uszkoreit* Google Research usz@google.com

Llion Jones* Google Research llion@google.com Aidan N. Gomez^{*†} University of Toronto aidan@cs.toronto.edu Łukasz Kaiser* Google Brain lukaszkaiser@google.com

Illia Polosukhin* [‡] illia.polosukhin@gmail.com

Attention vs convolution

Convolution

Global attention





Fully Connected layer







Source: https://colah.github.io/posts/2015-08-Understanding-LSTMs/



I moved to Germany ...

so I speak German fluently

I2DL: Prof. Dai
Attention: Architecture

• A decoder processes the information

- Decoders take as input:
 - Previous decoder hidden state
 - Previous output
 - Attention





Transformers

Deep Learning Revolution

	Deep Learning	Deep Learning 2.0
Main idea	Convolution	Attention
Field invented	Computer vision	NLP
Started	NeurIPS 2012	NeurIPS 2017
Paper	AlexNet	Transformers
Conquered vision	Around 2014-2015	Around 2020-2021
Replaced (Augmented)	Traditional ML/CV	CNNs, RNNs

Transformers





Intuition: Take the query Q, find the most similar key K, and then find the value V that corresponds to the key.

In other words, learn V, K, Q where: V – here is a bunch of interesting things. K – here is how we can index some things. Q – I would like to know this interesting thing.

Loosely connected to Neural Turing Machines (Graves et al.).



To train them well, divide by $\sqrt{d_k}$, "probably" because for large values of the key's dimension, the dot product grows large in magnitude, pushing the softmax function into regions where it has extremely small gradients.



Adapted from Y. Kilcher

I2DL: Prof. Dai





Values
\vee 1
V2
∨3
\vee 4
$\vee 5$



QK^T Essentially, dot product between (<Q,K1>), (<Q,K2>), (<Q,K3>), (<Q,K4>), (<Q,K5>).





softmax

I2DL: Prof. Dai

 $\begin{pmatrix} QK^T \\ \sqrt{d_k} \end{pmatrix}$ Is simply inducing a distribution over the values. The larger a value is, the higher is its softmax value. Can be interpreted as a differentiable soft indexing.



softmax $\left(\frac{\zeta}{-1}\right)$

I2DL: Prof. Dai

Is simply inducing a distribution over the values. The larger a value is, the higher is its softmax value. Can be interpreted as a differentiable soft indexing.





Selecting the value V where the network needs to attend..





Good old fullyconnected layers.





N layers of attention followed by FC





Same as multi-head attention, but masked. Ensures that the predictions for position i can depend only on the known outputs at positions less than i.





Multi-headed attention between encoder and the decoder.





Projection and prediction.





What is missing from self-attention?

- Convolution: a different linear transformation for each relative position. Allows you to distinguish what information came from where.
- Self-attention: a weighted average.



Uses fixed positional encoding based on trigonometric series, in order for the model to make use of the order of the sequence





Transformers – a final look



Self-attention: complexity

Layer Type	Complexity per Layer	Sequential	Maximum Path Length
		Operations	
Self-Attention	$O(n^2 \cdot d)$	O(1)	O(1)
Recurrent	$O(n \cdot d^2)$	O(n)	O(n)
Convolutional	$O(k\cdot n\cdot d^2)$	O(1)	$O(log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	O(1)	O(n/r)

where n is the sequence length, d is the representation dimension, k is the convolutional kernel size, and r is the size of the neighborhood.

Self-attention: complexity

Layer Type	Complexity per Layer	Sequential	Maximum Path Length
		Operations	
Self-Attention	$O(n^2 \cdot d)$	O(1)	O(1)
Recurrent	$O(n \cdot d^2)$	O(n)	O(n)
Convolutional	$O(k\cdot n\cdot d^2)$	O(1)	$O(log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	O(1)	O(n/r)

where n is the sequence length, d is the representation dimension, k is the convolutional kernel size, and r is the size of the neighborhood.

Considering that most sentences have a smaller dimension than the representation dimension (in the paper, it is 512), self-attention is very efficient.

I2DL: Prof. Dai

Transformers – training tricks

• ADAM optimizer with proportional learning rate:

 $lrate = d_{\text{model}}^{-0.5} \cdot \min(step_num^{-0.5}, step_num \cdot warmup_steps^{-1.5})$

- Residual dropout
- Label smoothing
- Checkpoint averaging

Transformers - results

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training C	Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR	
ByteNet [15]	23.75				
Deep-Att + PosUnk [32]		39.2		$1.0\cdot 10^{20}$	
GNMT + RL [31]	24.6	39.92	$2.3\cdot 10^{19}$	$1.4\cdot10^{20}$	
ConvS2S [8]	25.16	40.46	$9.6\cdot 10^{18}$	$1.5\cdot 10^{20}$	
MoE [26]	26.03	40.56	$2.0\cdot 10^{19}$	$1.2\cdot 10^{20}$	
Deep-Att + PosUnk Ensemble [32]		40.4		$8.0\cdot10^{20}$	
GNMT + RL Ensemble [31]	26.30	41.16	$1.8\cdot 10^{20}$	$1.1\cdot 10^{21}$	
ConvS2S Ensemble [8]	26.36	41.29	$7.7\cdot 10^{19}$	$1.2\cdot10^{21}$	
Transformer (base model)	27.3	38.1	3.3 •	$3.3\cdot 10^{18}$	
Transformer (big)	28.4	41.0	$2.3\cdot 10^{19}$		

Transformers - summary

- Significantly improved SOTA in machine translation
- Launched a new deep-learning revolution in MLP
- Building block of NLP models like BERT (Google) or GPT/ChatGPT (OpenAI)
- BERT has been heavily used in Google Search

• And eventually made its way to computer vision (and other related fields)



See you next time!