

Neural Networks

Recurrent Neural networks, LSDM

(P-ITEEA-0011)

Akos Zarandy Lecture 9 November 12, 2018



Administrative announces

- Next week 3rd of December at the lecture time
 - Replacement test paper (Pót ZH)
 - Email will be sent through Neptun
- No quizzes on Tuesday in the workclass.
 - It will be done later.
 - The Wednesday and Monday group will have quizzed.
- Consultation will be held in the usual consultation time and place
- Programming test is coming (18th of December)
 - Please practice!

Contents



- How to handle sequential signals with Neural Networks?
- Recurrent Networks
 - Training
 - Examples
 - Vanishing gradient problem
- Long Short Term Memory (LSTM)
 - LSTM versions

Static samples vs Data signal flow

AlexNet could recognize 1000s of images. ResNet could reach better then human performance.

- Though human can recognize
 - Single letters
 - Single sounds
 - Single tunes
 - Single pictures

- But in real life we handle
 - Texts
 - Speech
 - Music

Can feed forward neural networks (perceptrons,

conv. nets) solve these problems?

DATA MEMORY

Movies





Memory

- Our feed-forward nets had so far
 - Program memory (for the weights)
 - Registers
 - For store temporally due to implementation and not matematical resasons
 - Registers were not part of the networks
- After each inferences the net was reset
 - All registers were deleted
 - No information remained in the net after processing an input vector
 - Therefore the order of a test sequence made no difference

Recurrent networks (RNN)

- Unlike traditional neural networks, the output of the RNN depends on the previous inputs
 - State
- RNN contains feedback
- Theoretically:
 - Directed graph with cyclic loops
- From now, time has a role in execution
 - Time steps, delays

Jürgen lives in Berlin.

He speeks

Feedback loop





-







11/26/2018

Input Laver Hidden Laver Output Layer



Weights (multiple arrows) replaced with Input Vector vectors (single arrows)

Output Vector



11/26/2018







Activation function in feedback loop

- Activation function of the hidden layers is typically hyperbolic tangent
- It avoids large positive feedback
 - Keeps the output between
 -1 and +1
 - Avoids exploding the loop calculation
 - Gain should be smaller than 1 in the loop!

Positive feedback in a loop: A produces more of **B** which in turn produces more of **A**. It leeds to increase beyond any limit.







Timing of the RNN

How to calculate back propagation?

- Discrete time steps are used
- Input vector sequence to apply
- Signals are calculated in a node, when all inputs exist
- State machine

Time	Input	State	output
t=1	<i>x</i> (1)	$h(1) = f\big(h(0), x(1)\big)$	y(1) = g(h(1))
t=2	<i>x</i> (2)	$h(2) = f\left(\frac{h(1)}{x(2)}\right)$	y(2) = g(h(2))
t=3	<i>x</i> (3)	$h(3) = f\left(\frac{h(2)}{x(3)}\right)$	y(3) = g(h(3))
t=4	<i>x</i> (4)	$h(4) = f\left(\frac{h(3)}{x(4)}\right)$	$y(4) = g\big(h(4)\big)$
		• • •	





Unrolling



Unrolling



- Unrolling generates an acyclic directed graph from the original cyclic directed graph structure
- It generates a final impulse response (FIR) filter from the original infinite impulse response (IIR) filter
- Dynamic behavior

IIR filters may response to any finite length input with a infinite (usually decaying) response, due to their internal loop.



FIR filters response to any finite length input with a final response.

Weight matrix sharing

RNN re-uses the same weight matrix in every unrolled steps.





Example: Character-level Language Model

Vocabulary: One-hot [h,e,l,o] encoding

Example training sequence: "hello"



11/26/2018 http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture10.pdf





Example: Character-level Language Model

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







Example: Character-level Language Model Sampling

Vocabulary: [h,e,l,o]



Example: Character-level Language Model Sampling

Vocabulary: [h,e,l,o]



Example: Character-level Language Model Sampling

Vocabulary: [h,e,l,o]





Example: Character-level Language Model Sampling Backpropagation can be started using negative log Vocabulary: likelihood cost [h,e,l,o]

function



Back propagation through time

- Assuming that the length of the input vector sequence is limited
- It became a feedforward neural net
- Possible to apply back propagation
- We need multiple vector sequences to train!







Backpropagation through time

Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient



Truncated Backpropagation through time





Run forward and backward through chunks of the sequence instead of whole sequence

Truncated Backpropagation through time





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



Truncated Backpropagation through time





Image captioning example



Figure from Karpathy et a, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015; figure copyright IEEE, 2015. Reproduced for educational purposes.

Explain Images with Multimodal Recurrent Neural Networks, Mao et al. 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

Image captioning example Recurrent Neural Network





Convolutional Neural Network

test image



This image is CC0 public domain









Alexnet: scored 5 best guesses



žĸ



2 K 7 K









2 K 7 K



test image



**





2 K 2 K

Image captioning Example: Results





A cat sitting on a suitcase on the floor



A cat is sitting on a tree branch



Two people walking on the beach with surfboards

11/26/2018



A tennis player in action on the court



A dog is running in the grass with a frisbee



Two giraffes standing in a grassy field



A white teddy bear sitting in the grass



A man riding a dirt bike on a dirt track

41

Image captioning: Failure cases





A woman is holding a cat in her hand



A person holding a computer mouse on a desk



A woman standing on a beach holding a surfboard



A bird is perched on a tree branch



A man in a baseball uniform throwing a ball

Problem



• What happens if the input sequence is too long?

Vanishing gradient!



Vanishing Gradient Problem

- In case of long input vector sequencies, the old vectors has a strongly fading effect in inference phase
- In training phase, the stacked gradient functions will be very small



Practical problem of long term dependences

- Consider a network which predicts the next word in a text
 - If the information needed to predict is close, it can be successfully trained
 - If required information is far, the training will be difficult





Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994 Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



$$h_{t} = \tanh(W_{hh}h_{t-1} + W_{xh}x_{t})$$
$$= \tanh\left(\left(W_{hh} \quad W_{hx}\right) \begin{pmatrix}h_{t-1}\\x_{t}\end{pmatrix}\right)$$
$$= \tanh\left(W\begin{pmatrix}h_{t-1}\\x_{t}\end{pmatrix}\right)$$



Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994 Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

 $h_{t} = \tanh(W_{hh}h_{t-1} + W_{xh}x_{t})$ $= \tanh\left(\left(W_{hh} \quad W_{hx}\right) \begin{pmatrix}h_{t-1}\\x_{t}\end{pmatrix}\right)$ $= \tanh\left(W\begin{pmatrix}h_{t-1}\\x_{t}\end{pmatrix}\right)$

Backpropagation from h_t to h_{t-1} multiplies by W (actually W_{hh}^{T})





Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994 Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



Computing gradient of h₀ involves many factors of W (and repeated tanh)

Largest singular value > 1: Exploding gradients

Gradient clipping: Scale gradient if its norm is too big

Largest singular value < 1: Vanishing gradients

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994 Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



Computing gradient of h₀ involves many factors of W (and repeated tanh) Largest singular value > 1: **Exploding gradients**

Largest singular value < 1: Vanishing gradients Introduction of Long Short Term Memory (LSTM)

Change RNN architecture

11/26/2018 http://colah.github.io/posts/2015-08-Understanding-LSTMs/

Long Short Term Memory (LSTM)

- Was originally introduced Hochreiter & Schmidhuber (1997)
- Idea:
 - To be able to learn long term dependences
 - Collects data when the input is considered to be relevant
 - Keeps it as long as it considers to be important
 - Technique:
 - Handle the state as a memory with minor modifications
 - No matrix multiplication
 - No tanh
 - Apply memory handling kind signals
 - » data in, data out, write, enable





Derivation of LSTM

- Repeating module in Normal RNN
 - concatenates the input and the state
 - A neural network with tanh output and repeats the result
- LSTM
 - Uses the state as a memory
 - Uses 4 neural nets to control the memory
 - Forget, Input, W, Output



Components of LSTM |

- All wires represents vector
 - Vector transfer
 - Vector concatenation >>>
 - Vector copy
- Neural nets with (yellow boxes)
 - Multi-layer NN with tanh activation function used for update value tanh calculation
 - <u>Multi-layer NN</u> with *logistic* activation function (sigmoid) used for <u>value selection (kind of</u> <u>addressing)</u>
- Pointwise operation (pink circles)
 - Pointwise multifaction
 - Pointwise addition



Components of LSTM II

- State of the LSTM
 - This is the actual memory,
 - It can pass the previous values with or without update
 - Represented by the upper black line
 - Indicated with C_t
- Old content can be removed value-by-value
- New content can be ۲ added

 C_{t-1} h_{t-1} x_t Vector

11/26/2018

Neural Network Layer

Pointwise Operation Transfer

Concatenate

Copy





- Combines input and previous output (concatenation)
- Selects which values to forget
 - Sort of addressing
 - Done by the ٠ "Forget Gate"
 - Neural net with sigmoid output



$$f_t = \sigma \left(W_f \cdot [h_{t-1}, x_t] + b_f \right)$$



- Input: "James"
- **Forget Neural** network figures out:
 - Analyzes the concatenated vector
 - Name, Subject of a sentence, Male
- Selects which values to forget and how much
 - Position and weight
- Task:
 - Update gender of the subject (forget the old value)
 - Gender might be represented with a variable
 - c_1 : value proportional with the probability that the subject is a male
 - c₂: represents weather
 - Calculate the forget factor of the gender memories
 - 0 completely get rid of it
 - 1 keep the previous value
 - 0 .. 1 partial forget
 - Adressing and suppressing!!!



- Input: "James"
- Input Gate figures out:
 - Analyze the concatenated vector
 - Select which values to update (ENABLE!!!)
 - Calculate the update weights
- Cell Network calculates:
 - The update values
- Task:
 - Update gender of the subject (calculate the update value)
 - Gender might be represented with a variable
 - c₁: value proportional with the probability that the gender is male
 - c₂: represents weather
 - Calculate the update factor of the gender memories
 - 0 not to update
 - 1 fully update
 - 0 .. 1 partial update
 - 11/26/2018 ADRESSING!!!









ht /



Step 4

٠

٠

 h_t



reduced values

11/26/2018

LSTM network



• General form of an LSTM network







Gradient calculation in LSTM

tanh



Backpropagation from c_t to c_{t-1} only elementwise multiplication by f, no matrix multiply by W



11/26/2018

stack

C_{t-1}

h_{t-1}

64

Gradient calculation in LSTM



Uninterrupted gradient flow!



- Though we multiply the memory content with a smaller than 1 number
- And the W matrix is part of the memory update
- But it still preserves the content for longer time
- As it comes from the name: It is a elongated time <u>short term</u> memory

Achevements with LSTM networks

- Record results in natural language text compression
- Unsegmented connected handwriting recognition
- Natural speech recognition
- Smart voice assistants
 - Google Translate
 - Amazon Alexa
 - Microsoft Cortana
 - Apple Quicktype
- 95.1% recognition accuracy on the Switchboard corpus, incorporating a vocabulary of 165,000 words
 - Continuous spontaneous English native speech



Translate







Variants of LSTM I : Peephole connections h_t

- Introduced by Gers & Schmidhuber (2000)
- All the three gates receives input from the previous state and the input
- Since output can be sparse this version has more information for gating
 - addressing and weighting



11/26/2018

Variants of LSTM II : Joined forget and input ht A

- Input and forget gates has practically the same role
- Why not to join them?



Gated Recurrent Unit (GRU)

• Another variant of LSTM

- h_{t} h_{t-1} tanh x_t
- Introduced by Kyunghyun Cho (2014)
- There is no separate State and Output
- Only three neural nets
- At GRU the output will not be sparse (not gated)
- Similar performance in music and speech signal modelling and
- Learns faster for smaller data set

