

# Introductory Notes on Machine Translation and Deep Learning

February 20, 2017

Jindřich Libovický, Jindřich Helcl



Charles University in Prague  
Faculty of Mathematics and Physics  
Institute of Formal and Applied Linguistics



# What is machine translation?

---

Time for discussion.

## What we think...

---

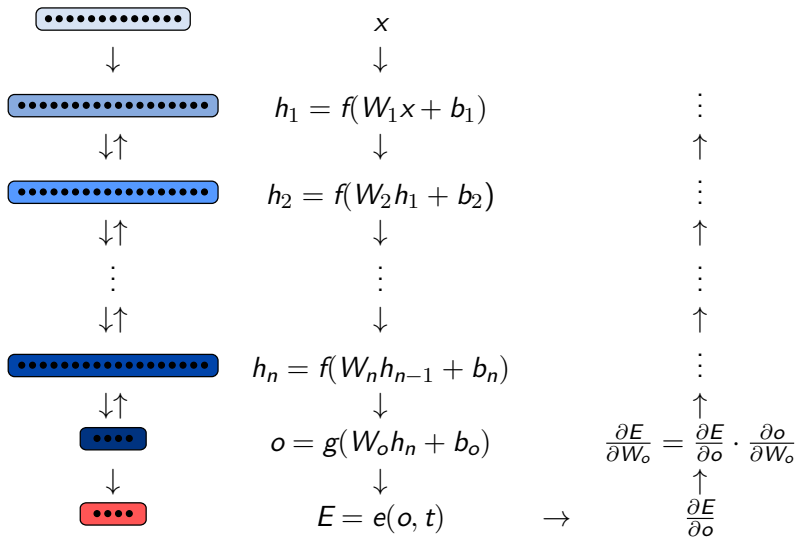
- MT does not care what translation is
- we believe people know what translation is and that it is captured in the data
- we evaluate how well we can mimic what humans do when they translate

# Deep Learning

---

- machine learning that hierarchically infers suitable data representation with the increasing level of complexity and abstraction (Goodfellow et al.)
- formulating end-to-end relation of a problems' raw inputs and raw outputs as parameterizable real-valued functions and finding good parameters for the functions (me)
- industrial/marketing buzzword for machine learning with neural networks (backpropaganda, ha, ha)

# Neural Network



## Building Blocks (1)

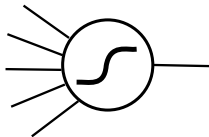
---

- individual neurons / more complex units like recurrent cells  
(allows innovations like inventing LSTM cells, ReLU activation)
- libraries like Keras, Lasagne, TFSlim conceptualize on layer-level  
(allows innovations like batch normalization, dropout)
- sometimes higher-level conceptualization, similar to functional programming concepts  
(allows innovations like attention)

## Building Blocks (2)

---

### Single Neuron



- computational model from 1940's
- adds weighted inputs and transforms to input

### Layer



$$f(Wx + b)$$

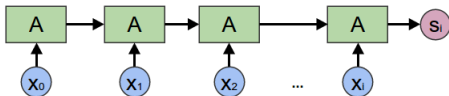
... $f$  nonlinearity,  $W$  ...weight matrix,  $b$  ...bias

- having the network in layers allows using matrix multiplication
- allows GPU acceleration
- vector space interpretations

# Encoder & Decoder

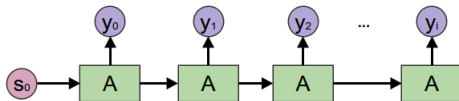
---

Encoder:



Functional fold (reduce) with function  
`foldl a s xs`

Decoder:

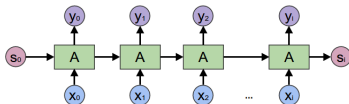


Inverse operation – functional unfold  
`unfoldr a s`



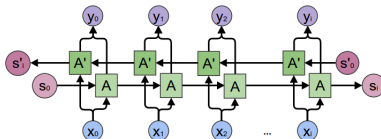
# RNNs & Convolutions

General RNN:



Map with accumulator  
`mapAccumR a s xs`

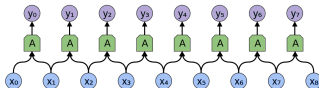
Bidirectional RNN:



Zip left and right accumulating map

`zip (mapAccumR a s xs)`  
`(mapAccumL a' s' xs)`

Convolution:



Zip neighbors and apply function  
`zipWith a xs (tail xs)`

Source: Colah's blog (<http://colah.github.io/posts/2015-09-NN-Types-FP/>)

# Optimization

---

- data is constant, treat the network as function of parameters
- the differentiable error is function of parameters as well
- clever variants of gradient descent algorithm

# Deep Learning as Alchemy

---

- there no rigorous manual how to develop a good deep learning model – just rules of thumb
- we don't know how to interpret the weights the network has learned
- there is no theory that is able to predict results of experiments (as in physics), there are only experiments

# Recoding in mathematics

Algebraic equations

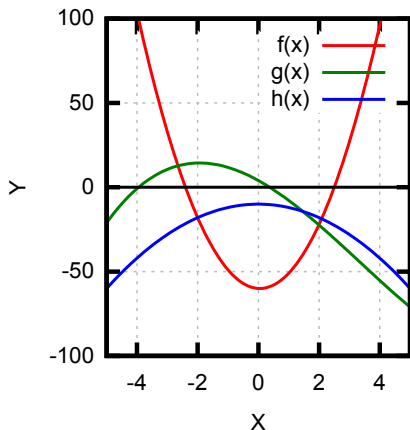
$$10x^2 - x - 60 = 0$$

$$0.2x^3 - 2x^2 - 10x + 4 = 0$$

$$-2x^2 - 10 = 0$$

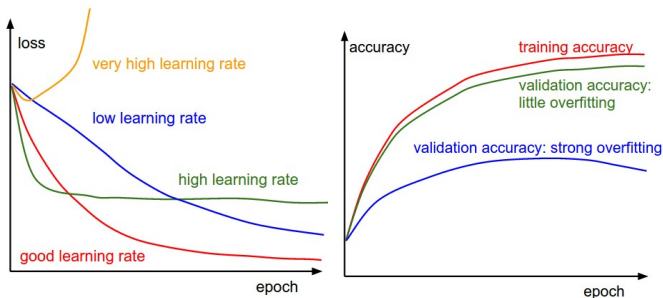


...became planar curves



# Watching Learning Curves

---

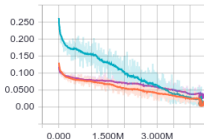


Source: Convolutional Neural Networks for Visual Recognition at Stanford University  
(<http://cs231n.github.io/neural-networks-3/>)

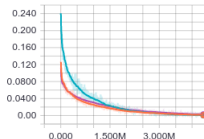
# Other Things to Watch During Training (1)

- train and validation loss

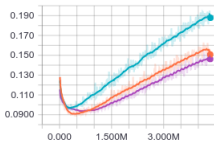
train\_target/runtime\_xent



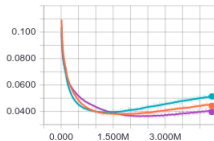
train\_target/train\_xent



val\_target/runtime\_xent



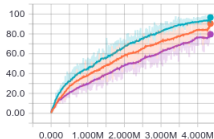
val\_target/train\_xent



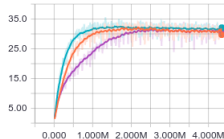
# Other Things to Watch During Training (2)

- target metric on training and validation data

train\_target/BLEU-4

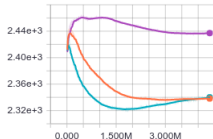


val\_target/BLEU-4

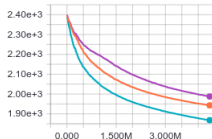


- L2 and L1 norm of parameters

train\_l1



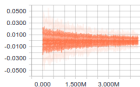
train\_l2



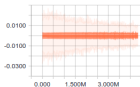
# Other Things to Watch During Training (3)

- gradients of the parameters

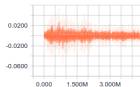
gr\_decoder/attention\_decoder/AttnOutputProje.  
clever/



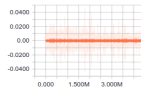
gr\_decoder/attention\_decoder/AttnOutputProje.  
clever/



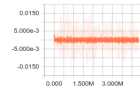
gr\_decoder/attention\_decoder/GRUCell/Candid.  
clever/



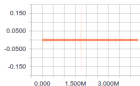
gr\_decoder/attention\_decoder/GRUCell/Candid.  
clever/



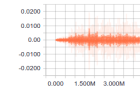
gr\_decoder/attention\_decoder/GRUCell/Gates/  
clever/



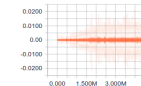
gr\_decoder/attention\_decoder/GRUCell/Gates/  
clever/



gr\_decoder/attention\_decoder/attention\_wrapp.  
clever/



gr\_decoder/attention\_decoder/attention\_wrapp.  
clever/



- non-linearities saturation



# What's Strange on Neural MT

---

- we naturally think of translation in terms of manipulating with symbols
- neural network represents everything as real-space vectors
- ignore pretty much everything we know about language

## Reading for the Next Week

---

LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton. "Deep learning." *Nature* 521.7553 (2015): 436.

<http://pages.cs.wisc.edu/~dyer/cs540/handouts/deep-learning-nature2015.pdf>

Question:

**Can you identify some implicit assumptions the authors make about sentence meaning while talking about NMT? Do you think they are correct? How do the properties that the authors attribute to LSTM networks correspond to your own ideas how should language be computationally processed?**