

Machine Learning (part II)

Recurrent Neural Networks

Angelo Ciaramella

RNNs

- family of neural networks for processing sequential data
- specialized for processing a sequence of values

 $x^{(1)}, \ldots, x^{(\tau)}$

- early ideas found in machine learning and statistical models of the 1980s
 - sharing parameters across different parts of a model



Related idea

use of convolution across a 1-D temporal sequence

time-delay neural networks

- RNNs
 - minibatches of sequences
 - may also be applied in two dimensions across spatial data such as images



Adaptive filters

- Adaptive filter
 - The parameters are estimated
 - learning algorithm
 - An error function is used
 - e.g., Linear Artificial Neural Network (Adaline)



Adaptive filters

Hospital

- ECG (electrocardiogram) corrupted by noise at 50 Hz (electricity)
- The current can vary between 47 Hz and 53 Hz
- A filter for the elimination of static noise at 50 Hz could give errors
- An adaptive filter can learn from the current shape of noise

Helicopter

- Pilot speaking with noise from rotating propeller
- The noise has not a spectrum well defined
- An adaptive filter learns the shape of the noise
- The noise can be subtracted from the signal for only the pilot's voice



Adaline





Adaptive filters





Related idea

- formalize the structure of a set of computations
- introduce the idea of an operation
 - an operation is a simple function of one or more variables





Graph using the x operation to compute z = xy





logistic regression prediction $\hat{y} = \sigma \left(\boldsymbol{x}^{\mathsf{T}} \boldsymbol{w} + b \right)$



Minibatch of inputs X





More than one operation of a linear regression model



RNNs

- Artificial Neural Networks
 - exhibit temporal dynamic behavior
 - can use their internal state (memory) to process sequences of inputs
 - models sequences
 - Time series
 - Natural Language
 - Speech
 - Convert non-sequences to sequences, eg: feed an image as a sequence of pixels!



Feed-forward NN



Temporal dipendencies



sequence of observations



Reber Grammar

Problem that can not be solved without memory



Jordan's sequential network

Limited short-term memory



M.I. Jordan NN (1986).



Jordan's sequential network



Jordan NN has been applied to categorize a class of English syllables.



Simple recurrent network





Elman RNN (1990).

Simple recurrent network



Elman RNN. Basic RNN structure called "Vanilla" RNN

Simple recurrent network



Elman SRN. A total of 60,000 randomly generated strings are used for training.

Applications of RNNs

A person riding a motorcycle on a dirt road.



Image Captioning

.. and Trump

12 12

Twitterbot

4 1

I'm a Neural Network trained on Trump's transcripts. Priming text in []s. Donate (http://www.gofundme.com/deepdrumpf) to interact! Created by @hayesbh.

17

RNN Generated Music **RNN** Generated Eminem rapper

... and more!

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;

Write like Shakespeare

 In reply to Thomas Paine DeepDrumpf @DeepDrumpf · Mar 20 There will be no amnesty. It is going to pass because the people are going to be gone. I'm giving a mandate. #ComeyHearing @Thomas1774Paine



Dynamical system (reccurrent expression)

$$s^{(t)} = f(s^{(t-1)}; \boldsymbol{\theta})$$

State of the system



23

Dynamical system driven by an external signal

$$s^{(t)} = f(s^{(t-1)}, x^{(t)}; \theta)$$



$$\boldsymbol{h}^{(t)} = f(\boldsymbol{h}^{(t-1)}, \boldsymbol{x}^{(t)}; \boldsymbol{\theta})$$

state of the hidden units of the network





This recurrent network just processes information from the input x by incorporating it into the state h that is passed forward through time

Representation of the unfolded recurrence

$$\begin{aligned} \boldsymbol{h}^{(t)} = & g^{(t)} \left(\boldsymbol{x}^{(t)}, \boldsymbol{x}^{(t-1)}, \boldsymbol{x}^{(t-2)}, \dots, \boldsymbol{x}^{(2)}, \boldsymbol{x}^{(1)} \right) \\ = & f(\boldsymbol{h}^{(t-1)}, \boldsymbol{x}^{(t)}; \boldsymbol{\theta}) \end{aligned}$$





universal – any function computable by a Turing machine can be computed by such a recurrent network of a finite size

Recurrent networks that produce an output at each time step and have recurrent connections between hidden units



ML – Recurrent NNs Recurrent networks that produce an output at each time step and have recurrent connections only from the output at one time step to the hidden units at the next time step



Recurrent networks with recurrent connections between hidden units, that read an entire sequence and then produce a single output

Vanilla RNN cell



$$h_t = f(W_h h_{t-1} + W_x x_t)$$





Unfolding





Deep RNN



Feedforward depth = 4



Recurrent depth = 3





Objective is to update the weight matrix:

$$\mathbf{W} \to \mathbf{W} - \alpha \frac{\partial L}{\partial \mathbf{W}}$$

Issue: W occurs each timestep Every path from W to L is one dependency

Find all paths from **W** to L!

(note: dropping subscript h from **W**_h for brevity)



How many paths exist from W to L through L_1 ?

Just 1. Originating at h_0 .





How many paths from **W** to L through L₂?

2. Originating at h_0 and h_1 .



35



And 3 in this case.

Origin of path = basis for Σ

$\frac{\partial L}{\partial \mathbf{W}}$

The gradient has two summations: 1: Over L_j 2: Over h_k


ML – Recurrent NNs

First summation over L

 $\frac{\partial L}{\partial \mathbf{W}} = \sum_{i=0}^{I-1} \frac{\partial L_j}{\partial \mathbf{W}}$



 Second summation over h: Each L_j depends on the weight matrices *before it*

$$\frac{\partial L_j}{\partial \mathbf{W}} = \sum_{k=1}^{j} \frac{\partial L_j}{\partial h_k} \frac{\partial h_k}{\partial \mathbf{W}}$$

$$\mathsf{L}_j \text{ depends on all } \mathsf{h}_k \text{ before it.}$$















$$\frac{\partial L}{\partial \mathbf{W}_{\mathbf{h}}} = \sum_{j=0}^{T-1} \sum_{k=1}^{j} \frac{\partial L_{j}}{\partial y_{j}} \frac{\partial y_{j}}{\partial h_{j}} \left(\prod_{m=k+1}^{j} \frac{\partial h_{m}}{\partial h_{m-1}}\right) \frac{\partial h_{k}}{\partial \mathbf{W}_{\mathbf{h}}}$$

The final Backpropagation equation







$$\frac{\partial L}{\mathbf{W}_{\mathbf{h}}} = \sum_{j=0}^{T-1} \sum_{k=1}^{j} \frac{\partial L_{j}}{\partial y_{j}} \frac{\partial y_{j}}{\partial h_{j}} \left(\prod_{m=k+1}^{j} \frac{\partial h_{m}}{\partial h_{m-1}} \right) \frac{\partial h_{k}}{\partial \mathbf{W}_{\mathbf{h}}}$$

- Often, to reduce memory requirement, we truncate the network
- Inner summation runs from *j-p* to *j* for some *p* ==> truncated BPTT

$$\frac{\partial L}{\partial W} = \sum_{j=0}^{T-1} \sum_{k=1}^{j} \frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \left(\prod_{m=k+1}^{j} \frac{\partial h_m}{\partial h_{m-1}} \right) \frac{\partial h_k}{\partial \mathbf{W}}$$

$$h_m = f(\mathbf{W}_h h_{m-1} + \mathbf{W}_x x_m)$$

$$\frac{\partial h_m}{\partial h_{m-1}} = \mathbf{W}_h^T \operatorname{diag}\left(f'(\mathbf{W}_h h_{m-1} + \mathbf{W}_x x_m)\right)$$

Expanding the Jacobian





Repeated matrix multiplications leads to **vanishing and exploding** gradients.



Vanishing gradients



$$\frac{\partial h_t}{\partial h_0} = \frac{\partial h_t}{\partial h_{t-1}} \cdot \dots \cdot \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial h_0}$$

Known problem for deep feed-forward networks. For recurrent networks (even shallow) makes **impossible to learn long-term** dependencies!

Considerations

- Smaller weight parameters lead to faster gradients vanishing
- Very big initial parameters make the gradient descent to diverge fast (explode)



Exploiding gradients



Large increase in the norm of the gradient during training

Diagnostics: NaNs; Cost function large fluctuations

Solutions:

- Use gradient clipping
- Try reduce learning rate
- Change loss function by setting constrains on weights (L1/L2 norms)



Eigenvalues and Stability

Consider identity activation function If Recurrent Matrix W_h is a diagonalizable:

Q matrix composed of eigenvectors of W_h

∧ is a diagonal matrix with eigenvalues placed on the diagonals

Computing powers of **W**_h is simple:

$$W_h^n = Q^{-1} * \Lambda^n * Q$$

 $W_h = Q^{-1} * \Lambda * Q$





Eigenvalues and Stability



Fundamental DL problem

- DNNs train difficulties
 - Vanishing gradient
 - Exploiding gradient
- Solutions
 - Previously proposed
 - Unsupervised pre-training
 - Improve network architecture



RNNs - forward propagation

- Assume the hyperbolic tangent activation function
- Initial state h⁽⁰⁾
- Update equation

$$m{a}^{(t)} = m{b} + m{W}m{h}^{(t-1)} + m{U}m{x}^{(t)}$$

$$\begin{aligned} \boldsymbol{h}^{(t)} &= \tanh(\boldsymbol{a}^{(t)}) \\ \boldsymbol{o}^{(t)} &= \boldsymbol{c} + \boldsymbol{V} \boldsymbol{h}^{(t)} \\ \hat{\boldsymbol{y}}^{(t)} &= \operatorname{softmax}(\boldsymbol{o}^{(t)}) \end{aligned}$$



ML – Recurrent NNs

RNNs - forward propagation

Total loss

$$L\left(\{x^{(1)}, \dots, x^{(\tau)}\}, \{y^{(1)}, \dots, y^{(\tau)}\}\right)$$

= $\sum_{t} L^{(t)}$
= $-\sum_{t} \log p_{\text{model}}\left(y^{(t)} \mid \{x^{(1)}, \dots, x^{(t)}\}\right)$

Negative log-likelihood



RNNs - Teacher forcing





Illustration of teacher forcing

RNNs - learning

- back-propagation through time (BPTT) algorithm
- For each node N we need to compute the gradient recursively
 - based on the gradient computed at nodes that follow it in the graph

$\nabla_{N}L$

Start the recursion

$$\frac{\partial L}{\partial L^{(t)}} = 1$$



Gradient on the outputs at time step t, for all i, t,

$$\left(\nabla_{\boldsymbol{o}^{(t)}}L\right)_{i} = \frac{\partial L}{\partial o_{i}^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial o_{i}^{(t)}}$$

$$\log \operatorname{softmax}(\boldsymbol{z})_i = z_i - \log \sum_j \exp(z_j)$$

$$\log \sum_{j} \exp(z_j) \approx \max_{j} z_j = z_i$$



Backwards starting from the end of the sequence

$$\nabla_{\boldsymbol{h}^{(\tau)}} L = \boldsymbol{V}^{\top} \nabla_{\boldsymbol{o}^{(\tau)}} L$$

Back-propagate gradients through time

$$\nabla_{\boldsymbol{h}^{(t)}} L = \left(\frac{\partial \boldsymbol{h}^{(t+1)}}{\partial \boldsymbol{h}^{(t)}}\right)^{\top} (\nabla_{\boldsymbol{h}^{(t+1)}} L) + \left(\frac{\partial \boldsymbol{o}^{(t)}}{\partial \boldsymbol{h}^{(t)}}\right)^{\top} (\nabla_{\boldsymbol{o}^{(t)}} L)$$
$$= \boldsymbol{W}^{\top} (\nabla_{\boldsymbol{h}^{(t+1)}} L) \operatorname{diag} \left(1 - \left(\boldsymbol{h}^{(t+1)}\right)^{2}\right) + \boldsymbol{V}^{\top} (\nabla_{\boldsymbol{o}^{(t)}} L)$$

Once the gradients on the internal nodes of the computational graph are obtained, we can obtain the gradients on the parameter nodes

RNNs - learning

For all the parameters

$$\nabla_{\mathbf{c}} L = \sum_{t} \left(\frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{c}} \right)^{\mathsf{T}} \nabla_{\mathbf{o}^{(t)}} L = \sum_{t} \nabla_{\mathbf{o}^{(t)}} L$$
$$\nabla_{\mathbf{b}} L = \sum_{t} \left(\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{b}^{(t)}} \right)^{\mathsf{T}} \nabla_{\mathbf{h}^{(t)}} L = \sum_{t} \operatorname{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^{2} \right) \nabla_{\mathbf{h}^{(t)}} L$$
$$\nabla_{\mathbf{V}} L = \sum_{t} \sum_{i} \left(\frac{\partial L}{\partial o_{i}^{(t)}} \right) \nabla_{\mathbf{V}} o_{i}^{(t)} = \sum_{t} \left(\nabla_{\mathbf{o}^{(t)}} L \right) \mathbf{h}^{(t)^{\mathsf{T}}}$$
$$\nabla_{\mathbf{W}} L = \sum_{t} \sum_{i} \left(\frac{\partial L}{\partial h_{i}^{(t)}} \right) \nabla_{\mathbf{W}^{(t)}} h_{i}^{(t)}$$
$$= \sum_{t} \operatorname{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^{2} \right) \left(\nabla_{\mathbf{h}^{(t)}} L \right) \mathbf{h}^{(t-1)^{\mathsf{T}}}$$
$$\nabla_{\mathbf{U}} L = \sum_{t} \sum_{i} \left(\frac{\partial L}{\partial h_{i}^{(t)}} \right) \nabla_{\mathbf{U}^{(t)}} h_{i}^{(t)}$$
$$= \sum_{t} \operatorname{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^{2} \right) \left(\nabla_{\mathbf{h}^{(t)}} L \right) \mathbf{x}^{(t)^{\mathsf{T}}}$$



RNNs - Bidirectional



ML – Recurrent NNs

prediction of y(t) which may depend on the whole input sequence e.g., speech recognition

58

RNNs - Encoder and Decoder



encoder-decoder or sequence-to-sequence RNN architecture

To map input sequence to an output sequence which is not necessarily the same lenght: NLP, speech recongition, ...

Attention meccanism to C could be added

Deep Recurrent Networks



The hidden recurrent state can be broken down into groups organized hierarchically



ML – Recurrent NNs

Deep Recurrent Networks



Deeper computation (e.g., an MLP) can be introduced in the input-to-hidden, hidden-to-hidden and hidden-to-output parts. This may lengthen the shortest path linking different time steps.

ML – Recurrent NNs

Deep Recurrent Networks



ML – Recurrent NNs

The path-lengthening effect can be mitigated by introducing skip connections

Recursive NNs



Generalization of recurrent networks Applied for structurated data

Long-Term dependendencies





Long-Term dependendencies

- Random W_h initialization of RNN has no constraint on eigenvalues
 - vanishing or exploding gradients in the initial epoch

- Careful initialization of W_h with suitable eigenvalues
 - allows the RNN to learn in the initial epochs
 - hence can generalize well for further iterations

ML – Recurrent NNs

Long-Term dependendencies

Trick #1(IRNN)

- W_h initialized to Identity
- Activation function: ReLU

- Trick# 2 (np-RNN)
 - Wh positive definite (+ve real eigenvalues)
 - At least one eigenvalue is 1, others all less than equal to one
 - Activation function: ReLU

ML – Recurrent NNs

Long Short-Term Memory



- Gated RNNs
 - Long Short-Term memory
 - Gated Recurrent Unit
- Idea
 - creating paths through time that have derivatives that neither vanish nor explode



Gated RNNs

ML – Recurrent NNs



69

LSTM cell



LSTM cell



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

$$i_t = \sigma \left(W_i \cdot [h_{t-1}, x_t] + b_i \right)$$
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C$$

i_t decides what component is to be updated. C'_t provides change contents



LSTM cell



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Updating the cell state

$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$
 Decide what part of
 $h_t = o_t * \tanh (C_t)$ the cell state to output


RNN vs LSTM



Contraction of the

ML – Recurrent NNs

Peephole LSTM



$$f_{t} = \sigma \left(W_{f} \cdot [C_{t-1}, h_{t-1}, x_{t}] + b_{f} \right)$$

$$i_{t} = \sigma \left(W_{i} \cdot [C_{t-1}, h_{t-1}, x_{t}] + b_{i} \right)$$

$$o_{t} = \sigma \left(W_{o} \cdot [C_{t}, h_{t-1}, x_{t}] + b_{o} \right)$$

Allows "peeping into the memory". Can learn the fine distinction between sequences of spikes separated by either 50 or 49 discrete time steps



Gated Recurrent Unit (GRU)



$$z_t = \sigma \left(W_z \cdot [h_{t-1}, x_t] \right)$$
$$r_t = \sigma \left(W_r \cdot [h_{t-1}, x_t] \right)$$
$$\tilde{h}_t = \tanh \left(W \cdot [r_t * h_{t-1}, x_t] \right)$$
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

It combines the forget and input into a single update gate. It also merges the cell state and hidden state. This is simpler than LSTM. There are many other variants too.



Clipping gradients



parameter gradient is very large



"landscape" in which one finds "cliffs"

 $\begin{array}{l} \text{if } ||\boldsymbol{g}|| > v \\ \boldsymbol{g} \leftarrow \frac{\boldsymbol{g} v}{||\boldsymbol{g}||} \end{array}$

Clipping the gradient





Saliency Heatmap



RNN vs LSTM





Sequence to sequence chat model



LSTM Encoder

LSTM Decoder



Speech recognition example (Deep Speech)



Andrew Ng

Reservoir computing

- The equivalent idea for RNNs
 - fix the input- hidden connections and the hidden-hidden connections at random values
 - only learn the hidden-output connections
- The learning is then very simple (assuming linear output units)
- Its important to set the random connections very carefully so the RNN does not explode or die
- See also Liquid State Machine



ML – Recurrent NNs

Reservoir computing

Herbert Jaeger, 2001





Reservoir computing



ML – Recurrent NNs

