

Recurrent Neural Networks

Samuel Cheng

School of ECE
University of Oklahoma

Spring, 2018

(Slides credit to Stanford CS231n and Hinton et al.)

Table of Contents

- 1 Motivation
- 2 Basic RNN
- 3 LSTM
- 4 Example: simple character-level language model
- 5 Example: image captioning
- 6 Overview of echo state networks
- 7 Conclusions

Review and Overview

- We looked into couple use cases of CNNs previously
 - Recognition and localization
 - Object detection
 - Some use of CNNs for arts
- Up to now, the network models we have studied are all memoryless
 - We will discuss a non-memoryless model—recurrent neural networks today

Why non-memoryless models

- Almost all natural signals are sequential if we take time into account (we just cannot escape time)
 - Memory is needed to remember the past
- They also offer a simplified solution for some problems (for example, number addition)
- They can treat some unsupervised problems as supervised problems
 - Consider prediction of a stock: unsupervised? Supervised?

Why non-memoryless models

- Almost all natural signals are sequential if we take time into account (we just cannot escape time)
 - Memory is needed to remember the past
- They also offer a simplified solution for some problems (for example, number addition)
- They can treat some unsupervised problems as supervised problems
 - Consider prediction of a stock: unsupervised? Supervised?

Why non-memoryless models

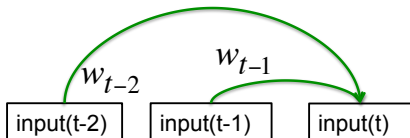
- Almost all natural signals are sequential if we take time into account (we just cannot escape time)
 - Memory is needed to remember the past
- They also offer a simplified solution for some problems (for example, number addition)
- They can treat some unsupervised problems as supervised problems
 - Consider prediction of a stock: unsupervised? Supervised?

[Hinton 2012, week 7]

Memoryless models for sequences

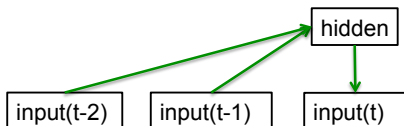
- Autoregressive models

Predict the next term in a sequence from a fixed number of previous terms using “delay taps”.



- Feed-forward neural nets

These generalize autoregressive models by using one or more layers of non-linear hidden units. *e.g. Bengio's first language model.*



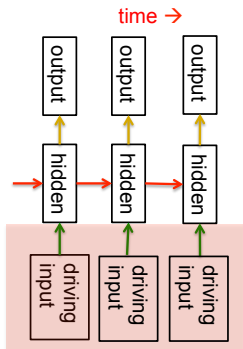
Non-memoryless models

- **Benefit: memories increase the expressive power of the model**
- Typically we do not know the exact values of the hidden states (that is why “hidden”). In many cases, the best we could do is just to infer a probability distribution over the hidden states
- Let's look at two classic examples

Non-memoryless models

- Benefit: memories increase the expressive power of the model
- Typically we do not know the exact values of the hidden states (that is why “hidden”). In many cases, the best we could do is just to infer a probability distribution over the hidden states
- Let's look at two classic examples

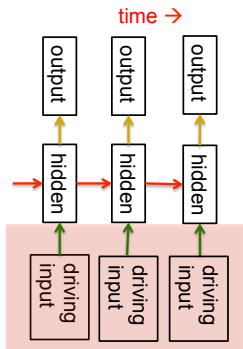
Linear dynamical systems (Engineers love them!)



- These are generative models with real **continuous** values as hidden states that cannot be observed directly
 - The hidden state has linear dynamics with Gaussian noise and produces the observations subjected to linear Gaussian noise
 - There can also be driving inputs
- To predict next output, we need to infer the hidden state

[Hinton 2012, Week 7]

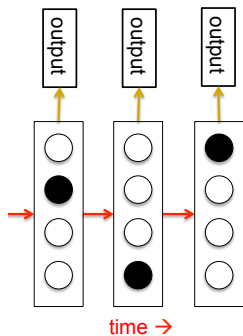
Linear dynamical systems (Engineers love them!)



- These are generative models with real **continuous** values as hidden states that cannot be observed directly
 - The hidden state has linear dynamics with Gaussian noise and produces the observations subjected to linear Gaussian noise
 - There can also be driving inputs
- To predict next output, we need to infer the hidden state

[Hinton 2012, Week 7]

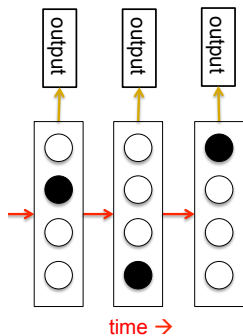
State-Space Models (Computer scientists love them!)



- State-Space Models or Hidden Markov Models (HMMs) have a **discrete** one-of- N hidden state. Transitions between states are stochastic and controlled by a transition matrix. The output produced by a state are also stochastic
 - We don't know which state produced a given output. So the state is "hidden"
 - We can represent the probability distribution across N states with N numbers
- To predict next output, we need to infer the probability distribution over the hidden state

[Hinton 2012, Week 7]

State-Space Models (Computer scientists love them!)



- State-Space Models or Hidden Markov Models (HMMs) have a **discrete** one-of- N hidden state. Transitions between states are stochastic and controlled by a transition matrix. The output produced by a state are also stochastic
 - We don't know which state produced a given output. So the state is "hidden"
 - We can represent the probability distribution across N states with N numbers
- To predict next output, we need to infer the probability distribution over the hidden state

[Hinton 2012, Week 7]

A fundamental limitation of state space models

- The only information stored in the model is which state the model currently is in
 - So with N hidden states it can only remember a maximum $\log(N)$ bits of information
- Consider the speech prediction of one half from earlier half
 - The syntax needs to fit (e.g. number and tense agreement)
 - The semantics needs to fit. The intonation needs to fit
 - The accent, rate, volume, and vocal tract characteristics must all fit
- All these aspects combined could be 100 bits of information that the first half of an utterance needs to convey to the second half
 2^{100} states

[Hinton 2012, week 7]

A fundamental limitation of state space models

- The only information stored in the model is which state the model currently is in
 - So with N hidden states it can only remember a maximum $\log(N)$ bits of information
- Consider the speech prediction of one half from earlier half
 - The syntax needs to fit (e.g. number and tense agreement)
 - The semantics needs to fit. The intonation needs to fit
 - The accent, rate, volume, and vocal tract characteristics must all fit
- All these aspects combined could be 100 bits of information that the first half of an utterance needs to convey to the second half
 2^{100} states

[Hinton 2012, week 7]

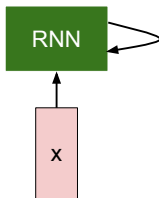
A fundamental limitation of state space models

- The only information stored in the model is which state the model currently is in
 - So with N hidden states it can only remember a maximum $\log(N)$ bits of information
- Consider the speech prediction of one half from earlier half
 - The syntax needs to fit (e.g. number and tense agreement)
 - The semantics needs to fit. The intonation needs to fit
 - The accent, rate, volume, and vocal tract characteristics must all fit
- All these aspects combined could be 100 bits of information that the first half of an utterance needs to convey to the second half
 2^{100} states

[Hinton 2012, week 7]

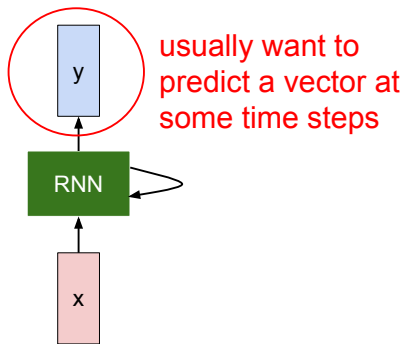
Recurrent neural networks (RNNs)

Recurrent Neural Network



Recurrent neural networks (RNNs)

Recurrent Neural Network



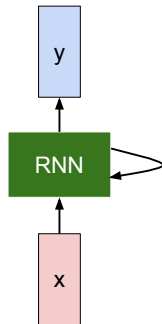
Recurrent neural networks (RNNs)

Recurrent Neural Network

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

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

new state some function with parameters W old state input vector at some time step



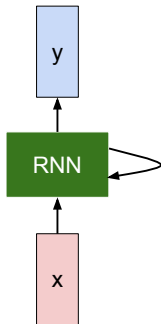
Recurrent neural networks (RNNs)

Recurrent Neural Network

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

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

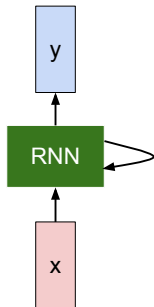
Notice: the same function and the same set of parameters are used at every time step.



Recurrent neural networks (RNNs)

(Vanilla) Recurrent Neural Network

The state consists of a single “hidden” vector h :



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

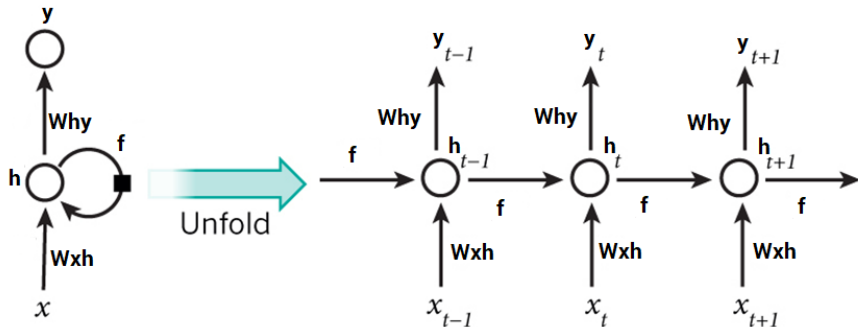


$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

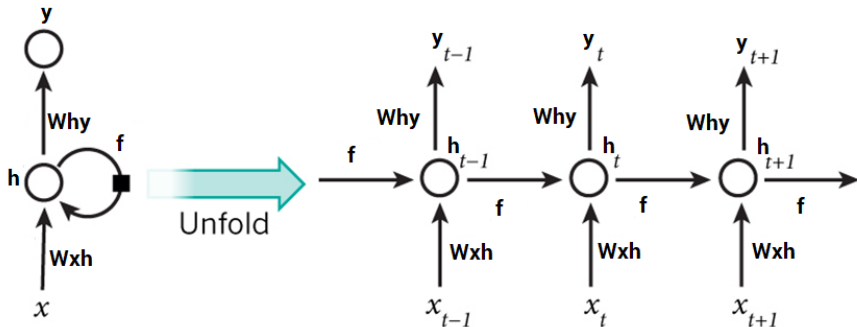
Back-Propagation Through Time (BPTT)

- For training, we can unroll all the time step to form a stack of activities and backprop will then similar to regular backprop
- The backward pass peels activities off the stack to compute the error derivatives at each time step
- After the backward pass we add together the derivatives at all the different times for each weight



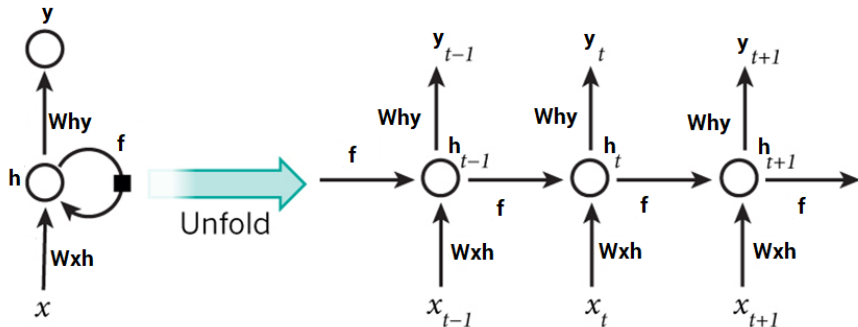
Back-Propagation Through Time (BPTT)

- For training, we can unroll all the time step to form a stack of activities and backprop will then similar to regular backprop
- The backward pass peels activities off the stack to compute the error derivatives at each time step
- After the backward pass we add together the derivatives at all the different times for each weight



Back-Propagation Through Time (BPTT)

- For training, we can unroll all the time step to form a stack of activities and backprop will then similar to regular backprop
- The backward pass peels activities off the stack to compute the error derivatives at each time step
- After the backward pass we add together the derivatives at all the different times for each weight



An irritative extra issue

- We need to specify the initial activity state of all the hidden and output units
- We could just fix these initial states to have some default value like 0.5
- But it is better to treat the initial states as learned parameters
- We learn them in the same way as we learn the weights
 - Start off with an initial random guess for the initial states
 - At the end of each training sequence, backpropagate through time all the way to the initial states to get the gradient of the error function with respect to each initial state
 - Adjust the initial states by following the negative gradient

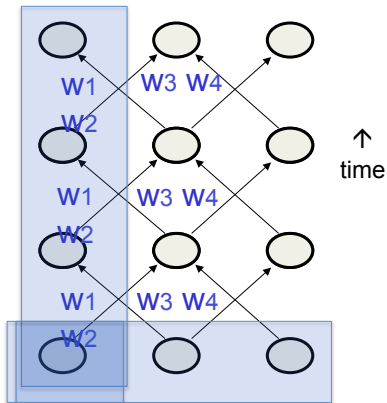
An irritative extra issue

- We need to specify the initial activity state of all the hidden and output units
- We could just fix these initial states to have some default value like 0.5
- But it is better to treat the initial states as learned parameters
- We learn them in the same way as we learn the weights
 - Start off with an initial random guess for the initial states
 - At the end of each training sequence, backpropagate through time all the way to the initial states to get the gradient of the error function with respect to each initial state
 - Adjust the initial states by following the negative gradient

An irritative extra issue

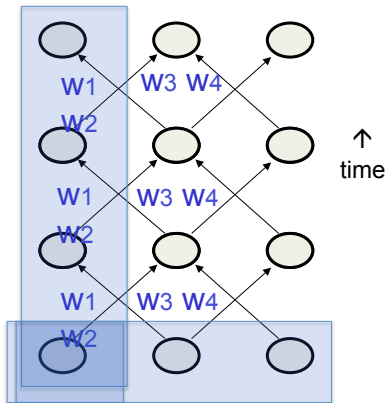
- We need to specify the initial activity state of all the hidden and output units
- We could just fix these initial states to have some default value like 0.5
- But it is better to treat the initial states as learned parameters
- We learn them in the same way as we learn the weights
 - Start off with an initial random guess for the initial states
 - At the end of each training sequence, backpropagate through time all the way to the initial states to get the gradient of the error function with respect to each initial state
 - Adjust the initial states by following the negative gradient

Providing inputs to recurrent networks



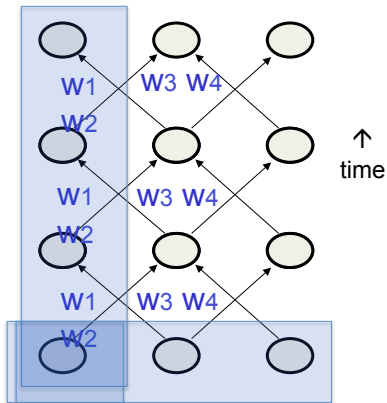
- We can specify inputs in several ways:
 - Specify the initial states of all the units
 - Specify the initial states of a subset of the units
 - Specify the states of the same subset of the units at every time step

Providing inputs to recurrent networks



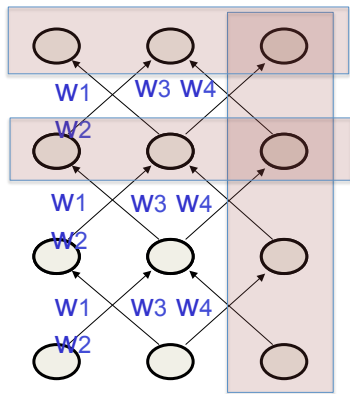
- We can specify inputs in several ways:
 - Specify the initial states of all the units
 - Specify the initial states of a subset of the units
 - Specify the states of the same subset of the units at every time step

Providing inputs to recurrent networks



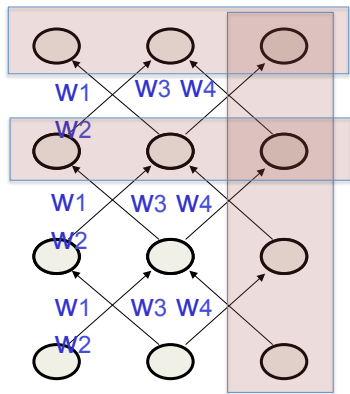
- We can specify inputs in several ways:
 - Specify the initial states of all the units
 - Specify the initial states of a subset of the units
 - Specify the states of the same subset of the units at every time step

Teaching recurrent networks to learn signals



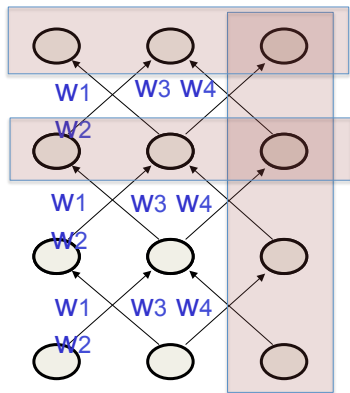
- We can specify targets in several ways:
 - Specify desired final activities of all the units
 - Specify desired activities of all units for the last few steps
 - Good for learning attractors
 - Specify the desired activity of a subset of the units.
 - The other units are input or hidden units.

Teaching recurrent networks to learn signals



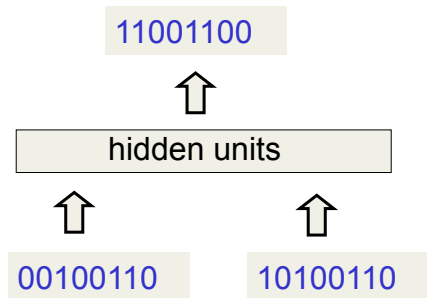
- We can specify targets in several ways:
 - Specify desired final activities of all the units
 - Specify desired activities of all units for the last few steps
 - Good for learning attractors
 - Specify the desired activity of a subset of the units.
 - The other units are input or hidden units.

Teaching recurrent networks to learn signals



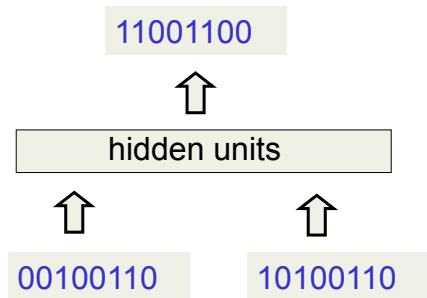
- We can specify targets in several ways:
 - Specify desired final activities of all the units
 - Specify desired activities of all units for the last few steps
 - Good for learning attractors
 - Specify the desired activity of a subset of the units.
 - The other units are input or hidden units.

Toy problem for RNN: binary addition



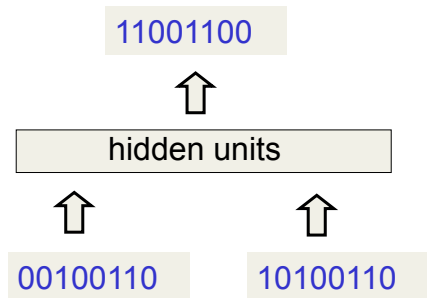
- We can train a feedforward net to do binary addition, but...
 - We must decide in advance the maximum number of digits in each number
 - We expect weights to process different bits to be the same, but it is tricky to enforce that
- As a result, feedforward nets do not generalize well for the binary addition task

Toy problem for RNN: binary addition



- We can train a feedforward net to do binary addition, but...
 - We must decide in advance the maximum number of digits in each number
 - We expect weights to process different bits to be the same, but it is tricky to enforce that
- As a result, feedforward nets do not generalize well for the binary addition task

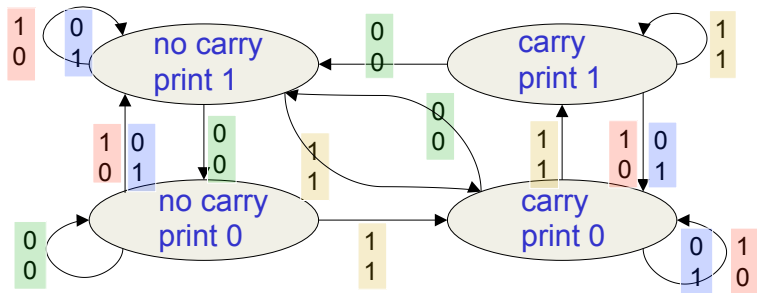
Toy problem for RNN: binary addition



- We can train a feedforward net to do binary addition, but...
 - We must decide in advance the maximum number of digits in each number
 - We expect weights to process different bits to be the same, but it is tricky to enforce that
- As a result, feedforward nets do not generalize well for the binary addition task

We are trying to learn this!

The algorithm for binary addition



This is a finite state automaton. It decides what transition to make by looking at the next column. It prints after making the transition. It moves from right to left over the two input numbers.

A little bit detail

$$x = [b_8, b_7, \dots, b_1]$$

$$y = [c_8, c_7, \dots, c_1]$$

$$z = x + y = [d_8, d_7, \dots, d_1]$$

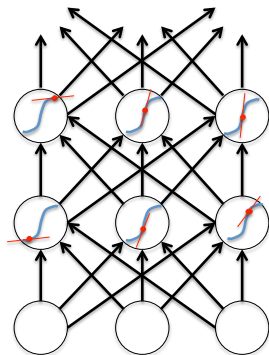
$$\hat{z} = [\hat{d}_8, \hat{d}_7, \dots, \hat{d}_1]$$

Hidden unit: $h_i = \text{sigm}(W_{x,h}[b_i, c_i]^T + W_{h,h}h_{i-1})$

Output: $\hat{d}_i = \text{sigm}(W_{h,z}h_i)$

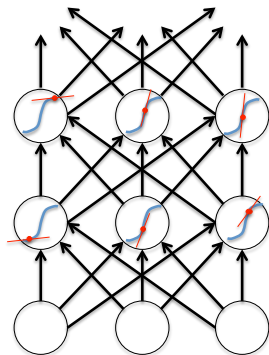
https://github.com/llSourcecell/recurrent_neural_net_demo

Why training RNN is difficult? The backward pass is linear



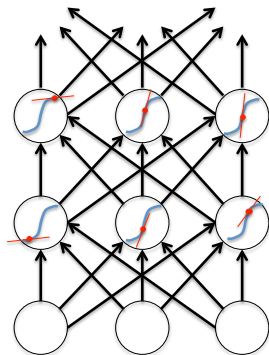
- There is a big difference between the forward and backward passes
- In the forward pass we use squashing functions (like the logistic) to prevent the activity vectors from exploding
- The backward pass, is completely linear. If you double the error derivatives at the final layer, all the error derivatives will double
 - The forward pass determines the slope of the linear function used for backpropagating through each neuron

Why training RNN is difficult? The backward pass is linear



- There is a big difference between the forward and backward passes
- In the forward pass we use squashing functions (like the logistic) to prevent the activity vectors from exploding
- The backward pass, is completely linear. If you double the error derivatives at the final layer, all the error derivatives will double
 - The forward pass determines the slope of the linear function used for backpropagating through each neuron

Why training RNN is difficult? The backward pass is linear



- There is a big difference between the forward and backward passes
- In the forward pass we use squashing functions (like the logistic) to prevent the activity vectors from exploding
- The backward pass, is completely linear. If you double the error derivatives at the final layer, all the error derivatives will double
 - The forward pass determines the slope of the linear function used for backpropagating through each neuron

The problem of exploding or vanishing gradients

- What happens to the magnitude of the gradients as we backpropagate through many layers?
 - If the weights are small, the gradients shrink exponentially.
 - If the weights are big the gradients grow exponentially
- Typical feed-forward neural nets can cope with these exponential effects when they only have a few hidden layers
- In an RNN trained on long sequences (e.g. 100 time steps) the gradients can easily explode or vanish
 - We could avoid this by initializing the weights very carefully
- Even with good initial weights, the dependency of the current target output from an input many time-steps ago tends to be numerically unstable
 - So RNNs have difficulty dealing with long-range dependencies

The problem of exploding or vanishing gradients

- What happens to the magnitude of the gradients as we backpropagate through many layers?
 - If the weights are small, the gradients shrink exponentially.
 - If the weights are big the gradients grow exponentially
- Typical feed-forward neural nets can cope with these exponential effects when they only have a few hidden layers
- In an RNN trained on long sequences (e.g. 100 time steps) the gradients can easily explode or vanish
 - We could avoid this by initializing the weights very carefully
- Even with good initial weights, the dependency of the current target output from an input many time-steps ago tends to be numerically unstable
 - So RNNs have difficulty dealing with long-range dependencies

The problem of exploding or vanishing gradients

- What happens to the magnitude of the gradients as we backpropagate through many layers?
 - If the weights are small, the gradients shrink exponentially.
 - If the weights are big the gradients grow exponentially
- Typical feed-forward neural nets can cope with these exponential effects when they only have a few hidden layers
- In an RNN trained on long sequences (e.g. 100 time steps) the gradients can easily explode or vanish
 - We could avoid this by initializing the weights very carefully
- Even with good initial weights, the dependency of the current target output from an input many time-steps ago tends to be numerically unstable
 - So RNNs have difficulty dealing with long-range dependencies

The problem of exploding or vanishing gradients

- What happens to the magnitude of the gradients as we backpropagate through many layers?
 - If the weights are small, the gradients shrink exponentially.
 - If the weights are big the gradients grow exponentially
- Typical feed-forward neural nets can cope with these exponential effects when they only have a few hidden layers
- In an RNN trained on long sequences (e.g. 100 time steps) the gradients can easily explode or vanish
 - We could avoid this by initializing the weights very carefully
- Even with good initial weights, the dependency of the current target output from an input many time-steps ago tends to be numerically unstable
 - So RNNs have difficulty dealing with long-range dependencies

Understanding gradient flow dynamics

Cute backprop signal video: <http://imgur.com/gallery/vaNahKE>

```
H = 5 # dimensionality of hidden state
T = 50 # number of time steps
Whh = np.random.randn(H,H)

# forward pass of an RNN (ignoring inputs x)
hs = {}
ss = {}
hs[-1] = np.random.randn(H)
for t in xrange(T):
    ss[t] = np.dot(Whh, hs[t-1])
    hs[t] = np.maximum(0, ss[t])

# backward pass of the RNN
dhs = {}
dss = {}
dhs[T-1] = np.random.randn(H) # start off the chain with random gradient
for t in reversed(xrange(T)):
    dss[t] = (hs[t] > 0) * dhs[t] # backprop through the nonlinearity
    dhs[t-1] = np.dot(Whh.T, dss[t]) # backprop into previous hidden state
```

Understanding gradient flow dynamics

```

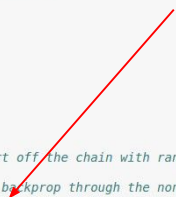
H = 5 # dimensionality of hidden state
T = 50 # number of time steps
Whh = np.random.randn(H,H)

# forward pass of an RNN (ignoring inputs x)
hs = {}
ss = {}
hs[-1] = np.random.randn(H)
for t in xrange(T):
    ss[t] = np.dot(Whh, hs[t-1])
    hs[t] = np.maximum(0, ss[t])

# backward pass of the RNN
dhs = {}
dss = {}
dhs[T-1] = np.random.randn(H) # start off the chain with random gradient
for t in reversed(xrange(T)):
    dss[t] = (hs[t] > 0) * dhs[t] # backprop through the nonlinearity
    dhs[t-1] = np.dot(Whh.T, dss[t]) # backprop into previous hidden state

```

if the largest eigenvalue is > 1 , gradient will explode
 if the largest eigenvalue is < 1 , gradient will vanish



[On the difficulty of training Recurrent Neural Networks, Pascanu et al., 2013]

Four effective ways to learn an RNN

- **Long Short Term Memory:**
Make the RNN out of little modules that are designed to remember values for a long time
- **Hessian Free Optimization:**
Deal with the vanishing gradients problem by using a fancy optimizer that can detect directions with a tiny gradient but even smaller curvature
 - The HF optimizer (Martens & Sutskever, 2011) is good at this
- **Echo State Networks:**
Initialize the input→ hidden and hidden→hidden and output→ hidden connections very carefully so that the hidden state has a huge reservoir of weakly coupled oscillators which can be selectively driven by the input
 - ESNs only need to learn the hidden→output connections
- **Good initialization with momentum:** Initialize like in Echo State Networks, but then learn all of the connections using momentum

Four effective ways to learn an RNN

- **Long Short Term Memory:**
Make the RNN out of little modules that are designed to remember values for a long time
- **Hessian Free Optimization:**
Deal with the vanishing gradients problem by using a fancy optimizer that can detect directions with a tiny gradient but even smaller curvature
 - The HF optimizer (Martens & Sutskever, 2011) is good at this
- **Echo State Networks:**
Initialize the input→ hidden and hidden→hidden and output→ hidden connections very carefully so that the hidden state has a huge reservoir of weakly coupled oscillators which can be selectively driven by the input
 - ESNs only need to learn the hidden→output connections
- **Good initialization with momentum:** Initialize like in Echo State Networks, but then learn all of the connections using momentum

Four effective ways to learn an RNN

- **Long Short Term Memory:**
Make the RNN out of little modules that are designed to remember values for a long time
- **Hessian Free Optimization:**
Deal with the vanishing gradients problem by using a fancy optimizer that can detect directions with a tiny gradient but even smaller curvature
 - The HF optimizer (Martens & Sutskever, 2011) is good at this
- **Echo State Networks:**
Initialize the input→ hidden and hidden→hidden and output→ hidden connections very carefully so that the hidden state has a huge reservoir of weakly coupled oscillators which can be selectively driven by the input
 - ESNs only need to learn the hidden→output connections
- **Good initialization with momentum:** Initialize like in Echo State Networks, but then learn all of the connections using momentum

Long Short Term Memory (LSTM)

- Hochreiter & Schmidhuber (1997) solved the problem of getting an RNN to remember things for a long time (like hundreds of time steps)
 - Keep short-term memory for a long period of time, thus the name
- They designed a memory cell using logistic and linear units with multiplicative interactions
- Information gets into the cell whenever its “write” gate is on
- The information stays in the cell so long as its “keep” gate is on
- Information can be read from the cell by turning on its “read” gate

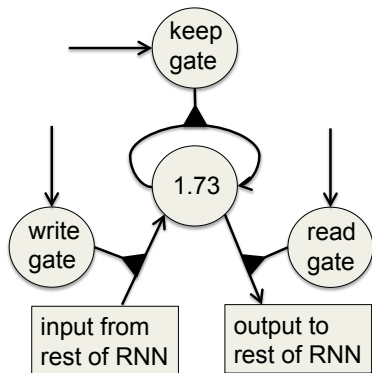
Long Short Term Memory (LSTM)

- Hochreiter & Schmidhuber (1997) solved the problem of getting an RNN to remember things for a long time (like hundreds of time steps)
 - Keep short-term memory for a long period of time, thus the name
- They designed a memory cell using logistic and linear units with multiplicative interactions
- Information gets into the cell whenever its “write” gate is on
- The information stays in the cell so long as its “keep” gate is on
- Information can be read from the cell by turning on its “read” gate

Long Short Term Memory (LSTM)

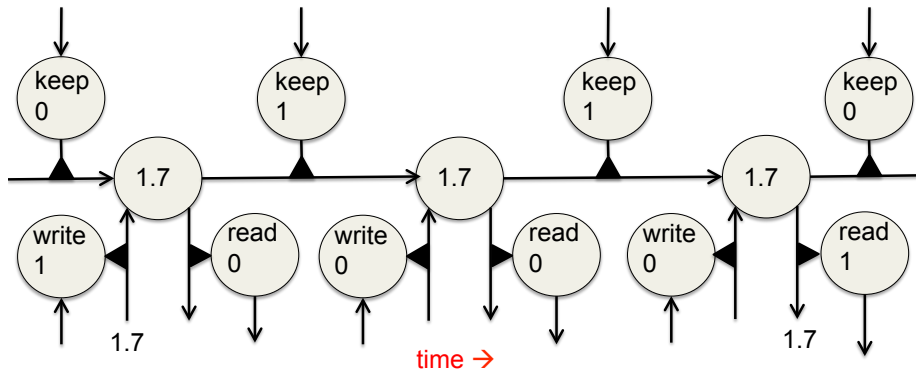
- Hochreiter & Schmidhuber (1997) solved the problem of getting an RNN to remember things for a long time (like hundreds of time steps)
 - Keep short-term memory for a long period of time, thus the name
- They designed a memory cell using logistic and linear units with multiplicative interactions
- Information gets into the cell whenever its “write” gate is on
- The information stays in the cell so long as its “keep” gate is on
- Information can be read from the cell by turning on its “read” gate

Implementing a memory cell in a neural network



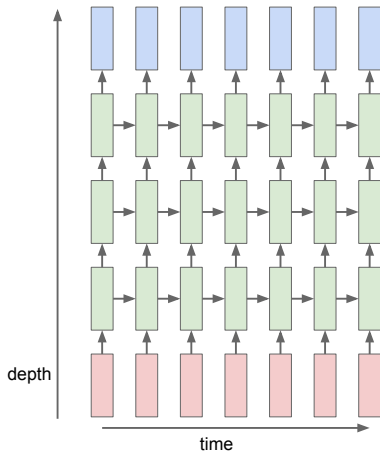
- To preserve information for a long time in the activities of an RNN, we use a circuit mimicking an analog memory cell
 - Information is kept in the cell when "keep" gate is on
 - Information is stored in the cell by activating its write gate
 - Information is retrieved by activating the read gate
 - We can backpropagate through this circuit because logistics are have nice derivatives

Backpropagation through a memory cell



RNN:

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

 $h \in \mathbb{R}^n, \quad W^l [n \times 2n]$


RNN:

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$h \in \mathbb{R}^n, \quad W^l [n \times 2n]$$

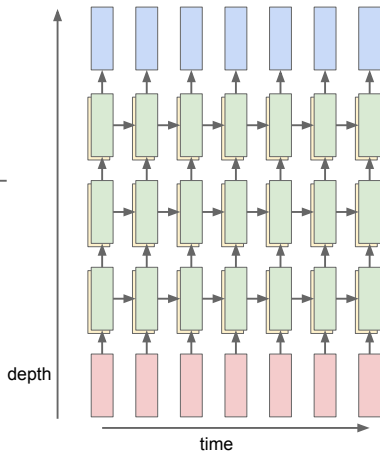
LSTM:

$$W^l [4n \times 2n]$$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

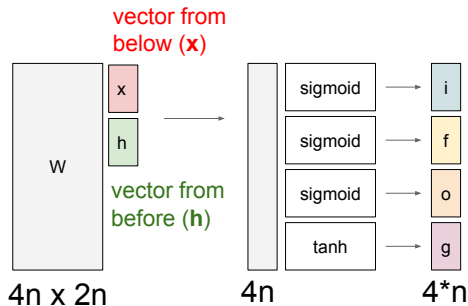
$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$



Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



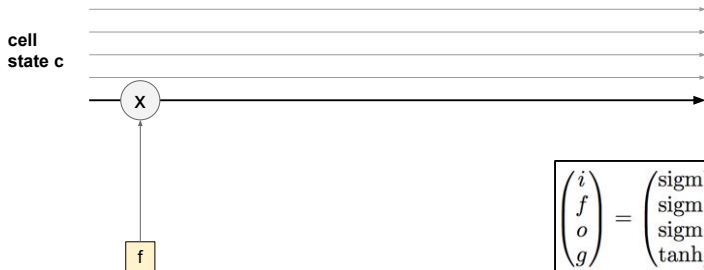
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



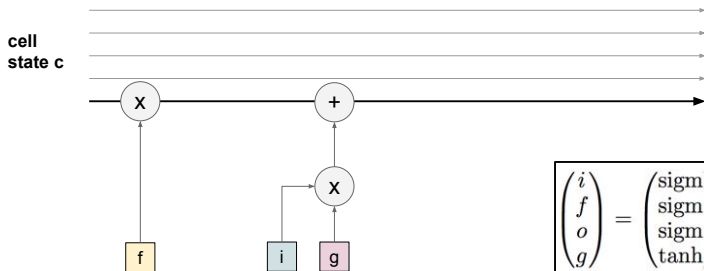
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_t^{l-1} \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



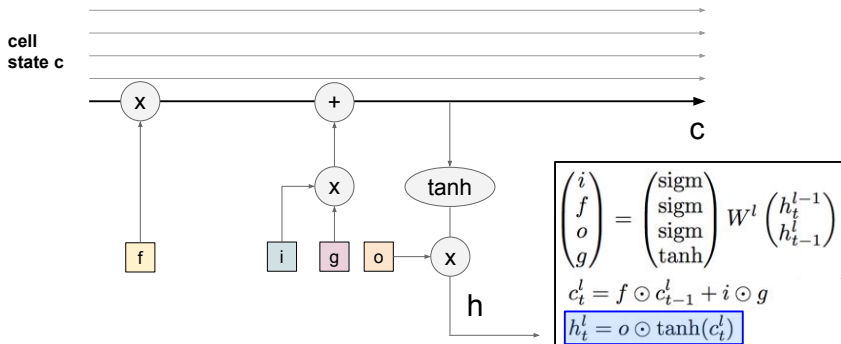
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_t^{l-1} \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



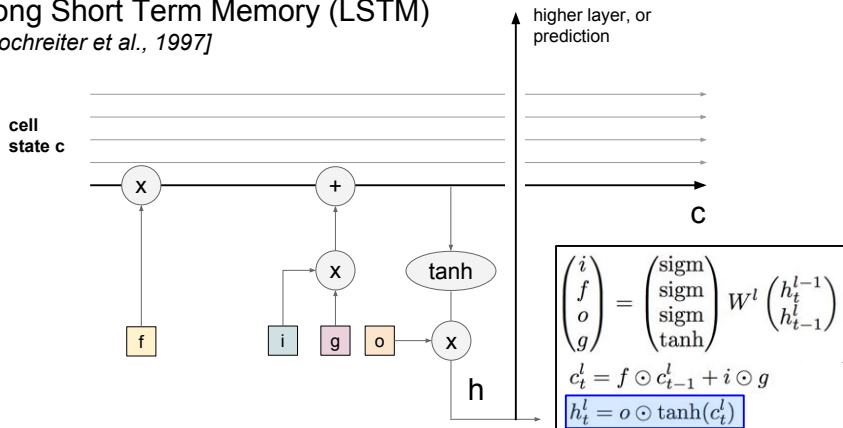
Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 10 - 73

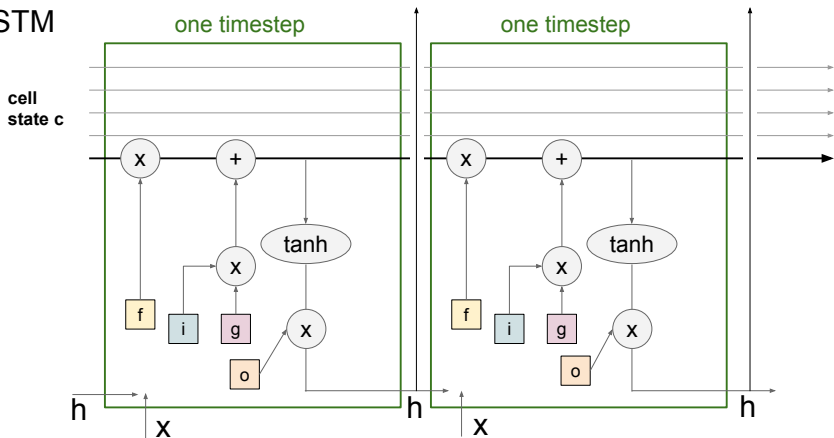
8 Feb 2016

Long Short Term Memory (LSTM)

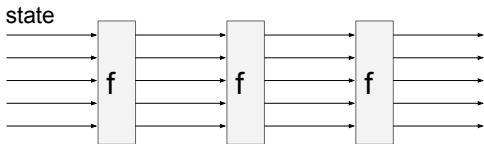
[Hochreiter et al., 1997]



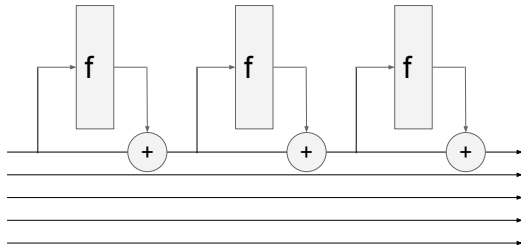
LSTM

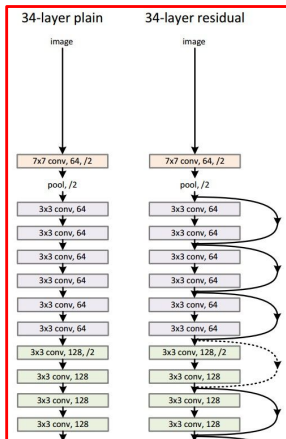


RNN



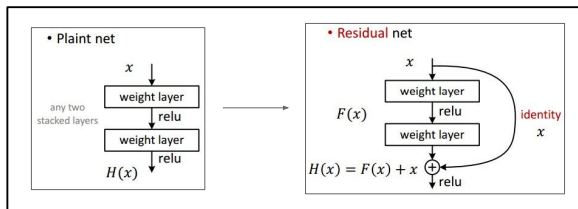
LSTM
(ignoring
forget gates)



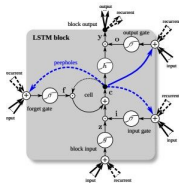


Recall: “PlainNets” vs. ResNets

ResNet is to PlainNet what LSTM is to RNN, kind of.



LSTM variants and friends



[LSTM: A Search Space Odyssey, Greff et al., 2015]

GRU [Learning phrase representations using rnn encoder-decoder for statistical machine translation, Cho et al. 2014]

$$\begin{aligned}
 r_t &= \text{sigm}(W_{xr}x_t + W_{hr}h_{t-1} + b_r) \\
 z_t &= \text{sigm}(W_{xz}x_t + W_{hz}h_{t-1} + b_z) \\
 \tilde{h}_t &= \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h) \\
 h_t &= z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t
 \end{aligned}$$

[An Empirical Exploration of Recurrent Network Architectures, Jozefowicz et al., 2015]

MUT1:

$$\begin{aligned}
 z &= \text{sigm}(W_{xz}x_t + b_z) \\
 r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\
 h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + \tanh(x_t) + b_h) \odot z \\
 &\quad + h_t \odot (1 - z)
 \end{aligned}$$

MUT2:

$$\begin{aligned}
 z &= \text{sigm}(W_{xz}x_t + W_{hz}h_t + b_z) \\
 r &= \text{sigm}(x_t + W_{hr}h_t + b_r) \\
 h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\
 &\quad + h_t \odot (1 - z)
 \end{aligned}$$

MUT3:

$$\begin{aligned}
 z &= \text{sigm}(W_{xz}x_t + W_{hz}\tanh(h_t) + b_z) \\
 r &= \text{sigm}(W_{xr}x_t + W_{hr}h_t + b_r) \\
 h_{t+1} &= \tanh(W_{hh}(r \odot h_t) + W_{xh}x_t + b_h) \odot z \\
 &\quad + h_t \odot (1 - z)
 \end{aligned}$$

Modelling text: why working with characters?

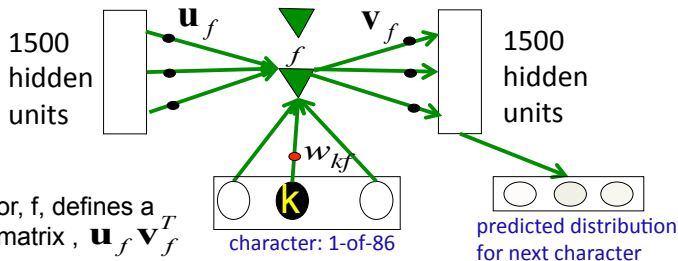
- The web is composed of character strings
- Any learning method powerful enough to understand the world by reading the web ought to find it trivial to learn which strings make words (this turns out to be true, as we shall see)
- Pre-processing text to get words is a big hassle
 - What about morphemes (prefixes, suffixes etc)
 - What about subtle effects like “sn” words?
 - What about New York vs new York Minster roof?
 - What about Finnish
 - ymmärtämättömyydellänsäkään

Modelling text: why working with characters?

- The web is composed of character strings
- Any learning method powerful enough to understand the world by reading the web ought to find it trivial to learn which strings make words (this turns out to be true, as we shall see)
- Pre-processing text to get words is a big hassle
 - What about morphemes (prefixes, suffixes etc)
 - What about subtle effects like “sn” words?
 - What about New York vs new York Minster roof?
 - What about Finnish
 - ymmärtämättömyydellänsäkään

Ideal model?

Using 3-way factors to allow a character to create a whole transition matrix



Each factor, f , defines a rank one matrix, $\mathbf{u}_f \mathbf{v}_f^T$

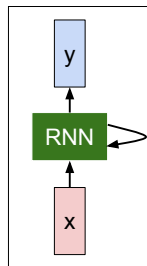
Each character, k , determines a gain w_{kf} for each of these matrices.

Simplest model: a first attempt

Character-level language model example

Vocabulary:
[h,e,l,o]

Example training sequence:
“hello”

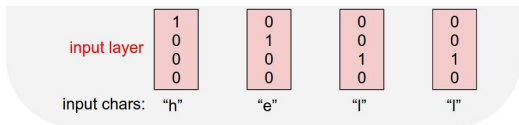


Simplest model: a first attempt

Character-level language model example

Vocabulary:
[h,e,l,o]

Example training sequence:
“hello”



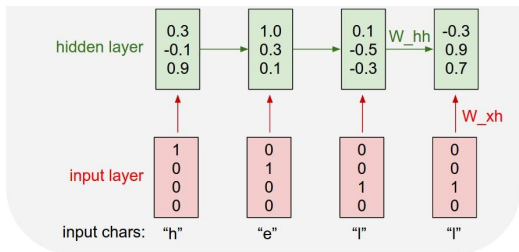
Simplest model: a first attempt

Character-level language model example

Vocabulary:
[h,e,l,o]

Example training sequence:
"hello"

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

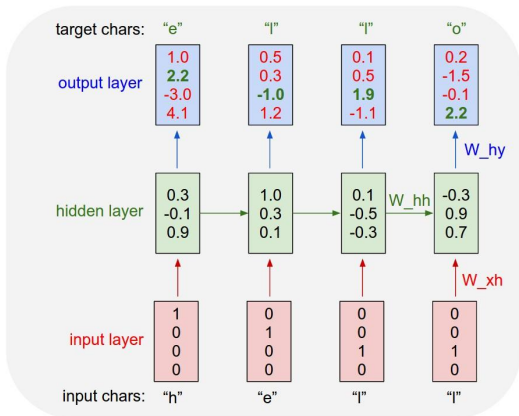


Simplest model: a first attempt

Character-level language model example

Vocabulary:
[h,e,l,o]

Example training sequence:
"hello"



Sampling

- Start the model with its default hidden state
- Give it a “burn-in” sequence of characters and let it update its hidden state after each character
- Then look at the probability distribution it predicts for the next character
- Pick a character randomly from that distribution and tell the net that this was the character that actually occurred
 - i.e. tell it that its guess was correct, whatever it guessed
- Continue to let it pick characters until bored

min-char-rnn.py gist: 112 lines of Python

```

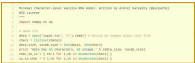
1 """
2 Minimal character-level vanilla RNN model. Written by Andrej Karpathy (2015/07/01)
3 BSD license
4 """
5
6 import numpy as np
7
8 # data loader
9
10 # data loader
11 # data loader
12 # data loader
13 # data loader
14 # data loader
15 # data loader
16 # data loader
17 # data loader
18 # data loader
19 # data loader
20 # data loader
21 # data loader
22 # data loader
23 # data loader
24 # data loader
25 # data loader
26 # data loader
27 # data loader
28 # data loader
29 # data loader
30 # data loader
31 # data loader
32 # data loader
33 # data loader
34 # data loader
35 # data loader
36 # data loader
37 # data loader
38 # data loader
39 # data loader
40 # data loader
41 # data loader
42 # data loader
43 # data loader
44 # data loader
45 # data loader
46 # data loader
47 # data loader
48 # data loader
49 # data loader
50 # data loader
51 # data loader
52 # data loader
53 # data loader
54 # data loader
55 # data loader
56 # data loader
57 # data loader
58 # data loader
59 # data loader
60 # data loader
61 # data loader
62 # data loader
63 # data loader
64 # data loader
65 # data loader
66 # data loader
67 # data loader
68 # data loader
69 # data loader
70 # data loader
71 # data loader
72 # data loader
73 # data loader
74 # data loader
75 # data loader
76 # data loader
77 # data loader
78 # data loader
79 # data loader
80 # data loader
81 # data loader
82 # data loader
83 # data loader
84 # data loader
85 # data loader
86 # data loader
87 # data loader
88 # data loader
89 # data loader
90 # data loader
91 # data loader
92 # data loader
93 # data loader
94 # data loader
95 # data loader
96 # data loader
97 # data loader
98 # data loader
99 # data loader
100 # data loader
101 # data loader
102 # data loader
103 # data loader
104 # data loader
105 # data loader
106 # data loader
107 # data loader
108 # data loader
109 # data loader
110 # data loader
111 # data loader
112 # data loader

```

<https://gist.github.com/karpathy/d4dee566867f8291f086>

Fei-Fei Li & Andrej Karpathy & Justin Johnson Lecture 10 - 22 8 Feb 2016

min-char-rnn.py gist



```

1 # Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
2 BSD License
3 """
4 import numpy as np
5
6 # data I/O
7 data = open('input.txt', 'r').read() # should be simple plain text file
8 chars = list(set(data))
9 data_size, vocab_size = len(data), len(chars)
10 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
11 char_to_ix = { ch:i for i,ch in enumerate(chars) }
12 ix_to_char = { i:ch for i,ch in enumerate(chars) }

```

Data I/O

```

1 """
2 Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
3 BSD License
4 """
5 import numpy as np
6
7 # data I/O
8 data = open('input.txt', 'r').read() # should be simple plain text file
9 chars = list(set(data))
10 data_size, vocab_size = len(data), len(chars)
11 print 'data has %d characters, %d unique.' % (data_size, vocab_size)
12 char_to_ix = { ch:i for i,ch in enumerate(chars) }
13 ix_to_char = { i:ch for i,ch in enumerate(chars) }

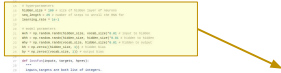
```

min-char-rnn.py gist

```

1 #!/usr/bin/env python3
2 # coding: utf-8
3 # min-char-rnn.py
4 # simple character-level language model
5 # based on the original code by Andre Kardecy (2015)
6 # see: https://gist.github.com/andrekardecy/564609
7 # see: https://github.com/andrekardecy/min-char-rnn
8
9 # hyperparameters
10 hidden_size = 100 # size of hidden layer of neurons
11 seq_length = 25 # number of steps to roll the RNN for
12 learning_rate = 1e-1
13
14 # model parameters
15 Wxh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
16 Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
17 Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
18 bh = np.zeros((hidden_size, 1)) # hidden bias
19 by = np.zeros((vocab_size, 1)) # output bias
20
21 # model
22 def init():
23     # initialize weights and biases
24     Wxh = np.random.randn(hidden_size, vocab_size)*0.01
25     Whh = np.random.randn(hidden_size, hidden_size)*0.01
26     Why = np.random.randn(vocab_size, hidden_size)*0.01
27     bh = np.zeros((hidden_size, 1))
28     by = np.zeros((vocab_size, 1))
29
30     # initialize hidden state
31     hidden = np.zeros((hidden_size, 1))
32
33     return Wxh, Whh, Why, bh, by, hidden
34
35 def forward(Wxh, Whh, Why, bh, by, hidden, x):
36     # calculate hidden state
37     hidden = np.tanh(Wxh.dot(x) + Whh.dot(hidden) + bh)
38
39     # calculate output
40     y = Why.dot(hidden) + by
41
42     return hidden, y
43
44 def backward(Wxh, Whh, Why, bh, by, hidden, y):
45     # calculate gradients
46     # ... (omitted for brevity)
47
48     return Wxh, Whh, Why, bh, by, hidden
49
50 def main():
51     # initialize model
52     Wxh, Whh, Why, bh, by, hidden = init()
53
54     # train model
55     # ... (omitted for brevity)
56
57 if __name__ == '__main__':
58     main()

```



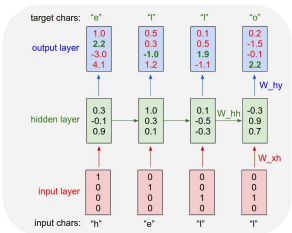
Initializations

```

15 # hyperparameters
16 hidden_size = 100 # size of hidden layer of neurons
17 seq_length = 25 # number of steps to unroll the RNN for
18 learning_rate = 1e-1
19
20 # model parameters
21 Wxh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
22 Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
23 Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
24 bh = np.zeros((hidden_size, 1)) # hidden bias
25 by = np.zeros((vocab_size, 1)) # output bias

```

recall:



min-char-rnn.py gist

```

1 from random import choice, randint, random, seed
2 import sys
3
4 # Model size
5 n_vocab = 27
6 n_hidden = 128
7 n_layers = 2
8
9 # Data
10 minchars = 'abcdefghijklmnopqrstuvwxyz '
11 maxchars = minchars * 100
12 data = minchars * 100000
13
14 # Parameters
15 seq_length = 200
16 batch_size = 128
17
18 # Training
19 # ... (omitted)
20
21 # Inference
22 # ... (omitted)
23
24 # Main loop
25 n, p = 0, 0
26 mw, mwh, mby = np.zeros_like(wxh), np.zeros_like(whh), np.zeros_like(why)
27 mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
28 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
29 while True:
30     # prepare inputs (we're sweeping from left to right in steps seq_length long)
31     if p+seq_length+1 >= len(data) or n == 0:
32         hprev = np.zeros((hidden_size,1)) # reset RNN memory
33         p = 0 # go from start of data
34     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
35     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
36
37     # sample from the model now and then
38     if n % 100 == 0:
39         sample_ix = sample(hprev, inputs[0], 200)
40         txt = ''.join(ix_to_char[ix] for ix in sample_ix)
41         print '----\n%s \n----' % (txt, )
42
43     # forward seq_length characters through the net and fetch gradient
44     loss, dwxh, dwhh, dwhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
45     smooth_loss = smooth_loss * 0.999 + loss * 0.001
46     if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
47
48     # perform parameter update with Adagrad
49     for param, dparam, mem in zip([wxh, whh, why, bh, by],
50                                 [dwxh, dwhh, dwhy, dbh, dby],
51                                 [mw, mwh, mwhy, mbh, mby]):
52         mem += dparam * dparam
53         param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
54
55     p += seq_length # move data pointer
56     n += 1 # iteration counter
57
58 # ... (omitted)

```

Main loop

```

80 n, p = 0, 0
81 mw, mwh, mwhy = np.zeros_like(wxh), np.zeros_like(whh), np.zeros_like(why)
82 mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
83 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
84 while True:
85     # prepare inputs (we're sweeping from left to right in steps seq_length long)
86     if p+seq_length+1 >= len(data) or n == 0:
87         hprev = np.zeros((hidden_size,1)) # reset RNN memory
88         p = 0 # go from start of data
89     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
90     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
91
92     # sample from the model now and then
93     if n % 100 == 0:
94         sample_ix = sample(hprev, inputs[0], 200)
95         txt = ''.join(ix_to_char[ix] for ix in sample_ix)
96         print '----\n%s \n----' % (txt, )
97
98     # forward seq_length characters through the net and fetch gradient
99     loss, dwxh, dwhh, dwhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
100     smooth_loss = smooth_loss * 0.999 + loss * 0.001
101     if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
102
103     # perform parameter update with Adagrad
104     for param, dparam, mem in zip([wxh, whh, why, bh, by],
105                                 [dwxh, dwhh, dwhy, dbh, dby],
106                                 [mw, mwh, mwhy, mbh, mby]):
107         mem += dparam * dparam
108         param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
109
110     p += seq_length # move data pointer
111     n += 1 # iteration counter

```

```

112 # ... (omitted)
113
114 # Inference
115 # ... (omitted)
116
117 # Main loop
118 n, p = 0, 0
119 mw, mwh, mwhy = np.zeros_like(wxh), np.zeros_like(whh), np.zeros_like(why)
120 mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
121 smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
122 while True:
123     # prepare inputs (we're sweeping from left to right in steps seq_length long)
124     if p+seq_length+1 >= len(data) or n == 0:
125         hprev = np.zeros((hidden_size,1)) # reset RNN memory
126         p = 0 # go from start of data
127     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
128     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
129
130     # sample from the model now and then
131     if n % 100 == 0:
132         sample_ix = sample(hprev, inputs[0], 200)
133         txt = ''.join(ix_to_char[ix] for ix in sample_ix)
134         print '----\n%s \n----' % (txt, )
135
136     # forward seq_length characters through the net and fetch gradient
137     loss, dwxh, dwhh, dwhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
138     smooth_loss = smooth_loss * 0.999 + loss * 0.001
139     if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
140
141     # perform parameter update with Adagrad
142     for param, dparam, mem in zip([wxh, whh, why, bh, by],
143                                 [dwxh, dwhh, dwhy, dbh, dby],
144                                 [mw, mwh, mwhy, mbh, mby]):
145         mem += dparam * dparam
146         param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
147
148     p += seq_length # move data pointer
149     n += 1 # iteration counter

```


min-char-rnn.py gist

```

1 from __future__ import division, unicode_literals, absolute_import, print_function
2 import sys
3 import random
4
5 # Model size
6 HIDDEN_DIM = 100  # hidden state dimensionality (number of hidden units)
7 EMBED_DIM = 26  # embedding dimensionality (number of characters)
8 INPUT_DIM = HIDDEN_DIM + EMBED_DIM  # input dimensionality
9
10 # Data loader
11 def load_data(filename):
12     """Load data from a file. Returns a list of (input, target) pairs.
13     Each pair is a list of integers representing characters. The first
14     element is the input sequence, and the second is the target sequence.
15     The length of the input sequence is one less than the length of the
16     target sequence, as the target sequence starts with the first character
17     of the input sequence.
18     """
19     with open(filename, 'r') as f:
20         lines = f.readlines()
21     data = []
22     for line in lines:
23         line = line.strip()
24         if not line:
25             continue
26         input_seq = []
27         target_seq = []
28         for i in range(len(line)):
29             input_seq.append(ord(line[i]) - ord('a'))
30             target_seq.append(ord(line[i+1]) - ord('a'))
31         data.append((input_seq, target_seq))
32     return data
33
34 # Utility functions
35 def softmax(x):
36     """Compute softmax probabilities from a list of raw scores.
37     """
38     exp_x = np.exp(x)
39     return exp_x / np.sum(exp_x)
40
41 def clip_by_value(arr, clip_min, clip_max):
42     """Apply a clip by value operation to an array.
43     """
44     arr[arr < clip_min] = clip_min
45     arr[arr > clip_max] = clip_max
46     return arr
47
48 def clip_by_norm(arr, clip_norm):
49     """Apply a clip by norm operation to an array.
50     """
51     norm = np.linalg.norm(arr)
52     if norm > clip_norm:
53         scale = clip_norm / norm
54         arr *= scale
55     return arr
56
57 def clip_by_value_and_norm(arr, clip_min, clip_max, clip_norm):
58     """Apply clip by value and clip by norm operations to an array.
59     """
60     arr = clip_by_value(arr, clip_min, clip_max)
61     arr = clip_by_norm(arr, clip_norm)
62     return arr
63
64 # Model
65 def init_model():
66     """Initialize the model parameters.
67     """
68     # Initialize hidden state
69     hidden_state = np.zeros((1, HIDDEN_DIM))
70     # Initialize embedding matrix
71     embed_matrix = np.random.randn(EMBED_DIM, EMBED_DIM)
72     # Initialize input-to-hidden weights
73     W_in = np.random.randn(HIDDEN_DIM, INPUT_DIM)
74     # Initialize hidden-to-hidden weights
75     W_hid = np.random.randn(HIDDEN_DIM, HIDDEN_DIM)
76     # Initialize hidden-to-output weights
77     W_out = np.random.randn(EMBED_DIM, HIDDEN_DIM)
78     # Initialize bias vectors
79     b_in = np.zeros(HIDDEN_DIM)
80     b_hid = np.zeros(HIDDEN_DIM)
81     b_out = np.zeros(EMBED_DIM)
82     # Initialize learning rate
83     lr = 0.01
84     # Initialize clip parameters
85     clip_min = -1.0
86     clip_max = 1.0
87     clip_norm = 5.0
88     return hidden_state, embed_matrix, W_in, W_hid, W_out, b_in, b_hid, b_out, lr, clip_min, clip_max, clip_norm
89
90 def forward_pass(inputs, targets, hidden_state, embed_matrix, W_in, W_hid, W_out, b_in, b_hid, b_out):
91     """Perform a forward pass through the model.
92     """
93     hidden_states = []
94     for t in range(len(inputs)):
95         input_vec = np.zeros(INPUT_DIM)
96         input_vec[:EMBED_DIM] = embed_matrix[inputs[t]]
97         input_vec[EMBED_DIM:] = hidden_state[t]
98         hidden_state = np.tanh(np.dot(W_in, input_vec) + b_in)
99         hidden_state = np.tanh(np.dot(W_hid, hidden_state) + b_hid)
100         output_vec = np.dot(W_out, hidden_state) + b_out
101         output_prob = softmax(output_vec)
102         hidden_states.append(hidden_state)
103     return hidden_states, output_prob
104
105 def backward_pass(hidden_states, output_prob, targets, lr, clip_min, clip_max, clip_norm):
106     """Perform a backward pass through the model to compute gradients.
107     """
108     # Compute gradients of output probabilities
109     d_output_prob = output_prob - targets
110     # Compute gradients of hidden states
111     d_hidden_state = np.zeros((len(hidden_states), HIDDEN_DIM))
112     for t in range(len(hidden_states)-1, -1, -1):
113         d_hidden_state[t] = d_output_prob[t] * W_out.T
114         d_hidden_state[t] = clip_by_value_and_norm(d_hidden_state[t], clip_min, clip_max, clip_norm)
115         d_hidden_state[t-1] = np.dot(W_hid.T, d_hidden_state[t]) + lr * d_hidden_state[t-1]
116         d_hidden_state[t-1] = clip_by_value_and_norm(d_hidden_state[t-1], clip_min, clip_max, clip_norm)
117     return d_hidden_state
118
119 def update_model(hidden_state, embed_matrix, W_in, W_hid, W_out, b_in, b_hid, b_out, lr, clip_min, clip_max, clip_norm):
120     """Update the model parameters based on the current batch.
121     """
122     # Compute gradients of model parameters
123     d_W_in = np.zeros(W_in.shape)
124     d_W_hid = np.zeros(W_hid.shape)
125     d_W_out = np.zeros(W_out.shape)
126     d_b_in = np.zeros(b_in.shape)
127     d_b_hid = np.zeros(b_hid.shape)
128     d_b_out = np.zeros(b_out.shape)
129     # Update model parameters
130     W_in = W_in + lr * d_W_in
131     W_hid = W_hid + lr * d_W_hid
132     W_out = W_out + lr * d_W_out
133     b_in = b_in + lr * d_b_in
134     b_hid = b_hid + lr * d_b_hid
135     b_out = b_out + lr * d_b_out
136     # Clip model parameters
137     W_in = clip_by_value_and_norm(W_in, clip_min, clip_max, clip_norm)
138     W_hid = clip_by_value_and_norm(W_hid, clip_min, clip_max, clip_norm)
139     W_out = clip_by_value_and_norm(W_out, clip_min, clip_max, clip_norm)
140     b_in = clip_by_value(b_in, clip_min, clip_max)
141     b_hid = clip_by_value(b_hid, clip_min, clip_max)
142     b_out = clip_by_value(b_out, clip_min, clip_max)
143     return W_in, W_hid, W_out, b_in, b_hid, b_out
144
145 def train_model(filename, num_epochs, num_batches_per_epoch):
146     """Train the model on a dataset.
147     """
148     # Load data
149     data = load_data(filename)
150     # Initialize model
151     hidden_state, embed_matrix, W_in, W_hid, W_out, b_in, b_hid, b_out, lr, clip_min, clip_max, clip_norm = init_model()
152     # Train model
153     for epoch in range(num_epochs):
154         for batch in range(num_batches_per_epoch):
155             # Sample a batch of data
156             input_seqs, target_seqs = [], []
157             for i in range(num_batches_per_epoch):
158                 input_seq, target_seq = data[random.randrange(len(data))]
159                 input_seqs.append(input_seq)
160                 target_seqs.append(target_seq)
161             # Perform forward and backward passes
162             hidden_states, output_prob = forward_pass(input_seqs, target_seqs, hidden_state, embed_matrix, W_in, W_hid, W_out, b_in, b_hid, b_out)
163             d_hidden_state = backward_pass(hidden_states, output_prob, target_seqs, lr, clip_min, clip_max, clip_norm)
164             # Update model parameters
165             W_in, W_hid, W_out, b_in, b_hid, b_out = update_model(hidden_state, embed_matrix, W_in, W_hid, W_out, b_in, b_hid, b_out, lr, clip_min, clip_max, clip_norm)
166     return hidden_state, embed_matrix, W_in, W_hid, W_out, b_in, b_hid, b_out
167
168 if __name__ == '__main__':
169     filename = 'data.txt'
170     num_epochs = 10
171     num_batches_per_epoch = 100
172     hidden_state, embed_matrix, W_in, W_hid, W_out, b_in, b_hid, b_out = train_model(filename, num_epochs, num_batches_per_epoch)
173     print('Training complete. Model parameters are stored in the global variables.')
174

```



Loss function

- forward pass (compute loss)
- backward pass (compute param gradient)

```

27 def lossFun(inputs, targets, hprev):
28     """
29     inputs, targets are both list of integers.
30     hprev is Hx1 array of initial hidden state.
31     returns the loss, gradients on model parameters, and last hidden state
32     """
33     xs, hs, ys, ps = [], [], [], {}
34     hs[-1] = np.copy(hprev)
35     loss = 0
36     # forward pass
37     for t in xrange(len(inputs)):
38         xs[t] = np.zeros((vocab_size, 1)) # encode in 1-of-k representation
39         xs[t][inputs[t]] = 1
40         hs[t] = np.tanh(np.dot(Wx, xs[t]) + np.dot(wh, hs[t-1]) + bh) # hidden state
41         ys[t] = np.dot(Wy, hs[t]) + by # unnormalized log probabilities for next chars
42         ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t])) # probabilities for next chars
43         loss += -np.log(ps[t][targets[t], 0]) # softmax (cross-entropy loss)
44     # backward pass: compute gradients going backwards
45     dwh, dwhh, dwhy = np.zeros_like(Wx), np.zeros_like(wh), np.zeros_like(Wy)
46     dbh, dby = np.zeros_like(bh), np.zeros_like(by)
47     dhnext = np.zeros_like(hs[0])
48     for t in reversed(xrange(len(inputs))):
49         dy = np.copy(ps[t])
50         dy[targets[t]] -= 1 # backprop into y
51         dwhy += np.dot(dy, hs[t].T)
52         dby = dy
53         dh = np.dot(Wy.T, dy) + dhnext # backprop into h
54         ddraw = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity
55         dbh = ddraw
56         dwh += np.dot(ddraw, xs[t].T)
57         dwhh += np.dot(ddraw, hs[t-1].T)
58         dhnext = np.dot(Wx.T, ddraw)
59     for dparam in [dwh, dwhh, dwhy, dbh, dby]:
60         np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61     return loss, dwh, dwhh, dwhy, dbh, dby, hs[len(inputs)-1]

```

min-char-rnn.py gist

```

1 #!/usr/bin/env python
2 """
3 Example: character-level language model, trained by simple backprop (BPTRT)
4 """
5 import numpy as np
6
7 # Model parameters
8 # Embedding matrix
9 Wxh = np.random.randn(EMBED_DIM, VOCAB_SIZE)
10 # Hidden-to-hidden weights
11 Whh = np.random.randn(HIDDEN_DIM, HIDDEN_DIM)
12 # Hidden-to-output weights
13 Why = np.random.randn(VOCAB_SIZE, HIDDEN_DIM)
14 # Bias for the hidden state
15 bh = np.random.randn(HIDDEN_DIM)
16 # Bias for the output
17 by = np.random.randn(VOCAB_SIZE)
18 # Learning rate
19 LR = 0.01
20 # Number of hidden units
21 HIDDEN_DIM = 100
22 # Number of embedding units
23 EMBED_DIM = 100
24 # Vocabulary size
25 VOCAB_SIZE = 26
26 # Number of training steps
27 STEPS = 100000
28 # Number of hidden states to store
29 NUM_HIDDEN_STATES_TO_STORE = 1000
30 # Initial hidden state
31 hprev = np.zeros(HIDDEN_DIM)
32 # Loss function
33 loss = 0
34 # Forward pass
35 for t in xrange(len(inputs)):
36     # Encode in 1-of-k representation
37     xs[t] = np.zeros((VOCAB_SIZE, 1))
38     xs[t][inputs[t]] = 1
39     # Hidden state
40     hs[t] = np.tanh(np.dot(Wxh, xs[t]) + np.dot(Whh, hs[t-1]) + bh)
41     # Unnormalized log probabilities for next chars
42     ys[t] = np.dot(Why, hs[t]) + by
43     # Probabilities for next chars
44     ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t]))
45     # Softmax (cross-entropy loss)
46     loss += -np.log(ps[t][targets[t], 0])
47
48 # Backward pass
49 # Gradients on model parameters
50 dWxh = np.zeros_like(Wxh)
51 dWhh = np.zeros_like(Whh)
52 dWhy = np.zeros_like(Why)
53 dbh = np.zeros_like(bh)
54 dby = np.zeros_like(by)
55 # Gradients on inputs
56 dxs = np.zeros_like(xs)
57 # Gradients on hidden states
58 dhs = np.zeros_like(hs)
59 # Gradients on outputs
60 dps = np.zeros_like(ps)
61 # Gradients on targets
62 dtargets = np.zeros_like(targets)
63 # Gradients on loss
64 dloss = 0
65 # Backward pass
66 for t in xrange(len(inputs)-1, -1, -1):
67     # Gradients on outputs
68     dps[t] = -1 / ps[t][targets[t], 0]
69     # Gradients on hidden states
70     dhs[t] = dWhy * ps[t] + dps[t]
71     # Gradients on inputs
72     dxs[t] = dWxh * ps[t] + dhs[t] * Wxh
73     # Gradients on model parameters
74     dWxh += dxs * xs[t]
75     dWhh += dhs[t] * hs[t-1]
76     dWhy += dps[t] * hs[t]
77     dbh += dhs[t]
78     dby += dps[t]
79     # Gradients on loss
80     dloss += dps[t][targets[t], 0]
81
82 # Print the results
83 print "Loss: %f" % loss
84 print "Gradients on model parameters: %f" % (dWxh + dWhh + dWhy + dbh + dby)
85 print "Gradients on inputs: %f" % dxs
86 print "Gradients on hidden states: %f" % dhs
87 print "Gradients on outputs: %f" % dps
88 print "Gradients on targets: %f" % dtargets
89 print "Gradients on loss: %f" % dloss
90
91 # Save the model parameters
92 np.save("min-char-rnn-params.npy", (Wxh, Whh, Why, bh, by))
93
94 # Load the model parameters
95 (Wxh, Whh, Why, bh, by) = np.load("min-char-rnn-params.npy")
96
97 # Test the model
98 # Inputs
99 inputs = "hello world"
100 # Targets
101 targets = "ello wrld"
102 # Initial hidden state
103 hprev = np.zeros(HIDDEN_DIM)
104 # Loss function
105 loss = 0
106 # Forward pass
107 for t in xrange(len(inputs)):
108     # Encode in 1-of-k representation
109     xs[t] = np.zeros((VOCAB_SIZE, 1))
110     xs[t][inputs[t]] = 1
111     # Hidden state
112     hs[t] = np.tanh(np.dot(Wxh, xs[t]) + np.dot(Whh, hs[t-1]) + bh)
113     # Unnormalized log probabilities for next chars
114     ys[t] = np.dot(Why, hs[t]) + by
115     # Probabilities for next chars
116     ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t]))
117     # Softmax (cross-entropy loss)
118     loss += -np.log(ps[t][targets[t], 0])
119
120 # Print the results
121 print "Loss: %f" % loss
122 print "Gradients on model parameters: %f" % (Wxh + Whh + Why + bh + by)
123 print "Gradients on inputs: %f" % xs
124 print "Gradients on hidden states: %f" % hs
125 print "Gradients on outputs: %f" % ps
126 print "Gradients on targets: %f" % targets
127 print "Gradients on loss: %f" % loss

```



```

27 def lossFun(inputs, targets, hprev):
28     """
29     inputs, targets are both list of integers.
30     hprev is Hx1 array of initial hidden state
31     returns the loss, gradients on model parameters, and last hidden state
32     """
33     xs, hs, ys, ps = {}, {}, {}, {}
34     hs[-1] = np.copy(hprev)
35     loss = 0
36     # forward pass
37     for t in xrange(len(inputs)):
38         xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
39         xs[t][inputs[t]] = 1
40         hs[t] = np.tanh(np.dot(wxh, xs[t]) + np.dot(whh, hs[t-1]) + bh) # hidden state
41         ys[t] = np.dot(why, hs[t]) + by # unnormalized log probabilities for next chars
42         ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t])) # probabilities for next chars
43         loss += -np.log(ps[t][targets[t],0]) # softmax (cross-entropy loss)

```

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

$$y_t = W_{hy}h_t$$

Softmax classifier

min-char-rnn.py gist

```

44 # backward pass: compute gradients going backwards
45 dwdx, dwhh, dwhy = np.zeros_like(wxh), np.zeros_like(whh), np.zeros_like(why)
46 dbh, dby = np.zeros_like(bh), np.zeros_like(by)
47 dhnext = np.zeros_like(hs[0])
48 for t in reversed(xrange(len(inputs))):
49     dy = np.copy(ps[t])
50     dy[targets[t]] -= 1 # backprop into y
51     dwhy += np.dot(dy, hs[t].T)
52     dby += dy
53     dh = np.dot(why.T, dy) + dhnext # backprop into h
54     dhraw = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity
55     dbh += dhraw
56     dwdx += np.dot(dhraw, xs[t].T)
57     dwhh += np.dot(dhraw, hs[t-1].T)
58     dhnext = np.dot(whh.T, dhraw)
59 for dparam in [dwdx, dwhh, dwhy, dbh, dby]:
60     np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61 return loss, dwdx, dwhh, dwhy, dbh, dby, hs[len(inputs)-1]

```

```

def _compute_grads_backward(self):
    """Compute gradients backwards"""
    # Compute gradients for each layer
    # ... (omitted code) ...
    # Compute gradients for the input layer
    # ... (omitted code) ...
    # Compute gradients for the hidden layer
    # ... (omitted code) ...
    # Compute gradients for the output layer
    # ... (omitted code) ...
    # Clip gradients to mitigate exploding gradients
    for dparam in [dwdx, dwhh, dwhy, dbh, dby]:
        np.clip(dparam, -5, 5, out=dparam)
    return loss, dwdx, dwhh, dwhy, dbh, dby, hs[len(inputs)-1]

```

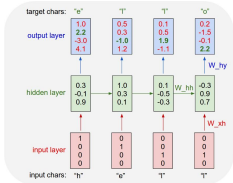


```

44 # backward pass: compute gradients going backwards
45 dwdx, dwhh, dwhy = np.zeros_like(wxh), np.zeros_like(whh), np.zeros_like(why)
46 dbh, dby = np.zeros_like(bh), np.zeros_like(by)
47 dhnext = np.zeros_like(hs[0])
48 for t in reversed(xrange(len(inputs))):
49     dy = np.copy(ps[t])
50     dy[targets[t]] -= 1 # backprop into y
51     dwhy += np.dot(dy, hs[t].T)
52     dby += dy
53     dh = np.dot(why.T, dy) + dhnext # backprop into h
54     dhraw = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity
55     dbh += dhraw
56     dwdx += np.dot(dhraw, xs[t].T)
57     dwhh += np.dot(dhraw, hs[t-1].T)
58     dhnext = np.dot(whh.T, dhraw)
59 for dparam in [dwdx, dwhh, dwhy, dbh, dby]:
60     np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61 return loss, dwdx, dwhh, dwhy, dbh, dby, hs[len(inputs)-1]

```

recall:



min-char-rnn.py gist

```

1 #
2 # LICENSE: MIT (https://creativecommons.org/licenses/by-sa/4.0/)
3 #
4 #
5 #
6 #
7 #
8 #
9 #
10 #
11 #
12 #
13 #
14 #
15 #
16 #
17 #
18 #
19 #
20 #
21 #
22 #
23 #
24 #
25 #
26 #
27 #
28 #
29 #
30 #
31 #
32 #
33 #
34 #
35 #
36 #
37 #
38 #
39 #
40 #
41 #
42 #
43 #
44 #
45 #
46 #
47 #
48 #
49 #
50 #
51 #
52 #
53 #
54 #
55 #
56 #
57 #
58 #
59 #
60 #
61 #
62 #
63 #
64 #
65 #
66 #
67 #
68 #
69 #
70 #
71 #
72 #
73 #
74 #
75 #
76 #
77 #
78 #
79 #
80 #
81 #
82 #
83 #
84 #
85 #
86 #
87 #
88 #
89 #
90 #
91 #
92 #
93 #
94 #
95 #
96 #
97 #
98 #
99 #
100 #

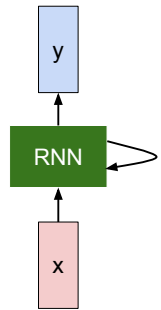
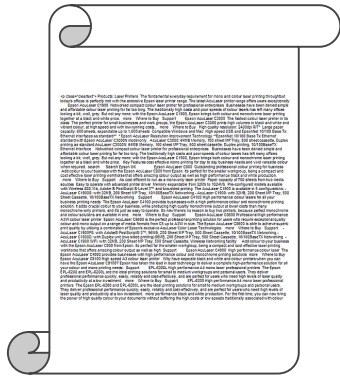
```



```

63 def sample(h, seed_ix, n):
64     """
65     sample a sequence of integers from the model
66     h is memory state, seed_ix is seed letter for first time step
67     """
68     x = np.zeros((vocab_size, 1))
69     x[seed_ix] = 1
70     ixes = []
71     for t in xrange(n):
72         h = np.tanh(np.dot(Wxh, x) + np.dot(Whh, h) + bh)
73         y = np.dot(Wyh, h) + by
74         p = np.exp(y) / np.sum(np.exp(y))
75         ix = np.random.choice(range(vocab_size), p=p.ravel())
76         x = np.zeros((vocab_size, 1))
77         x[ix] = 1
78         ixes.append(ix)
79     return ixes

```



Fei-Fei Li & Andrej Karpathy & Justin Johnson Lecture 10 - 34 8 Feb 2016

Sonnet 116 – Let me not ...

by William Shakespeare

Let me not to the marriage of true minds
Admit impediments. Love is not love
Which alters when it alteration finds,
Or bends with the remover to remove:
O no! it is an ever-fixed mark
That looks on tempests and is never shaken;
It is the star to every wandering bark,
Whose worth's unknown, although his height be taken.
Love's not Time's fool, though rosy lips and cheeks
Within his bending sickle's compass come:
Love alters not with his brief hours and weeks,
But bears it out even to the edge of doom.
If this be error and upon me proved,
I never writ, nor no man ever loved.

at first:

tyntd-iafhatawiaoihrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
 plia tklrqd t o idoe ns,smtt h ne etie h,hregtrs nigtkie,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. MeiDimorotio 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 offer.

↓
 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.

PANDARUS:

Alas, I think he shall be come approached and the day
When little strain 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.

open source textbook on algebraic geometry

 **The Stacks Project**

[home](#) [about](#) [tags explained](#) [tag lookup](#) [browse](#) [search](#) [bibliography](#) [recent comments](#) [blog](#) [add slogans](#)

Browse chapters

Part	Chapter	online	TeX source	view pdf
Preliminaries				
	1. Introduction	online	tex 	pdf 
	2. Conventions	online	tex 	pdf 
	3. Set Theory	online	tex 	pdf 
	4. Categories	online	tex 	pdf 
	5. Topology	online	tex 	pdf 
	6. Sheaves on Spaces	online	tex 	pdf 
	7. Sites and Sheaves	online	tex 	pdf 
	8. Stacks	online	tex 	pdf 
	9. Fields	online	tex 	pdf 
	10. Commutative Algebra	online	tex 	pdf 

Parts

- [Preliminaries](#)
- [Schemes](#)
- [Topics in Scheme Theory](#)
- [Algebraic Spaces](#)
- [Topics in Geometry](#)
- [Deformation Theory](#)
- [Algebraic Stacks](#)
- [Miscellany](#)

Statistics

The Stacks project now consists of

- 455910 lines of code
- 14221 tags (56 inactive tags)
- 2366 sections

Latex source



For $\bigoplus_{n=1, \dots, m}$ where $\mathcal{L}_{m_*} = 0$, 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 $U \rightarrow T$ is a separated algebraic space.

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

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

and the comparably in the fibre product covering we have to prove the lemma generated by $\prod Z \times_U U \rightarrow V$. Consider the maps M along the set of points Sch_{fppf} and $U \rightarrow U$ 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 $Sch(G)$ such that $\text{Spec}(R')$ $\rightarrow S$ 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 $\mathcal{O}_{X,x'} \rightarrow \mathcal{O}_{X',x'}$ is separated. By Algebra, Lemma ?? we can define a map of complexes $GL_{S'}(Z'/S')$ and we win. \square

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}_i 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

$$\tilde{M}^* = \mathcal{I}^* \otimes_{\text{Spec}(k)} \mathcal{O}_{S_n} - i_{X'}^{-1} \mathcal{F}$$

is a unique morphism of algebraic stacks. Note that

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

and

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

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. \square

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 (2) by the construction in the description.
 Suppose $X = \lim |X|$ (by the formal open covering X and a single map $\underline{Proj}_X(A) = \text{Spec}(B)$ over U compatible with the complex

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

When in this case of to show that $\mathcal{Q} \rightarrow \mathcal{C}_{Z/X}$ 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 $S = \text{Spec}(R)$ and $Y = \text{Spec}(R)$.

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

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

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

Lemma 0.3. In Situation ?? . Hence we may assume $\mathfrak{q}' = 0$.

Proof. We will use the property we see that \mathfrak{p} is the next 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 . \square

Proof. Omitted. □

Lemma 0.1. Let \mathcal{C} be a set of the construction.
 Let \mathcal{C} be a gerber covering. Let \mathcal{F} be a quasi-coherent sheaves of \mathcal{O} -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

.

Proof. This is an algebraic space with the composition of sheaves \mathcal{F} on $X_{\text{étale}}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{\text{morph}_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where \mathcal{G} defines an isomorphism $\mathcal{F} \rightarrow \mathcal{F}$ of \mathcal{O} -modules. □

Lemma 0.2. This is an integer Z is injective.
Proof. See Spaces, Lemma ?? □

Lemma 0.3. Let S be a scheme. Let X be a scheme and X is an affine open covering. Let $U \subset X$ be a canonical and locally of finite type. Let X be a scheme. Let X be a scheme which is equal to the formal complex.
 The following to the construction of the lemma follows.
 Let X be a scheme. Let X be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

be a morphism of algebraic spaces over S and Y .

Proof. Let X be a nonzero scheme of X . Let X be an algebraic space. Let \mathcal{F} be a quasi-coherent sheaf of \mathcal{O}_X -modules. The following are equivalent

- (1) \mathcal{F} is an algebraic space over S .
- (2) If X is an affine open covering.

Consider a common structure on X and X the functor $\mathcal{O}_X(U)$ which is locally of finite type. □

This since $\mathcal{F} \in \mathcal{F}$ and $x \in \mathcal{G}$ the diagram

is a limit. Then \mathcal{G} is a finite type and assume S is a flat and \mathcal{F} and \mathcal{G} is a finite type \mathcal{F} . This is of finite type diagrams, and

- the composition of \mathcal{G} is a regular sequence,
- $\mathcal{O}_{X'}$ is a sheaf of rings.

Proof. We have see that $X = \text{Spec}(R)$ and \mathcal{F} is a finite type representable by algebraic space. The property \mathcal{F} is a finite morphism of algebraic stacks. Then the cohomology of X is an open neighbourhood of U . □

Proof. This is clear that \mathcal{G} is a finite presentation, see Lemmas ??.
 A reduced above we conclude that U is an open covering of \mathcal{C} . The functor \mathcal{F} is a "field"

$$\mathcal{O}_{X, x} \rightarrow \mathcal{F}_x \rightarrow \mathcal{O}_{X, x} \rightarrow \mathcal{O}_{X, x} \rightarrow \mathcal{O}_{X, x} \rightarrow \mathcal{O}_{X, x}$$

is an isomorphism of covering of $\mathcal{O}_{X, x}$. If \mathcal{F} is the unique element of \mathcal{F} such that X is an isomorphism.
 The property \mathcal{F} is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme \mathcal{O}_X -algebra with \mathcal{F} are opens of finite type over S .
 If \mathcal{F} is a scheme theoretic image points. □

If \mathcal{F} is a finite direct sum $\mathcal{O}_{X, x}$ is a closed immersion, see Lemma ?? . This is a sequence of \mathcal{F} is a similar morphism.

GitHub repository interface for `torvalds / linux`. The repository has 520,037 commits, 1 branch, 420 releases, and 5,039 contributors. The current branch is `master`.

Recent commits:

Directory	Commit Message	Time Ago
Documentation	Merge git://git.kernel.org/pub/scm/linux/kernel/git/hab/target-pending	6 days ago
arch	Merge branch 'x86-urgent-for-linus' of git://git.kernel.org/pub/scm/...	a day ago
block	block: discard bdi_unregister() in favour of bdi_destroy()	9 days ago
crypto	Merge git://git.kernel.org/pub/scm/linux/kernel/git/herbert/crypto-2.6	10 days ago
drivers	Merge branch 'drm-fixes' of git://people.freedesktop.org/~airlied/linux	9 hours ago
firmware	firmware/ihex2hw.c: restore missing default in switch statement	2 months ago
fs	vfs: read file_handle only once in handle_to_path	4 days ago
include	Merge branch 'perf-urgent-for-linus' of git://git.kernel.org/pub/scm/...	a day ago
init	init: fix regression by supporting devices with major:minor:offset fo...	a month ago

Right sidebar options: Code, Pull requests (74), Pulse, Graphs, HTTPS clone URL (https://github.com), Clone in Desktop, Download ZIP.

```

static void do_command(struct seq_file *m, void *v)
{
    int column = 32 << (cmd[2] & 0x80);
    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 << i))
            pipe = (in_use & UMXTHREAD_UNCCA) +
                ((count & 0x00000000ffffffff) & 0x0000000f) << 8;
        if (count == 0)
            sub(pid, ppc_md.kexec_handle, 0x20000000);
        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++)
        seq_puts(s, "policy ");
}

```

Generated C code

```
/*
 * Copyright (c) 2006-2010, Intel Mobile Communications. All rights reserved.
 *
 * This program is free software; you can redistribute it and/or modify it
 * under the terms of the GNU General Public License version 2 as published by
 * the Free Software Foundation.
 *
 * This program is distributed in the hope that it will be useful,
 * but WITHOUT ANY WARRANTY; without even the implied warranty of
 * MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
 * GNU General Public License for more details.
 *
 * You should have received a copy of the GNU General Public License
 * along with this program; if not, write to the Free Software Foundation,
 * Inc., 675 Mass Ave, Cambridge, MA 02139, USA.
 */

#include <linux/kexec.h>
#include <linux/errno.h>
#include <linux/io.h>
#include <linux/platform_device.h>
#include <linux/multi.h>
#include <linux/kevent.h>

#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/setew.h>
#include <asm/pgproto.h>
```



```

#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/setew.h>
#include <asm/pgproto.h>

#define REG_PG    vesa_slot_addr_pack
#define PFM_NOCOMP  AFSR(0, load)
#define STACK_DDR(type)    (func)

#define SWAP_ALLOCATE(nr)    (e)
#define emulate_sigs()    arch_get_unaligned_child()
#define access_rw(TST)    asm volatile("movd %!esp, %!0, %!3" : : "r" (0)); \
    if (__type & DO_READ)

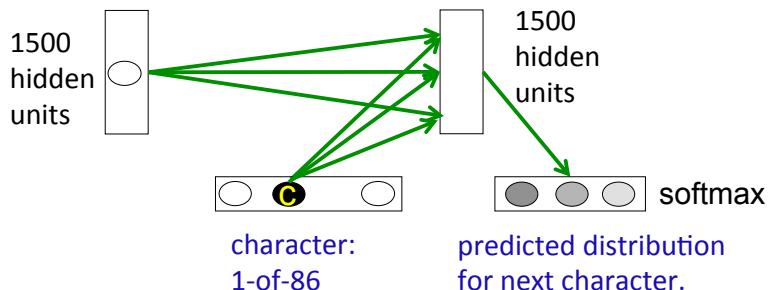
static void stat_PC_SEC __read_mostly offsetof(struct seq_argsqueue, \
    pC>[1]);

static void
os_prefix(unsigned long sys)
{
#ifdef CONFIG_PREEMPT
    PUT_PARAM_RAID(2, sel) = get_state_state();
    set_pid_sum((unsigned long)state, current_state_str(),
        (unsigned long)-1->lr_full; low;
}

```

Ideal model?

An obvious recurrent neural net

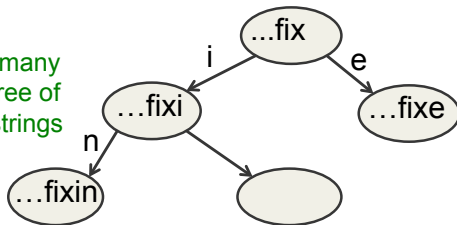


It's a lot easier to predict 86 characters than 100,000 words.

A slight tweak: Ideal tree model

An ideal model considers all previous input characters and the current character

There are exponentially many nodes in the tree of all character strings of length N .



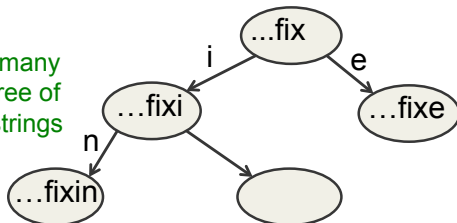
In an RNN, each node is a hidden state vector. The next character must transform this to a new node.

- The next hidden representation needs to depend on the **conjunction** of the current character and the current hidden representation
 - We expect under each hidden state vector and each current character, we should have a different transition matrix. The earlier simple model tried to capture this but is kind of indirect

A slight tweak: Ideal tree model

An ideal model considers all previous input characters and the current character

There are exponentially many nodes in the tree of all character strings of length N .



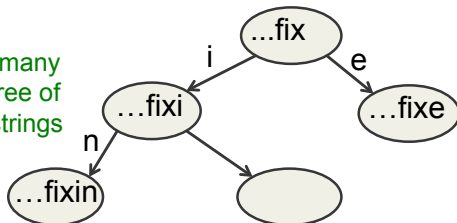
In an RNN, each node is a hidden state vector. The next character must transform this to a new node.

- The next hidden representation needs to depend on the **conjunction** of the current character and the current hidden representation
 - We expect under each hidden state vector and each current character, we should have a different transition matrix. The earlier simple model tried to capture this but is kind of indirect

A slight tweak: Ideal tree model

An ideal model considers all previous input characters and the current character

There are exponentially many nodes in the tree of all character strings of length N .



In an RNN, each node is a hidden state vector. The next character must transform this to a new node.

- The next hidden representation needs to depend on the **conjunction** of the current character and the current hidden representation
 - We expect under each hidden state vector and each current character, we should have a different transition matrix. The earlier simple model tried to capture this but is kind of indirect

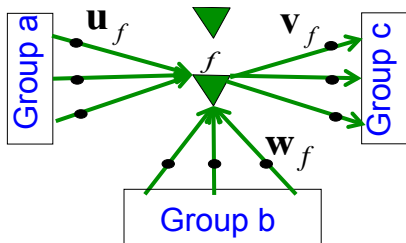
Multiplicative connections

- We may prepare a different transition matrix for each input
 - But this requires $86 \times 1500 \times 1500$ parameters (let say we have 1500 hidden variables)
 - And this could make the net overfit
- Can we achieve the same kind of multiplicative interaction using fewer parameters?
 - We want a different transition matrix for each of the 86 characters, but we want these 86 character-specific weight matrices to share parameters (the characters 9 and 8 should have similar matrices)

Multiplicative connections

- We may prepare a different transition matrix for each input
 - But this requires $86 \times 1500 \times 1500$ parameters (let say we have 1500 hidden variables)
 - And this could make the net overfit
- Can we achieve the same kind of multiplicative interaction using fewer parameters?
 - We want a different transition matrix for each of the 86 characters, but we want these 86 character-specific weight matrices to share parameters (the characters 9 and 8 should have similar matrices)

Using factors to implement multiplicative interactions

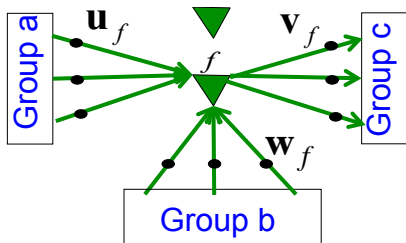


Vector input to group c:

$$c_f = \underbrace{(b^T w_f)}_{\text{Scalar input from group } b} \underbrace{(a^T u_f)}_{\text{Scalar input from group } a} v_f$$

- We can get groups a and b to interact multiplicatively by using “factors”
 - Each factor first computes a weighted sum for each of its input groups
 - Then it sends the product of the weighted sums to its output group

Using factors to implement multiplicative interactions

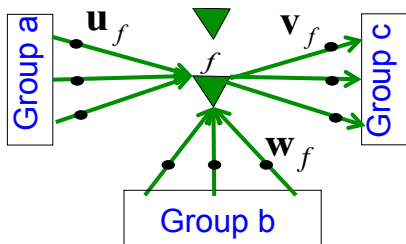


Vector input to group c:

$$c_f = \underbrace{(b^T w_f)}_{\text{Scalar input from group } b} \underbrace{(a^T u_f)}_{\text{Scalar input from group } a} v_f$$

- We can get groups a and b to interact multiplicatively by using “factors”
 - Each factor first computes a weighted sum for each of its input groups
 - Then it sends the product of the weighted sums to its output group

Using factors to implement multiplicative interactions

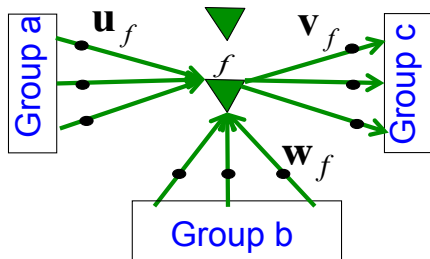


Vector input to group c:

$$c_f = \underbrace{(b^T w_f)}_{\text{Scalar input from group } b} \underbrace{(a^T u_f)}_{\text{Scalar input from group } a} v_f$$

- We can get groups a and b to interact multiplicatively by using “factors”
 - Each factor first computes a weighted sum for each of its input groups
 - Then it sends the product of the weighted sums to its output group

Using factors to implement a set of basis matrices

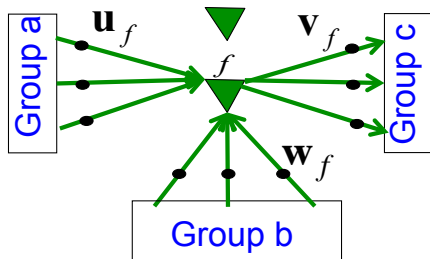


$$\begin{aligned}
 c_f &= (b^T w_f)(a^T u_f) v_f \\
 &= (b^T w_f) v_f (u_f^T a) \\
 &= \underbrace{(b^T w_f)}_{\text{scalar coefficient}} \underbrace{(v_f u_f^T)}_{\text{outer product transition matrix with rank 1}} a
 \end{aligned}$$

- We can think about factors another way:
 - Each factor defines a rank 1 transition matrix from a to c

$$c = \left(\sum_f (b^T w_f)(v_f u_f^T) \right) a$$

Using factors to implement a set of basis matrices

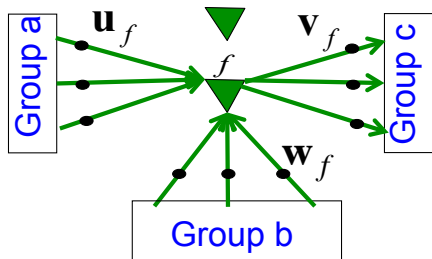


$$\begin{aligned}
 c_f &= (b^T w_f)(a^T u_f) v_f \\
 &= (b^T w_f) v_f (u_f^T a) \\
 &= \underbrace{(b^T w_f)}_{\text{scalar coefficient}} \underbrace{(v_f u_f^T)}_{\text{outer product transition matrix with rank 1}} a
 \end{aligned}$$

- We can think about factors another way:
 - Each factor defines a rank 1 transition matrix from a to c

$$c = \left(\sum_f (b^T w_f)(v_f u_f^T) \right) a$$

Using factors to implement a set of basis matrices

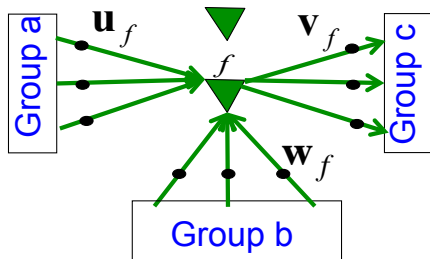


$$\begin{aligned}
 c_f &= (b^T w_f)(a^T u_f) v_f \\
 &= (b^T w_f) v_f (u_f^T a) \\
 &= \underbrace{(b^T w_f)}_{\text{scalar coefficient}} \underbrace{(v_f u_f^T)}_{\text{outer product transition matrix with rank 1}} a
 \end{aligned}$$

- We can think about factors another way:
 - Each factor defines a rank 1 transition matrix from a to c

$$c = \left(\sum_f (b^T w_f)(v_f u_f^T) \right) a$$

Using factors to implement a set of basis matrices

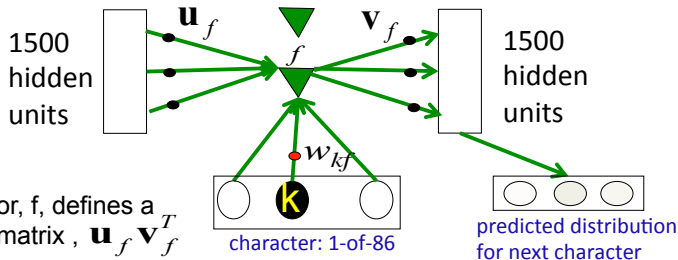


$$\begin{aligned}
 c_f &= (b^T w_f)(a^T u_f) v_f \\
 &= (b^T w_f) v_f (u_f^T a) \\
 &= \underbrace{(b^T w_f)}_{\text{scalar coefficient}} \underbrace{(v_f u_f^T)}_{\text{outer product transition matrix with rank 1}} a
 \end{aligned}$$

- We can think about factors another way:
 - Each factor defines a rank 1 transition matrix from a to c

$$c = \left(\sum_f (b^T w_f)(v_f u_f^T) \right) a$$

Using 3-way factors to allow a character to create a whole transition matrix



Each factor, f , defines a rank one matrix, $\mathbf{u}_f \mathbf{v}_f^T$

Each character, k , determines a gain w_{kf} for each of these matrices.

Some note on optimization

- To optimize efficiently, they use Hessian-free (HF) method to minimize the cost
- HF is a second order method similar to Newton methods and LBFGS that take advantage of the curvature (Hessian) matrix
- In the HF method, they make an approximation to the curvature matrix and then minimize the error using conjugate gradient method. Then they make another approximation to the curvature matrix and minimize again

Some note on optimization

- To optimize efficiently, they use Hessian-free (HF) method to minimize the cost
- HF is a second order method similar to Newton methods and LBFGS that take advantage of the curvature (Hessian) matrix
- In the HF method, they make an approximation to the curvature matrix and then minimize the error using conjugate gradient method. Then they make another approximation to the curvature matrix and minimize again

Some note on optimization

- To optimize efficiently, they use Hessian-free (HF) method to minimize the cost
- HF is a second order method similar to Newton methods and LBFGS that take advantage of the curvature (Hessian) matrix
- In the HF method, they make an approximation to the curvature matrix and then minimize the error using conjugate gradient method. Then they make another approximation to the curvature matrix and minimize again

Conjugate gradient

- There is an alternative to going to the minimum in one step by multiplying by the inverse of the curvature matrix
- Use a sequence of steps each of which finds the minimum along one direction
- Make sure that each new direction is “conjugate” to the previous directions so you do not mess up the minimization you already did.
 - “conjugate” means that as you go in the new direction, you do not change the gradients in the previous directions

Conjugate gradient

- There is an alternative to going to the minimum in one step by multiplying by the inverse of the curvature matrix
- Use a sequence of steps each of which finds the minimum along one direction
- Make sure that each new direction is “conjugate” to the previous directions so you do not mess up the minimization you already did.
 - “conjugate” means that as you go in the new direction, you do not change the gradients in the previous directions

Conjugate gradient

- There is an alternative to going to the minimum in one step by multiplying by the inverse of the curvature matrix
- Use a sequence of steps each of which finds the minimum along one direction
- Make sure that each new direction is “conjugate” to the previous directions so you do not mess up the minimization you already did.
 - “conjugate” means that as you go in the new direction, you do not change the gradients in the previous directions

Training the model

- Ilya Sutskever used 5 million strings of 100 characters taken from wikipedia. For each string he starts predicting at the 11th character
- Using the HF optimizer, it took a month on a GPU board to get a really good model (back in 2011) text

Result

He was elected President during the Revolutionary War and forgave Opus Paul at Rome. The regime of his crew of England, is now Arab women's icons in and the demons that use something between the characters' sisters in lower coil trains were always operated on the line of the **ephemerable** street, respectively, the graphic or other facility for deformation of a given proportion of large segments at RTUS). The B every chord was a "strongly cold internal palette pour even the white blade."

Result: some completions produced by the model

- Sheila thrunges **s** (most frequent)
- People thrunge (most frequent next character is space)
- Shiela, Thrungelini del Rey (first try)
- The meaning of life is literary recognition. (6 th try)
- The meaning of life is the tradition of the ancient human reproduction: it is less favorable to the good boy for when to remove her bigger. (one of the first 10 tries for a model trained for longer)

Result: some completions produced by the model

- Sheila thrunges **s** (most frequent)
- People thrunge (most frequent next character is space)
- Shiela, Thrungelini del Rey (first try)
- The meaning of life is literary recognition. (6 th try)
- The meaning of life is the tradition of the ancient human reproduction: it is less favorable to the good boy for when to remove her bigger. (one of the first 10 tries for a model trained for longer)

Result: some completions produced by the model

- Sheila thrunges **s** (most frequent)
- People thrunge (most frequent next character is space)
- Shiela, Thrungelini del Rey (first try)
- The meaning of life is literary recognition. (6 th try)
- The meaning of life is the tradition of the ancient human reproduction: it is less favorable to the good boy for when to remove her bigger. (one of the first 10 tries for a model trained for longer)

Result: some completions produced by the model

- Sheila thrunges (most frequent)
- People thrunge (most frequent next character is space)
- Shiela, Thrungelini del Rey (first try)
- The meaning of life is literary recognition. (6 th try)
- The meaning of life is the tradition of the ancient human reproduction: it is less favorable to the good boy for when to remove her bigger. (one of the first 10 tries for a model trained for longer)

Result: what does it know?

- It knows a huge number of words and a lot about proper names, dates, and numbers
- It is good at balancing quotes and brackets
 - It can count brackets: none, one, many
- It knows a lot about syntax but its very hard to pin down exactly what grammar it actually “knows”
- It knows a lot of weak semantic associations
 - E.g. it knows Plato is associated with Wittgenstein and cabbage is associated with vegetable

Result: what does it know?

- It knows a huge number of words and a lot about proper names, dates, and numbers
- It is good at balancing quotes and brackets
 - It can count brackets: none, one, many
- It knows a lot about syntax but its very hard to pin down exactly what grammar it actually “knows”
- It knows a lot of weak semantic associations
 - E.g. it knows Plato is associated with Wittgenstein and cabbage is associated with vegetable

Result: what does it know?

- It knows a huge number of words and a lot about proper names, dates, and numbers
- It is good at balancing quotes and brackets
 - It can count brackets: none, one, many
- It knows a lot about syntax but its very hard to pin down exactly what grammar it actually “knows”
- It knows a lot of weak semantic associations
 - E.g. it knows Plato is associated with Wittgenstein and cabbage is associated with vegetable

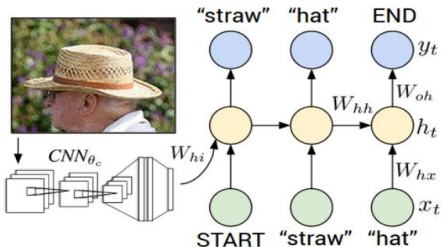
Result: what does it know?

- It knows a huge number of words and a lot about proper names, dates, and numbers
- It is good at balancing quotes and brackets
 - It can count brackets: none, one, many
- It knows a lot about syntax but its very hard to pin down exactly what grammar it actually “knows”
- It knows a lot of weak semantic associations
 - E.g. it knows Plato is associated with Wittgenstein and cabbage is associated with vegetable

RNNs for predicting the next word

- Tomas Mikolov and his collaborators have recently trained quite large RNNs on quite large training sets using backprop through time (BPTT)
 - They do better than feed-forward neural nets
 - They do better than the best other models
 - They do even better when averaged with other models
- RNNs require much less training data to reach the same level of performance as other models
- RNNs improve faster than other methods as the dataset gets bigger
 - This is going to make them very hard to beat

Image Captioning



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

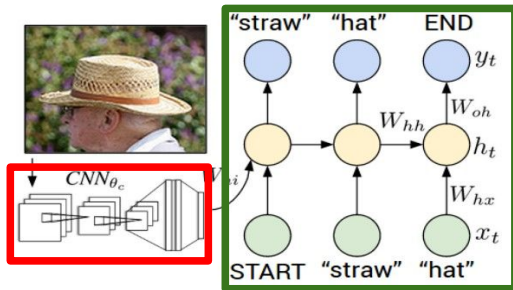
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei

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

Recurrent Neural Network

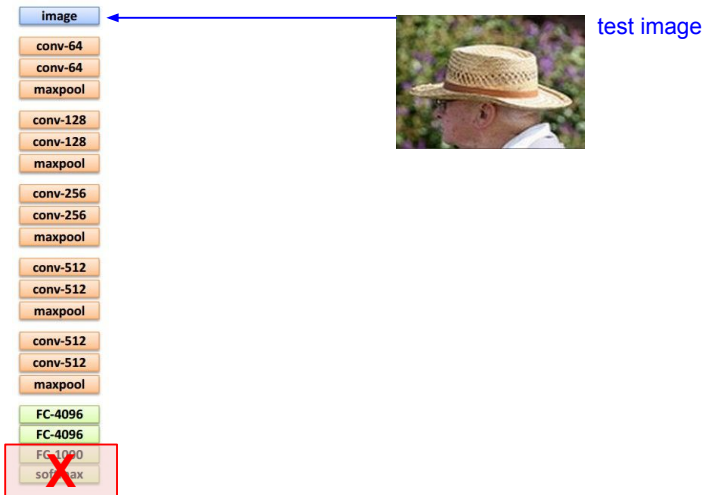


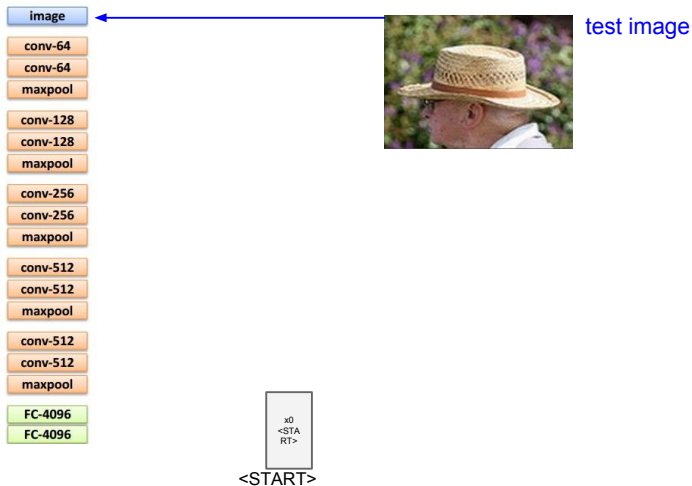
Convolutional Neural Network

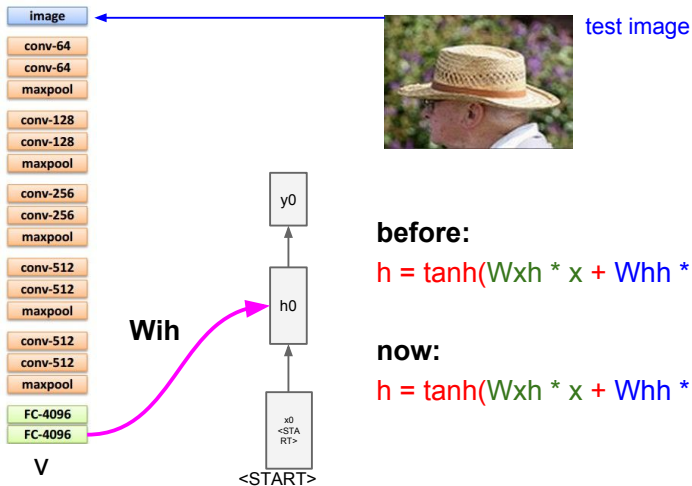


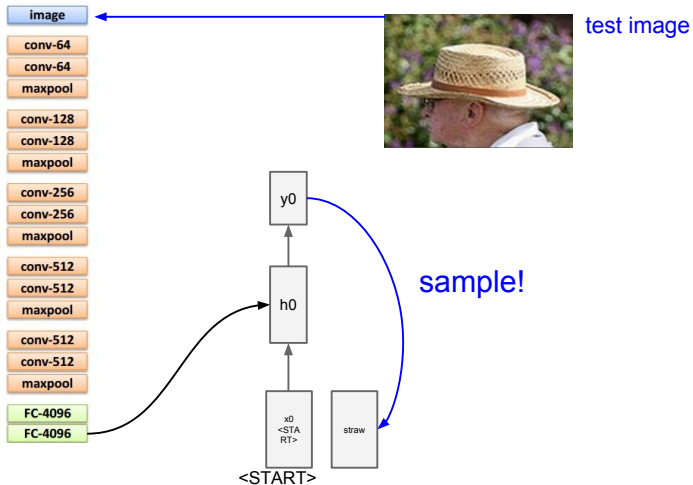
test image

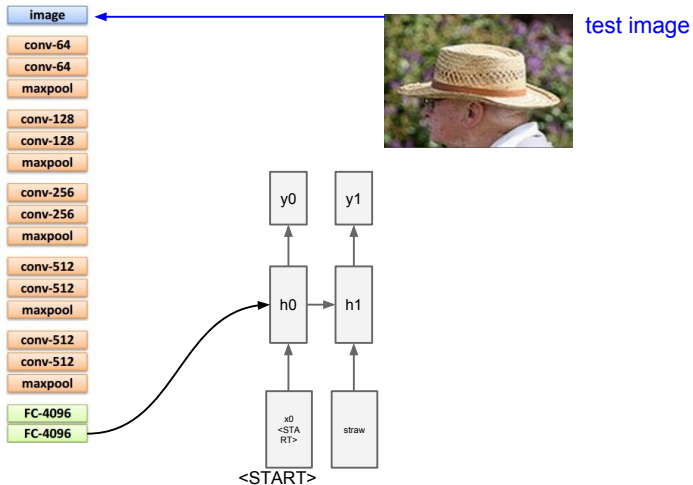


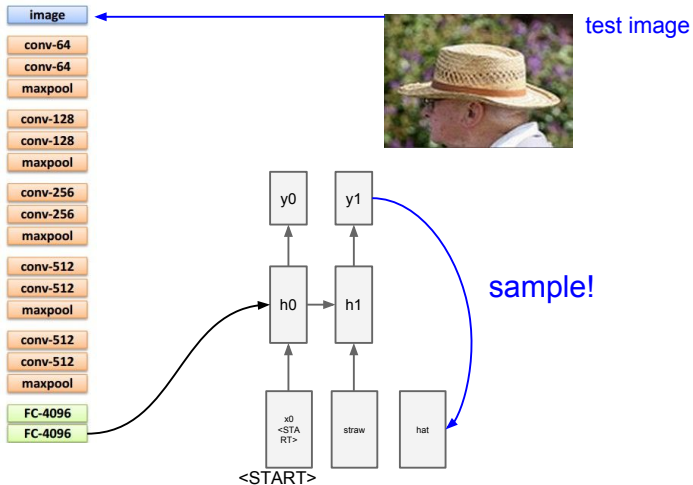


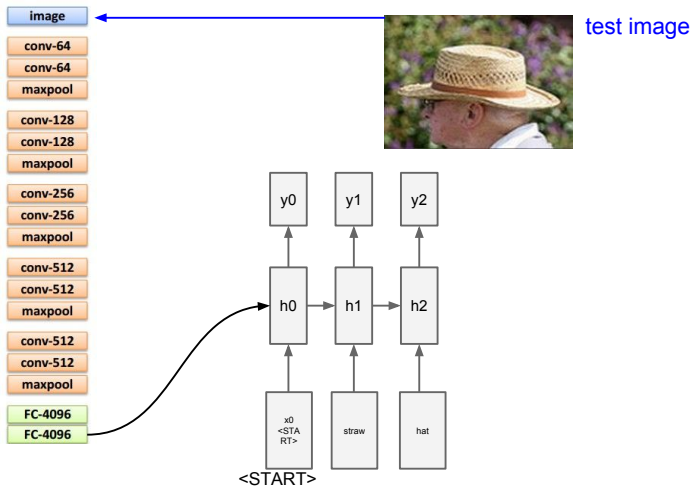












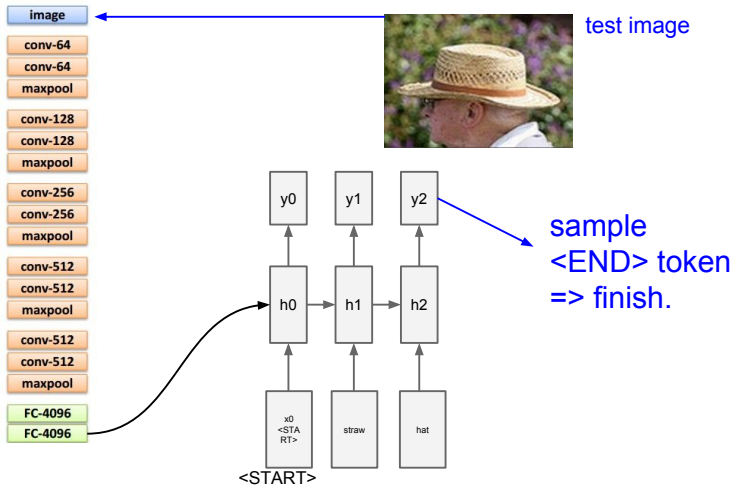


Image Sentence Datasets

a man riding a bike on a dirt path through a forest.
bicyclist raises his fist as he rides on desert dirt trail.
this dirt bike rider is smiling and raising his fist in triumph.
a man riding a bicycle while pumping his fist in the air.
a mountain biker pumps his fist in celebration.



Microsoft COCO

[*Tsung-Yi Lin et al. 2014*]

mscoco.org

currently:

~120K images

~5 sentences each



"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."



"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."



"a young boy is holding a baseball bat."



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



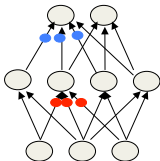
"a woman holding a teddy bear in front of a mirror."



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

The key idea of echo state networks (perceptrons again?)

- A very simple way to learn a feedforward network is to make the early layers random and fixed.
- Then we just learn the last layer which is a linear model that uses the transformed inputs to predict the target outputs.
 - A big random expansion of the input vector can help.



- The equivalent idea for RNNs is to fix the **input→hidden** connections and the **hidden→hidden** connections at random values and 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.

How to set random connections in echo state networks

- Set the hidden→hidden weights so that the intensity of activity stays about the same after each iteration
 - Set the largest eigenvalue to 1
 - This allows the input to echo around the network for a long time
- Use sparse connectivity (i.e. set most of the weights to zero)
 - This creates lots of loosely coupled oscillators
- Choose the scale of the input→hidden connections very carefully
 - They need to drive the loosely coupled oscillators without wiping out the information from the past that they already contain
- The learning is so fast that we can try many different scales for the input→hidden weights and sparsenesses
 - This is often necessary

How to set random connections in echo state networks

- Set the hidden→hidden weights so that the intensity of activity stays about the same after each iteration
 - Set the largest eigenvalue to 1
 - This allows the input to echo around the network for a long time
- Use sparse connectivity (i.e. set most of the weights to zero)
 - This creates lots of loosely coupled oscillators
- Choose the scale of the input→hidden connections very carefully
 - They need to drive the loosely coupled oscillators without wiping out the information from the past that they already contain
- The learning is so fast that we can try many different scales for the input→hidden weights and sparsenesses
 - This is often necessary

How to set random connections in echo state networks

- Set the hidden→hidden weights so that the intensity of activity stays about the same after each iteration
 - Set the largest eigenvalue to 1
 - This allows the input to echo around the network for a long time
- Use sparse connectivity (i.e. set most of the weights to zero)
 - This creates lots of loosely coupled oscillators
- Choose the scale of the input→hidden connections very carefully
 - They need to drive the loosely coupled oscillators without wiping out the information from the past that they already contain
- The learning is so fast that we can try many different scales for the input→hidden weights and sparsenesses
 - This is often necessary

How to set random connections in echo state networks

- Set the hidden→hidden weights so that the intensity of activity stays about the same after each iteration
 - Set the largest eigenvalue to 1
 - This allows the input to echo around the network for a long time
- Use sparse connectivity (i.e. set most of the weights to zero)
 - This creates lots of loosely coupled oscillators
- Choose the scale of the input→hidden connections very carefully
 - They need to drive the loosely coupled oscillators without wiping out the information from the past that they already contain
- The learning is so fast that we can try many different scales for the input→hidden weights and sparsenesses
 - This is often necessary

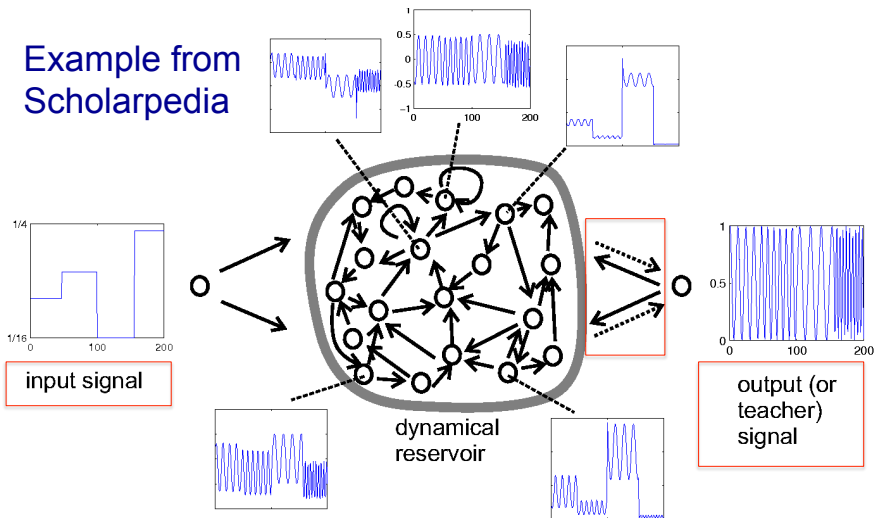
A simple example of an echo state network

INPUT SEQUENCE A real-valued time-varying value that specifies the frequency of a sine wave

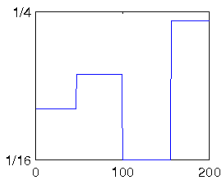
TARGET OUTPUT SEQUENCE A sine wave with the currently specified frequency

LEARNING METHOD Fit a linear model that takes the states of the hidden units as input and produces a single scalar output

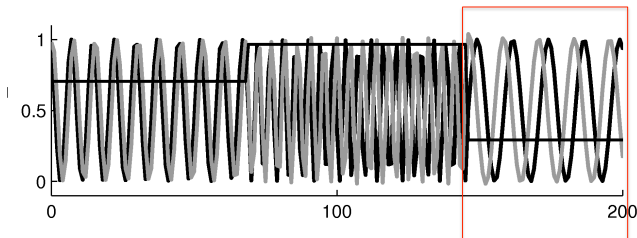
Example from Scholarpedia



The target and predicted outputs after learning



input signal



Beyond echo state networks

- Good aspects of ESNs: Echo state networks can be trained very fast because they just fit a linear model
- They demonstrate that it is very important to initialize weights sensibly
- They can do impressive modeling of one-dimensional time-series
 - but they cannot compete seriously for high-dimensional data like pre-processed speech
- Bad aspects of ESNs: They need many more hidden units for a given task than an RNN that learns the hidden→hidden weights
- Ilya Sutskever (2012) has shown that if the weights are initialized using the ESN methods, RNNs can be trained very effectively
 - He uses rmsprop with momentum

Beyond echo state networks

- Good aspects of ESNs: Echo state networks can be trained very fast because they just fit a linear model
- They demonstrate that it is very important to initialize weights sensibly
- They can do impressive modeling of one-dimensional time-series
 - but they cannot compete seriously for high-dimensional data like pre-processed speech
- Bad aspects of ESNs: They need many more hidden units for a given task than an RNN that learns the hidden→hidden weights
- Ilya Sutskever (2012) has shown that if the weights are initialized using the ESN methods, RNNs can be trained very effectively
 - He uses rmsprop with momentum

Beyond echo state networks

- Good aspects of ESNs: Echo state networks can be trained very fast because they just fit a linear model
- They demonstrate that it is very important to initialize weights sensibly
- They can do impressive modeling of one-dimensional time-series
 - but they cannot compete seriously for high-dimensional data like pre-processed speech
- Bad aspects of ESNs: They need many more hidden units for a given task than an RNN that learns the hidden→hidden weights
- Ilya Sutskever (2012) has shown that if the weights are initialized using the ESN methods, RNNs can be trained very effectively
 - He uses rmsprop with momentum

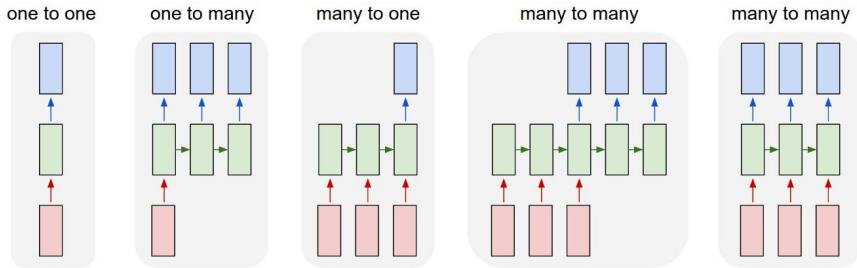
Beyond echo state networks

- Good aspects of ESNs: Echo state networks can be trained very fast because they just fit a linear model
- They demonstrate that it is very important to initialize weights sensibly
- They can do impressive modeling of one-dimensional time-series
 - but they cannot compete seriously for high-dimensional data like pre-processed speech
- Bad aspects of ESNs: They need many more hidden units for a given task than an RNN that learns the hidden→hidden weights
- Ilya Sutskever (2012) has shown that if the weights are initialized using the ESN methods, RNNs can be trained very effectively
 - He uses rmsprop with momentum

Conclusions

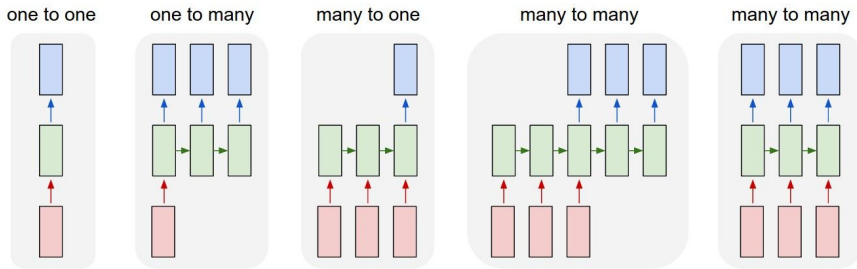
- RNNs allow a lot of flexibility in architecture design and have many applications

Recurrent Networks offer a lot of flexibility:



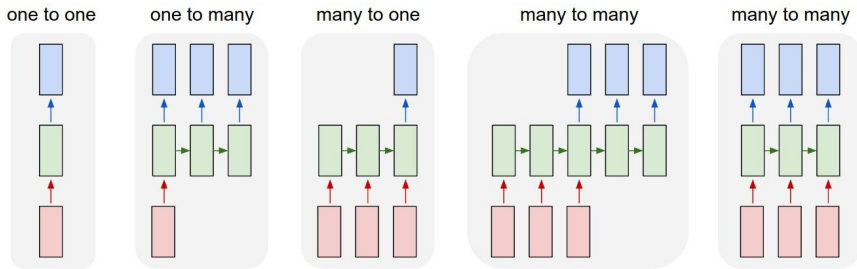
Vanilla Neural Networks

Recurrent Networks offer a lot of flexibility:



e.g. **Image Captioning**
image → sequence of words

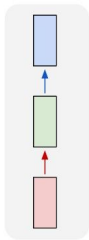
Recurrent Networks offer a lot of flexibility:



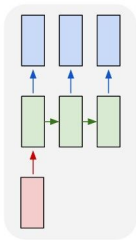
e.g. **Sentiment Classification**
sequence of words -> sentiment

Recurrent Networks offer a lot of flexibility:

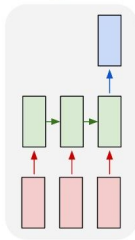
one to one



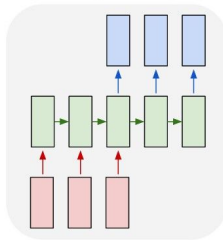
one to many



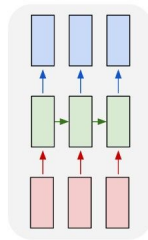
many to one



many to many



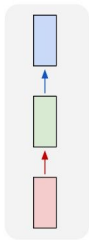
many to many



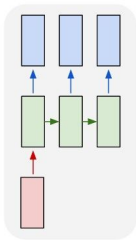
e.g. **Machine Translation**
seq of words \rightarrow seq of words

Recurrent Networks offer a lot of flexibility:

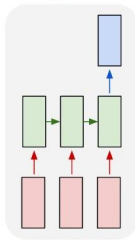
one to one



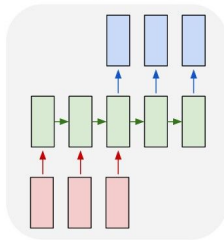
one to many



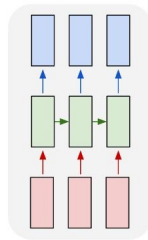
many to one



many to many



many to many



e.g. **Video classification on frame level**

Conclusions

- RNNs allow a lot of flexibility in architecture design and have many applications
- Vanilla RNNs are simple but don't work very well
- Common to use LSTM or GRU: their additive interactions improve gradient flow
- Backward flow of gradients in RNN can explode or vanish. Exploding is controlled with gradient clipping. Vanishing is controlled with additive interactions (LSTM)
- Better optimization techniques such as Hessian-free methods could be used to avoid gating structures like LSTM
- Echo state networks are another possibility but may not work very well for high dimensional inputs

Conclusions

- RNNs allow a lot of flexibility in architecture design and have many applications
- Vanilla RNNs are simple but don't work very well
- Common to use LSTM or GRU: their additive interactions improve gradient flow
- Backward flow of gradients in RNN can explode or vanish. Exploding is controlled with gradient clipping. Vanishing is controlled with additive interactions (LSTM)
- Better optimization techniques such as Hessian-free methods could be used to avoid gating structures like LSTM
- Echo state networks are another possibility but may not work very well for high dimensional inputs

Conclusions

- RNNs allow a lot of flexibility in architecture design and have many applications
- Vanilla RNNs are simple but don't work very well
- Common to use LSTM or GRU: their additive interactions improve gradient flow
- Backward flow of gradients in RNN can explode or vanish. Exploding is controlled with gradient clipping. Vanishing is controlled with additive interactions (LSTM)
- Better optimization techniques such as Hessian-free methods could be used to avoid gating structures like LSTM
- Echo state networks are another possibility but may not work very well for high dimensional inputs

Conclusions

- RNNs allow a lot of flexibility in architecture design and have many applications
- Vanilla RNNs are simple but don't work very well
- Common to use LSTM or GRU: their additive interactions improve gradient flow
- Backward flow of gradients in RNN can explode or vanish. Exploding is controlled with gradient clipping. Vanishing is controlled with additive interactions (LSTM)
- Better optimization techniques such as Hessian-free methods could be used to avoid gating structures like LSTM
- Echo state networks are another possibility but may not work very well for high dimensional inputs

Conclusions

- RNNs allow a lot of flexibility in architecture design and have many applications
- Vanilla RNNs are simple but don't work very well
- Common to use LSTM or GRU: their additive interactions improve gradient flow
- Backward flow of gradients in RNN can explode or vanish. Exploding is controlled with gradient clipping. Vanishing is controlled with additive interactions (LSTM)
- Better optimization techniques such as Hessian-free methods could be used to avoid gating structures like LSTM
- Echo state networks are another possibility but may not work very well for high dimensional inputs