



10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

RNNs + PAC Learning

Matt Gormley Lecture 15 Mar. 4, 2020

Reminders

- Homework 5: Neural Networks
 - Out: Fri, Feb. 28
 - Due: Wed, Mar. 18 at 11:59pm
- Today's In-Class Poll
 - http://poll.mlcourse.org (latest version works on mobile!)
- Exam Viewing

Q&A

RECURRENT NEURAL NETWORKS

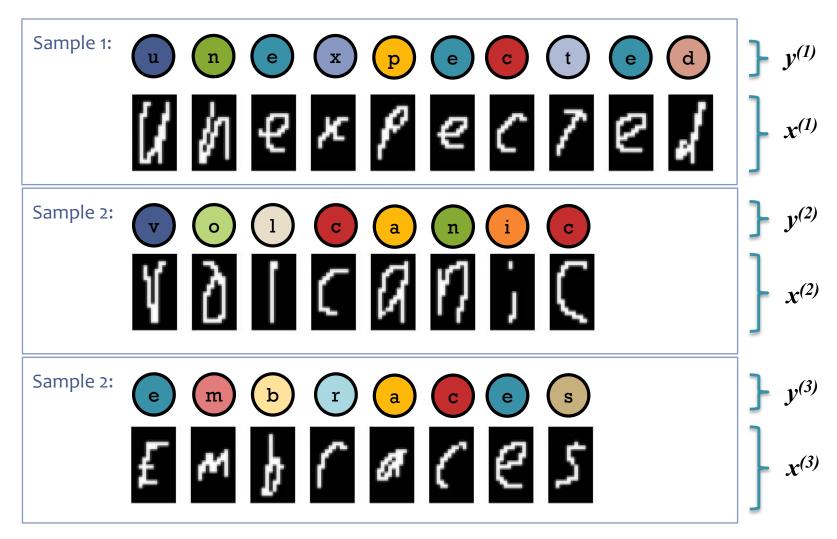
Dataset for Supervised Part-of-Speech (POS) Tagging

Data: $\mathcal{D} = \{oldsymbol{x}^{(n)}, oldsymbol{y}^{(n)}\}_{n=1}^N$

Sample 1:	n	(flies)	p like	an	$ \begin{array}{c c} & & \\ & $
Sample 2:	n	n	like	d	$\begin{array}{c c} & & \\ & & \\ \hline & & \\ &$
Sample 3:	n	fly	with	n	$\begin{cases} \mathbf{n} \\ \mathbf{vings} \end{cases} \mathbf{y}^{(3)}$
Sample 4:	with	n	you	will	$\begin{cases} \mathbf{v} \\ \mathbf{see} \end{cases} y^{(4)}$

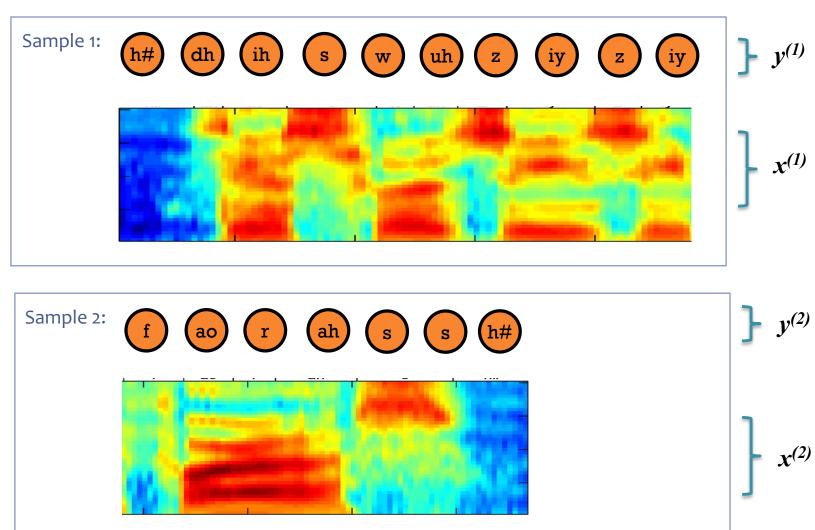
Dataset for Supervised Handwriting Recognition

Data: $\mathcal{D} = \{oldsymbol{x}^{(n)}, oldsymbol{y}^{(n)}\}_{n=1}^N$



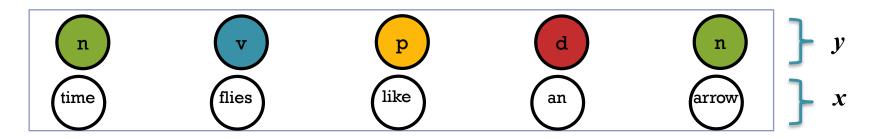
Dataset for Supervised Phoneme (Speech) Recognition

Data: $\mathcal{D} = \{oldsymbol{x}^{(n)}, oldsymbol{y}^{(n)}\}_{n=1}^N$



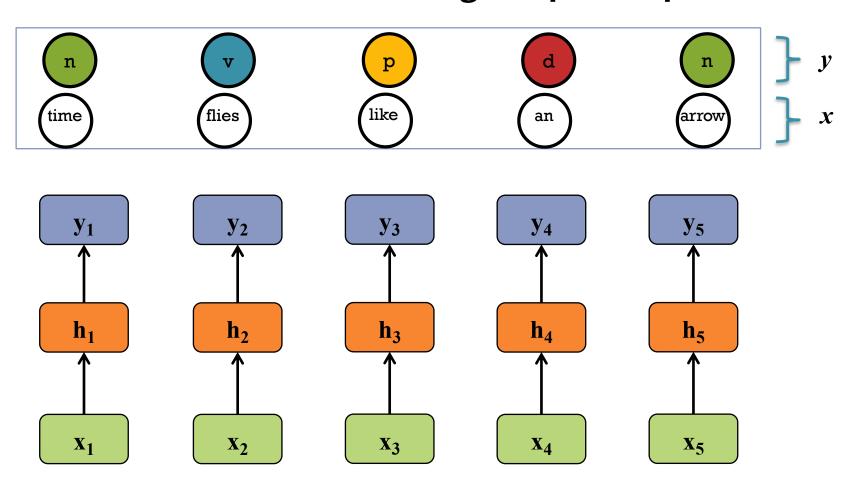
Time Series Data

Question 1: How could we apply the neural networks we've seen so far (which expect **fixed size input/output**) to a prediction task with **variable length input/output**?



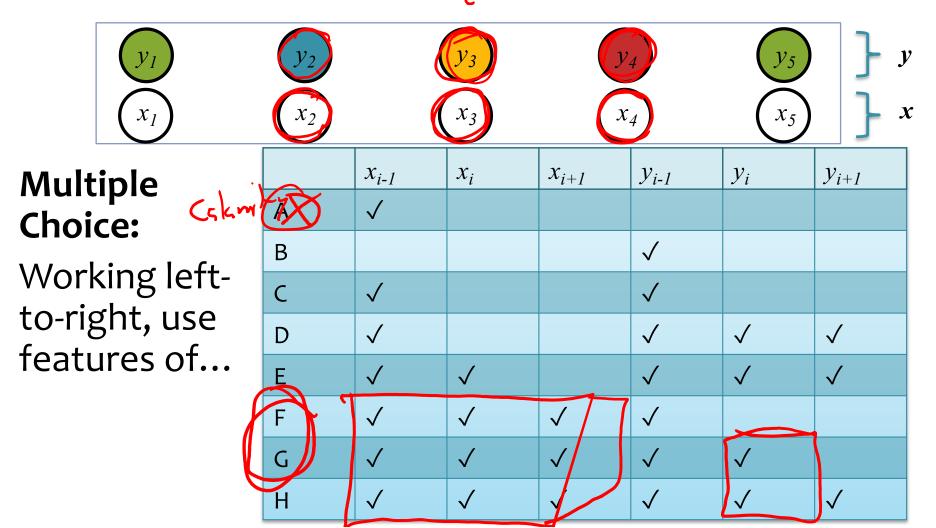
Time Series Data

Question 1: How could we apply the neural networks we've seen so far (which expect **fixed size input/output**) to a prediction task with **variable length input/output**?



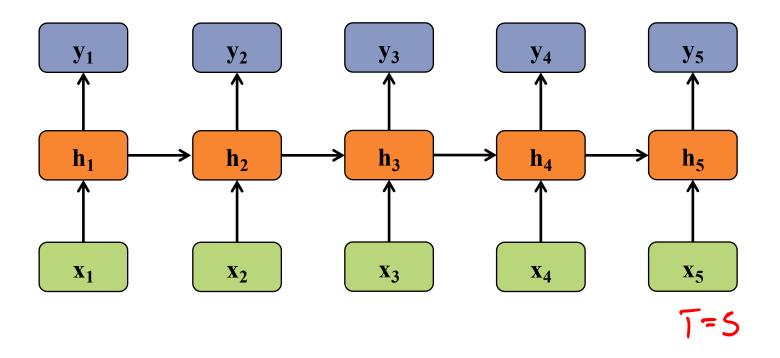
Time Series Data

Question 2: How could we incorporate context (e.g. words to the left/right, or tags to the left/right) into our solution?



inputs:
$$\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$$
hidden units: $\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$
outputs: $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$
nonlinearity: \mathcal{H}

$$h_t = \mathcal{H}(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$$
$$y_t = W_{hy}h_t + b_y$$



inputs: $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

hidden units: $\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$

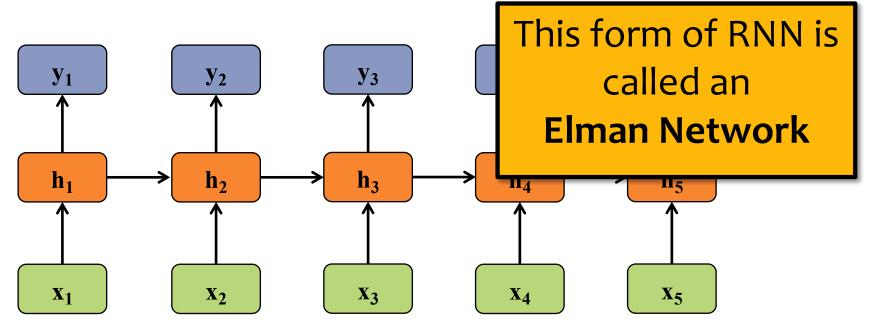
outputs: $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

nonlinearity: \mathcal{H}

$$h_t = \mathcal{H}\left(W_{xh}x_t + W_{hh}h_{t-1} + b_h\right)$$

$$y_t = W_{hy}h_t + b_y$$





inputs:
$$\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$$

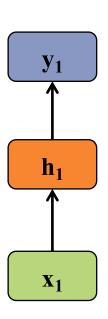
hidden units: $\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$

outputs:
$$\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$$

nonlinearity: \mathcal{H}

$$h_t = \mathcal{H}\left(W_{xh}x_t + W_{hh}h_{t-1} + b_h\right)$$

$$y_t = W_{hy}h_t + b_y$$



- If T=1, then we have a standard feed-forward neural net with one hidden layer
- All of the deep nets from last lecture required fixed size inputs/outputs

Background

A Recipe for Machine Learning

1. Given training data:

$$\{oldsymbol{x}_i, oldsymbol{y}_i\}_{i=1}^N$$

3. Define goal:

$$oldsymbol{ heta}^* = rg\min_{oldsymbol{ heta}} \sum_{i=1}^N \ell(f_{oldsymbol{ heta}}(oldsymbol{x}_i), oldsymbol{y}_i)$$

- 2. Choose each of these:
 - Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Loss function

$$\ell(\hat{m{y}},m{y}_i)\in\mathbb{R}$$

4. Train with SGD:

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

Background

A Recipe for Machine Learning

- Recurrent Neural Networks (RNNs) provide another form of decision function
 - An RNN is just another differential function

Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Train with SGD:

(take small steps opposite the gradient)

- We'll just need a method of computing the gradient efficiently
- Let's use Backpropagation Through Time...



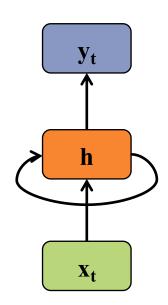
inputs: $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

outputs: $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K \mid y_t = W_{hy} h_t + b_y$

nonlinearity: \mathcal{H}

hidden units:
$$\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$$
 $h_t = \mathcal{H}(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$

$$y_t = W_{hy}h_t + b_y$$



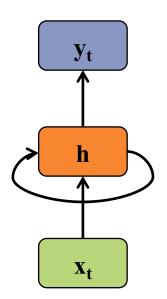
inputs:
$$\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$$

hidden units: $\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$
outputs: $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

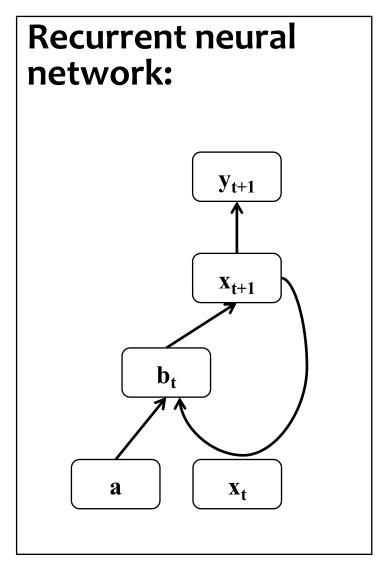
nonlinearity: \mathcal{H}

$$h_t = \mathcal{H}(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$$
$$y_t = W_{hy}h_t + b_y$$

- By unrolling the RNN through time, we can share parameters and accommodate arbitrary length input/output pairs
- Applications: time-series data such as sentences, speech, stock-market, signal data, etc.

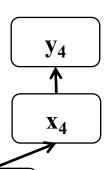


Background: Backprop through time



BPTT:

1. Unroll the computation over time





 $\mathbf{b_3}$ \mathbf{X}_{3} $\mathbf{b_2}$ $\mathbf{X_2}$ $\mathbf{b_1}$ $\mathbf{X_1}$ a

2. Run backprop through the resulting feed-forward network

Bidirectional RNN

inputs:
$$\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$$

hidden units: $\overrightarrow{\mathbf{h}}$ and $\overleftarrow{\mathbf{h}}$

outputs:
$$\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$$

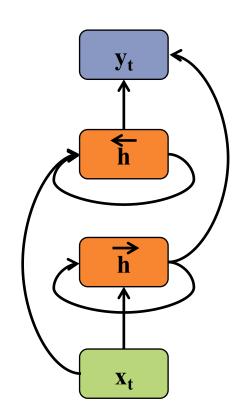
nonlinearity: \mathcal{H}

inputs:
$$\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$$
 en units: $\overrightarrow{\mathbf{h}}$ and $\overleftarrow{\mathbf{h}}$ outputs: $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$ linearity: \mathcal{H}

Recursive Definition:
$$\overrightarrow{h}_t = \mathcal{H}\left(W_x \overrightarrow{h} x_t + W_{\overrightarrow{h}} \overrightarrow{h} \overrightarrow{h}_{t-1} + b_{\overrightarrow{h}}\right)$$

$$\overleftarrow{h}_t = \mathcal{H}\left(W_x \overleftarrow{h} x_t + W_{\overleftarrow{h}} \overleftarrow{h} \overleftarrow{h}_{t+1} + b_{\overleftarrow{h}}\right)$$

$$y_t = W_{\overrightarrow{h}y} \overrightarrow{h}_t + W_{\overleftarrow{h}y} \overleftarrow{h}_t + b_y$$



Bidirectional RNN

inputs: $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

hidden units: $\overrightarrow{\mathbf{h}}$ and $\overleftarrow{\mathbf{h}}$

outputs: $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

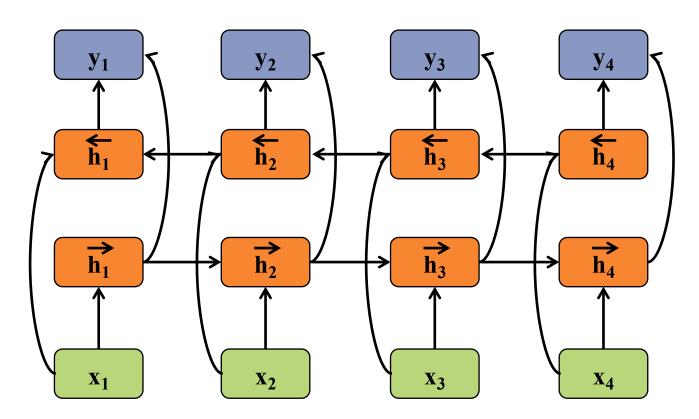
nonlinearity: \mathcal{H}

Recursive Definition:

$$\overrightarrow{h}_{t} = \mathcal{H}\left(W_{x\overrightarrow{h}}x_{t} + W_{\overrightarrow{h}}\overrightarrow{h}\overrightarrow{h}\overrightarrow{h}_{t-1} + b_{\overrightarrow{h}}\right)$$

$$\overleftarrow{h}_{t} = \mathcal{H}\left(W_{x\overleftarrow{h}}x_{t} + W_{\overleftarrow{h}}\overleftarrow{h}\overleftarrow{h}_{t+1} + b_{\overleftarrow{h}}\right)$$

$$y_{t} = W_{\overrightarrow{h}}\overrightarrow{h}_{y}\overrightarrow{h}_{t} + W_{\overleftarrow{h}}\overrightarrow{h}_{y}\overleftarrow{h}_{t} + b_{y}$$



Bidirectional RNN

inputs: $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

hidden units: $\overrightarrow{\mathbf{h}}$ and $\overleftarrow{\mathbf{h}}$

outputs: $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

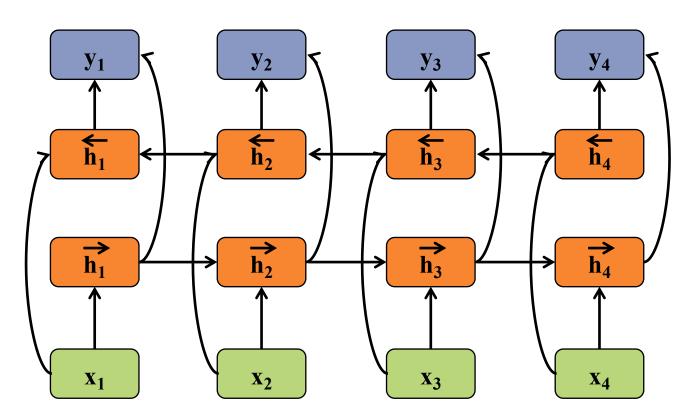
nonlinearity: \mathcal{H}

Recursive Definition:

$$\overrightarrow{h}_{t} = \mathcal{H}\left(W_{x\overrightarrow{h}}x_{t} + W_{\overrightarrow{h}}\overrightarrow{h}\overrightarrow{h}\overrightarrow{h}_{t-1} + b_{\overrightarrow{h}}\right)$$

$$\overleftarrow{h}_{t} = \mathcal{H}\left(W_{x\overleftarrow{h}}x_{t} + W_{\overleftarrow{h}}\overleftarrow{h}\overleftarrow{h}_{t+1} + b_{\overleftarrow{h}}\right)$$

$$y_{t} = W_{\overrightarrow{h}}\overrightarrow{h}_{y}\overrightarrow{h}_{t} + W_{\overleftarrow{h}}\overrightarrow{h}_{y}\overleftarrow{h}_{t} + b_{y}$$



Deep RNNs

inputs: $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

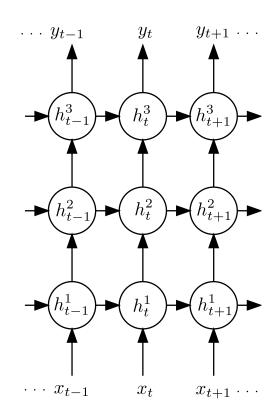
outputs: $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

nonlinearity: \mathcal{H}

Recursive Definition:

$$h_t^n = \mathcal{H}\left(W_{h^{n-1}h^n}h_t^{n-1} + W_{h^nh^n}h_{t-1}^n + b_h^n\right)$$

$$y_t = W_{h^N y} h_t^N + b_y$$



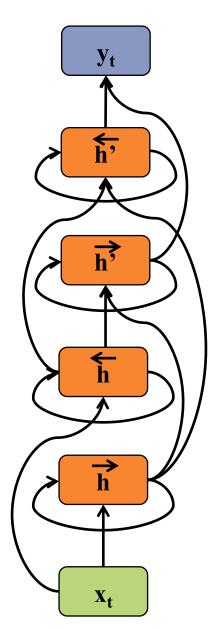
Deep Bidirectional RNNs

inputs: $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

outputs: $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

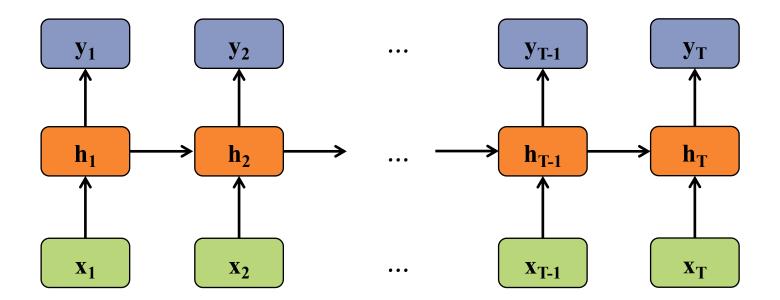
nonlinearity: \mathcal{H}

- Notice that the upper level hidden units have input from two previous layers (i.e. wider input)
- Likewise for the output layer
- What analogy can we draw to DNNs, DBNs, DBMs?



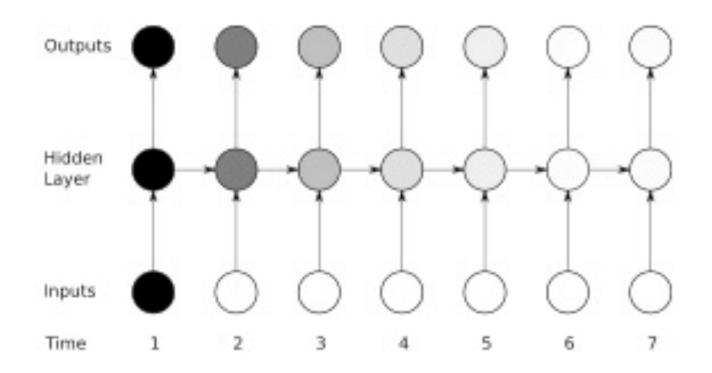
Motivation:

- Standard RNNs have trouble learning long distance dependencies
- LSTMs combat this issue



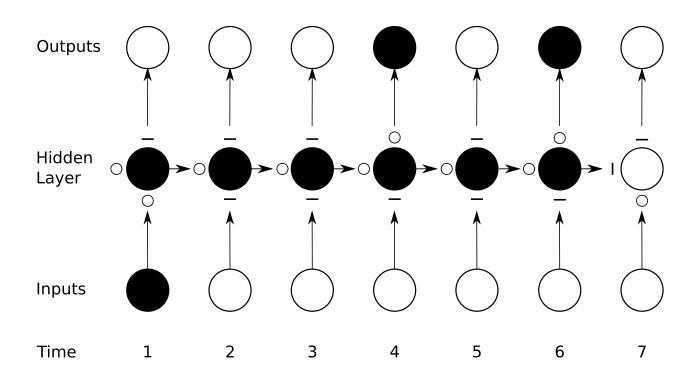
Motivation:

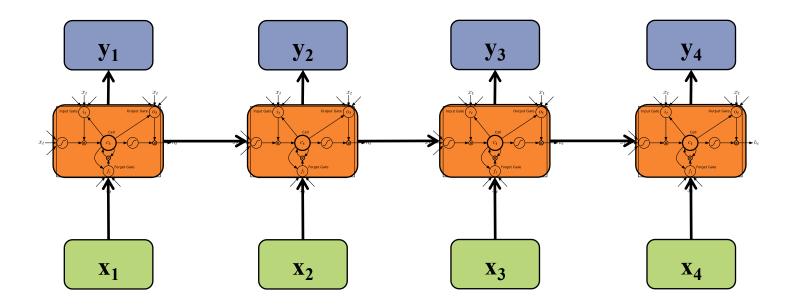
- Vanishing gradient problem for Standard RNNs
- Figure shows sensitivity (darker = more sensitive) to the input at time t=1



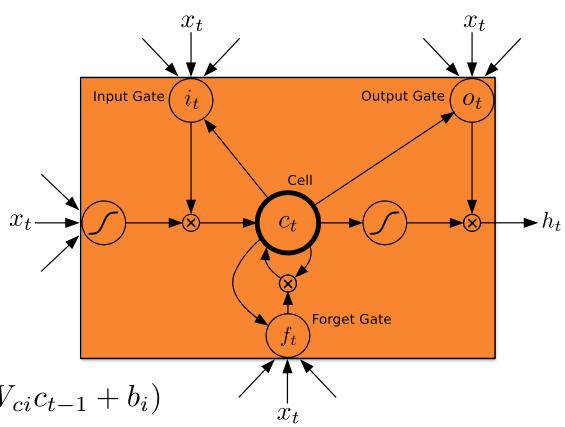
Motivation:

- LSTM units have a rich internal structure
- The various "gates" determine the propagation of information and can choose to "remember" or "forget" information





- Input gate: masks out the standard RNN inputs
- Forget gate: masks out the previous cell
- Cell: stores the input/forget mixture
- Output gate: masks out the values of the next hidden



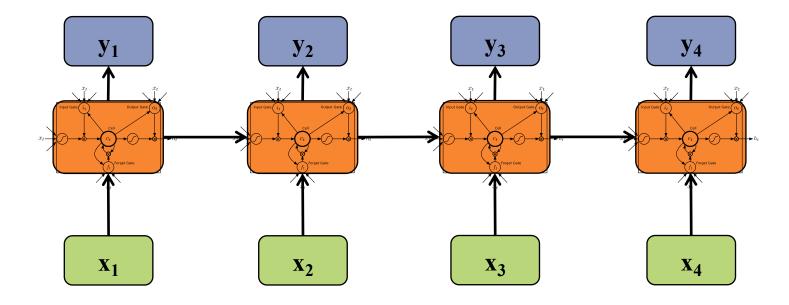
$$i_{t} = \sigma \left(W_{xi} x_{t} + W_{hi} h_{t-1} + W_{ci} c_{t-1} + b_{i} \right)$$

$$f_{t} = \sigma \left(W_{xf} x_{t} + W_{hf} h_{t-1} + W_{cf} c_{t-1} + b_{f} \right)$$

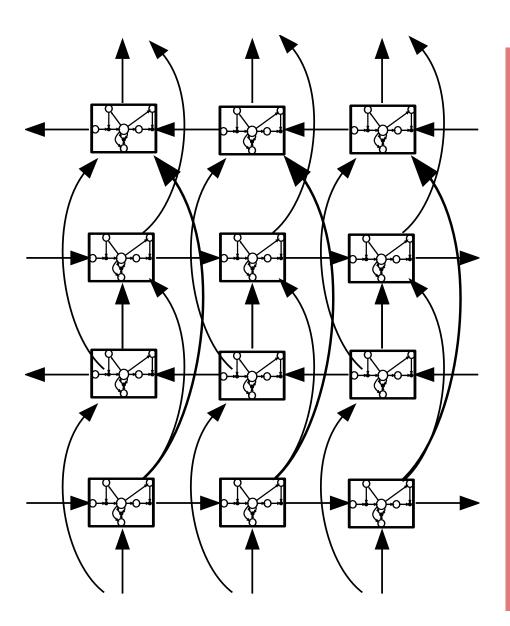
$$c_{t} = f_{t} c_{t-1} + i_{t} \tanh \left(W_{xc} x_{t} + W_{hc} h_{t-1} + b_{c} \right)$$

$$o_{t} = \sigma \left(W_{xo} x_{t} + W_{ho} h_{t-1} + W_{co} c_{t} + b_{o} \right)$$

$$h_{t} = o_{t} \tanh (c_{t})$$
Figure from (Graves et al., 2013)

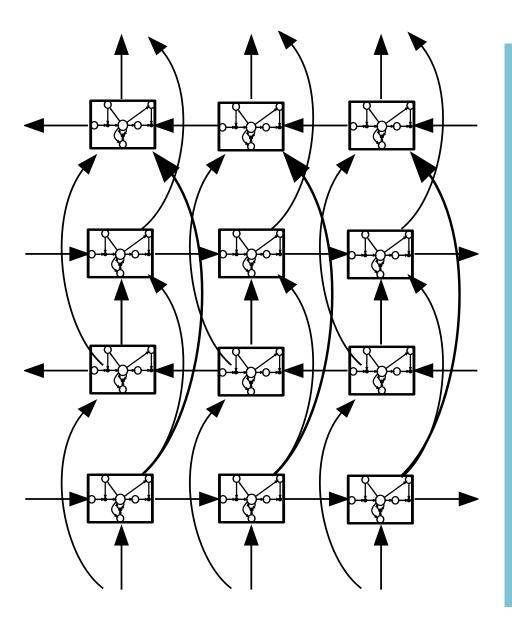


Deep Bidirectional LSTM (DBLSTM)



- Figure: input/output layers not shown
- Same general topology as a Deep Bidirectional RNN, but with LSTM units in the hidden layers
- No additional representational power over DBRNN, but easier to learn in practice

Deep Bidirectional LSTM (DBLSTM)



How important is this particular architecture?

Jozefowicz et al. (2015)
evaluated 10,000
different LSTM-like
architectures and
found several variants
that worked just as
well on several tasks.

RNN Training Tricks

- Deep Learning models tend to consist largely of matrix multiplications
- Training tricks:
 - mini-batching with masking

	Metric	DyC++	DyPy	Chainer	DyC++Seq	Theano	TF
RNNLM (MB=1)	words/sec	190	190	114	494	189	298
RNNLM (MB=4)	words/sec	830	825	295	1510	567	473
RNNLM (MB=16)	words/sec	1820	1880	794	2400	1100	606
RNNLM (MB=64)	words/sec	2440	2470	1340	2820	1260	636

- sorting into buckets of similar-length sequences, so that mini-batches have same length sentences
- truncated BPTT, when sequences are too long, divide sequences into chunks and use the final vector of the previous chunk as the initial vector for the next chunk (but don't backprop from next chunk to previous chunk)

RNN Summary

RNNs

- Applicable to tasks such as sequence labeling, speech recognition, machine translation, etc.
- Able to learn context features for time series data
- Vanishing gradients are still a problem but
 LSTM units can help

Other Resources

 Christopher Olah's blog post on LSTMs <u>http://colah.github.io/posts/2015-08-</u> <u>Understanding-LSTMs/</u>

LEARNING THEORY

PAC-MAN Learning For some hypothesis $h \in \mathcal{H}$:

1. True Error

2. Training Error

$$\hat{R}(h)$$

Question 2:

What is the expected number of PAC-MAN levels Matt will complete before a **Game**-

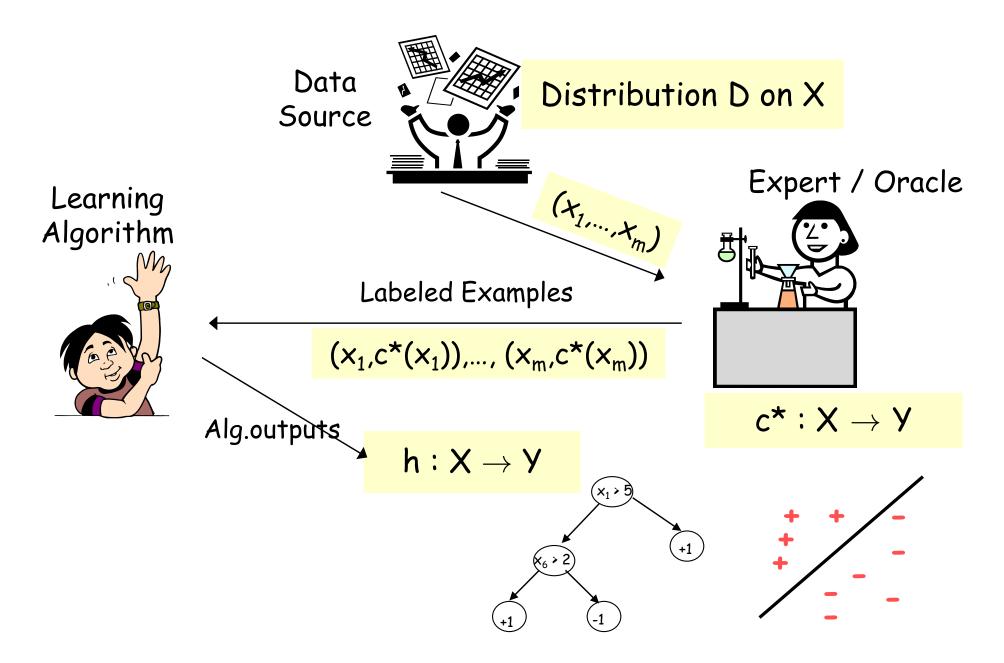
Over?

- A. 1-10
- B. 11-20
- C. 21-30

Questions For Today

- Given a classifier with zero training error, what can we say about true error (aka. generalization error)?
 (Sample Complexity, Realizable Case)
- Given a classifier with low training error, what can we say about true error (aka. generalization error)?
 (Sample Complexity, Agnostic Case)
- 3. Is there a theoretical justification for regularization to avoid overfitting? (Structural Risk Minimization)

PAC/SLT models for Supervised Learning



Two Types of Error

1. True Error (aka. expected risk)

$$R(h) = P_{\mathbf{x} \sim p^*(\mathbf{x})}(c^*(\mathbf{x}) \neq h(\mathbf{x}))$$

2. Train Error (aka. empirical risk)

$$\hat{R}(h) = P_{\mathbf{x} \sim \mathcal{S}}(c^*(\mathbf{x}) \neq h(\mathbf{x}))$$

$$= \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}(c^*(\mathbf{x}^{(i)}) \neq h(\mathbf{x}^{(i)}))$$

$$= \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}(y^{(i)} \neq h(\mathbf{x}^{(i)}))$$

This quantity is always unknown

We can measure this on the training data

where $\mathcal{S} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}_{i=1}^N$ is the training data set, and $\mathbf{x} \sim \mathcal{S}$ denotes that \mathbf{x} is sampled from the empirical distribution.

PAC / SLT Model

We've also referred to this as the "Function View"

1. Generate instances from unknown distribution p^*

$$\mathbf{x}^{(i)} \sim p^*(\mathbf{x}), \, \forall i$$
 (1)

2. Oracle labels each instance with unknown function c^{st}

$$y^{(i)} = c^*(\mathbf{x}^{(i)}), \forall i$$
 (2)

3. Learning algorithm chooses hypothesis $h \in \mathcal{H}$ with low(est) training error, $\hat{R}(h)$

$$\hat{h} = \underset{h}{\operatorname{argmin}} \hat{R}(h) \tag{3}$$

4. Goal: Choose an h with low generalization error R(h)

Three Hypotheses of Interest

The **true function** c^* is the one we are trying to learn and that labeled the training data:

$$y^{(i)} = c^*(\mathbf{x}^{(i)}), \,\forall i \tag{1}$$

The **expected risk minimizer** has lowest true error:

$$h^* = \operatorname*{argmin}_{h \in \mathcal{H}} R(h)$$

Question:

True or False: h* and c* are always equal.

The **empirical risk minimizer** has lowest training error:

$$\hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \, \hat{R}(h) \tag{3}$$

PAC LEARNING

Probably Approximately Correct (PAC) Learning

Whiteboard:

- PAC Criterion
- Meaning of "Probably Approximately Correct"
- Def: PAC Learner
- Sample Complexity
- Consistent Learner

PAC Learning

The **PAC criterion** is that our learner produces a high accuracy learner with high probability:

$$P(|R(h) - \hat{R}(h)| \le \epsilon) \ge 1 - \delta \tag{1}$$

Suppose we have a learner that produces a hypothesis $h \in \mathcal{H}$ given a sample of N training examples. The algorithm is called **consistent** if for every ϵ and δ , there exists a positive number of training examples N such that for any distribution p^* , we have that:

$$P(|R(h) - \hat{R}(h)| > \epsilon) < \delta \tag{2}$$

The **sample complexity** is the minimum value of N for which this statement holds. If N is finite for some learning algorithm, then $\mathcal H$ is said to be **learnable**. If N is a polynomial function of $\frac{1}{\epsilon}$ and $\frac{1}{\delta}$ for some learning algorithm, then $\mathcal H$ is said to be **PAC learnable**.

SAMPLE COMPLEXITY RESULTS

Sample Complexity Results

Definition 0.1. The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

We'll start with the Four Cases we care about... finite case... Realizable Agnosti Finite $|\mathcal{H}|$ Infinite $|\mathcal{H}|$

Generalization and Overfitting

Whiteboard:

- Realizable vs. Agnostic Cases
- Finite vs. Infinite Hypothesis Spaces
- Theorem 1: Realizable Case, Finite |H|
- Proof of Theorem 1

Sample Complexity Results

Definition 0.1. The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

Four Cases we care about...

	Realizable	Agnostic
Finite $ \mathcal{H} $	Thm. 1 $N \geq \frac{1}{\epsilon} \left[\log(\mathcal{H}) + \log(\frac{1}{\delta}) \right]$ labeled examples are sufficient so that with probability $(1-\delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$.	
Infinite $ \mathcal{H} $		