

# Graves 2013, “Generating Sequences with Recurrent Neural Networks”

Johannes Bausch and Jack Kamm

14 November 2017

# Outline

- 1 RNN and LSTM
- 2 Generating text
- 3 Generating handwriting

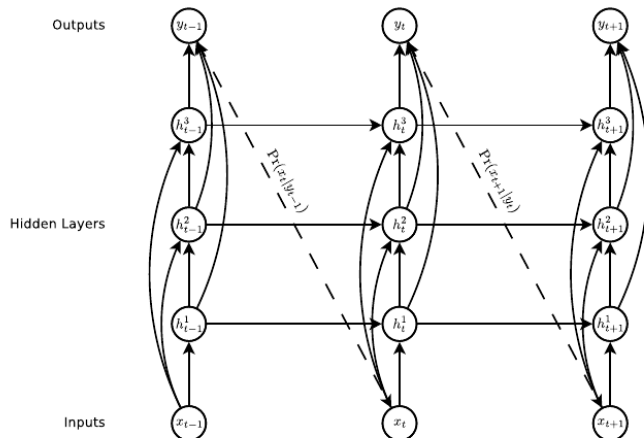
# Table of Contents

1 RNN and LSTM

2 Generating text

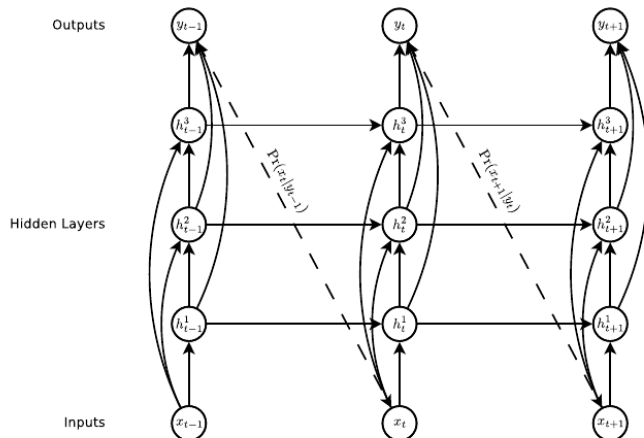
3 Generating handwriting

# Recurrent neural networks



Input  $x_t$ , hidden layers  $h_t^n$ , output  $y_t$ , generative model  $\mathbb{P}(x_{t+1} | y_t)$

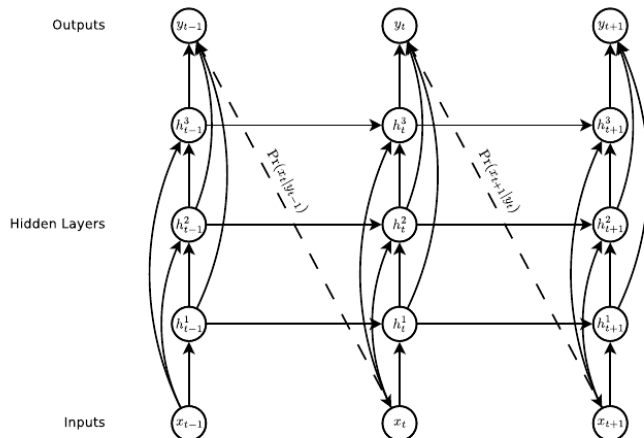
# Recurrent neural networks



$h_t^n$  = nonlinear link  $\circ$  affine combo of  $x_t$ ,  $h_{t-1}^n$ ,  $h_t^{n-1}$

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

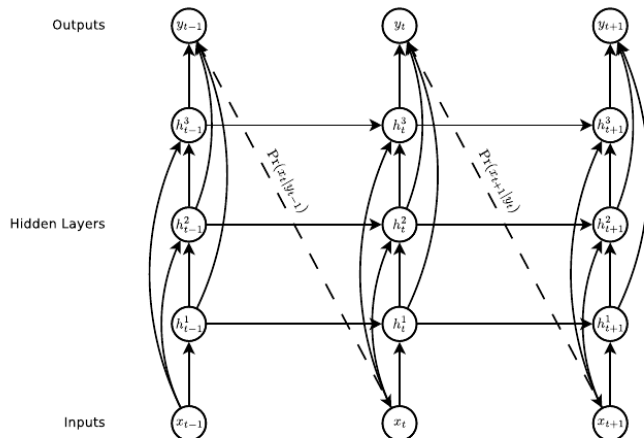
# Recurrent neural networks



$y_t = \text{nonlinear link} \circ \text{affine combo of } h_t^n$

$$y_t = \mathcal{Y}(b_y + \sum_{n=1}^N W_{h^n y} h_t^n)$$

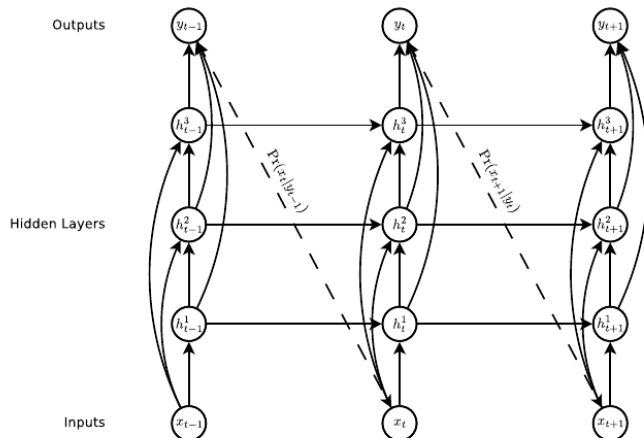
# Recurrent neural networks



Train by maximizing likelihood of generative model:

$$\mathbb{P}(\mathbf{x}) = \prod_{t=1}^T \mathbb{P}(x_t | y_{t-1})$$

# Recurrent neural networks

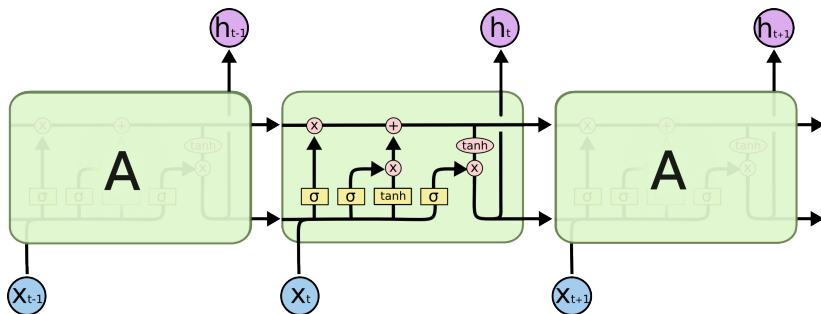


Compute  $\nabla_{\Theta} \log \mathbb{P}_{\Theta}(\mathbf{x})$  by "truncated backpropagation through time"

- i.e., reverse chain-rule + "clip" exploding derivatives

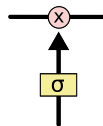


# Long short term memory<sup>1</sup>

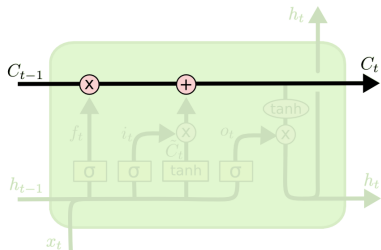


Information passes through a series of “gates”

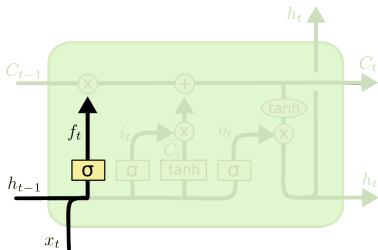
- “Gate” = multiplication with sigmoid
  - $\sigma = 0 \Rightarrow$  “let nothing thru”
  - $\sigma = 1 \Rightarrow$  “let all thru”



<sup>1</sup>graphics have slight differences with Graves 2013; they are taken from:

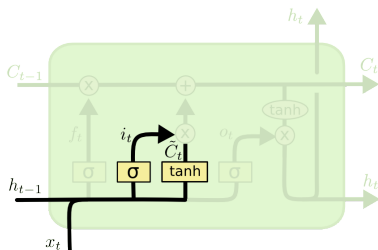


$C_t$  = “cell state” = flows horizontally across LSTM units



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

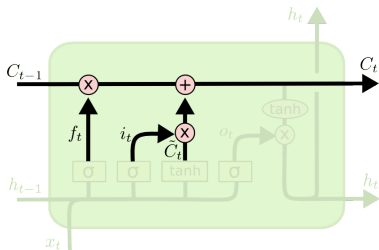
$f_t$  = “forget gate” = gate to forget information from  $C_{t-1}$



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

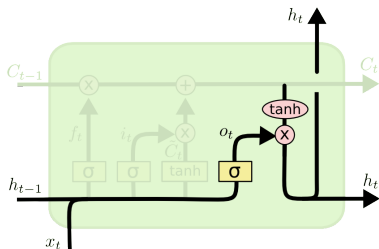
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

$i_t$  = “input gate” = gate to add information from  $h_{t-1}$  and  $x_t$  to  $C_t$



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

Update  $C_t$  using  $f_t$  and  $i_t$

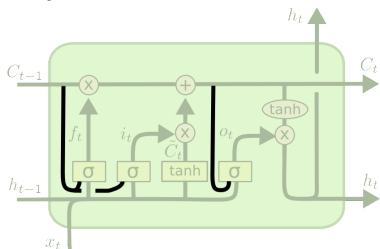


$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$

$o_t$  = “output gate” = gate to output information to cell above/right

Many variations on LSTM exist. Here is the one used in Graves 2013:



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

$$i_t = \sigma(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$

$$o_t = \sigma(W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$

Notice the extra “peephole” connections from  $C$  to  $f, i, o$

# Table of Contents

1 RNN and LSTM

2 Generating text

3 Generating handwriting



# Text Prediction

Basic framework for text prediction

$$y_t = \mathbb{P}(x_{t+1})$$

Per-character vs per-word prediction

# Penn Treebank Test Set

Penn Treebank  
Perplexity? BPC?  
Regularization schemes

# Wikipedia Experiments

- Wikipedia experiments
- `karpathy.github.io/2015/05/21/rnn-effectiveness` has some nice visualizations for understanding what's going on, e.g. what the inner neurons represent

# Table of Contents

- 1 RNN and LSTM
- 2 Generating text
- 3 Generating handwriting

would find the bus safe and sound  
As for Mark, unless it were a  
cancer at the age of fifty-five  
Editorial. Dilemma of  
the the tides in the affairs of men;

Figure: Training samples from IAM handwriting database

Input  $x_t = (\text{pos}_t - \text{pos}_{t-1}) \times (\text{is-end-of-stroke}) \in \mathbb{R}^2 \times \{0, 1\}$

# Mixture density network

$$\mathbf{x}_t \in \mathbb{R} \times \mathbb{R} \times \{0, 1\}$$

$$\mathbf{y}_t = (\mathbf{e}_t, \{\pi_t^j, \mu_t^j, \sigma_t^j, \rho_t^j\}_{j=1}^M)$$

$$\mathbb{P}(\mathbf{x}_{t+1} \mid \mathbf{y}_t) = \sum_{j=1}^M \pi_t^j \mathcal{N}(\mathbf{x}_{t+1} \mid \mu_t^j, \sigma_t^j, \rho_t^j) \begin{cases} \mathbf{e}_t & \text{if } (\mathbf{x}_{t+1})_3 = 1 \\ 1 - \mathbf{e}_t & \text{else} \end{cases}$$

- $\mathbf{e}_t$  = end-of-stroke prob =  $\frac{1}{1 + \exp(\hat{e}_t)}$
- $\pi_t^j$  = mixture prob =  $\frac{\exp(\hat{\pi}_t^j)}{\sum_{j'} \exp(\hat{\pi}_t^{j'})}$
- $\mu_t^j$  = component mean =  $\hat{\mu}_t^j$
- $\sigma_t^j$  = component variance =  $\exp(\hat{\sigma}_t^j)$
- $\rho_t^j$  = component x/y correlation =  $\tanh(\hat{\rho}_t^j)$

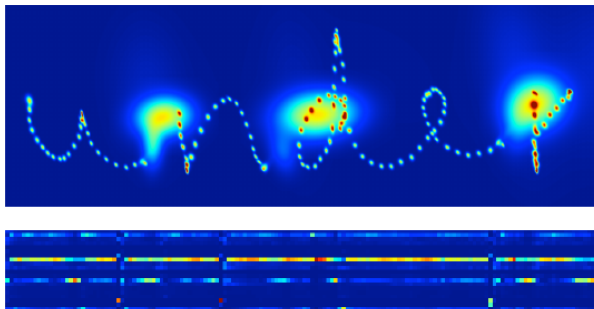


Figure: Mixture density outputs as word “under” is written

- Blobs = Predictions at end of strokes for first point in next stroke
- Lower panel = component weights

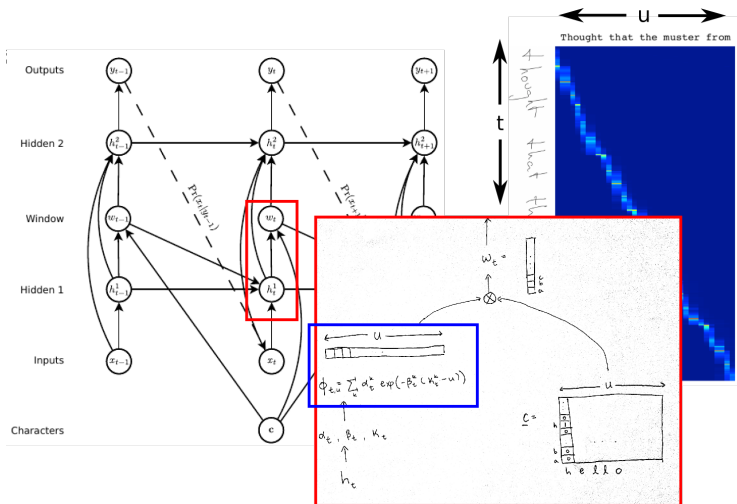




# Handwriting synthesis

- Generate handwriting for a given text
- Main challenge: aligning continuous handwriting sequence with discrete text
- Main idea: bottom hidden layer maintains a distribution of the current position in the text (“soft window”)

# Synthesis Network



**Figure:**  $\phi_t(h_t^1)$  = distn over positions;  $w_t = c\phi_t$  = distn over characters

# Unbiased sampling

from his travels it might have been  
 from his travels it might have been  
 from his travels it might have been  
 from his travels it might have been  
 from his travels it might have been  
 from his travels it might have been

more of national temperament  
 more of national temperament  
 more of national temperament  
 more of national temperament  
 more of national temperament  
 more of national temperament

- Sample from  $\mathbb{P}(\mathbf{x} \mid \mathbf{c})$  by iteratively sampling from  $\mathbb{P}(x_{t+1} \mid y_t)$
- Stop when  $\phi(t, U + 1) > \max_{u \leq U} \phi(t, u)$
- First line of each block is real; subsequent lines generated by network

# Biased sampling

0 when the samples are biased  
0.1 towards more probable sequences  
0.5 they get easier to read  
2 but less diverse  
5 until they all look  
10 exactly the same  
10 exactly the same  
10 exactly the same

- Decreasing the variance of the output sequences leads to neater handwriting

**Figure:** what is the bias term mean here

# Primed sampling