# Sequence Modeling: Recurrent and Recursive Networks

Markus Dumke

27th January 2016

# Contents

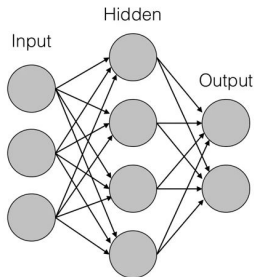# Why RNN's?



Input

Hidden

Output

- Independence
- Fixed Length

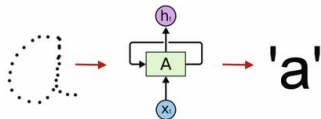https://www.nervanasys.com/recurrent-

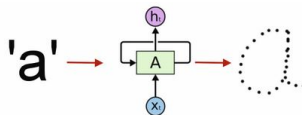neural-networks/

He went to Germany in 2010.

In 2010 he went to Germany.

- sequential data:
  texts, speech, time series
- variable length
- long-term dependencies
- memory

# Applications



Handwriting recognition

Handwriting generation

https://greydanus.github.io/2016/08/21/handwriting/

# Applications



"man in black shirt is playing guitar."

Image Captioning



Smart reply

cs.stanford.edu/people/karpathy/deepimagesent/

https://research.googleblog.com/2016/05/chat-smarter-with-allo.html

# Applications

- Machine translation
- Sentiment analysis
- Text summaries
- Speech recognition and generation
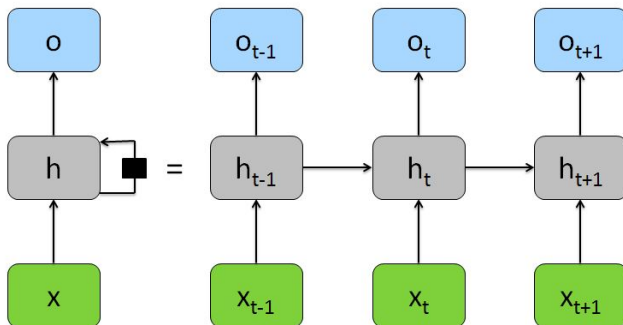- Time series
- Deep Reinforcement Learning

# Contents

# Recurrent Neural Network

# Recurrent Neural Network
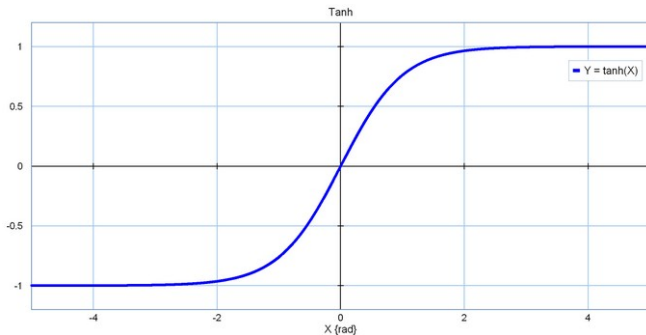


for $t = 1$ to $T$:

$$h_t = f(b + W\,h_{t-1} + U\,x_t)$$
$$o_t = c + V\,h_t$$
$$\hat{y}_t = softmax(o_t)$$
$$= \frac{exp(o_t^{k'})}{\sum_k exp(o_t^k)} \quad \forall k'$$

# Which activation function?

$$f(x) = tanh(x) = \frac{sinh(x)}{cosh(x)} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

# Language Modeling

- Input: word/character encoded as one-hot vector
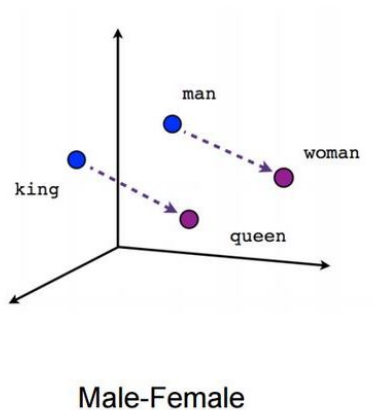- Output: Probability distribution over words given previous words

$$P(y_1, ..., y_T) = \prod_{i=1}^{T} P(y_i | y_1, ..., y_{i-1})$$
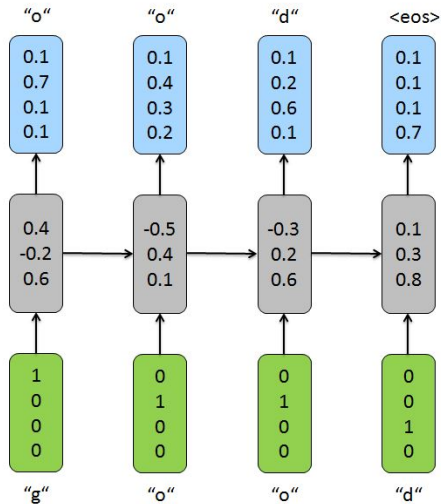
- scoring candidates

# Word embeddings (Word2vec)

- Data sparsity

$$man \rightarrow \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0.35 \\ 0.83 \\ \vdots \\ 0.34 \\ 0.66 \end{bmatrix}$$



Male-Female

https://www.tensorflow.org/tutorials/word2vec/

# Recurrent Neural Network

# Sampling from an RNN

- Sample from conditional distribution at each time step

- How to generate sequence length?
  - special end symbol
  - Bernoulli random variable
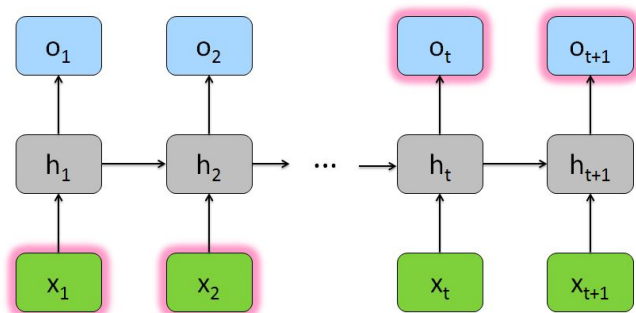  - integer value $\tau$

# Contents

# Optimization

- Forward Propagation:
  - compute hidden states, outputs and loss
  - Loss function, e.g. Bernoulli loss, MSE

$$L = \sum_t L_t$$

- Backward Propagation through time (BPTT):
  - compute gradients

- Stochastic Gradient Descent
  - Minibatch

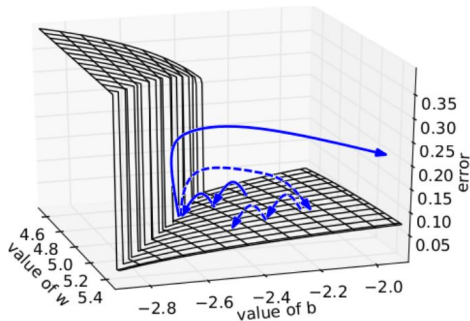# Vanishing (and Exploding) Gradient Problem

# How to deal with exploding gradients?

Gradient Clipping

$$\text{if } ||\nabla W|| > \text{threshold}:$$

$$\nabla W \leftarrow \frac{\text{threshold}}{||\nabla W||} \nabla W$$



http://www.jmlr.org/proceedings/papers/v28/pascanu13.pdf

# How to deal with vanishing gradients?

- Regularization $\nabla_{h_t} L \approx (\nabla_{h_t} L) \dfrac{\partial h_t}{\partial h_{t-1}}$

- skip-connections over time

- Leaky units $\mu = \alpha \, \mu_{t-1} + (1 - \alpha) \, \nu_t$

- remove short-term connections
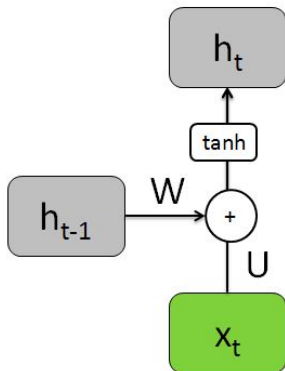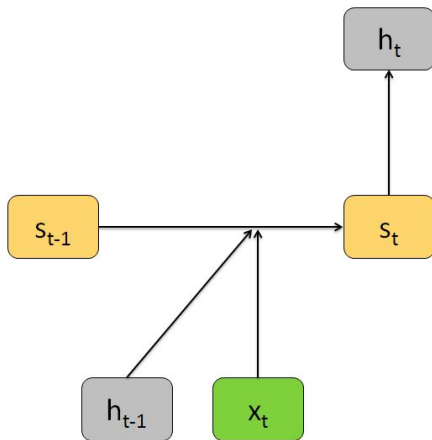
- **LSTM**, GRU and other gated RNNs

# Contents

# Vanilla RNN

$$h_t = \tanh(b + W\, h_{t-1} + U\, x_t)$$

# LSTM

# LSTM

# LSTM

# LSTM



$$\tilde{s}_t = \tanh(b + W\, h_{t-1} + U\, x_t)$$

# LSTM



$$s_t = f_t \, s_{t-1} + i_t \, \tilde{s}_t$$

$$i_t = \sigma \left( b^i + U^i \, x_t + W^i \, h_{t-1} \right)$$

# LSTM



$$h_t = q_t \, \tanh(s_t)$$

$$q_t = \sigma \left( b^q + U^q \, x_t + W^q \, h_{t-1} \right)$$

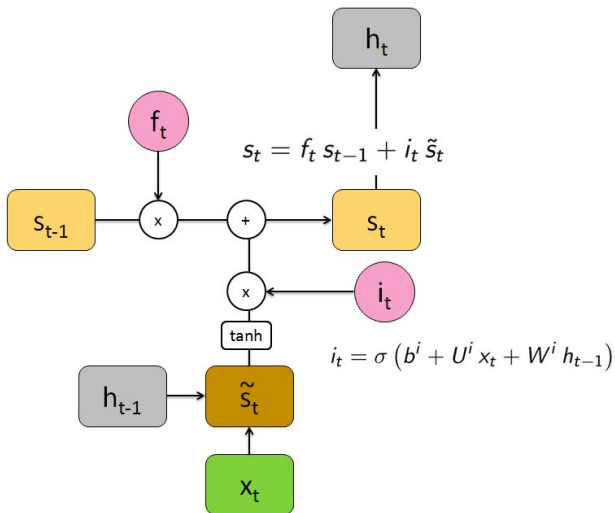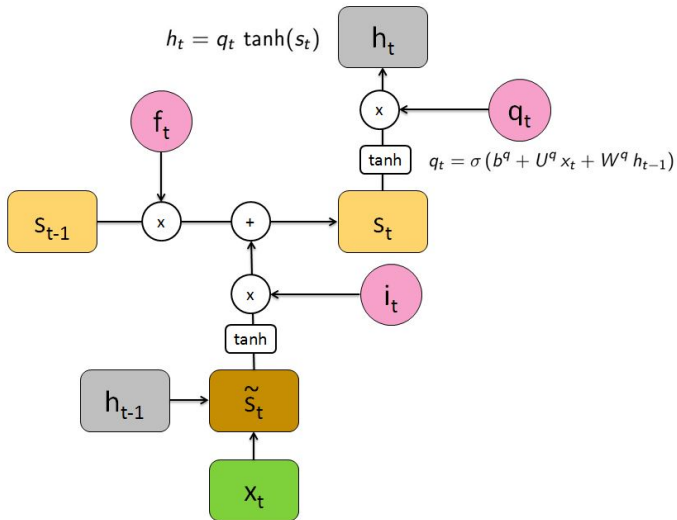# LSTM in R (mxnet)
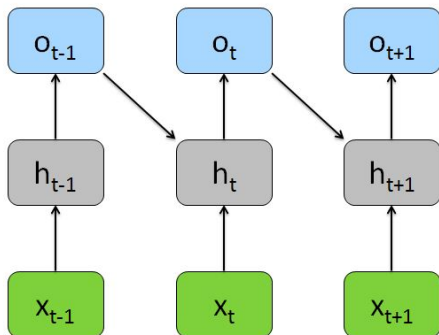
```
 1  model <- mx.lstm(X.train, X.val,
 2                   ctx = mx.cpu(),
 3                   num.round = 100,
 4                   update.period = 1,
 5                   num.lstm.layer = 1,
 6                   seq.len = 32,
 7                   num.hidden = 16,
 8                   num.embed = 16,
 9                   num.label = 100,
10                   batch.size = 32,
11                   input.size = 100,
12                   initializer = mx.init.uniform(0.1),
13                   learning.rate = 0.1,
14                   wd = 0.00001,
15                   clip_gradient = 1)
```

# Generating Text

```
1  infer.model <- mx.lstm.inference(num.lstm.layer=num.lstm.layer,
2                                    input.size=vocab,
3                                    num.hidden=num.hidden,
4                                    num.embed=num.embed,
5                                    num.label=vocab,
6                                    arg.params=model\$arg.params,
7                                    ctx=mx.cpu())
8
9  mx.lstm.forward(infer.model, input, FALSE)
```
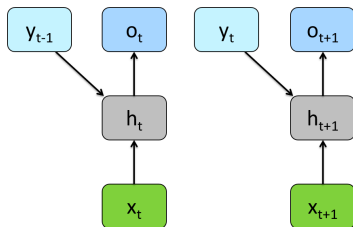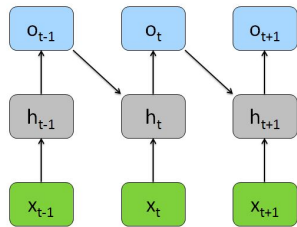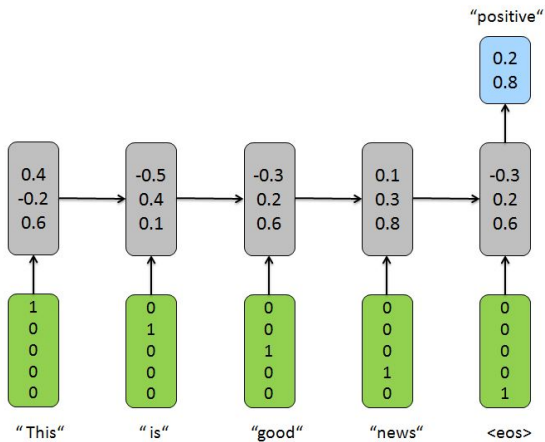
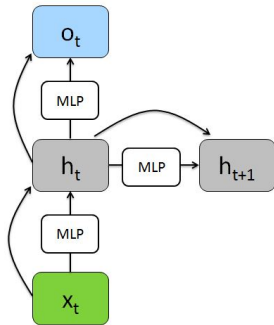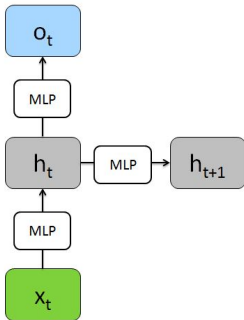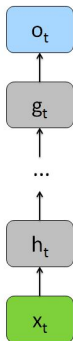# RNN with output recurrence

# Teacher Forcing



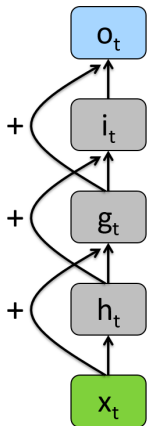At train time

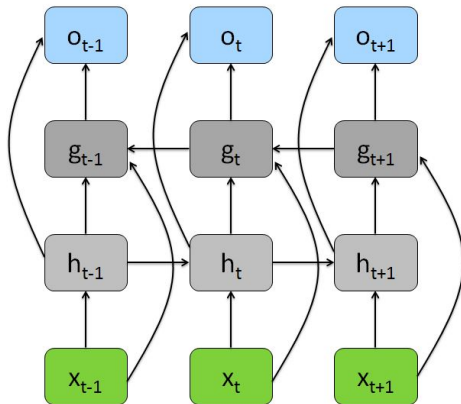At test time

# One-output RNN

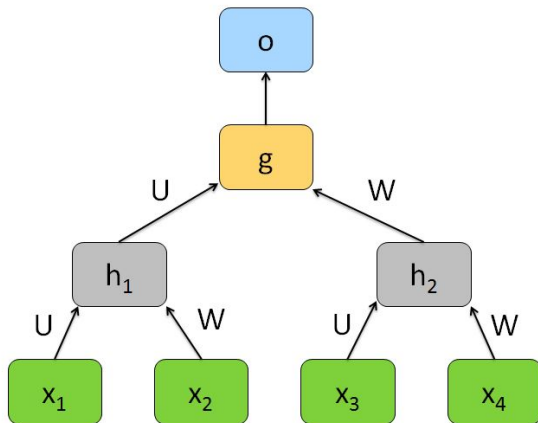# Deep RNNs

# Residual Networks (Res-Nets)



- training of very deep models possible
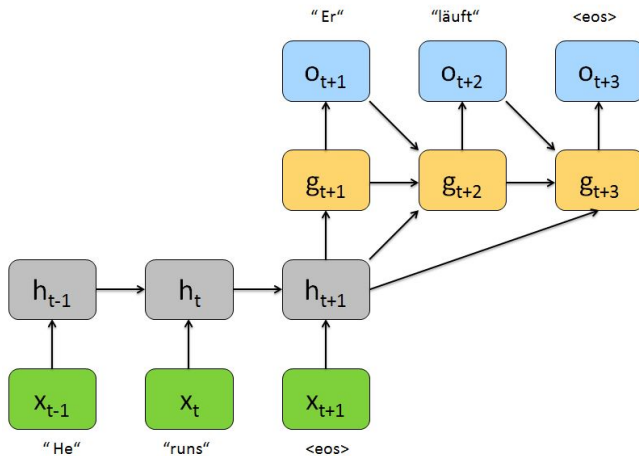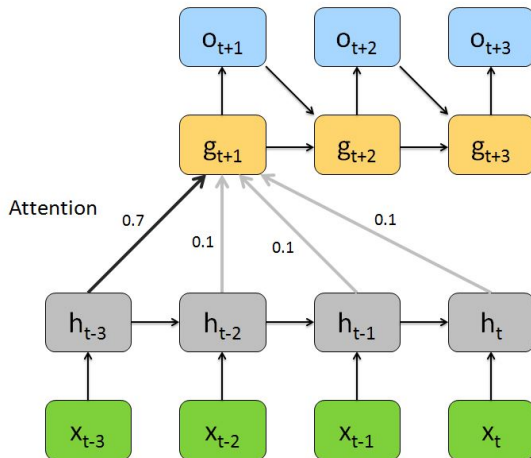- like an ensemble of shallow architectures
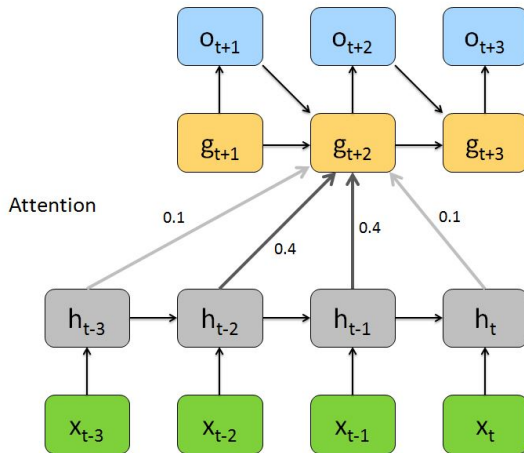
# Bidirectional RNN

# Recursive Neural Network

# Encoder-Decoder Network

# Attention

# Attention

# Contents

# Machine Translation

- last decades: phrase-based systems

- neural networks as part of phrase-based systems

- Encoder-decoder RNNs:
  - Sutskever et al. (2014), Bahdanau et al. (2015)

- Google's Neural Machine Translation (September/November 2016)

# Google's Neural Machine Translation System



Wu et al. (2016): Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation
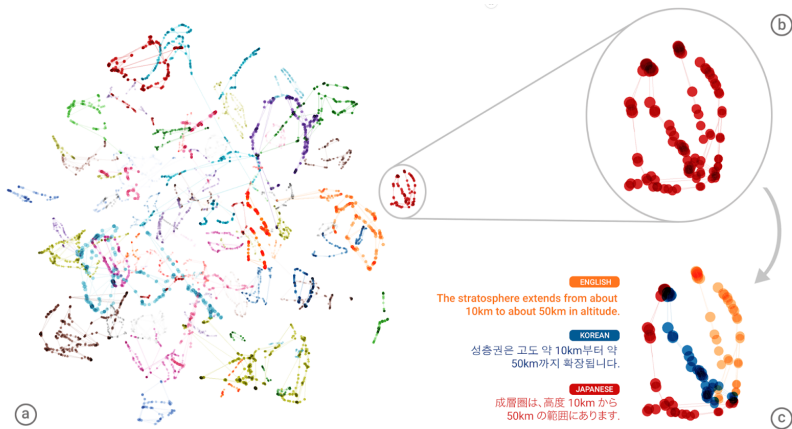
# Details

Main challenges:

- speed
- handling of rare words
- not translating all words (coverage)

Solutions:

- GPU training
- sub-word units (wordpieces)
- coverage penalty
- length-normalization

# Language Embeddings



ENGLISH
The stratosphere extends from about
10km to about 50km in altitude.

KOREAN
성층권은 고도 약 10km부터 약
50km까지 확장됩니다.

JAPANESE
成層圏は、高度 10km から
50km の範囲にあります。

Johnson et al. (2016): Google's Multilingual Neural Machine Translation System: Enabling Zero-Shot Translation