

Recurrent neural nets.

Victor Kitov

v.v.kitov@yandex.ru

Intro

- Sequences
 - words: sequences of symbols
 - sentences: sequences of words
 - documents: sequences of words
- Need fixed vector representation for prediction!
- Bag-of-words allows to do that:
 - one-hot encoding: indicator, TF, TF-IDF models
 - embeddings: get average embedding for sentence/document.

Problem of bag-of-words approach

- Problem: bag-of-words completely ignores word order.
 - information loss!
- Recurrent neural nets account for positions of all elements in sequence!
 - output fixed size sequence representation
 - this feature representation is feature extraction for later model.
 - e.g. MLP.

Recurrent neural net (RNN)

- Consider input sequence $\mathbf{x}_{i:j} := \mathbf{x}_i, \dots \mathbf{x}_j$, $\mathbf{x}_i \in \mathbb{R}^{d_{in}}$.
- RNN outputs single vector $\hat{\mathbf{y}}_n \in \mathbb{R}^{d_{out}}$:

$$\hat{\mathbf{y}}_n = RNN(\mathbf{x}_{1:n})$$

- This implicitly defines RNN^* with sequential output:

$$\hat{\mathbf{y}}_{1:n} = RNN^*(\mathbf{x}_{1:n})$$

$$\hat{\mathbf{y}}_i = RNN(\mathbf{x}_{1:i})$$

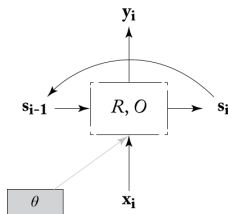
- Comments:
 - RNN shrinks history $\mathbf{x}_{1:n}$ to fixed size vector \mathbf{y}_n .
 - No Markov assumption: all info is aggregated!

Technical details of RNN

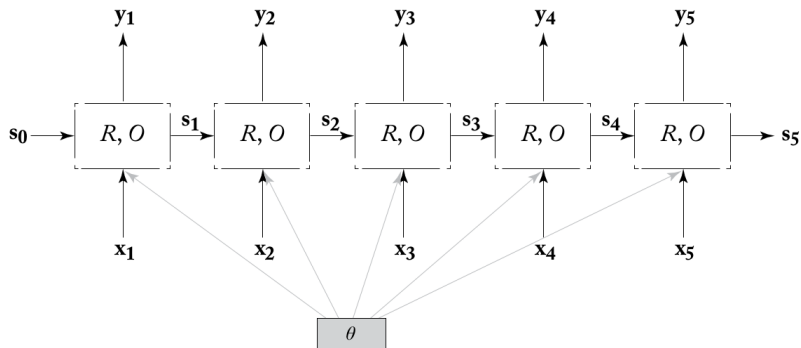
$$\begin{aligned}RNN^*(\mathbf{x}_{1:n}, \mathbf{s}_0) &= \mathbf{y}_{1:n} \\ \hat{\mathbf{y}}_i &= O(\mathbf{s}_i) \\ \mathbf{s}_i &= R(\mathbf{s}_{i-1}, \mathbf{x}_i)\end{aligned}$$

$$\mathbf{x}_i \in \mathbb{R}^{d_{in}}, \mathbf{y}_i \in \mathbb{R}^{d_{out}}, \mathbf{s}_i \in \mathbb{R}^{d_{state}}$$

Typical usage: $O(\mathbf{s}) \equiv \mathbf{s}$, $d_{state} = d_{out}$, $\mathbf{s}_0 = \mathbf{0}$.



Unrolled RNN



$$\begin{aligned}
 s_4 &= R(s_3, x_4) = R(R(s_2, x_3), x_4) \\
 &= R(R(R(s_1, x_2), x_3), x_4) = R(R(R(R(s_0, x_1), x_2), x_3), x_4)
 \end{aligned}$$

Training

Training: unroll RNN and use parameter sharing.

- called **backpropagation through time** (BPTT)
- Variant: unroll RNN for all non-intersecting subsequences of given sequence of given length.

```

init  $s_0$ 

for  $i$  in  $0, 1, \dots, n/k - 1$ :
     $\hat{y}_{ki+1:ki+k} = RNN^*(x_{ki+1:ki+k}, s_{ki})$ 
    calculate loss  $\sum_{j=ki+1}^{ki+k} L(\hat{y}_j, y_j)$ 
    backpropagate gradients, update weights
  
```

Common use-cases of RNN

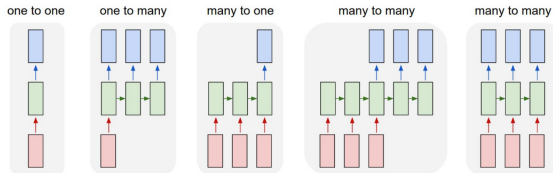
- **Acceptor:** output prediction in $\hat{\mathbf{y}}_n$.
 - e.g. read sentence and output its polarity probabilities.
- **Transducer:** tag sequence x_1, \dots, x_n with RNN outputs y_1, \dots, y_n .
Loss function:

$$\mathcal{L}(\hat{\mathbf{y}}_{1:n}, \mathbf{y}_{1:n}) = \sum_{i=1}^n L(\hat{\mathbf{y}}_i, \mathbf{y}_i)$$

- e.g. POS tagging, language modelling.
 - summarization: for each sentence classify whether to include it into summary or not.
- **Encoder:** encode input sequence representation as $\hat{\mathbf{y}}_n$, e.g.:
 - machine translation: translation generation done with another "decoding" RNN
decode starting from $\mathbf{s}_0 = \hat{\mathbf{y}}_n$.

NN & RNN architectures

NN & RNN architectures:



Examples, where these architectures arise:

- **one to one:** classical classification, image classification.
- **one to many:** image captioning, story generation based on topic.
- **many to one:** text classification, sentiment analysis.
- **many to many:** machine translation, summarization.
- **syncd many to many:** POS tagging, activity detection on video.

Table of Contents

- 1 RNN extensions
- 2 Concrete RNN Architectures
- 3 LSTM model

Bidirectional RNN

- Bidirectional RNN consists of 2 RNNs
 - forward RNN (R^f, O^f) with state $\mathbf{s}_i^f, i = \overline{1, n}$
 - backward RNN (R^b, O^b) with state $\mathbf{s}_i^b, i = \overline{1, n}$
- Forward RNN goes in forward direction $\mathbf{x}_1, \mathbf{x}_2 \dots \mathbf{x}_n$.
- Backward RNN goes in backward direction $\mathbf{x}_n, \mathbf{x}_{n-1} \dots \mathbf{x}_1$.
- At each moment i we have 2 states:
 - 1 $\mathbf{s}_i^f = F_1(\mathbf{x}_1, \mathbf{x}_2 \dots \mathbf{x}_i)$
 - 2 $\mathbf{s}_i^b = F_2(\mathbf{x}_i, \mathbf{x}_{i+1}, \dots \mathbf{x}_n)$

Bidirectional RNN

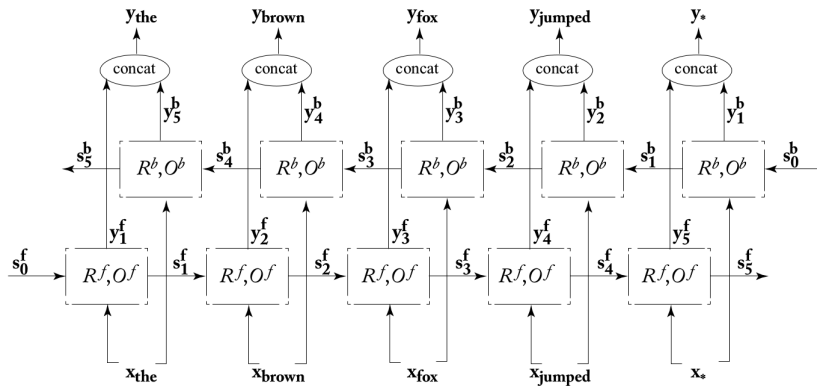
- So we can output

$$biRNN(\mathbf{x}_{1:n}, i) = \hat{\mathbf{y}}_i = [\hat{\mathbf{y}}_i^f; \hat{\mathbf{y}}_i^b] = [RNN^f(\mathbf{x}_{1:i}); RNN^b(\mathbf{x}_{n:i})]$$

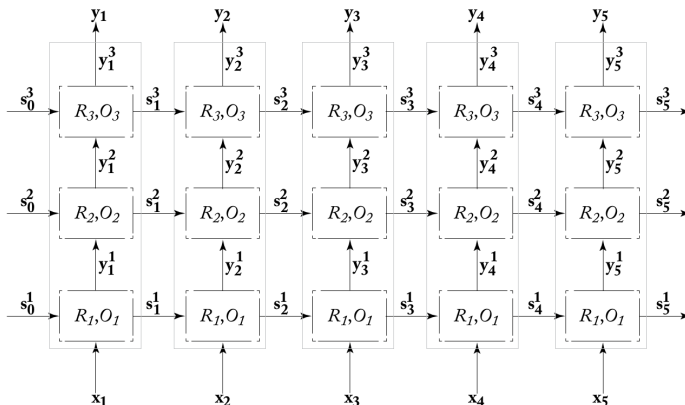
$$biRNN^*(\mathbf{x}_{1:n}) = \mathbf{y}_{1:n} = [biRNN(\mathbf{x}_{1:n}, 1); \dots; biRNN(\mathbf{x}_{1:n}, n)]$$

- encoding takes into account past and future!
- biRNN is very effective for tagging sequences (e.g. POS tagging).

biRNN illustration



Stacked RNN



- Output of previous layer RNN is input to next layer.
- Empirically stacked RNNs work better than single layer RNNs.
- biRNNs can also be stacked.

Table of Contents

- 1 RNN extensions
- 2 Concrete RNN Architectures
- 3 LSTM model

Bag-of-words RNN

Bag-of-words RNN:

$$\mathbf{s}_i = \mathbf{s}_{i-1} + \mathbf{x}_i$$

$$\mathbf{y}_i = \mathbf{s}_i$$

- \mathbf{x}_i : input vector
- \mathbf{s}_i : hidden layer
- \mathbf{y}_i : output vector

Order of words does not matter, not very informative.

Simple RNN (S-RNN)

Simple RNN (S-RNN)¹:

$$\begin{aligned}\mathbf{s}_i &= g_s(W_s \mathbf{s}_{i-1} + V_s \mathbf{x}_i + \mathbf{b}_s) \\ \mathbf{y}_i &= g_y(W_y \mathbf{s}_i + \mathbf{b}_y)\end{aligned}$$

- \mathbf{x}_i : input vector
- \mathbf{s}_i : hidden layer
- \mathbf{y}_i : output vector
- W_s, V_s, W_y : parameter matrices
- $\mathbf{b}_s, \mathbf{b}_y$: parameter vectors
- $g_s(\cdot), g_y(\cdot)$: activation functions

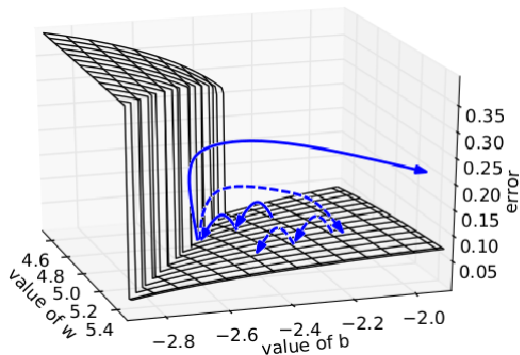
¹also called Elman network

Properties of S-RNN

- S-RNN is sensitive to the order of the inputs.
- Due to recurrent multiplications by W_s is subject to:
 - exploding gradient problem
 - solved by gradient clipping
 - vanishing gradient problem
 - solved by gated models and memory models

Exploding gradient problem

Exploding gradient problem:



Exploding gradient problem

Solutions:

- add regularization
- gradient clipping: clip norm of gradient by threshold.
 - if $\|\nabla_{\theta} L(\hat{\mathbf{y}}_i, \mathbf{y}_i)\| < t$

$$\theta \rightarrow \theta - \varepsilon \nabla_{\theta} L(\hat{\mathbf{y}}_i, \mathbf{y}_i)$$

- else

$$\theta \rightarrow \theta - \varepsilon \frac{t}{\|\nabla_{\theta} L(\hat{\mathbf{y}}_i, \mathbf{y}_i)\|} \nabla_{\theta} L(\hat{\mathbf{y}}_i, \mathbf{y}_i)$$

Vanishing gradients

- Repetitive multiplication of state by the same matrix W_s and saturating non-linearities also cause net to forget past quickly due to vanishing gradients.
- Ways to combat this:
 - 1 Initialize $W_s = I, b_s = 0, g_s = ReLu$.
 - so initially network sums information
 - will change behavior after training if needed
 - 2 Better solution: use LSTM model.

Table of Contents

- 1 RNN extensions
- 2 Concrete RNN Architectures
- 3 LSTM model**

Gates

- Consider n dimensional vectors:
 - old state \mathbf{s} , update \mathbf{x} and new state \mathbf{s}' .
- Gate $g \in \{0, 1\} \in \mathbb{R}^n$ controls state positions where change is applied.
- Example (\odot defines point-wise multiplication):

$$\begin{array}{ccc}
 \begin{bmatrix} 8 \\ 11 \\ 3 \\ 7 \\ 5 \\ 15 \end{bmatrix} & \leftarrow & \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \odot \begin{bmatrix} 10 \\ 11 \\ 12 \\ 13 \\ 14 \\ 15 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \end{bmatrix} \odot \begin{bmatrix} 8 \\ 9 \\ 3 \\ 7 \\ 5 \\ 8 \end{bmatrix} \\
 \mathbf{s}' & & \mathbf{g} \quad \mathbf{x} \quad (1-\mathbf{g}) \quad \mathbf{s}
 \end{array}$$

- Problems:
 - gates need to be learned
 - piece-wise constant gates cannot be optimized.
- Solution: use sigmoid gate $\mathbf{g} = \sigma(f(\mathbf{x}, \mathbf{s}, \theta))$
 - θ : learned parameters
 - f : any differentiable function

Long short-term memory (LSTM) model

$$\mathbf{f}_t = \sigma(W_f \mathbf{x}_t + U_f \mathbf{h}_{t-1} + \mathbf{b}_f)$$

forget gate

$$\mathbf{i}_t = \sigma(W_i \mathbf{x}_t + U_i \mathbf{h}_{t-1} + \mathbf{b}_i)$$

input gate

$$\mathbf{o}_t = \sigma(W_o \mathbf{x}_t + U_o \mathbf{h}_{t-1} + \mathbf{b}_o)$$

output gate

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tanh(W_c \mathbf{x}_t + U_c \mathbf{h}_{t-1} + \mathbf{b}_c)$$

inner state

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$

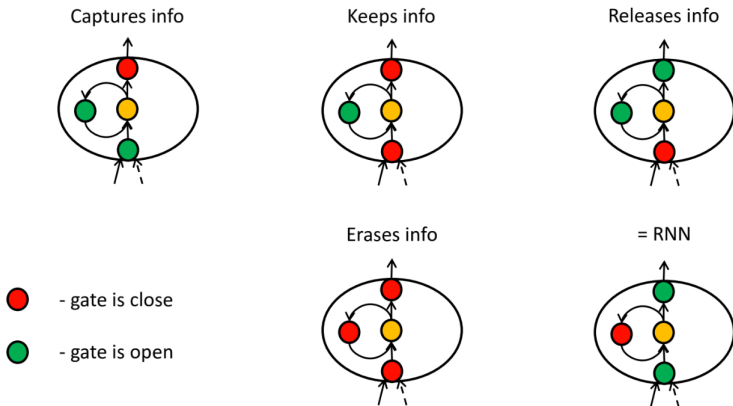
observed output

\mathbf{x}_t -input, parameters:

- matrices: $W_f, U_f, W_i, U_i, W_o, U_o, W_c, U_c$
- vectors: b_f, b_i, b_o, b_c
- initialization: c_0, h_0

Illustration²

Input from below, output above, memory (yellow) in the middle.



²Illustration by Lobacheva Julia.

Comments

- Architecture excluded repetitive multiplication of state by the same matrix W_s (which cause vanishing and exploding gradients)
- Gating mechanisms allow for gradients related to c_t to stay high across long time ranges.
- It's recommended to initialize $\mathbf{b}_f = \mathbf{1}$
 - so initially neural net tries to remember everything