

Рекуррентные нейронные сети

Виктор Китов

v.v.kitov@yandex.ru

Работа с последовательностями

- Данные в виде последовательности
 - динамика цен на акции
 - динамика действий посетителя веб-сайта
 - динамика погоды
 - предложения - последовательности слов
 - речь - последовательность звуков
 - видео - последовательность кадров
- Для текстов: необходимо представление входов небольшой фикс. длины
 - Word2Vec, glove, и др.
- Возможно использование сверточных нейросетей над последовательностями
 - но свертка имеет ограниченную область видимости в историю
- Решение: рекуррентные нейросети (Recurrent neural net, RNN)
 - помнят (в теории) всю историю

Рекуррентные нейросети

- Входная последовательность $\mathbf{x}_{1:n} := \mathbf{x}_1, \dots, \mathbf{x}_n$, $\mathbf{x}_i \in \mathbb{R}^{d_{in}}$.
- RNN выдает вектор фикс. размера $\hat{\mathbf{y}}_n \in \mathbb{R}^{d_{out}}$:

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

- Варьируя n получаем отображения RNN^* из последовательности в последовательность:

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

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

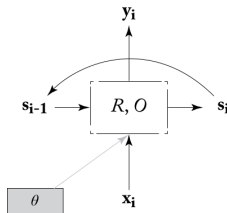
- Т.к. RNN сжимает всю историю $\mathbf{x}_{1:n}$ в вектор фикс. размера \mathbf{y}_n , его можно подавать как вектор признаков др. модели
 - например MLP

Модель рекуррентной нейросети

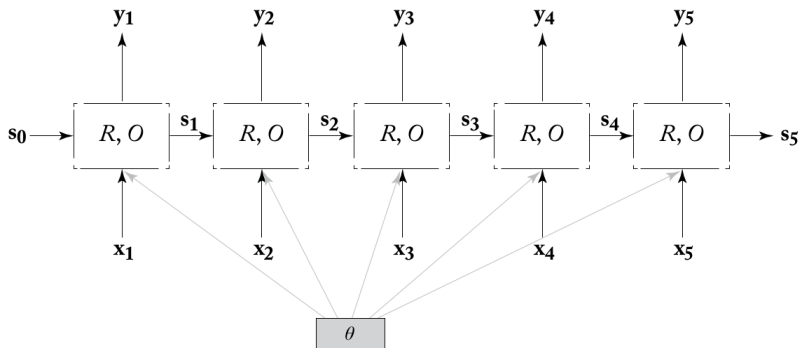
$$\begin{aligned}RNN^*(\mathbf{x}_{1:n}, \mathbf{s}_0) &= \mathbf{y}_{1:n} \\ \hat{\mathbf{y}}_i &= O(\mathbf{s}_i) \\ \mathbf{s}_i &= R(\mathbf{s}_{i-1}, \mathbf{x}_i)\end{aligned}$$

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

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



Развернутая нейросеть (unrolled RNN)



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

Обучение

Обучение: развернуть RNN и использовать неизменность весов для разных t .

- называется **backpropagation through time (BPTT)**
- на практике: развернуть RNN для всех не пересекающихся подпоследовательностей фикс. длины из длинной последовательности.

```

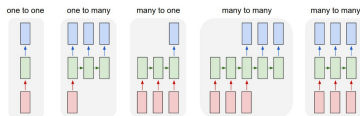
init  $s_0$ 

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
  
```

- Распараллелить вычисления в рамках последовательности нельзя - последующее состояние зависит от предыдущего.

Архитектуры обычной и рекуррентных нейросетей

Архитектуры обычной и рекуррентных нейросетей:



- **one to one:** классическая классификация и регрессия в ML.
- **one to many:** описание изображений (image captioning), генерация текстов по заданной теме.
- **many to one:** классификация текстов, например определение тональности (sentiment analysis).
- **many to many:** машинный перевод, суммаризация длинного текста.
- **syncd many to many:** разметка частей речи, определение событий на видео, распознавание речи.

Варианты применения RNN

- **Acceptor:** выдает итоговый $\hat{\mathbf{y}}_n$.
 - пример: прочитать комментарий и указать его полярность.
- **Encoder:** закодировать последовательность в виде высокоразмерного $\hat{\mathbf{y}}_n$
 - машинный перевод: перевод выполняется декодирующей RNN, стартующей из $\mathbf{s}_0 = \hat{\mathbf{y}}_n$.
 - описание текста вектором фикс. длины $\hat{\mathbf{y}}_n$, на основе которого будет производиться суммаризация текста
- **Transducer:** по x_1, \dots, x_n выдать y_1, \dots, y_n .
 - разметка частей речи, предсказание следующего слова (языковое моделирование), текст- \rightarrow речь, речь- \rightarrow текст
 - Ф-ция потерь:

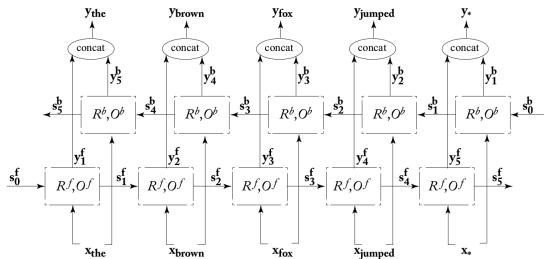
$$\mathcal{L}(\hat{\mathbf{y}}_{1:n}, \mathbf{y}_{1:n}) = \sum_{i=1}^n L(\hat{\mathbf{y}}_i, \mathbf{y}_i)$$

Содержание

- 1 Расширения RNN
- 2 Простая модель и основные проблемы
- 3 Рекуррентные сети с вентилями

Двунаправленная RNN

- Двунаправленная рекуррентная нейросеть (bidirectional RNN) состоит из 2х RNN:
 - forward RNN (R^f, O^f) с состоянием $s_i^f, i = \overline{1, n}$
 - backward RNN (R^b, O^b) с состоянием $s_i^b, i = \overline{1, n}$
- Forward RNN идет слева-направо $x_1, x_2 \dots x_n$.
- Backward RNN идет справа-налево $x_n, x_{n-1} \dots x_1$.



Двунаправленная RNN

- В каждый момент i имеем 2 состояния:

- 1 $s_i^f = F_1(x_1, x_2 \dots x_i)$

- 2 $s_i^b = F_2(x_i, x_{i+1}, \dots x_n)$

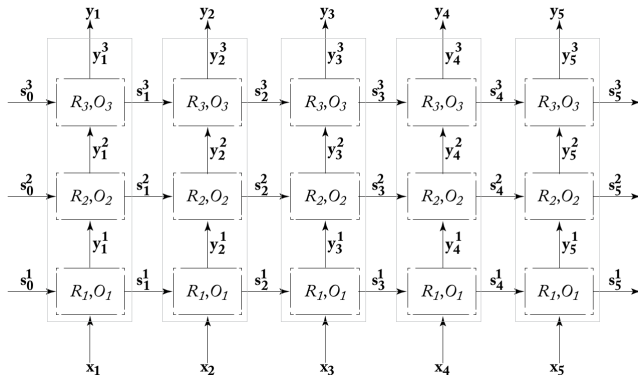
- Для каждого i выдаем их конкатенацию, которую подаем на вход др. модели (MLP)

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

$$biRNN^*(x_{1:n}) = y_{1:n} = [biRNN(x_{1:n}, 1); \dots; biRNN(x_{1:n}, n)]$$

- Двунаправленная RNN эффективно размечает последовательность (режим Transducer), например по тексту разметить части речи.
 - т.к. выход учитывает контекст и справа, и слева.

Многослойная RNN (stacked RNN)



- Выход предыдущего слоя - вход для следующего.
- На практике работают точнее однослойных RNN.
- Можно использовать многослойную двунаправленную RNN.

Содержание

- 1 Расширения RNN
- 2 Простая модель и основные проблемы
- 3 Рекуррентные сети с вентилями

Bag-of-words RNN

Bag-of-words RNN:

$$s_i = s_{i-1} + x_i$$

$$y_i = s_i$$

- x_i : вход
- s_i : скрытый слой
- y_i : выход

Порядок входов не имеет значения, поэтому на практике не применяется.

Сеть Элмана

Сеть Элмана (Elman net, simple RNN):

$$\mathbf{s}_i = g_s (W_s \mathbf{s}_{i-1} + V_s \mathbf{x}_i + \mathbf{b}_s)$$

$$\mathbf{y}_i = g_y (W_y \mathbf{s}_i + \mathbf{b}_y)$$

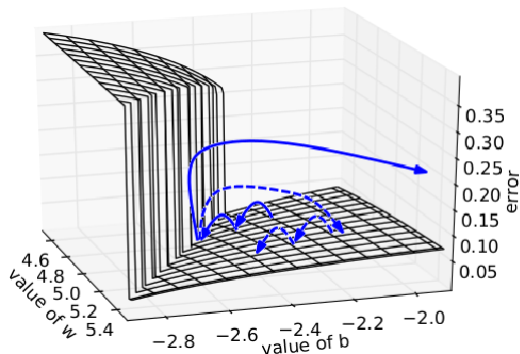
- \mathbf{x}_i : вход
- \mathbf{s}_i : скрытый слой
- \mathbf{y}_i : выход
- W_s, V_s, W_y : матрицы параметров
- $\mathbf{b}_s, \mathbf{b}_y$: векторы параметров
- $g_s(\cdot), g_y(\cdot)$: функции активации

Свойства сети Элмана

- Сеть Элмана чувствительна к порядку входов.
- Зависимость $s_i \leftarrow s_{i-1} \leftarrow s_{i-1} \leftarrow \dots \leftarrow s_0$ через W_s .
- Из-за рекуррентной зависимости от W_s возможны:
 - проблема взрывающегося градиента (exploding gradient) при $|\lambda| \gg 1$
 - проблема затухающего градиента (vanishing gradient) при $|\lambda| \approx 0$

Проблема взрывающегося градиента

Exploding gradient problem:



Проблема взрывающегося градиента

Решение проблемы взрывающегося градиента:

- добавить регуляризацию
- обрезать норму градиента по порогу

$$\text{если } \|\nabla_w L(\hat{y}_i, y_i)\| \leq t : \quad w \rightarrow w - \varepsilon \nabla_w L(\hat{y}_i, y_i)$$

$$\text{если } \|\nabla_w L(\hat{y}_i, y_i)\| > t : \quad w \rightarrow w - \varepsilon \frac{t}{\|\nabla_w L(\hat{y}_i, y_i)\|} \nabla_w L(\hat{y}_i, y_i)$$

Проблема затухающего градиента

- Инициализировать $W_s = I, b_s = 0, g_s = \text{ReLU}, \text{LeakyReLU}$.
- Добавить регуляризацию на несильное отклонение от ортогональности

$$\left\| W_s^T W_s - I \right\|_F^2 \text{ либо } \left\| W_s^T W_s - I \right\|_F^2 + \left\| W_s W_s^T - I \right\|_F^2$$

т.к. у ортогональных матриц все $|\lambda_i| = 1^1$.

- Разделить состояние на быстро и медленно меняющиеся компоненты: $s_t = [s_t^{slow}; s_t^{fast}]$
 - $s_t^{slow} = \alpha s_{t-1}^{slow} + (1 - \alpha) W_{slow} x_t, \alpha \lesssim 1$
 - $s_t^{fast} = \sigma(V_x x_t + V_{fast} s_{t-1}^{fast} + V_{slow} s_t^{slow})$

¹Докажите.

Содержание

- 1 Расширения RNN
- 2 Простая модель и основные проблемы
- 3 Рекуррентные сети с вентилями

Проблема обычной RNN

- Проблема обычной RNN: вся память перезаписывается на каждой итерации.
 - поэтому RNN быстро забывает прошлое (или градиент взрывается)
- Важно помнить прошлое глубоко в истории:
 - автоматические ответы на вопросы
 - машинный перевод
 - суммаризация текстов

Проблема обычной RNN

- Проблема обычной RNN: вся память перезаписывается на каждой итерации.
 - поэтому RNN быстро забывает прошлое (или градиент взрывается)
- Важно помнить прошлое глубоко в истории:
 - автоматические ответы на вопросы
 - машинный перевод
 - суммаризация текстов
- Решение: использование вентиляей



Вентили

- Пусть s -старое состояние, x -новый вход, s' -новое состояние, $s, s', x \in \mathbb{R}^n$.
- Вентиль $g \in \{0, 1\} \in \mathbb{R}^n$ контролирует позиции, которые нужно обновить:
- Пример (\odot - поэлементное умножение):

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

- Какие вентили перекрывать? нужно их настраивать.
- Кусочно-постоянный вентиль не сможем оптимизировать.
 - поэтому используем гладкий вентиль $g = \sigma(f(x, s, \theta))$
 - θ : настраиваемые параметры
 - f : любая гладкая функция

Популярные RNN с вентилями

Популярные RNN с вентилями:

- Long short-term memory (LSTM)
 - "сеть долгой кратковременной памяти"
- Gated recurrent unit (GRU)

RNN с вентилями работают лучше на больших данных чем простые RNN.

Long short-term memory (LSTM)

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

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

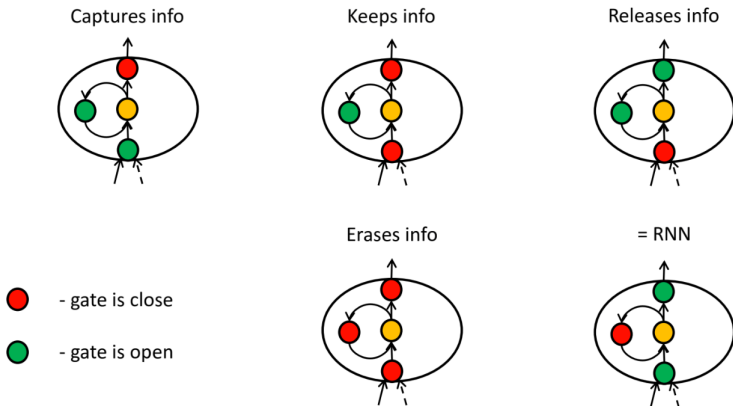
\mathbf{x}_t -ВХОД, \mathbf{h}_t -ВЫХОД.

Параметры:

- матрицы: $W_f, U_f, W_i, U_i, W_o, U_o, W_c, U_c$
- вектора: b_f, b_i, b_o, b_c
- инициализация: c_0, h_0

Иллюстрация

Вход-внизу, память-в середине, выход-наверху.



Комментарии

- $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \dots$ обеспечивает более долгую память (но все равно не бесконечную - см. memory networks)
- Рекомендуется инициализация $\mathbf{b}_f \geq 1$
 - вначале сеть пытается запоминать все

Gated recurrent unit (GRU)

$$\mathbf{z}_t = \sigma(W_z \mathbf{x}_t + U_z \mathbf{h}_{t-1} + \mathbf{b}_z)$$

$$\mathbf{r}_t = \sigma(W_r \mathbf{x}_t + U_r \mathbf{h}_{t-1} + \mathbf{b}_r)$$

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tanh(W_h \mathbf{x}_t + U_h (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_h)$$

\mathbf{x}_t - вход, \mathbf{z}_t - forget gate, \mathbf{r}_t - input gate, \mathbf{h}_t - выход.

Параметры:

- матрицы: $W_z, U_z, W_r, U_r, W_h, U_h$
- вектора: $\mathbf{b}_z, \mathbf{b}_r, \mathbf{b}_h$
- инициализация: \mathbf{h}_0

Комментарии

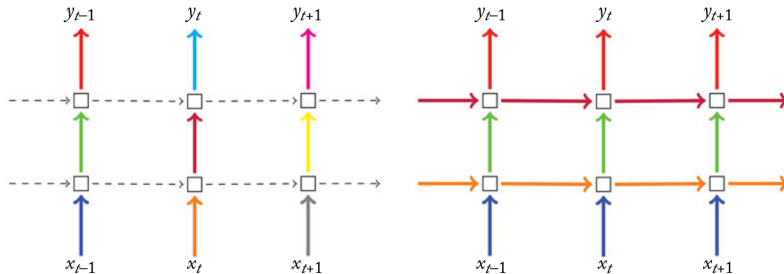
- По сравнению с LSTM GRU:
 - имеет 2, а не 3 вентиля
 - память и выход - то же самое
- GRU содержит меньше параметров и меньше переобучается.

Dropout в RNN

- Проблема: dropout во времени нарушает временные связи.
- Решения
 - 1 (Pham et al. [2013]): применять dropout к нереккуррентным связям.
 - 2 (Gal [2015], variational RNN dropout): для каждой последовательности $x_{1:T}$ регенерировать и фиксировать маску.
 - маска не меняется от времени.
 - маска регенерируется для другой последовательности.

Иллюстрация вариантов dropout в RNN

Иллюстрация dropout в решениях 1 и 2.



Заключение

- Рекуррентные нейросети - удобная модель для обработки и генерации последовательных данных.
- Специфические проблемы настройки рекуррентных нейросетей:
 - взрывающийся градиент
 - решение: обрезка градиента по норме
 - затухающий градиент
 - решение: использование архитектур с вентилями.