

# Обработка последовательностей и модели внимания

К. В. Воронцов  
vokov@forecsys.ru

Этот курс доступен на странице вики-ресурса  
<http://www.MachineLearning.ru/wiki>  
«Машинное обучение (курс лекций, К.В.Воронцов)»

20 октября 2020 • ШАД Яндекс

- 1 Задачи обработки последовательностей**
  - Рекуррентная сеть
  - Рекуррентная сеть с моделью внимания
  - Прикладные задачи
- 2 Разновидности моделей внимания**
  - Разновидности функций сравнения
  - Многомерное внимание (multi-head attention)
  - Внимание к себе ;) (self-attention)
- 3 Модели внимания на графах**
  - Модель внимания GAT
  - Многомерное обобщение GAT

## Напоминание. Рекуррентная сеть (RNN)

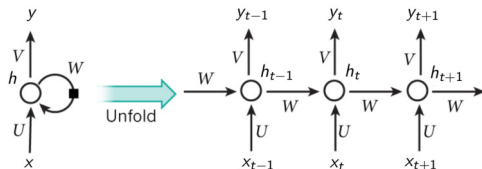
$x_t$  — входной вектор в момент  $t = 1, \dots, T$

$y_t$  — выходной вектор (в некоторых приложениях  $y_t \equiv h_t$ )

$h_t$  — вектор скрытого состояния в момент  $t$

$$h_t = \sigma_h(Ux_t + Wh_{t-1})$$

$$y_t = \sigma_y(Vh_t)$$



Обучение рекуррентной сети:  $\sum_{t=0}^T \mathcal{L}_t(U, V, W) \rightarrow \min_{U, V, W}$

- длины входного и выходного сигнала обязаны совпадать
- невозможно заглядывание вперёд
- не подходит для многих задач (MT, QA и др.)

## Рекуррентная сеть для обработки последовательностей (seq2seq)

$\{x_i: i = 1, \dots, n\}$  — входная последовательность

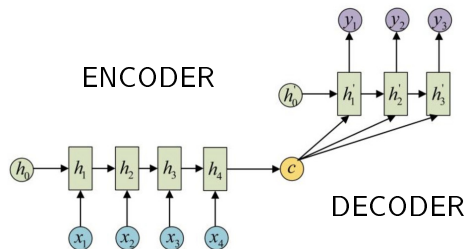
$\{y_t: t = 1, \dots, m\}$  — выходная последовательность

$c \equiv h_n$  кодирует всю информацию про  $\{x_i\}$  для синтеза  $\{y_t\}$

$$h_i = f_{in}(x_i, h_{i-1})$$

$$h'_t = f_{out}(h'_{t-1}, y_{t-1}, c)$$

$$y_t = f_y(h'_t, y_{t-1})$$



- $h_n$  лучше помнит конец последовательности, чем начало
- чем больше  $n$ , тем труднее упаковать всю информацию в  $c$
- придётся контролировать затухание/взрывы градиента
- RNN трудно распараллеливается

## Рекуррентная сеть с вниманием (attention mechanism)

$a(h_i, h'_t)$  — функция сходства состояний входа  $i$  и выхода  $t$

$\alpha_{ti}$  — важность входа  $i$  для выхода  $t$  (attention score),  $\sum_i \alpha_{ti} = 1$

$c_t$  — вектор входного контекста для выхода  $t$  (context vector)

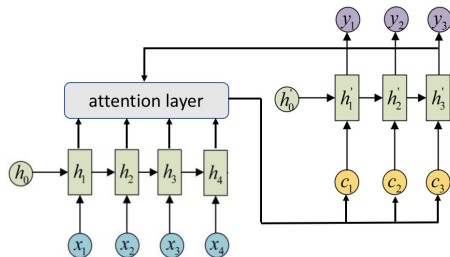
$$h_i = f_{in}(x_i, h_{i-1})$$

$$h'_t = f_{out}(h'_{t-1}, y_{t-1})$$

$$\alpha_{ti} = \frac{a(h_i, h'_t)}{\sum_{k=1}^n a(h_k, h'_t)}$$

$$c_t = \sum_i \alpha_{ti} h_i$$

$$y_t = f_y(h'_t, y_{t-1}, c_t)$$



- можно отказаться от рекуррентности как по  $h_i$ , так и по  $h'_t$
- можно вводить обучаемые параметры в  $a$  и  $c$

## Применения моделей внимания

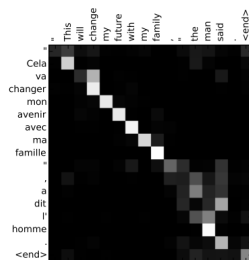
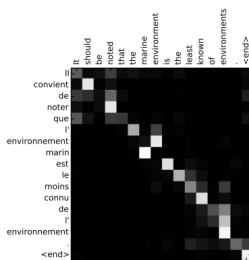
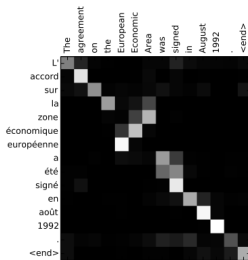
Преобразование одной последовательности в другую, seq2seq:

- Машинный перевод (machine translation)
- Ответы на вопросы (question answering)
- Суммаризация текста (text summarization)
- Описание изображений, аудио, видео (multimedia description)
- Распознавание речи (speech recognition)
- Синтез речи (speech synthesis)

Обработка последовательности:

- Классификация текстовых документов
- Анализ тональности документа / предложений / аспектов

## Применения моделей внимания в машинном переводе



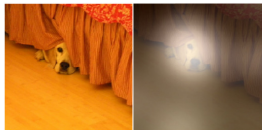
Интерпретируемость моделей внимания:

При обработке конкретной последовательности  $x$  визуализация матрицы  $\alpha_{tj}$  показывает, на какие слова  $x_j$  модель обращает внимание, генерируя слово перевода  $y_t$

## Применения моделей внимания в описании изображений



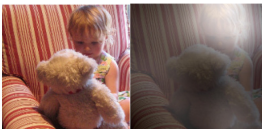
A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

При генерации каждого слова в описании изображения визуализация показывает, на какие области изображения модель обращает внимание, генерируя данное слово

---

*Kelvin Xu et al.* Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. 2016



## Разновидности функций сравнения

$a(h_i, h'_t) = h_i^T h'_t$  — скалярное произведение

$a(h_i, h'_t) = \exp(h_i^T h'_t)$  — тогда norm превращается в SoftMax

$a(h_i, h'_t) = h_i^T W h'_t$  — обобщение, с матрицей параметров  $W$

$a(h_i, h'_t) = w^T \text{th}(U h_i + V h'_t)$  — аддитивное внимание ( $w, U, V$ )

Обобщение с тремя матрицами Query, Key, Value:

$a(h_i, h'_t) = (K h_i)^T (Q h'_t)$

$\alpha_{ti} = \text{norm}_j a(h_i, h'_t)$

$c_t = \sum_j \alpha_{tj} V h_j$

$y_t = f_y(h'_t, y_{t-1}, c_t)$

Возможно упрощение:  $K \equiv V$

Возможно преобразование размерности:  $K, V \in \mathbb{R}^{\dim(h') \times \dim(h)}$

---

*Vaswani et al.* Attention is all you need. 2017.

*Dichao Hu.* An Introductory Survey on Attention Mechanisms in NLP Problems. 2018.

*Sneha Chaudhari et al.* An Attentive Survey of Attention Models. 2019.

## Многомерное внимание (multi-head attention)

**Идея:** несколько разных моделей совместно обучаются  
обращать внимание на разные аспекты входной информации

Вычисляется  $K$  функций сходства вершин  $i, t$ :

$$a^k(h_i, h'_t) = h_i^\top W^k h'_t, \quad k = 1, \dots, K$$

$$\alpha_{ti} = \text{norm}_i a^k(h_i, h'_t)$$

$$c_t^k = \sum_i \alpha_{ti} V^k h_i$$

Два варианта агрегирования выходного вектора:

$$c_t = \text{concat} [c_t^k]_{k=1}^K \text{ — конкатенация}$$

$$c_t = \frac{1}{K} \sum_{k=1}^K c_t^k \text{ — усреднение}$$

$$y_t = f_y(h'_t, y_{t-1}, c_t) \text{ — предсказание по агрегированному вектору}$$

---

*Vaswani et al.* Attention is all you need. 2017.

*Dichao Hu.* An Introductory Survey on Attention Mechanisms in NLP Problems. 2018.

*Sneha Chaudhari et al.* An Attentive Survey of Attention Models. 2019.

## Self-attention для обработки одной последовательности

$\{x_t: t = 1, \dots, n\}$  — входная последовательность токенов

$\{y_t: t = 1, \dots, n\}$  — выходная последовательность

### Идея:

модель обращает внимание на схожие токены из контекста;  
не столь важно, генерируется новая последовательность или  
генерируются новые эмбединги исходной последовательности

Теперь  $h_i$  и  $h_t$  — эмбединги из одной последовательности

$\alpha_{ti} = \text{norm}_i a(h_i, h_t)$  — важность токена  $i$  в контексте токена  $t$

$c_t = \sum_j \alpha_{tj} \mathbf{V} h_j$  — эмбединг контекста токена  $t$  с обучаемым  $\mathbf{V}$

$y_t = f_y(h_t, y_{t-1}, c_t)$  — предсказание для токена  $t$

---

*Vaswani et al.* Attention is all you need. 2017.

*Dichao Hu.* An Introductory Survey on Attention Mechanisms in NLP Problems. 2018.

*Sneha Chaudhari et al.* An Attentive Survey of Attention Models. 2019.

## Модель внимания Graph Attention Network (GAT)

Дано: граф  $\langle V, E \rangle$

$h_i, i \in V$  — входные векторы признаков (или эмбединги) вершин

$h'_i, i \in V$  — выходные векторы вершин

$\mathcal{N}(t)$  — множество вершин  $i \in V$  в окрестности вершины  $t$

Функция сходства вершин  $i, t$  с параметрами  $u, v, W$ :

$$a(h_i, h_t) = \exp(\text{LeakyReLU}(uWh_i + vWh_t))$$

$\alpha_{ti} = \text{norm}_{i \in \mathcal{N}(t)} a(h_i, h_t)$  — важность вершины  $i$  в контексте  $t$

$c_t = \sum_{i \in \mathcal{N}(t)} \alpha_{ti} Wh_i$  — эмбединг контекста вершины  $t$

$h'_t = \sigma(c_t)$  — выходной вектор для вершины  $t$

Функция потерь определяется решаемой на графе задачей.

## Многомерное обобщение Multi-Head Attention для GAT

Дано: граф  $\langle V, E \rangle$

$h_i, i \in V$  — входные векторы признаков (или эмбединги) вершин

$h'_i, i \in V$  — выходные векторы вершин

$\mathcal{N}(t)$  — множество вершин  $i \in V$  в окрестности вершины  $t$

$K$  функций сходства вершин  $i, t$  с параметрами  $u^k, v^k, W^k$ :

$$a(h_i, h_t) = \exp(\text{LeakyReLU}(u^k W^k h_i + v^k W^k h_t))$$

$$\alpha_{ti} = \text{norm}_{i \in \mathcal{N}(t)} a(h_i, h_t) \text{ — важность вершины } i \text{ в контексте } t$$

$$c_t^k = \sum_{i \in \mathcal{N}(t)} \alpha_{ti} W^k h_i \text{ — эмбединг контекста вершины } t$$

Два варианта выходного вектора для вершины  $t$ :

$$h'_t = \text{concat}[\sigma(c_t^k)]_{k=1}^K \text{ — конкатенация}$$

$$h'_t = \sigma\left(\frac{1}{K} \sum_{k=1}^K c_t^k\right) \text{ — усреднение}$$

---

*Petar Veličković et al.* Graph Attention Networks. ICLR-2018.

## Многомерное обобщение Multi-Head Attention для GAT

Дано: граф  $\langle V, E \rangle$

$h_i, i \in V$  — входные векторы признаков (или эмбединги) вершин

$h'_i, i \in V$  — выходные векторы вершин

Пример.  $K = 3$  моделей внимания для преобразования  $h_1 \rightarrow h'_1$

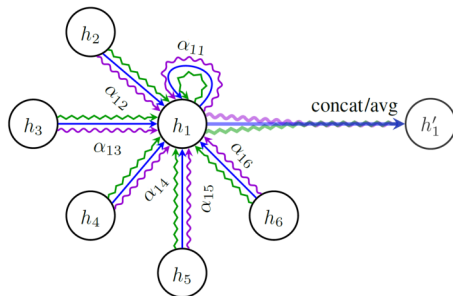
$$\alpha_{ti} = \text{norm}_{i \in \mathcal{N}(t)} a(h_i, h_t)$$

$$c_t^k = \sum_{i \in \mathcal{N}(t)} \alpha_{ti} W^k h_i$$

concat / average:

$$h'_t = \text{concat}[\sigma(c_t^k)]$$

$$h'_t = \sigma\left(\frac{1}{K} \sum_{k=1}^K c_t^k\right)$$



- Модели внимания сначала встраивались в RNN или CNN, но оказалось, что они самодостаточны
- Модель внимания работает точнее и быстрее RNN
- Легко обобщается на тексты, графы, изображения
- Доказано, что модель внимания multi-head self-attention (MHSA) эквивалентна свёрточной сети [Cordonnier, 2020]
- Модель внимания используются в наиболее продвинутых нейросетевых моделях BERT, GPT-2/3

---

*Vaswani et al.* Attention is all you need. 2017.

*Dichao Hu.* An Introductory Survey on Attention Mechanisms in NLP Problems. 2018.

*Sneha Chaudhari et al.* An Attentive Survey of Attention Models. 2019.

*Cordonnier et al.* On the relationship between self-attention and convolutional layers. 2020