

Глубокие нейронные сети

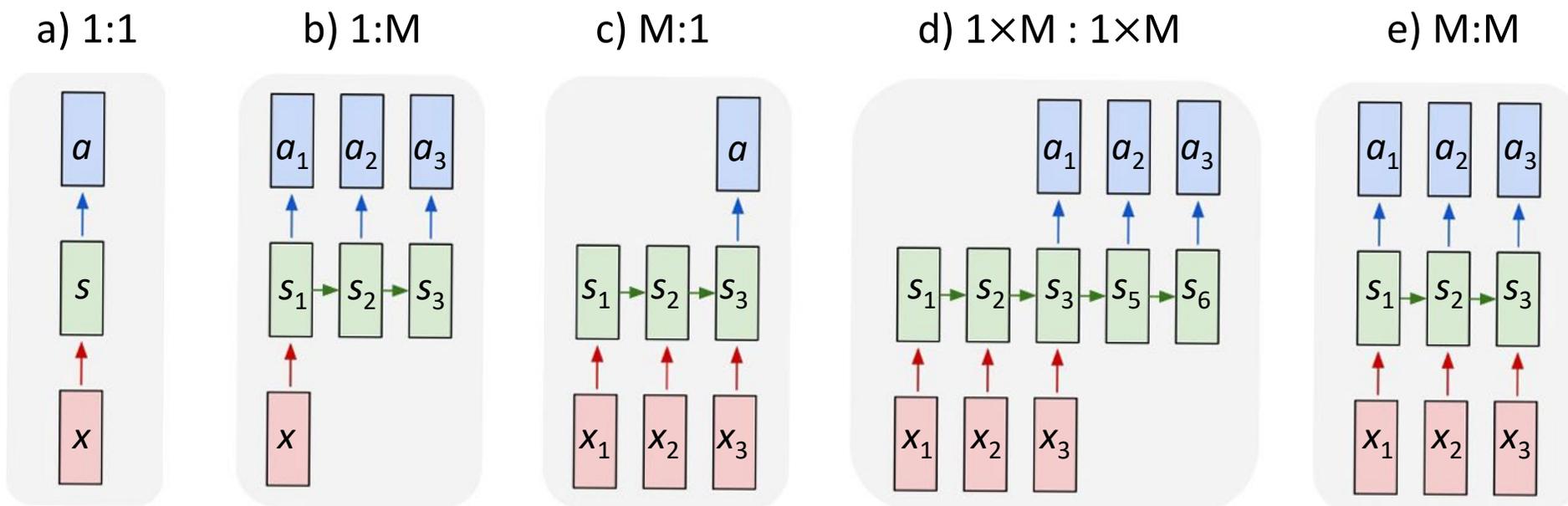
Рекуррентные нейронные сети (Recurrent neural network - RNN)

Лекция 10

Отличие RNN от CNN

- CNN ориентированы на анализ изображений, являющихся моментальными снимками состояний группы объектов в некоторый момент времени (нет анализа предыстории)
- RNN ориентированы на анализ временной последовательности состояний группы объектов

Типы приложений



- a) Один вход, один выход без запоминания состояния (распознавание изображений)
 b) Последовательность на выходе (вход – картинка, выход – ее словесное описание)
 c) Последовательность на входе (вход – рецензия на фильм, выход – оценка по десятибалльной шкале)
 d) Последовательность на входе, последовательность на выходе (перевод текста с одного языка на другой)
 e) Синхронизированные последовательности на входе и выходе (разметка смены кадра на видео; очистка аудиопотока от посторонних шумов)



Область применения рекуррентных нейронных сетей

- Анализ временных рядов
 - Изменения цен акций
 - Показания датчиков
- Понимание (письменных) текстов на естественном языке
- Машинный перевод
- Понимание человеческой речи (перевод речи в текст)
- Обнаружение аномалий (определение предынфарктного состояния по кардиограмме)
- Классификация и аннотация текстов
- Предсказательная аналитика

Applications of RNNs

Speech recognition



“The quick brown fox jumped
over the lazy dog.”

Music generation

∅



Sentiment classification

“There is nothing to like
in this movie.”



DNA sequence analysis

AGCCCCTGTGAGGAACTAG



AG**CCCCTGTGAGGAACTAG**

Machine translation

Voulez-vous chanter avec
moi?



Do you want to sing with
me?

Video activity recognition



Running

Name entity recognition

Yesterday, Harry Potter
met Hermione Granger.



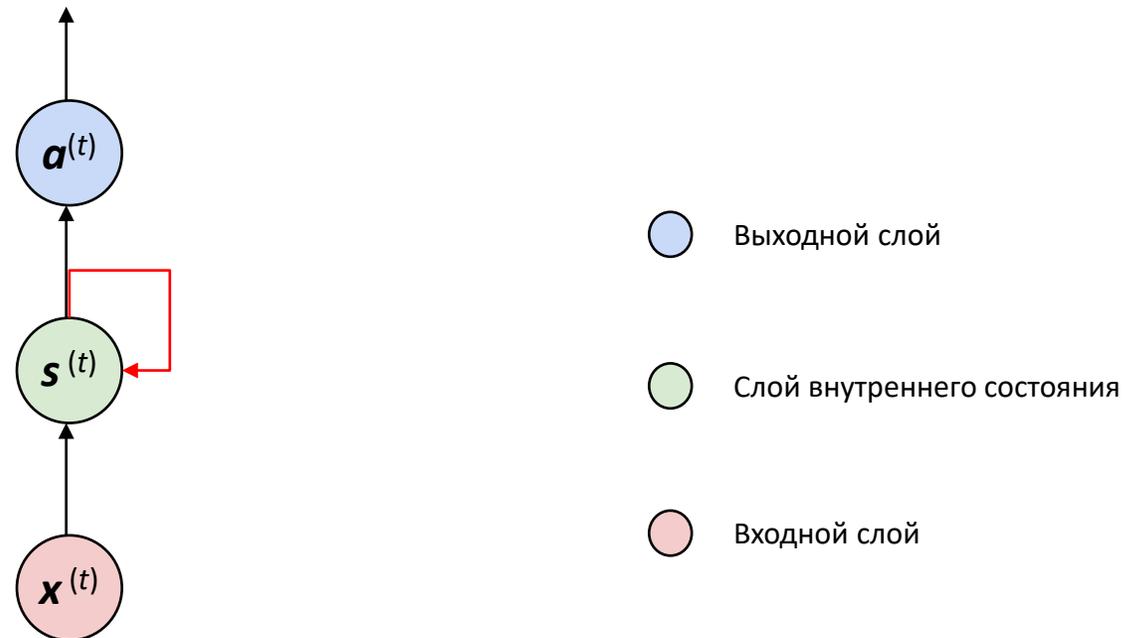
Yesterday, **Harry Potter**
met **Hermione Granger**.

Andrew Ng

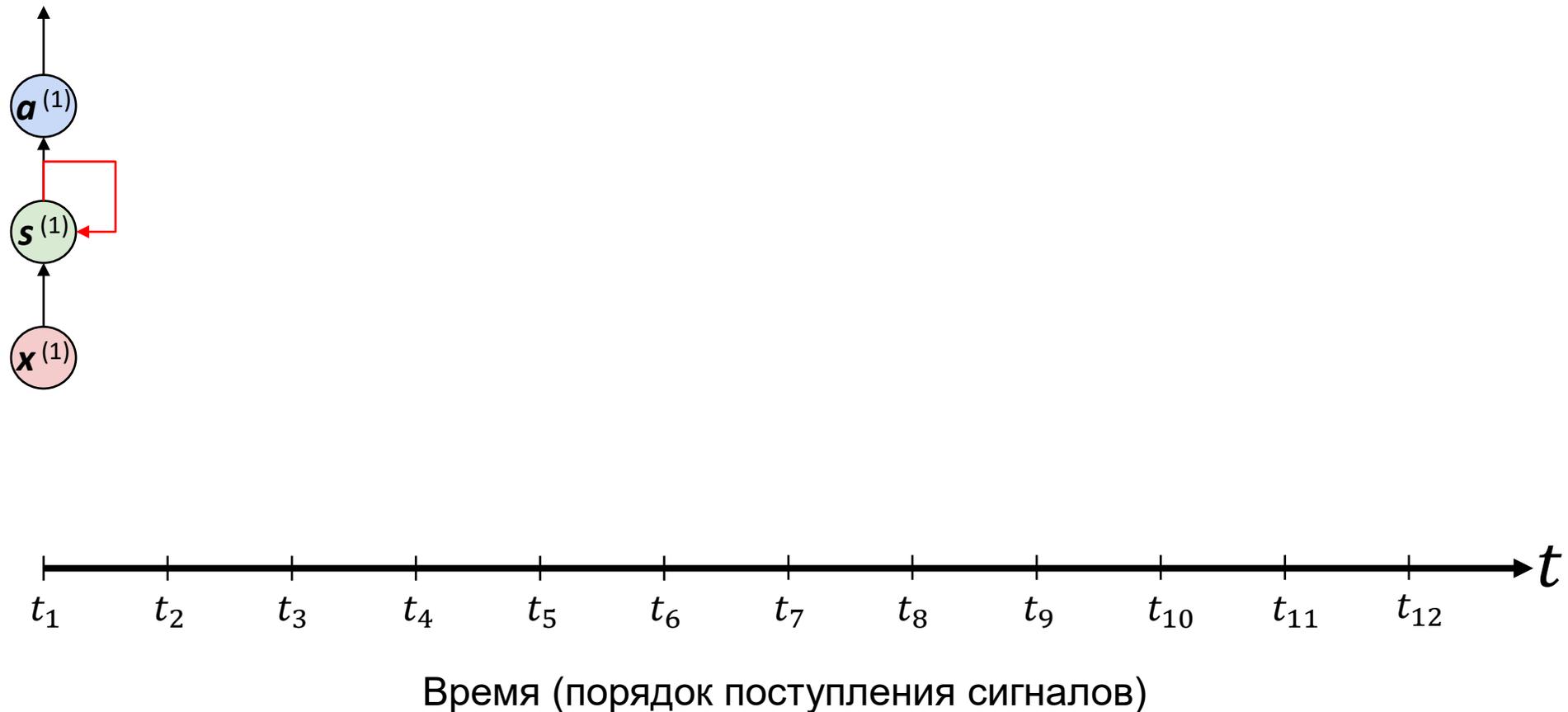
Идея

Запоминать что-то из истории приходящих на вход данных, сохраняя некое внутреннее состояние, которое можно было бы потом использовать для анализа текущих данных

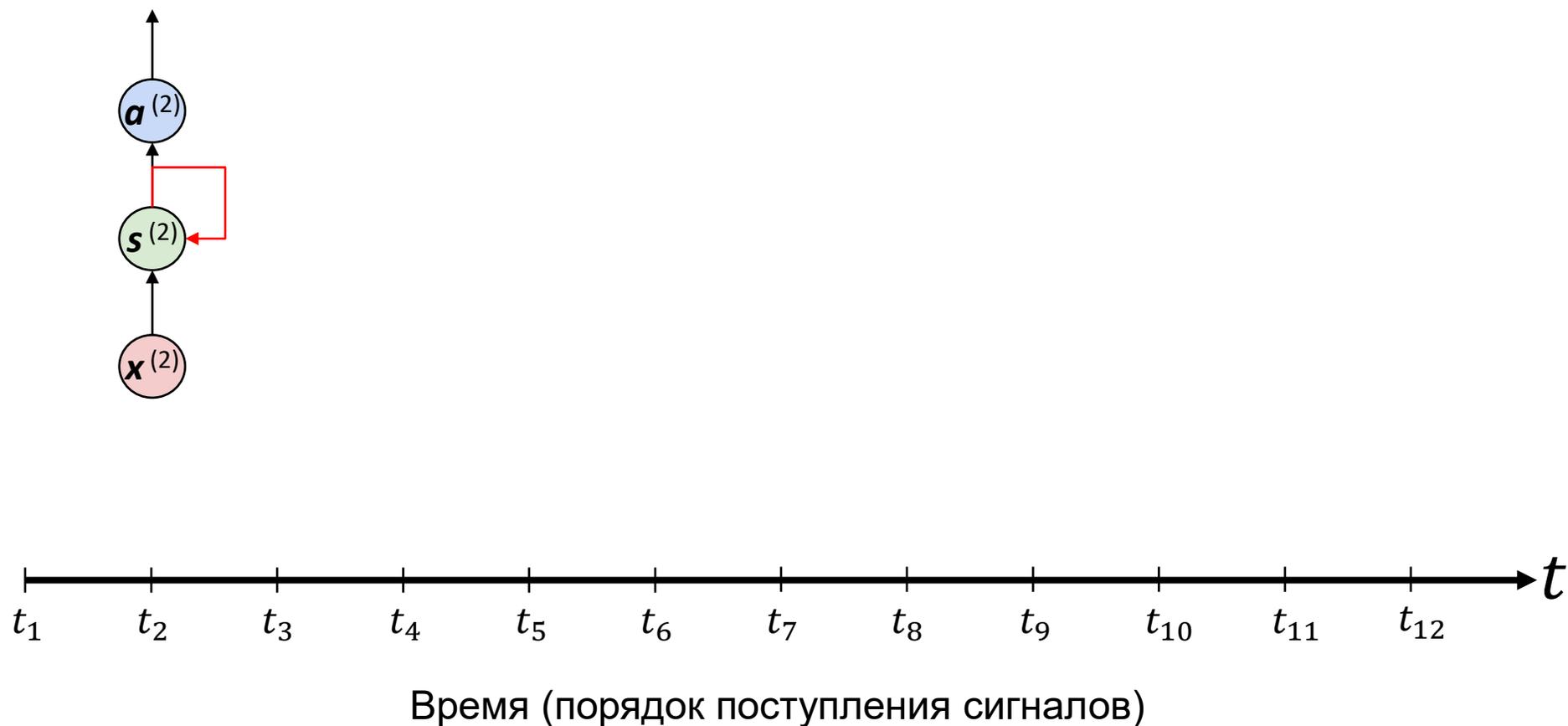
Простейшая рекуррентная сеть



Работа рекуррентной нейронной сети



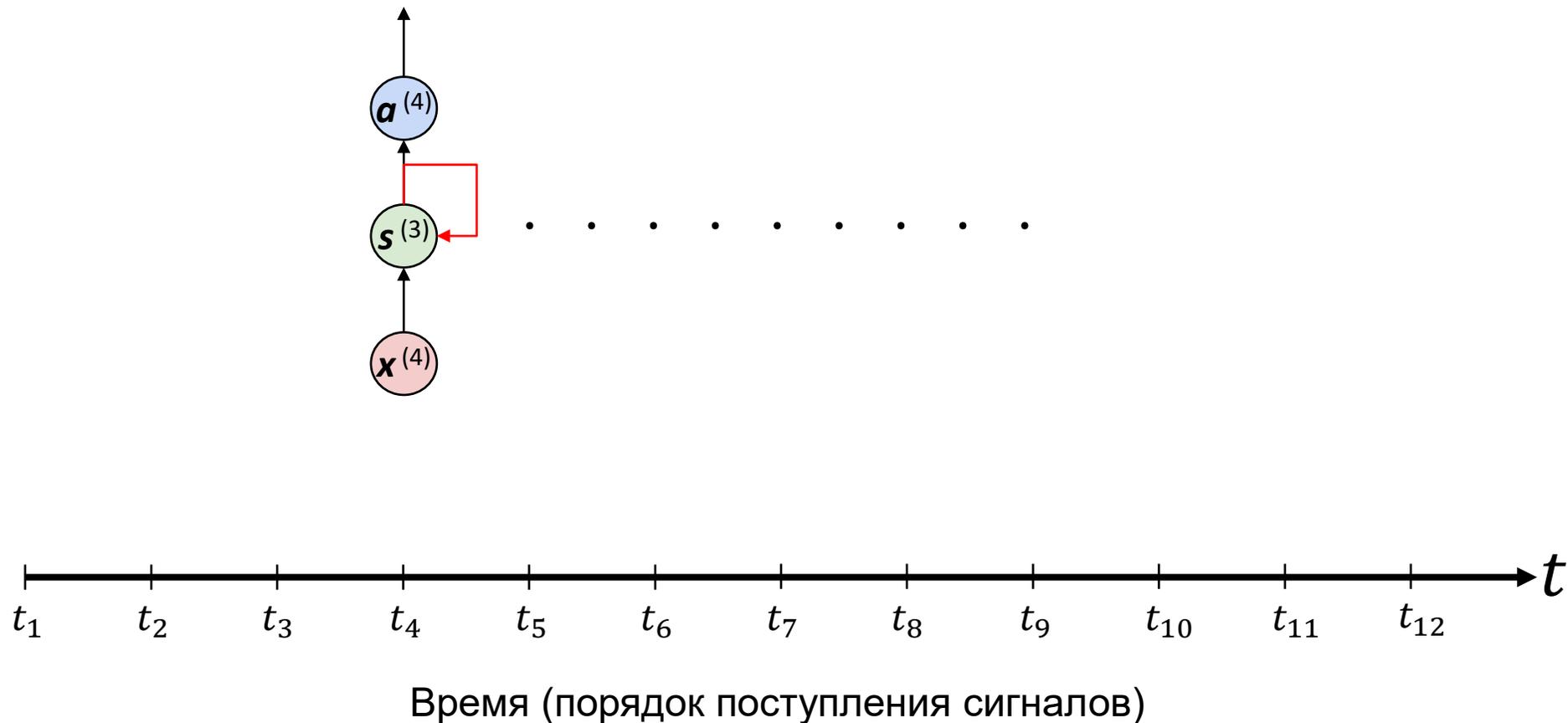
Работа рекуррентной нейронной сети



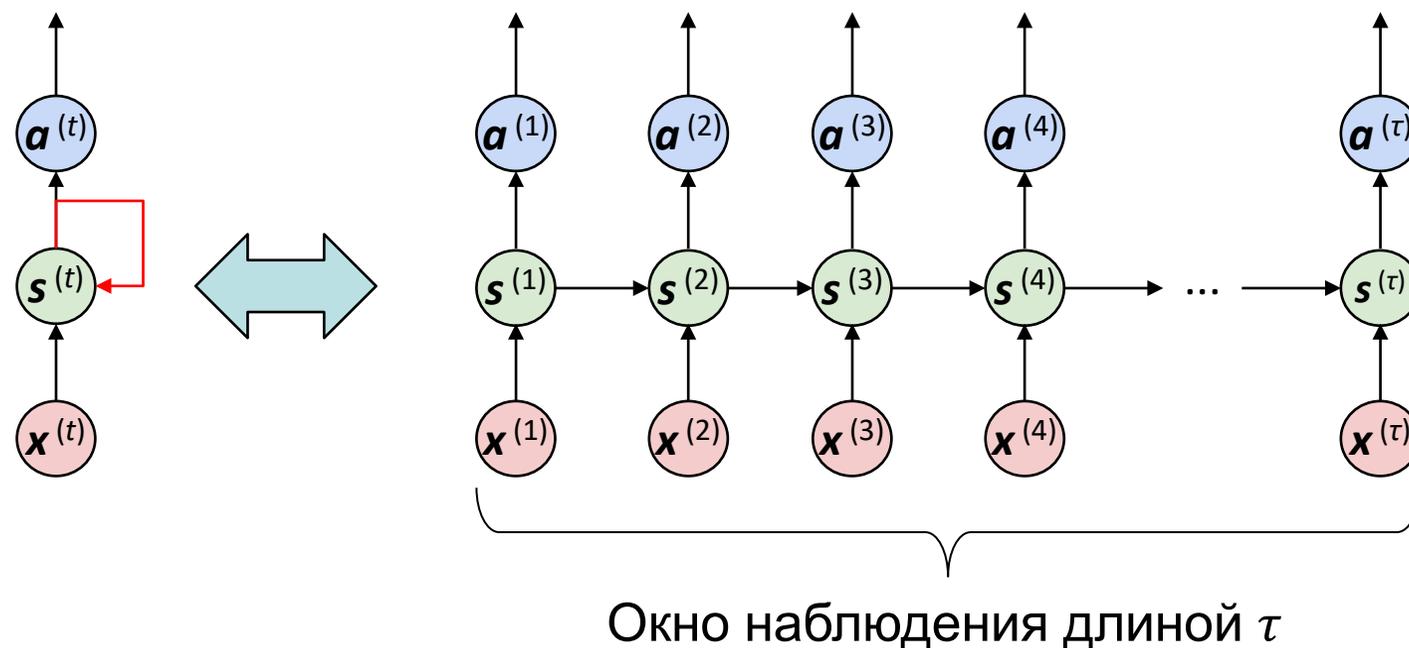
Работа рекуррентной нейронной сети



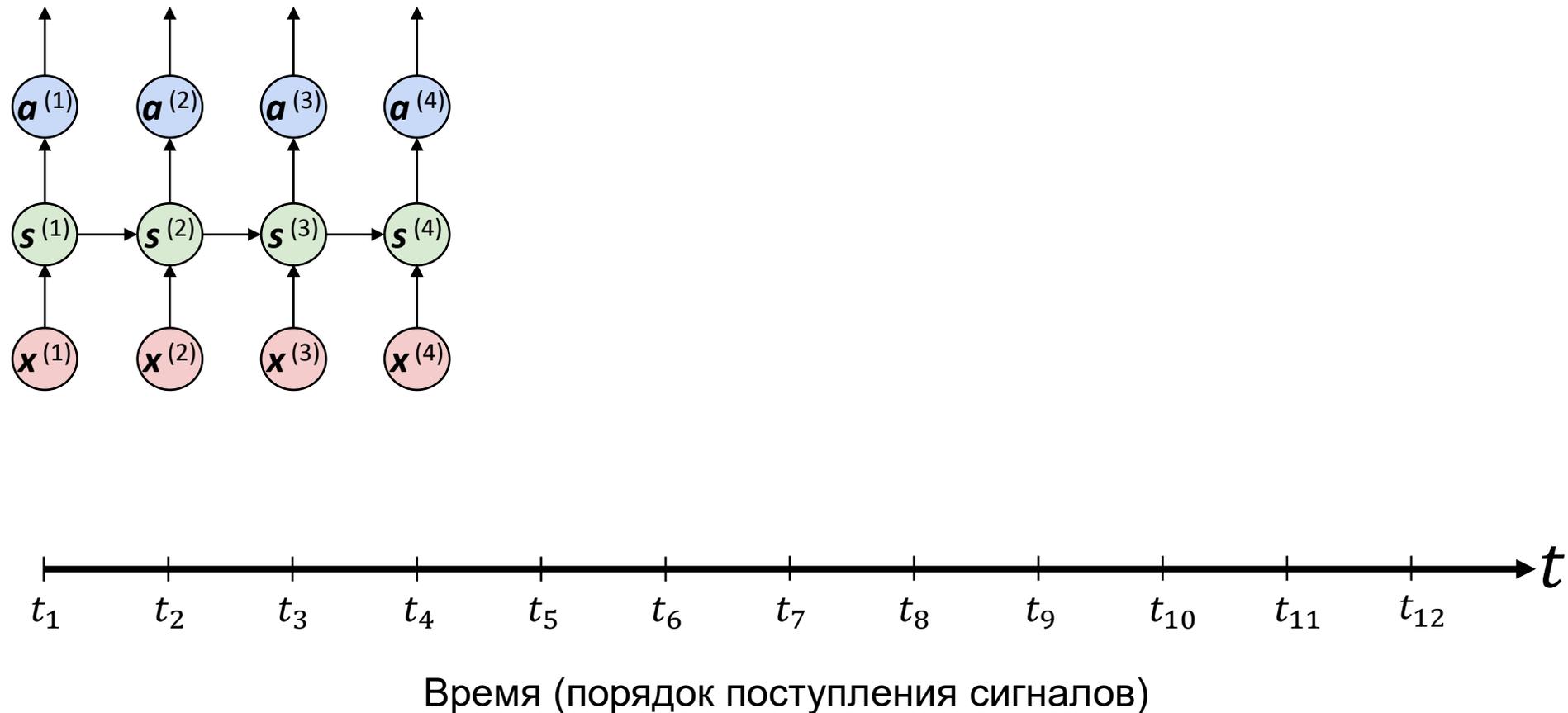
Работа рекуррентной нейронной сети



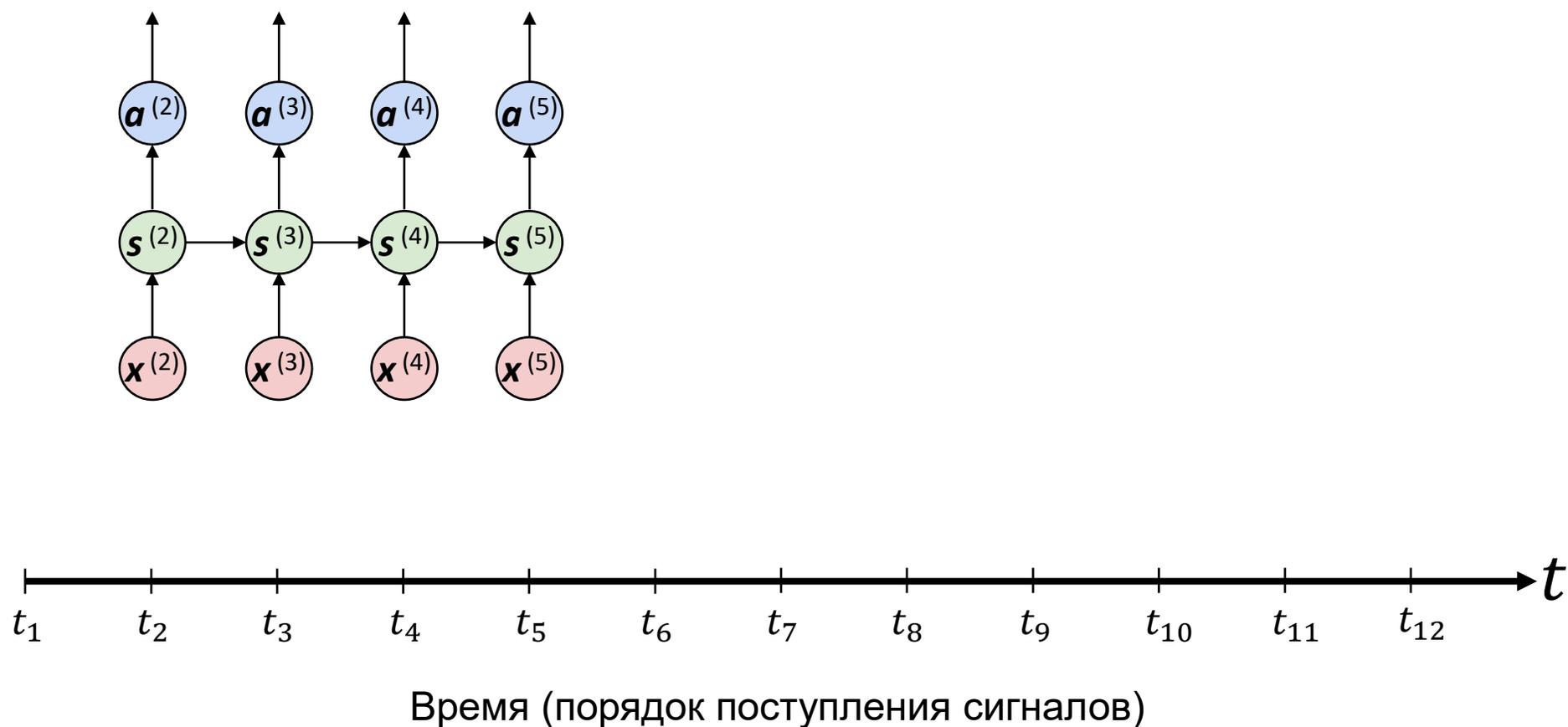
Развертка рекуррентной сети



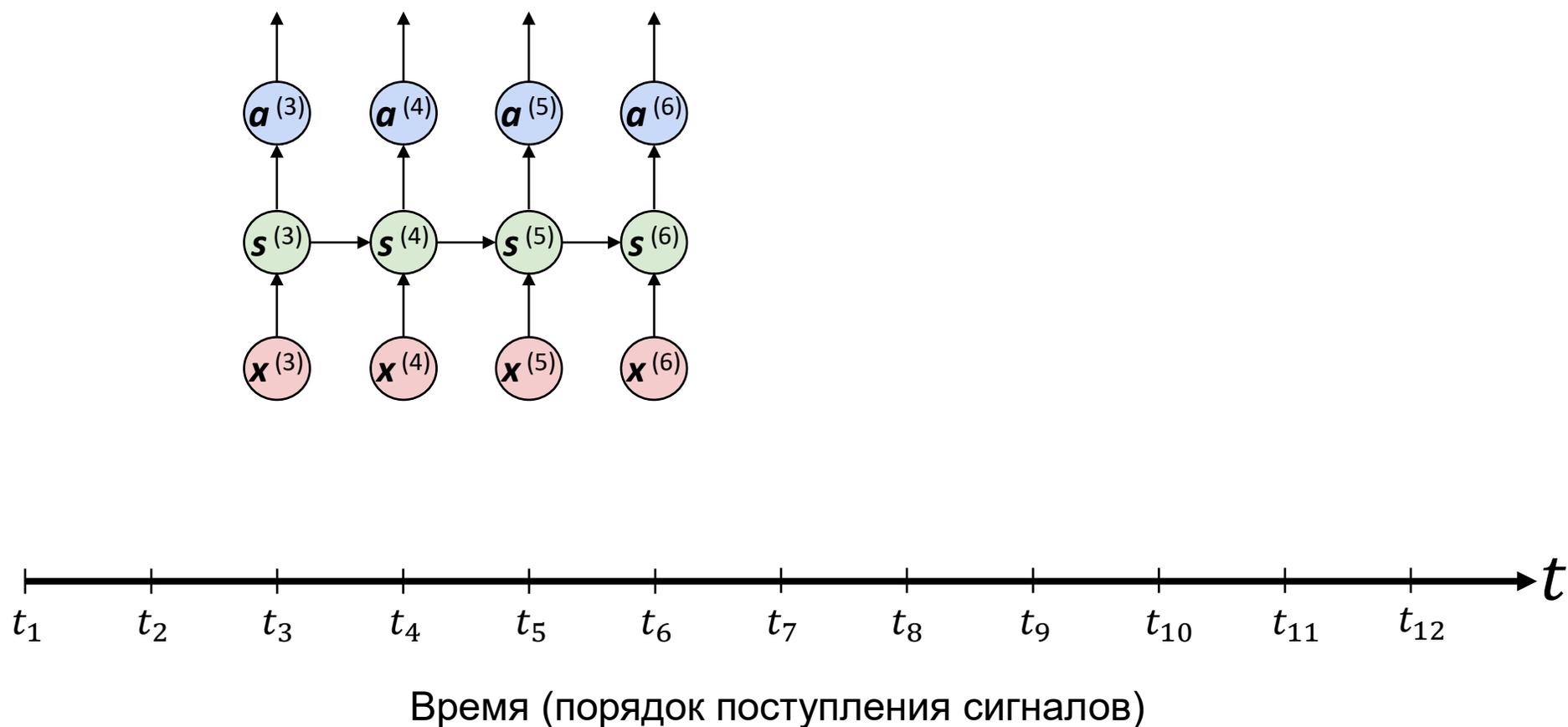
Работа развертки ($\tau = 4$)



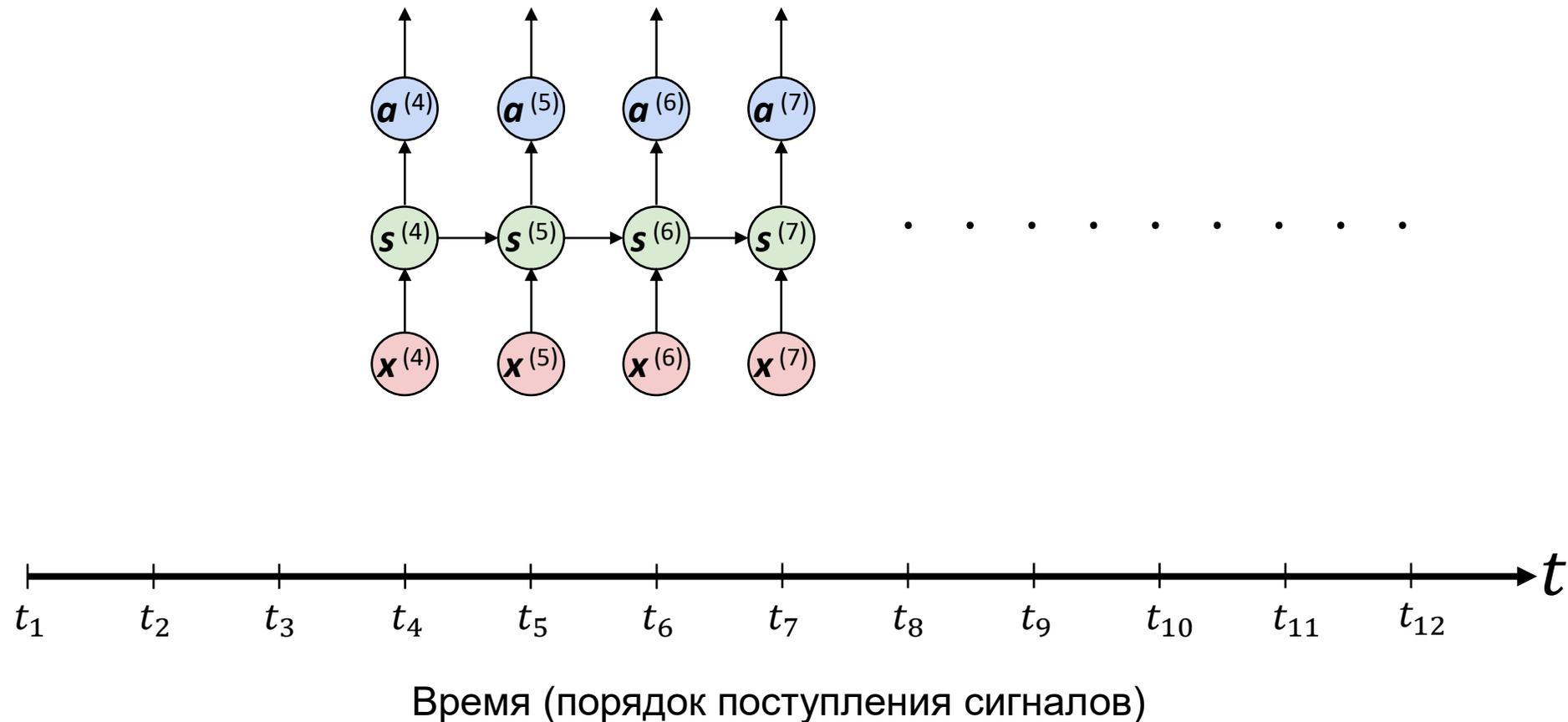
Работа развертки ($\tau = 4$)



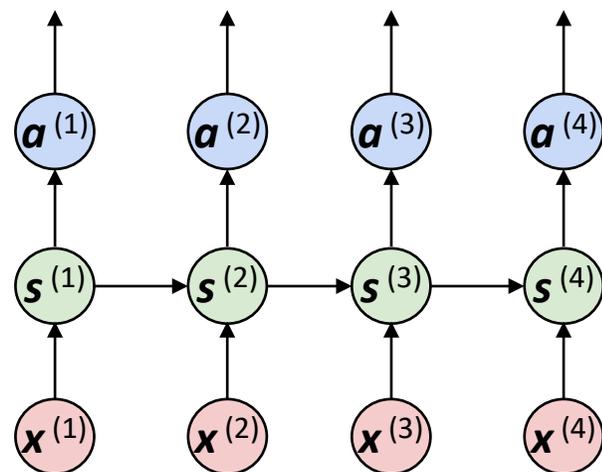
Работа развертки ($\tau = 4$)



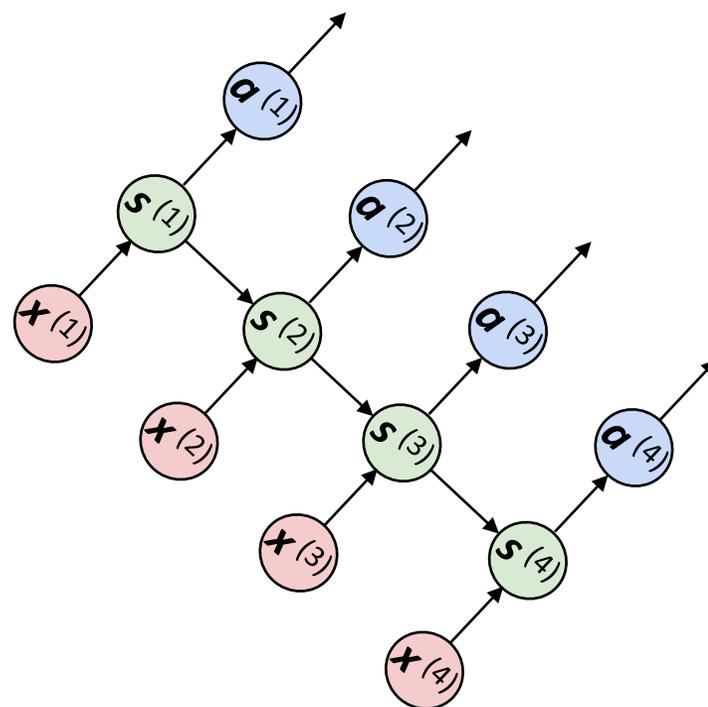
Работа развертки ($\tau = 4$)



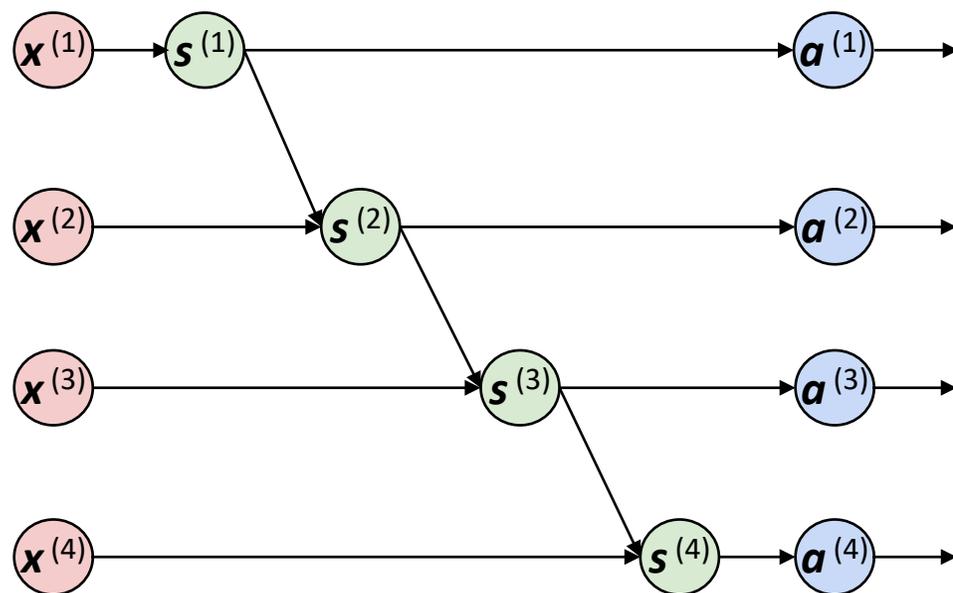
Развертка => сеть прямого распространения



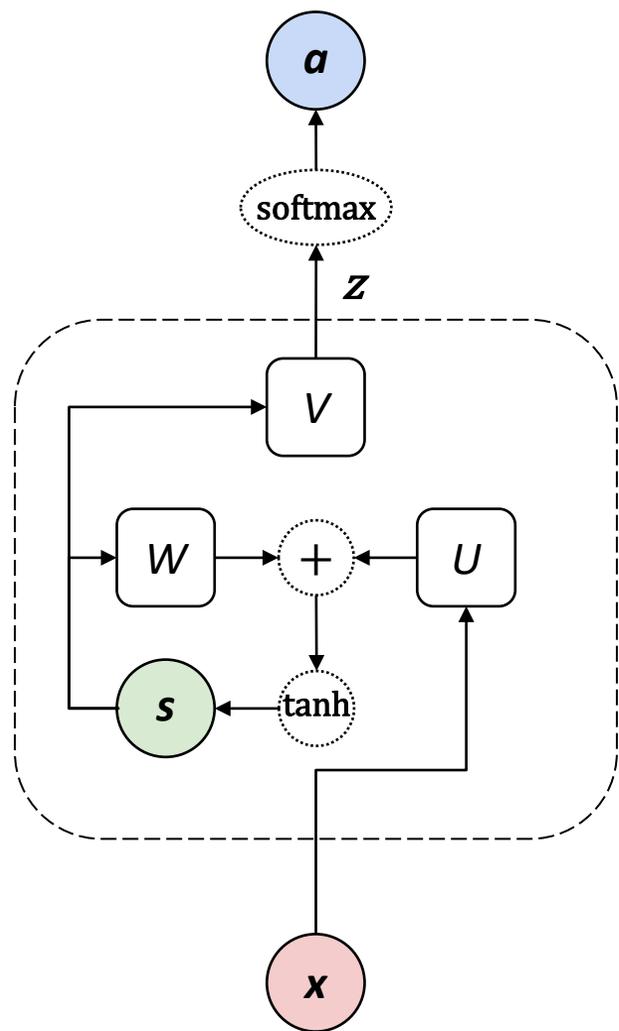
Развертка => сеть прямого распространения



Развертка => сеть прямого распространения



Модель Simple RNN



$$1) \mathbf{s}^{(t)} = \mathit{tanh}(W\mathbf{s}^{(t-1)} + U\mathbf{x}^{(t)} + \mathbf{b})$$

$$2) \mathbf{z}^{(t)} = V\mathbf{s}^{(t)} + \mathbf{d}$$

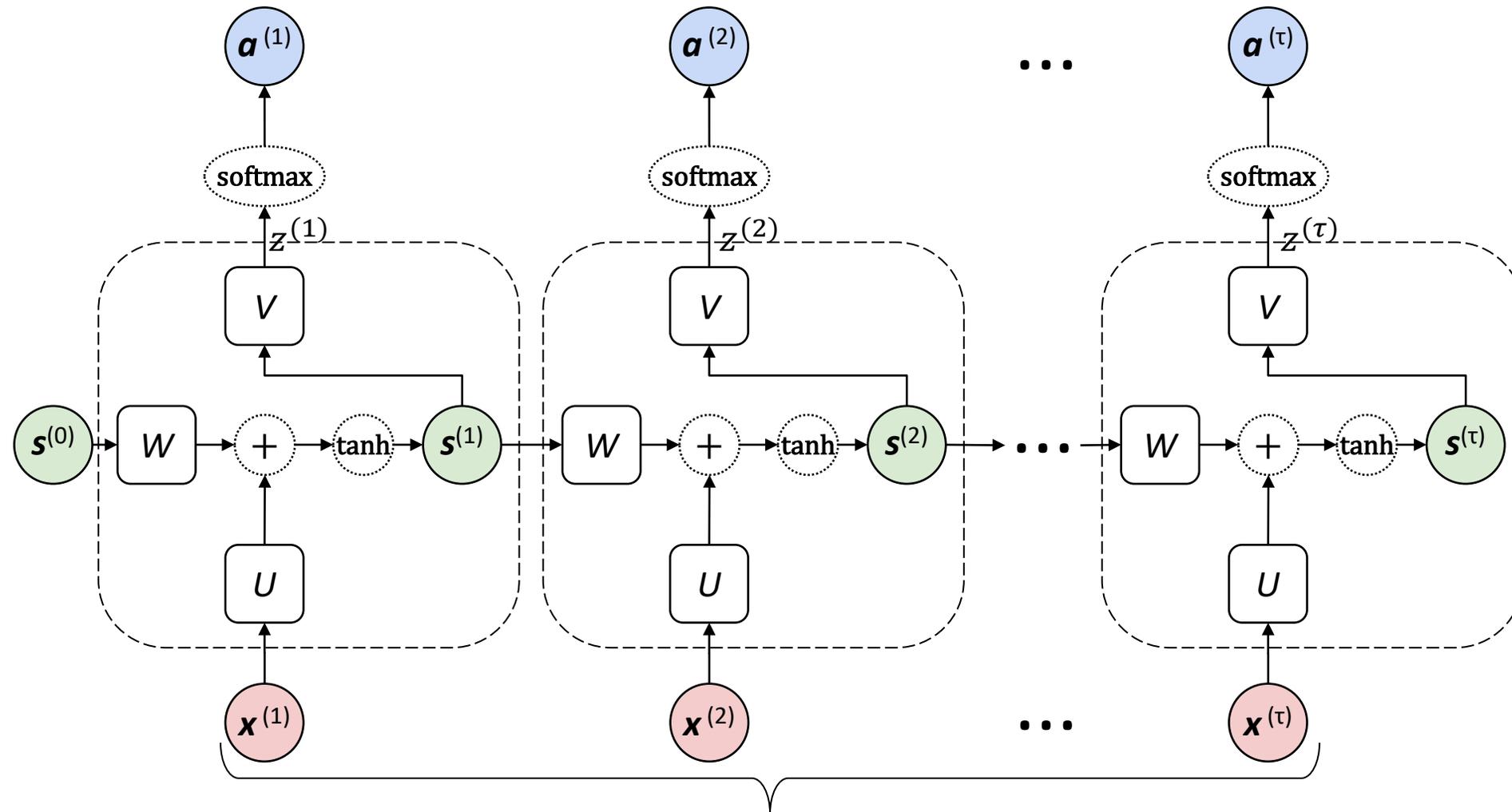
$$3) \mathbf{a}^{(t)} = \mathit{softmax}(\mathbf{z}^{(t)})$$

$$t = 1, 2, 3, \dots$$

$\mathbf{s}^{(t)}$ – вектор начального состояния

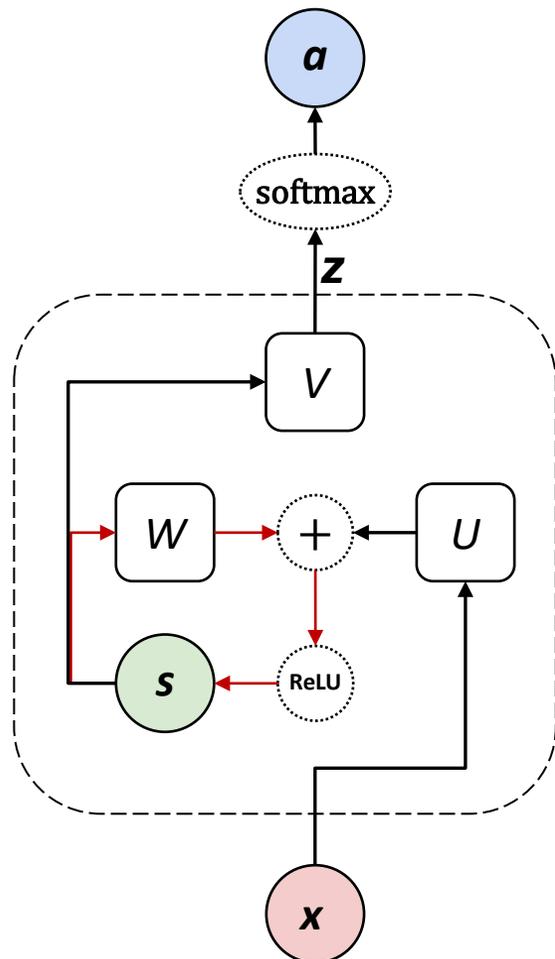
\mathbf{b}, \mathbf{d} – векторы смещений

Развертка Simple RNN



Окно наблюдения длиной τ

Задача: вычислить $a^{(2)}$



$$b = 0$$

$$d = 0$$

$$s^{(0)} = (0, 0)$$

$$x^{(1)} = (1, 2); x^{(2)} = (2, 1)$$

$$U = \begin{bmatrix} 0.9 & 0.2 \\ 0.4 & 0.6 \end{bmatrix}; W = \begin{bmatrix} 0.5 & 0.8 \\ 0.4 & 0.2 \end{bmatrix}; V = \begin{bmatrix} 0.3 & 0.8 \\ 0.9 & 0.4 \end{bmatrix}$$

$$s^{(1)} = Ux^{(1)} + Ws^{(0)} = (1.3, 1.6)$$

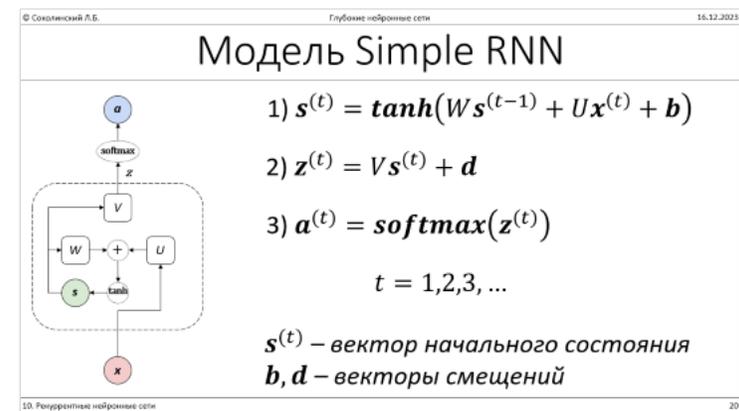
$$s^{(2)} = Ux^{(2)} + Ws^{(1)} = (3.93, 2.24)$$

$$z^{(2)} = Vs^{(2)} = (2.971, 4.433)$$

$$e^{z_1^{(2)}} = 19.51142; e^{z_2^{(2)}} = 84.18359$$

$$e^{z_1^{(2)}} + e^{z_2^{(2)}} = 103.69501$$

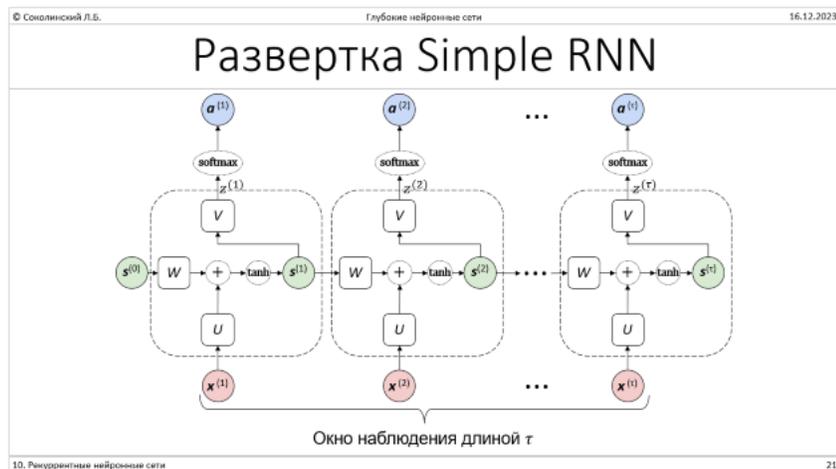
$$a^{(2)} = \left(\frac{e^{z_1^{(2)}}}{e^{z_1^{(2)}} + e^{z_2^{(2)}}}, \frac{e^{z_2^{(2)}}}{e^{z_1^{(2)}} + e^{z_2^{(2)}}} \right) = (0.18816, 0.81184)$$



Обратное распространение во времени (Back propagation through time)

Ошибка в момент времени t :

$$C_{(x^{(t)}, y^{(t)})} = \sum_j y_j^{(t)} \ln a_j^{(t)}$$



© Соколинский Л.Б. Глубокие нейронные сети 16.12.2023

Использование в качестве функции потерь функции перекрестной энтропии

$$C = - \sum_j y_j \ln a_j^t$$

Свойства функции потерь:

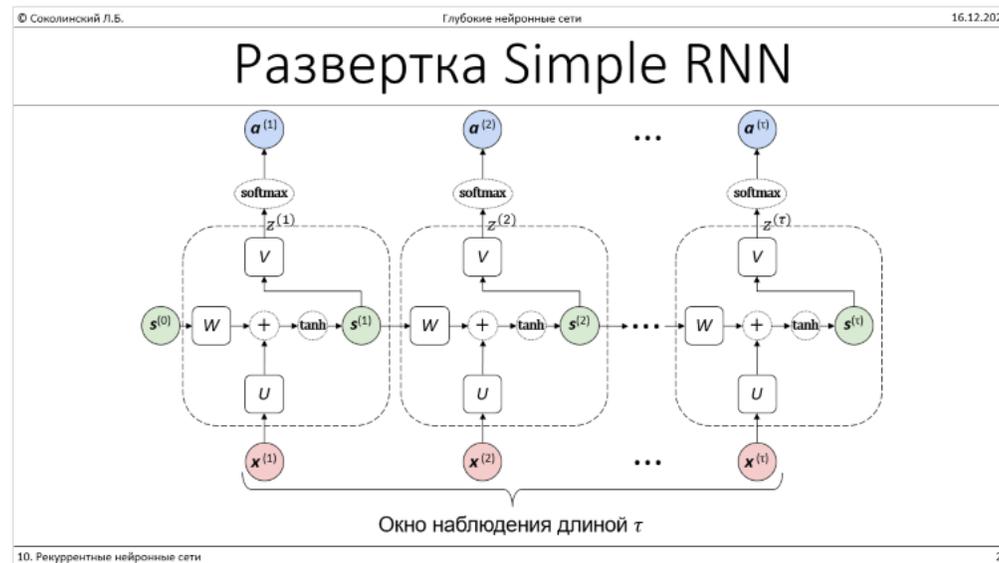
- 1) $C \geq 0$
- 2) Минимум C достигается при $a^L = y$:

$$\min_{a^L} C = - \sum_j y_j \ln y_j$$

10. Рекуррентные нейронные сети 48

Суммарная ошибка

$$C(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}) = \sum_{t=1}^{\tau} C(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})$$



Производная для суммарной ошибки

Обозначим

$$\mathbb{C} = \mathbb{C}(\{x^{(1)}, \dots, x^{(\tau)}\}, \{y^{(1)}, \dots, y^{(\tau)}\})$$

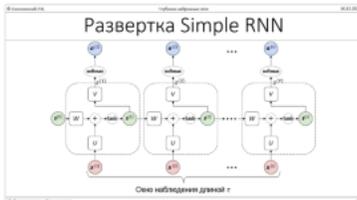
$$C^{(t)} = C(x^{(t)}, y^{(t)})$$

Имеем

$$\frac{\partial \mathbb{C}}{\partial C^{(t)}} = 1$$

© Соколинский Л.Б. Глубокие нейронные сети 16.12.2023

Суммарная ошибка

$$\mathbb{C}(\{x^{(1)}, \dots, x^{(\tau)}\}, \{y^{(1)}, \dots, y^{(\tau)}\}) = \sum_{t=1}^{\tau} C(x^{(t)}, y^{(t)})$$


Развертка Simple RNN

Схема наблюдения данных τ

10. Рекуррентные нейронные сети 24

Вычисление градиента для выходного слоя

$$\left(\nabla_{\mathbf{z}^{(t)}} \mathbb{C}\right)_j = \frac{\partial \mathbb{C}}{\partial z_j^{(t)}} = \frac{\partial \mathbb{C}}{\partial C^{(t)}} \frac{\partial C^{(t)}}{\partial z_j^{(t)}} = a_j^{(t)} - y_j$$

© Соколинский Л.Б. Глубокие нейронные сети 16.12.2023

Производная для суммарной ошибки

Обозначим $\mathbb{C} = \mathbb{C}_{(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n)}\})}$

Имеем $C^{(t)} = C_{(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})}$

$$\frac{\partial \mathbb{C}}{\partial C^{(t)}} = 1$$

Суммарная ошибка

$$\mathbb{C}_{(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n)}\})} = \sum_{t=1}^n C_{(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})}$$


10. Рекуррентные нейронные сети 25

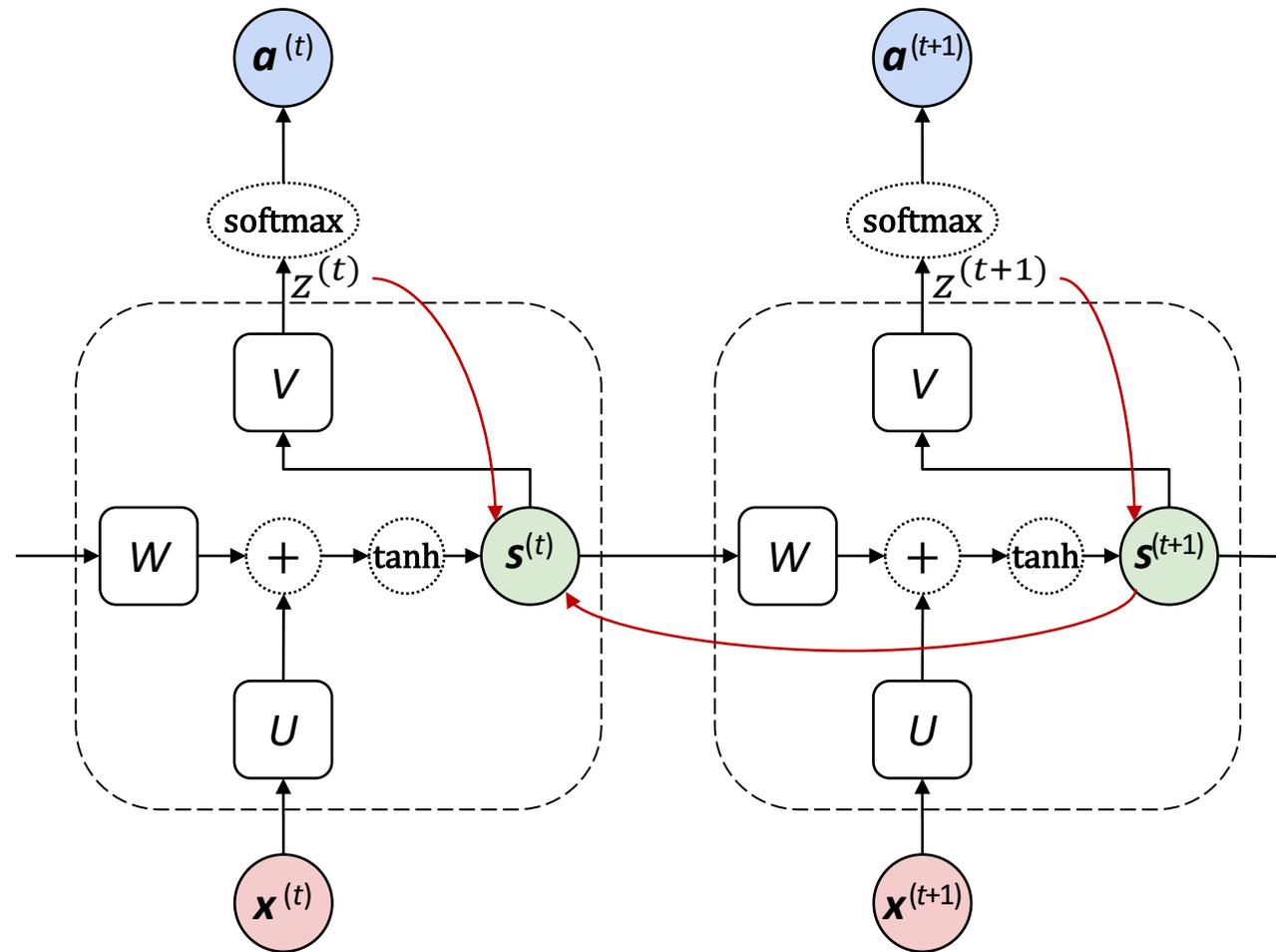
© Соколинский Л.Б. Глубокие нейронные сети 16.12.2023

Вычисление ошибки δ_i^L для выходного слоя

$$\begin{aligned} \delta_i^L &= \frac{\partial \mathbb{C}}{\partial z_i^L} = \sum_{j=1}^n \frac{\partial \mathbb{C}}{\partial a_j^L} \cdot \frac{\partial a_j^L}{\partial z_i^L} = \sum_{j=1}^n -\frac{y_j}{a_j^L} \cdot \frac{\partial a_j^L}{\partial z_i^L} = -\sum_{j=1}^n \frac{y_j}{a_j^L} \cdot \frac{\partial a_j^L}{\partial z_i^L} \\ &= -\frac{y_i}{a_i^L} \cdot \frac{\partial a_i^L}{\partial z_i^L} - \sum_{j \neq i}^n \frac{y_j}{a_j^L} \cdot \frac{\partial a_j^L}{\partial z_i^L} \\ &= -\frac{y_i}{a_i^L} a_i^L (1 - a_i^L) - \sum_{j \neq i}^n -\frac{y_j}{a_j^L} \cdot a_i^L a_j^L = -y_i + y_i a_i^L + \sum_{j \neq i}^n y_j a_i^L \\ &= -y_i + \sum_{j=1}^n y_j a_i^L = -y_i + a_i^L \sum_{j=1}^n y_j \\ &= -y_i + a_i^L = a_i^L - y_i \end{aligned}$$

10. Рекуррентные нейронные сети 49

Обратные зависимости для слоя состояний



Вычисление градиентов для слоя состояний

Для $t = \tau$:

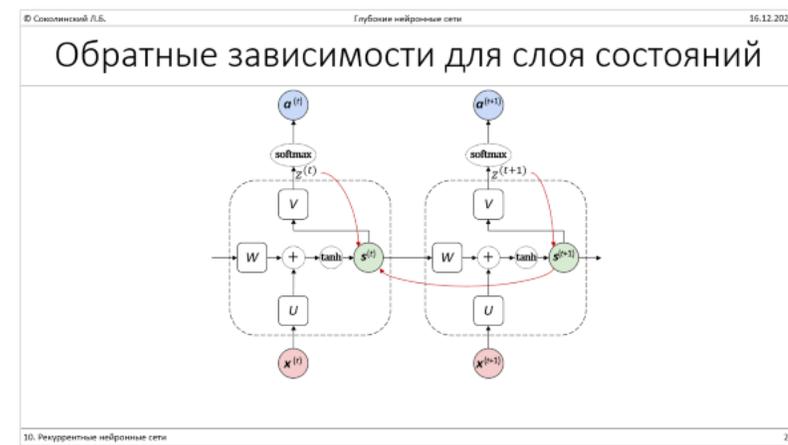
$$\nabla_{\mathbf{s}^{(\tau)}} \mathbb{C} = \frac{\partial \mathbf{z}^{(\tau)}}{\partial \mathbf{s}^{(\tau)}} \circ \nabla_{\mathbf{z}^{(\tau)}} \mathbb{C} = V^T \nabla_{\mathbf{z}^{(\tau)}} \mathbb{C}$$

Для $t = \tau-1, \dots, 1$:

$$\begin{aligned} \nabla_{\mathbf{s}^{(t)}} \mathbb{C} &= \frac{\partial \mathbf{s}^{(t+1)}}{\partial \mathbf{s}^{(t)}} \circ \nabla_{\mathbf{s}^{(t+1)}} \mathbb{C} + \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{s}^{(t)}} \circ \nabla_{\mathbf{z}^{(t)}} \mathbb{C} \\ &= W^T (\nabla_{\mathbf{s}^{(t+1)}} \mathbb{C}) \cdot \text{diag} \left(1 - (\mathbf{s}^{(t+1)})^2 \right) + V^T (\nabla_{\mathbf{z}^{(t)}} \mathbb{C}) \end{aligned}$$

$\text{diag} \left(1 - (\mathbf{s}^{(t+1)})^2 \right)$ - диагональная матрица с элементами $1 - (s_j^{(t+1)})^2$

$\tanh'(x) = 1 - \tanh^2(x)$ (см. <https://socratic.org/questions/what-is-the-derivative-of-tanh-x>)



Вычисление градиента для выходного слоя

$$(\nabla_{\mathbf{z}^{(t)}} \mathbb{C})_j = \frac{\partial \mathbb{C}}{\partial z_j^{(t)}} = \frac{\partial \mathbb{C}}{C^{(t)}} \frac{\partial C^{(t)}}{\partial z_j^{(t)}} = a_j^{(t)} - y_j$$

Производная для суммарной ошибки

Обозначим $C = C_1 + C_2 + \dots + C_n$

Итак $\frac{\partial C}{\partial C^{(t)}} = 1$

Вычисление ошибки δ_j^t для выходного слоя

$\delta_j^t = \frac{\partial C}{\partial z_j^{(t)}} = a_j^{(t)} - y_j$

10. Рекуррентные нейронные сети 26

Вычисление градиентов для весов и смещений

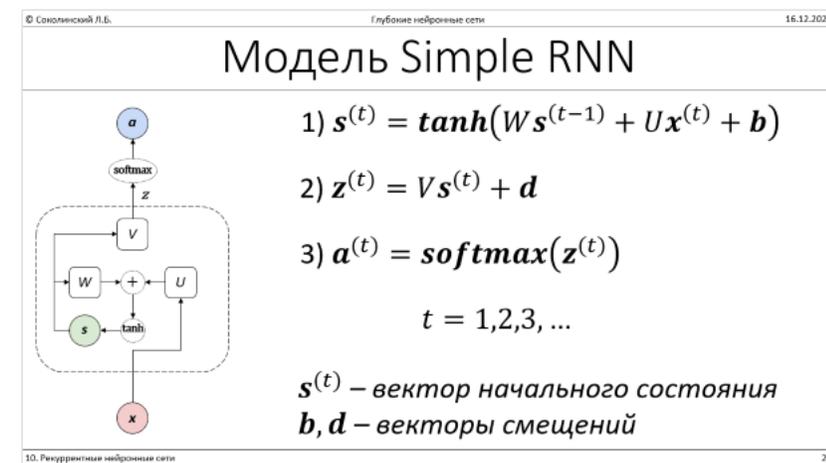
$$\nabla_{\mathbf{d}} \mathbb{C} = \sum_t \nabla_{\mathbf{z}^{(t)}} \mathbb{C}$$

$$\nabla_{\mathbf{b}} \mathbb{C} = \sum_t \text{diag} \left(1 - (\mathbf{s}^{(t)})^2 \right) \nabla_{\mathbf{s}^{(t)}} \mathbb{C}$$

$$\nabla_{\mathbf{V}} \mathbb{C} = \sum_t \left(\nabla_{\mathbf{z}^{(t)}} \mathbb{C} \right) \circ \mathbf{s}^{(t)}$$

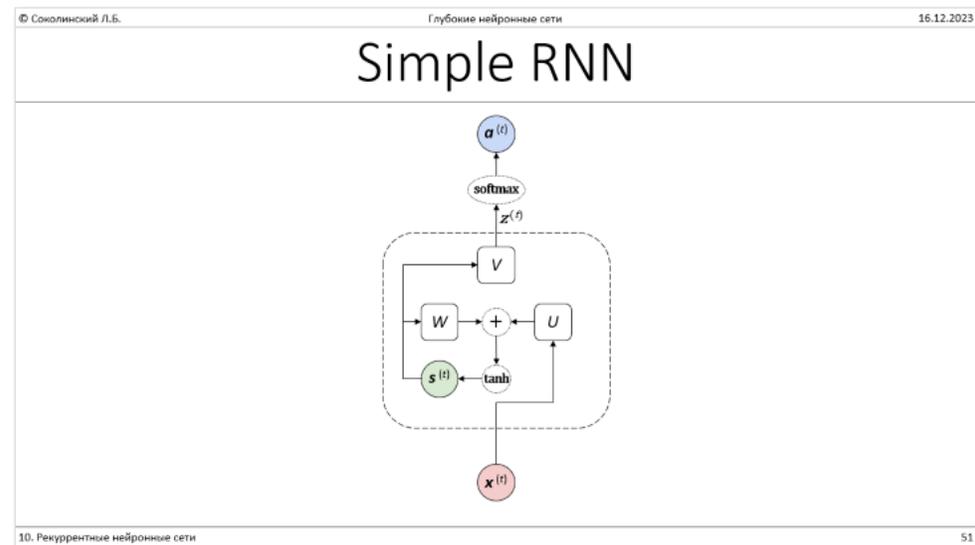
$$\nabla_{\mathbf{W}} \mathbb{C} = \sum_t \text{diag} \left(1 - (\mathbf{s}^{(t)})^2 \right) \left(\nabla_{\mathbf{s}^{(t)}} \mathbb{C} \right) \circ \mathbf{s}^{(t-1)}$$

$$\nabla_{\mathbf{U}} \mathbb{C} = \sum_t \text{diag} \left(1 - (\mathbf{s}^{(t)})^2 \right) \left(\nabla_{\mathbf{s}^{(t)}} \mathbb{C} \right) \circ \mathbf{x}^{(t)}$$

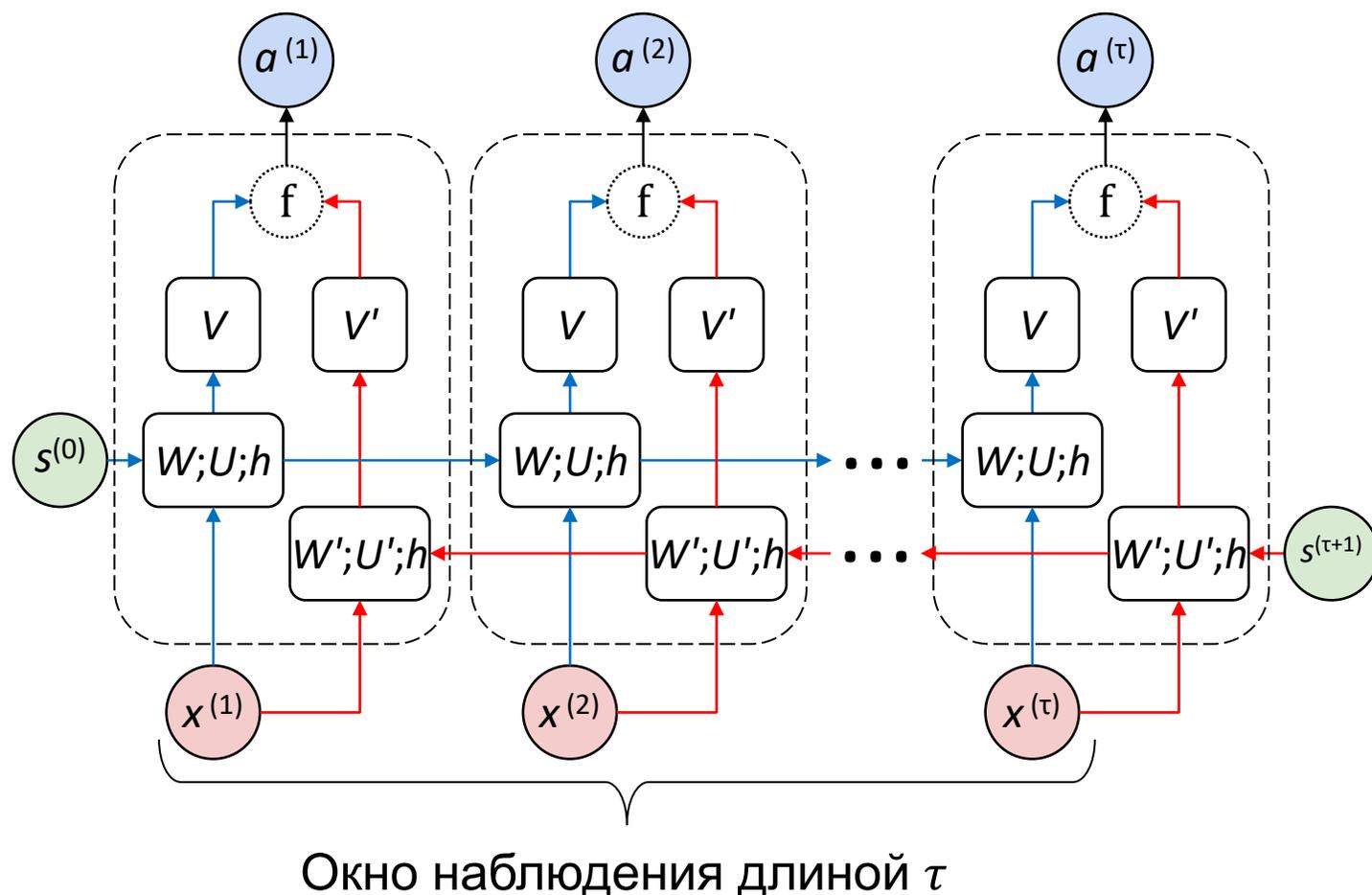


Глубокие RNN

1. Вставить глубокую сеть перед слоем внутреннего состояния (позволяет понимать временную структуру информации, например, выделять слова при распознавании речи)
2. Вставить глубокую сеть после слоя внутреннего состояния (позволяет выявлять закономерности, не зависящие от времени)
3. Использовать глубокую сеть вместо \tanh (фрагментация изображений)
4. Использовать выходы одной RNN в качестве входов другой RNN (каждый слой действует в своем масштабе времени)



Двунаправленные RNN



Позволяют получить два состояния, отражающие контекст как слева, так и справа для каждого элемента последовательности

Пример: разметка слов в предложении по частям речи

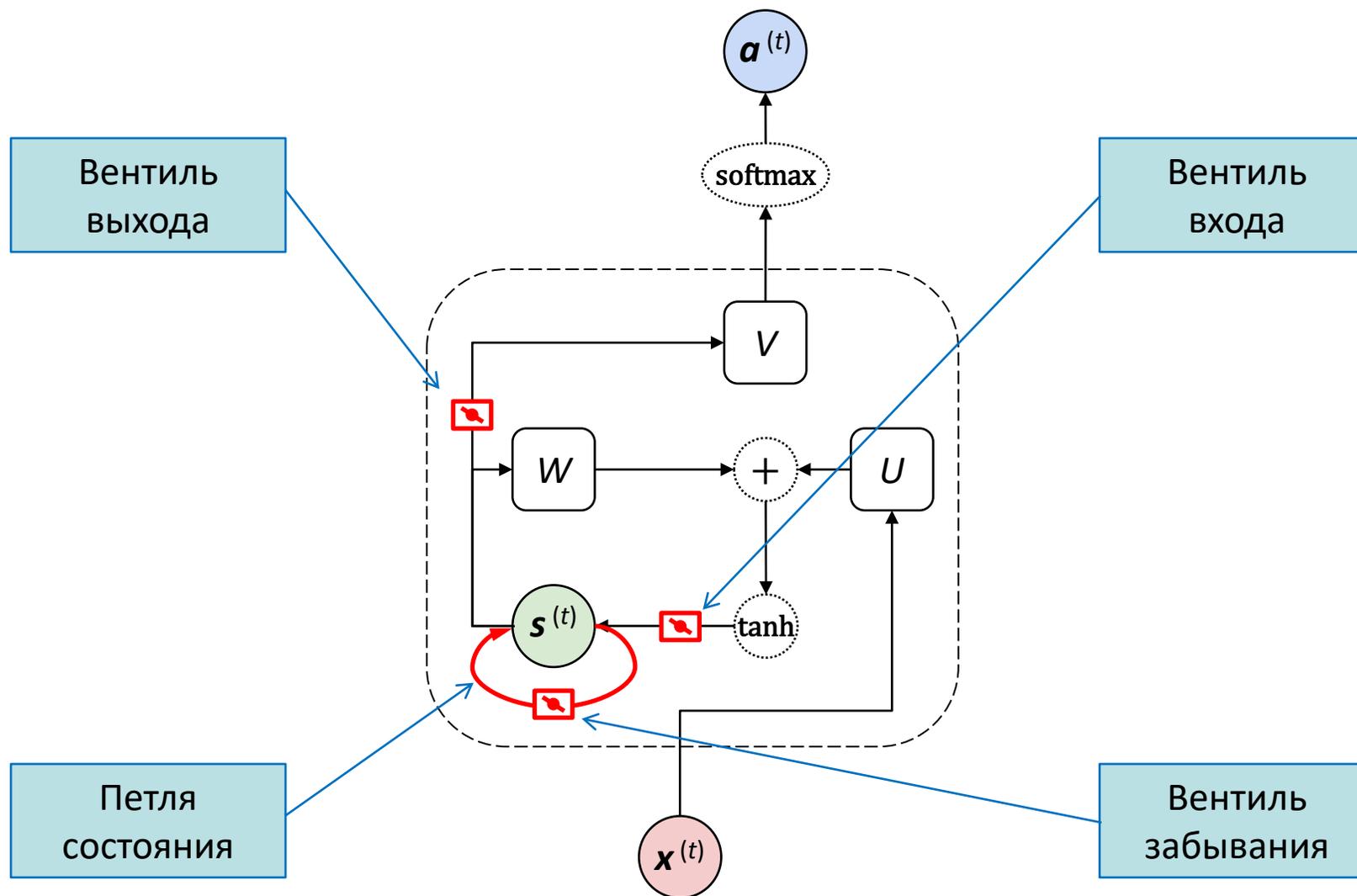
Недостатки RNN

- **Не позволяют надолго запоминать информацию во временной последовательности: влияние текущего состояния на будущее состояния экспоненциально затухает**
- Градиенты, распространяющиеся через много одинаковых слоев, либо исчезают (в большинстве случаев), либо начинают взрывообразно расти (редко, но с большим уроном для обучения)
- Высокая вычислительная сложность: не удастся обучать RNN с окном наблюдения размером 10 и более

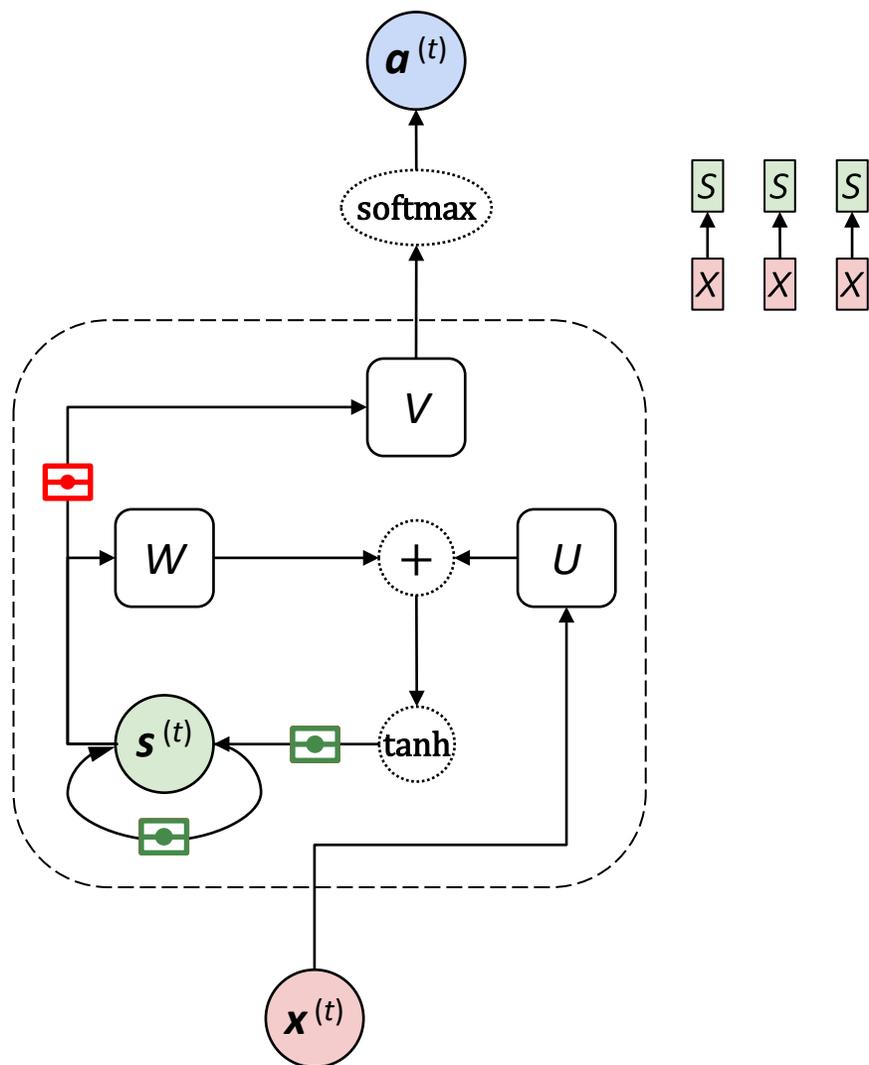


Как преодолеть забывчивость RNN?

Идея! Вентильные РНС (Gated RNN)



Пример: синхронный машинный перевод



© Соколинский Л.Б. Глубокие нейронные сети 16.12.2023

Типы приложений

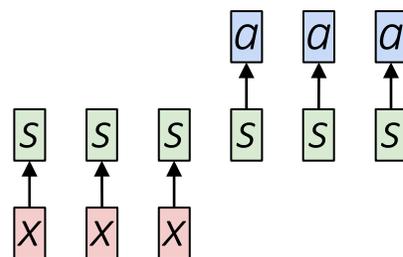
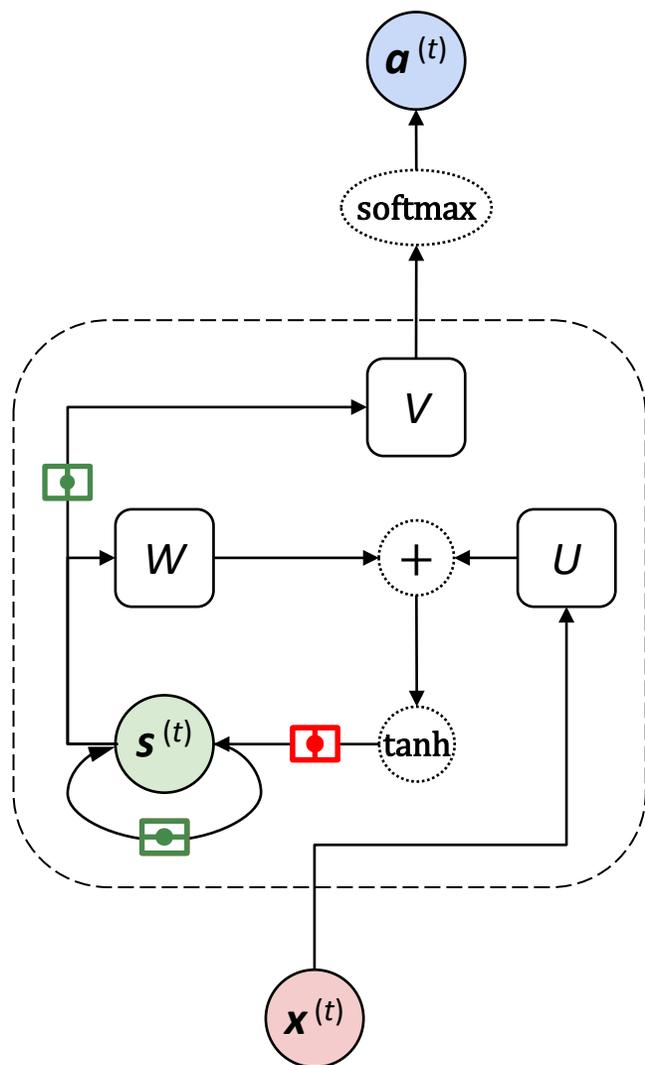
a) 1:1 b) 1:M c) M:1 d) 1xM : 1xM e) M:M

а) Один вход, один выход без запоминания состояния (распознавание изображений)
 б) Последовательность на выходе (вход – картинка, выход – ее словесное описание)
 в) Последовательность на входе (вход – рецензия на фильм, выход – оценка по десятибалльной шкале)
 г) Последовательность на входе, последовательность на выходе (перевод текста с одного языка на другой)
 е) Синхронизированные последовательности на входе и выходе (разметка смены кадров на видео; очистка аудиопотока от посторонних шумов)

Выходной вектор
 Вектор состояния
 Входной вектор

10. Рекуррентные нейронные сети 3

Пример: синхронный машинный перевод



© Соколинский Л.Б. Глубокие нейронные сети 16.12.2023

Типы приложений

a) 1:1

b) 1:M

c) M:1

d) 1xM : 1xM

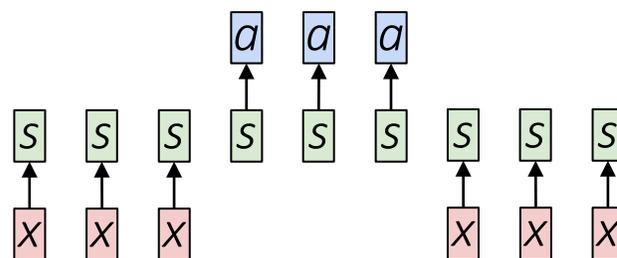
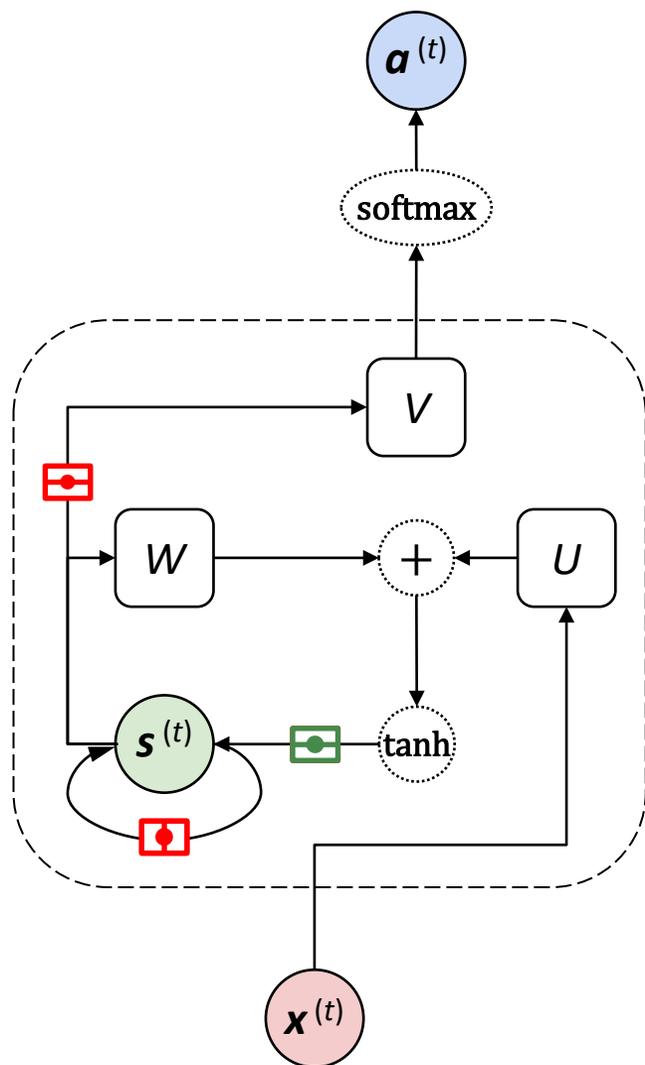
e) M:M

■ Выходной вектор
■ Вектор состояния
■ Входной вектор

а) Один вход, один выход без запоминания состояния (распознавание изображений)
 б) Последовательность на выходе (вход – картинка, выход – ее словесное описание)
 в) Последовательность на входе (вход – рецензия на фильм, выход – оценка по десятибалльной шкале)
 г) Последовательность на входе, последовательность на выходе (перевод текста с одного языка на другой)
 д) Синхронизированные последовательности на входе и выходе (разметка смены кадров на видео; очистка аудиопотока от посторонних шумов)

10. Рекуррентные нейронные сети 3

Пример: синхронный машинный перевод



© Соколинский Л.Б. Глубокие нейронные сети 16.12.2023

Типы приложений

a) 1:1

b) 1:M

c) M:1

d) 1xM : 1xM

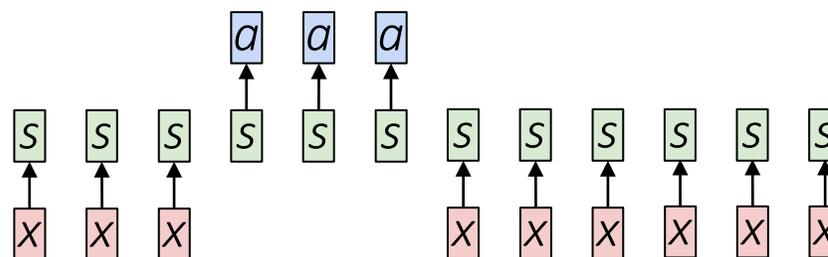
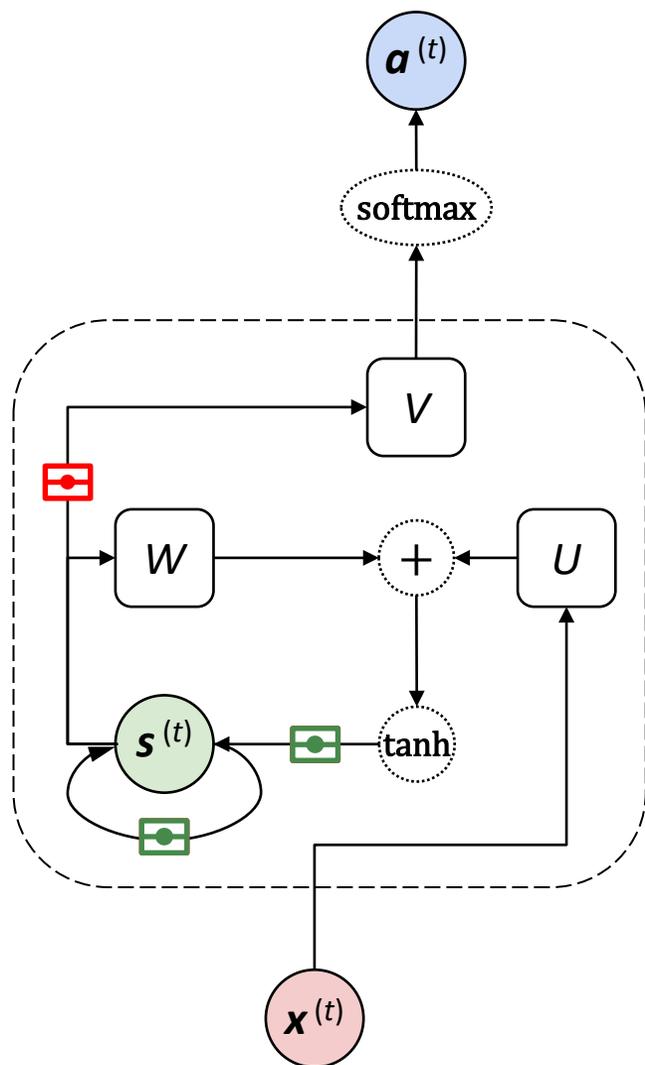
e) M:M

■ Выходной вектор
■ Вектор состояния
■ Входной вектор

а) Один вход, один выход без запоминания состояния (распознавание изображений)
 б) Последовательность на выходе (вход – картинка, выход – ее словесное описание)
 в) Последовательность на входе (вход – рецензия на фильм, выход – оценка по десятибалльной шкале)
 г) Последовательность на входе, последовательность на выходе (перевод текста с одного языка на другой)
 е) Синхронизированные последовательности на входе и выходе (разметка смены кадров на видео; очистка аудиопотока от посторонних шумов)

10. Рекуррентные нейронные сети 3

Пример: синхронный машинный перевод



© Соколинский Л.Б. Глубокие нейронные сети 16.12.2023

Типы приложений

a) 1:1

b) 1:M

c) M:1

d) 1xM : 1xM

e) M:M

a) Один вход, один выход без запоминания состояния (распознавание изображений)

b) Последовательность на выходе (вход – картинка, выход – ее словесное описание)

c) Последовательность на входе (вход – рецензия на фильм, выход – оценка по десятибалльной шкале)

d) Последовательность на входе, последовательность на выходе (перевод текста с одного языка на другой)

e) Синхронизированные последовательности на входе и выходе (разметка смены кадра на видео; очистка аудиопотока от посторонних шумов)

■ Выходной вектор

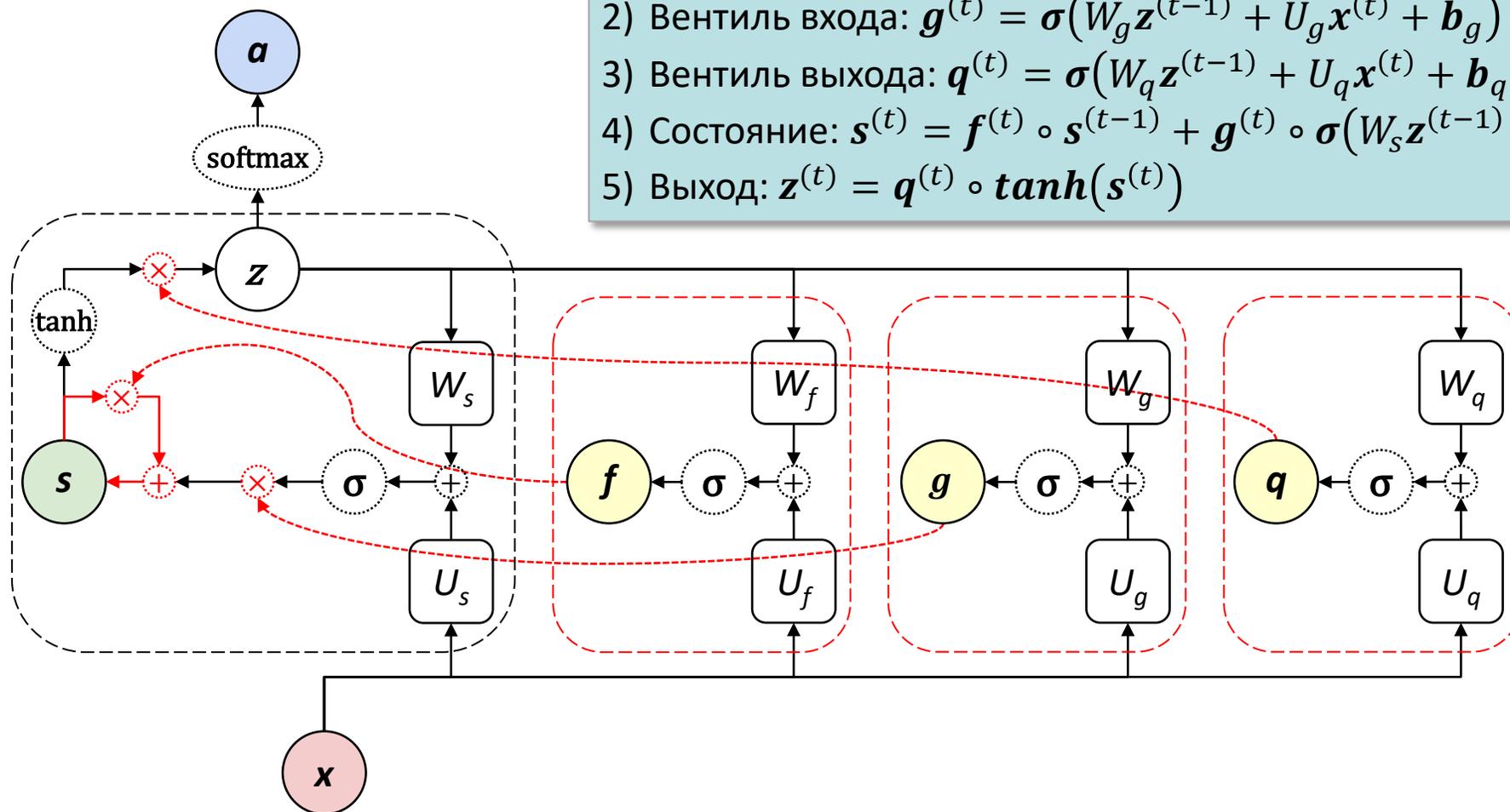
■ Вектор состояния

■ Входной вектор

3

Модель долгой краткосрочной памяти (long short-term memory – LSTM)

- 1) Вентиль забывания: $f^{(t)} = \sigma(W_f z^{(t-1)} + U_f x^{(t)} + b_f)$
- 2) Вентиль входа: $g^{(t)} = \sigma(W_g z^{(t-1)} + U_g x^{(t)} + b_g)$
- 3) Вентиль выхода: $q^{(t)} = \sigma(W_q z^{(t-1)} + U_q x^{(t)} + b_q)$
- 4) Состояние: $s^{(t)} = f^{(t)} \circ s^{(t-1)} + g^{(t)} \circ \sigma(W_s z^{(t-1)} + U_s x^{(t)} + b_s)$
- 5) Выход: $z^{(t)} = q^{(t)} \circ \tanh(s^{(t)})$

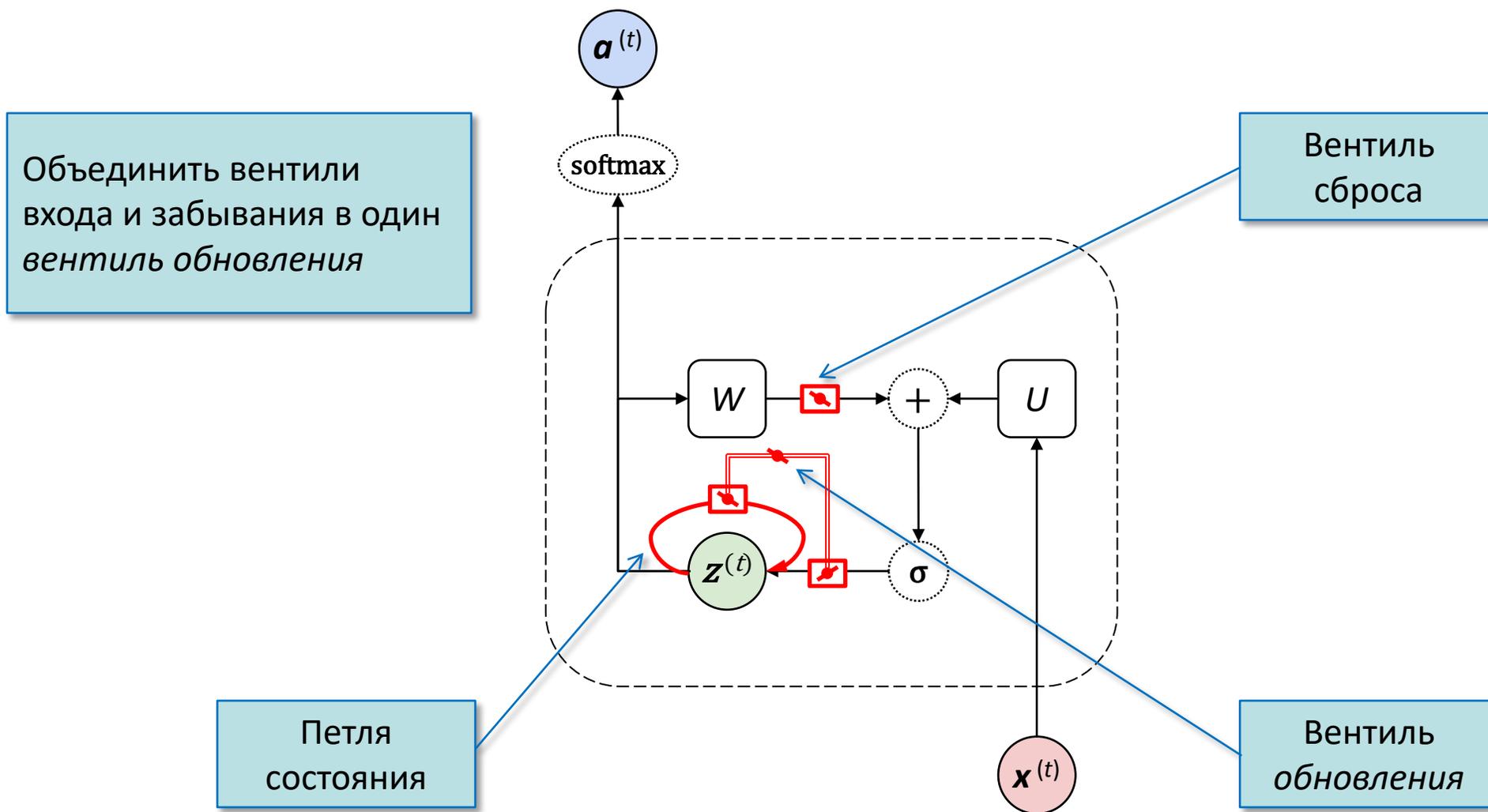


Применения LSTM

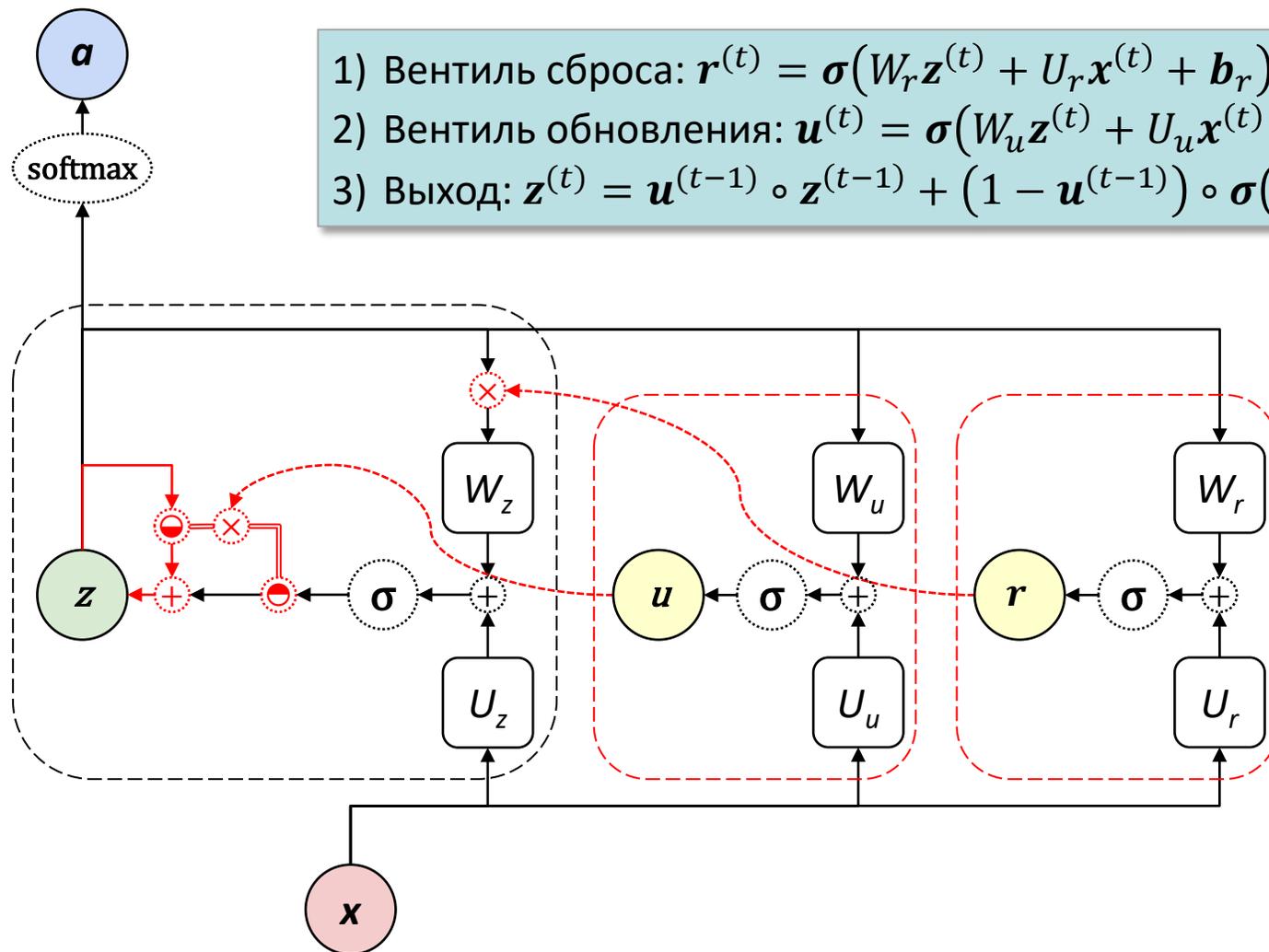
- Большинство применений RNN используют модель LSTM
- Ключевыми компонентами LSTM являются вентиль забывания и функция активации \tanh перед вентилем **ВЫХОДА** (без вентиля забывания качество сразу сильно падает, а без функции активации на выходе этот самый выход теоретически может расти неограниченно, что приводит к нежелательным эффектам)
- Основной недостаток LSTM – большое количество матриц весов (8 против 3 у Simple RNN)

Идея!

Вентильный рекуррентный модуль (Gated recurrent unit – GRU)



Модель GRU



- 1) Вентиль сброса: $r^{(t)} = \sigma(W_r z^{(t)} + U_r x^{(t)} + b_r)$
- 2) Вентиль обновления: $u^{(t)} = \sigma(W_u z^{(t)} + U_u x^{(t)} + b_u)$
- 3) Выход: $z^{(t)} = u^{(t-1)} \circ z^{(t-1)} + (1 - u^{(t-1)}) \circ \sigma(W_z(r^{(t-1)} \circ z^{(t-1)} + U_z x^{(t-1)} + b_z))$

Применения модели GRU

- Практика показывает, что модель GRU практически всегда работает так же хорошо, как LSTM, а параметров у нее гораздо меньше: 6 матриц весов против 8 в LSTM
- GRU быстрее обучается
- GRU требует меньше памяти, что позволяет расширить окно наблюдения

Задача: вычислить $\mathbf{a}^{(2)}$

$$\mathbf{s}^{(0)} = (0, 0)$$

$$\mathbf{x}^{(1)} = (4, 6);$$

$$\mathbf{x}^{(2)} = (9, 7)$$

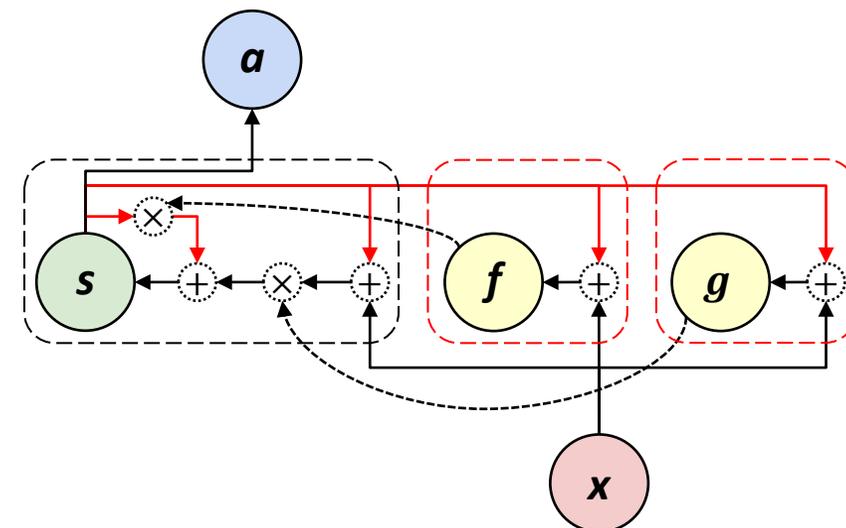
$$\mathbf{f}^{(1)} = \mathbf{g}^{(1)} = \mathbf{x}^{(1)} + \mathbf{s}^{(0)} = (4, 6)$$

$$\mathbf{s}^{(1)} = (\mathbf{x}^{(1)} + \mathbf{s}^{(0)}) \circ \mathbf{g}^{(1)} + \mathbf{s}^{(0)} \circ \mathbf{f}^{(1)} = (16, 36)$$

$$\mathbf{f}^{(2)} = \mathbf{g}^{(2)} = \mathbf{x}^{(2)} + \mathbf{s}^{(1)} = (25, 43)$$

$$\mathbf{s}^{(2)} = (\mathbf{x}^{(2)} + \mathbf{s}^{(1)}) \circ \mathbf{g}^{(2)} + \mathbf{s}^{(1)} \circ \mathbf{f}^{(2)} = (1025, 3397)$$

$$\mathbf{a}^{(2)} = \mathbf{s}^{(2)} = (1025, 3397)$$



Конец лекции 10

Вспомогательные слайды

Использование в качестве функции потерь функции перекрестной энтропии

$$C = - \sum_j y_j \ln a_j^L$$

Свойства функции потерь:

1) $C \geq 0$

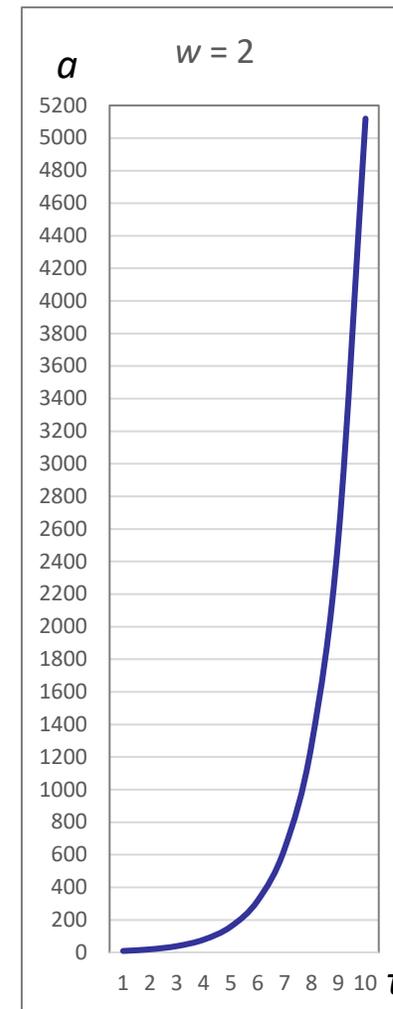
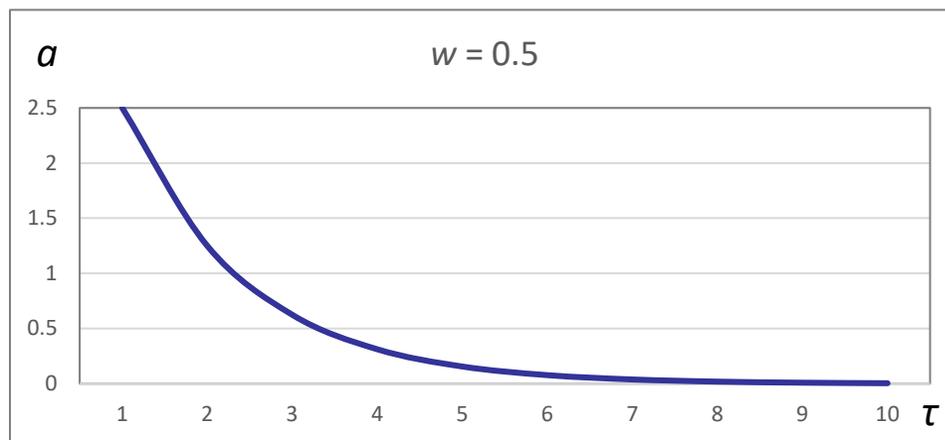
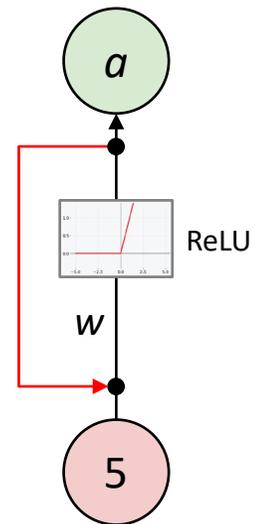
2) Минимум C достигается при $\mathbf{a}^L = \mathbf{y}$:

$$\min_{\mathbf{a}^L} C = - \sum_j y_j \ln y_j$$

Вычисление ошибки δ_i^L для выходного слоя

$$\begin{aligned}
 \delta_i^L &= \frac{\partial C}{\partial z_i^L} = \sum_{j=1}^n \frac{\partial C}{\partial a_j^L} \cdot \frac{\partial a_j^L}{\partial z_i^L} = \sum_{j=1}^n -\frac{y_j}{a_j^L} \cdot \frac{\partial a_j^L}{\partial z_i^L} = -\sum_{j=1}^n \frac{y_j}{a_j^L} \cdot \frac{\partial a_j^L}{\partial z_i^L} \\
 &= -\frac{y_i}{a_i^L} \cdot \frac{\partial a_i^L}{\partial z_i^L} - \sum_{\substack{j=1 \\ j \neq i}}^n \frac{y_j}{a_j^L} \cdot \frac{\partial a_j^L}{\partial z_i^L} \\
 &= -\frac{y_i}{a_i^L} a_i^L (1 - a_i^L) - \sum_{\substack{j=1 \\ j \neq i}}^n -\frac{y_j}{a_j^L} \cdot a_i^L a_j^L = -y_i + y_i a_i^L + \sum_{\substack{j=1 \\ j \neq i}}^n y_j a_i^L \\
 &= -y_i + \sum_{j=1}^n y_j a_i^L = -y_i + a_i^L \sum_{j=1}^n y_j \\
 &= -y_i + a_i^L = a_i^L - y_i
 \end{aligned}$$

Исчезающие и взрывающиеся градиенты



Simple RNN

