

Нейронные сети

Обработка текстов, лекция 6

18.10.2017

Константин Архипенко

Содержание

- 1 Устройство
- 2 Обучение
 - Стохастический градиентный спуск (SGD)
 - Вычисление градиента в графе
 - Модификации SGD
- 3 Продвинутое архитектуры
 - Сверточные нейронные сети – кратко
 - Рекуррентные нейронные сети
- 4 Программное обеспечение и пример кода

1 Устройство

2 Обучение

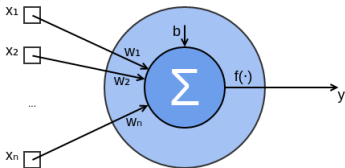
- Стохастический градиентный спуск (SGD)
- Вычисление градиента в графе
- Модификации SGD

3 Продвинутое архитектуры

- Сверточные нейронные сети – кратко
- Рекуррентные нейронные сети

4 Программное обеспечение и пример кода

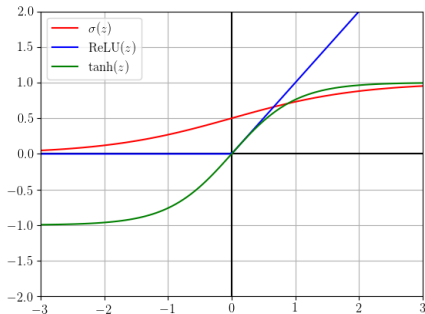
Нейрон



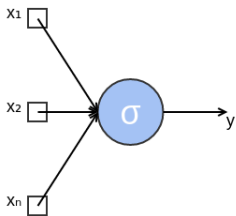
- $y(\mathbf{x}) = f(\mathbf{x}\mathbf{w}^T + b) = f(\sum_i x_i w_i + b)$
- $f(\cdot)$ – функция активации, чаще нелинейная
- Обучение нейрона есть настройка его параметров \mathbf{w} и b

Функции активации

- $\sigma(z) = \frac{1}{1+e^{-z}}$
- $\text{ReLU}(z) = \max(0, z)$
- $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = 2\sigma(2z) - 1$
- Бывают и другие

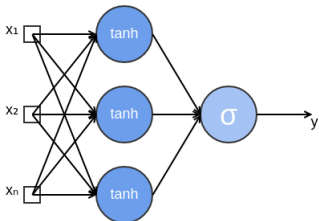


Нейрон с функцией активации $\sigma(z)$



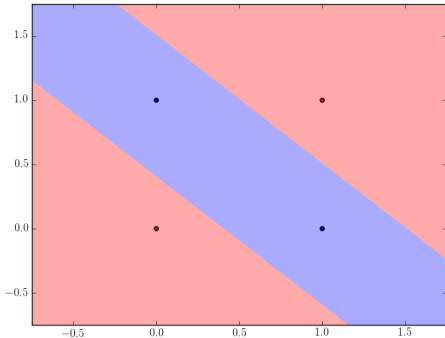
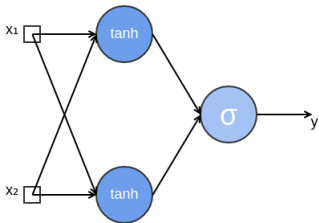
- $y(\mathbf{x}) = \sigma(\mathbf{x}\mathbf{w}^T + b)$
- Получается не что иное, как **логистическая регрессия**

Добавим слой нейронов



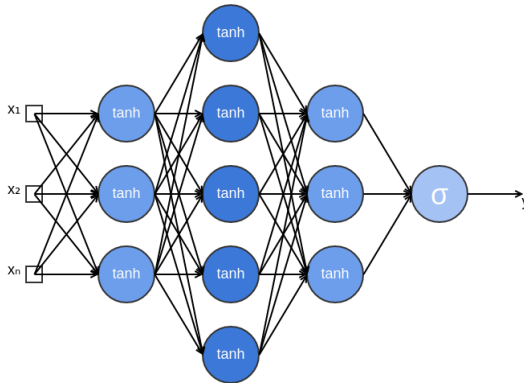
- $\mathbf{h} = \tanh(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)$
 $y = \sigma(\mathbf{h}\mathbf{w}_2^T + b_2)$
- Заведомо более мощная модель, чем логистическая регрессия

XOR

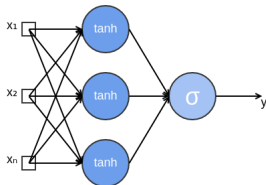


- Линейные классификаторы на такое не способны

Multilayer perceptron (MLP)

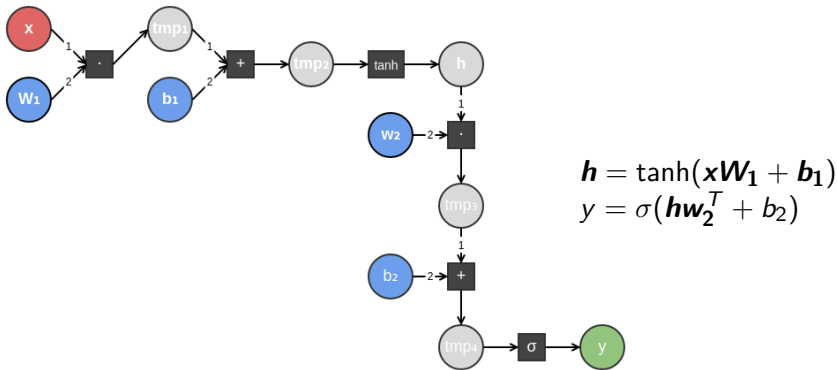


- Вместо подобных картинок



для визуализации архитектуры нейросетей можно нужно использовать более компактные и удобные **графы вычислений**, максимально отражающие факт векторизации

Пример графа вычислений



- **Входные**, промежуточные и **выходные** вершины, операции и **настраиваемые параметры**

Мультиномиальная логистическая регрессия

- Вероятность класса i :

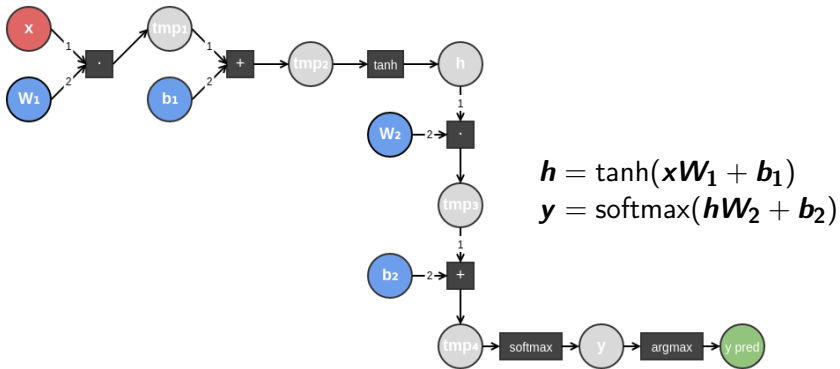
$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

- Аналогичным образом можно усилить модель добавлением слоя:

$$\mathbf{h} = \tanh(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)$$

$$\mathbf{y} = \text{softmax}(\mathbf{h}\mathbf{W}_2 + \mathbf{b}_2)$$

Граф вычислений



- После предсказания вектора y вероятностей классов выбирается наиболее подходящий класс $y_{pred} = \arg \max_i y_i$

Содержание

1 Устройство

2 Обучение

- Стохастический градиентный спуск (SGD)
- Вычисление градиента в графе
- Модификации SGD

3 Продвинутое архитектуры

- Сверточные нейронные сети – кратко
- Рекуррентные нейронные сети

4 Программное обеспечение и пример кода

- Метод максимального правдоподобия:

$$\mathcal{L}(\mathbf{X}, \mathbf{y}_{true}, \mathbf{w}, b) = -\frac{1}{N} \sum_{n=1}^N \log P(y_{true}^{[n]} | \mathbf{x}^{[n]}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w}, b}$$

$$P(1 | \mathbf{x}^{[n]}) = 1 - P(0 | \mathbf{x}^{[n]}) = y^{[n]} = \sigma(\mathbf{x}^{[n]} \mathbf{w}^T + b)$$

Бинарная кросс-энтропия

- Перепишем $-\log P(y_{true}^{[n]} | \mathbf{x}^{[n]})$:

$$-\log P(y_{true}^{[n]} | \mathbf{x}^{[n]}) = -(y_{true}^{[n]} \log y^{[n]} + (1 - y_{true}^{[n]}) \log(1 - y^{[n]}))$$

- В правой части находится бинарная кросс-энтропия $\text{BCE}(y_{true}^{[n]}, y^{[n]})$
- Таким образом, максимизация правдоподобия есть то же самое, что минимизация кросс-энтропии:

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N \text{BCE}(y_{true}^{[n]}, y^{[n]}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w}, b}$$

$$y^{[n]} = \sigma(\mathbf{x}^{[n]} \mathbf{w}^T + b)$$

Мультиномиальная логистическая регрессия

- Для многоклассовой классификации бинарная кросс-энтропия заменяется на обычную:

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N \text{CrossEntropy}(\mathbf{y}_{true}^{[n]}, \mathbf{y}^{[n]}) + \frac{\lambda}{2} \|\mathbf{W}\|_F^2 \rightarrow \min_{\mathbf{W}, \mathbf{b}}$$

$$\mathbf{y}^{[n]} = \text{softmax}(\mathbf{x}^{[n]}\mathbf{W} + \mathbf{b})$$

$$\text{CrossEntropy}(\mathbf{y}_{true}^{[n]}, \mathbf{y}^{[n]}) = -\mathbf{y}_{true}^{[n]} \log(\mathbf{y}^{[n]})^T$$

- Здесь $\mathbf{y}_{true}^{[n]}$ – вектор из нулей, кроме единицы в позиции i , соответствующей истинному классу объекта $\mathbf{x}^{[n]}$ (**one-hot encoding**)

Градиентный спуск

- Рассматриваем бинарную логистическую регрессию
- Для минимизации \mathcal{L} есть много методов (включая методы второго порядка), наиболее простым и пригодным для больших данных является **градиентный спуск**
- На каждой итерации t обновляем параметры \mathbf{w} и b ($\alpha_t > 0$):

$$w_i \leftarrow w_i - \alpha_t \nabla_{w_i} \mathcal{L}$$

$$b \leftarrow b - \alpha_t \nabla_b \mathcal{L}$$

- Градиент \mathcal{L} выписывается в явном виде, но найти стационарную точку аналитически не получится

Стохастический градиентный спуск (SGD)

- Вычисление градиента \mathcal{L} на каждой итерации t дорого, поскольку в него входит сумма по всем объектам выборки
- Вместо этого предлагается на каждой итерации брать случайный объект $(\mathbf{x}^{[n]}, y_{true}^{[n]})$ выборки и считать градиент потери $\mathcal{L}^{[n]}$ только на этом объекте:

$$\mathcal{L}^{[n]} = \text{BCE}(y_{true}^{[n]}, y^{[n]}) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

$$w_i \leftarrow w_i - \alpha_t \nabla_{w_i} \mathcal{L}^{[n]}$$

$$b \leftarrow b - \alpha_t \nabla_b \mathcal{L}^{[n]}$$

- Есть стратегии уменьшения α_t с теоретическими гарантиями сходимости, но им следуют далеко не всегда

SGD с минибатчами

- Разобьем выборку на минибатчи (группы обычно одинакового размера), а затем будем на каждой итерации брать минибатч b и считать градиент потери $\mathcal{L}^{[b]}$:

$$\mathcal{L}^{[b]} = \frac{1}{|b|} \sum_{n \in b} \text{BCE}(y_{true}^{[n]}, y^{[n]}) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- По сравнению с простым SGD появляется возможность распараллеливания вычислений
- Для обучения нейронных сетей на GPU минибатчи используются почти всегда

Обучение нейронной сети

- Рассмотрим подробно обучение нейронной сети для бинарной классификации:

$$\mathbf{h} = \tanh(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)$$
$$y = \sigma(\mathbf{h}\mathbf{w}_2^T + b_2)$$

- Используем бинарную кросс-энтропию в качестве функции потерь и SGD
- Опустим для простоты L_2 -регуляризацию (но она очень желательна!) и пока забудем о минибатчах
- Потребуется ввести понятие якобиана и вспомнить цепное правило

- Пусть $\mathbf{v} = f(\mathbf{u})$. Тогда $\nabla_{\mathbf{u}} \mathbf{v}$ – матрица:

$$(\nabla_{\mathbf{u}} \mathbf{v})_{ij} = \nabla_{u_j} v_i$$

- Обобщим понятие якобиана

- Пусть $f: \mathbb{R}^m \rightarrow \mathbb{R}$, $v = f(\mathbf{u})$. Тогда $\nabla_{\mathbf{u}} v$ – вектор:

$$(\nabla_{\mathbf{u}} v)_i = \nabla_{u_i} v$$

- Пусть $f: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$, $v = f(\mathbf{U})$. Тогда $\nabla_{\mathbf{U}} v$ – матрица:

$$(\nabla_{\mathbf{U}} v)_{ij} = \nabla_{u_{ij}} v$$

- Пусть $f: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^n$, $\mathbf{v} = f(\mathbf{U})$. Тогда $\nabla_{\mathbf{U}} \mathbf{v}$ – 3-мерный тензор:

$$(\nabla_{\mathbf{U}} \mathbf{v})_{ijk} = \nabla_{u_{jk}} v_i$$

Цепное правило

- Если $w = g(v)$, $v = f(u)$, а функции дифференцируемы, то

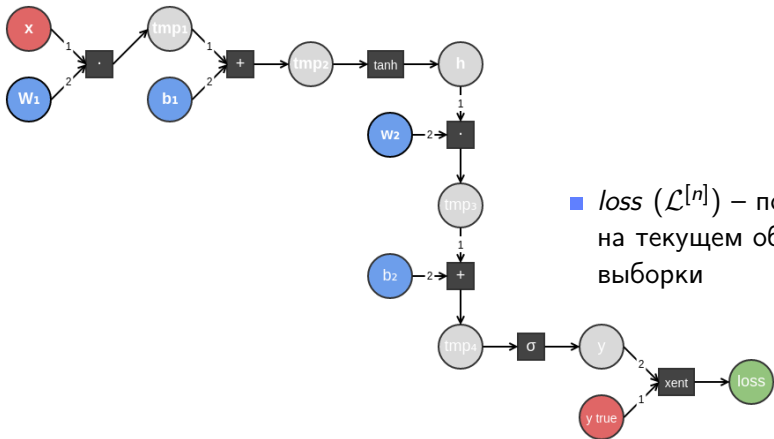
$$\nabla_u w = \nabla_v w \nabla_u v$$

- Справедливо и для наших якобианов
- Например, если $w = g(\mathbf{v})$, $\mathbf{v} = f(\mathbf{U})$, то

$$\nabla_{\mathbf{U}} w = \nabla_{\mathbf{v}} w \nabla_{\mathbf{U}} \mathbf{v}$$

- В последнем равенстве вектор-строка умножается на 3-мерный тензор – получается матрица

Граф вычислений нейронной сети

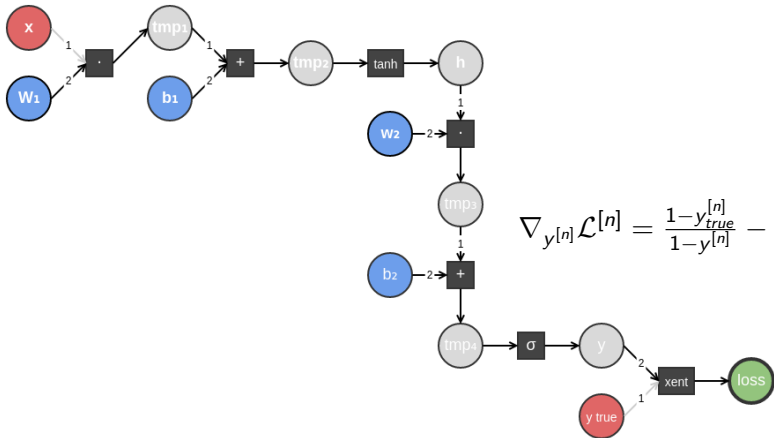


Backpropagation

- Пусть мы вычислили значения вершин от **входных** к **выходным** на текущем объекте выборки
- Наша цель – вычислить градиент $\mathcal{L}^{[n]}$ по всем **настраиваемым параметрам** нейросети, чтобы затем сделать шаг SGD
- Для этого будем идти назад от вершины *loss* к вершинам-параметрам, вычисляя по цепному правилу градиенты $\mathcal{L}^{[n]}$ по круглым вершинам на нашем пути

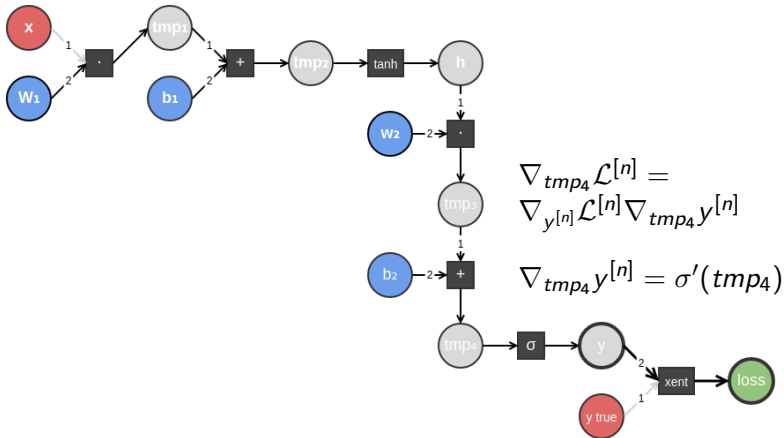
Backpropagation

■ $\mathcal{L}^{[n]} = -(y_{true}^{[n]} \log y^{[n]} + (1 - y_{true}^{[n]}) \log(1 - y^{[n]}))$



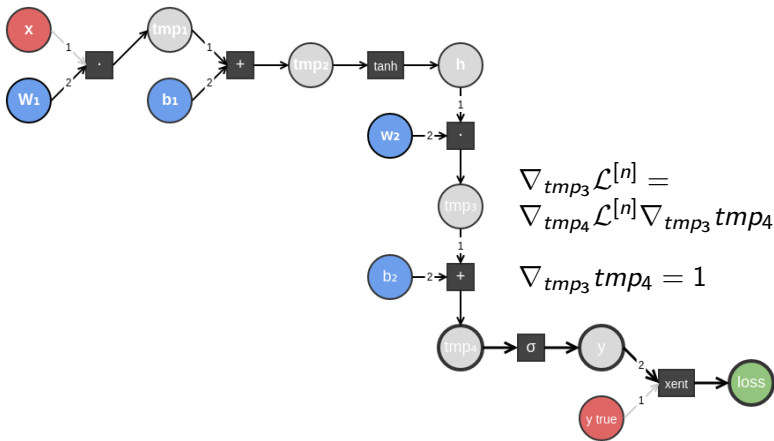
Backpropagation

■ $y = \sigma(tmp_4)$



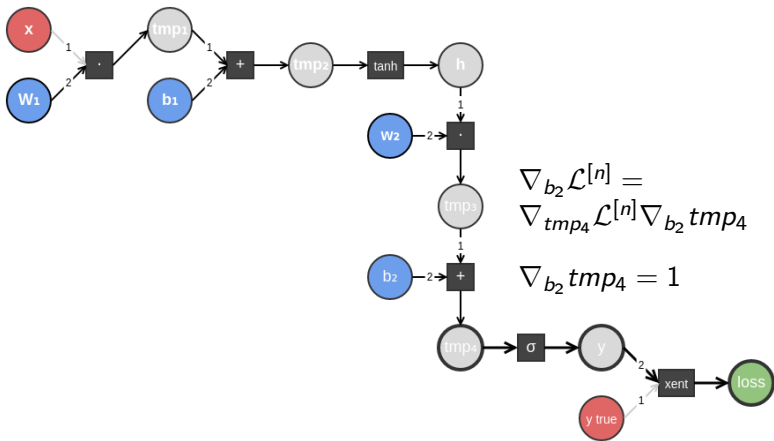
Backpropagation

■ $tmp_4 = tmp_3 + b_2$



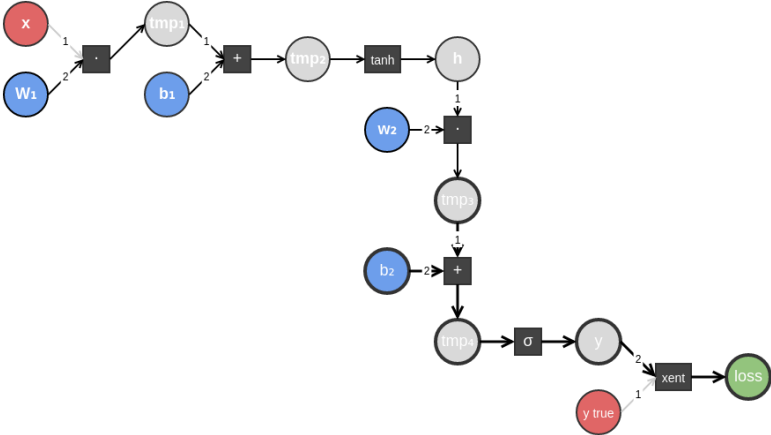
Backpropagation

■ $tmp_4 = tmp_3 + b_2$



Backpropagation

■ И так далее



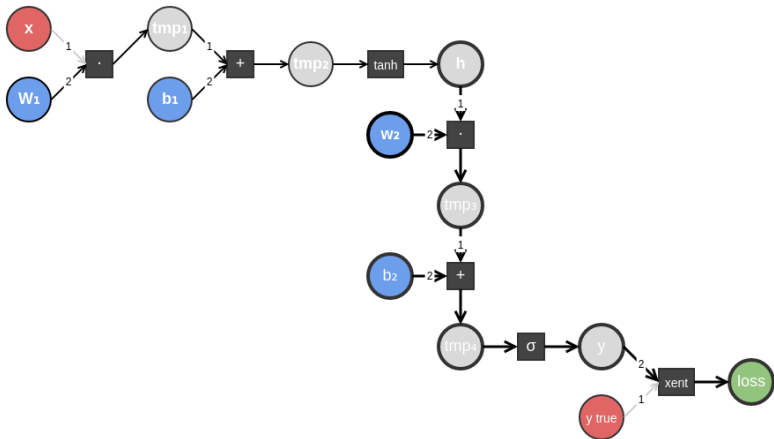
Автоматическое дифференцирование

- Вершина tmp_3 знает градиент $\nabla_{tmp_3} \mathcal{L}^{[n]}$ и отправляет его операции-предку
- Прямоугольник-операция, получив градиент по tmp_3 , умножает его на свои якобианы $\nabla_{\mathbf{h}} tmp_3$ и $\nabla_{\mathbf{w}_2} tmp_3$ (подставляя текущие значения вершин) и отправляет результаты вершинам \mathbf{h} и \mathbf{w}_2 соответственно:

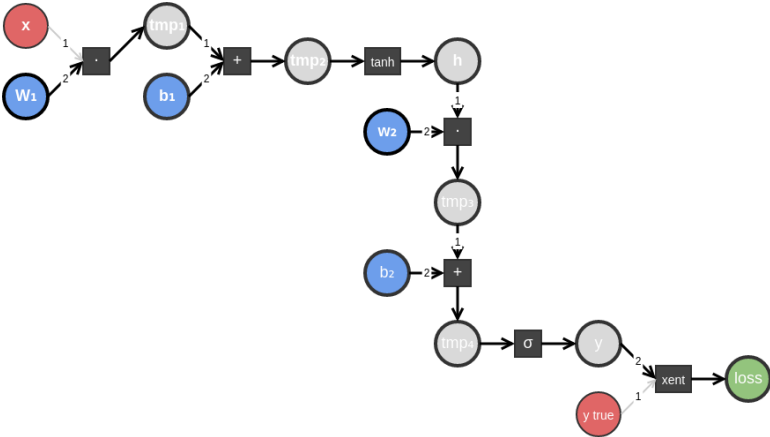
$$\nabla_{\mathbf{h}} \mathcal{L}^{[n]} = \nabla_{tmp_3} \mathcal{L}^{[n]} \nabla_{\mathbf{h}} tmp_3 = \nabla_{tmp_3} \mathcal{L}^{[n]} \mathbf{w}_2$$

$$\nabla_{\mathbf{w}_2} \mathcal{L}^{[n]} = \nabla_{tmp_3} \mathcal{L}^{[n]} \nabla_{\mathbf{w}_2} tmp_3 = \nabla_{tmp_3} \mathcal{L}^{[n]} \mathbf{h}$$

Backpropagation



Backpropagation



Минибатчи

- На входе теперь не вектор $\mathbf{x}^{[n]}$, а матрица $\mathbf{X}^{[b]}$ с равным размером минибатча количеством строк
- \mathbf{tmp}_1 , \mathbf{tmp}_2 , \mathbf{h} превращаются в матрицы
- \mathbf{tmp}_3 , \mathbf{tmp}_4 – в векторы
- $y^{[n]}$ и $y_{true}^{[n]}$ – в векторы $\mathbf{y}^{[b]}$ и $\mathbf{y}_{true}^{[b]}$
- К графу добавляется новая **выходная** вершина – усредненная потеря $\mathcal{L}^{[b]}$ на минибатче
- В итоге распараллеливается не только предсказание нейросети, но и backpropagation
- Минибатчи используются и на тестовых данных после обучения

Общая схема обучения нейросети

- Параметры нейросети инициализируются случайными малыми значениями (например, Glorot uniform¹), смещения ***b*** – обычно нулями
- Каждую эпоху обучающая выборка случайным образом разбивается на минибатчи
- Для каждого минибатча вычисляют градиент потери и обновляют параметры
- Обучение заканчивается по достижении определенного числа эпох или при прекращении снижения функции потерь на валидационной выборке

¹Glorot & Bengio “Understanding the difficulty of training deep feedforward neural networks” (2010),
<http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf>

Оптимизация функции потерь

- В отличие от логистической регрессии, обучение нейросети есть оптимизация очень сложной функции с большим количеством локальных минимумов
- Применять SGD становится трудно
- Рассмотрим улучшения SGD для более быстрого нахождения более оптимальных параметров

Momentum

- Пусть θ – совокупность настраиваемых параметров нейросети
- $\mathbf{v}_t = \gamma \mathbf{v}_{t-1} - \alpha_t \nabla_{\theta} \mathcal{L}^{[b]}|_{\theta=\theta_{t-1}}$
 $\theta_t = \theta_{t-1} + \mathbf{v}_t$
- По функции потерь катится мяч, ускоряясь на спусках
- Чуть более умная версия мяча (Nesterov accelerated gradient):
 $\mathbf{v}_t = \gamma \mathbf{v}_{t-1} - \alpha_t \nabla_{\theta} \mathcal{L}^{[b]}|_{\theta=\theta_{t-1} + \gamma \mathbf{v}_{t-1}}$

- SGD



- SGD + momentum



Sebastian Ruder "An overview of gradient descent optimization algorithms", <http://ruder.io/optimizing-gradient-descent>

- Хотелось бы, чтобы скорость \mathbf{v} для каждого из направлений была примерно одинаковой
- Решение $\mathbf{v}_t = -\alpha_t \frac{\nabla_{\theta} \mathcal{L}^{[b]}}{|\nabla_{\theta} \mathcal{L}^{[b]}| + \epsilon} \Big|_{\theta = \theta_{t-1}}$ не очень подходит: представим два минибатча с противоположно направленными градиентами разной величины
- Возьмем **скользящее среднее**:

$$\mathbf{a}_t = 0.9\mathbf{a}_{t-1} + 0.1|\nabla_{\theta} \mathcal{L}^{[b]}|_{\theta = \theta_{t-1}} \quad \mathbf{v}_t = -\alpha_t \frac{\nabla_{\theta} \mathcal{L}^{[b]}}{\mathbf{a}_t + \epsilon} \Big|_{\theta = \theta_{t-1}}$$

- Еще лучше: **корень (R)** из скользящего **среднего (M) квадратов (S)**:

$$\mathbf{a}_t = 0.9\mathbf{a}_{t-1} + 0.1(\nabla_{\theta} \mathcal{L}^{[b]} \Big|_{\theta = \theta_{t-1}})^2 \quad \mathbf{v}_t = -\alpha_t \frac{\nabla_{\theta} \mathcal{L}^{[b]}}{\sqrt{\mathbf{a}_t + \epsilon}} \Big|_{\theta = \theta_{t-1}}$$

- Использование корня из среднего квадратов имеет связь с кривизной функции потерь по каждому из направлений
- Проблема разной величины градиентов по разным направлениям решается изящнее, скорость не становится слишком большой
- Небольшой **недостаток**: для хранения скользящего среднего квадратов нужна дополнительная память (на GPU)
- Если заменить градиент в числителе на его скользящее среднее, получим Adam². Это **не** RMSprop + momentum

²Kingma & Ba "Adam: A Method for Stochastic Optimization" (2014), <https://arxiv.org/abs/1412.6980>

Проблема выбора α_t

- Частично решается методами RMSprop и Adam (SGD + momentum – в меньшей степени)
- Некоторые исследователи³ рекомендуют обучать нейросеть некоторое количество эпох (фиксированное или до прекращения снижения значения функции потерь) с фиксированным α_t , а затем **многokrатно уменьшать** (до 10 раз)

³He et al. “Deep Residual Learning for Image Recognition” (2015), <https://arxiv.org/abs/1512.03385>

Dropout

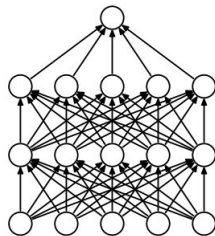
- Каждый элемент тензора обнуляется с вероятностью p либо умножается на $\frac{1}{p}$ с вероятностью $1 - p$
- Пример:

$$\tilde{\mathbf{h}} = \tanh(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)$$

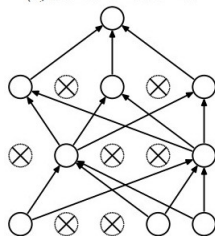
$$\mathbf{h} = \text{dropout}(\tilde{\mathbf{h}})$$

$$y = \sigma(\mathbf{h}\mathbf{w}_2^T + b_2)$$

- Более эффективным обобщением является **Gaussian dropout** (с умножением на нормальный шум с матожиданием 1 вместо Бернулли)



(a) Standard Neural Net



(b) After applying dropout.

Dropout

- Помогает от переобучения
- При выводе (работе на тестовых данных) dropout выключают
- Сейчас на практике чаще используют различные виды **нормализации** (batch/layer/weight) и другие изоциренные способы регуляризации (**zoneout**⁴ и т. п.)

⁴Krueger et al. "Zoneout: Regularizing RNNs by Randomly Preserving Hidden Activations" (2016), <https://arxiv.org/abs/1606.01305>

Содержание

1 Устройство

2 Обучение

- Стохастический градиентный спуск (SGD)
- Вычисление градиента в графе
- Модификации SGD

3 Продвинутое архитектуры

- Сверточные нейронные сети – кратко
- Рекуррентные нейронные сети

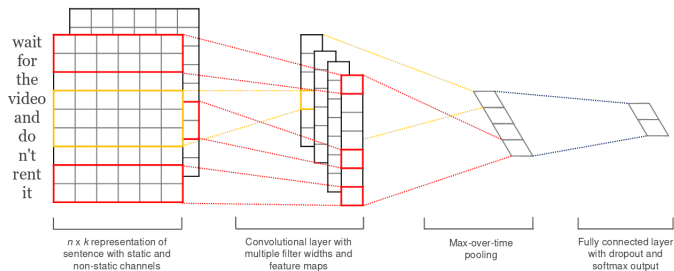
4 Программное обеспечение и пример кода

Сверточные нейронные сети

- Изначально использовались в задачах компьютерного зрения и широко используются сейчас (например, многочисленные варианты ResNet)
- Иногда помогают и при обработке текстов
- Рассмотрим одно из возможных применений сверточных сетей для классификации коротких текстов (предложений)⁵

⁵Kim “Convolutional Neural Networks for Sentence Classification” (2014), <http://www.aclweb.org/anthology/D14-1181>

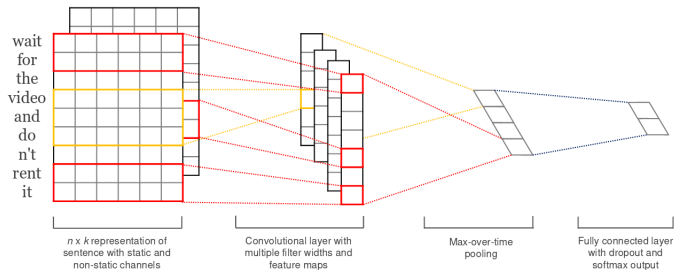
Классификация текстов с помощью CNN



- На вход подается матрица **векторных представлений** (например, **word2vec**⁶, где векторные представления могут строиться автоматически по большому неразмеченному корпусу текстов) слов предложения

⁶Mikolov et al. "Distributed Representations of Words and Phrases and their Compositionality" (2013), <https://arxiv.org/abs/1310.4546>

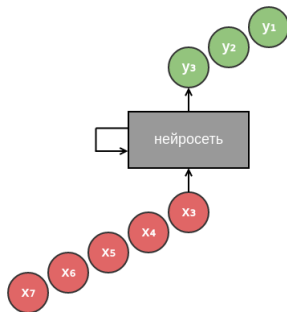
Классификация текстов с помощью CNN



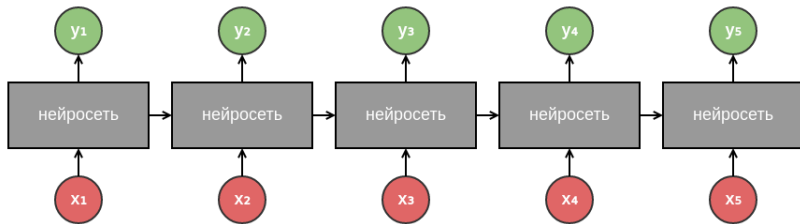
- Свертка – скалярное умножение фрагмента входа на обучаемый (тем же backpropagation) **фильтр**
- Используются фильтры с разной шириной по словам, pooling осуществляется выбором максимального значения

Рекуррентные нейронные сети: мотивация

- Обработка **последовательных данных**: предложений на естественном языке, временных рядов и т. п.
- Различная длина последовательностей
- Взаимосвязи между элементами последовательностей, возможно, далеко отстоящими
- Рекуррентная нейросеть (RNN) **аккумулирует** информацию о предыдущих элементах подаваемой на вход последовательности

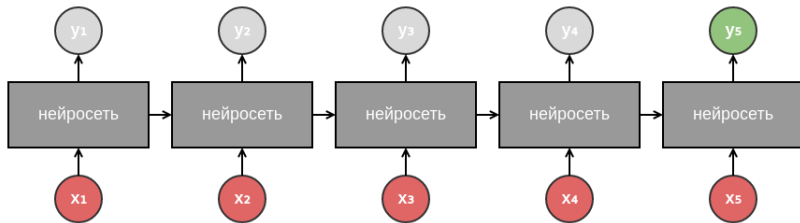


RNN, представленная в другом виде



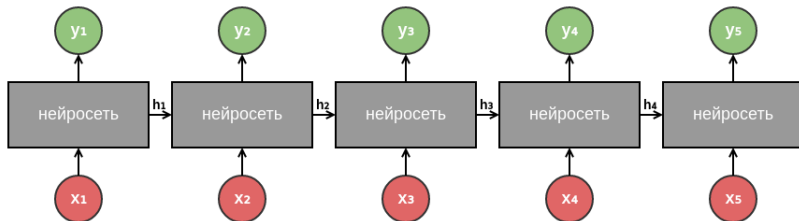
- Можно использовать как все элементы последовательности выходов (например, в задаче определения частей речи)...

RNN, представленная в другом виде



- так и только последний, игнорируя и не вычисляя градиент потери по остальным (например, в задаче классификации предложений)

Простая RNN⁷



- $h_t = f_h(x_t \mathbf{W} + h_{t-1} \mathbf{U} + \mathbf{b})$ – скрытое состояние
- $y_t = f_y(h_t \mathbf{V} + \mathbf{c})$
- $h_0 = 0$

⁷Elman "Finding Structure in Time" (1990),
<http://psych.colorado.edu/~kimlab/Elman1990.pdf>

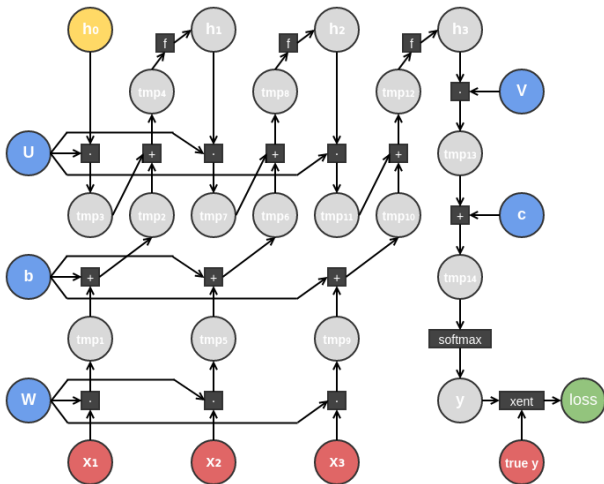
Вычисление градиента в RNN

- Пусть мы хотим классифицировать короткие тексты
- Возьмем обученный word2vec и представим каждый текст как матрицу $T \times D$, где T – количество слов, D – размерность word2vec
- Используем RNN:

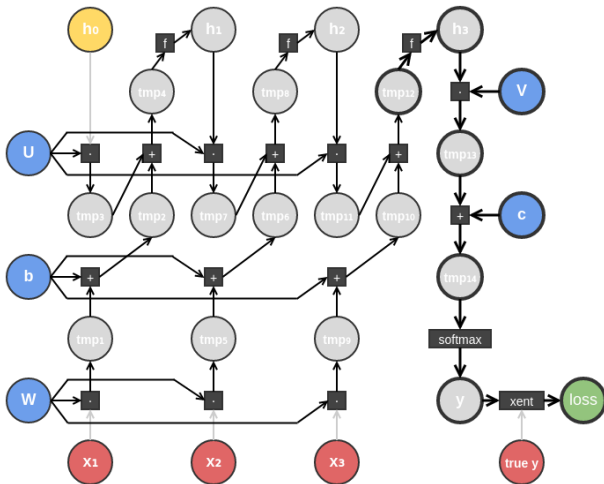
$$\begin{aligned}h_t &= \tanh(\mathbf{x}_t \mathbf{W} + \mathbf{h}_{t-1} \mathbf{U} + \mathbf{b}), t = 1 \dots T \\ \mathbf{y} &= \text{softmax}(\mathbf{h}_T \mathbf{V} + \mathbf{c}) - \text{вероятности классов} \\ h_0 &= \mathbf{0}\end{aligned}$$

- И кросс-энтропию в соответствии с принципом максимального правдоподобия

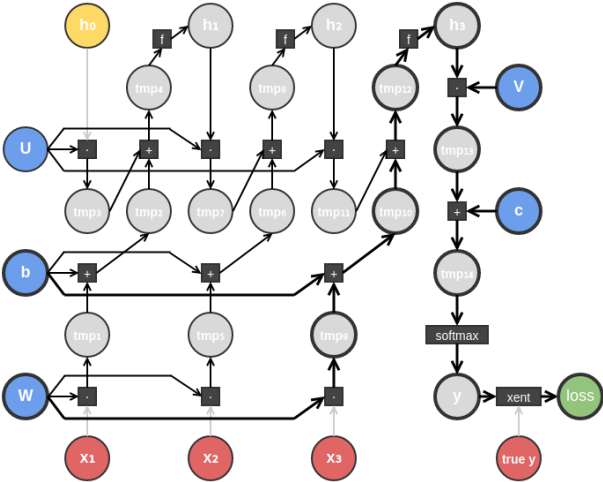
Вычисление градиента в RNN



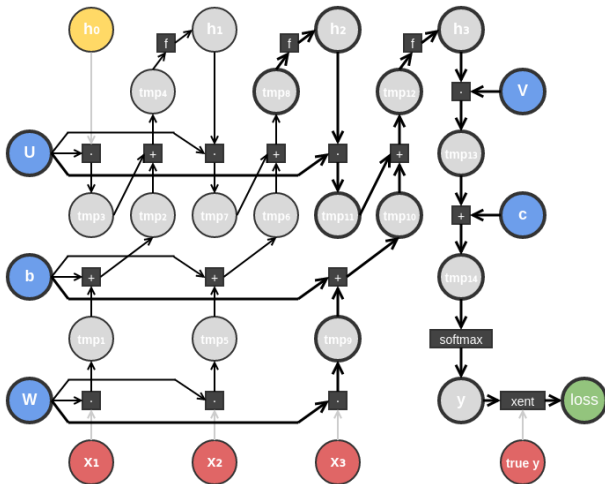
Вычисление градиента в RNN



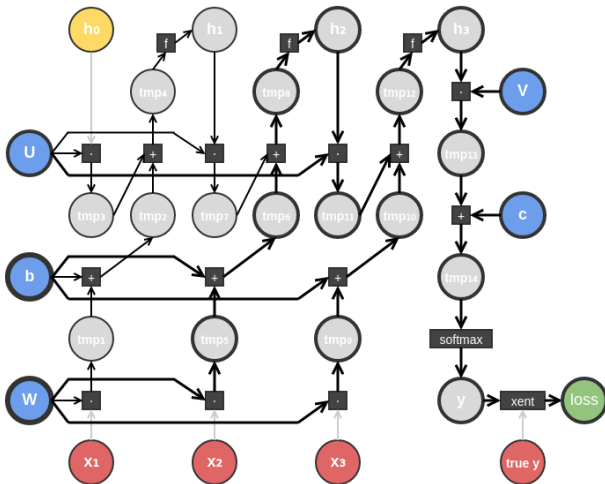
Вычисление градиента в RNN



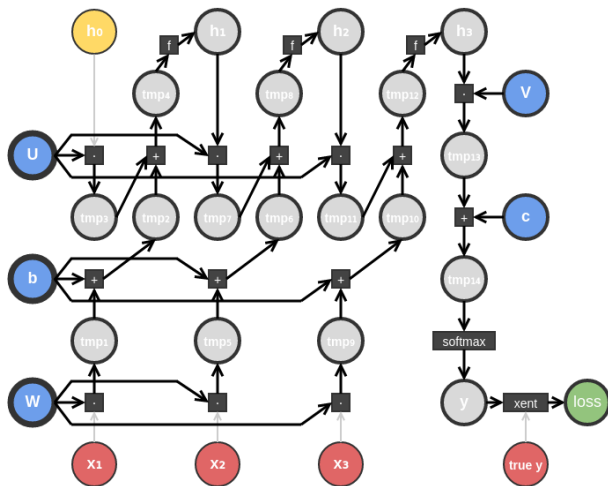
Вычисление градиента в RNN



Вычисление градиента в RNN



Вычисление градиента в RNN



Вычисление градиента в RNN

- Вершина W получает части градиента потери по себе от трех вершин tmp_1 , tmp_5 и tmp_9 – эти части суммируются
- Аналогично для вершин U и b

RNN и минибатчи

- Последовательности на входе RNN имеют разную длину
- Чтобы сформировать минибатч (3-мерный тензор), придется выровнять (`padding`) последовательности среди всех объектов минибатча
- При этом выходы RNN для частей входов, относящихся к выравниванию, следует проигнорировать

Затухание (vanishing) градиента в RNN

- $\mathbf{h}_t = f_h(\mathbf{x}_t \mathbf{W} + \mathbf{h}_{t-1} \mathbf{U} + \mathbf{b})$
 $\mathbf{y}_t = f_y(\mathbf{h}_t \mathbf{V} + \mathbf{c})$
 $\mathbf{h}_0 = \mathbf{0}$

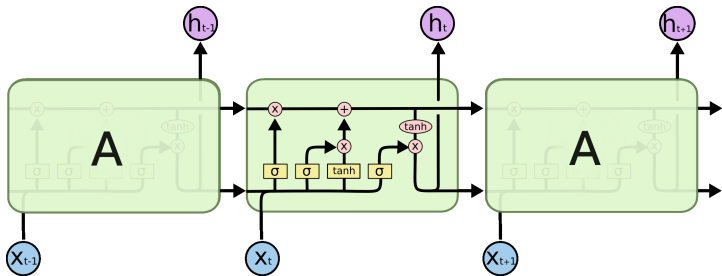
$$\nabla_{\mathbf{x}_5} \mathbf{h}_5 = \text{diag}(f'_h(\mathbf{x}_5 \mathbf{W} + \mathbf{h}_4 \mathbf{U} + \mathbf{b})) \mathbf{W}^T$$

$$\nabla_{\mathbf{x}_2} \mathbf{h}_5 = \nabla_{\mathbf{h}_4} \mathbf{h}_5 \nabla_{\mathbf{h}_3} \mathbf{h}_4 \nabla_{\mathbf{h}_2} \mathbf{h}_3 \cdot \text{diag}(f'_h(\mathbf{x}_2 \mathbf{W} + \mathbf{h}_1 \mathbf{U} + \mathbf{b})) \mathbf{W}^T$$

$$\nabla_{\mathbf{h}_{t-1}} \mathbf{h}_t = \text{diag}(f'_h(\mathbf{x}_t \mathbf{W} + \mathbf{h}_{t-1} \mathbf{U} + \mathbf{b})) \mathbf{U}^T$$

- Если $|f'_h(\cdot)| \|\mathbf{U}\|$ существенно меньше 1, то произойдет **затухание градиента**, и RNN не сможет запоминать информацию на долгое время (градиент по давним входам почти нулевой)
- Также опасна обратная ситуация (**взрыв градиента**)

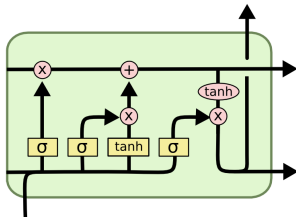
Long Short-Term Memory



Christopher Olah "Understanding LSTM Networks",
<http://colah.github.io/posts/2015-08-Understanding-LSTMs>

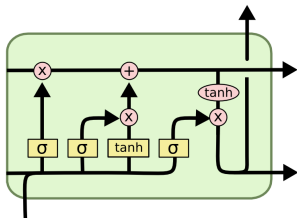
Long Short-Term Memory

- $f_t = \sigma(x_t W_f + h_{t-1} U_f + b_f)$
 $i_t = \sigma(x_t W_i + h_{t-1} U_i + b_i)$
 $\tilde{c}_t = \tanh(x_t W_c + h_{t-1} U_c + b_c)$
 $o_t = \sigma(x_t W_o + h_{t-1} U_o + b_o)$
 $c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$
 $h_t = o_t \odot \tanh(c_t)$
- LSTM обладает дополнительным скрытым состоянием c , а состояние h одновременно является выходом



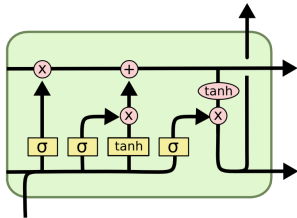
Long Short-Term Memory

- $f_t = \sigma(x_t W_f + h_{t-1} U_f + b_f)$
 $i_t = \sigma(x_t W_i + h_{t-1} U_i + b_i)$
 $\tilde{c}_t = \tanh(x_t W_c + h_{t-1} U_c + b_c)$
 $o_t = \sigma(x_t W_o + h_{t-1} U_o + b_o)$
 $c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$
 $h_t = o_t \odot \tanh(c_t)$
- **Forget gate:** с учетом x_t определяет, какую информацию из c_{t-1} стоит сохранить



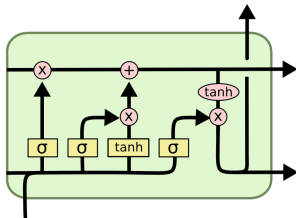
Long Short-Term Memory

- $f_t = \sigma(x_t W_f + h_{t-1} U_f + b_f)$
 $i_t = \sigma(x_t W_i + h_{t-1} U_i + b_i)$
 $\tilde{c}_t = \tanh(x_t W_c + h_{t-1} U_c + b_c)$
 $o_t = \sigma(x_t W_o + h_{t-1} U_o + b_o)$
 $c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$
 $h_t = o_t \odot \tanh(c_t)$
- **Input gate** и **cell gate**: определяют важность информации из x_t и преобразуют важную информацию для сохранения в состоянии c



Long Short-Term Memory

- $f_t = \sigma(x_t W_f + h_{t-1} U_f + b_f)$
 $i_t = \sigma(x_t W_i + h_{t-1} U_i + b_i)$
 $\tilde{c}_t = \tanh(x_t W_c + h_{t-1} U_c + b_c)$
 $o_t = \sigma(x_t W_o + h_{t-1} U_o + b_o)$
 $c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$
 $h_t = o_t \odot \tanh(c_t)$
- **Output gate:** определяет, какую информацию из обновленного состояния c выдать в качестве ответа



LSTM и проблема нестабильного градиента

$$\nabla_{\mathbf{c}_{t-1}} \mathbf{c}_t = \text{diag}(\mathbf{f}_t)$$

- Теперь нейросеть может запоминать информацию надолго в состоянии \mathbf{c}
- Смещение \mathbf{b}_f forget gate желательно инициализировать не вектором из нулей, а вектором из единиц
- Для уменьшения вероятности дестабилизации градиента можно инициализировать матрицы \mathbf{U} случайными ортогональными матрицами, также очень полезно делать [gradient clipping](#)

Содержание

1 Устройство

2 Обучение

- Стохастический градиентный спуск (SGD)
- Вычисление градиента в графе
- Модификации SGD

3 Продвинутое архитектуры

- Сверточные нейронные сети – кратко
- Рекуррентные нейронные сети

4 Программное обеспечение и пример кода

Библиотеки глупинного глубокого обучения

- Theano – некогда популярная, дала развитие современным библиотекам. Поддержка прекращена
- TensorFlow (Google) – сейчас наиболее известная, пригодная для промышленного использования, сложна в освоении и поддержке
- PyTorch – более простая и гибкая, чем TF, активно развивается, пока есть проблемы с промышленным использованием и распределенным обучением очень больших моделей



PYTORCH

Недостатки TensorFlow

- В TensorFlow граф вычислений компилируется целиком в эффективный машинный код (использующий GPU)
- Появляются проблемы⁸ с реализацией некоторых моделей, например, TreeLSTM над деревьями синтаксического разбора произвольной глубины
- Запущенная модель кажется отрезанной от пользователя, с ней трудно контактировать на языке Python
- PyTorch с самого начала лишена этих недостатков

⁸Решаются использованием далеко не самого изящного решения TensorFlow Fold

Пример кода на PyTorch

- Задача с SentiRuEval-2016: классифицируем твиты о мобильных операторах на **негативные**, **нейтральные** и **позитивные**
- Имеем матрицу X , содержащую выровненные признаки (индексы слов в словаре) для всех объектов обучающей выборки
- Также имеем вектор y меток классов этих объектов и вектор $lengths$ с длиной признаков (количеством слов без выравнивания) для каждого объекта
- Пусть у нас уже есть массив $numpy$ $w2v$ векторных представлений $word2vec$

Класс LSTM

```
import torch
import torch.nn as nn
from torch.autograd import Variable
from torch.nn.utils.rnn import pack_padded_sequence, pad_packed_sequence

class LSTM(nn.Module):
    def __init__(self, input_size, output_size, num_layers):
        super().__init__()
        self.lstm = nn.LSTM(input_size, output_size, num_layers)
        self.output_size, self.num_layers = output_size, num_layers

    def forward(self, inputs, lengths):
        """
        :param inputs: torch Variable of size (seq_len, batch, input_size)
        :param lengths: torch LongTensor of size (batch,)
        :return: torch Variable of size (batch, output_size)
        """
        inputs_packed = pack_padded_sequence(inputs, list(lengths))
        h0 = Variable(torch.zeros(
            self.num_layers, inputs.size(1), self.output_size)).cuda()
        c0 = Variable(torch.zeros(
            self.num_layers, inputs.size(1), self.output_size)).cuda()
        outputs_packed, _ = self.lstm(inputs_packed, (h0, c0))
        outputs = pad_packed_sequence(outputs_packed)[0]
        return outputs[lengths - 1, range(outputs.size(1))]
```


Класс SentiNet

```
class SentiNet(nn.Module):
    def __init__(self, w2v, rnn, num_classes):
        super().__init__()
        self.rnn = rnn
        self.fc = nn.Linear(rnn.output_size, num_classes)
        self.softmax = nn.Softmax()
        self.cuda()
        self.emb = nn.Embedding(w2v.size(0), w2v.size(1))
        self.emb.weight = nn.Parameter(w2v)
        self.emb.weight.requires_grad = False

    def forward(self, inputs, lengths):
        """
        :param inputs: torch Variable of size (seq_len, batch)
        :param lengths: torch LongTensor of size (batch,)
        :return: torch Variable of size (batch, num_classes)
        """
        inputs_emb = self.emb(inputs).cuda()
        return self.fc(self.rnn(inputs_emb, lengths.cuda()))

    def predict(self, inputs, lengths):
        return self.softmax(self.forward(inputs, lengths)).cpu()

    def parameters(self):
        return filter(lambda p: p.requires_grad, super().parameters())
```

Обучение нейросети

```
vocab_size, vector_size = w2v.shape
seq_len, num_objects = X.shape
num_lstm_layers = 2
num_classes = 3
batch_size = 64

rnn = LSTM(vector_size, vector_size, num_lstm_layers)
net = SentiNet(torch.from_numpy(w2v), rnn, num_classes)

criterion = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(net.parameters(), lr=.001)

for epoch in range(10):
    optimizer.zero_grad()
    for i in range(num_objects // batch_size):
        X_batch = Variable(torch.from_numpy(X[:, i * batch_size: (i + 1) * batch_size]))
        y_batch = Variable(torch.from_numpy(y[i * batch_size: (i + 1) * batch_size])).cuda()
        l_batch = torch.LongTensor(lengths[i * batch_size: (i + 1) * batch_size])
        outputs = net(X_batch, l_batch)
        loss = criterion(outputs, y_batch)
        loss.backward()
        optimizer.step()

y_pred = net.predict(Variable(torch.from_numpy(X_test)), torch.LongTensor(l_test))\
    .data.numpy().argmax(axis=1)
from sklearn.metrics import f1_score
print(f1_score(y_test, y_pred, average='macro'))
```

Приложения deep learning в обработке языка

- Определение эмоциональной окраски, сарказма и т. п.
- Синтаксический анализ текстов
- Более продвинутый информационный поиск (Яндекс “Королев”)
- Машинный перевод (Google Translate)
- Чат-боты и голосовые помощники (Replika, Samsung Bixby, Алиса)
- ?

Ключевые слова для поиска в Google

pytorch

pytorch tutorials

torchvision

tensorflow

fasttext

syntaxnet

acl 2017

deep learning book

arxiv sanity

attention mechanism

variational autoencoder

layer normalization

simple recurrent unit

selu activation