

ГЛУБОКОЕ ОБУЧЕНИЕ III: ИНИЦИАЛИЗАЦИЯ, НОРМАЛИЗАЦИЯ, ГРАДИЕНТНЫЙ СПУСК

Сергей Николенко

uData School — Киев
28 августа 2018 г.

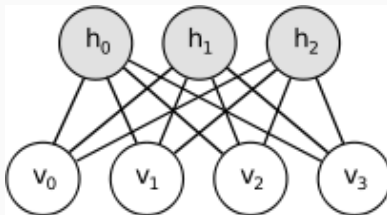
Random facts:

- 28 августа 1851 г. открыто железнодорожное движение между Санкт-Петербургом и Москвой
- 28 августа 1898 г. Кaleb Брэдхем переименовал свой «Brad's Drink» в «Pepsi-Cola»
- 28 августа 1789 г. Вильям Гершель открыл спутник Сатурна Энцелад, а 28 августа 1993 г. аппарат «Галилео» открыл Дактиль, вращающийся вокруг 243 Ида — первый известный спутник астероида

ИНИЦИАЛИЗАЦІЯ ВЕСОВ

ПРЕДОБУЧЕНИЕ БЕЗ УЧИТЕЛЯ

- Революция глубокого обучения началась с *предобучением без учителя* (unsupervised pretraining).
- Главная идея: добраться до хорошей области пространства весов, затем уже сделать fine-tuning градиентным спуском.
- Ограниченные машины Больцмана (restricted Boltzmann machines):



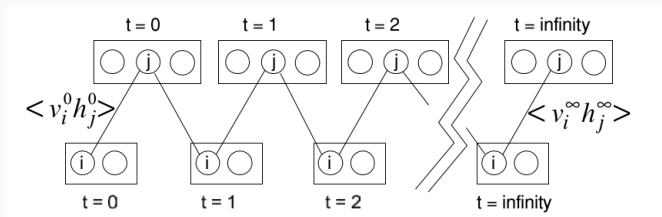
ПРЕДОБУЧЕНИЕ БЕЗ УЧИТЕЛЯ

- Это ненаправленная графическая модель, задающая распределение

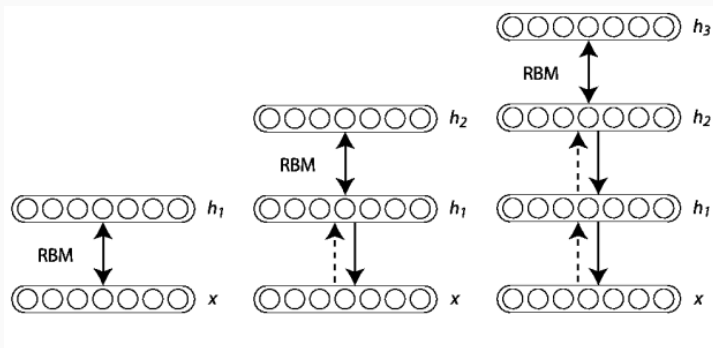
$$p(\mathbf{v}) = \sum_h p(\mathbf{v}, h) = \frac{1}{Z} \sum_h e^{-E(\mathbf{v}, h)}, \text{ где}$$

$$E(\mathbf{v}, h) = -\mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top h - h^\top W \mathbf{v}.$$

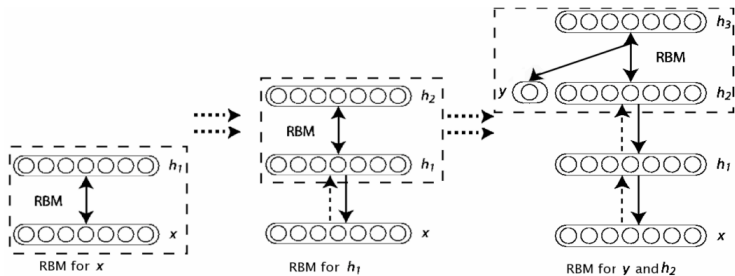
- Обучают алгоритмом Contrastive Divergence (приближение к сэмплингованию по Гиббсу).



- Из RBM можно сделать глубокие сети, поставив одну на другую:

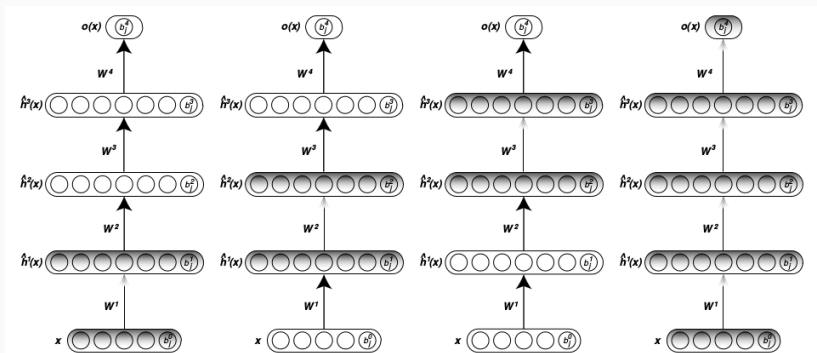


- И вывод можно вести последовательно, уровень за уровнем:

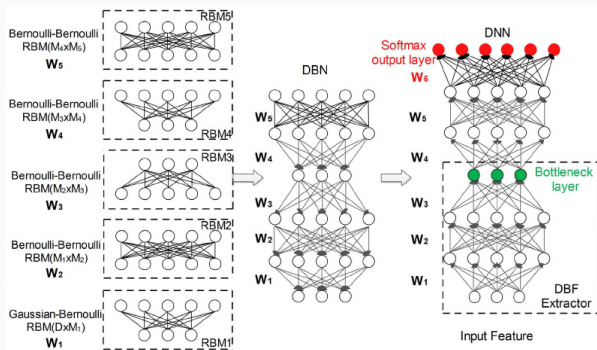


ПРЕДОБУЧЕНИЕ БЕЗ УЧИТЕЛЯ

- А потом уже дообучать градиентным спуском (fine-tuning).



- Этот подход привёл к прорыву в распознавании речи.



- Но обучать глубокие сети из RBM довольно сложно, они хрупкие, и вычислительно тоже нелегко.
- И сейчас уже не очень-то и нужны сложные модели вроде RBM для того, чтобы попасть в хорошую начальную область.
- Инициализация весов — важная часть этого.

- *Xavier initialization* (Glorot, Bengio, 2010).
- Рассмотрим простой линейный нейрон:

$$y = \mathbf{w}^\top \mathbf{x} + b = \sum_i w_i x_i + b.$$

- Его дисперсия равна

$$\begin{aligned} \text{Var}[y_i] &= \text{Var}[w_i x_i] = \mathbb{E}[w_i^2 x_i^2] - (\mathbb{E}[w_i x_i])^2 = \\ &= \mathbb{E}[x_i]^2 \text{Var}[w_i] + \mathbb{E}[w_i]^2 \text{Var}[x_i] + \text{Var}[w_i] \text{Var}[x_i]. \end{aligned}$$

ИНИЦИАЛИЗАЦИЯ ХАВЬЕРА

- Его дисперсия равна

$$\begin{aligned}\text{Var}[y_i] &= \text{Var}[w_i x_i] = \mathbb{E}[w_i^2 x_i^2] - (\mathbb{E}[w_i x_i])^2 = \\ &= \mathbb{E}[x_i]^2 \text{Var}[w_i] + \mathbb{E}[w_i]^2 \text{Var}[x_i] + \text{Var}[w_i] \text{Var}[x_i].\end{aligned}$$

- Для нулевого среднего весов

$$\text{Var}[y_i] = \text{Var}[w_i] \text{Var}[x_i].$$

- И если w_i и x_i инициализированы независимо из одного и того же распределения,

$$\text{Var}[y] = \text{Var}\left[\sum_{i=1}^{n_{\text{out}}} y_i\right] = \sum_{i=1}^{n_{\text{out}}} \text{Var}[w_i x_i] = n_{\text{out}} \text{Var}[w_i] \text{Var}[x_i].$$

- Иначе говоря, дисперсия на выходе пропорциональна дисперсии на входе с коэффициентом $n_{\text{out}} \text{Var}[w_i]$.

ИНИЦИАЛИЗАЦИЯ ХАВЬЕРА

- До (Glorot, Bengio, 2010) стандартным способом инициализации было

$$w_i \sim U \left[-\frac{1}{\sqrt{n_{\text{out}}}}, \frac{1}{\sqrt{n_{\text{out}}}} \right].$$

- См., например, *Neural Networks: Tricks of the Trade*.
- Так что с дисперсиями получается

$$\text{Var} [w_i] = \frac{1}{12} \left(\frac{1}{\sqrt{n_{\text{out}}}} + \frac{1}{\sqrt{n_{\text{out}}}} \right)^2 = \frac{1}{3n_{\text{out}}}, \text{ и}$$

$$n_{\text{out}} \text{Var} [w_i] = \frac{1}{3},$$

и после нескольких уровней сигнал совсем умирает; аналогичный эффект происходит и в backprop.

ИНИЦИАЛИЗАЦИЯ ХАВЬЕРА

- Инициализация Хавьера — давайте попробуем уменьшить изменение дисперсии, т.е. взять

$$\text{Var}[w_i] = \frac{2}{n_{\text{in}} + n_{\text{out}}};$$

для равномерного распределения это

$$w_i \sim U \left[-\frac{\sqrt{6}}{\sqrt{n_{\text{in}} + n_{\text{out}}}}, \frac{\sqrt{6}}{\sqrt{n_{\text{in}} + n_{\text{out}}}} \right].$$

- Но это работает только для симметричных активаций, т.е. не для ReLU...

- ...до работы (He et al., 2015). Вернёмся к

$$\text{Var} [w_i x_i] = \mathbb{E} [x_i]^2 \text{Var} [w_i] + \mathbb{E} [w_i]^2 \text{Var} [x_i] + \text{Var} [w_i] \text{Var} [x_i]$$

- Мы теперь можем обнулить только второе слагаемое:

$$\text{Var} [w_i x_i] = \mathbb{E} [x_i]^2 \text{Var} [w_i] + \text{Var} [w_i] \text{Var} [x_i] = \text{Var} [w_i] \mathbb{E} [x_i^2], \text{ и}$$

$$\text{Var} [y^{(l)}] = n_{\text{in}}^{(l)} \text{Var} [w^{(l)}] \mathbb{E} [(x^{(l)})^2].$$

- Мы теперь можем обнулить только второе слагаемое:

$$\text{Var} [y^{(l)}] = n_{\text{in}}^{(l)} \text{Var} [w^{(l)}] \mathbb{E} [(x^{(l)})^2].$$

- Предположим, что $x^{(l)} = \max(0, y^{(l-1)})$, и у $y^{(l-1)}$ симметричное распределение вокруг нуля. Тогда

$$\mathbb{E} [(x^{(l)})^2] = \frac{1}{2} \text{Var} [y^{(l-1)}], \quad \text{Var} [y^{(l)}] = \frac{n_{\text{in}}^{(l)}}{2} \text{Var} [w^{(l)}] \text{Var} [y^{(l-1)}].$$

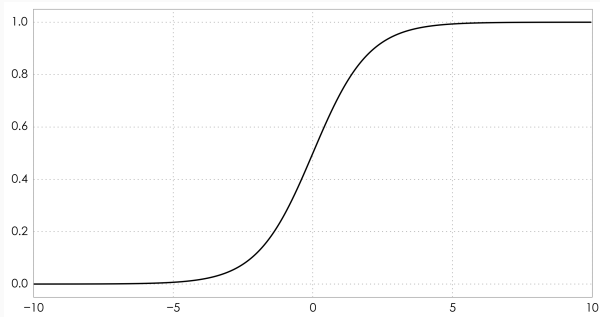
- И это приводит к формуле для дисперсии активации ReLU; теперь нет никакого n_{out} :

$$\text{Var}[w_i] = 2/n_{\text{in}}^{(l)}.$$

- Кстати, равномерную инициализацию делать не обязательно, можно и нормальное распределение:

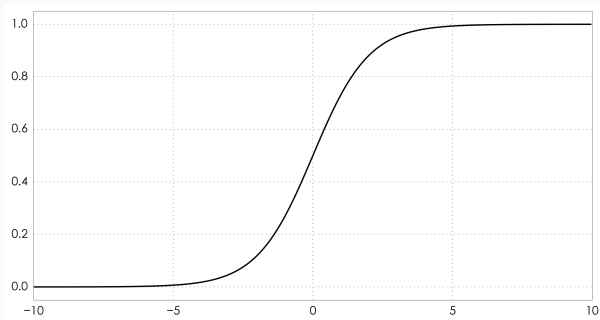
$$w_i \sim \mathcal{N}\left(0, \sqrt{2/n_{\text{in}}^{(l)}}\right).$$

- Кстати, о (Glorot, Bengio, 2010) – ещё одна важная идея.
- Эксперименты показали, что $\sigma(x) = \frac{1}{1+e^{-x}}$ работает в глубоких сетях довольно плохо.



О СИГМОИДАХ

- *Насыщение*: если $\sigma(x)$ уже «обучилась», т.е. даёт большие по модулю значения, то её производная близка к нулю и «поменять мнение» трудно.

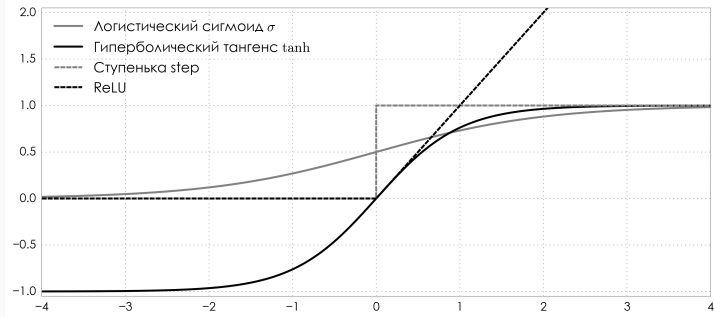


- Но ведь другие тоже насыщаются? В чём разница?

- Рассмотрим последний слой сети $h(W\mathbf{a} + \mathbf{b})$, где \mathbf{a} — выходы предыдущего слоя, \mathbf{b} — свободные члены, h — функция активации последнего уровня, обычно `softmax`.
- Когда мы начинаем оптимизировать сложную функцию потерь, поначалу выходы h не несут полезной информации о входах, ведь первые уровни ещё не обучены.
- Тогда неплохим приближением будет константная функция, выдающая средние значения выходов.
- Это значит, что $h(W\mathbf{a} + \mathbf{b})$ подберёт подходящие свободные члены \mathbf{b} и постарается обнулить слагаемое Wh , которое поначалу скорее шум, чем сигнал.

О СИГМОИДАХ

- Иначе говоря, в процессе обучения мы постараемся привести выходы предыдущего слоя к нулю.
- Здесь и проявляется разница: у $\sigma(x) = \frac{1}{1+e^{-x}}$ область значений $(0, 1)$ при среднем $\frac{1}{2}$, и при $\sigma(x) \rightarrow 0$ будет и $\sigma'(x) \rightarrow 0$.
- А у \tanh наоборот: когда $\tanh(x) \rightarrow 0$, $\tanh'(x)$ максимальна.

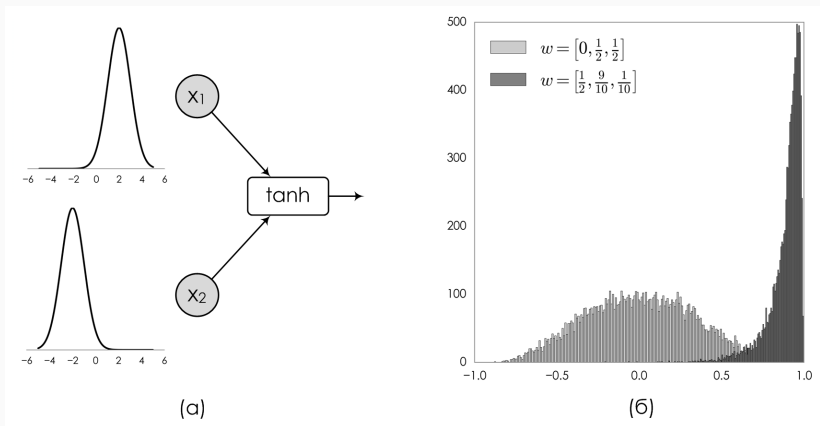


НОРМАЛИЗАЦИЯ ПО МИНИ-БАТЧАМ

- Ещё одна важная проблема в глубоких сетях: *внутренний сдвиг переменных* (internal covariate shift).
- Когда меняются веса слоя, меняется распределение его выходов.
- Это значит, что следующему уровню придётся всё начинать заново, он же не ожидал таких входов, не видел их раньше!
- Более того, нейроны следующего уровня могли уже и насытиться, и им теперь сложно быстро обучиться заново.
- Это серьёзно мешает обучению.

НОРМАЛИЗАЦИЯ ПО МИНИ-БАТЧАМ

- Вот характерный пример:



- Что делать?

НОРМАЛИЗАЦИЯ ПО МИНИ-БАТЧАМ

- Можно пытаться нормализовать входы каждого уровня.
- Не работает: рассмотрим для простоты уровень с одним только bias b и входами u :

$$\hat{\mathbf{x}} = \mathbf{x} - \mathbb{E}[\mathbf{x}], \text{ где } \mathbf{x} = u + b.$$

- На следующем шаге градиентного спуска получится $b := b + \Delta b \dots$
- ...но $\hat{\mathbf{x}}$ не изменится:

$$u + b + \Delta b - \mathbb{E}[u + b + \Delta b] = u + b - \mathbb{E}[u + b].$$

- Так что всё обучение сведётся к тому, что b будет неограниченно расти — не очень хорошо.

- Можно пытаться добавить нормализацию как отдельный слой:

$$\hat{\mathbf{x}} = \text{Norm}(\mathbf{x}, \mathcal{X}).$$

- Это лучше, но теперь этому слою на входе нужен весь датасет \mathcal{X} !
- И на шаге градиентного спуска придётся вычислить $\frac{\partial \text{Norm}}{\partial \mathbf{x}}$ и $\frac{\partial \text{Norm}}{\partial \mathcal{X}}$, да ещё и матрицу ковариаций

$$\text{Cov}[\mathbf{x}] = \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\mathbf{x}\mathbf{x}^\top] - \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}]^\top.$$

- Это точно не сработает.

НОРМАЛИЗАЦИЯ ПО МИНИ-БАТЧАМ

- Решение в том, чтобы нормализовать каждый вход отдельно, и не по всему датасету, а по текущему кусочку; это и есть *нормализация по мини-батчам* (batch normalization).
- После нормализации по мини-батчам получим

$$\hat{x}_k = \frac{x_k - \mathbb{E}[x_k]}{\sqrt{\text{Var}[x_k]}}$$

где статистики подсчитаны по текущему мини-батчу.

- Ещё одна проблема: теперь пропадают нелинейности!
- Например, σ теперь практически всегда близка к линейной.

НОРМАЛИЗАЦИЯ ПО МИНИ-БАТЧАМ

- Чтобы это исправить, нужно добавить гибкости уровню batchnorm.
- В частности, нужно разрешить ему обучаться иногда *ничего не делать* со входами.
- Так что вводим дополнительные параметры (сдвиг и растяжение):

$$y_k = \gamma_k \hat{x}_k + \beta_k = \gamma_k \frac{x_k - \mathbb{E}[x_k]}{\sqrt{\text{Var}[x_k]}} + \beta_k.$$

- γ_k и β_k — это новые переменные, тоже будут обучаться градиентным спуском, как веса.

НОРМАЛИЗАЦИЯ ПО МИНИ-БАТЧАМ

- И ещё добавим ϵ в знаменатель, чтобы на ноль не делить.
- Теперь мы можем формально описать слой батч-нормализации — для очередного мини-батча $B = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$:

- вычислить базовые статистики по мини-батчу

$$\mu_B = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i, \quad \sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu_B)^2,$$

- нормализовать входы

$$\hat{\mathbf{x}}_i = \frac{\mathbf{x}_i - \mu_b}{\sqrt{\sigma_B^2 + \epsilon}},$$

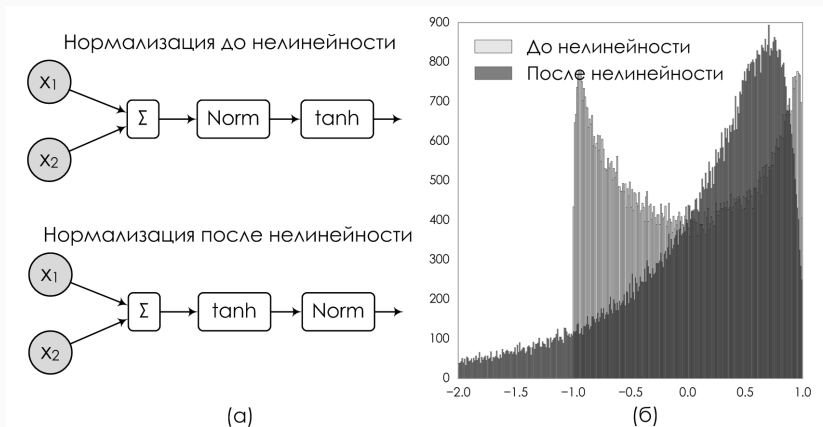
- вычислить результат

$$\mathbf{y}_i = \gamma \mathbf{x}_i + \beta.$$

- Через всё это совершенно стандартным образом пропускаются градиенты, в том числе по γ и β .

НОРМАЛИЗАЦИЯ ПО МИНИ-БАТЧАМ

- Последнее замечание: важно, куда поместить слой batchnorm.
- Можно до, а можно после нелинейности.



- Нормализация по мини-батчам сейчас стала фактически стандартом.
- Очередная очень крутая история, примерно как дропаут.
- Но мысль идёт и дальше:
 - (Laurent et al., 2016): BN не помогает рекуррентным сетям;
 - (Cooijmans et al., 2016): recurrent batch normalization;
 - (Salimans and Kingma, 2016): нормализация весов для улучшения обучения — давайте добавим веса как

$$h_i = f \left(\frac{\gamma}{\|\mathbf{w}_i\|} \mathbf{w}_i^\top \mathbf{x} + b_i \right);$$

тогда мы будем перемасштабировать градиент, стабилизировать его норму и приближать его матрицу ковариаций к единичной, что улучшает обучение.

ВАРИАНТЫ ГРАДИЕНТНОГО СПУСКА

- «Ванильный» стохастический градиентный спуск:

$$\theta_t = \theta_{t-1} - \eta \nabla E(\mathbf{x}_t, \theta_{t-1}, y_t).$$

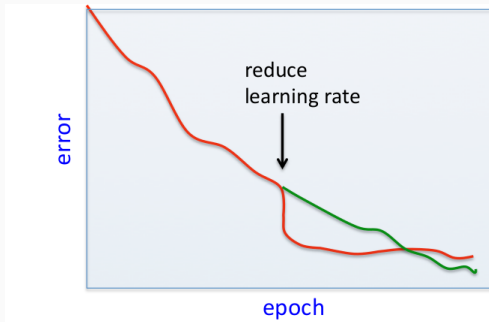
- Всё зависит от скорости обучения η .
- Первая мысль — пусть η уменьшается со временем:
 - линейно (linear decay):

$$\eta = \eta_0 \left(1 - \frac{t}{T}\right);$$

- или экспоненциально (exponential decay):

$$\eta = \eta_0 e^{-\frac{t}{T}}.$$

- Скорость обучения лучше не уменьшать слишком быстро.



- Но это в любом случае никак не учитывает собственно E ; лучше быть *адаптивным*.

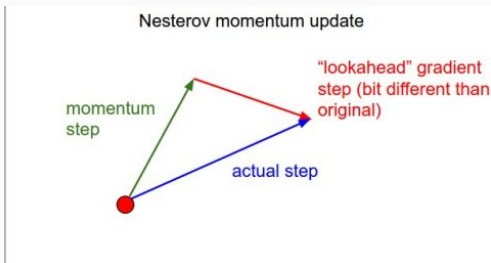
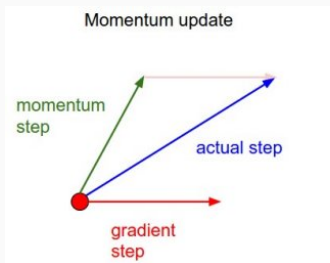
- *Метод моментов* (momentum): сохраним часть скорости, как у материальной точки.
- С инерцией получается

$$u_t = \gamma u_{t-1} + \eta \nabla_{\theta} E(\theta),$$
$$\theta = \theta - u_t.$$

- И теперь мы сохраняем γu_{t-1} .

МЕТОД МОМЕНТОВ

- Но на самом деле мы уже знаем, что попадём в γu_{t-1} на промежуточном шаге.
- Давайте прямо там, на полпути, и вычислим градиент!



- *Метод Нестерова* (Nesterov's momentum):

$$u_t = \gamma u_{t-1} + \eta \nabla_{\theta} E(\theta - \gamma u_{t-1})$$



- Можно ли ещё лучше?..

- ...ну, можно попробовать методы второго порядка.
- Метод Ньютона:

$$E(\theta) \approx E(\theta_0) + \nabla_{\theta} E(\theta_0)(\theta - \theta_0) + \frac{1}{2}(\theta - \theta_0)^{\top} H(E(\theta))(\theta - \theta_0).$$

- Когда работает, это обычно гораздо быстрее, и нет никакой η , ничего настраивать не надо.
- Но нужно считать гессиан $H(E(\theta))$, и это нереально.

- Есть, правда, приближения.
- L-BFGS (limited memory Broyden–Fletcher–Goldfarb–Shanno):
 - строим аппроксимацию к H^{-1} ;
 - для этого сохраняем последовательно апдейты аргументов функции и градиентов и выражаем через них H^{-1} .
- Интересный открытый вопрос: можно ли заставить L-BFGS работать для deep learning?
- Но пока не получается («в лоб» было бы нужно считать градиент по всему датасету, а по мини-батчам непонятно как).

- Но дальше улучшить всё равно можно.
- Заметим, что до сих пор скорость обучения была одна во всех направлениях.
- Идея: давайте быстрее двигаться по тем параметрам, которые не сильно меняются, и медленнее по быстро меняющимся параметрам.

- *Adagrad*: давайте накапливать историю этой скорости изменений и учитывать её.
- Обозначая $g_{t,i} = \nabla_{\theta_i} L(\theta)$, получим

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot g_{t,i},$$

где G_t – диагональная матрица с $G_{t,ii} = G_{t-1,ii} + g_{t,i}^2$, которая накапливает общее значение градиента по всей истории обучения.

- Так что скорость обучения всё время уменьшается, но с разной скоростью для разных θ_i .

- Проблема: G всё увеличивается и увеличивается, и скорость обучения иногда уменьшается слишком быстро.
- *Adadelta* (Zeiler, 2012) – та же идея, но две новых модификации.
- Во-первых, историю градиентов мы теперь считаем с затуханием:

$$G_{t,ii} = \rho G_{t-1,ii} + (1 - \rho) g_{t,i}^2.$$

- А всё остальное здесь точно так же:

$$u_t = -\frac{\eta}{\sqrt{G_{t-1} + \epsilon}} \mathbf{g}_{t-1}.$$

- Во-вторых, надо бы «единицы измерения» привести в соответствие.
- В предыдущих методах была проблема:
 - в обычном градиентном спуске или методе моментов «единицы измерения» обновления параметров $\Delta\theta$ — это единицы измерения градиента, т.е. если веса в секундах, а целевая функция в метрах, то градиент будет иметь размерность «метр в секунду», и мы вычитаем метры в секунду из секунд;
 - а в Adagrad получалось, что значения обновлений $\Delta\theta$ зависели от отношений градиентов, и величина обновлений вовсе безразмерная.

АДАПТИВНЫЕ МЕТОДЫ ГРАДИЕНТНОГО СПУСКА

- Эта проблема решается в методе второго порядка: обновление параметров $\Delta\theta$ пропорционально $H^{-1}\nabla_{\theta}f$, то есть размерность будет

$$\Delta\theta \propto H^{-1}\nabla_{\theta}f \propto \frac{\frac{\partial f}{\partial\theta}}{\frac{\partial^2 f}{\partial\theta^2}} \propto \text{размерность } \theta.$$

- Чтобы привести *Adadelta* в соответствие, нужно домножить на ещё одно экспоненциальное среднее, но теперь уже от квадратов обновлений параметров, а не от градиента.
- Настоящее среднее мы не знаем, аппроксимируем предыдущими шагами:

$$\mathbb{E}[\Delta\theta^2]_t = \rho\mathbb{E}[\Delta\theta^2]_{t-1} + (1-\rho)\Delta\theta^2, \text{ где}$$
$$u_t = -\frac{\sqrt{\mathbb{E}[\Delta\theta^2]_{t-1} + \epsilon}}{\sqrt{G_{t-1} + \epsilon}} \cdot g_{t-1}.$$

- Следующий вариант – *RMSprop* из курса Хинтона.
- Практически то же, что *Adadelta*, только *RMSprop* не делает вторую поправку с изменением единиц и хранением истории самих обновлений, а просто использует корень из среднего от квадратов (вот он где, RMS) от градиентов:

$$u_t = -\frac{\eta}{\sqrt{G_{t-1} + \epsilon}} \cdot g_{t-1}.$$

- И последний алгоритм – *Adam* (Kingma, Ba, 2014).
- Модификация *Adagrad* со сглаженными версиями среднего и среднеквадратичного градиентов:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t,$$

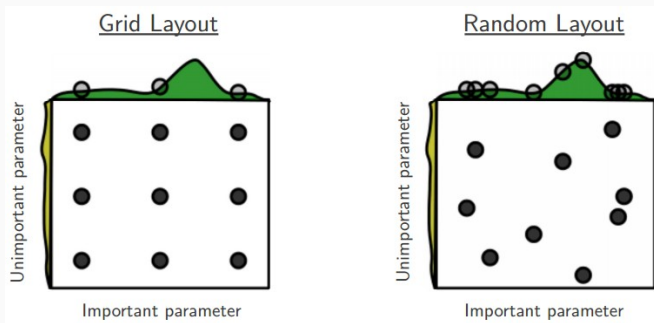
$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2,$$

$$u_t = \frac{\eta}{\sqrt{v_t + \epsilon}} m_t.$$

- (Kingma, Ba, 2014) рекомендуют $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$.
- *Adam* практически не требует настройки, используется на практике очень часто.
- ...[анимированные примеры]...

ПРАКТИЧЕСКИЕ ЗАМЕЧАНИЯ

- Ещё практические замечания об оптимизации гиперпараметров:
 - одного валидационного множества достаточно;
 - лучше гиперпараметры искать на логарифмической шкале;
 - и лучше случайным поиском, а не по сетке (Bergstra and Bengio).



Спасибо за внимание!