

Инициализация весов и регуляризация

ФКН

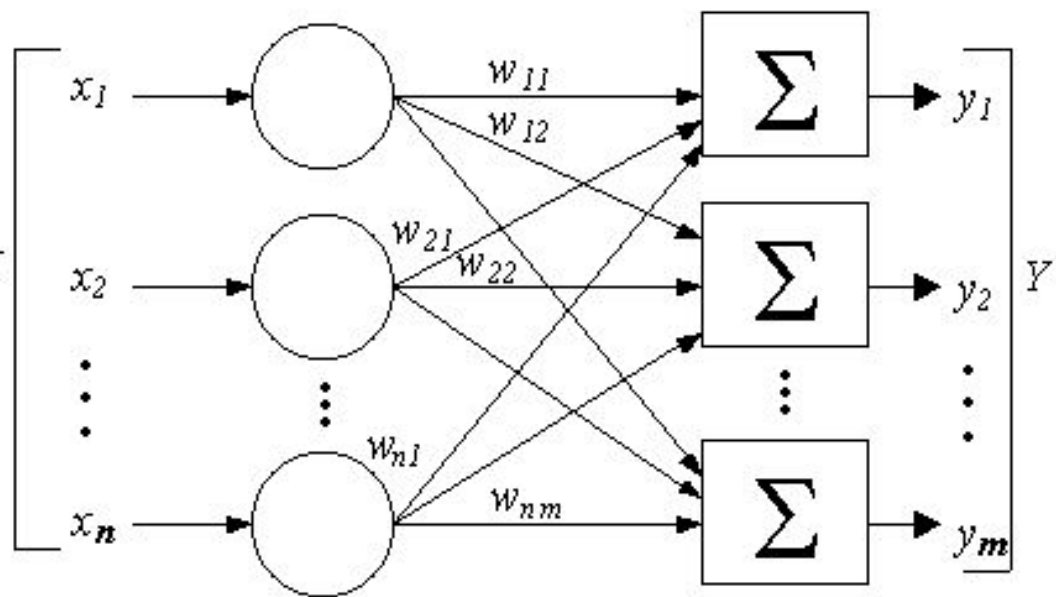
План лекции

- Повторение линейного слоя и нелинейностей
- ML common knowledge
- Инициализация весов сети
 - Инициализация константой
 - Нормальное распределение
 - Xavier
 - Kaiming
- Методы регуляризации сети
 - Модификации функций ошибки
 - Слои нейросетей
 - Оптимизация
 - Работа с данными

Повторение*

Линейный слой

- Матрица весов: $W = \begin{bmatrix} w_{0;1} & \dots & w_{0;m} \\ \vdots & \ddots & \vdots \\ w_{d;1} & \dots & w_{d;m} \end{bmatrix} \in \mathbb{R}^{(d+1) \times m}$.
- $\sigma(z_1, \dots, z_m) \sim (\sigma(z_1), \dots, \sigma(z_m))$.
- Тогда $y = \sigma(x^T * W)$.



Два линейных слоя подряд

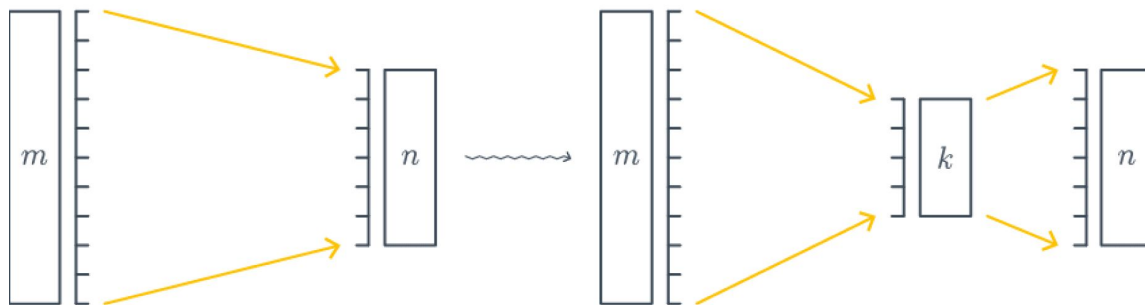
Линейная комбинация линейных отображений есть линейное отображение, то есть два последовательных линейных слоя эквивалентны одному линейному слою.

Но на самом деле бывают ситуации, когда два линейных слоя подряд — это полезно.



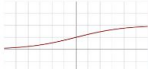
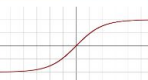




Два линейных слоя подряд

Например, если вы понимаете, что у вас очень много параметров, а информации в данных не так много, вы можете заменить линейный слой, превращающий m -мерные векторы в n -мерные, на два, вставив посередине k -мерное представление, где $k \ll m, n$:

С точки зрения линейной алгебры это примерно то же самое, что потребовать, чтобы матрица исходного линейного слоя имела ранг не выше k . И с точки зрения сужения «информационного канала» это иногда может сработать. Два линейных слоя подряд стоит ставить, только если вы хорошо понимаете, чего хотите добиться.



Нелийности

ACTIVATION FUNCTION	PLOT	EQUATION	DERIVATIVE	RANGE
Linear		$f(x) = x$	$f'(x) = 1$	$(-\infty, \infty)$
Binary Step		$f(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{if } x \neq 0 \\ \text{undefined} & \text{if } x = 0 \end{cases}$	$\{0, 1\}$
Sigmoid		$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$	$(0, 1)$
Hyperbolic Tangent(tanh)		$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$	$f'(x) = 1 - f(x)^2$	$(-1, 1)$
Rectified Linear Unit(ReLU)		$f(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ \text{undefined} & \text{if } x = 0 \end{cases}$	$[0, \infty)$
Softplus		$f(x) = \ln(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$	$(0, 1)$
Leaky ReLU		$f(x) = \begin{cases} 0.01x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0.01 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$	$(-1, 1)$
Exponential Linear Unit(ELU)		$f(x) = \begin{cases} \alpha(e^x - 1) & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$	$f'(x) = \begin{cases} \alpha e^x & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ 1 & \text{if } x = 0 \text{ and } \alpha = 1 \end{cases}$	$[0, \infty)$

SoftMax

Output
layer

$$\begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ 0.7 \\ 1.1 \end{bmatrix}$$

Softmax
activation function

$$\frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

Probabilities

$$\begin{bmatrix} 0.02 \\ 0.90 \\ 0.05 \\ 0.01 \\ 0.02 \end{bmatrix}$$

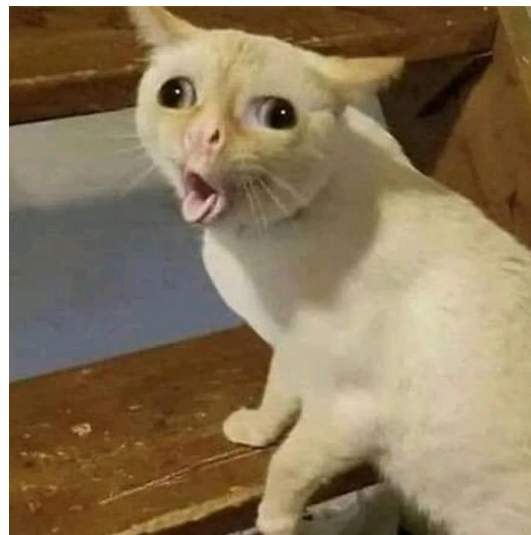
ML common knowledge

Робастность

Робастность — это устойчивость метода или модели к шуму, ошибкам и неожиданным ситуациям

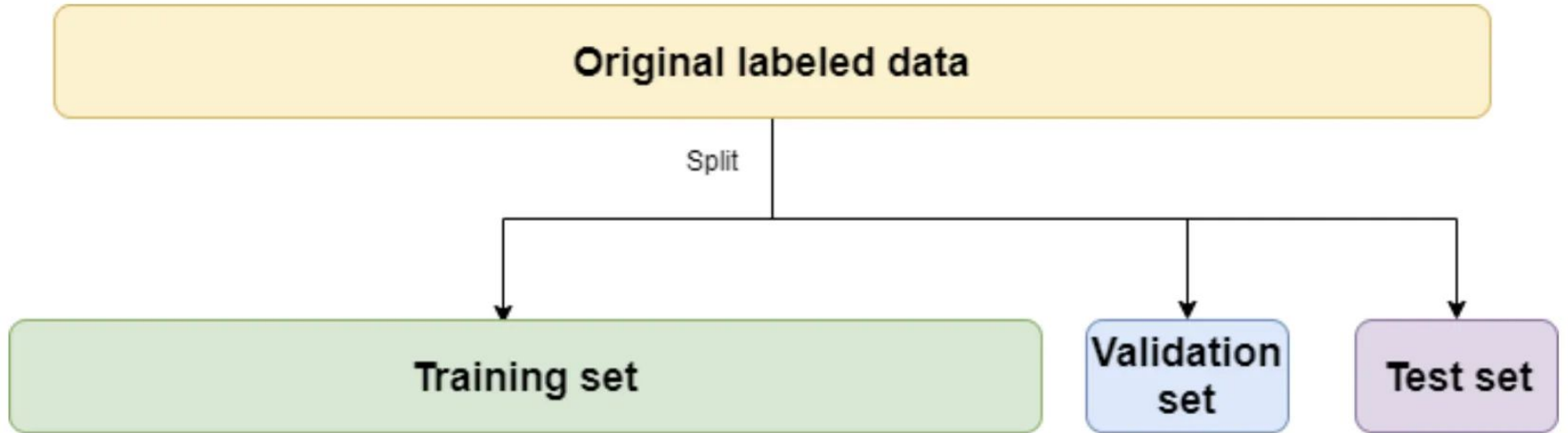


neural net → Котик

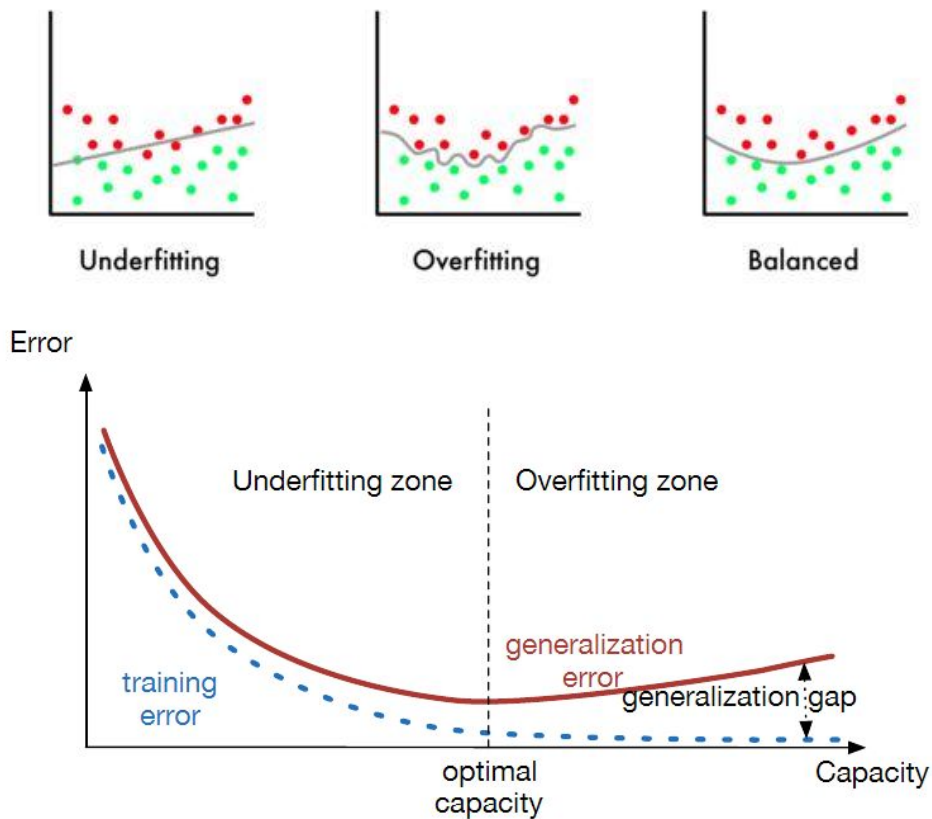


neural net → ???

Train, Validation, Test

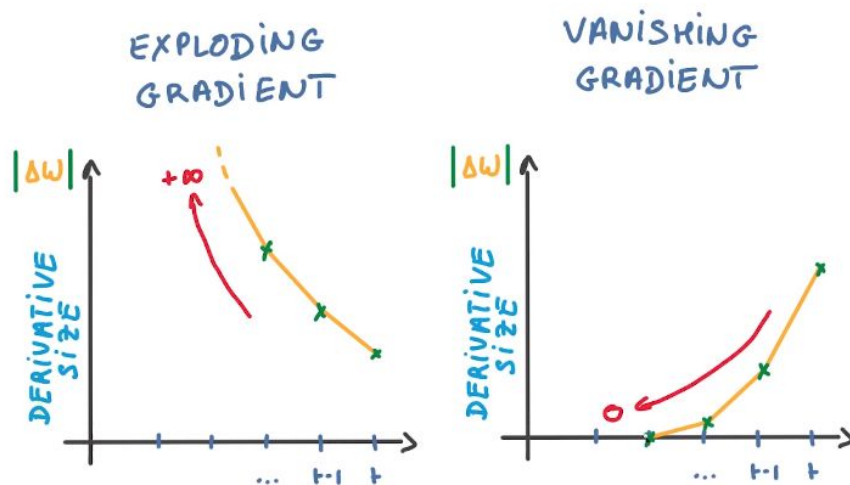


Переобучение



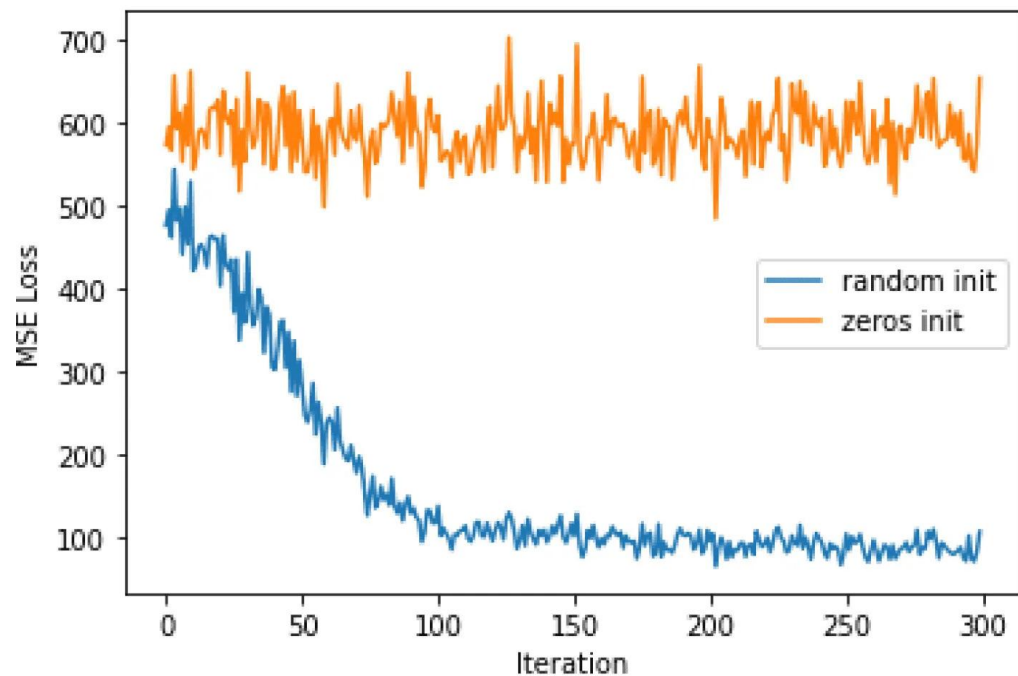
Взрыв и затухание градиента

При обратном распространении ошибки при прохождении через слои нейронной сети в элементах градиента могут накапливаться большие значения, что будет приводить к сильным изменениям весов, что сильно “разболтает” веса – взрыв градиента. Обратная ситуация – затухание.



Инициализация весов

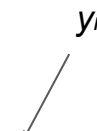
Инициализация константой



Веса будут меняться одинаково внутри слоёв

Xavier

Рассмотрим одну активацию:

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



Тогда её дисперсия:

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

Используем симметричную функцию активации (значит мат. ож. $x = 0$) и берём распределение w с мат. ожиданием в 0

Xavier

Рассмотрим одну активацию:

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


Тогда её дисперсия:

$$\begin{aligned} \text{Var}(y_i) &= \text{Var}(w_i x_i) = \mathbb{E}[x_i^2 w_i^2] - (\mathbb{E}[x_i w_i])^2 = \leftarrow x \text{ и } w \text{ независимы} \\ &= \cancel{\mathbb{E}[x_i]^2 \text{Var}(w_i)} + \cancel{\mathbb{E}[w_i]^2 \text{Var}(x_i)} + \text{Var}(w_i) \text{Var}(x_i) \end{aligned}$$

Используем симметричную функцию активации (значит мат. ож. $x = 0$) и берём распределение w с мат. ожиданием в 0

Xavier

Итого получаем: $\text{Var}(y_i) = \text{Var}(w_i) \text{Var}(x_i)$

Если предположить, что x и w сэмплируются независимо из одного распределения (сильное утверждение), то получаем:

$$\text{Var}(y) = \text{Var}\left(\sum_{i=1}^{n_{\text{out}}} y_i\right) = \sum_{i=1}^{n_{\text{out}}} \text{Var}(w_i x_i) = n_{\text{out}} \text{Var}(w_i) \text{Var}(x_i)$$

← число нейронов на выходе слоя

Т.е. дисперсия выходов пропорциональна дисперсии входов с коэффициентом $n_{\text{out}} \text{Var}(w_i)$

Xavier

Ранее часто веса инициализировали из распределения:

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


Из-за чего дисперсия активации падала втрое с каждым слоем, т.е. сигнал “затухал”:

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

Xavier (backward pass)

На обратном проходе похожая ситуация. Если $z^{(l+1)} = f(y^{(l)})$, где l – номер слоя, а f – функция активации, то обратное распространение ошибки выглядит следующим образом:

$$\frac{\partial L}{\partial y_i^{(l)}} = f'(y_i^{(l)}) \sum_j w_{i,j}^{(l+1)} \frac{\partial L}{\partial y_j^{(l+1)}}$$

 ~ 1 для симметричной f с единичной производной в нуле (напр., \tanh)

Т.е. дисперсия градиентов пропорциональна дисперсии выходов с коэффициентом $n_{in} \text{Var}(w_i)$

Xavier (idea)

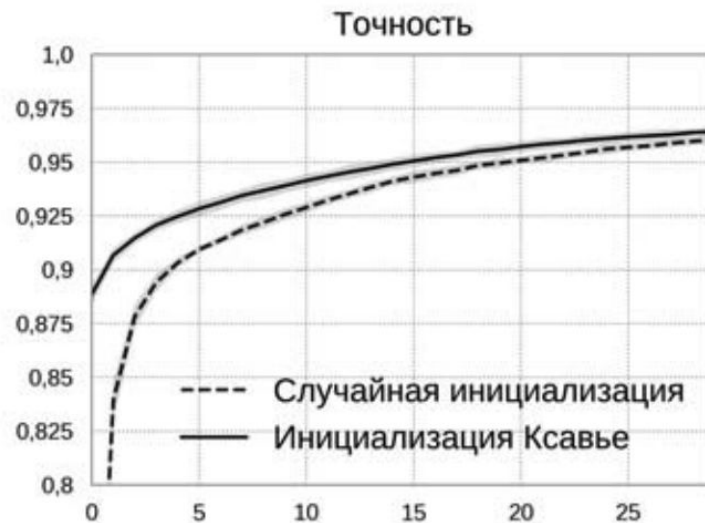
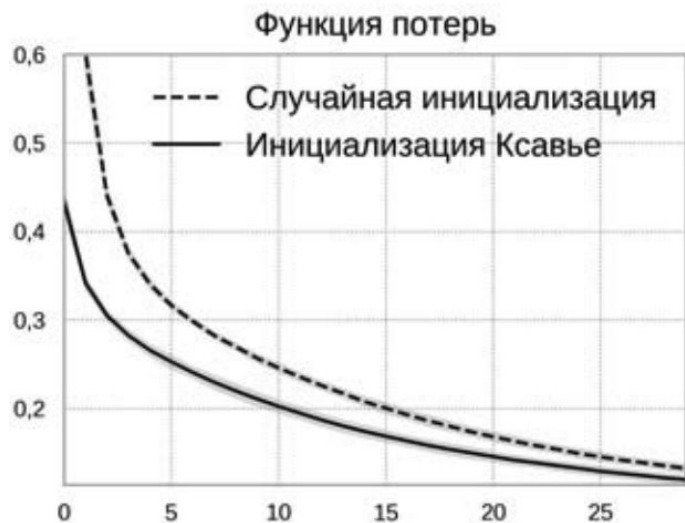
Чтобы не затухал ни сигнал, ни градиент, надо сэмплировать веса из распределения с дисперсией:

$$\text{Var}(w_i) = \frac{2}{n_{\text{in}} + n_{\text{out}}}$$

Тогда для равномерного распределения в центре с нулём, границы будут такие:

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

Инициализация просто нормальным распределением



Kaiming (for ReLU)

Но функция активации не всегда симметрично относительно нуля. Есть популярный ReLU, для отбросить первое слагаемое не получится:

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

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

Дальше выкладки похожи на Xavier, но явно используется формула ReLU. Для него выводится следующее оптимальное распределение:

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

Регуляризация нейросетей

Регуляризация нейросетей

Техники для борьбы с переобучением, повышения робастности и для получения более подходящего решения с точки зрения эксперта. В нейронных сетях можно добиться изменением

- функции потерь
- структуры сети
- процесса оптимизации
- данных

Функции потерь

Модифицируем loss, по которому пускаем расчёт градиента

$$Loss_{result} = L_{original} + L_{regularization}$$

Добавка может быть любой, подходящей под задачу, но пара наиболее распространённых примеров: L1, L2 -регуляризации (a.k.a.weight decay)

Регуляризатор может быть применён как к весам модели, так и к активациям/выходам, в зависимости от желаемых целей

L2-регуляризация

- L2-норма на параметры

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \ell(f(x_i, \theta), y_i) + \frac{\lambda}{2} \|\theta\|_2^2$$

- Препятствует росту параметров
- Обычно встроена в оптимизатор (в SGD – это weight decay)
$$\theta_{t+1} \leftarrow \theta_t(1 - \lambda) - \gamma \nabla_{\theta} \ell_t(\theta_t)$$
- Обычные значения 10^{-3} , 10^{-4} (по умолчанию – 0)
- L1 – примерно такой же эффект

Модификации архитектуры

- Специальные слои (нормализации и DropOut)
- Дистилляция
- Квантизация
- Пруннинг

Нормализации

- BatchNorm
- LayerNorm
- InstanceNorm

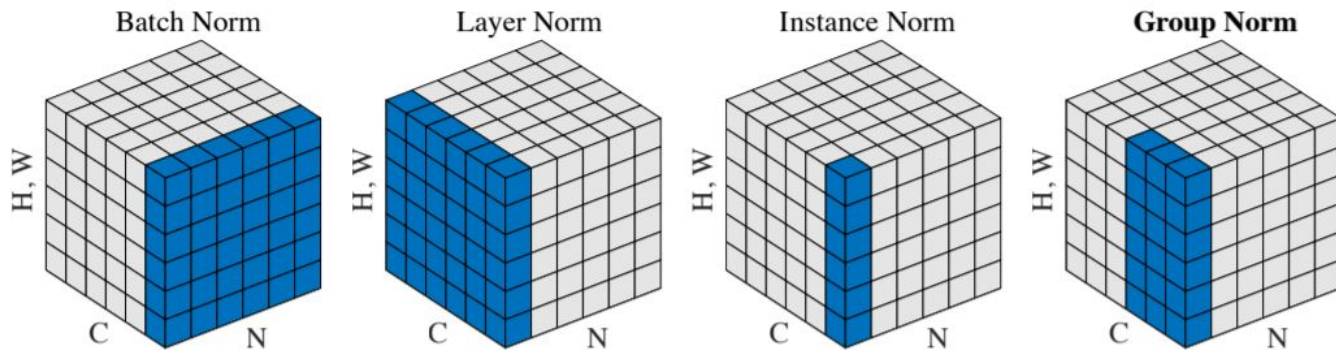
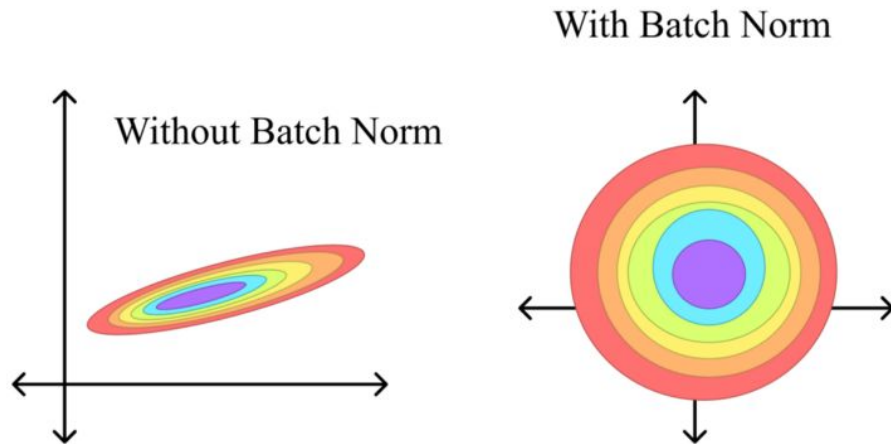


Figure 2. **Normalization methods.** Each subplot shows a feature map tensor, with N as the batch axis, C as the channel axis, and (H, W) as the spatial axes. The pixels in blue are normalized by the same mean and variance, computed by aggregating the values of these pixels.

BatchNorm

Когда веса нижележащих слоёв обновляются, распределение активаций на входе последующих слоёв "дрейфует", из-за чего обучение становится менее стабильным и требует более аккуратного подбора learning rate и инициализации. Это называется Internal Covariate Shift, с ним борется Batch Normalization и его аналоги.



BatchNorm

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots x_m\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

[Ioffe&Szegedy, 2015]

На тесте работает по другому!

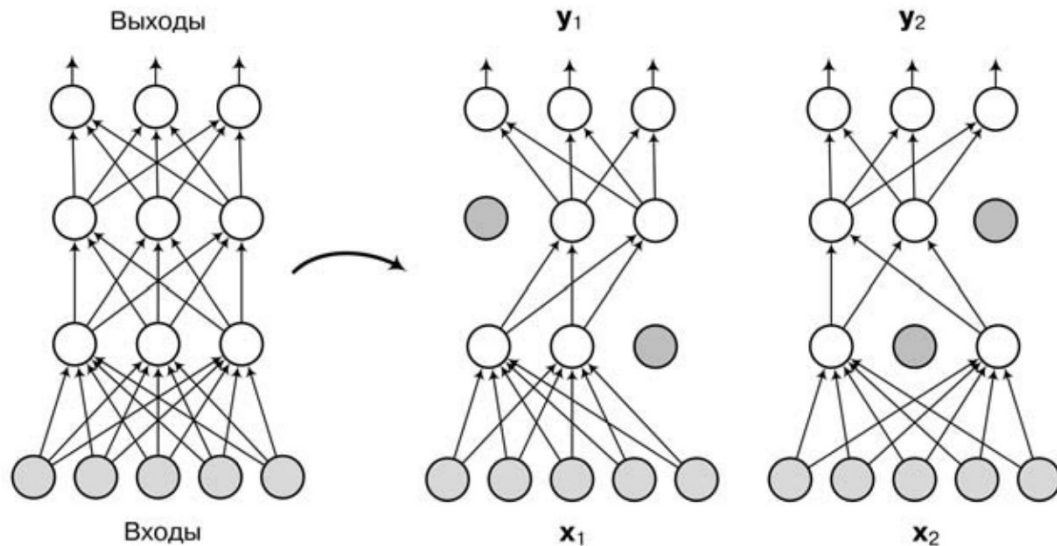
Используются предсчитанные оценки средних и дисперсий.

Оценки средних и дисперсий вычисляются скользящим средним во время обучения.

μ , σ – статистики, накапливаемые по батчу
 γ , β – параметры, обучаемые градиентным спуском

DropOut

Оставляем каждую активацию с вероятностью $1 - p$, а с вероятностью p приравниваем её 0



На последнем слое DropOut обычно не делают, т.к. нужен вектор определённой размерности

DropOut. А как применять?

Сэмплировать “прореженные” DropOut’ом сети, применять их все, а затем усреднять

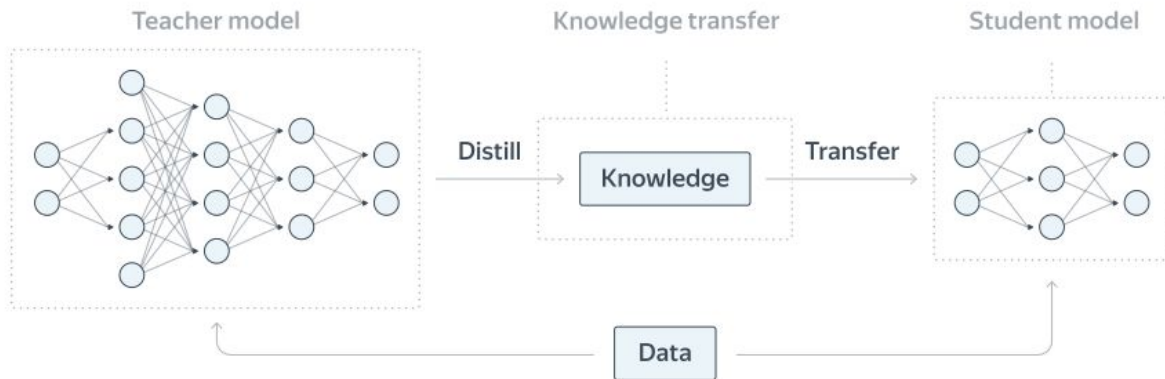
Оказывается это примерно эквивалентно применению сети без “прореживаний”, в которой каждая активация умножается на $1 - p$, благодаря чему сохраняем их математическое ожидание

Это объясняет и мотивацию: DropOut позволяет нам в некотором смысле усреднить все возможные под-архитектуры нашей сети, сделать активацию каждого нейрона более информативной.

DropOut. Техническая реализация

	DropOut	Inverted DropOut
Обучение	$h_{\text{train}} = m \odot h, \quad m \sim \text{Bernoulli}(1 - p)$	$h_{\text{train}} = \frac{m}{1 - p} \odot h, \quad m \sim \text{Bernoulli}(1 - p)$
Инференс	$h_{\text{test}} = (1 - p) \cdot h$	$h_{\text{test}} = h$

Дистилляция



Учим маленькую модель повторять предсказания большой модели. Её возможности более ограничены, поэтому ей не удаётся переобучаться, как большой модели

Подробнее об этом методе через одну лекцию 😊

Квантизация

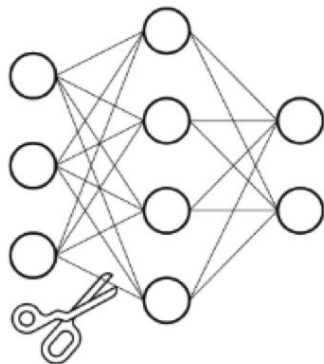


Веса модели переводим с высокой точности (fp32) к более низкой (fp16/fp8/fp4/int8). Модель снова не может так сильно переобучаться под данные и вынуждена описывать закономерности проще

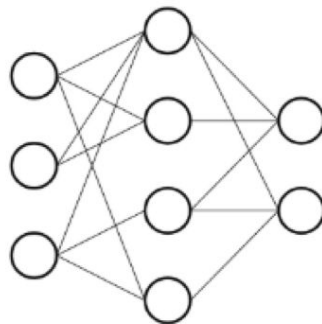
Подробнее об этом методе через одну лекцию 😊

Пруннинг

До:



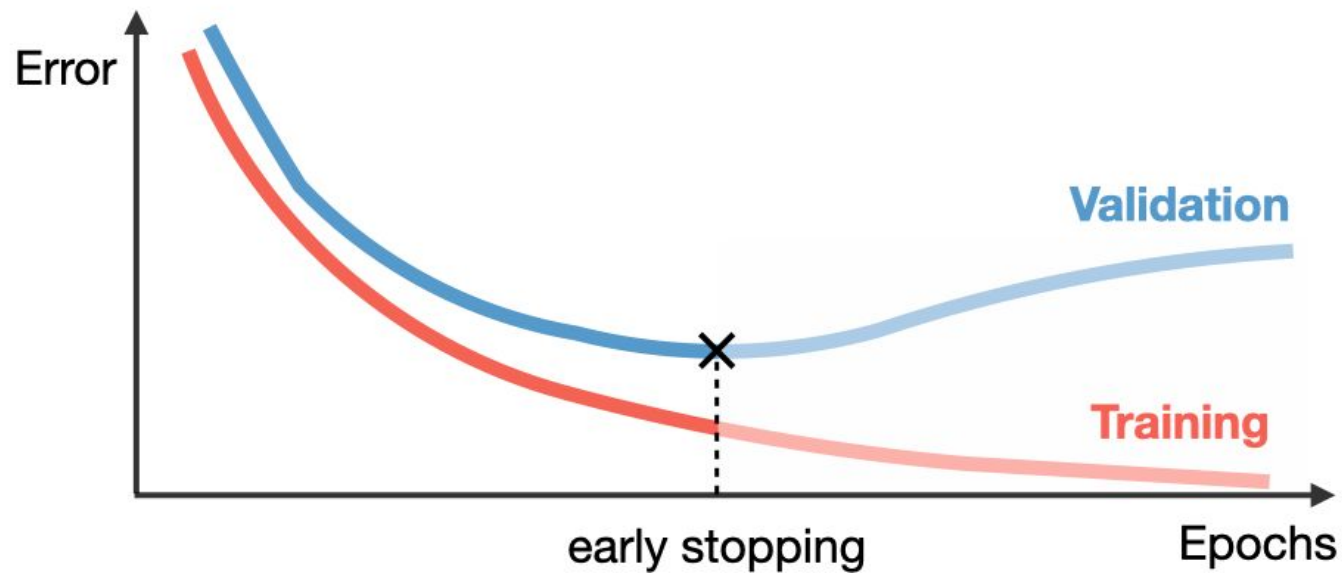
После:



Выбираем критерий важности для весов, наименее важные веса убираем и дообучаем модель. Архитектура модели становится всё более специфичной, модель держится за “важные” признаки и игнорирует “шумные”

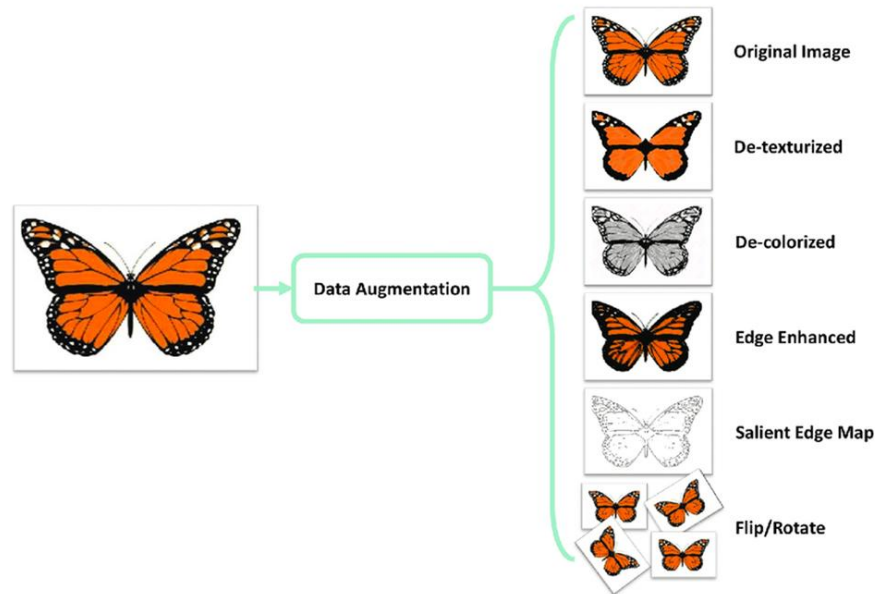
Подробнее об этом методе через одну лекцию 😊

Оптимизация. Early stopping



Изменение данных

- Аугментации – устойчивость к некоторым преобразованиям
 - повороты и отражения картинок
 - замена случайных слов в тексте на синонимы



- Data mixup – создание промежуточных примеров через линейную комбинацию входов и меток
- Label smoothing – заменяем строгие метки классификации 0/1 на “вероятности”, позволяя модели “сомневаться в разметке”

Источники

- <https://education.yandex.ru/handbook/ml>
- https://github.com/aosokin/dl_cshse_ami/tree/master/2020-fall/lectures
- Глубокое обучение. Погружение в мир нейронных сетей – Николенко С., Кадурын А.
- Глубокое обучение – Гудфеллоу Я., Бенджио И., Курвилль А.