

## Batch Normalization

## Weights Initialization

Инициализация генератора инициализации ZMC

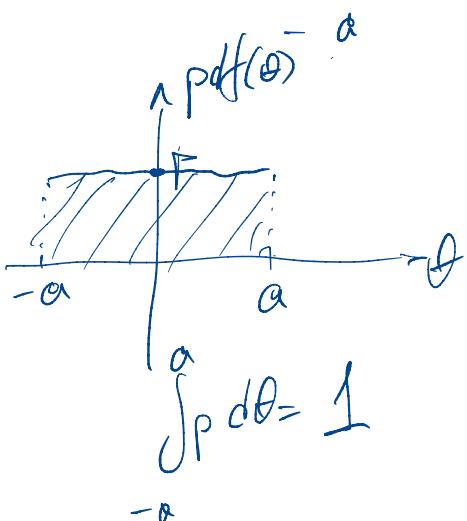
$$\text{Var } \theta = \frac{2}{m+k}$$

$m$  - кон-60 выходов единиц  
 $k$  - кон-60 выходов единиц

$\theta \sim U[-\alpha; \alpha]$   $\alpha = ?$

$$\theta \sim \mathcal{N}(0, \sigma^2) \quad \sigma^2 = ? = \frac{2}{m+k}$$

$$\sigma^2(\theta) = \int_{-\alpha}^{\alpha} (\theta - \bar{\theta})^2 \cdot \text{pdf}(\theta) d\theta = \int_{-\alpha}^{\alpha} \theta^2 = \sigma^2(\theta) = \frac{2}{m+k}$$

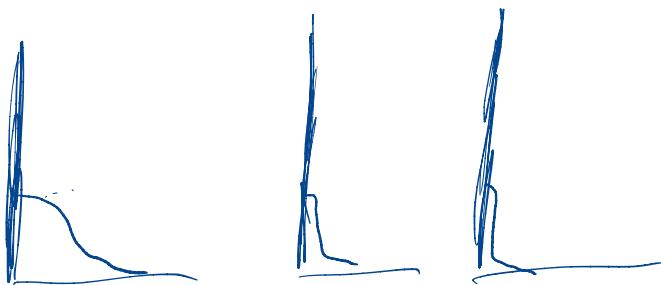
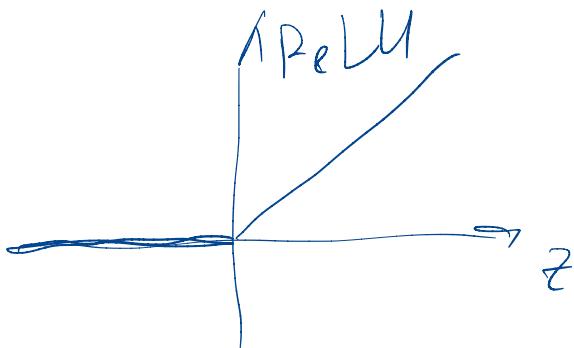
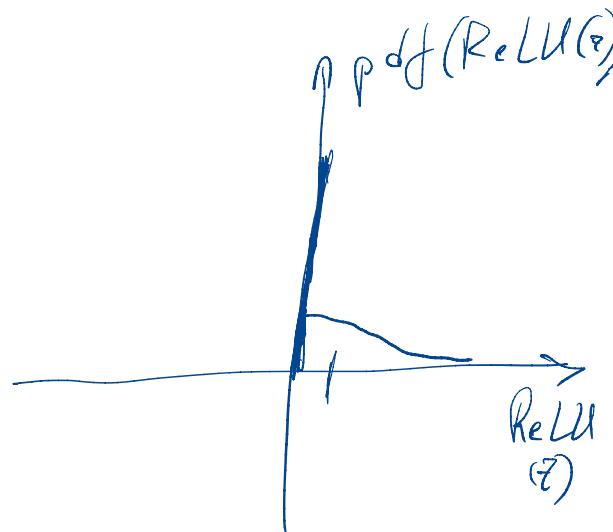
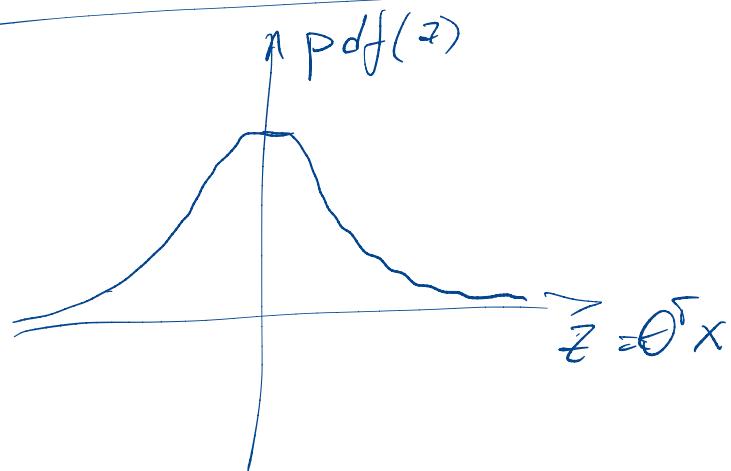


$$\alpha = \frac{\sqrt{6}}{\sqrt{m+k}}$$

$$P \int_{-\alpha}^{\alpha} d\theta = 1 \quad 2\alpha \cdot P = 1 \quad P = \frac{1}{2\alpha}$$

Xavier Glorot

## Batch Normalization



$$z^{(l)} = \theta^T h^{(l-1)}$$

$$\mu_t^* = \mu_{t-1}^* \cdot \beta + E_{B_t}(z)$$

$$z^* = \frac{z - \mu_B(z)}{\sigma_B(z)}$$

$$\sigma_t^* = \sigma_{t-1}^* \cdot \lambda + G_{B_t}(z)$$

$$z^* = \frac{z - \mu_t^*(z)}{\sigma_t^*(z) + \epsilon} \cdot \gamma + \xi \approx 0$$

# Batch Norm

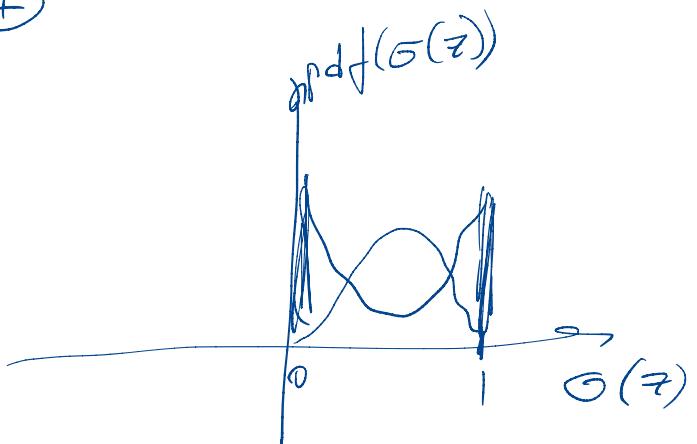
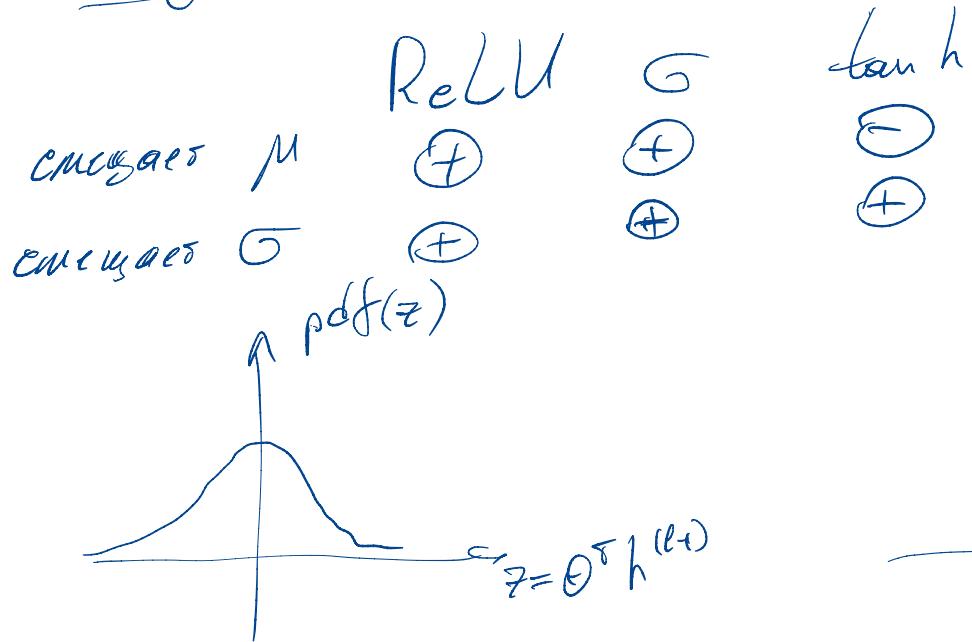
## Pros

- ① Ускоряет сходимость обучения
- ② Тренировка не так сильно зависит от  $\theta$ .
- ③ Сложности эффективны для нейронных сетей.
- ④ Работает как перенормализация.

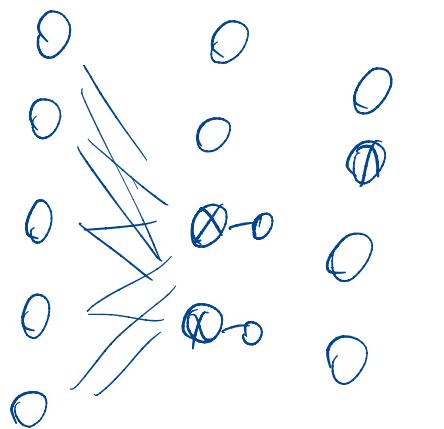
## Cons

- ① Плохо работает при небольших batch-size
- ② Небольшое непредсказуемое изменение

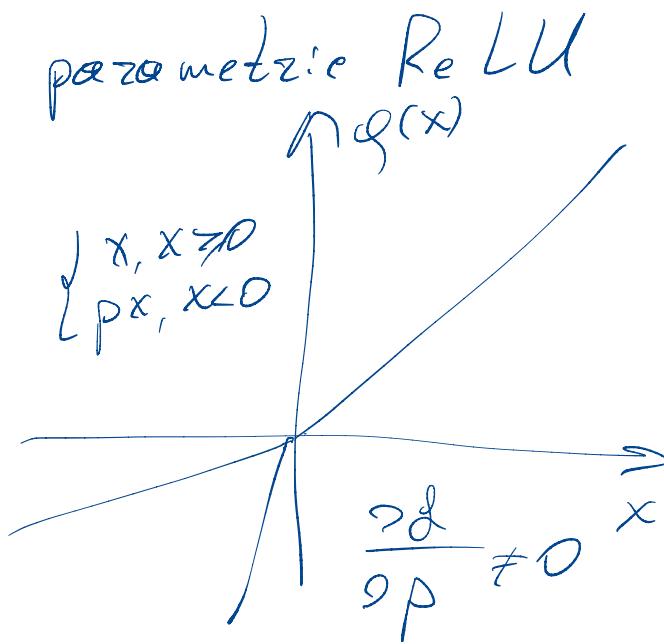
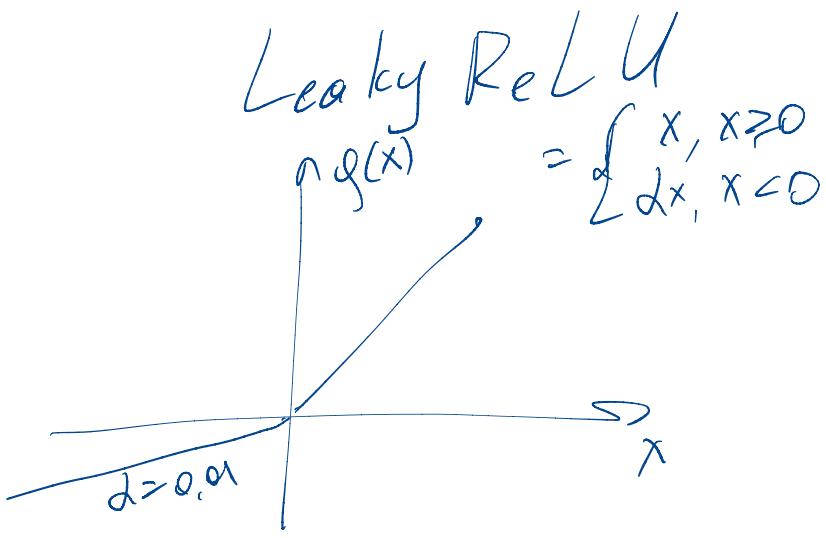
# Physics and Logits



ReLU:



dead neurons



$$\underline{E_{lu}}(x) = \begin{cases} x, & x \geq 0 \\ \alpha(e^x - 1), & x < 0 \end{cases}$$

$$\underline{\text{Если}} \quad \frac{x \sim \mathcal{N}(0, 1)}{\theta_0 \sim \mathcal{N}(0, \frac{\sqrt{6}}{\sqrt{m+k}})}$$

$$\stackrel{\text{TO}}{\underline{\mu(E_{lu}(x)) = \mu(x)}}$$

Таким образом  $\underline{E_{lu}}$   $\mu(E_{lu}(x))$  схожа с  $\mu(x)$   
но xоду обусловлено не так просто, как в   
ReLU или Leaky ReLU