

$$\hat{y} = \psi(\dots \psi(\theta^{(2)} g(\theta^{(1)} \psi(\theta^{(1)} \phi(\theta^{(0)} x))))$$

$h^{(i)}$ — векторы затягивания
hidden representations
embeddings

коэффициентов

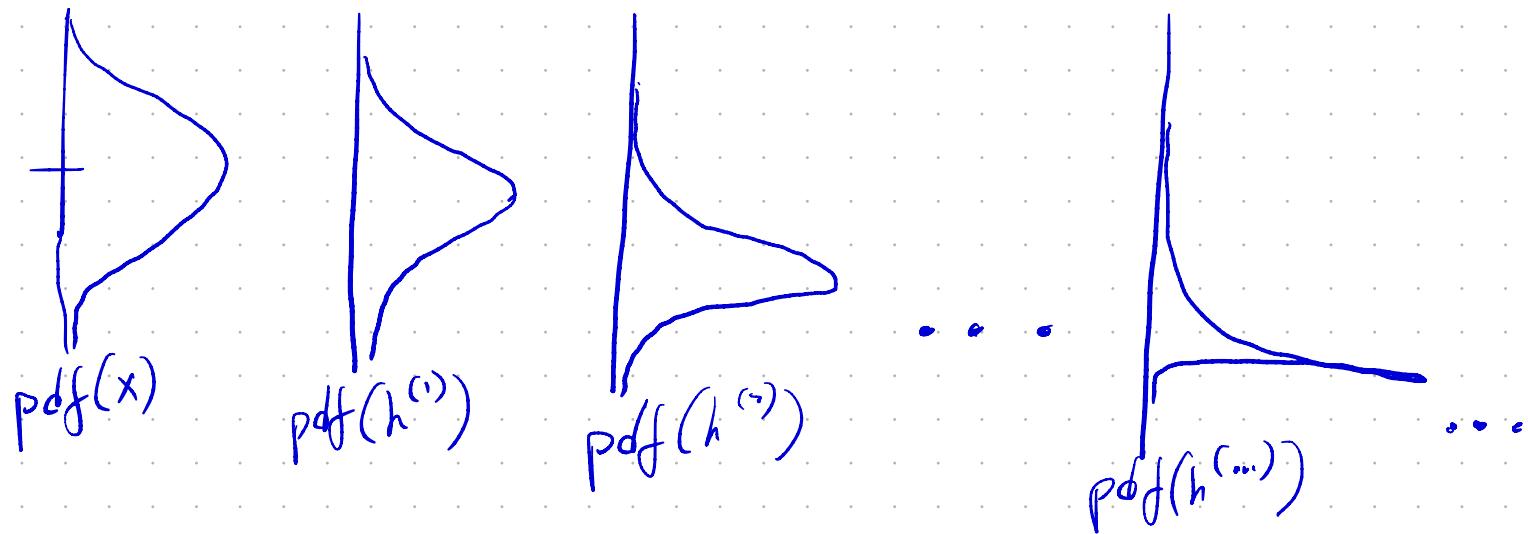
$$\hat{y} \in \mathbb{R}^1 \quad \mathbb{R}^d$$

$$\psi: I(z^{(L)})$$

$$\sigma(z^{(L)})$$

$$\text{Softmax}(z^{(L)})$$

$$\underline{\text{Var}}(h^{(e)}) \approx k \text{Var}(\theta^{(e+1)}) \text{Var}(h^{(e+1)})$$



$$k \text{Var}(\theta^{(\dots)}) < 1$$

$$\text{Var}(\hat{y}) \approx 0 \\ \hat{y} = \text{const}$$

$$k \text{Var}(\theta^{(\dots)}) > 1$$

$$\text{Var}(\hat{y}) \rightarrow \infty \\ \frac{\partial \hat{y}}{\partial x}_{ij} \rightarrow \infty$$

$$k \cdot \text{Var}(\theta^{(\dots)}) \approx 1$$

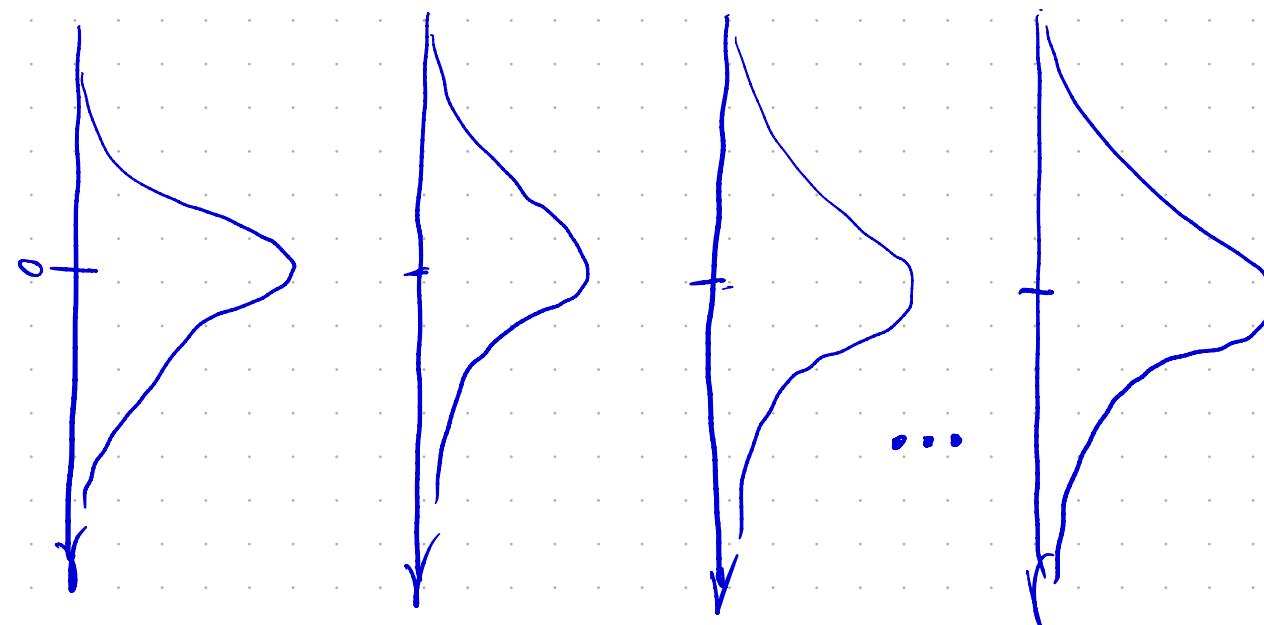
$$\text{Var}(\theta^{(\dots)}) \approx \frac{1}{k}$$

Выведение, применение UHC, inference

$$\hat{y} = F_{MLP}(x)$$

Обучение, оптимизация (from noise), training, learning

$$\theta^* = \operatorname{argmin} L(F_{MLP}(X), Y)$$

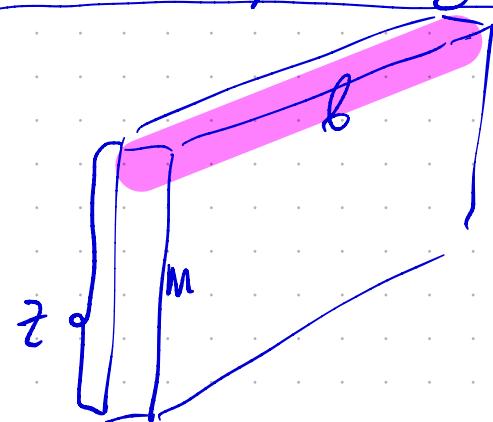


Batch Normalization

Батчевое нормализование

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

$$z^* = \frac{z - \mu_b(z)}{\sigma_b(z)}$$



$$z \in \mathbb{R}^m$$

$$\mu \in \mathbb{R}^m$$

$$\sigma_z \in \mathbb{R}^m$$

$$\mu(z^*) = 0$$

$$\sigma(z^*) = 1$$

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

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

$$z^* = \frac{z - \mu_t^*}{G_t^* + \varepsilon} \cdot \gamma + \xi \quad \begin{cases} \gamma = 1 \\ \xi = 0 \end{cases}$$

Такое обозначение оптимального ряда μ^*, ξ
"наиоптимальнее" μ^*, G^*

Таким образом: μ^*, G^*, γ, ξ функционалы
и не являются