

w - model parameters (weights)

$$\text{Loss}(w, x) = \frac{1}{N} \sum_{i=1}^N (y_{\text{pred}}^i(w, x) - y_{\text{true}}^i)^2 \Rightarrow \min_w \text{WER}^P$$

P is extremely large

Adam or SGD or SGD with momentum

$w, \nabla_w L, v, s$

$w, \nabla_w L, v$

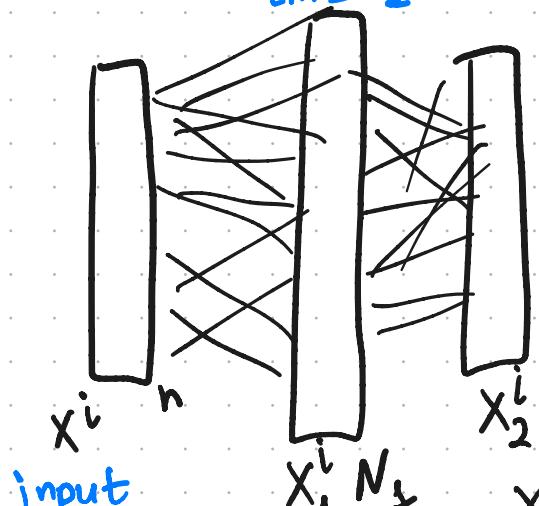
$w, \nabla_w L, v$

3p

4p

2p

$y_{\text{pred}}(w, x)$ - NEURAL NETWORK LAYERS

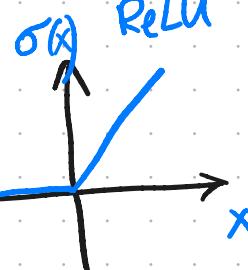


$$X_2^i = W_2 \cdot W_1 \cdot X_1^i$$

$n \times 3$ $n \times n_1$ $n_1 \times n$

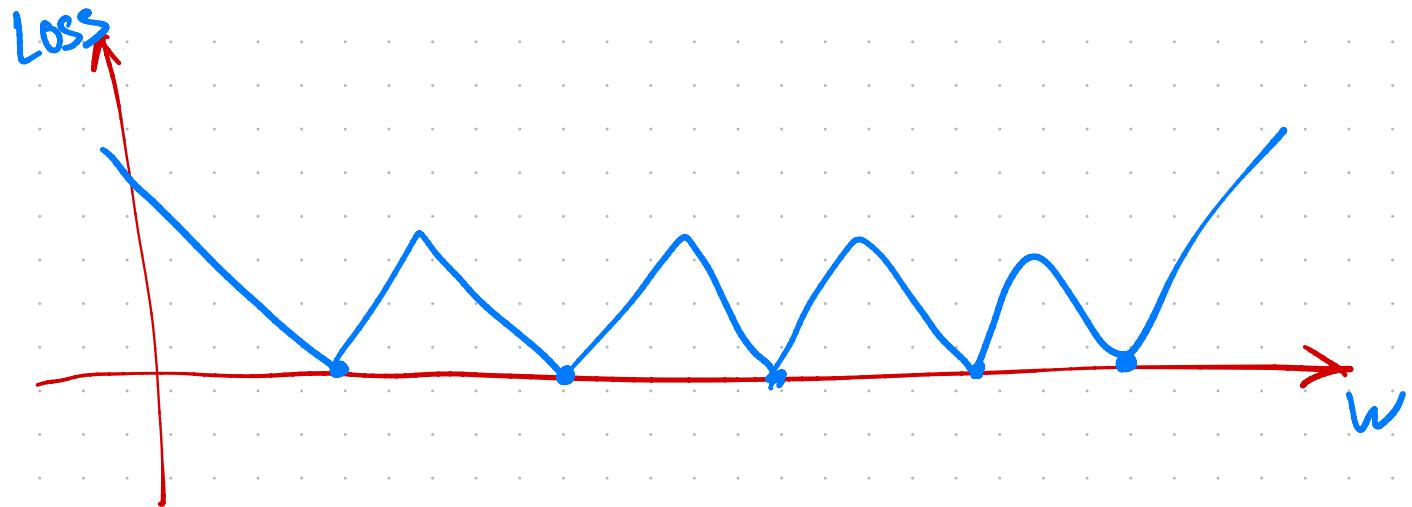
$$x_1^i = \sigma_1(W_1 \cdot x_1^i + b_1)$$

activation function



GPT-3, transformers and etc.
1 GPU?
(255 years)

PROBLEM IS NOT CONVEX.

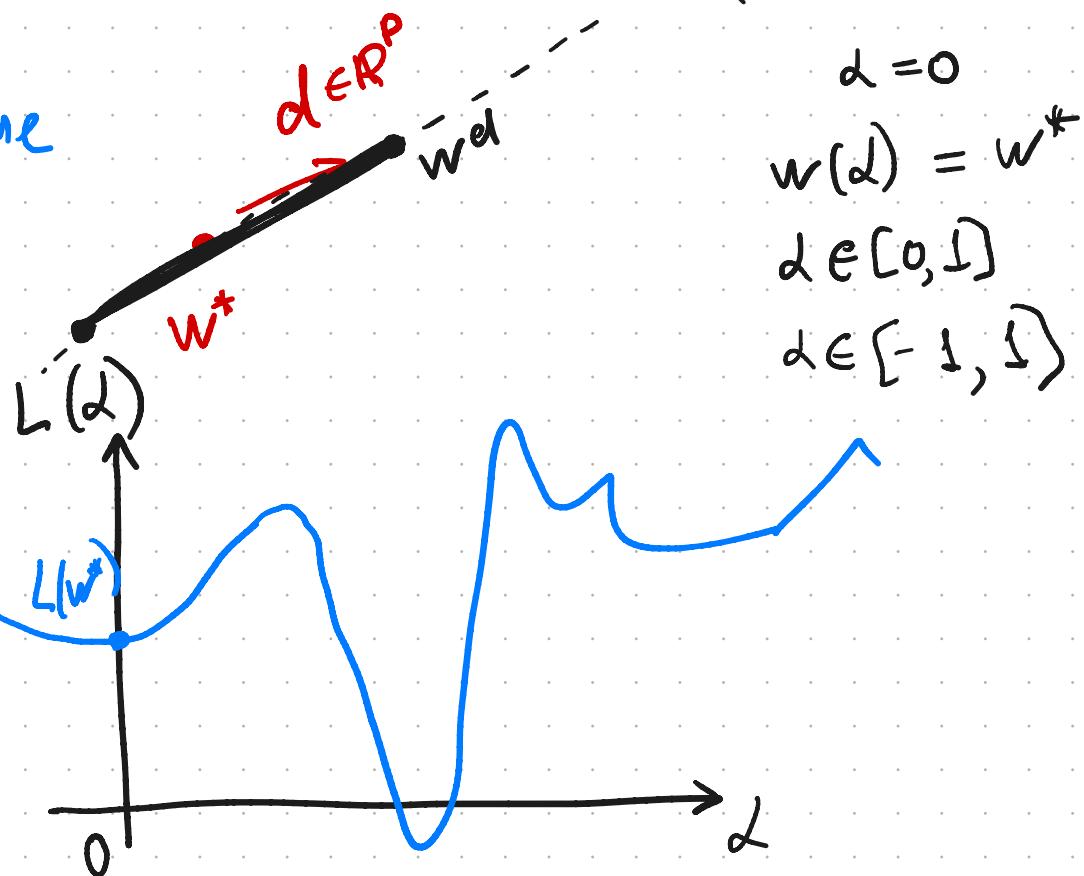


$$L(w) \rightarrow \min_{w \in \mathbb{R}^P}$$

$$w^* \\ w(\lambda) = w^* + \lambda \cdot d$$

$$w(\lambda) = (1-\lambda)w^* + \lambda \cdot W^d$$

projection
on a line



w^* w_1, w_2 $w^*, w, w_1, w_2 \in \mathbb{R}^p$

$$w(\alpha, \beta) = w^* + \alpha \cdot w_1 + \beta w_2$$

$$\alpha \in [-1, 1]$$

$$\beta \in [-1, 1]$$

$$\underline{\alpha, \beta \in \mathbb{R}}$$

$$L(w(\alpha, \beta)) = \dots$$

