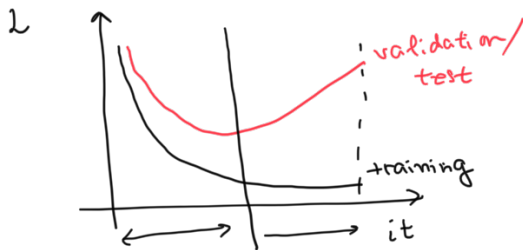


Overfitting and regularization :

$$\mathcal{L}(\theta) := \frac{1}{2} \sum_{i=1}^n (f_{\theta}(x_i) - y_i)^2$$

$$\theta := \{w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}, w^{(0)}, b^{(0)}\}$$

1.) Early stopping :



Idea : Monitor the validation loss and stop the training early (as soon as the validation loss starts to diverge).

2.) ℓ_1 / ℓ_2 - parameter regularization (MAP estimation)

$$\mathcal{L}(\theta) = \underbrace{\frac{1}{2} \sum_{i=1}^n [f_{\theta}(x_i) - y_i]^2}_{\text{data fit}} + \underbrace{\left[\frac{\lambda_1}{2} \|W\|_1 + \frac{\lambda_2}{2} \|W\|_2 \right]}_{\text{model complexity}} \quad \lambda_1, \lambda_2 \in \mathcal{O}(10^{-1} - 10^{-2})$$

$$p(\theta|x, y) \propto \underbrace{p(y|x, \theta)}_{\text{Gaussian likelihood}} \underbrace{p(\theta)}_{\text{prior}}$$

Laplace prior Gaussian prior :
 $\underline{w \sim \mathcal{N}(0, I)}$

3.) Dropout :

Recall : $H^{(e)} = f(H^{(e-1)} W^{(e)} + b^{(e)})$

$$r_j^{(e)} \sim \text{Bernoulli}(\underline{p}) \quad , \quad j = 1, \dots, Q^{(e)} : \# \text{ of neurons in the } e\text{-th layer}$$

$$Z^{(e)} = r^{(e)} \odot H^{(e-1)}$$

$$H^{(e)} = f(z^{(e)} \cdot W^{(e)} + b^{(e)})$$

4.) Data augmentation (Classification tasks)

$$\tilde{x} = \frac{x - E[x]}{\text{std}[x]}$$

Network initialization :

$\mathcal{D} := \{x, y\}$ is a "standardized" data-set : $E[x] = 0, \text{Var} = 1$
 $\xrightarrow{\text{dim}\{d_{in}\}} \xrightarrow{\text{dim}\{d_{out}\}}$
 i.i.d. $E[y] = 0, \text{Var} = 1$

Linear regression : $y = w_1 x_1 + \dots + w_d x_{d_{in}}$, w_i are zero-mean

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

$$\xRightarrow{\text{i.i.d.}} \text{Var}[y] = d_{in} \cdot \text{Var}[x_i] \text{Var}[w_i]$$

$$\rightarrow \text{Var}[w_i] = \frac{1}{d_{in}} \quad (1)$$

Repeat this analysis for the back-propagated gradient signal

$$\rightarrow \text{Var}[w_i] = \frac{1}{d_{out}} \quad (2)$$

Empirical rule for initializing w :

$$\text{Var}[w_i] = \frac{2}{d_{in} + d_{out}} \quad \checkmark$$

Glorot initialization :

$$W \sim \mathcal{U}\left[-\frac{\sqrt{6}}{\sqrt{d_{in}+d_{out}}}, \frac{\sqrt{6}}{\sqrt{d_{in}+d_{out}}}\right]$$

Glorot
Uniform
initialization

$$W \sim \mathcal{N}\left(0, \frac{2}{d_{in}+d_{out}}\right)$$

Glorot
normal
initialization