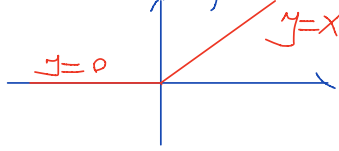


- o Linear Model

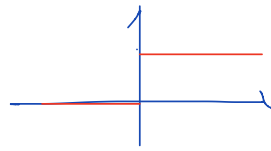
- Limited: no cross term $x_1 \cdot x_2 \rightarrow$ introduce nonlinearity
- Efficient: GPU
- Stable

- Rectified Linear Unit (ReLU)

Simplest non-linear function



Derivative

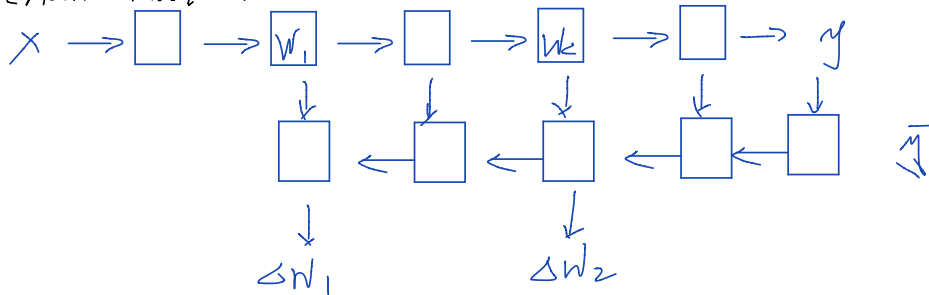


- Network of ReLUs

$$X \times \boxed{w} + \boxed{b} \rightarrow y \rightarrow \text{Sigmoid} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$X \times [w_1] + [b_1] \rightarrow \begin{matrix} \text{ReLU} \\ \text{ReLU} \\ \text{ReLU} \end{matrix} \times [w_1] + [b_1] \rightarrow y \rightarrow S(y)$$

- The chain Rule :



• Training a Deep Neural Network

- Going deeper rather than wider:

wider (more neurons in single layer): not efficient, hard to train.

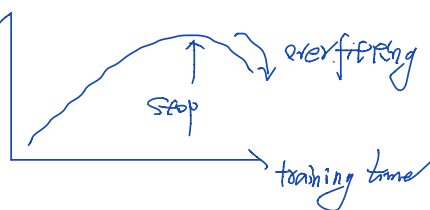
deeper (more layers): parameter efficiency. Capture hierarchical structure in deep model.

• Regularization

- Avoid overfit:

- Early Termination

validation set performance



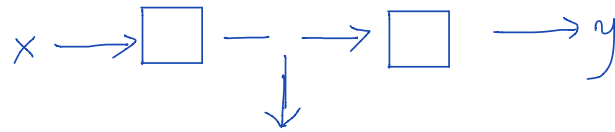
- Regularization

Artificially add constraints on the parameters, implicitly reduce free parameters

• L2 Regularization

$$\mathcal{L}' = \mathcal{L} + \underbrace{\beta \frac{1}{2} \|W\|_2^2}_{\frac{1}{2} (W_1^2 + W_2^2 + \dots + W_n^2)} \quad \mathcal{L}: \text{loss} \quad \beta: \text{hyperparameter}$$

• Dropout (regularization)



Take activations from one layer next layer, randomly drop half of them (set to 0), then repeat the process.

Intuition: Your network should never rely on any given activation to be present, because they may get squashed at any given moment so it's a way to force to learn the redundant representation, prevent overfitting.

Introduced by Geoffrey Hinton

Training and Evaluation:

• Training

randomly drop half of activations, at the same time double the value of remaining activations get output \hat{y}_t .

repeat this for many times, get the average of prediction $\hat{y}_e = \overline{\hat{y}_t}$ then on average this is equivalent to not dropping the activations. Where \hat{y}_e is the evaluation of dropout

