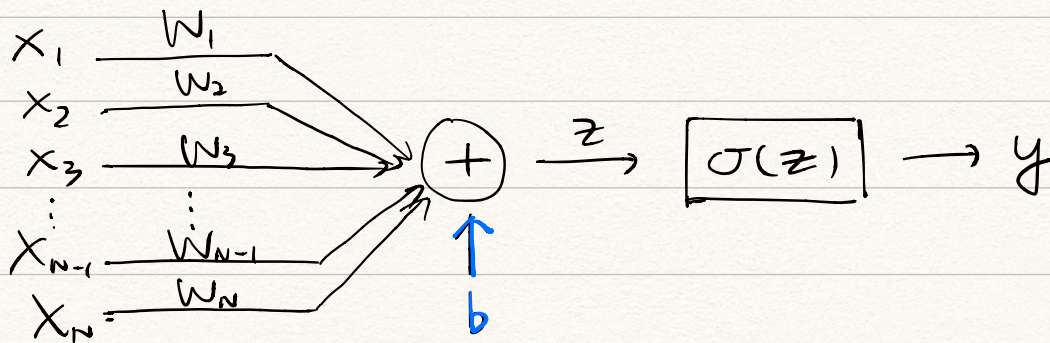
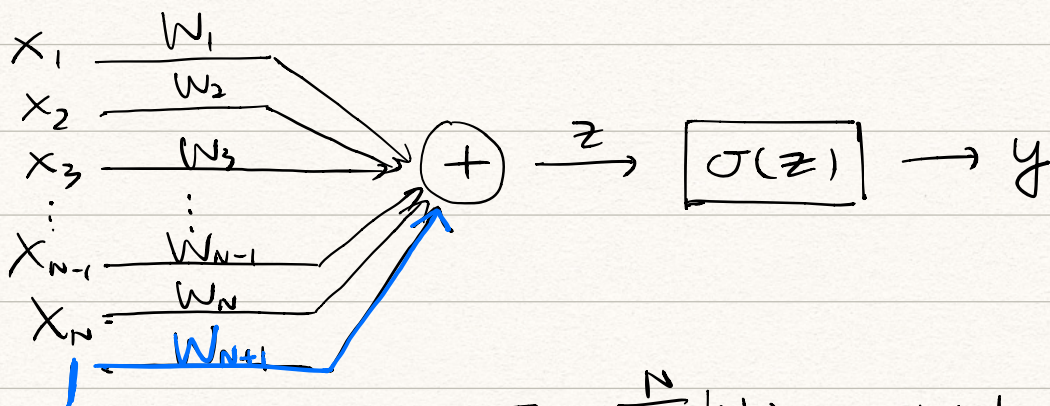


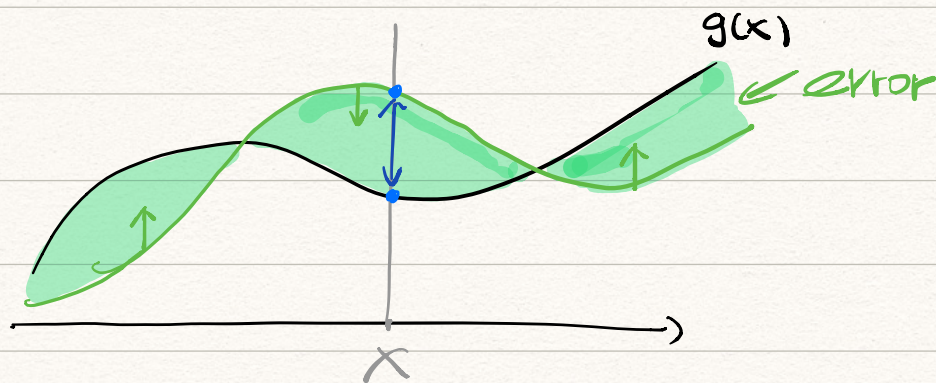
Preliminaries: The units in the network



$$z = \sum_{i=1}^N w_i x_i + b$$



$$z = \sum_{i=1}^N w_i x_i + 1 \cdot w_{n+1}$$



Divergence:

$$\hat{w} = \underset{w}{\operatorname{argmin}} \int_x \operatorname{div}(f(x; w), g(x)) dx$$

**NOTICE:** Learning a neural network  $\Leftrightarrow$

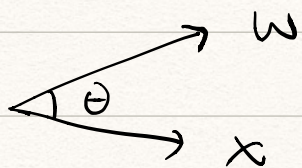


determining the parameters of the network (weights and biases) required for it to model a desired function.

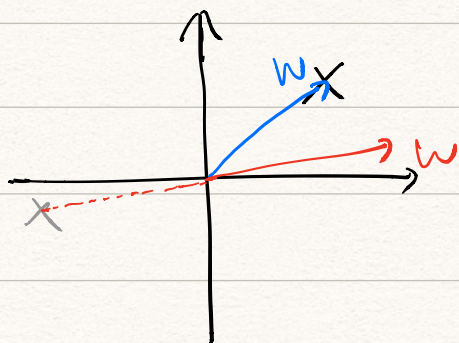
$$W^T X = 0 \Rightarrow \underline{X} \text{ orthogonal } \underline{W}$$

$$y = \text{sign}(W^T X) \quad \text{"inner product"}$$

$$W^T X = |W| |X| \cos \theta$$



$$-\frac{\pi}{2} < \theta < \frac{\pi}{2}$$



1.  $W$  is orthogonal to the plane

2. positive instances:  $W \xrightarrow[\mathbb{R}]{\mathbb{R}} x$

Perceptron learning algorithm:

Given  $N$  training instances  $(x_1, y_1), (x_2, y_2), \dots$



$$\star Y_i = +1 \text{ or } -1$$

② Initialize  $W$

③ Cycle through the training instances

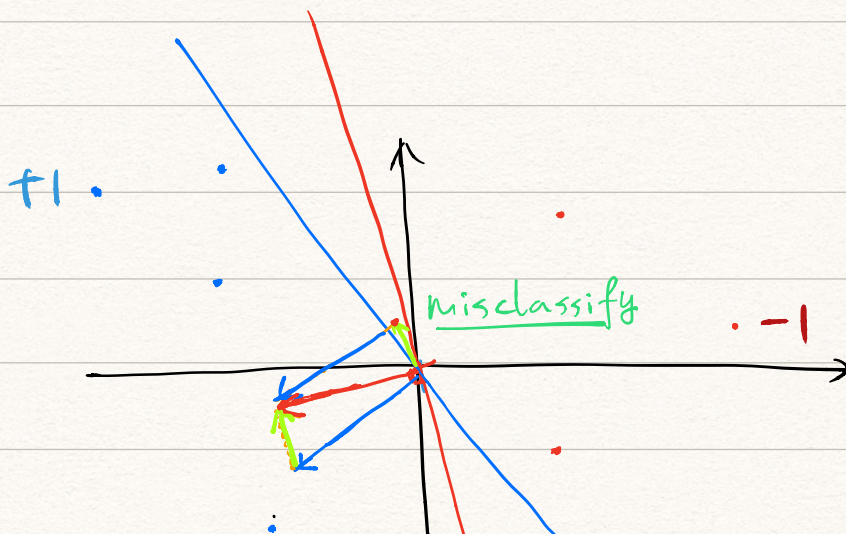
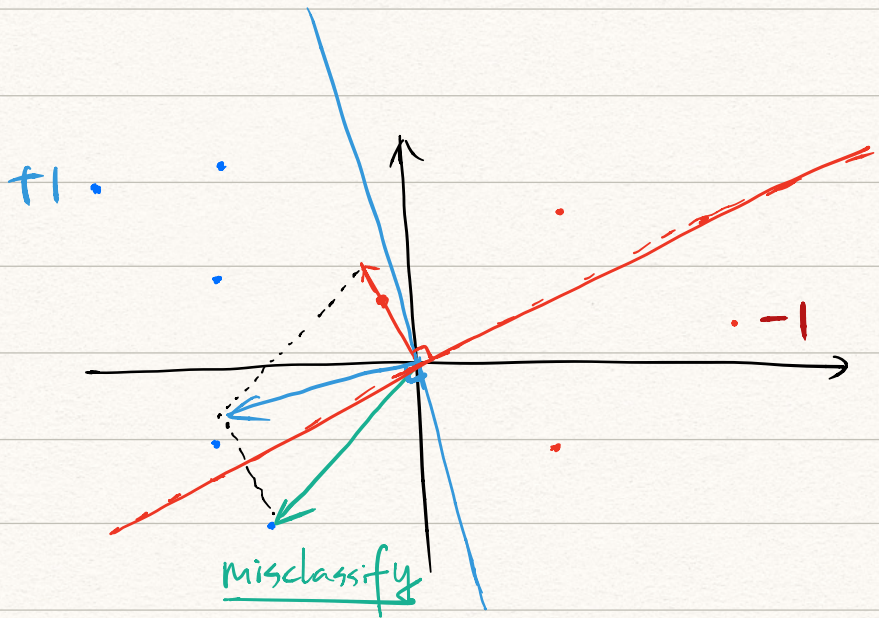
④ While more classification errors:

For  $i = 1, 2, \dots, N_{\text{train}}$

$$O(X_i) = \text{sign}(W^T X_i)$$

if  $O(X_i) \neq Y_i$

$$W = W + Y_i X_i$$





1972 Paul Werbos MIT PhD "back propagation"

differentiable activation function:

$$\frac{dy}{dw_i} = \frac{dy}{dz} \cdot \frac{dz}{dw_i} = \sigma'(z) x_i$$

$$\frac{dy}{dx_i} = \frac{dy}{dz} \cdot \frac{dz}{dx_i} = \sigma'(z) w_i$$

input-output pairs:  $(x_1, d_1), (x_2, d_2), \dots, (x_N, d_N)$

$$\hat{W} = \arg \min_w \int_X \text{div}(f(x; w), g(x)) dx$$

$$\begin{aligned} \hat{W} &= \arg \min_w \int_X \text{div}(f(x; w), g(x)) p(x) dx \\ &= \arg \min_w E[\text{div}(f(x; w), g(x))] \end{aligned}$$

$$E[\text{div}(f(x; w), g(x))] \approx \frac{1}{N} \sum_{i=1}^N \text{div}(f(x_i, w), d_i)$$



## function Minimization

$$\left\{ \begin{array}{l} \text{Err}(W) = \frac{1}{N} \sum_i \text{div}(f(x_i; W), d_i) \\ \hat{W} = \underset{W}{\operatorname{argmin}} \text{Err}(W) \end{array} \right.$$