

True label		
TP	FP	Precision = $\frac{TP}{TP+FP}$
FN	TN	Recall = $\frac{TP}{TP+FN}$

$$\text{Accuracy} = \frac{\text{TP+TN}}{\text{TP+FP+TN+FN}}$$

$$F_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{TP rate} = \frac{\text{TP}}{\text{P}} = \text{Recall}$$

$$\text{FP rate} = \frac{\text{FP}}{\text{P}}$$

Learning curve → identify how much data needed
 increase rapidly → get more data
 Completely flat → more data can help
 Gradual increase → need a lot more data

Train Acc → decrease with more data

Validation Acc → increase with more data

→ don't close overfitting
 → close too fast underfitting

ROC Graph is TP rate (y) vs FP rate (x), used to depict trade offs between benefits (TP) and costs (FP). Note goes through (0,0) (1,1)

Bias → inherent error
 Variance → how over specialized

$$\text{Regularized Loss: } J_{\lambda}(\theta) = J(\theta) + \lambda R(\theta)$$

- $J(\theta)$ → training loss / cost function
- $R(\theta)$ → regularizer
- λ Regularization parameter

$$\min \frac{1}{2} \sum_{i=1}^m (y_i - \theta^T \cdot x_i)^2 + \frac{\lambda}{2} \|\theta\|^2$$

Linear Regression:

- $\theta = \emptyset$ → no regularization → maybe overfitting
- if λ too big → forcing all parameters to be close to 0

Sign unit
 1 when > 0
 -1 otherwise

$$g = g(w_0 + w^T X)$$

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \quad W = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{pmatrix}$$

linear comb inputs
 w_0 is bias
 $g = g\left(w_0 + \sum_{i=1}^m w_i \cdot x_i\right)$ non activation linear function

Sigmoid function
 $g(z) = \sigma(z) = \frac{1}{1+e^{-z}}$
 $g'(z) = \sigma'(z) (1 - \sigma(z))$

$$\text{Hyperbolic Tangent}$$

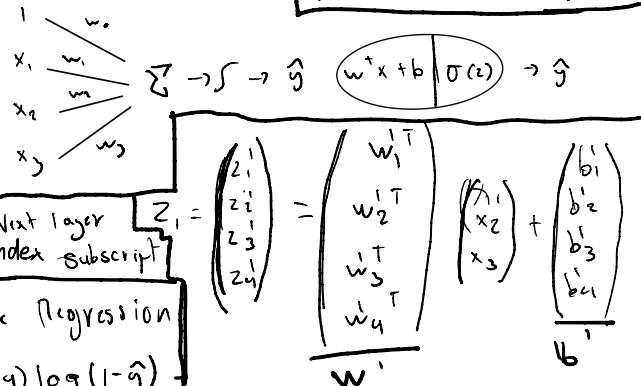
$$g(z) = \sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

$$\text{Rectified Linear Unit (ReLU)}$$

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$



Gradient Rule

$$w^1 \leftarrow w^1 - \alpha dW^1$$

Forward Prop
 $\rightarrow z$ values

Backward Prop
 \rightarrow derivative

Hyperparameters

- learning rate α

- # hidden layer L

- # hidden units n_1, n_2

- choice of activation functions

Build ML System:

Start with idea,

implement to code,

experiment to test idea worked

$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

- $y=1 \rightarrow$ want $-\log \hat{y} \rightarrow$ want \hat{y} large
- $y=0 \rightarrow$ want $\log (1-\hat{y}) \rightarrow$ want \hat{y} small

Reinforcement learning
 key input → Reward Func
 State S → Action A
 $\Omega(s)$ is input for the next second layer
 $\text{Return} = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$

Discount Factor

$$\gamma: 0 < \gamma < 1$$

$\Omega(s)$ is input for the next second layer

$$z^1 = w_3^T x + b_3, a_3^1 = \sigma(z_3^1)$$

$$z_4^1 = w_4^T x + b_4, a_4^1 = \sigma(z_4^1)$$

$$z_1^1 = w_1^T x + b_1, a_1^1 = \sigma(z_1^1)$$

$$z_2^1 = w_2^T x + b_2, a_2^1 = \sigma(z_2^1)$$

Bellman Equation

$$\Omega(s, a) = R(s) + \gamma \max_{a'} (\Omega^1(a'))$$

Bellman max is \max of
 \rightarrow or \leftarrow input on state

$$z = w^T x + b \rightarrow a = \sigma(z) \rightarrow L(a, y)$$

$$\frac{da}{d\alpha} = \frac{dL(a, y)}{d\alpha} = -\frac{y}{\alpha} + \frac{(1-y)}{1-\alpha} = -\frac{y}{\alpha} - \frac{(1-y)}{1-\alpha}$$

$$dz = da \cdot \frac{da}{dz} = da \cdot \sigma'(z) = -\frac{y}{\alpha} - \frac{(1-y)}{1-\alpha} \times \alpha(1-\alpha) = \alpha - y$$

$$dw = dz \cdot \frac{dz}{dw} = dz \cdot x \quad \frac{dH}{dV} = 3$$

$$db = dz \cdot \frac{dz}{db} = dz \quad \frac{dH}{da} = \frac{dH}{dV} \cdot \frac{dV}{da} = 3$$

$$\frac{dH}{dv} = \frac{dH}{db} \cdot \frac{db}{dv} = 3 \times 2 = 6$$

$x_1 \xrightarrow{3} x_2 \rightarrow \dots \rightarrow y$

 $y = g(w_0 + w^T x)$
 $= g(1 + \begin{pmatrix} 3 \\ -2 \end{pmatrix}^T \begin{pmatrix} x_1 \\ x_2 \end{pmatrix})$
 $= g(1 + 3x_1 - 2x_2)$

given input $\begin{pmatrix} 1 \\ 2 \end{pmatrix}$

 $y = g(1 + 3(-1) - 2(2))$
 $= g(-6) = \frac{1}{1+e^{-6}} = 0.002$

$a = 5$

 $b = 3 \rightarrow u = bc$
 $c = 2$
 $\frac{dH}{dc} = \frac{dH}{du} \cdot \frac{du}{dc} = 3 \times 3 = 9$

$v = a + u \rightarrow H = 3v$

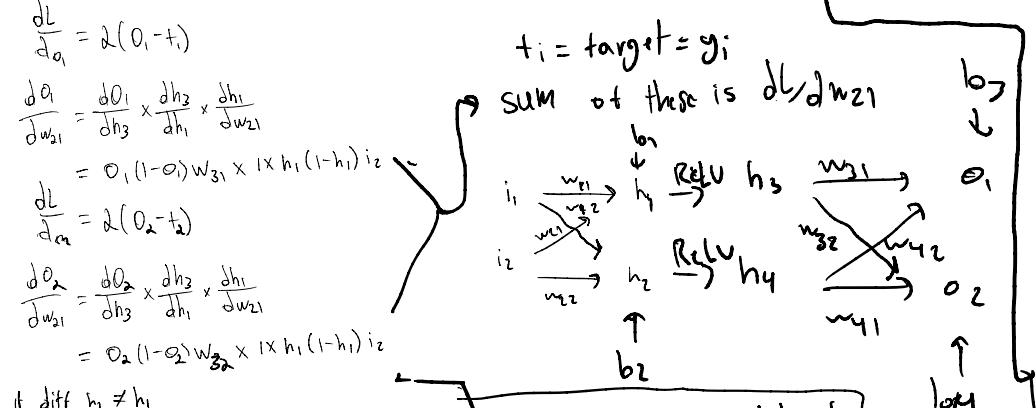
$m = \# \text{ inputs pair or sight}$

Gradient Decent

Perception tree use values to scale branch then match then multiply resulting to node and add bias \rightarrow the perception formula

 $AUC = \frac{1}{2} \sum_{i=1}^m (FPR_i \cdot TPR_i) (TPR_i + FPR_i)$
 $L = \frac{1}{m} \sum_{i=1}^m (o_i - y_i)^2$
 $\min J(\theta_1, \theta_2)$
 $\theta_1, \theta_2 \in \mathbb{R}$

$\frac{dL}{dw_{21}} = \frac{dL}{d\theta_1} \times \frac{d\theta_1}{dw_{21}} + \frac{dL}{d\theta_2} \times \frac{d\theta_2}{dw_{21}}$
 $\leftarrow w_{21} \leftarrow w_{21} - \alpha \frac{dL}{dw_{21}}$
 $\frac{dL}{d\theta_1} = 2(o_1 - t_1)$
 $\frac{d\theta_1}{dw_{21}} = \frac{d\theta_1}{dh_3} \times \frac{dh_3}{dh_1} \times \frac{dh_1}{dw_{21}}$
 $= o_1(1-o_1)w_{31} \times 1 \times h_1(1-h_1)t_1$
 $\frac{dL}{d\theta_2} = 2(o_2 - t_2)$
 $\frac{d\theta_2}{dw_{21}} = \frac{d\theta_2}{dh_3} \times \frac{dh_3}{dh_1} \times \frac{dh_1}{dw_{21}}$
 $= o_2(1-o_2)w_{32} \times 1 \times h_1(1-h_1)t_2$



Plug into Activation derivative

 $\frac{dh_3}{dz_3} = h_3 \cdot (1-h_3)$
 $z_3 = w_{31}h_1 + w_{32}h_2 + b_3$
 $\frac{dz_3}{dh_1} = w_{31}$

in gradient descent if weights set zero only need to compute one set of hidden units they're symmetric