RECALL THE DISCRIMINATIVE MODELING
APPROACH FROM Lecture 03:

$$X \to f(X;\theta) \to \boxed{Z} \to \sigma_{SOFTMAX}(Z) \to \hat{Y}$$

$$\nabla_\theta NLL \leftarrow NLL(\theta; Y, \hat{Y}, X)$$

$$Y$$

SOFTMAX REGRESSION : $f(X;\theta) : XW^T + b$

- CONVEX
- LINEAR IN X
- DECISION Boundary Lies DIRECTLY IN THE SPACE at the INPUT X

XOR PROBLEM

$X_1$

- $y = +1$    ○ $y = 0$

$\Rightarrow$ A LINEAR manifold
fails to Classify
X CORRECTLY!

- $y = 0$    ● $y = +1$

$------ X_D$

ANN with Single Hidden Layer

$$f(X; \theta) = \left[ \sigma_a(X W^{(1)T} + b^{(1)}) W^{(2)T} + b^{(2)} \right]$$

where

$$\theta = \{ W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)} \}$$

$\sigma_a(\cdot)$ is an "Activation function"

↳ Apply this to the XOR problem

LET  $W^{(1)} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$  $\in \mathbb{R}^{2 \times 2}$

$b^{(1)} = [0 \ -1]$  $\in \mathbb{R}^2$

$W^{(2)} = [1 \ -2]$  $\in \mathbb{R}^2$

$b^{(2)} = 0$  $\in \mathbb{R}$

$y \in \{0, 1\}$

$X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$

$\sigma_a(z_j) = \max(0, z_j) \Rightarrow$ "Rectified Linear Unit"

"RELU"

ANN / XOR cont...

$$\to XW^{(1)T} + b^{(1)} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} + [0 \ -1] = \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \to Z^{(1)}$$

$$\to \sigma_a(Z^{(1)}) = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \to a^{(1)}$$

$$\to a^{(1)} W^{(2)T} + b^{(2)} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} + 0 = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} \overset{\hat{y}}{\Longleftrightarrow} \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \overset{X}{}$$

$\hookrightarrow$ CORRECT XOR
class ASSIGNMENTS!

SLIDE 10

UNIVERSAL APPROXIMATION Theorem

- LINEAR MODELS SUCH AS $XW^T + b$ CAN ONLY MODEL LINEAR functions

- The ACTIVATION functions Allow ANNs to learn NON-linear manifolds in the input SPACE.

$\longrightarrow$ "UNIVERSAL APPROXIMATION Theorem"
   - AN ANN with AT least one activation layer can APPROXIMATE "ANY" Borel-measurable function in finite dimensions when enough hidden units. This extends to the derivatives of the function.
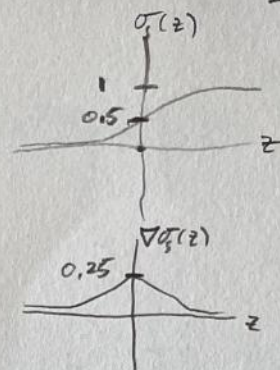
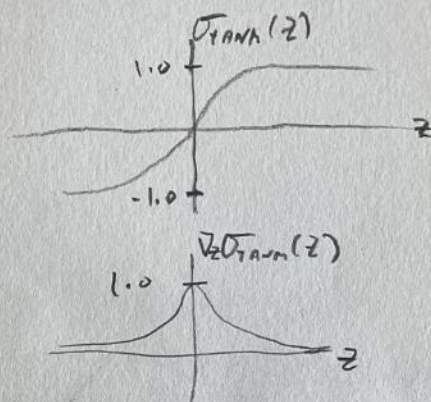# ACTIVATION FUNCTIONS

SIGMOID
$$\sigma_s(z) = \frac{1}{1 + e^{-z}}$$

$$\nabla_z \sigma_s(z) = \sigma_s(z)(1 - \sigma_s(z))$$

Hyperbolic
Tangent
$$\sigma_{TANH}(z) = \frac{2}{1 + e^{-2z}} - 1$$

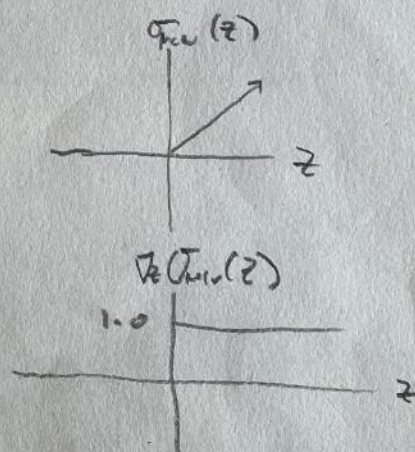$$\nabla_z \sigma_{TANH}(z) = 1 - \sigma_{TANH}(z)^2$$

RECTified
LINEAR
UNIT
$$\sigma_{relu}(z) = \begin{cases} z & z \geq 0 \\ 0 & else \end{cases}$$

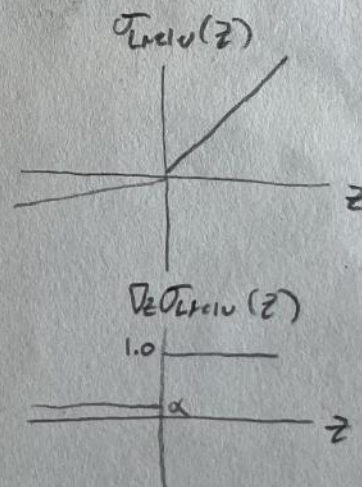$$\nabla_z \sigma_{relu}(z) = \begin{cases} 1 & z \geq 0 \\ 0 & else \end{cases}$$

Leaky
Relu
$$\sigma_{Lrelu}(z) = \begin{cases} z & z \geq 0 \\ \alpha & else \end{cases}$$

$$\nabla_z \sigma_{Lrelu}(z) = \begin{cases} 1 & z \geq 0 \\ \alpha & else \end{cases}$$

# ANN Learning

RECALL from lecture 03:

$$\hat{\theta} = \underset{\theta}{\text{ARGMIN}} \; NLL(\theta; DATA)$$

$$= \underset{\theta}{\text{ARGMIN}} \; -\sum_{i=1}^{M} f(x^{(i)};\theta)_{y^{(i)}} - \lg \sum_{k=0}^{k-1} e^{f(x^{(i)};\theta)_k}$$

$\hat{\theta}$ found VIA GRADIENT DECSENT:

$$\nabla_\theta NLL = -\sum_{i=1}^{M} \nabla_\theta f(x^{(i)};\theta)_{y^{(i)}} - \nabla_\theta \sum_{k=0}^{k-1} e^{f(x^{(i)};\theta)_k}$$

⟩ SOFTMAX REGRESSOR : $f(x;w) = xw^T + b$

  - NLL IS CONVEX
  - ANY LOCAL MINIMUM IS A GLOBAL MINIMUM

⟩ ANN

  - ~~NLL IS NOT CONVEX~~

  ~~MANY LOCAL MINIMA~~

  - $\nabla_\theta f(x;\theta)$ REQUIRES SUCCESSIVE APPLICATION OF THE CHAIN RULE, AKA BACK PROPAGATION

      ⤷ LECTURE 08

  - NLL IS NON-CONVEX, MANY LOCAL MINIMA

# STATISTICAL FOUNDATION for "REGULARIZATION"

RECALL THAT MLE finds $\hat{\theta}_{MLE}$ that maximizes
the likelihood of the assumed data:

$$\hat{\theta}_{MLE} = \underset{\theta}{Argmax} \sum_{D} P(D; \theta)$$

-D LETS USE BAYES Theorem TO EXPRESS AN
RELATED OBJECTIVE:

$$\underset{R.V.}{\underline{P(\theta \mid D^{(i)})}} \leftarrow L(D; \theta) \blacktriangleleft P(X,Y) \rightarrow L(\theta; D) + \underset{NOT\ R.V.}{\underline{P(D; \theta)}}$$

NEW
⌐D OBJECTIVE

$$\hat{\theta}_{MAP} = \underset{\theta}{Argmin} - \sum_{D} \lg P(\theta \mid D)$$

$$= \underset{\theta}{Argmin} - \sum_{D} \log \left[ \frac{P(D \mid \theta) P(\theta)}{P(D)} \right]$$

$$= \underset{\theta}{Argmin} - \sum_{D} \log P(D; \theta) + \log P(\theta) - \log P(D)$$

$$= \underset{\theta}{Argmin} - \sum_{D} \log P(D; \theta) + \log P(\theta)$$

## 2 IMPORTANT THINGS:

① WE STARTED BY DESCRIBING $\theta$ AS A R.V.
IN $L(D; \theta)$, BUT HAVE REDUCED THE
PROBLEM TO WHERE $\theta$ REPRESENTS A POINT
ESTIMATE!

② THE DIFFERENCE BETWEEN $\hat{\theta}_{MAP}$ AND $\hat{\theta}_{MLE}$ LIES IN

REGULARIZATION CONT...

Look AT Common PARAMETERIZATIONS for $P(\theta)$

① $\theta \sim \text{Unif}(\lambda)$

$\hat{\theta}_{MAP} = \text{ARGMIN}_{\theta} - \sum_{D} \log P(D;\theta) + \log \text{Unif}(\lambda)$

$= \text{ARGMIN}_{\theta} - \sum_{D} \log P(D;\theta)$

$= \hat{\theta}_{MLE}$

THUS, MAP IS A GENERALIZATION of MLE

② $\theta \sim N(0, \sigma_{\theta}^2)$ - GAUSSIAN PRIOR ON $\theta$

$\hat{\theta}_{MAP} = \text{ARGMIN}_{\theta} - \sum_{D} \log P(D;\theta) + \log\left[\frac{1}{\sqrt{2\pi\sigma_{\theta}^2}} e^{-\frac{\theta^2}{2\sigma_{\theta}}}\right]$

$= \text{ARGMIN}_{\theta} - \sum_{D} \log P(D;\theta) + \frac{\theta^2}{2\sigma_{\theta}}$

$\boxed{L_2 \text{ REGULARIZATION}}$

③ $\theta \sim \text{Laplace}(0, \sigma_{\theta})$

$\hat{\theta}_{MAP} = \text{ARGMIN}_{\theta} - \sum_{D} \log P(D;\theta) + \log\left[\frac{1}{2\sigma} e^{-\frac{|\theta|}{\sigma_{\theta}}}\right]$

$= \text{ARGMIN}_{\theta} - \sum_{D} \log P(D;\theta) + \frac{|\theta|}{\sigma_{\theta}}$

$\boxed{L_1 \text{ REGULARIZATION}}$

# REGULARIZATION Cont...

## DROPOUT

Randomly selects nodes in each hidden layer, $\ell_i$ and sets them to zero during training.

This is an element-wise operation:

$$a'^{(\ell)} = \gamma \; m \odot a^{(\ell)}$$

where

$$M_j \sim \begin{cases} 0 & \text{with } P = P_{dropout} \\ 1 & \text{with } P = 1 - P_{dropout} \end{cases},$$

$\gamma$ is a const function of $P_{dropout}$

Thus:

During training: $\quad Z_j^{(\ell+1)} = a'^{(\ell)} W_j^{T(\ell+1)}$

During inference: $\quad Z_j^{(\ell+1)} = a^{(\ell)} W_j^{T(\ell+1)}$

## GRADIENT MOMENTUM

### SGD w/ simple momentum

Set momentum $\alpha$
Set learning rate $h$
Initialize $\theta$, velocity $V$

Repeat:
$\quad X, y \sim$ data minibatch size $m$
$\quad \nabla_\theta$ from $\frac{1}{m} \sum_{i=1}^{m} LL(x^{(s)}, y^{(s)}) \theta)$
$\quad V = \alpha V - h \nabla_\theta$
$\quad \theta = \theta + V$

Until: stopping condition

### POPULAR VARIANTS

- NESTEROV
- RMS PROP
- ADAM