

①

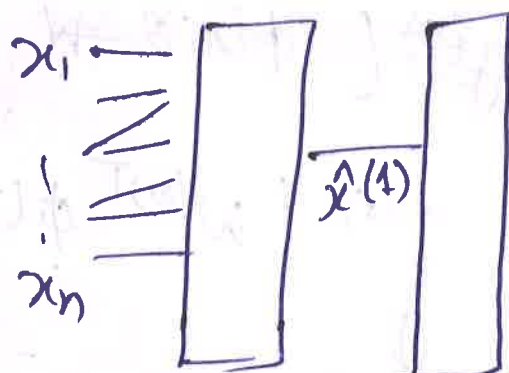
COL 865
Deep Learning
July 31, 2017

Last class -

- (I) Motivating Applications
- (II) Standard Architectures
↳ 5 of them
CNNs, RNNs, Autoencoders, GANs, DeepRL
- (III) Implementation: PyTorch
- (IV) Feedforward Networks

L : - # of layers

Notation:



m : - # of examples

n : - Dimensionality of x

K Hidden Units

$$x^{(0)} = \begin{bmatrix} 1 \\ x \end{bmatrix}$$

$n^{(k)}$: - Dimensionality of output coming out of layer k .

~~$W^{(k)}$~~ $n^{(k)} \times n^{(k+1)}$

\Downarrow
Weights in layer k

Feature representation
 $\phi(x; \theta)$

$$R^{n^{(k-1)} \times n^{(k)}} \quad (k \geq 1)$$

$$\theta = (\theta^{(1)} \dots \theta^{(k)})$$

$$\theta = (\theta, \theta^{(1)})$$

$$\phi = \phi^{(k)}(\phi^{(k-1)} \dots \phi^{(1)}(x, \theta^{(1)}), \dots)$$

$$f(x; \theta) = \hat{y}$$

$$\hat{y} \in n(\theta) \rightarrow \text{Output Layer}$$

$$\hat{y} \in \{0, \dots, L\}$$

of classes

$$y = f^{(0)}(\phi(x; \theta); \theta^{(0)})$$

$$\hat{x}^{(k+1)}$$

$$= g^{(k)}$$

Elementwise Application

$$\hat{x}^{(k)} = g^{(k)}[\hat{x}^{(k-1)} W^{(k)T} + b^{(k)}]$$

ReLU

$$h^{(k)}(x) \equiv g^{(k)}(W^{(k)T} x + b^{(k)}) \quad \text{Rectified Linear Unit}$$

$$\hat{y} = g^{(0)}(W^{(0)T} \phi(x; \theta) + b^{(0)})$$

$$y = g^{(0)}(W^{(0)T} \phi(x; \theta) + b^{(0)})$$

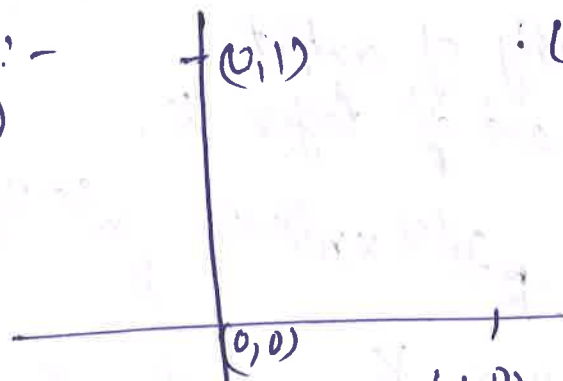
↓

Logistic / Softmax Function

Now, - Example (Non-linearity).

XOR:-

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



(1,1)

Not linearly separable

Transform:-
ReLU

$$g(W^T x + b) \quad W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

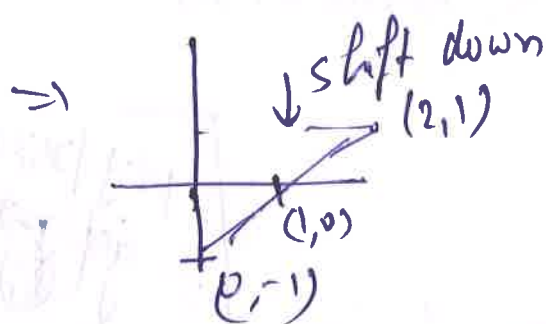
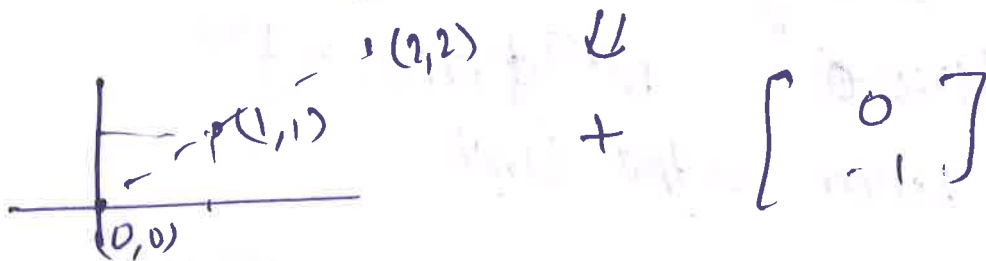
$$b = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

②

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

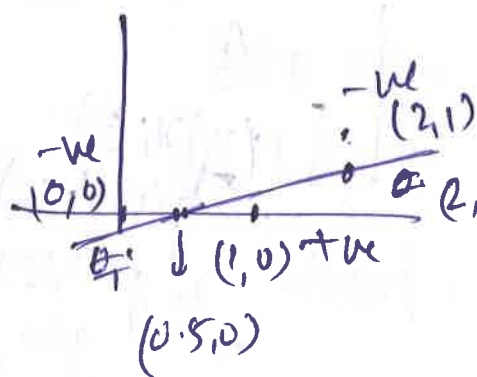
$$W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$W^T X + b = \begin{bmatrix} 1 & 0 & 1 & 2 & 2 \\ 0 & 1 & 1 & 2 & 2 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$



$$\begin{bmatrix} 0 & 1 & 1 & 2 \\ -1 & 0 & 0 & 1 \end{bmatrix}$$

$$g(W^T X + b) = \begin{bmatrix} 0 & 1 & 1 & 2 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$



linearly separable
After the transformation
Assume: $y \in \{0, 1\}$

Output Unit: $y = f^{(0)}(x; \theta) = g^{(0)} \left(W^{(0)T} \phi(x; \theta) + b^{(0)} \right)$

what is its form?

① $g^{(0)} =$ ~~identity~~ $f^{(0)}(x, \theta) =$

$W^{(0)T} \phi(x; \theta) + b^{(0)}$

Not a good idea.
For 0/1 loss

Understand the loss function:

part of input
or can be termed (6.2.4)

~~$L(y, \hat{y}; \theta)$~~

Assume $(x, y) \in \mathcal{N}(\mu, \sigma^2)$

$\phi(x; \theta) = \hat{x}(k)$

① $y \in \mathbb{R}$

$g(w) \equiv \text{Identity}$

$y | x \sim \phi(x; \theta) = w^{(0)T} \phi(x; \theta) + b^{(0)}$
Linear Output Unit

② $y \in \{0, 1\}$ Bernoulli $y \sim \text{Bernoulli}(\phi)$

$g(w) \equiv \text{identity}$

$\neq g(w) = \text{sigmoid}$

$g^{(0)}(z) = g^{(0)}(z) = \frac{1}{1 + e^{-z}}$

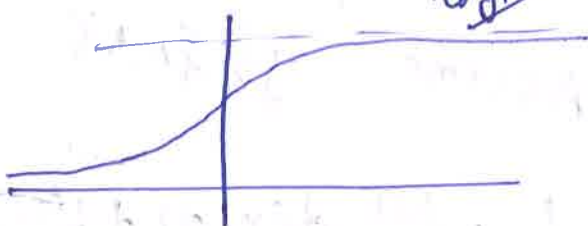
$P(y | x; \theta) = g^{(0)}(z)$

~~$g^{(0)}(w^{(0)T}$~~

$z = w^{(0)T} \hat{x}(k) + b^{(0)}$

$g^{(0)}(z) = \frac{1}{1 + e^{-z}}$
 \downarrow
logit

$P(y | x; \theta) = g^{(0)}(z)$
 $P(y=0 | x; \theta) = 1 - g^{(0)}(z) = g^{(0)}(-z)$



③ ~~$y \in \{0, 1\}$~~ $y \in \{1, \dots, L\}$

$g(w) \equiv \text{Softmax function}$

$z = w^{(0)T} \hat{x}(k) + b^{(0)}$
 \downarrow
vector

3

$$g^{(0)}(z)_l =$$

$$\frac{e^{z_l}}{\sum_{l=1}^L e^{z_l}}$$

softmax function
(generalization of logistic).

Defines a multinoulli distribution.
over $y \in \{0, 1, \dots, L\}$.

$$P(y = l | x; \theta) = g^{(0)}(z_l)$$

Now, what is the right loss function?

$$L(y, \hat{y}, \theta) = ?$$

$$= \|y - f(x, \theta)\|^2$$

squared difference

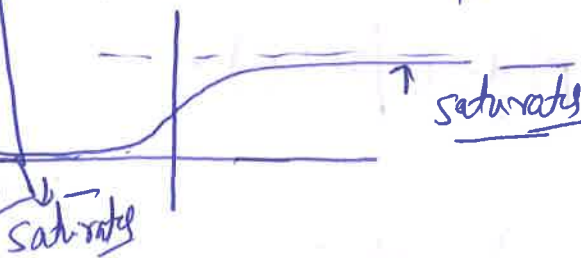
$$J(\theta) = \sum_{x \in \mathcal{X}} \|y - f(x, \theta)\|^2$$

↓
All instances

$$\hat{J}(\theta) = \sum_{\{x^{(i)}, y^{(i)}\}_{i=1}^m} \|y^{(i)} - f(x^{(i)}, \theta)\|^2$$

↓ Not a good cost function

$$\begin{aligned} & \frac{\partial}{\partial w} g(w^T \tilde{x}^{(k)} + b^{(0)}) \\ &= g'(w^T \tilde{x}^{(k)} + b^{(0)}) \\ &= (1 - g(w^T \tilde{x}^{(k)} + b^{(0)})) g(w^T \tilde{x}^{(k)} + b^{(0)}) \\ &\neq \tilde{x}^{(k)} \end{aligned}$$



$$g'(z) = g(z)(1-g(z)) \rightarrow 0$$

$$\text{as } z \rightarrow \pm \infty$$

$$z \rightarrow -\infty$$

Another idea -

$$J(\theta) =$$

$$L(y, \hat{y}; \theta) = -\log P(y|x; \theta)$$

$$z' = \frac{(1-2y)z}{(2y+1)z}$$

Bernoulli:-

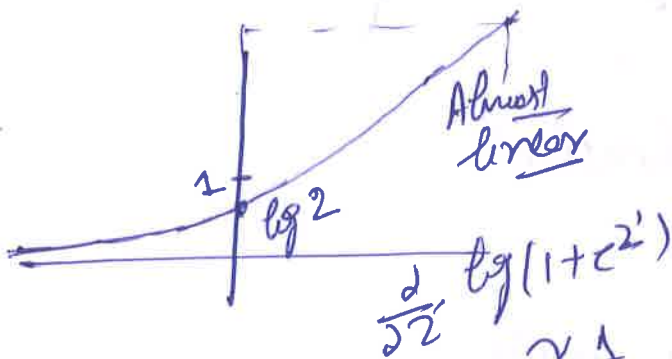
$$= -\log [g(z')] =$$

$$= \log \left[\frac{1}{g(z')} \right]$$

$$= \log (1 + e^{z'})$$

$$\downarrow \text{softmax function}$$

$$z \gg 1 \quad \log(1 + e^z) \sim z$$



$$J(\theta) = E_{(x,y) \sim \text{Dist.}} [-\log P(y|x; \theta)]$$

Does not saturate when $z' \geq 0$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m -\log P(y_i|x_i; \theta)$$

$y \neq \hat{y}$ & $z > 0$
 $y = 1$ & $z < 0$ } High gradient

Good.

Multinoulli:-

$$L(y, \hat{y}; \theta) = -\log P(y|x; \theta)$$

$$P(y=c|x; \theta) = \frac{e^{z_c}}{\sum_{c=1}^C e^{z_c}}$$

$$\log P(y=c|x; \theta) = z_c - \log \left[\sum_{c=1}^C e^{z_c} \right]$$

if c is the label with highest z_c , then

④

Note:-

if $e^{z_e^*} \geq e^{z_e}$ st $z_e^* \geq z_e$ $\forall e \neq e^*$

then $\log \sum_{e=1}^L e^{z_e} \approx z_e^* + \log \left\{ \sum_{e=1}^L \frac{e^{z_e}}{e^{z_e^*}} \right\}$

$\overline{z_e^*} = z_e^* + \log \sum_{e=1}^L e^{z_e - z_e^*}$

\Rightarrow if $e = e^*$

then $\log p(y = e | x; \theta)$ saturates

ex,

$-\log p(y = e | x; \theta) = -[z_e - z_e^*]$ Saturates near 0.

$J(\theta) = E_{(x,y) \sim \text{dist}} -\log p(y = e | x; \theta)$

$J(\theta) = \frac{1}{m} \sum_{i=1}^m -\log [p(y^{(i)} | x^{(i)}; \theta)]$

Note:-

$\text{Softmax}(z) = \text{Softmax}(z - \max_i z_i)$

Stable version

Gaussian

Finally,

$y \sim \mathcal{N}(\hat{y}, \sigma^2)$

$z = w^{(0)T} \hat{x}^{(0)} + b^{(0)} = \hat{y}$

$-\log p(y | x; \theta) = C + \frac{1}{2\sigma^2} \| \hat{y} - y \|^2$

squared error.

Summary:-

Output unit

Output:-

① Logistic / Softmax / Identity

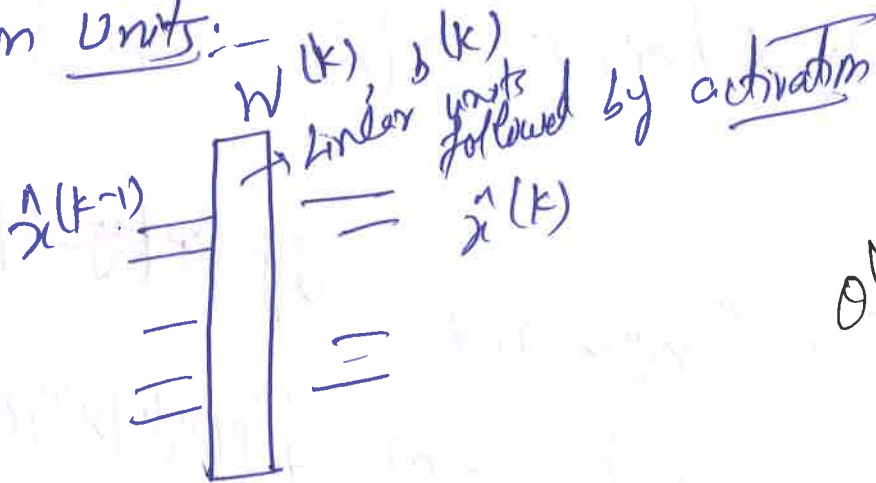
$$J(\theta) = \text{cost} = \sum_{i=1}^m -\log P(y^{(i)} | x^{(i)}; \theta)$$

~~Hidden unit~~ - Book:- Mixture of components

Mixture density networks

More advanced! Enter

Hidden units:-

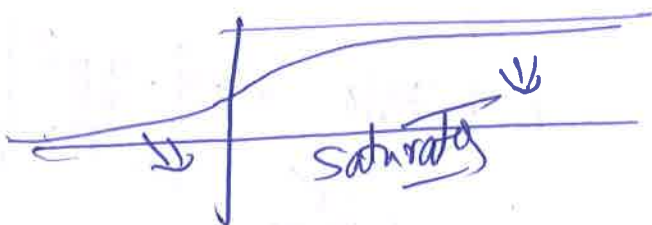


$$\theta^{(k)} = (w^{(k)}, b^{(k)})$$

$$x^{(k)} = g^{(k)} [w^{(k)T} x^{(k-1)} + b^{(k)}] = h^{(k)}_{x^{(k-1)}; \theta^{(k)}}$$

what is the right form for $g^{(k)}$?

① $g^{(k)}(z) \equiv \sigma(z) = \frac{1}{1+e^{-z}}$



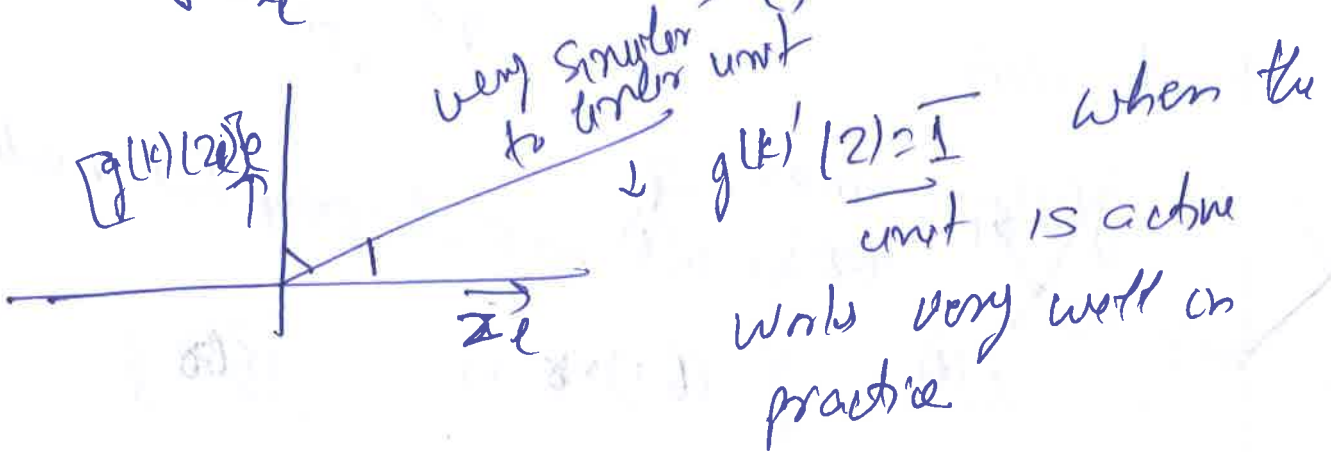
5

①

Relu:- Rectified Linear Unit:-

$$z = w^{(k)T} \hat{x}^{(k-1)} + b^{(k)}$$

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



$\frac{d}{dz} [g'(z)_k]_{z=0}$ is not defined
not an issue:-

(I). Typically may not reach exact
minima (early stopping)

(II) Use sub gradient.

(III) Initialize: - $w^{(k)}$ to be slightly
two numbers. \Rightarrow unit is active
& $w^{(k)}$ to close to zero.

Variations:-

Features
invariant to
sign
reversal

$$g(z, d)_k = \max(0, z) + d \cdot \min(0, z)$$

$$\text{when } d = -1, \quad g(z)_k = |z|$$

Absolute value rectification

④

Leaky Relu:-

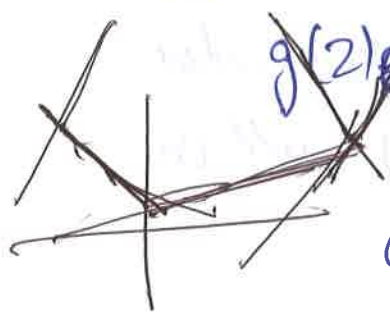
$$\boxed{\alpha \approx .01}$$

Parametric & Leaky Relu:-

α :- learnable

- Data dependent non-linearities

Maxout Units:-



$$g(z)_i = \max_j z_j$$

$\forall j \in G(i) \rightarrow$ Each group has r values

$$G(i) = \{ (i-1)r + 1, \dots, ir \}$$

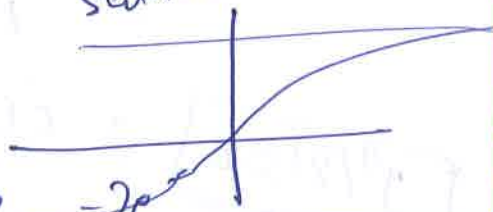
Can learn k piecewise linear (convex) function with upto k pieces.

~~$g(\tanh(z)) =$~~

$$g(z)_e = \tanh(2e) =$$

$$\frac{e^{2e} - e^{-2e}}{e^{2e} + e^{-2e}} = \tanh(2e)$$

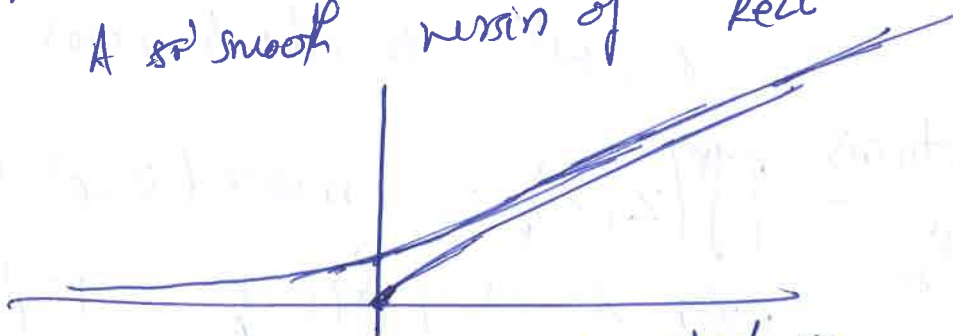
saturation issue



Softplus:-

$$g(z)_e = \log(1 + e^{2e})$$

A smooth version of ReLU



Some other units described in the book