# Lecture 2

1. Activation Functions



$$z = wx + b \quad \underline{\quad\quad} \text{ vector}$$
$$\quad\quad\uparrow\quad\llcorner\text{vector}$$
$$\quad\quad\text{matrix}$$

sigmod



act

$$\Rightarrow$$

$$g(x) = \frac{e^x}{1+e^x} \quad \text{or} \quad \frac{1}{1+e^{-x}}$$

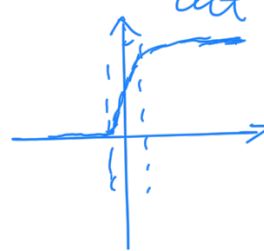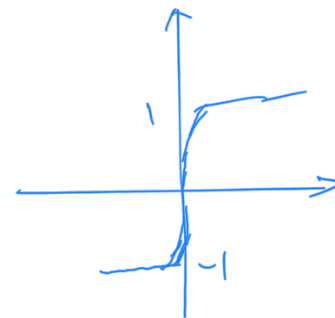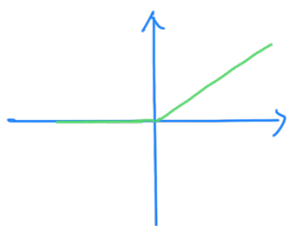tanh

$$\tanh = 2g(2x) - 1$$



$$\tanh = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

ReLU — Rectified Linear unit



$$ReLU(x) = Max(0, x)$$

$$\begin{cases} x & \text{when } x > 0 \\ 0 & \text{else} \end{cases}$$

Leaky ReLU



$$max(x, 0.1x)$$

$$max(W_1 x, W_2 x)$$



## 2. Output function

For classification

$$\begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix}$$   ground truth  $$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

softmax



output function



$$y_i = \frac{e^{z_i}}{\sum e^{z_i}}$$

$$z = \begin{bmatrix} 1 \\ 3 \end{bmatrix} \implies y = \begin{bmatrix} \frac{e^1}{e^1 + e^3 + e^7} \\ \frac{e^3}{} \end{bmatrix} = \begin{bmatrix} 0.002 \\ 0.02 \end{bmatrix}$$

$$(7) \qquad p \uparrow \quad \frac{e^2}{8} \quad (0.98)$$

## 3. Feed forward

input —□—○— output
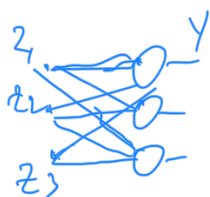
$x$ —□$\overset{z}{\downarrow}$○□○— ··· □—○— output $y$

$$z^1 = \underline{w}\, x + b$$

$$h^1 = \text{sigmoid}(z^1)$$

$$z^2 = w\, h^1 + b$$

$$h^2 = \text{sigmoid}(z^2)$$

$$\vdots$$

$w$ — matrix

$b, x$ — vector

$h^i, z^i$ — vector

$x$ —□— $\overset{z.}{\phantom{.}}$ [s]$\underset{3}{\phantom{.}}$ — $y$

$x$ is a vector of five dimensions $\qquad x = \begin{bmatrix} x \\ \vdots \\ x \end{bmatrix}$

$y$ — — — — three — —

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$\begin{matrix} x \\ 5 \end{matrix} \qquad \begin{matrix} z \\ 3 \end{matrix}$

$$z = w\, x + b$$
$$\underset{3\times1}{\phantom{z}} \quad \underset{3\times5}{\phantom{w}} \underset{5\times1}{\phantom{x}} \underset{3\times1}{\phantom{b}}$$

—□—○$\overset{h}{\downarrow}$□—○— ··· $\overset{hn}{\downarrow}$□○— $\overset{y}{\phantom{.}}$ $3\times1$

$$z = w\, x + \underline{b} \qquad h = \delta(z)$$

$$5 \times 1$$

1: $\underline{W_1}$    (10 neuron)    $50 + 10$
$$10 \times 5$$

2: $W_2$              $200 + 20$
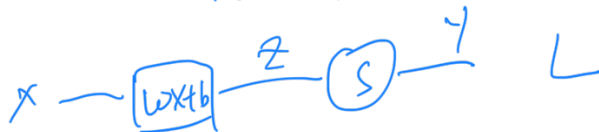$$20 \times 10$$

3: $W_3$              $60 + 3$
$$3 \times 20$$

## 4. Back propagation

$$y = a_0 x + a_1$$

$$L = (y - \hat{y})^2$$

$$a_0' = a_0^0 - \lambda \frac{\partial L}{\partial a_0} \quad \longrightarrow \text{derivative}$$

↑ ↑ step size

└ initial

$$x - \boxed{wx + b} \overset{z}{-} \bigcirc{s} \overset{y}{-} \quad L$$

$$\dot{w} = \dot{w} - \lambda \frac{\partial L}{\partial w}$$

↑ └─── step size

initial

$$x - \boxed{w_1, b} \overset{z_1}{-} \bigcirc{s} \overset{h_1}{-} \boxed{w_2, b} \overset{z_2}{-} \bigcirc{s} \overset{h_2}{-} ..$$

$$\Downarrow$$

$$[ w_1, w_2, \cdots w_n ]$$

$$\begin{cases} w_1^i = w_1^{i-1} - \lambda \frac{\partial L(w_1^{i-1})}{\partial w_1} \\ \vdots \\ w_n^i = w_n^{i-1} - \lambda \frac{\partial L(w_n^{i-1})}{\partial w_n} \end{cases}$$

$n$ layers

$$\frac{\partial L(W_n^{[i]})}{\partial W_n} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial W}$$

$$= \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_n} \frac{\partial z_n}{\partial W_n}$$

$$X - \boxed{W_1} \xrightarrow{z_1} \bigcirc \xrightarrow{h_1} \boxed{W_2} - \bigcirc \xrightarrow{h_2} \cdots \xrightarrow{z_{n-1}} \boxed{W_{n-1}} \bigcirc \xrightarrow{h_{n-1}} \boxed{W_n} \xrightarrow{z_n} \bigcirc - y \qquad L$$

no $W_n$

$$\frac{\partial L}{\partial W_n} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_n} \frac{\partial z_n}{\partial W_n}$$
$$\frac{\partial z_n = \partial(W_n h_{n-1} + b_n)}{\partial W_n} \quad \frac{}{\partial W_n}$$

$h_{n-1}$

$$= h_{n-1}$$

$$L = (y - \hat{y})^2 \qquad y =$$

$$\frac{\partial L}{\partial W_{n-1}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_n} \frac{\partial z_n}{\partial h_{n-1}} \frac{\partial h_{n-1}}{\partial z_{n-1}} \frac{\partial z_{n-1}}{\partial W_n}$$

$$W_n^{i-1} \qquad \qquad h_{n-1}$$

training data set
$$\{ (x_1, y_1), (x_2, y_2), \cdots, (x_m, y_m) \} \quad m \text{ samples.}$$

~~it~~ initialization
n layers
$$\{ w_1, w_2, \cdots w_n \}$$

it = 1

to get $\{W_1^1, W_2^1, \cdots W_n^1\}$

$$W_n^1 = W_n^0 - \lambda \frac{\partial L}{\partial W_n}(W_n^0)$$

$$\frac{\partial L}{\partial W_n} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_n} \frac{\partial z_n}{\partial W_n}$$

$$W_{n-1}^1 = W_{n-1}^0 - \lambda \frac{\partial L}{\partial W_{n-1}}(W_1^0, W_2^0, \cdots W_n^0)$$

$$\vdots$$

$$W_1^1 = W_1^0 - \lambda \frac{\partial L}{\partial W_1}(W_1^0, \cdots W_n^0)$$

## 5. Derivatives of Actuation Functions

Sigmoid $\quad y = \dfrac{1}{1+e^{-z}} \quad$ or $\quad \dfrac{e^z}{1+e^z}$

$$\frac{\partial y}{\partial z} = \frac{-1(-e^{-z})}{(1+e^{-z})^2}$$

$$= \underbrace{\frac{1}{1+e^{-z}}}_{y} \cdot \underbrace{\frac{e^{-z}}{1+e^{-z}}}_{(1-y)}$$

$$= y(1-y)$$

Softmax $\quad y_i = \dfrac{e^{z_i}}{\sum_i e^{z_i}} = \dfrac{e^{z_i}}{c + e^{z_i}}$

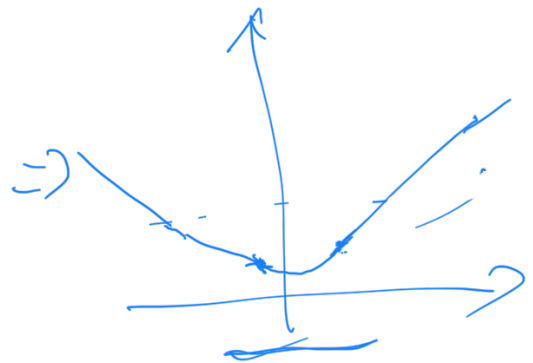$$\frac{\partial Y_i}{\partial Z_i} = Y_i(1-Y_i) \qquad \frac{\uparrow}{\sum_{k \in i}} e^{z_k}$$

6. Loss functions

For regression
$$L = \frac{1}{2}(Y - \hat{y})^2$$
$$= (Y - \hat{y})^T (Y - \hat{y})$$
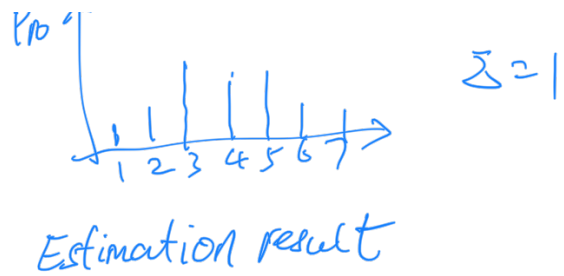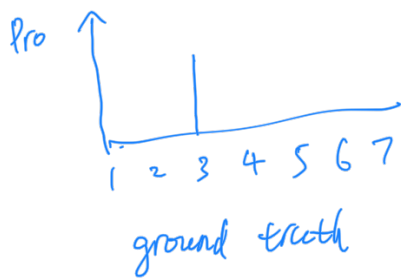$$\frac{\partial L}{\partial Y} = -(Y - \hat{y})$$

Huber
$$L = \begin{cases} \frac{1}{2}(Y - \hat{y})^2 & \text{if } |Y - \hat{y}| < c \\ c(|Y - \hat{y}| - \frac{1}{2}c) \end{cases}$$

Pseudo Huber
$$L = c^2 \left( \sqrt{1 + \left(\frac{Y - \hat{y}}{c}\right)^2} - 1 \right)$$

cross entropy Loss

Pro

| 1 2 3 4 5 6 7

ground truth

Pro'

1 2 3 4 5 6 7

Estimation result

$\sum = 1$

$$L = -\sum_i t_i \log(y_i)$$

target/label      estimate

ground truth $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$ , estimated $\begin{bmatrix} 0.002 \\ 0.02 \\ 0.98 \end{bmatrix}$

$$L = -\left( 0 + 0 + \log(0.98) \right) = 0.02$$

log $\longrightarrow$ natural log

if estimated is $\begin{bmatrix} 0.33 \\ 0.33 \\ 0.33 \end{bmatrix}$

$$L = -\left( 0 + 0 + \log(0.33) \right) = 1.1$$