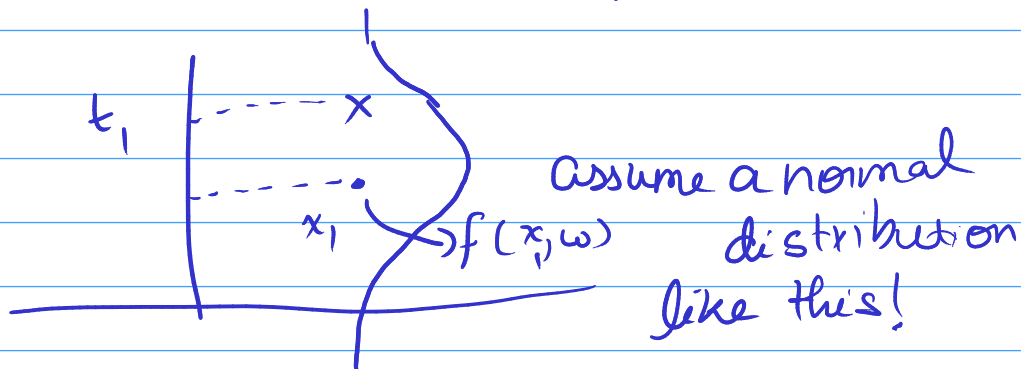
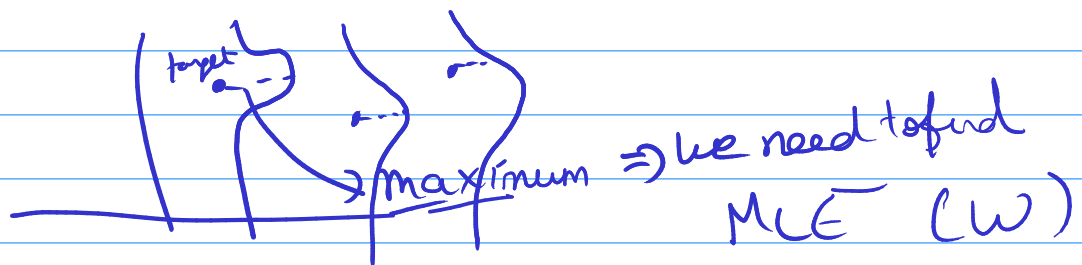


$$y/x \sim \mathcal{N}(y | f(x, w), \Sigma)$$



Consider covariance matrix to be  $I$  for now  
 now  $f(x, w)$  is quite away from target

now lets have another  $w_2$  such that



Covariance matrix can also be determined  
 positive, semidefinite  
 you need a more complex neural network

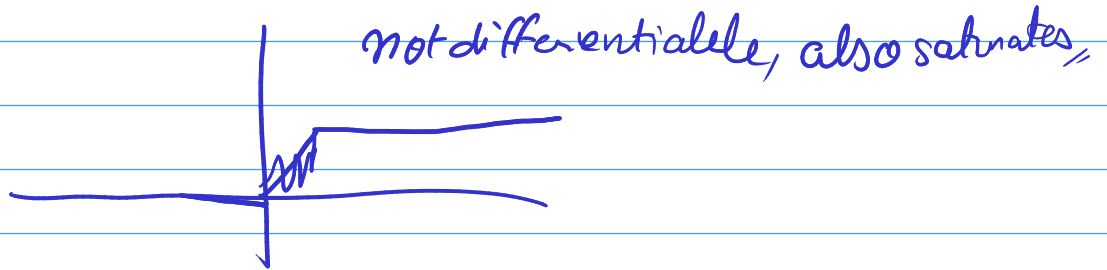
What cost functions are good??

$$\rightarrow P(y=1) = \sigma(w_0 + w_1 t_1 + w_2 t_2 \dots + w_n t_n)$$

what if

$$p(y=1) = \max\{0, \min\{1, w^T h + b\}\}$$

Gradients are zero outside



we need a better cost function !!

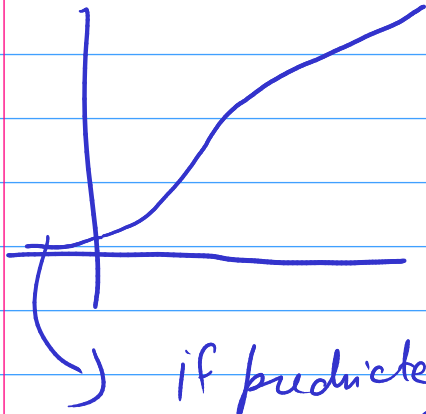
$$-\log(p(\dots)) \dots$$

So we need a cost function that saturates  
at correct classification  
& doesn't saturate @ incorrect  
classifications

$y \setminus x$	1	0
$p(y x)$	$p$	$1-p$

$p$  is a bernoulli

Cost functions aren't forced, it's a consequence  
of an assumption  
of a normal  
distribution



if predicted output is low  
(-ve) } saturates  
but target is positive

—  $(Z)$ , if there is small change  $Z \approx 0$   
then function won't change

Ex.  $Z \rightarrow \infty$  or  $Z \rightarrow -\infty$

and in an ~~iter~~ interval, it changes a lot!

My error function should change when I'm not  
near a global or local minima

in SGD, since you're doing it example by  
example

you want to improve vs drastically

$\Rightarrow$  Function shouldn't  
saturate

# Deriving Cost

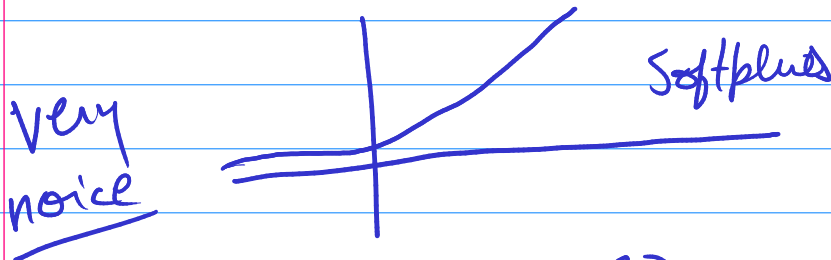
$$\sigma(z) = \frac{1}{1+e^{-z}} = \frac{e^z}{1+e^z}$$

$$\frac{d\sigma(z)}{dz} = \sigma(z)(1-\sigma(z))$$

assume  $\leftarrow \boxed{\ell(z) = \log(1+e^z)}$   $\rightarrow$  softplus function

$z \rightarrow +\infty \quad \log(1+e^z) \rightarrow z$   
 $z \rightarrow -\infty \quad \log(1+e^{-z}) \rightarrow 0$

$x^+ = \max\{0, x\}$   
 Softer version of max function

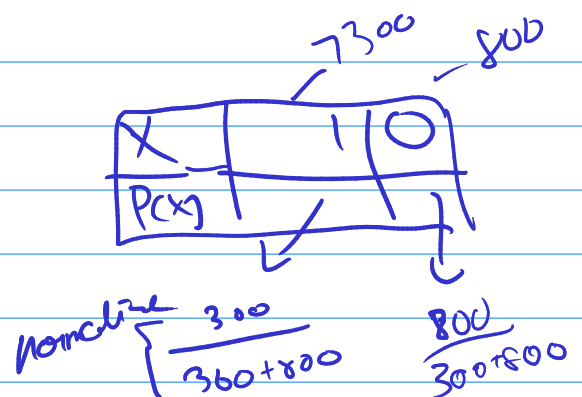
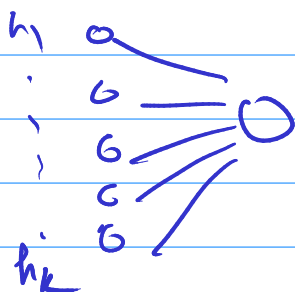


$$\ell(1+e^{-z}) = \log\left(\frac{e^z+1}{e^z}\right) = \log\left(\frac{1}{\sigma(z)}\right)$$

$$= -\log(\sigma(z))$$

$$\log(\sigma(z)) = -\ell(-z)$$

$$z = w_0 + w_1 h_1 + \dots + w_k h_k$$



I want to combine  $z$  with target to arrive @ cost function that doesn't saturate

So they choose  $yz$  as that combination such that

$$\log(\hat{p}(y)) = yz = \hat{p}(y) = e^{yz}$$

$$P(y) = \frac{e^{yz}}{\sum_{y'} e^{y'z}} = \frac{e^{yz}}{1 + e^z} \quad (y \text{ takes values } 0 \text{ or } 1)$$

A "kind of" probability distribution

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

$$\boxed{P(y) = \sigma((2y-1)z)} \quad \left[ \begin{array}{l} \text{when } y=1 \\ 2y-1=1 \\ y=0, \sigma(-z) \\ = \frac{1}{1+e^{-(-z)}} = \frac{1}{1+e^z} \end{array} \right]$$

Cost function of regular logistic regression also saturates

$$P(y_1, \dots, y_m | x_1, \dots, x_m) = \prod_{i=1}^m P(y_i | x_i)$$

$$\max \log \left( \prod_{i=1}^m p(y_i | x_i) \right)$$

$$\max \sum_{i=1}^m \log(p(y_i | x_i))$$

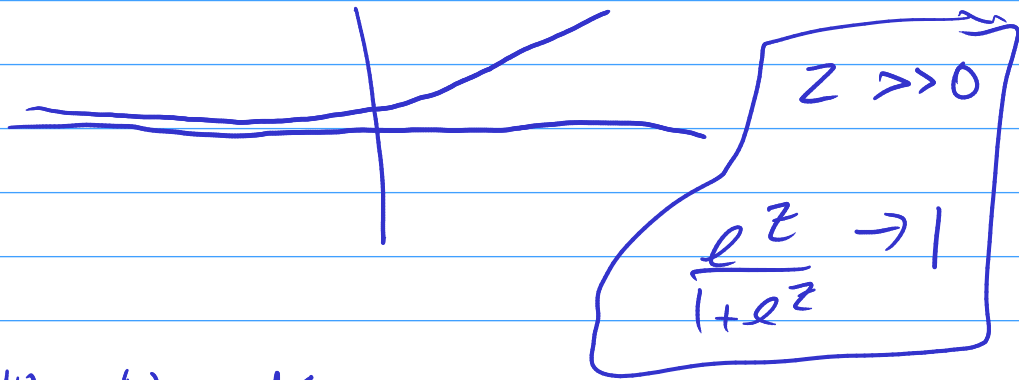
$$\min \sum_{i=1}^m -\log(p(y_i | x_i))$$

$$\text{Cost func: } \ell(y) = -\log p(y|x)$$

$$= -\log(\sigma((2y-1)z)) = \ell(-(2y-1)z)$$

$$= e^{((1-2y)z)}$$

remember  $e$  is the softplus function



$$w^{(k+1)} = w^{(k)} - \eta \underbrace{\frac{d\mathcal{L}}{dw}}$$

gradient should be small,  $z \gg 0$ , & prediction is 1 & is correct

This  $e$  function achieves that as, if  $y = 1$

$$e^{((1-2y)z)}$$

$$e^{-z} \approx 0 \text{ since } z \gg 0$$

$z \gg 0$ ,  $y = 0$ , you want high gradients

then  $e(z) \rightarrow z$   
you get

high  
for

$$z \ll 0 \quad y = 0$$

then you get

$$e(z) \rightarrow 0$$

$$z \ll 0 \quad y = 1$$

$$e(-z) \rightarrow (-z)$$

Why the  $\sigma$  function  
WTF is wrong with what I have?!

$$P(y_i=1/x_i) = \sigma(z_i)$$

$$P(y_i=0/x_i) = 1 - \sigma(z_i)$$

$$P(y_i) = (\sigma(z_i))^{y_i} (1 - \sigma(z_i))^{1-y_i}$$

$$P(y_1, y_2, \dots, y_m | x_1, x_2, \dots, x_m) = \prod_{i=1}^m P(y_i/x_i)$$

$$= \prod_{i=1}^m (\sigma(z_i))^{y_i} (1 - \sigma(z_i))^{1-y_i}$$

948

$$= \log(P(y_1, y_2, \dots, y_m | x_1, x_2, \dots, x_m))$$

$$\max \sum_{i=1}^m \left[ y_i \log(\sigma(z_i)) + (1-y_i) \log(1 - \sigma(z_i)) \right]$$

$$= \min \sum_{i=1}^m - \left( y_i \log \sigma(z_i) + (1-y_i) \log(1 - \sigma(z_i)) \right)$$

$z \gg 0 \quad y=1 \rightarrow \text{saturates}$

$z \gg 0 \quad y=0 \quad \text{also saturates ; - ;}$

$\sigma(z)$  saturates

$\Rightarrow \log(1 - \sigma(z))$  saturates (ultimately,

when  $\sigma(z)$   
doesn't change  
much at all)

$\Rightarrow$  problem,

Up next: softmax, Orelu

& how to deal with them!

Alright, let's say  $m$  outputs are there  
then how do we predict

$(w^T h + b)$  is input

$$(w^T h + b) = z_i = \log \hat{p}(y=i/x)$$

unnormalized  
log probabilities

$$\Rightarrow e^{w^T h + b} = \hat{p}(y=i/x)$$

$$p(y=i/x) = \frac{e^{z_i}}{\sum_{i=1}^m e^{z_i}} \quad \text{for } m \text{ classes}$$

Why distribution? This is a multinomial dist.



→ so we can "turn"  $z$  into a probability this way.

→ we don't normally maximize the probabilities  
so we ~~not~~ maximize the (unnormalized) log probabilities  
then we normalize



Softmax = activation

$$\sum_{i=1}^n \frac{e^{z_i}}{e^{z_i}}$$

1) one hot encoding

$$\min_{\theta} -\log \prod_{i=1}^n (p(y=2|x_i))^{\theta} (p(y=1|x_i))^{1-\theta}$$

It is however not for sure that  
you'll hit local minima after  
some iterations

✓ 0, since  
~~the~~ this  
is class  
✓ 1, since this  
is class 2

we might need early stopping

## Other activation fncs

Relu: differentiability problem

@ '0' take either LHS, or RHS derivative

## Leaky ReLU

Leaky ReLU

$$h_i = g(z_i) = \max(0, z_i) + \alpha_i \min(0, z_i)$$

$\alpha = -1 \Rightarrow$  absolute value activation function

$$\begin{array}{r} 0 \\ \hline 000 \\ \hline 000 \\ \hline 000 \end{array}$$

take max of each  $w^T x + b$  in each division  
and all neurons in that division  
have max of division in their (maxout)  
output

This is maxout

Don't use sigmoid,  $\tanh(z)$  saturates lesser, use that instead  
(make sure you shift origin of the graph)

radial basis

softplus  $\log(1 + e^z)$

hard tanh  $\max(-1, \min(1, a))$