

Day 4:

MULTI-CLASS CROSS ENTROPY

$i \rightarrow$		$j \rightarrow$		
		Door 1	Door 2	Door 3
$i \downarrow$	1 Duck	0.7 P_{11}	0.3 P_{12}	0.1 P_{13}
	2 Bear	0.2 P_{21}	0.4 P_{22}	0.5 P_{23}
	3 walrus	0.1 P_{31}	0.3 P_{32}	0.4 P_{33}

Support assume.

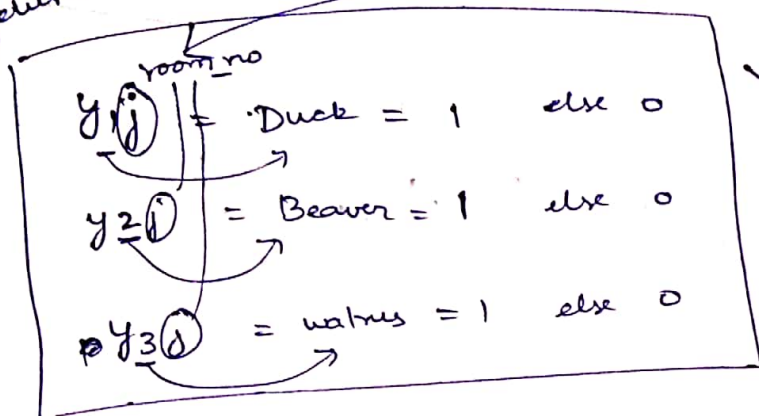
Door 1	Door 2	Door 3
Duck	walrus	walrus
0.7	0.3	0.4

$$P = 0.7 \times 0.3 \times 0.4 = 0.084$$

$$CE = -\ln(0.7) - \ln(0.3) - \ln(0.4)$$

$$CE = 2.48$$

Conclusion:



$$\text{Cross Entropy} = - \sum_{i=1}^n \sum_{j=1}^m P_{y_{ij}} \log(P_{ij})$$

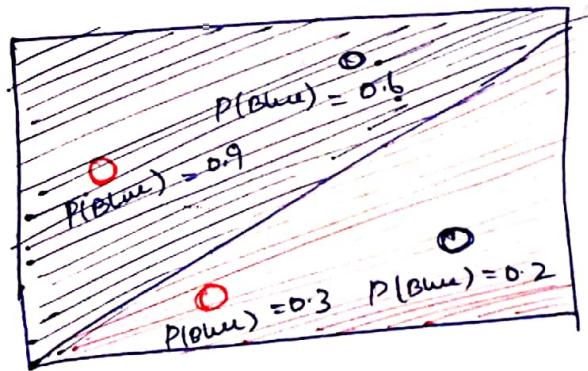
adding the log of prob of events that actually have occurred

[m → no of classes animals]

n → door no.

LOGISTIC REGRESSION

Calculating the Error function



If $y=1$ (our intention is to minimize the error function)

$$\begin{aligned} P(\text{Blue}) &= \hat{y} \\ \text{Error} &= -\log(\hat{y}) \\ P(\text{red}) &= 1 - P(\text{Blue}) \\ &= 1 - \hat{y} \\ &= \log(\hat{y}) \\ \text{Error} &= -\log(1 - \hat{y}) \end{aligned}$$

Combining them to produce a generalized one

$$\text{Error} = -(1-y) \log(1-\hat{y}) + y \log(\hat{y})$$

If $y=0$
If $y=1$

$$\text{Error} = -\frac{1}{m} \sum_{i=1}^m (1-y_i) \log(1-\hat{y}_i) + y_i \log(\hat{y}_i)$$

↑
generalized for all data points

$$\text{If } \hat{y}_i \Rightarrow \sigma(w x + b)$$

$$\text{Error} = -\frac{1}{m} \sum_{i=1}^m (1-y_i) \log(1-\sigma(w x_i + b)) + y_i \log(\sigma(w x_i + b))$$

Binary ↗

If we have multi class prob we have to go for multi-class Entropy

$$\Rightarrow -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n y_{ij} \log(\hat{y}_{ij})$$

After minimizing the error function we are proceeding with Gradient Descent with smaller heights making it as w', b'
 new weights
 $\sigma(w x_i' + b')$

Gradient Descent.

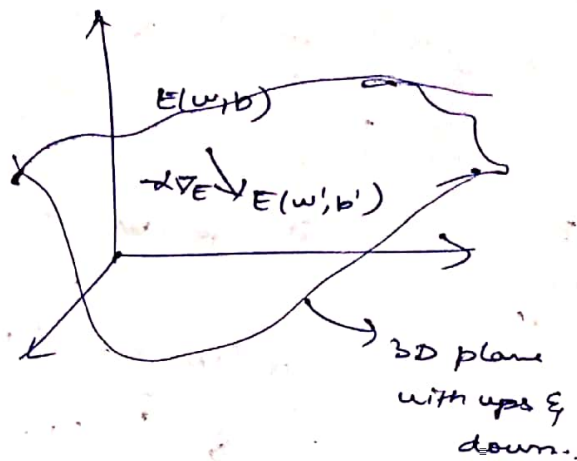
$$\hat{y} = \sigma(wx + b) \leftarrow \text{bad}$$

$$\hat{y} = \sigma(w_1 x_1 + \dots + w_n x_n + b)$$

$$\nabla E = \left(\frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n}, \frac{\partial E}{\partial b} \right)$$

$$\alpha = 0.1$$

$$\begin{aligned} w_i' &\leftarrow w_i - \alpha \left(\frac{\partial E}{\partial w_i} \right) \\ b_i' &\leftarrow b - \alpha \left(\frac{\partial E}{\partial b} \right) \\ \hat{y} &\leftarrow \sigma(w'x + b') \end{aligned}$$



Gradient Calc for sigmoid function:

$$\sigma'(x) = \frac{\partial}{\partial x} \frac{1}{1 + e^{-x}}$$

$$= \frac{e^{-x}}{(1 + e^{-x})^2}$$

$$= \frac{1}{1 + e^{-x}} \times \frac{e^{-x}}{(1 + e^{-x})}$$

$$\sigma'(x) = \sigma(x) * (1 - \sigma(x))$$

If we have m points for binary classification.

$$E = \frac{1}{m} \sum_{i=1}^m \left(y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right)$$

$$\left[\begin{array}{l} \text{WKT:} \\ \hat{y}_i = w x_i^{(i)} + b \end{array} \right]$$

Our goal is to calc

Error at a point $x_i [x_1, x_2, \dots, x_n]$.

$$\nabla E = \left(\frac{\partial}{\partial w_1} E, \dots, \frac{\partial}{\partial w_n} E, \frac{\partial}{\partial b} E \right)$$

• Error each pt produces
& calculate derivative of the error
& average of error at all pts

$$E = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

In order to calculate the derivative of this error with respect to weights, we'll first calculate $\frac{\partial}{\partial w_j} \hat{y}$.

$$\hat{y} = \sigma(w \cdot x + b)$$

$$\frac{\partial}{\partial w_j} \hat{y} = \frac{\partial}{\partial w_j} \sigma(w \cdot x + b)$$

$$= \frac{\sigma(w \cdot x + b) (1 - \sigma(w \cdot x + b))}{\sigma(x) (1 - x)} \times \frac{\partial}{\partial w_j} (w \cdot x + b)$$

format

$$\Rightarrow \hat{y} (1 - \hat{y}) \times \frac{\partial}{\partial w_j} (w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b)$$

constant

$$\Rightarrow \hat{y} (1 - \hat{y}) \times x_j$$

now we can calc the error derivative:

$$\frac{\partial}{\partial w_j} E = \frac{\partial}{\partial w_j} [-y \log(\hat{y}) - (1-y) \log(1-\hat{y})]$$

$$= -y \cdot \frac{\partial}{\partial w_j} \log(\hat{y}) - (1-y) \frac{\partial}{\partial w_j} \log(1-\hat{y})$$

$$= -y \cdot \frac{1}{\hat{y}} \left[\frac{\partial}{\partial w_j} (\hat{y}) \right] - (1-y) \cdot \frac{-1}{(1-\hat{y})} \left[\frac{\partial}{\partial w_j} (1-\hat{y}) \right]$$

$$\Rightarrow -y \cdot \frac{1}{\hat{y}} [\cancel{\hat{y}} (1-\hat{y}) \cdot x_j] - (1-y) \cdot \frac{1}{(1-\hat{y})} (-1) \hat{y} (1-\hat{y}) x_j$$

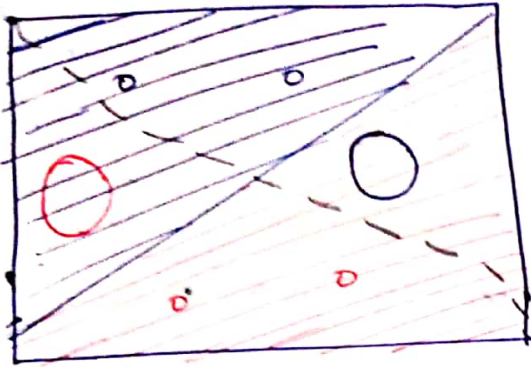
$$= -y [1-\hat{y}] \cdot x_j + (1-y) \hat{y} x_j$$

$$\Rightarrow \cancel{-y x_j + y \hat{y} x_j} + \hat{y} x_j - y \hat{y} x_j$$

$$\Rightarrow -(y - \hat{y}) x_j$$

$$\frac{\partial}{\partial b} E = -(y - \hat{y})$$

Gradient Descent :



1. Start with random weights.

$$w_1, \dots, w_n, b$$

2. For every point (x_i, \dots, x_n) :

2.1 for $i \leftarrow 1$ to n :

$$2.1.1 \text{ update } w_i' \leftarrow w_i - \alpha \left(\frac{\partial E}{\partial w_i} \right)$$

this is what we done before

$$w_i \leftarrow w_i - \alpha (\hat{y} - y) x_i$$

$$2.1.2 \text{ update } b' \leftarrow b - \alpha \left(\frac{\partial E}{\partial b} \right)$$

$$\leftarrow b - \alpha (\hat{y} - y)$$

3. repeat until error is small.

Perceptron vs Gradient Descent;

If a point is misclassified.

$$w_i \rightarrow \begin{cases} w_i + \alpha x_i & \hat{y} \neq y \\ w_i - \alpha x_i & \hat{y} = y \end{cases}$$

$$w_i \leftarrow w_i + \alpha (y - \hat{y}_i) x_i$$

for all points.

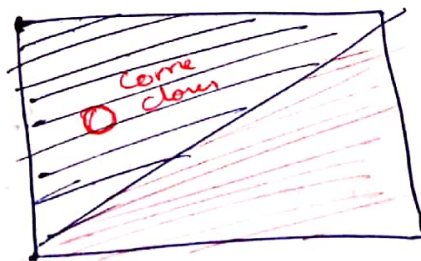
are those same?

Comparative study

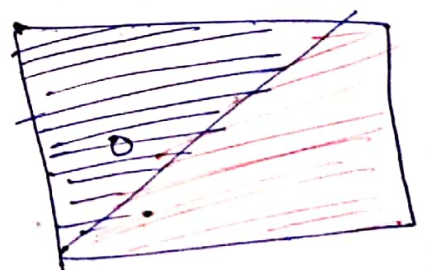
labels: 1 & 0



Correctly classified $y - \hat{y} = 0 \Rightarrow \sqrt{y = \hat{y}}$
 misclassified $y - \hat{y} = 1$ if true
 labeled as blue $y - \hat{y} = -1$ if false



(a) MISCLASSIFIED



(b) correctly classified

[well if it was a case in perceptron algo it does nothing]

BUT

IN GD change weights the point is telling the line to move farther away

HYPOTHESIS:

(a) \Rightarrow come closer
 (b) \Rightarrow move farther away

CONCLUSION

Atlast the line listens to all pts.
 It takes steps to a good soln!

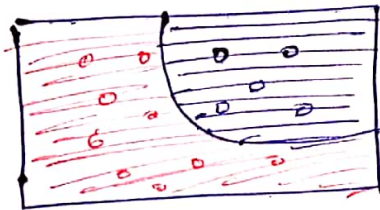
CONTINUOUS PERCEPTION!

Gradient Descent:

correct classified points	\Rightarrow move farther away	\Rightarrow more likely to be in my <u>area</u>
misclassified	\Rightarrow come closer	\Rightarrow So that I can also be in my <u>area</u>

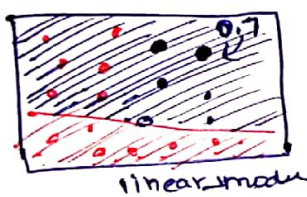
It's all happening with straightening line.

But it's be more comfort when ^{mostly} dealing with curve ~~expecting~~ the expected needs.

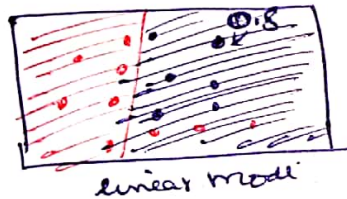


non-linear
arise of Neural network

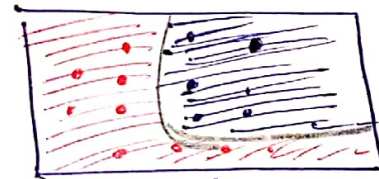
Neural Network:



+



=



add the 1st blue point $\rightarrow 0.7 + 0.8 = 1.5$

$0 < x < 1$

Squish $\rightarrow x \rightarrow$ sigmoid fn

\downarrow
0.82

$$\sigma \left(\begin{array}{l} \text{Probability of} \\ \text{Point} \\ \text{1st linear} \\ \text{mod} \end{array} + \begin{array}{l} \text{prob of point} \\ \text{in 2nd} \\ \text{linear model} \end{array} \right)$$

weight the sum

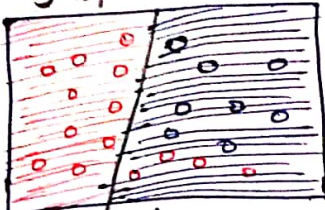
I need to have 7 times of model 1
than 5 times of model 2

$$\Rightarrow \text{so, do this } 7 \times \left(\frac{0.7}{\text{model 1}} \right) + 5 \times (0.8) = \underline{\underline{6 \text{ bias}}}$$

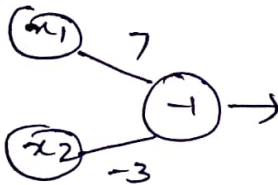
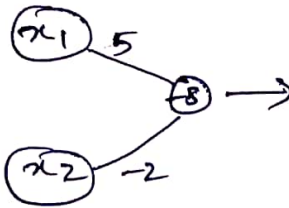
$$\Rightarrow \sigma(2.9) = \underline{\underline{0.95}}$$

Neural Network

$$5x_1 - 2x_2 \rightarrow 8$$



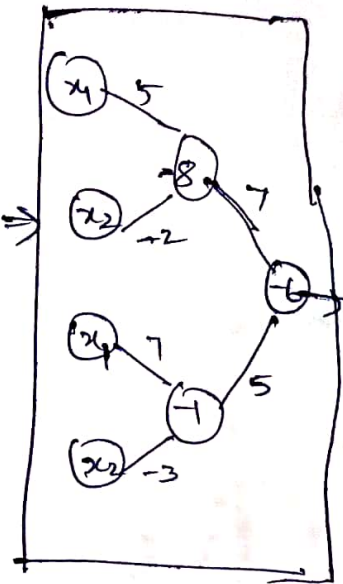
$$7x_1 - 3x_2 - 1$$



二



也



↓ make
Compass

\Downarrow whenever
 $7x_1 + 5x_2 - 6$ we combine
we get non-linear region

