

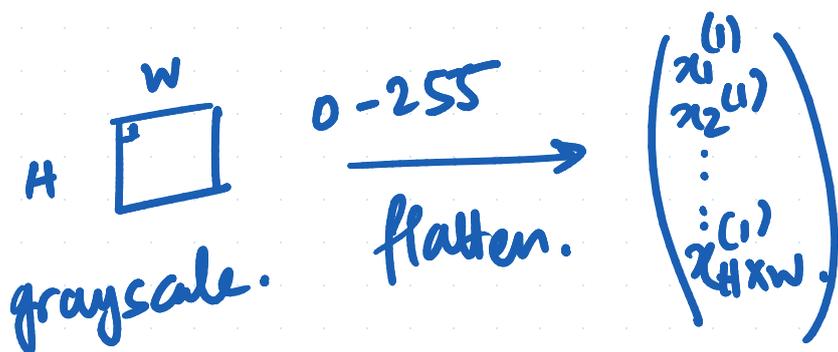
Lecture - 4

vectors \rightarrow capital letters V, X, Y

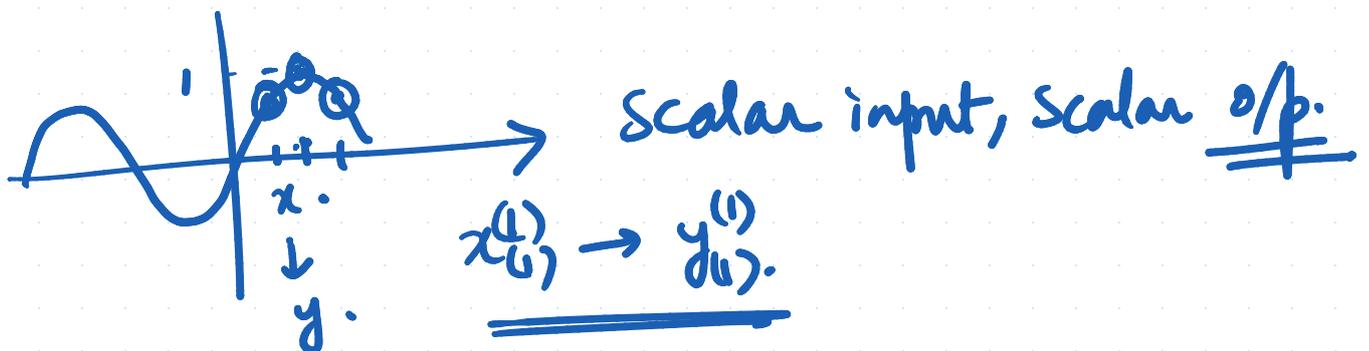
Scalars \rightarrow small letters v, x, y .

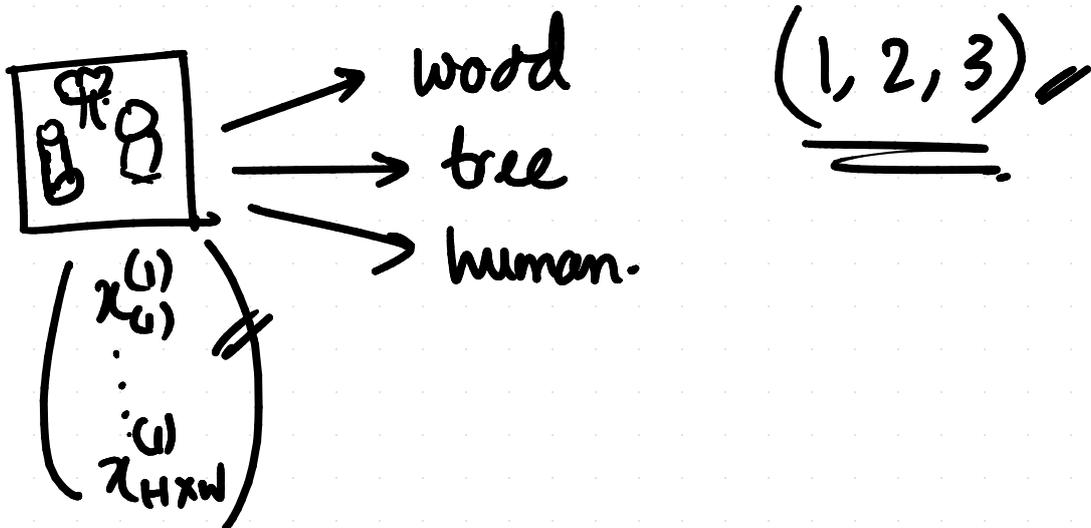
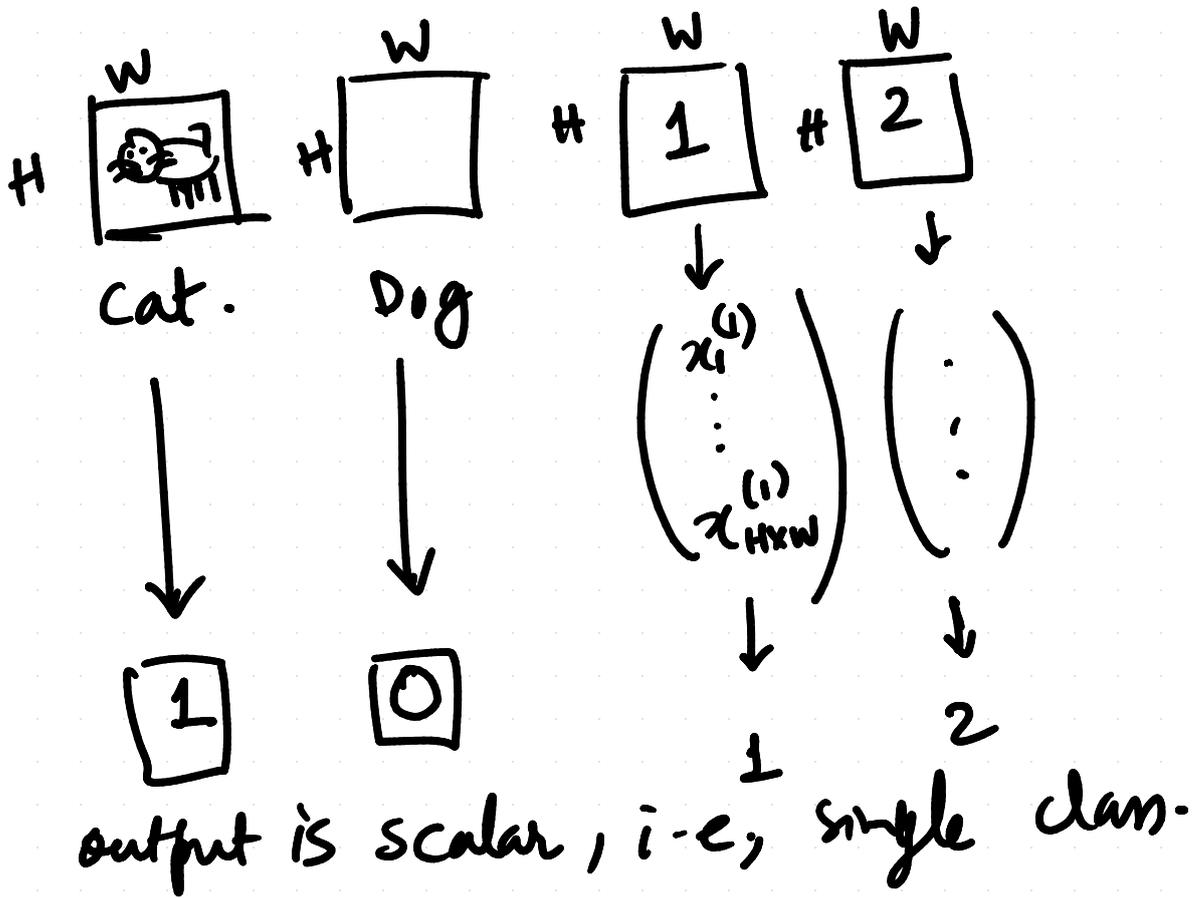
$$x_1^{(1)}, \begin{pmatrix} x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \\ \vdots^{(1)} \\ x_n^{(1)} \end{pmatrix} \rightarrow y_1^{(1)}, \begin{pmatrix} y_1^{(1)} \\ y_2^{(1)} \\ y_3^{(1)} \\ \vdots^{(1)} \\ y_s^{(1)} \end{pmatrix}$$

$$\mathcal{D} = \{x, y\}_{i=1}^n$$



input using vectors.





$x^i \rightarrow$ i th input vector.

$o^i \rightarrow$ i th output scalar

$net^i \rightarrow W \cdot x^i$ (n -input/output).

$W \rightarrow$ weight vector. Observation:

W & x^i has $m+1$ components.

$$W = \langle w_m, w_{m-1}, \dots, w_2, w_1, w_0 \rangle$$

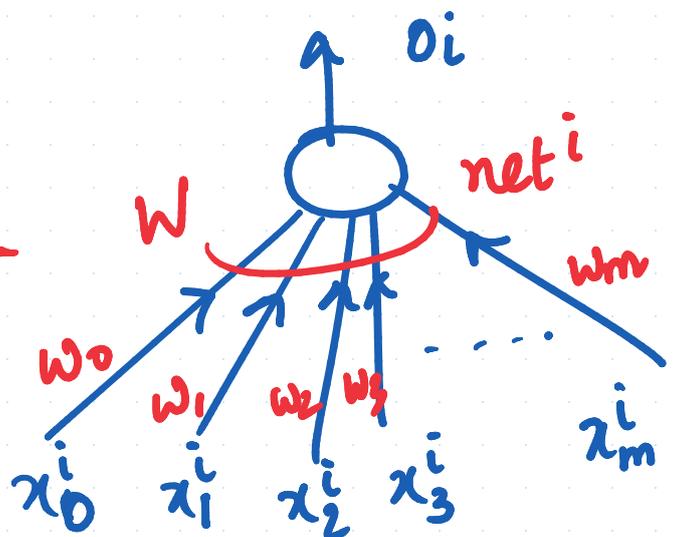
$$x^i = \langle x_m^i, x_{m-1}^i, \dots, x_2^i, x_0^i \rangle.$$

$$net_i = net^i$$

Derivative of Sigmoid.

$$o^i = \frac{1}{1 + e^{-net^i}}$$

$$net^i = \sum_j x_j^i w_j$$



$$o^i = \frac{1}{1 + e^{-net^i}} \rightarrow \text{for } i\text{th input.}$$

$$\ln o^i = - \ln (1 + e^{-net^i}) \quad \boxed{\ln = \log_e}$$

$\frac{\partial}{\partial net^i}$ on both sides

$$\frac{1}{o^i} \frac{\partial o^i}{\partial net^i} = - \frac{1}{(1 + e^{-net^i})} (-e^{-net^i})$$

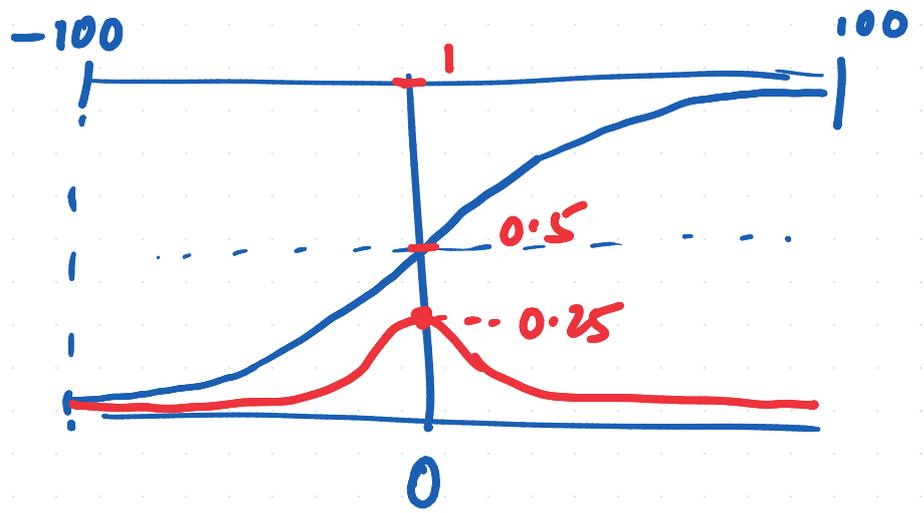
$$= \frac{e^{-net^i}}{1 + e^{-net^i}} = (1 - o^i)$$

$$\frac{1}{o^i} \frac{\partial o^i}{\partial net^i} = (1 - o^i)$$

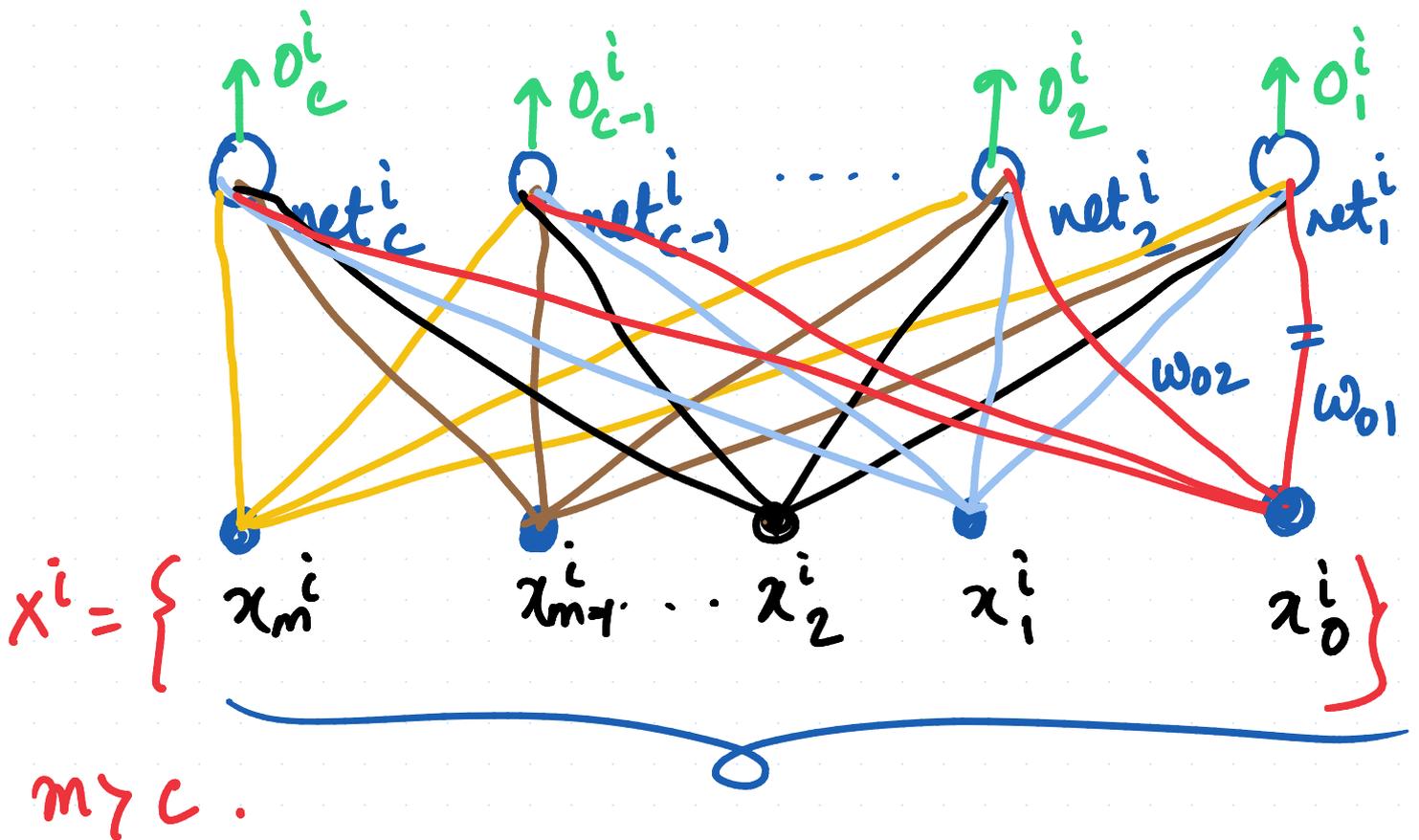
$$\boxed{\frac{\partial o^i}{\partial net^i} = o^i (1 - o^i)}$$

Derivative
of sigmoid

Maximum value of derivative of sigmoid = 0.25



Softmax using FCNN, DNN, Layer Perse



MNIST \rightarrow 28×28 \rightarrow 784 m

classes \rightarrow (0-10) \rightarrow 10 c

$m \gg c$

Output for class c , $c: 1$ to C .

$$o_c^i = S(\text{NET}^i)_c = \frac{e^{\text{net}^i_c}}{\sum_{k=1}^C e^{\text{net}^i_k}}$$

②

$$\left\{ \begin{array}{l} [x_0^0, x_1^0, \dots, x_m^0] \rightarrow [o_0^0, o_1^0, \dots, o_c^0] \\ [x_0^1, x_1^1, \dots, x_m^1] \rightarrow [o_0^1, o_1^1, \dots, o_c^1] \\ \vdots \\ [x_0^N, x_1^N, \dots, x_m^N] \rightarrow [o_0^N, o_1^N, \dots, o_c^N] \end{array} \right.$$

One-hot vectors / encoding

Cat Dog ... Elephant

↓
1

↓
2

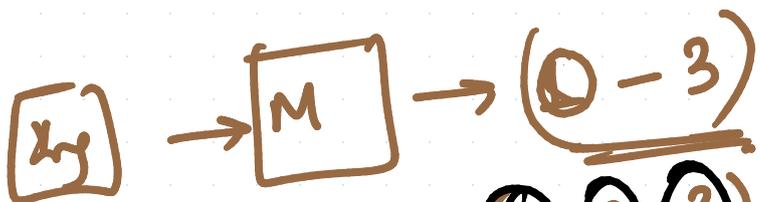
↓
3

Softmax.

→ (1, 0, 0) → Cat

→ (0, 1, 0) → Dog

→ (0, 0, 1) → Elephant



$\|y - \hat{y}\|_2^2$

(1, 2, 3)

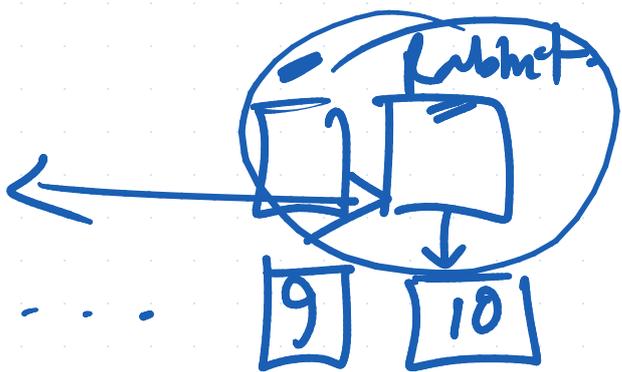
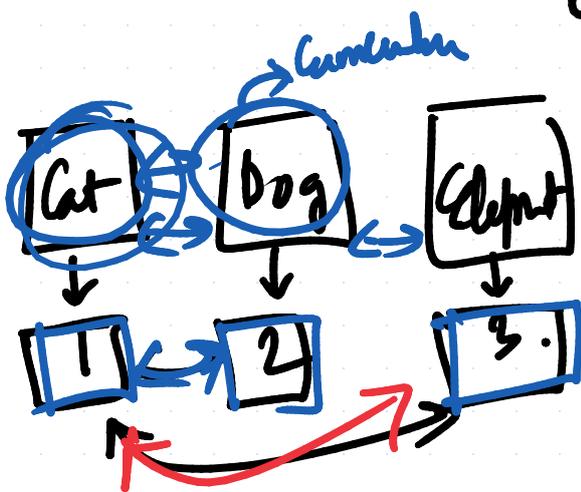
Cat

dog

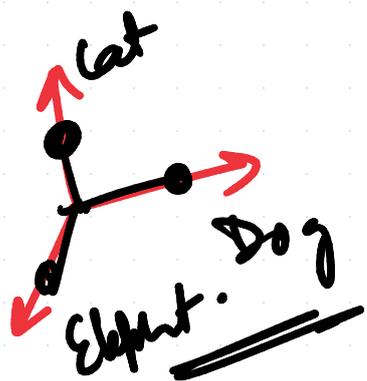
elephant.

[M]

(0.4, 0.3, 0.6)



How to compare kilometers with Kgs.



Softmax

$$\begin{Bmatrix} -1 \\ 1 \\ 0 \end{Bmatrix} \rightarrow \begin{Bmatrix} 0.02 \\ 0.5 \\ 0.43 \end{Bmatrix} \text{ argmax } \begin{Bmatrix} 0 \\ 1 \\ 0 \end{Bmatrix}$$

values are between 0 & 1
& sum of all the values = 1.

Complex

$$o_c^i = \frac{e^{\text{net}_c^i}}{\sum_{k=1}^c e^{\text{net}_k^i}}, \text{ i-th ip pattern.}$$

$$\Rightarrow \ln o_c^i = \underbrace{\log_{\frac{1}{e}}}_{\text{log}} \left(\frac{e^{\text{net}_c^i}}{\sum_{k=1}^c e^{\text{net}_k^i}} \right)$$

Derivative:- w.r.t. net^i_c we get:-

$$\frac{1}{o^i_c} \frac{\partial o^i_c}{\partial net^i_c} = 1 - \frac{e^{net^i_c}}{\sum_{k=1}^c e^{net^i_k}}$$

Case 1:- when $c = 0$ & net are the same

$$\frac{1}{o^i_c} \frac{\partial o^i_c}{\partial net^i_c} = 1 - o^i_c.$$

$$\frac{\partial o^i_c}{\partial net^i_c} = o^i_c (1 - o^i_c)$$

Case 2:- when c' and net^i_c are different from $c = 0$

$$\frac{1}{o^i_c} \frac{\partial o^i_c}{\partial net^i_c} = 0 - \frac{e^{net^i_{c'}}}{\sum_{k=1}^c e^{net^i_k}} = -o^i_{c'}$$

$$\frac{\partial o_c^i}{\partial \text{net}_c^i} = -o_c^i o_c^i$$

Maximum Likelihood & Cross Entropy loss-

Q) How did cross entropy loss come to existence?

Input \rightarrow R.V.

We are interested in modelling probability.

$p(o_c^i | x^i)$

here $o^i =$ output vector
 $x^i =$ input vector.

$o_c^i \rightarrow$ component of $o^i \rightarrow$ probability of x^i belonging to the class $C = \{1, 2, \dots, C\}$

C -components, where 1 is redundant.

prob of class $c = 1 - \sum \text{prob. (class } \neq c)$

$\langle \underline{0.1}, \underline{0.6}, \textcircled{0.3} \rangle$

$$1 - (0.6 + 0.1) \\ \Rightarrow 1 - 0.7$$

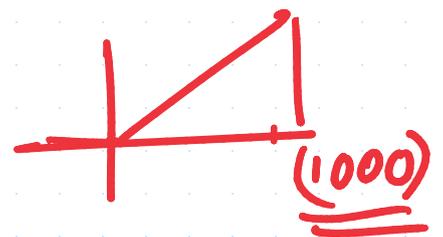
(C-1)

2-class \rightarrow sigmoid } activation functions
 > 2 \rightarrow softmax }

(0-100) \rightarrow what activation?



Sigmoid



(0-1) \rightarrow $\times 100$

(0-100)

(5, 95)

$(0 - 1)$ \rightarrow Sigmoid.

$$(0 - 1) \times 90 + 5.$$

$(0 - 90)$

$(5, 95)$

Scaling & translation.

BCE

Likelihood of observation for two classes.

For N o/p - i/p pairs.

likelihood.

$$L = \prod_{i=1}^N (o_i)^{t_i} (1 - o_i)^{(1 - t_i)}$$

where target
 $t^i = 1$ or 0 .

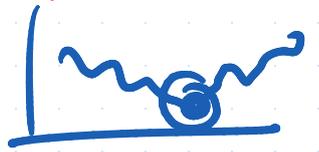
Binomial

① How likely an output
is belonging to target
L or 0.

$$\underline{\log L} = \log \left(\prod_{i=1}^N (o^i)^{t^i} (1-o^i)^{(1-t^i)} \right)$$

0.01 x 0.7 x 0.3 x 0.2

likelihood - maximize -



ML \rightarrow minimization algorithms.

$$-LL = \min \sum_{i=1}^N \left\{ t^i \log o^i + (1-t^i) \log(1-o^i) \right\}$$

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

Kullback-Leibler Divergence.

$$D_{KL}(p \parallel q) = \sum_{k=1}^n p_k \log \left(\frac{p_k}{q_k} \right)$$

$$D_{KL}(p \parallel q) = H(p, q) - H(p)$$

\swarrow
 \searrow

cross entropy
entropy

KL-divergence not KL distance

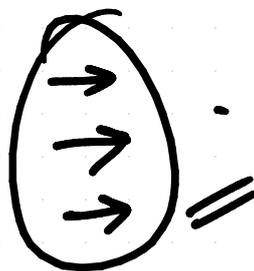
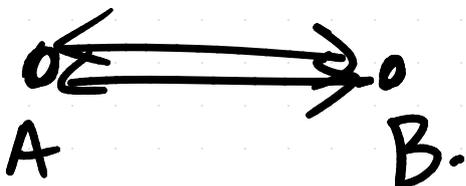
$$p_k \log\left(\frac{p_k}{q_k}\right)$$

$$p_k = 0.1$$

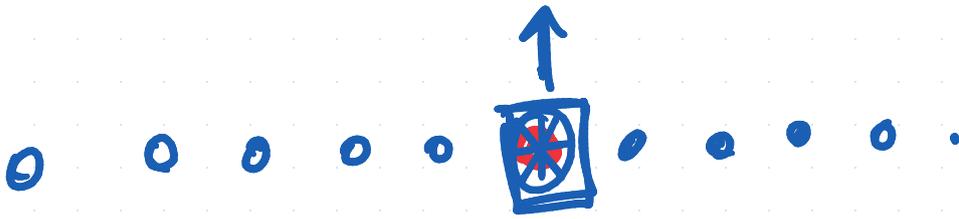
$$q_k = 0.7$$

$$0.1 \log\left(\frac{0.1}{0.7}\right) \neq 0.7 \log\left(\frac{0.7}{0.1}\right)$$

Not symmetric ; hence this is not a distance.



Entropy ↑ I.G. ↓



highly likely.

Categorical Cross entropy derivation.

For N no. of i/p o/p pairs:-

$$\mathcal{L} = \prod_{i=1}^N \prod_{k=1}^C (o_k^i)^{t_k^i} ; t_k^i = 1/0.$$

For a pattern i , only one of t_k^i 's is 1,

rest are 0:-

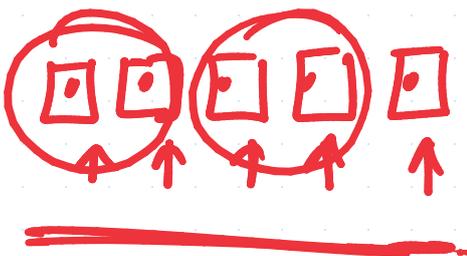
$$\max L = \sum_{i=1}^N \sum_{k=1}^C t_k^i \log o_k^i.$$

$$\min -L = - \sum_{i=1}^N \sum_{k=1}^C t_k^i \log o_k^i$$

for softmax as activation function,
we maximize the likelihood =

minimizing the $-L = \text{cross entropy}$

loss / error.



Backprop & Gradient Descent.

(greedy algorithm)

Always in the direction of reducing errors.

$$\Delta w_{ji} = -\eta \frac{\partial E}{\partial w_{ji}}$$

change in weight.

η = learning rate.

E = loss function.

w_{ji} = weight of connection from the i th neuron to the j th neuron.

