

$$L(\theta, P_a)$$

$$MAE = \frac{1}{N} \sum |e|$$

$$MSE = \frac{1}{N} \sum e^2$$

$$\frac{G.T}{P} \quad 0.8, 0.9$$

learn fast

Optimizing
Parameters
 $\alpha, \frac{\partial J(\theta)}{\partial w}$

* Choose loss function wisely

Entropy

Randomness / chaos (Measurement of these)

~~Information~~ Information

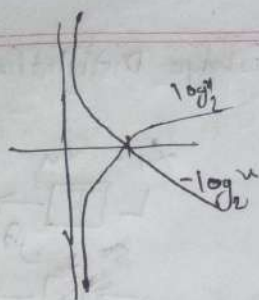
$P(E) = 1$ [where we know the output]

→ least Entropy

$P(E) < 1$ [where we are not sure about output]

→ Higher Entropy

		Bit length
A	0.3	$-\log_2 0.3$
B	0.6	$-\log_2 0.6$
C	0.05	$-\log_2 0.05$
D	0	$-\log_2 0$
F	0.05	$-\log_2 0.05$



(- to avoid get positive results)

Average bit length

$$\text{Bit length} = (0.3 \times -\log_2 0.3) \text{ take } 10 \text{ coke } 3 \text{ Sandwich } 20$$

$$+ (0.6 \times -\log_2 0.6)$$

$$+ (0.05 \times -\log_2 0.05)$$

$$\text{cost}_{\text{avg}} = \frac{10 + 3 \times 20}{3} = 14$$

$$= \frac{2 \times 10}{6} + \frac{3}{6} + \frac{3 \times 20}{6}$$

$$= \left(\frac{2}{6}\right) \times 10 + \left(\frac{1}{6}\right) \times 3 + \left(\frac{3}{6}\right) \times 20$$

equivalent to

$$\text{Average bit length} = - \sum_{i=1}^n p_i \log_2 p_i$$

Entropy \Rightarrow Measurement of cost, Information, Randomness

$$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \rightarrow \text{Entropy } 0$$

$$\begin{pmatrix} 0.2 \\ 0.3 \\ 0.3 \\ 0.2 \end{pmatrix}$$

Entropy ?

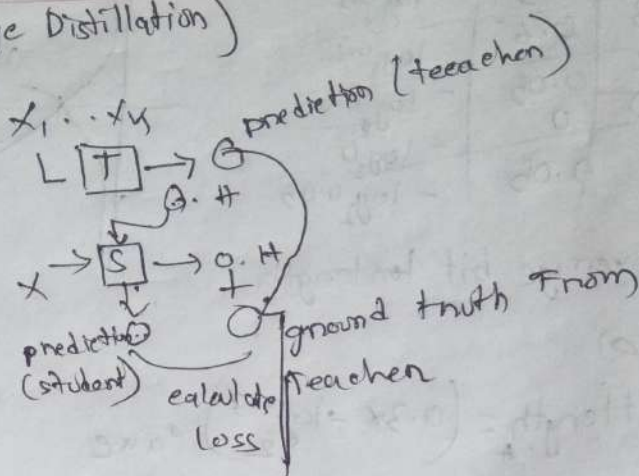
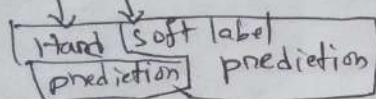
So we developed

cross-Entropy

What is the distance between these entropy, an amount of Avg bit length

Dark knowledge (knowledge Distillation)

Q	0	0.1
B	1	0.3
P	0	0.2
F	0	0.1
C	0	0.3



$$\alpha L_{0.4} + (1-\alpha) L_T$$

↓ cross Entropy

↓ KL Divergence

Dr. Sifat Momen (Sfm1)

Kacchi	0.1
Burger	0.3
Pizza	0.2
Fried chicken	0.4
$\Sigma = 1$	

$H(\text{Sifat Bhai})$

$$= 1.85$$

[checks Avg bit length equation]

Dr. Nabeel Mohammed (Nabeel)

Kacchi	0.25
Burger	0.1
Pizza	0.1
Fried	0.3
chicken	0.2
$\Sigma = 1$	

$H(\text{Nabeel})$

$$= 1.67 \text{ (Slightly less Diversity)}$$

What is the avg code length if Dr. Sifert uses my codebook.

$$0.1 * 2 + 0.3 * 63.32 + 0.2 * 3.32 + 0.4 * 1.73 = 2.435$$

(cross Entropy)

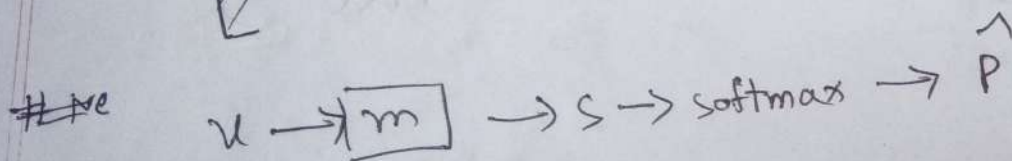
Equation 8

$$-\sum_{i=1}^L \frac{P_i^S}{\sum_{i=1}^L P_i^S} \log_{\frac{P_i^N}{\sum_{i=1}^L P_i^N}} \left[\text{cross Entropy} \right] \Rightarrow 2.435 \quad [\text{Extra 0.6}]$$

$$-\sum_{i=1}^L \frac{P_i^S}{\sum_{i=1}^L P_i^S} \log_{\frac{P_i^S}{\sum_{i=1}^L P_i^S}} \left[\text{Not cross Entropy} \right] \Rightarrow 1.85$$

$$KL(S,N) = \underbrace{-\sum_{i=1}^L P_i^S \log P_i^N}_{\text{cross Entropy}} - \underbrace{\left(-\sum_{i=1}^L P_i^S \log P_i^S \right)}_{\text{Entropy}}$$

$KL(S,N)$ not necessarily equals to $KL(N,S)$



$$\hat{P} = \begin{bmatrix} \hat{P}_0^N \\ \hat{P}_e \\ \hat{P}_{TB} \\ \hat{P}_N \end{bmatrix} \begin{matrix} \hat{P}_1 \\ \hat{P}_2 \\ \hat{P}_3 \\ \hat{P}_4 \end{matrix}$$

$$P = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \begin{matrix} P_1 \\ P_2 \\ P_3 \\ P_4 \end{matrix}$$

$P_N = 1$

$e = 2$

$TB = 3$

$N = 4$

$$KL(P, \hat{P}) = - \sum_{i=1}^4 P_i \log \hat{P}_i - \left(- \sum_{i=1}^4 P_i \log P_i \right)$$

$$\Downarrow$$

$$(-1 \log 0 + 0 + 1 \log 1 + \dots) = 0$$

[This will be eliminated]

Simplified (When Ground Truth is One Hot Encoding)

$$KL(P, \hat{P}) = - \sum_{i=1}^4 P_i \log \hat{P}_i$$

$$= -1 \log \hat{P}_3$$

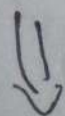
$$= -\log \hat{P}_3$$

log loss = Cross Entropy loss = logarithmic loss
= Negative log loss

$$\vec{x} \rightarrow \boxed{m} \rightarrow \vec{S} \rightarrow \text{Softmax} \rightarrow \hat{P} \rightarrow L(P, \hat{P})$$

$$\vec{S} = \begin{bmatrix} s_1 \\ s_2 \\ s_3 \\ s_4 \end{bmatrix}$$

$$\text{Softmax}(\vec{S}) = \begin{bmatrix} \hat{P}_1 \\ \hat{P}_2 \\ \hat{P}_3 \\ \hat{P}_4 \end{bmatrix}$$



$$\hat{P}_i = \frac{e^{s_i}}{e^{s_1} + e^{s_2} + e^{s_3} + e^{s_4}}$$

$$\sum_{i=1}^4 \hat{P}_i = 1$$

$$0 \leq \hat{P}_i \leq 1$$

$$= \frac{e^{s_i}}{\sum_{j=1}^4 e^{s_j}}$$

Alternatively

$$P_i = \frac{|s_i|}{\sum_{j=1}^4 |s_j|}$$

But for $-7, 7$ both it will get same result $(-7) = 7, (7) = 7$

$$\hat{P}_i = \frac{s_i^2}{\sum_{j=1}^4 s_j^2} \quad \left[\text{Same } 7^2 = 49, (-7)^2 = 49 \right]$$

Regulation \rightarrow (It reduces complexity of model, prevents overfitting)

$y = ax^3 + bx + c$

$$L(t, P) = \text{mae}(\quad) + R(\theta) \uparrow \left[\text{If overfit, } R(\theta) \uparrow \right]$$

then $a \rightarrow 0, b \rightarrow 0$
close to zero
 $y = ax^3 + bx + c$ becomes simple model

Also Dropout \rightarrow does not same things (off some neuron)

softmax function with Temperature;

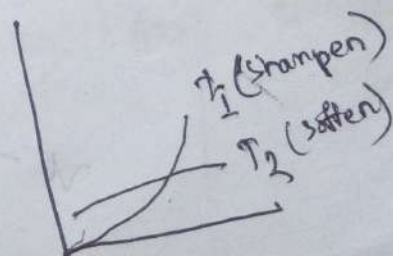
$$P_i = \frac{e^{s_i}}{\sum_{j=1}^4 e^{s_j}} \Rightarrow P_i = \frac{e^{(s_i/T)}}{\sum_{j=1}^4 e^{(s_j/T)}}$$

s_1	1
s_2	2
s_3	3
s_4	4

$$e' = e$$

	$T=1$	$T=2$
P_1	0.03	0.10
P_2	0.07	0.17
P_3	0.24	0.28
P_4	0.64	0.44
Σ	1	1

For $T=1, 2$



x_1	x_2	x_3	\dots	x_L
t_1	t_2	t_3	\dots	t_L
\downarrow	\downarrow			
$o_H(t_1)$	$o_H(t_2)$			

$$x \rightarrow [m] \rightarrow \vec{S} \rightarrow \text{softmax} \rightarrow \vec{P}$$

$$\mathcal{L} = -\log(P_{\text{argmax}}(o_H(t_i)))$$

$$\text{argmax} \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \end{pmatrix} = 0$$

$$m_{\text{Teacher}} = m$$

x_1	x_2	x_3	\dots	x_i
-------	-------	-------	---------	-------

$$\vec{S}_L = m_{\text{teacher}}(x_i)$$

x_1	x_2	x_3	\dots
-------	-------	-------	---------

$$\vec{S}_L = m_{\text{teacher}}(x_i)$$

Temperature controlled softmax

$$\vec{P}_L \rightarrow \text{Teacher}$$

$$x \rightarrow m_{\text{student}} \rightarrow \vec{S}^{\text{student}} \rightarrow \text{softmax}(\vec{S}^{\text{student}}) \rightarrow \vec{P}^{\text{student}}$$