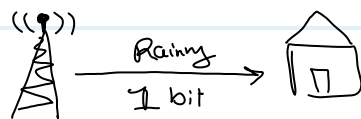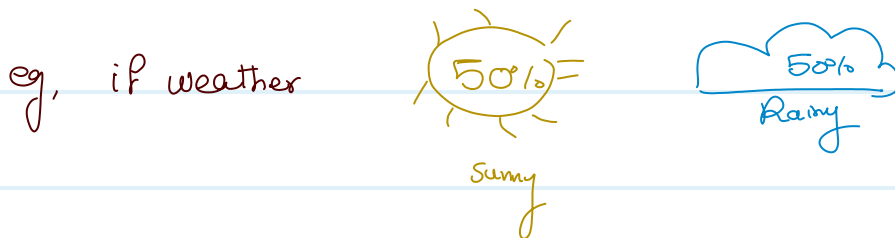in Machine learning ⇒

Cross-Entropy mostly used as cost function

in training the classifier

Come from Information theory (Claude Shannon) :

in theorem :

to Transfer 1 bit of information = Uncertainty divide by 2

eg, if weather



Sunny    Rainy



Rainy
1 bit

if send 1 bit it divide uncertainty by 2

& there's two options, now there's just 2 one

& mean : it send you in 1 bit of useful info even if it send 8 bits

eg, what if there's 8 states of weather equally likely

( Sunny , Cloudy , Rainy ------ )

When station send you information ⇒ it divide your uncertainty

by factor of 8=2^? = 2^3

Meaning ③ bits are the actual useful information

it's easy to calculate  #useful bits = $\log_2$( uncertainty Reduction factor )

$$= \log_2(8)$$

eg, what if probability not equally likely

75%   25%

if ((•)) —Rainny→ 🏠 ⟹ this mean your uncertainty drop By factor of $\underline{\underline{4}} = 2^{②}$

∞ Notes:  uncertainty Reduction $=$ inverse of probability of event

$$= \frac{1}{P}$$

so  #useful bits $= \log_2($ uncertainty Reduction $) = \log_2(\frac{1}{P}) = -\log_2(P)$

if ((•)) —sunny Raining→ 🏠    uncertainty Reduction $= \frac{1}{0.75} = 1.33$

#useful bits $= -\log(0.75) = 0.41$

Not so much Useful ∼ as I am already 75%

sure that's sunny Before you till me

⟹ Note | #useful bits | can be translated also as | useful info. |

What the average of useful information ?

↳ as there's 75% it's sunny ↗ it'll send (0.41) useful info

↳ also there's 25% it's Raing ↗ it'll send (2) useful info.

$$\& \#avg\ info = 0.75\ (0.41) + 0.25(2) = 0.81\ bits$$

$$= -p_S log(p_S) - p_R log(p_R)$$

$$= \text{Entropy} \implies \text{nice Measure of } \underline{uncertain} \\ \overline{Events\ are}$$

$$\text{Entopy} = H(\vec{P}) = -\sum_i p_i log(p_i) = \text{how much on avg information} \\ \text{you get} \longrightarrow \text{when you Sample} \\ \text{event from distributi} \\ \vec{P}$$

$$= \text{how unpredictable the probability} \\ \text{distribution}$$

Cross Entropy = avg message length ( that actually sent)

(1)



**Entropy**
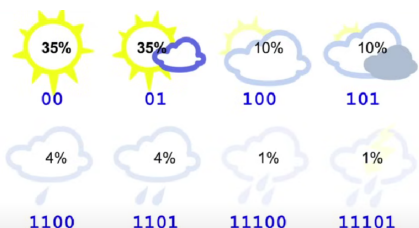$= 0.35 \log_2(0.35) + \ldots + 0.01 \log_2(0.01)$
$= 2.23\ bits$

Cross Entropy $= \underline{0.35} * 3 + \underline{0.35} * 3 + 0.\underline{1} * 3 + \ldots$

$$= 3\ bits$$

Other words: Wei sent 3 bits on avg

But also 2.23 on avg are useful

② 

35%×2 + 35%×2 + 10%×3
+ 10%×3 + 4%×4 + 4%×4
+ 1%×5 + 1%×5 = **2.42 bits**

⇐ We can do Better By changing Encoding

↑
Cross-Entropy

③ What if we use Encoding in different location with different distribut...



Cross-Entropy = 4.58 bits

Why? as we send ② bits for Sunny wheather ≃ we predicted that uncertainty Reduct-
knowing it's Sunny
= 4
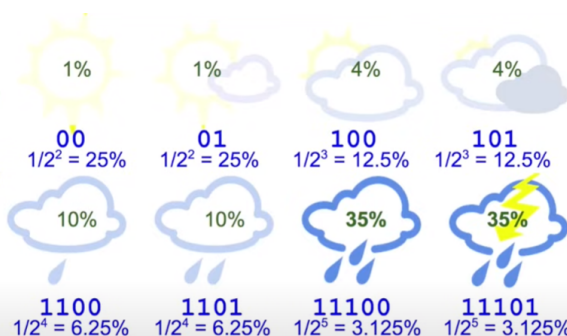
≃ meaning we predicted 25% Sunny

Very WRong assumption/predication

So our predicated probability = $\frac{1}{\text{uncertainty Red --}}$ = $\frac{1}{2^{\#bits}}$ = q

p = true distribution

q = predicted distribution



∴ Cross Entropy = $H(\vec{p}, \vec{q})$

= $-\sum_{i=1} p_i \log_2(q_i)$

Notes : _ if  Cross Entropy = Entropy ⟹ very efficient

_ if  Cross Entropy >>> Entropy ⟹ Make very wrong Assumption

_ "Cross Entropy" _ "Entropy" = Relative Entropy

$$= \text{Kull back\_leibler divergence}$$

$$= \text{KL Divergence}$$

OR Better :

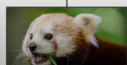$$\text{Cross Entropy} = \text{Entropy} + \text{KL Diverge}$$

OR :

$$D_{KL}(P \| q) = H(P, q) - H(P)$$

Cross Entropy as Cost Function :

True distribution:

| | Cat | Dog | Fox | Cow | Red Panda | Bear | Dolphin |
|---|---|---|---|---|---|---|---|
| True distribution: | 0% | 0% | 0% | 0% | 100% | 0% | 0% |
| Predicted distribution: | 2% | 30% | 45% | 0% | 25% | 5% | 0% |

Classifier

natural log →

**Cross-Entropy Loss:**
$H(\mathbf{p}, \mathbf{q}) = -\Sigma_i p_i \log(q_i)$
$= -\log(0.25) = 1.386$

Since it's One hot Encoding

then "Cross Entropy loss" $= -\log(0.25)$

if it predicated "1" ⟹ loss = 0

if it predicated "0" ⟹ Loss = 1