

Cross-Entropy & KL Divergence

Entropy is a measure of how surprised we are.

⇒ The more unpredictable / uncertain a system is, the more potential for surprise (e.g. surprise of correctly predicting lottery >> coin flip)

Principles for information content of event, $I(E)$.

① Continuous function

② Surprise \uparrow as probability of event \downarrow \rightarrow relate to $\frac{1}{p}$

③ Info from 2 independent events should be additive

↓

$$I(E) = \log\left(\frac{1}{p}\right) = -\log(p)$$

E.g. Flip 2 coins:

- Prob. ([Head, Head]) = $\frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$
 - Info ([Head, Head]) = $-\log\left(\frac{1}{4}\right) = 2$
- } we gain 2 bits of info.
 $E \in \{00, 01, 10, 11\}$

Def. One bit of info = info that reduces uncertainty by $\frac{1}{2}$

- Case 1: 25% HH, 25% HT, 25% TH, 25% TT

⇒ Knowing 1st flip is H is 1 bit
("Surprise" decreases by $\cdot \frac{-\log 4}{-\log 2} = \frac{2}{1} = 2$)

- Case 2: 75% Sunny, 25% Rain

⇒ Getting "Sunny" = $-\log_2(0.75) = 0.415$ bits

⇒ Getting "Rainy" = $-\log_2(0.25) = 2$ bits

} A more surprising outcome gives us more info to update our existing beliefs!

Entropy. Average uncertainty / expected info gain of an RV.

$$H(S) = \underbrace{-\sum_{i \in S}}_{\text{across all events}} \underbrace{P(i)}_{\text{prob. of event}} \underbrace{\log P(i)}_{\text{info gain from event}}$$

} expected value!!

Greater uncertainty in system \rightarrow higher entropy \rightarrow higher expected info gain from each obs

Know nothing abt coin

↓
Full 1 bit info

Know coin is 99% H

↓
Little info gain, v. likely H

Cross-Entropy: How surprised are we when we use an estimated distribution q to predict a true distribution p ?

$$L = - \sum_{k=1}^{K \text{ classes}} \underbrace{p(k)}_{\text{actual prob of observing } k} \underbrace{\log q(k)}_{\text{our surprise when observing } k}$$

\Rightarrow Larger divergence between p and $q \rightarrow$ obs are more surprising on avg
 \rightarrow higher entropy \rightarrow more info gain when updating our beliefs !!

[Perfect for loss function — captures how surprised I am by the correct ans']

$$\begin{aligned} \text{KL Divergence: } KL(P \parallel Q) &= \overbrace{-\sum_{x \in X} P(x) \log Q(x)}^{\text{surprise using } Q} - \overbrace{\sum_{x \in X} P(x) \log P(x)}^{\text{surprise using } P} \\ &= -\sum_{x \in X} P(x) \cdot [\log Q(x) - \log P(x)] \\ &= -\sum_{x \in X} P(x) \cdot \log \left(\frac{Q(x)}{P(x)} \right) \end{aligned}$$

\downarrow
Measure of how diff 2 prob dists are'

\hookrightarrow Cross-Entropy: (surprise due to divergence of P & Q) + (inherent surprise in p)

\Downarrow
KL divergence isolates this !!