

CROSS-ENTROPY

$$H_c(X_1, X_2) = - \sum_{x \in D} p_1(x) \log(p_2(x))$$

↓
Both take value x from D

$$= - \int_{x \in D} p_1(x) \log(p_2(x)) dx \quad \begin{matrix} \text{continuous} \\ \text{random var} \end{matrix}$$

$$= - \int_{\bar{x} \in D} p_1(\bar{x}) \log(p_2(\bar{x})) d\bar{x} \quad \text{vector}$$

majority contribution from dissimilarity

- $P_{gt} = 1, P_{pred} \approx 1 \Rightarrow$ does not contribute to H_c
- $P_{gt} \neq 1, P_{pred} \neq 1 \Rightarrow$ high H_c
- $P_{gt} \downarrow \Rightarrow$ will not contribute a lot

uniform distri → high entropy
steep peaks + low elsewhere → low entropy

minimum no. of bits needed to communicate a message if we know where message is drawn from

K-L DIVERGENCE

Cross entropy = entropy + divergence in distri

$$H_c(X_1, X_2) = H(X_1) + D(X_1, X_2)$$

$$D(X_1, X_2) = H_c(X_1, X_2) - H(X_1)$$

$$D(X_1, X_2) = - \sum_{x \in D} p_1(x) \log(p_2(x)) + \sum_{x \in D} p_1(x) \log p_1(x)$$

$$= - \sum_{x \in D} p_1(x) \log(p_1(x) - p_2(x))$$

$$KL = - \sum_{x \in D} p(x) \log \left(\frac{p_1(x)}{p_2(x)} \right)$$

$$= - \int_{x \in D} p_1(x) \log \left(\frac{p_1(x)}{p_2(x)} \right)$$

→ multi-class classification: $\min(H_c) \geq \min(KL)$

→ KL not symmetric, depends on gt

→ always positive

→ KL divergence b/w identical random vars = 0

$$\log \left(\frac{p_1(x)}{p_1(x)} \right) = \log(1) = 0$$

JS DIVERGENCE

- KL asymm → depends on which distri is considered "reference"

- JS = symm

$$JSD(P||Q) = \frac{KL(P||M) + KL(Q||M)}{2}$$

order matters!

$$D_{JS}(P||Q) = \frac{1}{2} D_{KL}(P, \frac{P+Q}{2}) + \frac{1}{2} D_{KL}(Q, \frac{P+Q}{2})$$

where $M = \text{avg var} = \frac{P+Q}{2}$

→ $JSD > 0 \Rightarrow$ diff distri

$$P = [0.7, 0.2, 0.1] \quad Q = [0.3, 0.4, 0.3]$$

$$M = [0.5, 0.3, 0.2]$$

$$D_{KL}(P||M) = - \sum_{i=1}^n P_i \left(\frac{\log P_i}{\log Q_i} \right)$$

$$= - \left[0.7 \left(\frac{\log 0.7}{\log 0.5} \right) \right] - \left[0.2 \left(\frac{\log 0.2}{\log 0.3} \right) \right] - \left[0.1 \left(\frac{\log 0.1}{\log 0.2} \right) \right]$$

$$= 0.1228$$

$$D_{KL}(Q||M) = 0.1204$$

$$D_{JS} = \frac{D_{KL}(P||M) + D_{KL}(Q||M)}{2} = 0.125 \text{ bits}$$

P, Q are similar

CONDITIONAL ENTROPY

Entropy of one variable given that another is a constant

$$H(w|H=h_i) = - \sum_{j=1}^n P(w_j|h_i) \log(P(w_j|h_i))$$

$$H(w|H) = \sum_{i=1}^n P(h_i) \left[- \sum_{j=1}^n P(w_j|h_i) \log(P(w_j|h_i)) \right]$$

w_1	...	w_n
h_1	...	h_n
:	...	:
one row		