

Optimization

Mahdi Roozbahani Georgia Tech

Outline

Motivation

Entropy

Conditional Entropy and Mutual Information

Cross-Entropy and KL-Divergence



Let's work on this subject in our Optimization lecture

Cross Entropy

Cross Entropy: The expected number of bits when a wrong distribution Q is assumed while the data actually follows a distribution P

$$H(p,q) = -\sum_{x \in \mathcal{X}} \overline{p(x)} \, \log \overline{q(x)} = H(P) + KL[P][Q]$$

This is because:

$$egin{align} H(p,q) &= \mathrm{E}_p[l_i] = \mathrm{E}_p\left[\lograc{1}{q(x_i)}
ight] \ H(p,q) &= \sum_{x_i} p(x_i)\,\lograc{1}{q(x_i)} \ H(p,q) &= -\sum_{x} p(x)\,\log q(x). \end{gathered}$$

Labeling target values Label encoding (ordinal) and One-hot encoding

$$X = \begin{bmatrix} h & \omega & age=a & \cdots \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & \\ & & & \\$$

Why Cross entropy and not simply use dot product?

$$CE = -\sum_{i=1}^{n} P(x) \log_{i} Q(x)$$

$$CE = -\left(1x\log_{i} 0.8 + 0x\log_{i} 0.1 + 0x\log_{i} 0.1\right) - \left(0x\log_{i} 0.3 + \cdots\right) - \frac{\log_{i} x}{\log_{i} x}$$

$$\frac{\partial_{i}}{\partial x} \log_{i} x$$

Kullback-Leibler Divergence

Another useful information theoretic quantity measures the difference between two distributions.

$$\begin{aligned} \mathbf{KL}[P(S) \| Q(S)] &= \sum_{s} P(s) \log \frac{P(s)}{Q(s)} \\ &= \underbrace{\sum_{s} P(s) \log \frac{1}{Q(s)}}_{\mathbf{Cross\ entropy}} - \mathbf{H}[P] = H(P,Q) - H(P) \end{aligned}$$
 KL Divergence is

Excess cost in bits paid by encoding according to Q instead of P.

a **KIND OF**distance
measurement

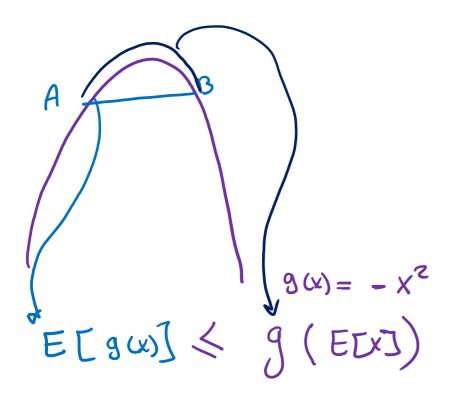
$$-\mathbf{KL}[P\|Q] = \sum_{s} P(s) \log \frac{Q(s)}{P(s)}$$
 log function is concave or convex?
$$\sum_{s} P(s) \log \frac{Q(s)}{P(s)} \leq \log \sum_{s} P(s) \frac{Q(s)}{P(s)} \quad \text{By Jensen Inequality}$$

$$= \log \sum_{s} Q(s) = \log 1 = 0$$

So $KL[P||Q] \ge 0$. Equality iff P = Q

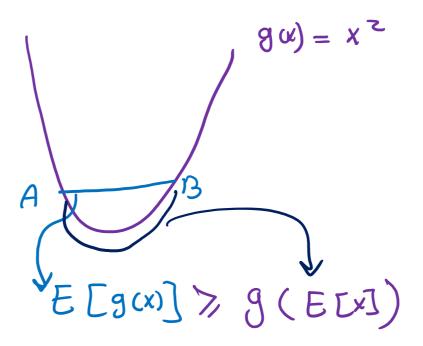
When P = Q, KL[P||Q] = 0

Concave



Jensen Inequality

Convex



$$g'(x) = 2x$$
$$g''(x) = 2$$

$$-\text{KL}\left[P\right]\left[Q\right] = \sum_{\alpha} P(x) \log \left(\frac{\alpha(x)}{P(x)}\right)^{\alpha(x)} = g(x) = \sum_{\alpha} P(x) \log g(x) = E\left[\log g(x)\right]$$

$$-\text{KL}\left[P\right]\left[Q\right] = E\left[\log g(x)\right] \leq \log\left(E\left[g(x)\right]\right)$$

$$\leq \log\left(\sum_{\alpha} P(x) g(x)\right)$$

$$\leq \log\left(\sum_{\alpha} P(x)\right)$$

$$\leq \log\left(\sum_{\alpha} P(x$$

-> Probabilistic models -> Gaussian distribution DE [1,5] MLE ~> derivative OPtimization > Non- Probabilistic models ~> No-constraint ~> derivative Equality Constraint In equality constraint Lagrange function -> 4 conditions derivative derivative

$$f(M,S) = 6M^2 + 3S^2$$

M # hours you study ML perday

S # hours you sleep per day

$$\frac{\partial f(M,s)}{\partial M} = 0 \implies 12M = 0 \implies M = 0$$

$$\frac{\partial f(M_3S)}{\partial S} = 0 \implies \delta S = 0 \implies S = 0$$

$$f(M,S) = 6M^2 + 3S^2 \rightarrow \text{Objective function}$$

$$S.t \quad M+S = 24 \Rightarrow g(M,S) = M+S-24$$

$$S.t \quad M-S=8$$

$$L(M,S,S) = f(M,S) - Sg(M,S)$$

$$\sum_{s=1}^{\infty} f(M,S) = \sum_{s=1}^{\infty} f(M,S) - Sg(M,S) = 0$$

$$\sum_{s=1}^{\infty} f(M,S) - Sg(M,S) = 0$$

$$\nabla f(M,S) = S \nabla g(M,S)$$

$$\nabla + (M,S) \approx \nabla g (M,S)$$

5 f(Ms) - S19, (Ms) - S29, (Ms)

 $\lfloor (M_3S_3S_1,S_2) = 1$

$$L(M, S, S) = 6M^{2} + 3S^{2} - S(M+S-24)$$

$$\frac{\partial L(M, S, S)}{\partial S} = 0 \Rightarrow M+S - 24 = 0 \Rightarrow M+S = 24 \Rightarrow \frac{S}{12} + \frac{S}{6} = 24$$

$$\frac{\partial L(M, S, S)}{\partial M} = 0 \Rightarrow 12M - S = 0 \Rightarrow M+S = 24$$

$$\frac{\partial L(M, S, S)}{\partial M} = 0 \Rightarrow 12M - S = 0 \Rightarrow M+S = 24$$

$$\frac{\partial L(M, S, S)}{\partial S} = 0 \Rightarrow 6S - S = 0 \Rightarrow S = \frac{S}{6} = 16$$

$$8 + 16 = 24$$

$$f(w,s) = 6w_5 + 3s_5$$

M+5-24 50

$$\nabla L(M,S,S) = 0$$

$$g(M,S) \leqslant 9$$

$$g(M,S) \neq 0$$

$$\delta = 0$$

$$g(M,S)=0$$

$$\int (x) = x^2$$

$$\int \frac{\partial f(x)}{\partial x} = 0 \implies 2x = 0 \implies x = 0$$

$$\text{Bradient direction}$$

$$\begin{cases} \{+1\} \\ \times \\ = \\ \times \end{cases} = \frac{\lambda f(x)}{\lambda x}$$

$$\begin{cases} earning \ rate \\ \alpha = 0.02 \end{cases}$$

$$X = X - 92X$$

M, M, S ~ 1) initialize randomly M, M, S

Example 1:

https://www.geogebra.org/3d/srzmv8uh

Example 2:

https://www.geogebra.org/3d/syhkqpk7