

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

$$H(p, q) = - \sum_{x \in \mathcal{X}} \overbrace{p(x)}^{\text{actual pdf}} \log \overbrace{q(x)}^{\text{predicted pdf}} = H(P) + KL[P][Q]$$

This is because:

$$H(p, q) = \mathbb{E}_p[l_i] = \mathbb{E}_p \left[\log \frac{1}{q(x_i)} \right]$$

$$H(p, q) = \sum_{x_i} p(x_i) \log \frac{1}{q(x_i)}$$

$$H(p, q) = - \sum_x p(x) \log q(x).$$

Labeling target values

Label encoding (ordinal) and One-hot encoding

$$\begin{array}{l}
 X = \begin{bmatrix} & h & w & a & \dots \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix}_{n \times d}
 \end{array}$$

$$\begin{array}{l}
 \text{Cat} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \\
 \text{dog} = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \\
 \text{fish} = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}
 \end{array}$$

$$\begin{array}{l}
 Y = \begin{bmatrix} \text{cat} \\ \text{dog} \\ \text{fish} \\ \text{cat} \\ \vdots \end{bmatrix}_{n \times 1} = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 1 \\ \vdots \end{bmatrix} \\
 \downarrow \\
 \text{actual}
 \end{array}$$

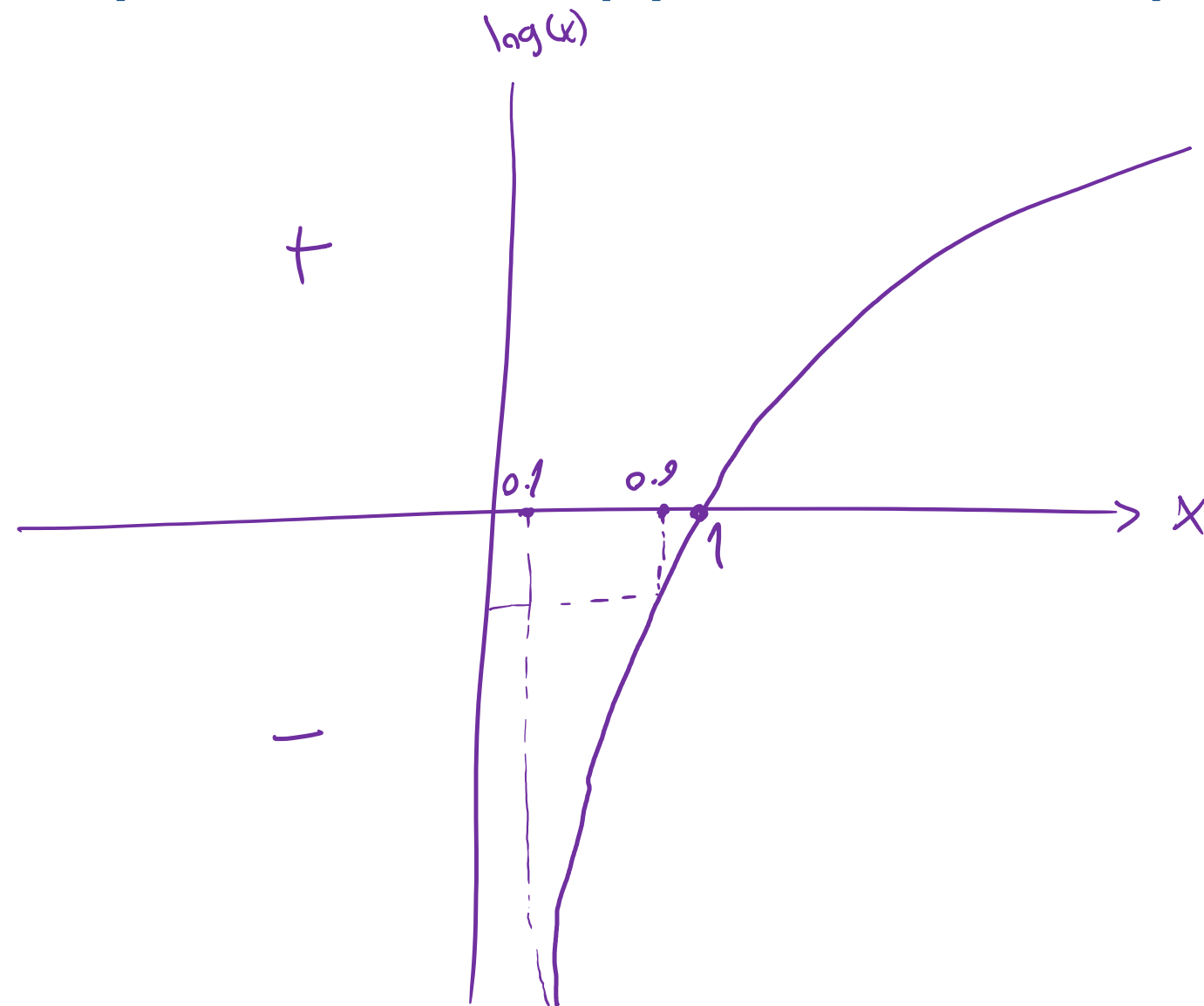
$$\begin{array}{l}
 \hat{Y}_{\text{Predicted}} = \begin{bmatrix} 1 \\ 2 \\ 2.5 \\ 1.5 \\ \vdots \end{bmatrix}
 \end{array}$$

$$\begin{array}{l}
 Y_{\text{actual}} = \begin{bmatrix} \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \\ \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \\ \begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \\ \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \\ \vdots \end{bmatrix} \\
 \hat{Y}_{\text{Predicted}} = \begin{bmatrix} \begin{bmatrix} 0.8 & 0.1 & 0.1 \end{bmatrix} \\ \begin{bmatrix} 0.3 & 0.6 & 0.1 \end{bmatrix} \\ \begin{bmatrix} 0.1 & 0.2 & 0.7 \end{bmatrix} \\ \vdots \end{bmatrix}
 \end{array}$$

$$\text{Loss (dot product)} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0.8 & 0.1 & 0.1 \end{bmatrix} + \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0.3 & 0.6 & 0.1 \end{bmatrix} + \dots$$

$$CE = H(p, q) = - \sum p(x) \log_2 q(x) = - [1 \log_2 0.8 + 0 \log_2 0.1 + 0 \log_2 0.1] - [0 \log_2 0.3 + 1 \log_2 0.6 + \dots]$$

Why Cross entropy and not simply use dot product?



$$\log 1 = 0$$

$$\log 0.9$$

$$\log 0.1$$

Kullback-Leibler Divergence

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

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

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

KL Divergence is
a **KIND OF**
distance
measurement

$$-\mathbf{KL}[P||Q] = \sum_s P(s) \log \frac{Q(s)}{P(s)}$$

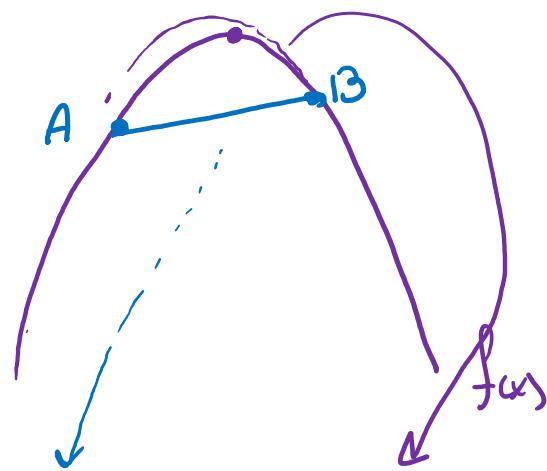
log function is
concave or
convex?

$$\begin{aligned}\sum_s P(s) \log \frac{Q(s)}{P(s)} &\leq \log \sum_s P(s) \frac{Q(s)}{P(s)} && \text{By Jensen Inequality} \\ &= \log \sum_s Q(s) = \log 1 = 0\end{aligned}$$

So $\mathbf{KL}[P||Q] \geq 0$. Equality iff $P = Q$

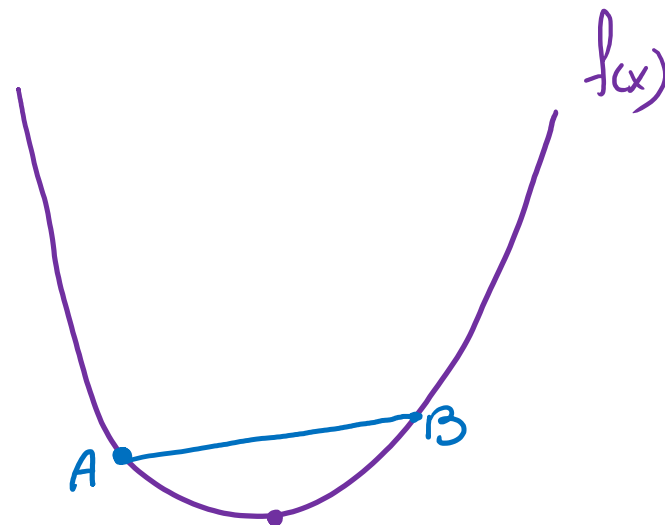
When $P = Q$, $KL[P||Q] = 0$

Concave



$$E[f(x)] < f(E[x])$$

Convex

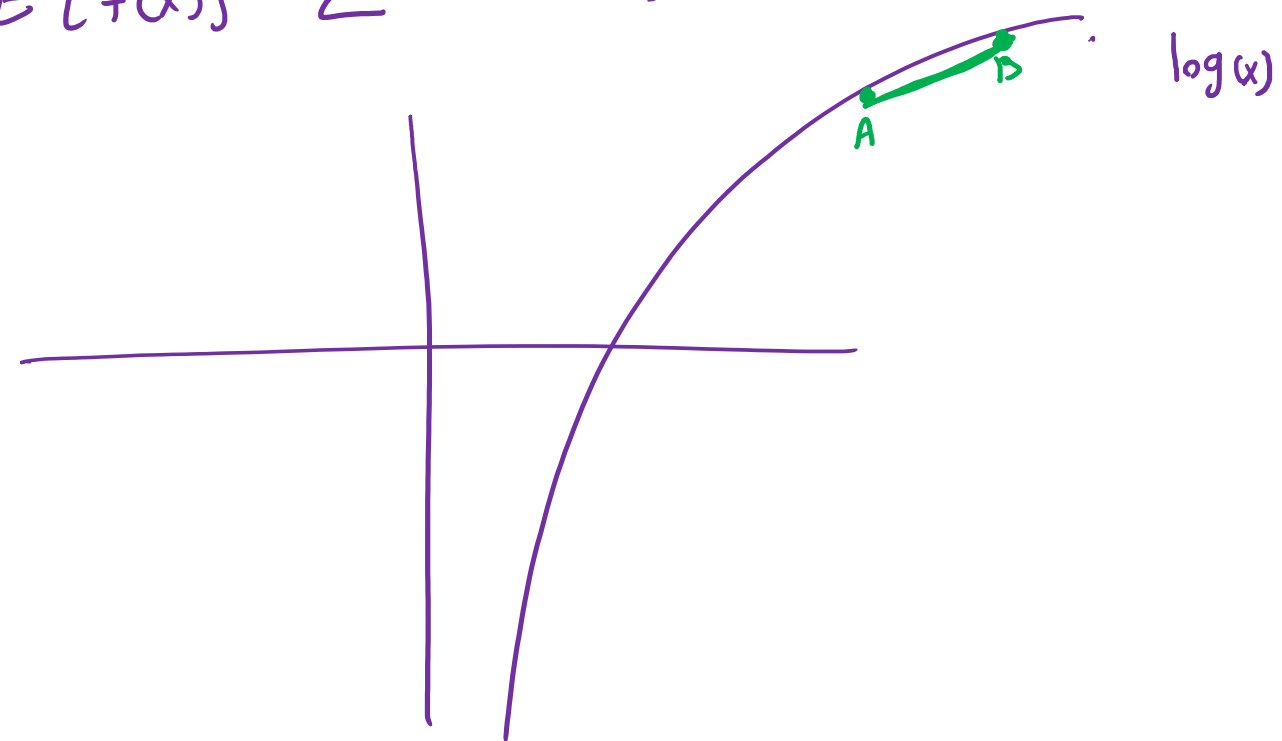
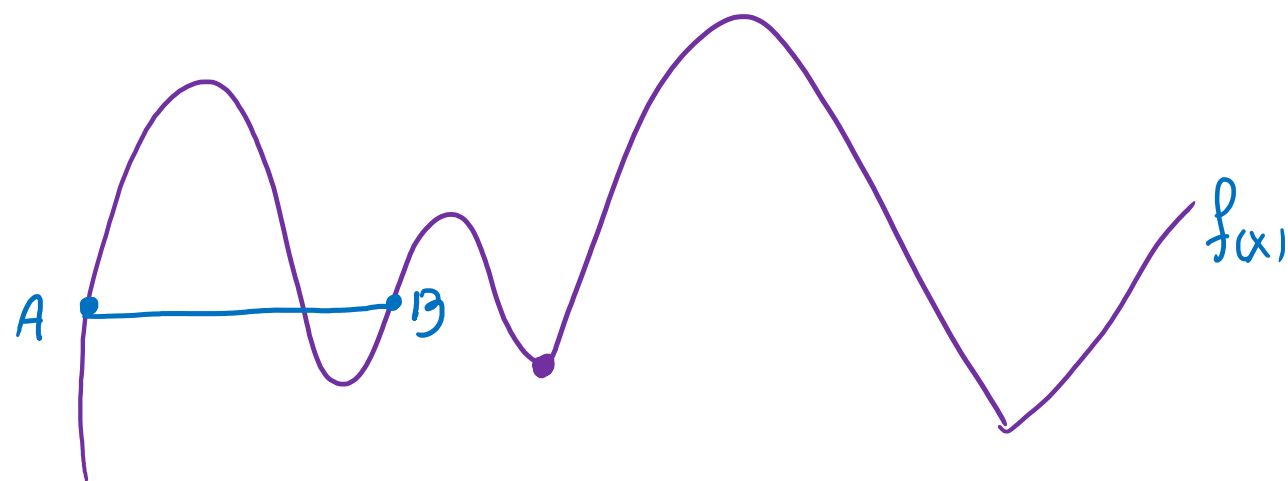


$$f(x) = x^2$$

$$f''(x) = 2$$

$$E[f(x)] > f(E[x])$$

$$E[f(x)] = \sum p(x) f(x)$$



$$E[\log(x)] < \log(E[x])$$

$$-KL[P][Q] = \sum P(x) \log \underbrace{\left(\frac{Q(x)}{P(x)} \right)}_{g(x)} = \sum P(x) \log g(x) = E[\log g(x)]$$

$$-KL[P][Q] = E[\log g(x)] \leq \log(E[g(x)])$$

$$\dots \leq \log(\sum P(x) g(x))$$

$$\dots \leq \log\left(\sum \cancel{P(x)} \frac{Q(x)}{\cancel{P(x)}}\right)$$

$$\dots \leq \log(\sum Q(x))$$

$$-KL[P][Q] \leq \log(1) = 0$$

$$KL[P][Q] \geq 0$$

Optimization

- ① Probabilistic models \rightsquigarrow Gaussian distribution $\begin{cases} \mu \\ \Sigma \end{cases} \Rightarrow \frac{MLE}{\downarrow}$
take a derivative
- ② Non-probabilistic models \rightsquigarrow constraints $\begin{cases} \text{Equality} \rightsquigarrow \text{Lagrange multipliers} \\ \text{Inequality} \rightsquigarrow \text{KKT} \end{cases}$

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

M : # hours you study ML per day
 S : # " you sleep " "

$$\frac{\partial f(M, S)}{\partial M} = 0 \Rightarrow 12M = 0 \Rightarrow M = 0$$

$$\frac{\partial f(M, S)}{\partial S} = 0 \Rightarrow 6S = 0 \Rightarrow S = 0$$

$f(m,s) = 6m^2 + 3s^2 \leadsto$ objective function

$$\text{s.t. } m+s=24 \leadsto g(m,s) = m+s-24$$

\Downarrow

$$L(m,s,\lambda) = f(m,s) - \lambda g(m,s)$$

$$L(m,s,\lambda) = 6m^2 + 3s^2 - \lambda (m+s-24)$$

$$\frac{\partial L(m,s,\lambda)}{\partial \lambda} = 0 \Rightarrow m+s=24 \Rightarrow \frac{\lambda}{12} + \frac{\lambda}{6} = 24 \Rightarrow \lambda = 96$$

$$\frac{\partial L(m,s,\lambda)}{\partial m} = 0 \Rightarrow 12m - \lambda = 0 \Rightarrow m = \frac{\lambda}{12} = \frac{96}{12} = 8$$

$$\frac{\partial L(m,s,\lambda)}{\partial s} = 0 \Rightarrow 6s - \lambda = 0 \Rightarrow s = \frac{\lambda}{6} = \frac{96}{6} = 16$$

$$m+s = 8+16 = 24 \quad \checkmark$$

$f(m,s)$

$\lambda, \alpha, \beta, \dots$

$$g(m,s) = m+s-24$$

$$h(m,s) = m-s-8$$

what is you have
two constraints

$$L = (m,s,\lambda_1,\lambda_2) = f(m,s) - \lambda_1 g(m,s) - \lambda_2 h(m,s)$$

$$\nabla f(m,s) = \nabla g(m,s)$$

$$\nabla f(m,s) - \nabla g(m,s) = 0$$

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

$$\text{s.t. } M + S \leq 24 \Rightarrow g(M, S) = M + S - 24$$

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

$$L(M, S, \lambda) = f(M, S) - \lambda g(M, S)$$

We need to satisfy 4 conditions

① Stationary condition

② Primal feasibility $g(M, S) < 0$

③ Dual feasibility $\lambda \geq 0$

④ Complementary slackness

$$g(M, S) \lambda = 0 \begin{cases} \lambda = 0 \text{ then } g(M, S) \neq 0 \\ \lambda \neq 0 \text{ then } g(M, S) = 0 \end{cases}$$

Example 1:

<https://www.geogebra.org/3d/srzm8uh>

Example 2:

<https://www.geogebra.org/3d/sy8kpk7>