Machine Learning CS 4641-7641



## 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

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

This is because:

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



$$E\left[-\log P + \log P - \log \varphi\right] \quad \log \alpha - \log \beta = \log \alpha \\ = -\log \beta \\ = -\log \beta \\ = -\log \beta \\ = \sum_{l=1}^{N} P(x) g(k)$$

$$E\left[E - \log \varphi\right] + E\left[-\log \varphi\right] \\E\left[E - \log \varphi\right] = \sum_{l=1}^{N} P(x) g(k)$$

$$E\left[E - \log \varphi\right] \\E\left[E - \log \varphi\right] = \sum_{l=1}^{N} P(x) g(k)$$

$$E\left[E - \log \varphi\right] \\E\left[E - \log \varphi\right$$





## **Kullback-Leibler Divergence**

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

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

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

a **KIND OF** distance measurement

log function is concave or convex?

$$\begin{aligned} -\mathsf{KL}[P||Q] &= \sum_{s} P(s) \log \frac{Q(s)}{P(s)} \\ \sum_{s} P(s) \log \frac{Q(s)}{P(s)} &\leq \log \sum_{s} P(s) \frac{Q(s)}{P(s)} \end{aligned} \qquad \begin{array}{l} & \text{By Jensen Inequality} \\ &= \log \sum_{s} Q(s) = \log 1 = 0 \end{aligned}$$

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

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







$$EE^{hog_{X}} EEE^{hog_{X}} \leq f(ECX) \Rightarrow EUo_{X} \leq bg(ECX)$$

$$EEg(x) = \sum P(x) \log (Q(x))^{2}g(x) = (\sum P(x) \log g(x))$$

$$EE(\log g(x)) = EE(\log g(x))$$

$$EE\log g(x) = EE\log g(x)$$

$$EE\log g(x) \leq \log (EEg(x))$$

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

$$\leq \log g \sum P(x) \frac{Q(x)}{P(x)}$$

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

$$-kLEP EQ \leq Q(x)$$

$$-kLEP EQ \leq Q(x)$$

Optimization MIN 
$$f(x) \rightarrow objective function
 $s.t g(x) = c \rightarrow Equality constraint
 $g(x) \leq c \qquad j \Rightarrow Inequality constraint \rightarrow KKT conditions$   
 $Lagrange function$   
 $f(x,y) = c \qquad H= \begin{bmatrix} \frac{\partial^2 f(y_x)}{\partial x^2} & \frac{\partial^2 f(x_y)}{\partial x \partial y} \\ \frac{\partial^2 f(y_y)}{\partial y \partial x} & \frac{\partial^2 f(y_y)}{\partial y^2} \end{bmatrix}$   
(D Linear Programming  $f(x,y) = x+y$   
(2) Quadratic Programming  $f(x,y) = x^2+y$   
 $s.t x+y=2$   
(3) Non linear Programming  $f(x,y) = x^3+y^2$   
 $s.t x^2+y=2$$$$

$$Minp = 6M^2 + 35^2$$

$$\frac{\partial f(M,s)}{\delta M} = 0 \implies 12M+0 = 0 \implies M=0$$
  
$$\frac{\partial f(M,s)}{\delta S} = 0 \implies 0+6S = 0 \implies S=0$$
  
$$\frac{\partial s}{\delta S}$$

S, W, B f(M,S) = 6M<sup>2</sup> + 3S<sup>2</sup> ~> objective function S.t.  $M+S=24 \rightarrow g(M,S)=M+S-24 \rightarrow Constraint function$ L(M,S,S) = f(M,S) - Sg(M,S) $L(M, S, S) = (6M^2 + 3S^2 - S(M+S - 24))$  $\frac{\partial L}{\partial S} = 0 \implies 0 + 0 - (M + S - 2Y) = 0 \implies M + S = 24 \implies \frac{S}{12} + \frac{S}{6} = 24 \implies$  S = 96  $\frac{\partial L}{\partial M} = 0 \implies 12 M + 0 - S(1) = 0 \implies M = \frac{S}{12} = \frac{96}{12} = 8$  $\frac{\partial L}{\partial s} = 0 \implies 6s - s = 0 \implies (s) = \frac{5}{6} = \frac{96}{6} = 16$ M + S = 29 8 + 16 = 29

$$f(M_{3}S) = 6M^{2} + 3S^{2}$$
St.  $M_{+}S = 2Y \implies g(M_{3}S) = M_{+}S - 2Y \implies S_{1}$ 
S.t.  $M_{-}S = IZ \implies h(M_{3}S) = M_{-}S - IZ \sim S_{2}$ 

$$L(M,S,S_1,S_2) = f(M,S) - S_1 g(M,S) - S_2 h(M,S)$$

L(M, S, S) = I(M, S) - Sg(M,S) $\nabla L(M,S,S) = 0 \implies \nabla (f(M,S) - Sg(M,S)) = 0$ 

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



$$f(M_{SS}) = \frac{1}{2}M^{2} + \frac{1}{2}S^{2}$$

$$M + S = 24$$

$$L(M_{SS},S) = \frac{1}{2}M^{2} + \frac{1}{2}S^{2} - S(M + S - 24)$$

$$\frac{\partial L}{\partial M} = 0 \Rightarrow M - S = 0 \Rightarrow M = S$$

$$\frac{\partial L}{\partial S} = 0 \Rightarrow S = S$$

$$L(S) = \frac{S^2}{2} + \frac{S^2}{2} - S(S+S-2Y) = S^2 - 2S^2 + 24S = -S^2 + 24S$$
  
Dual form  $\rightarrow$  concave  $\rightarrow$  Max  
 $\frac{\delta L(S)}{\delta S} = 0 \Rightarrow -2S + 24 = 0 \Rightarrow S = 12$ 





Derivation rules:

Quotient rule: 
$$\left(\frac{f}{g}\right)^{r} = \frac{fg - fg^{r}}{g^{2}}$$
  
Product rule:  $\left(fg\right)^{r} = fg + fg^{r}$   
Choin rule:  $\left(\psi^{n}\right)^{r} = n\psi^{r}\psi^{n-1}$