

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

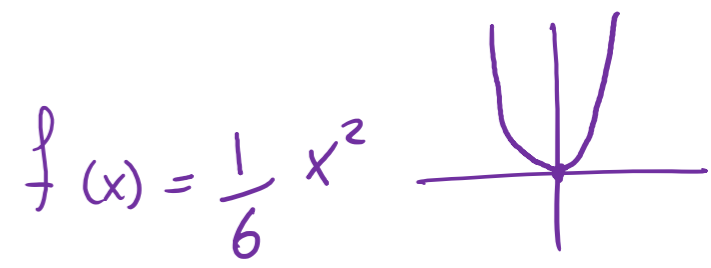
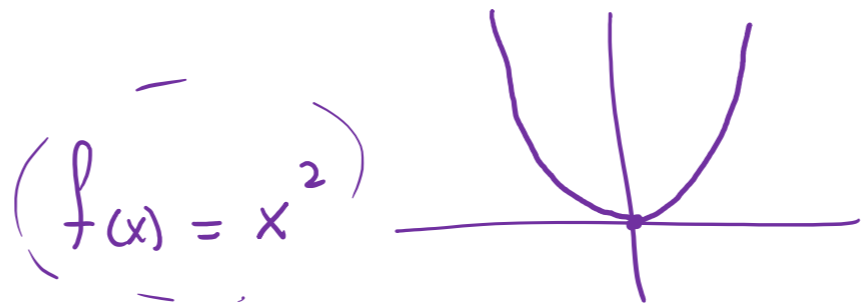
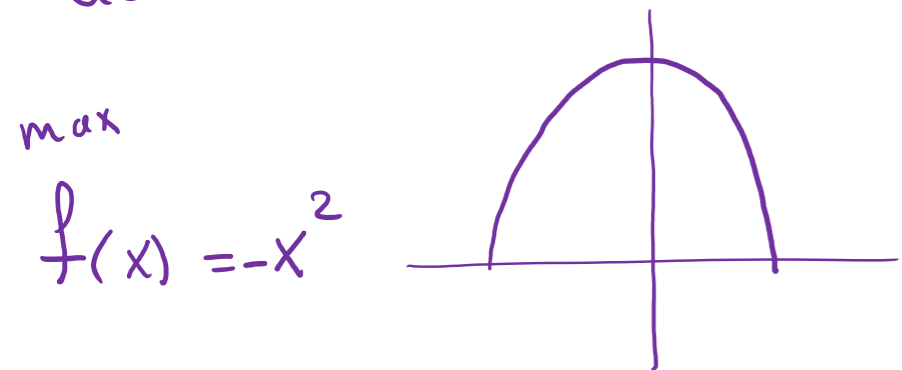
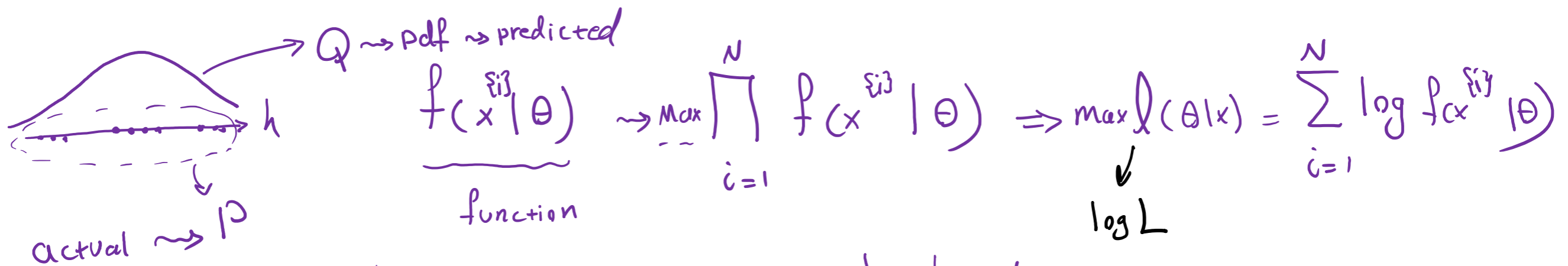
$$H(p, q) = - \sum_{x \in \mathcal{X}} p(x) \log q(x) = H(P) + KL[P][Q]$$

This is because:

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

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

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



$\left(\frac{1}{N} \sum \right) (\cdot) = E[\cdot] \rightsquigarrow \frac{1}{N} \sum g(x) = E[g(x)]$

$\min \log L(\theta|x) = \sum -\log f(x^{i3}|\theta) \Rightarrow \left(\frac{1}{N} \sum \right) -\log f(x^{i3}|\theta) = E[-\log f(x^{i3}|\theta)]$
 $= E[-\log Q]$

$$E[-\log P + \log P - \log Q]$$

$$E[a+b+c] = E[a] + E[b+c]$$

$$E[-\log P] + E\left[\log \frac{P}{Q}\right] \xrightarrow{E[g(x)] = \sum p(x) g(x)}$$

$$\sum p * (-\log p) + \sum p * \log \frac{p}{q} = \sum p (-\log p) - \sum p \log \left(\frac{q}{p}\right)$$

$\sum p \log \frac{1}{p}$ $-KL[P][Q]$

$$H(P) + \left(-\sum p \log \frac{q}{p}\right) = H(P) + KL[P][Q] = CE$$

CE = negative average log likelihood

$$E[g(x)] = \frac{1}{N} \sum (g(x)) = \sum p(x) g(x)$$

Labeling target values

Label encoding (ordinal) and One-hot encoding

$$X = \begin{matrix} \text{height} = h & \text{weight} = w \\ \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} & \begin{bmatrix} 5 \\ 6 \\ 7 \\ 8 \end{bmatrix} \end{matrix} \quad n \times d$$

$$y_a = \begin{bmatrix} \text{cat} \\ \text{fish} \\ \text{dog} \\ \text{cat} \end{bmatrix} \quad n \times 1 = \begin{bmatrix} 0 \\ 1 \\ 2 \\ 0 \end{bmatrix}$$

$$\xrightarrow{ML} \hat{y}_p = \begin{bmatrix} 0.1 \\ 0.5 \\ 2 \\ 0.2 \end{bmatrix}$$

$$\text{cat} = [1 \ 0 \ 0] \quad \text{fish} = [0 \ 1 \ 0] \quad \text{dog} = [0 \ 0 \ 1]$$

$$y_a = \begin{bmatrix} [1 \ 0 \ 0] \\ [0 \ 1 \ 0] \\ [0 \ 0 \ 1] \\ [1 \ 0 \ 0] \end{bmatrix} \xrightarrow{ML} \hat{y}_p = \begin{bmatrix} [23 \ 2 \ 1] \\ [72 \ 120 \ 23] \\ \vdots \\ \{i\} \end{bmatrix} \xrightarrow{\text{Softmax}} \begin{bmatrix} [0.7 \ 0.2 \ 0.1] \\ [0.1 \ 0.8 \ 0.1] \\ \vdots \end{bmatrix}$$

$$y_a * \hat{y}_p = \begin{bmatrix} 0.7 \\ 0.2 \\ \vdots \end{bmatrix}$$

$$\Rightarrow \max \sum_{i=1}^N (y_a * \hat{y}_p) \quad \text{best result } 1$$

$$\min \|y_a - \hat{y}_p\|_2^2 = \sum_{i=1}^N [(1-0.7)^2 + (0-0.2)^2 + (0-0.1)^2 + \dots + \dots] \quad \text{best result is } \underline{\underline{\text{zero}}}$$

$$Y_a = \begin{bmatrix} [1 & 0 & 0] \\ [0 & 1 & 0] \\ [0 & 0 & 1] \\ [1 & 0 & 0] \end{bmatrix}$$

$$\hat{Y}_p = \begin{bmatrix} [0.7 & 0.2 & 0.1] \\ [0.1 & 0.8 & 0.1] \\ [& &] \\ [& &] \end{bmatrix}$$

$$\text{Max} \sum_{i=1}^N (y_a^{(i)} \cdot \hat{y}_p^{(i)T}) = \left(\underbrace{[1 \ 0 \ 0]}_{x^{(1)}} \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix} + \underbrace{[0 \ 1 \ 0]}_{x^{(2)}} \begin{bmatrix} 0.1 \\ 0.8 \\ 0.1 \end{bmatrix} + \dots + \underbrace{\quad}_{x^{(N)}} \right)$$

$$\text{Min} \sum_{i=1}^N \| y_a^{(i)} - \hat{y}_p^{(i)} \|_2$$

Why Cross entropy and not simply use dot product?

Min

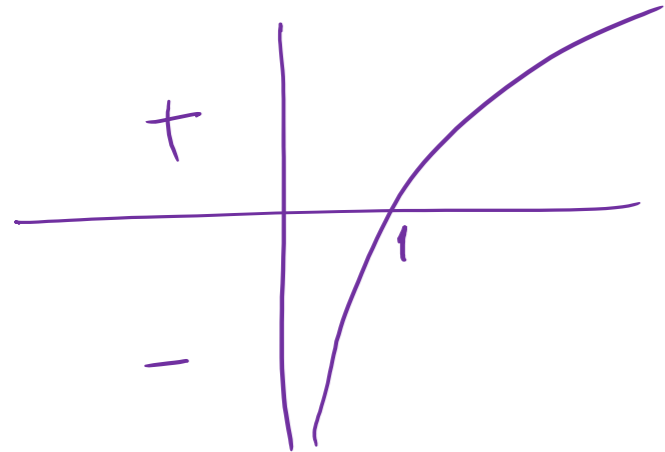
$$CE = - \sum p(x) \log q(x)$$

$$= - \sum y_a \log \hat{y}_p$$

$$p(x) = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$$

$$\log q(x) = \log(\hat{y}_p) =$$

$$= \begin{bmatrix} \log .7 & \log .2 & \log .1 \\ \vdots & \vdots & \vdots \end{bmatrix}$$

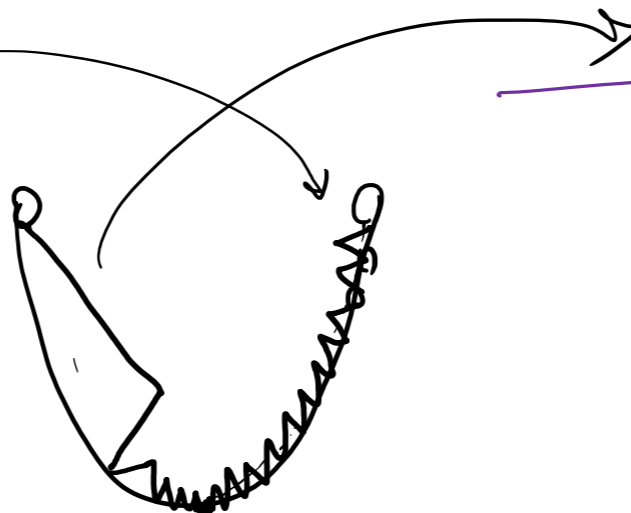
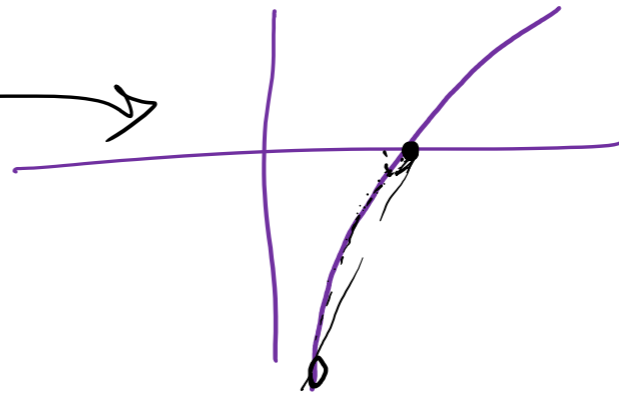
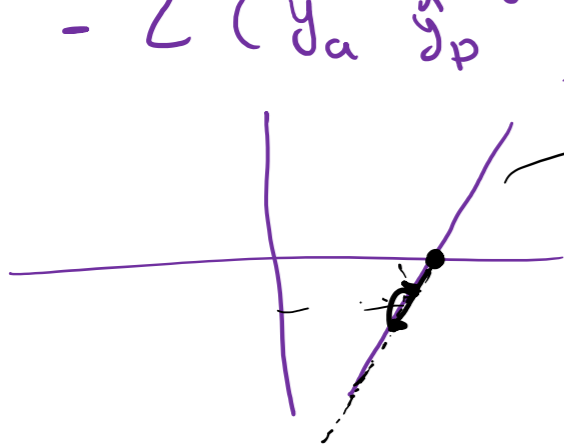


$$CE = - ([1 \log .7 + 0 \log .2 + 0 \log .1] + \dots + \dots)$$

$$CE = - \sum_{i=1}^N y_a^{i} \log \hat{y}_p^{i}$$

$$\text{Min} - \sum (y_a^{i} \hat{y}_p^{i})^T$$

$$CE = - \sum y_a^{i} \log \hat{y}_p^{i}$$



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

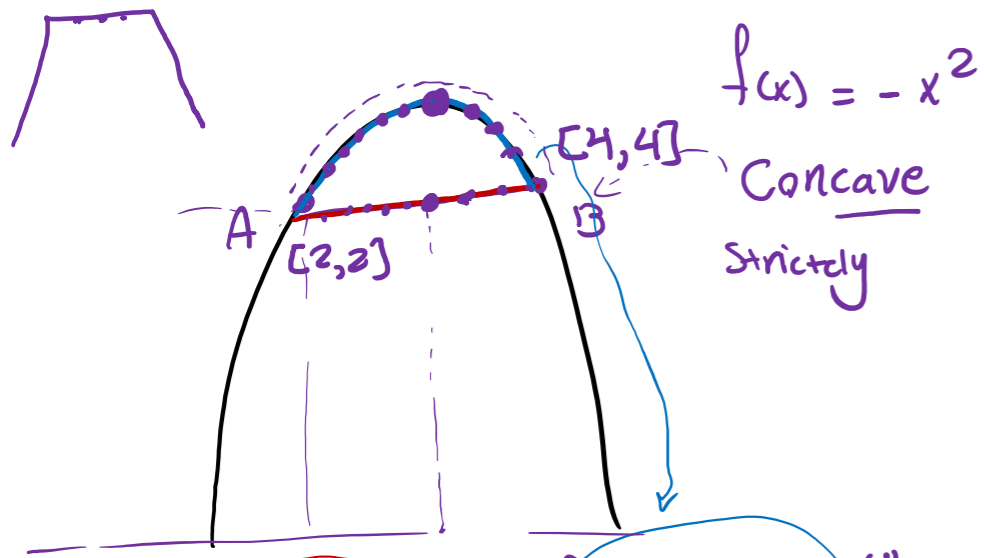
$$-\mathbf{KL}[P||Q] = \sum_s P(s) \log \frac{Q(s)}{P(s)}$$

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

log function is
concave or
convex?

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

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

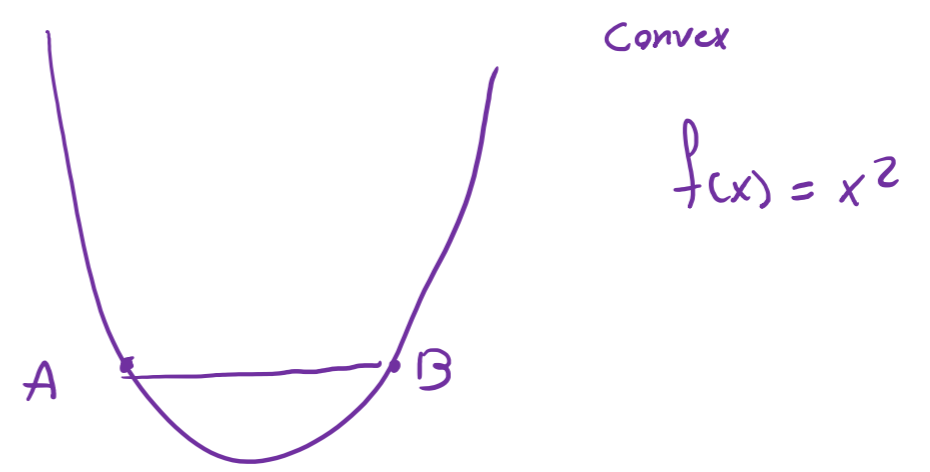


$$E[f(x)] \leq f(E[x])$$

Jensen inequality

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

$$= \underbrace{\frac{1}{2}}_{0.4} f(x=A) + \underbrace{\frac{1}{2}}_{0.6} f(x=B)$$

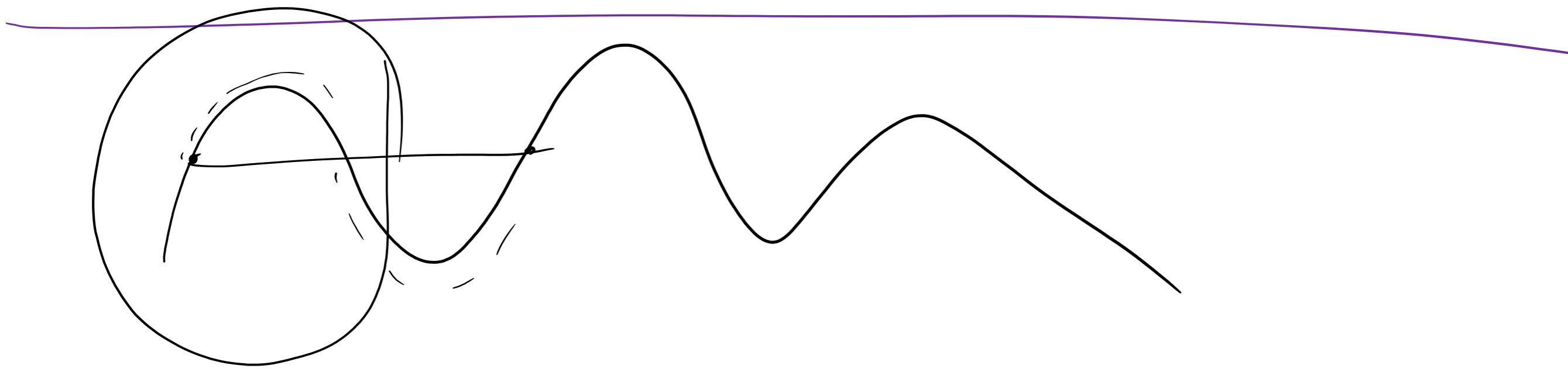


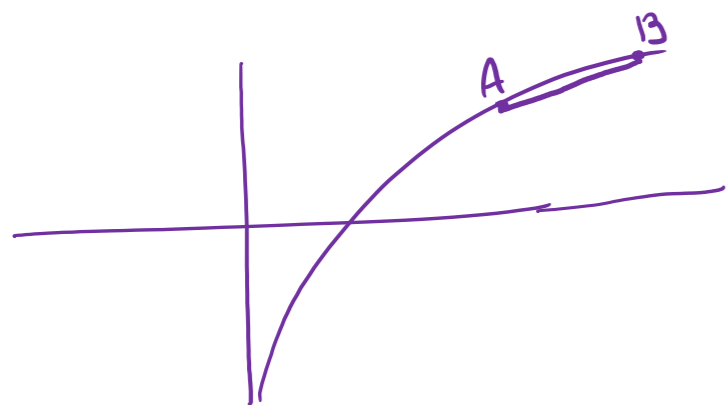
$$E[f(x)] \geq f(E[x])$$

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

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

$$\frac{1}{2} * 2 + \frac{1}{2} * 4 = 3$$





Concave

$$E[f(x)] \leq f(E[x])$$

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

$$f(x) = \log x$$

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

$$\sum p \log g(x) = E[\log g(x)] \leq \log(E[g(x)])$$


$$\leq \log \left(\sum Q(x) \right) = 1$$

$$\leq \log 1$$

$$-KL[P][Q] \leq 0$$

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

$$\begin{aligned} E[g(x)] &= \sum P(x) g(x) \\ &= \sum P(x) \frac{Q(x)}{P(x)} \end{aligned}$$

Objective function $f(x, y)$ 

s.t. $g(x, y) = 0$ Equality constraint \rightsquigarrow Lagrange function

s.t. $h(x, y) > 2$ Inequality constraint \rightsquigarrow KKT conditions

① Linear Programming

$$\rightsquigarrow f(x, y) = x + y$$

s.t. $x + y = 20$

② Quadratic Programming

$$\rightsquigarrow f(x) = x^2 + y$$

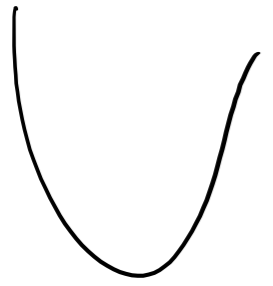
s.t. $x + y = 30$

③ Non linear Programming


$$\rightsquigarrow f(x) = x^3 + y$$

s.t. $x^2 + y = 20$

$$H = \begin{bmatrix} \frac{\partial^2 f(x, y)}{\partial x^2} & \frac{\partial^2 f(x, y)}{\partial x \partial y} \\ \frac{\partial^2 f(x, y)}{\partial y \partial x} & \frac{\partial^2 f(x, y)}{\partial y^2} \end{bmatrix}$$



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



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

$M = \#$ hours you study ML / day

$S = \#$ " " sleep

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

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

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

λ, α, β → Lagrange multiplier

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

Objective

$$L(M, S, \lambda) = f(M, S) - \lambda g(M, S) = 6M^2 + 3S^2 - \lambda (M + S - 24)$$

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

$$\frac{\lambda}{12} + \frac{\lambda}{6} = 24 \Rightarrow \lambda = 96$$

$$\frac{\partial L}{\partial M} = 0 \Rightarrow 12M + 0 - \lambda = 0$$

$$\Rightarrow M = \frac{\lambda}{12} = \frac{96}{12} = 8$$

$$\frac{\partial L}{\partial S} = 0 \Rightarrow 6S - \lambda = 0$$

$$\Rightarrow S = \frac{\lambda}{6} = \frac{96}{6} = 16$$

$$M + S = 24$$

$$8 + 16$$

$$\min f(m, s)$$

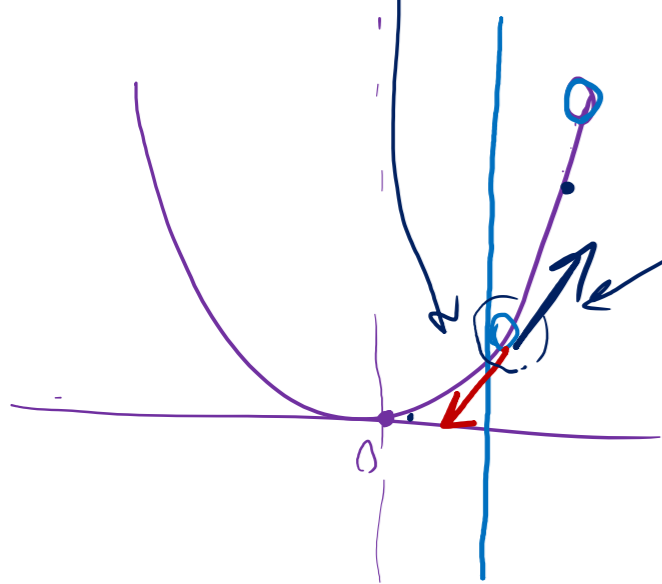
$$\text{s.t. } m + s = 20 \rightsquigarrow g(m, s) = m + s - 20 \rightsquigarrow \lambda_1$$

$$\text{s.t. } m - s = 10 \rightsquigarrow h(m, s) = m - s - 10 \rightsquigarrow \lambda_2$$

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

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

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



Min $f(M, S) = 6M^2 + 3S^2$
 $M + S \leq 24$ \rightsquigarrow KKT condition \rightsquigarrow
 $\rightsquigarrow g(M, S) = M + S - 24 \leq 0$

$g(M, S) = 0$ Active Solution
 $g(M, S) < 0$ Inactive Solution

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

① Stationary condition $\frac{\partial L}{\partial \cdot} = 0$

② Primal feasibility $g(M, S) \leq 0$

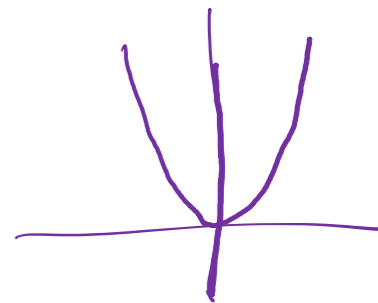
③ Dual feasibility $\lambda \geq 0$

④ Complementary slackness

$\lambda g(M, S) = 0 \implies$

$g(M, S) = 0 \rightsquigarrow \lambda \geq 0$

$\lambda = 0 \rightsquigarrow g(M, S) \neq 0$



$$\min f(M, S) = \frac{1}{2}M^2 + \frac{1}{2}S^2 \rightsquigarrow \text{Convex} \rightsquigarrow \text{minimize}$$

Primal form

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

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


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

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

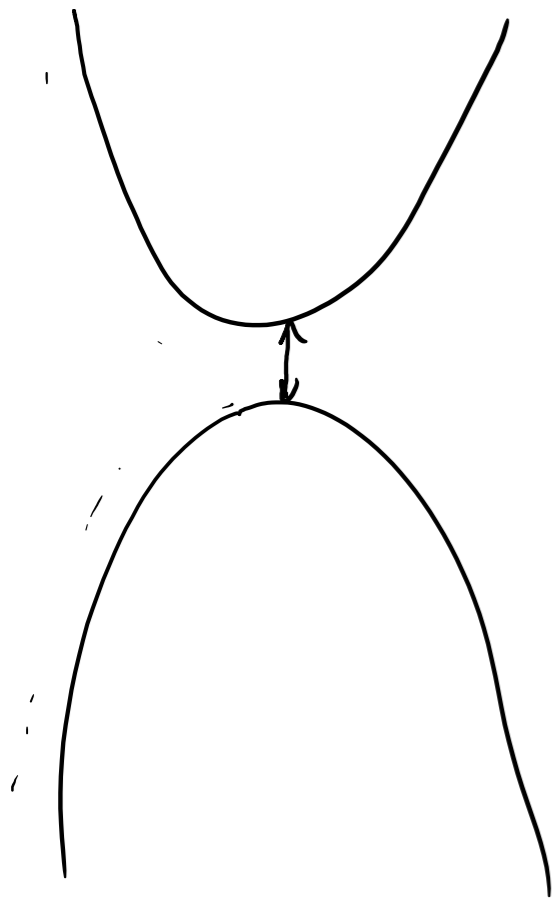
x^2 

$$L(\lambda) = \frac{\lambda^2}{2} + \frac{\lambda^2}{2} - \lambda(2\lambda - 24) = \lambda^2 - 2\lambda^2 + 24\lambda = -\lambda^2 + 24\lambda$$

Dual form

$-x^2$ 

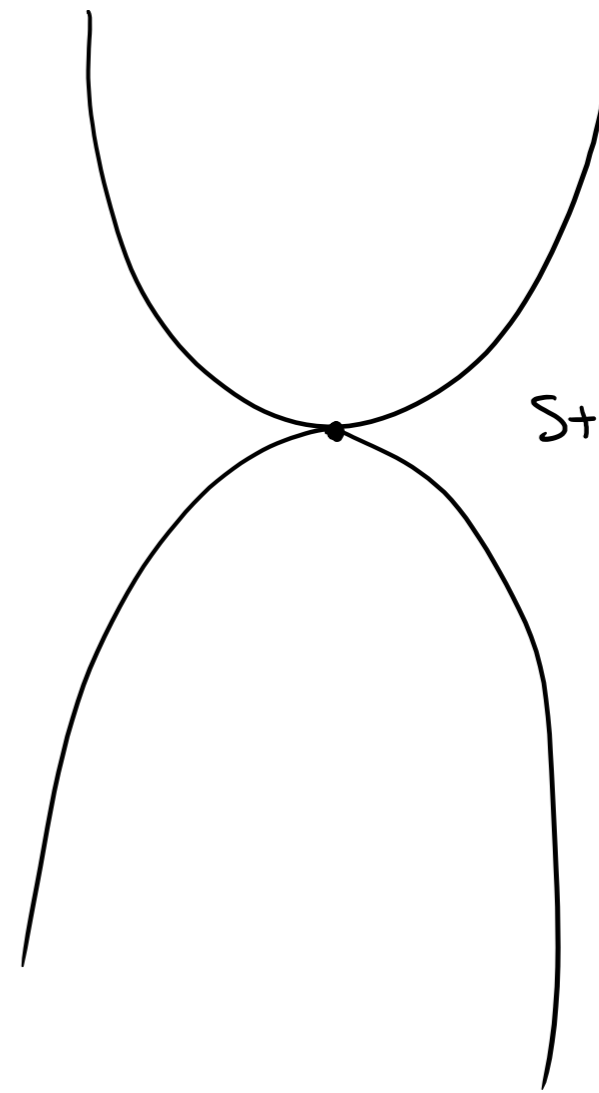
Concave



Convex
Primal

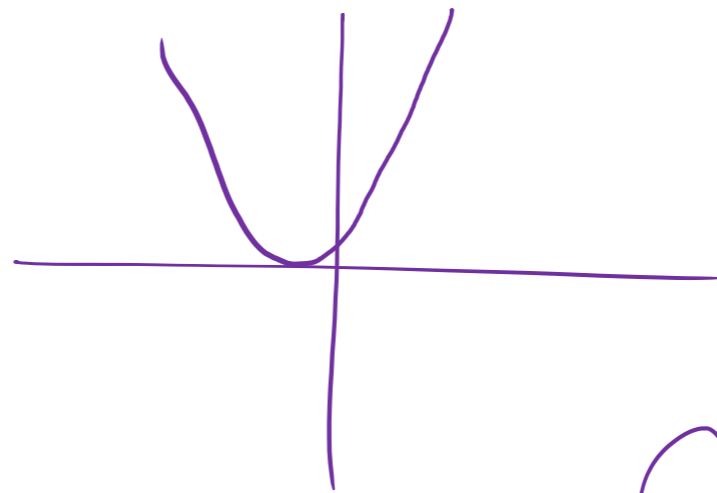
weak duality

Dual form



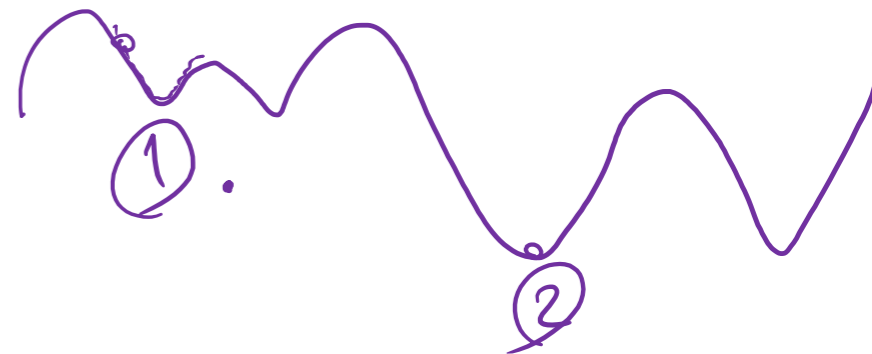
Strong duality

$$f(x) = \exp(x) + x^2$$

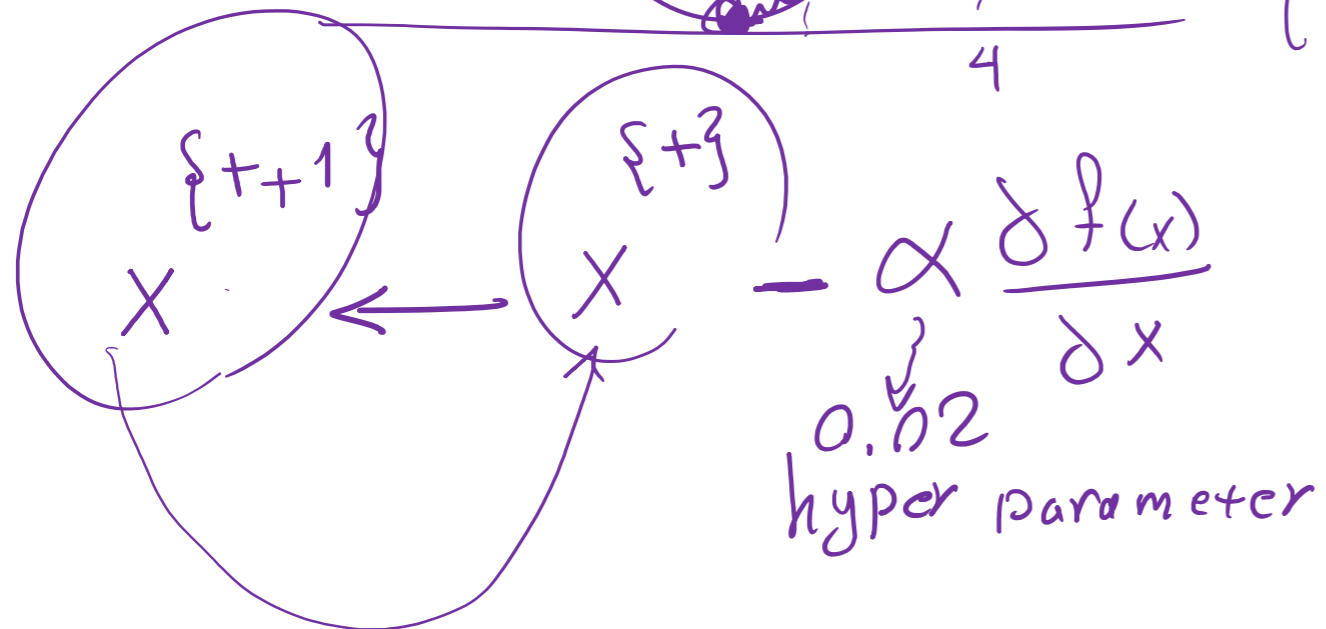
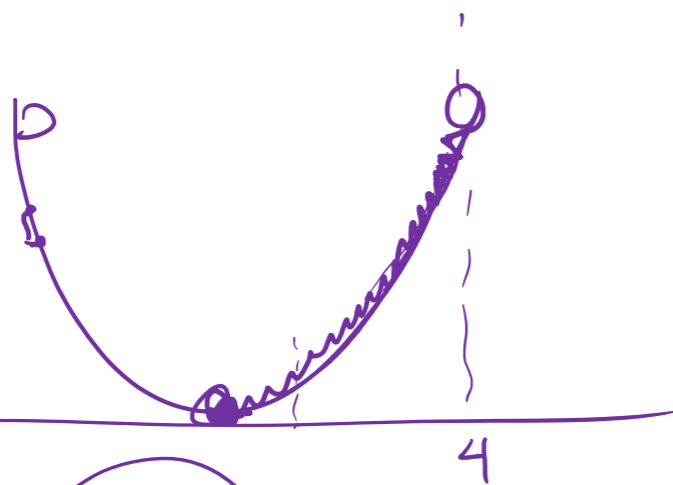


$$\frac{\partial f(x)}{\partial x} = 0 = \exp(x) + 2x = 0$$

$$\Rightarrow x = .$$



GD



GA

$$4 - 0.02 (\exp(4) + 2 \times 4)$$

① Quotient rule: $\left(\frac{f}{g}\right)' = \frac{f'g - fg'}{g^2}$

② Chain rule: $(u)^n' = nu' u^{n-1}$

③ Product rule: $(fg)' = f'g + fg'$

1	0	-1
1	0	-1
1	0	-1

$k = \text{kernel}$

*

0	1	2
2	1	3
1	2	1

$S = \text{sub matrix}$

=

0	0	-2
2	0	-3
1	0	-1

Sum (

) =

$$0 + 0 + -2 + 2 + \dots = -3$$

0	1	2	4	5	6	0	2
2	1	3	8	9	255	72	83
1	2	1	4	79	65	53	33
43	97	15	67	104	77	163	43
36	173	13	76	205	89	179	34
63	163	113	86	209	91	185	98
84	153	123	96	134	101	196	121
96	143	133	79	135	103	216	211

	-3						

$k = [1 \ 0 \ -1 \ 1 \ 0 \ -1 \ 1 \ 0 \ -1]$
 $S = [0 \ 1 \ 2 \ 2 \ 1 \ 3 \ 1 \ 2 \ 1]$
 $k \cdot S$

