

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

$$X = \begin{bmatrix} h & w & \text{age} = a & \dots \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}_{n \times d}$$

$$Y = \begin{bmatrix} \text{Cat} \\ \text{dog} \\ \text{fish} \\ \text{Cat} \\ \vdots \end{bmatrix}_{n \times 1} \xrightarrow{\text{actual}} \begin{bmatrix} 1 \\ 2 \\ 3 \\ 1 \\ \vdots \end{bmatrix}$$

$$\xrightarrow{\text{ML}} \hat{Y} = \begin{bmatrix} 1 \\ 2 \\ 2.5 \\ 1 \\ \vdots \end{bmatrix}_{\text{Predicted}}$$

One-hot encoding \rightarrow Pdf

$$Y = \begin{bmatrix} [1 & 0 & 0] \\ [0 & 1 & 0] \\ [0 & 0 & 1] \\ [1 & 0 & 0] \\ \vdots \end{bmatrix}_{\text{actual}} \xrightarrow{\text{ML}} \text{Softmax} \rightarrow \hat{Y} = \begin{bmatrix} [0.8 & 0.1 & 0.1] \\ [0.3 & 0.6 & 0.1] \\ \vdots \end{bmatrix}_{\text{Pdf}}$$

①

$$\|Y - \hat{Y}\|_2 = \sqrt{(1-0.8)^2 + (0-0.1)^2 + (0-0.1)^2} + \sqrt{(0-0.3)^2 + (1-0.6)^2 + \dots} + \dots \quad \text{Minimize}$$

②

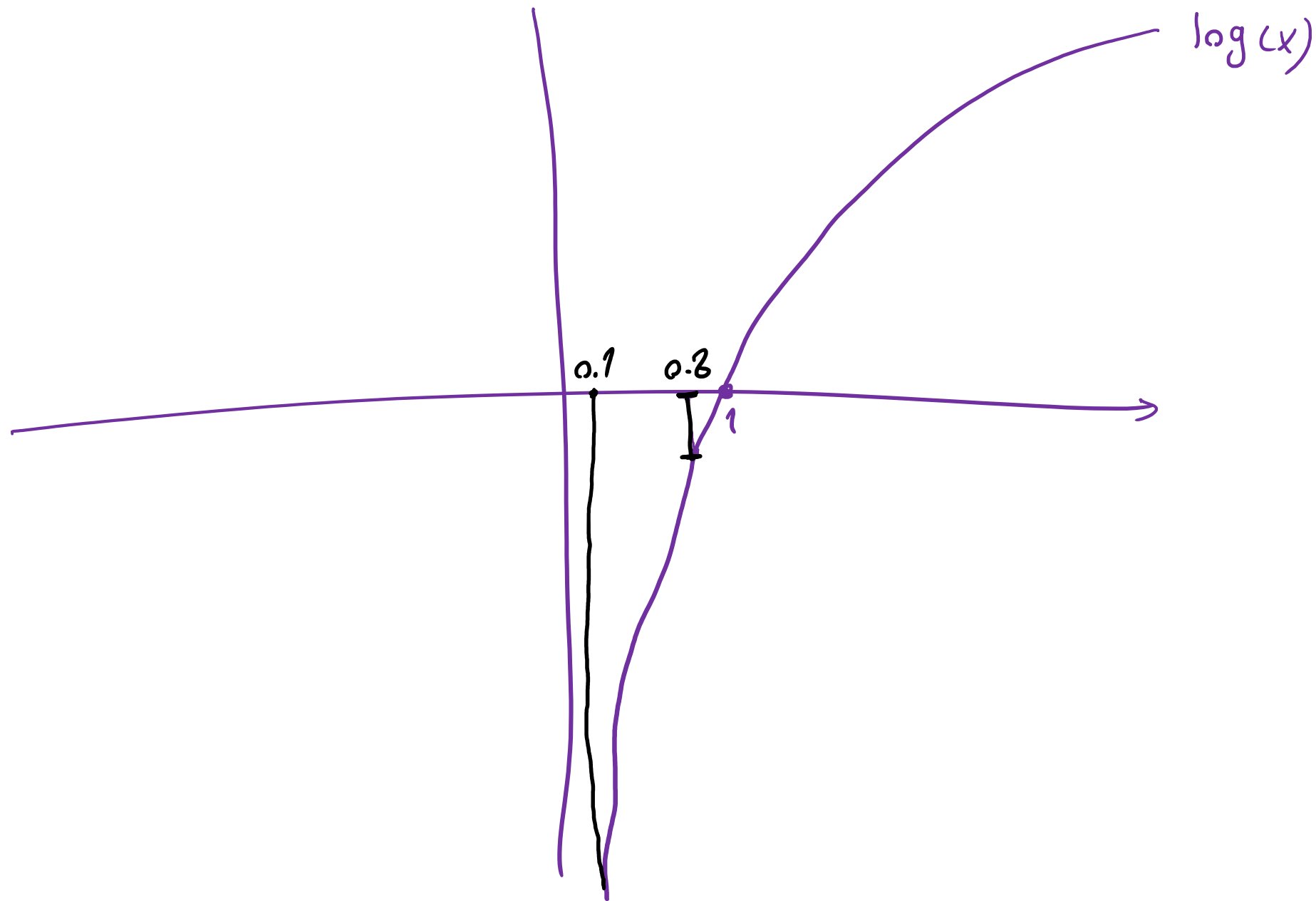
$$\sum |Y \cdot \hat{Y}| = \begin{bmatrix} 0.8 \\ 0.6 \\ \vdots \end{bmatrix} = \dots \checkmark \quad \text{maximize}$$

Why Cross entropy and not simply use dot product?

$$CE = - \sum P(x) \log Q(x)$$

Minimization

$$CE = - (1 \times \log 0.8 + 0 \times \log 0.1 + 0 \times \log 0.1) - (0 \times \log 0.3 + \dots) - \dots$$



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

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

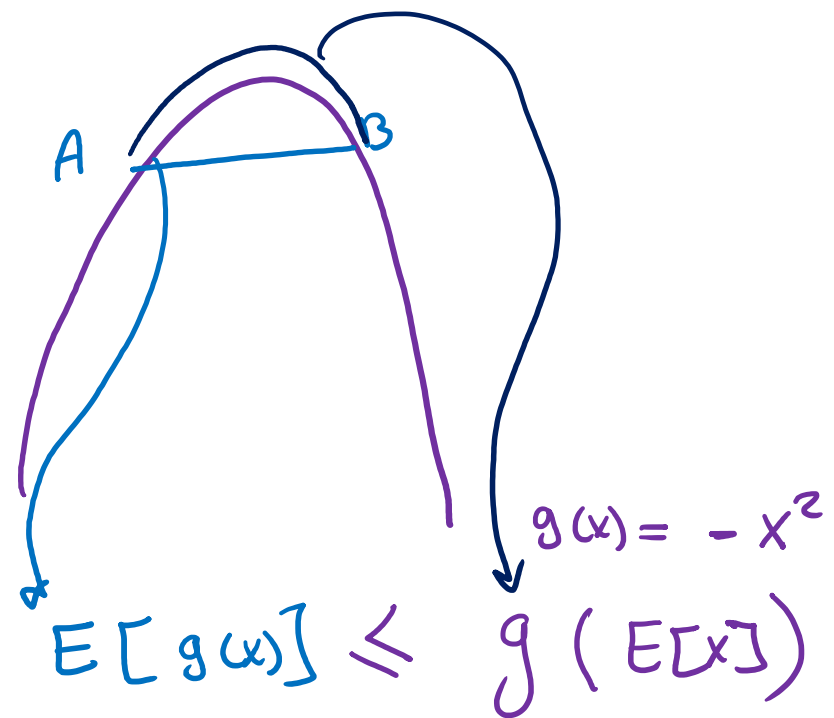
log function is
concave or
convex?

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

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

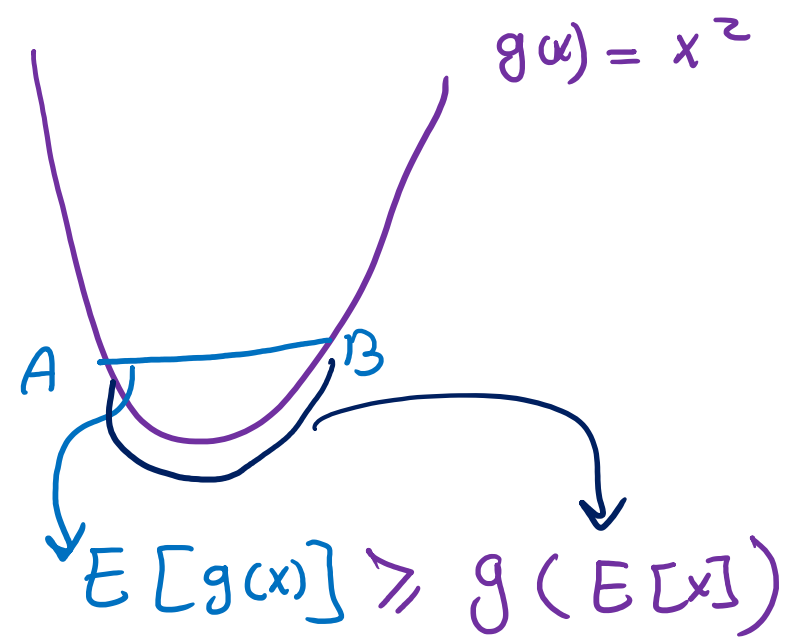
When $P = Q$, $KL[P||Q] = 0$

Concave



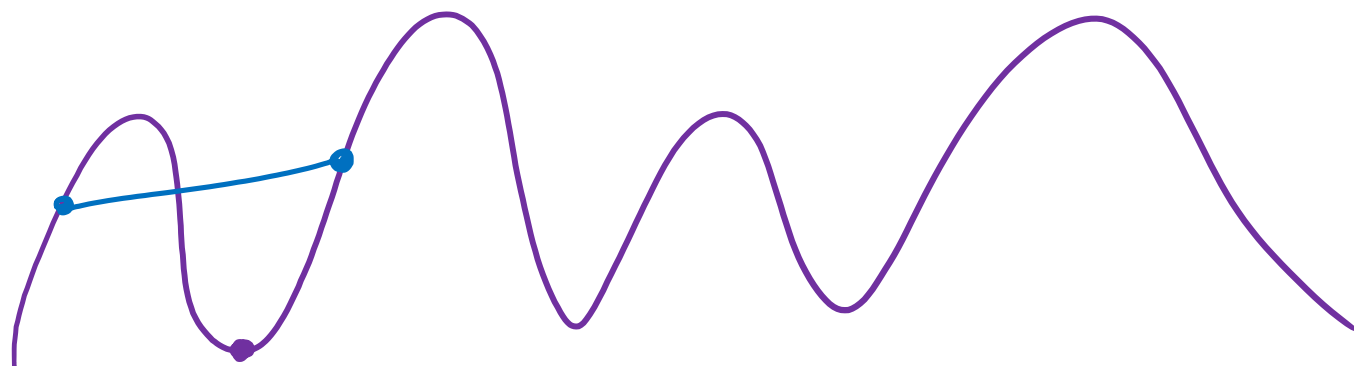
Jensen Inequality

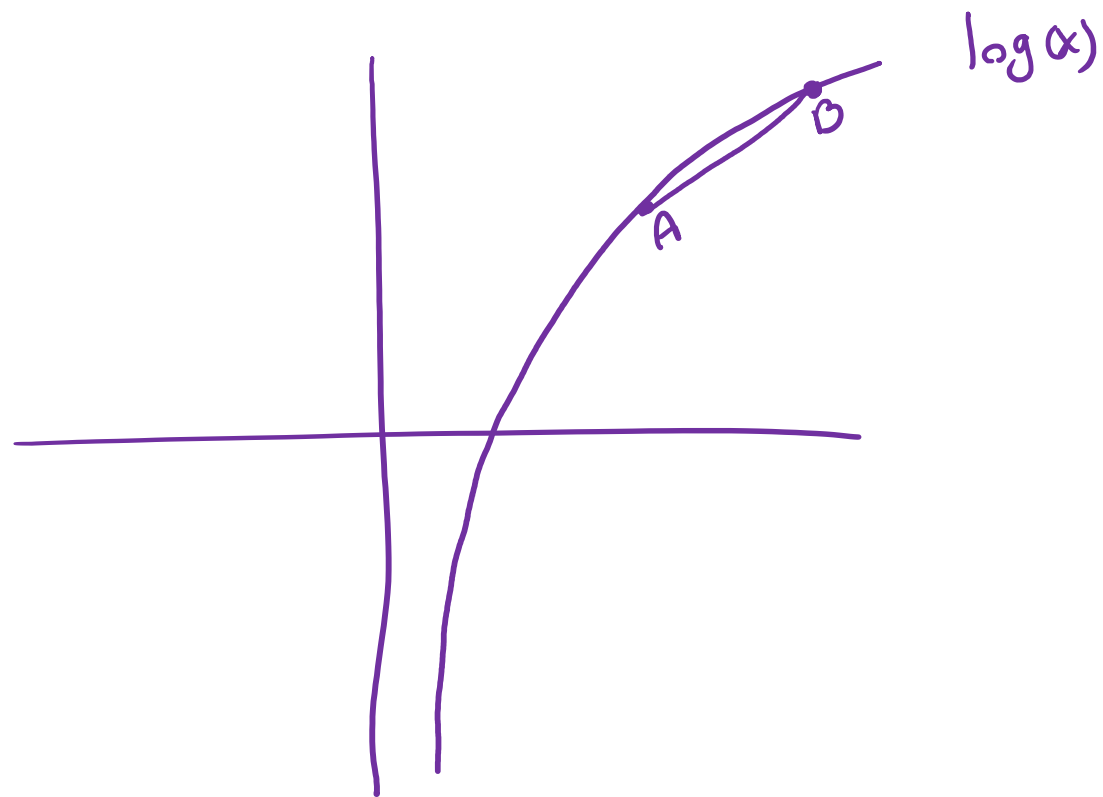
Convex



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

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



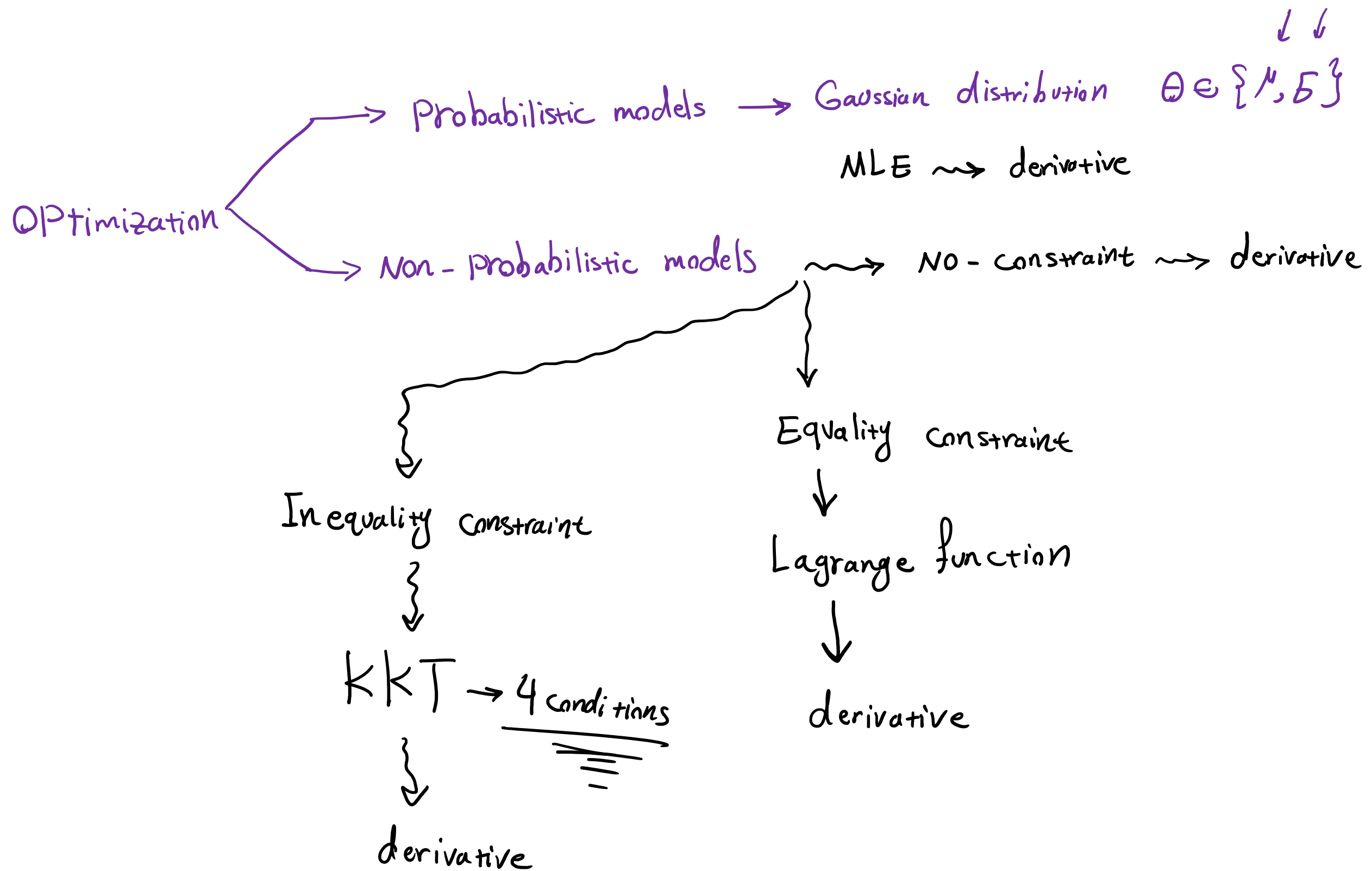


$$E[\log(x)] \leq \log(E[x])$$

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

$$\begin{aligned} -KL[P][Q] &= E[\log g(x)] \leq \log(E[g(x)]) \\ &\leq \log\left(\sum P(x) g(x)\right) \\ &\leq \log\left(\sum \cancel{P(x)} \frac{Q(x)}{\cancel{P(x)}}\right) \\ &\leq \log \sum Q(x) = \log 1 = 0 \end{aligned}$$

$$-KL[P][Q] \leq 0 \text{ or } KL[P][Q] \geq 0$$



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

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

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

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

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

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

$$\text{s.t. } m - s = 8$$

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

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

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

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

$$L(m, s, \lambda_1, \lambda_2) =$$

$$f(m, s) - \lambda_1 g_1(m, s) - \lambda_2 g_2(m, s)$$

$$\lambda, \alpha, \beta, m$$

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

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

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

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

$$\lambda = 96$$

$$m + s = 24$$

$$8 + 16 = 24$$

$$f(m,s) = 6m^2 + 3s^2$$

$$m+s \leq 24$$

KKT conditions

$$m+s-24 \leq 0$$

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

$$g(m,s) \leq 0$$

$$L(m,s,\delta) = f(m,s) + \delta g(m,s)$$

① Stationarity condition

$$\nabla L(m,s,\delta) = 0$$

② Primal feasibility

$$g(m,s) \leq 0$$

③ Dual feasibility

$$\delta \geq 0$$

④ Complementary slackness

$$g(m,s) \neq 0$$

$$\delta = 0$$

$$g(m,s) = 0$$

$$\delta \geq 0$$

Example 1:

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

Example 2:

<https://www.geogebra.org/3d/sy8qpk7>