

# 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  $P$

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

This is because:

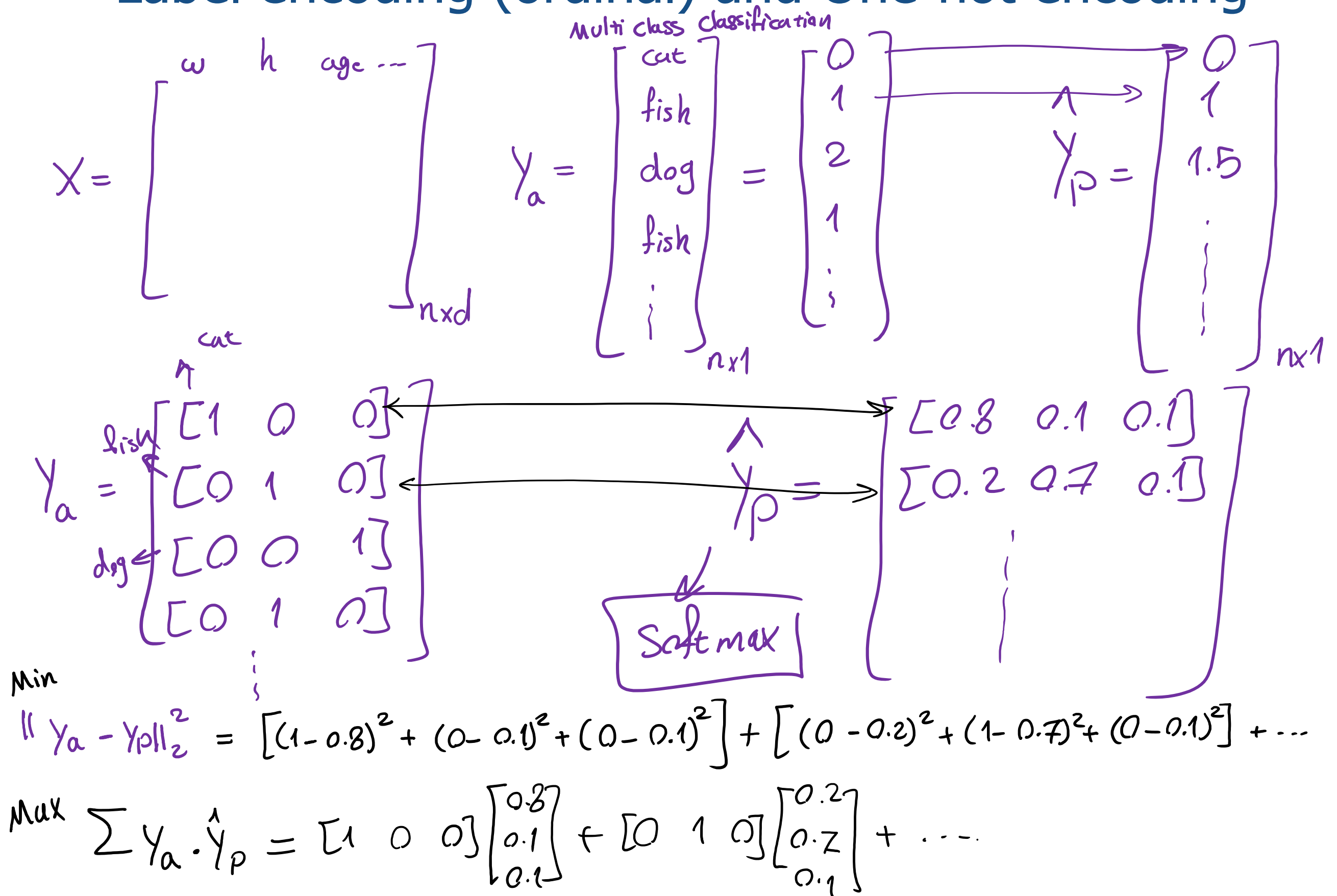
$$H(p, q) = \mathbb{E}_p[l_i] = \mathbb{E}_p \left[ \log \frac{1}{q(x_i)} \right]$$

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

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

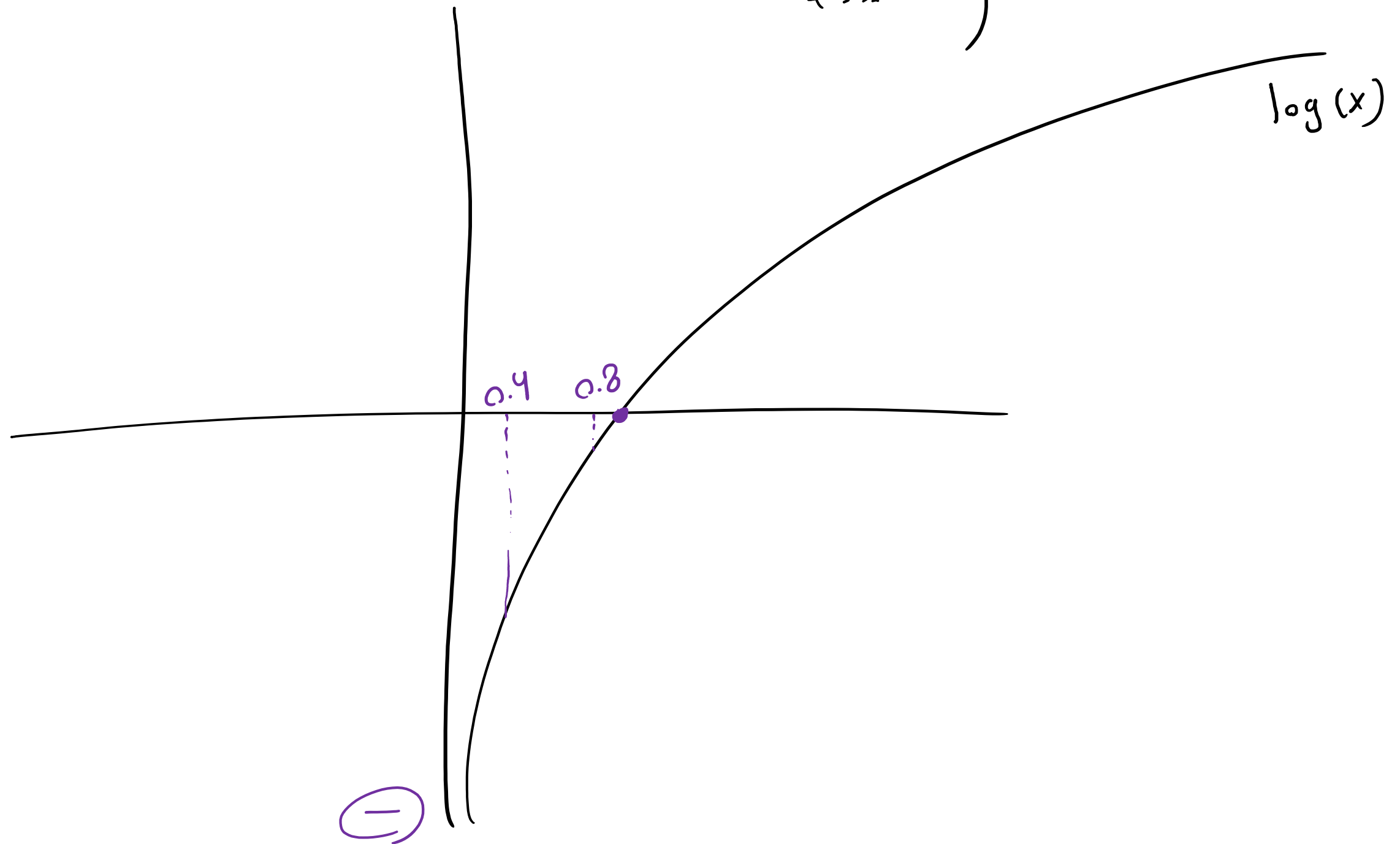
# Labeling target values

## Label encoding (ordinal) and One-hot encoding



# Why Cross entropy and not simply use dot product?

$$H(p, q) = - \sum p(x) \log q(x) = - \left( [1 \ 0 \ 0] \begin{bmatrix} \log 0.8 \\ \log 0.1 \\ \log 0.1 \end{bmatrix} + [0 \ 1 \ 0] \begin{bmatrix} \log 0.2 \\ \log 0.7 \\ \log 0.1 \end{bmatrix} + \dots \right)$$



# 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?

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

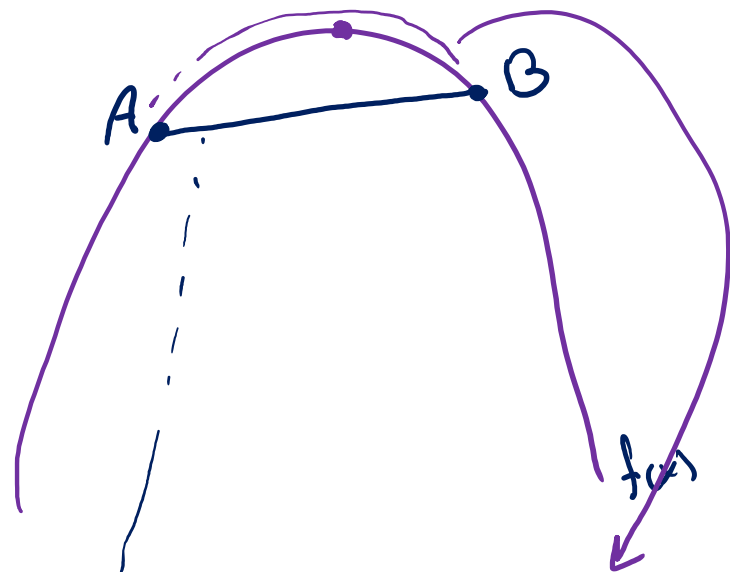
So  $\mathbf{KL}[P||Q] \geq 0$ . Equality iff  $P = Q$

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

Concave

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

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

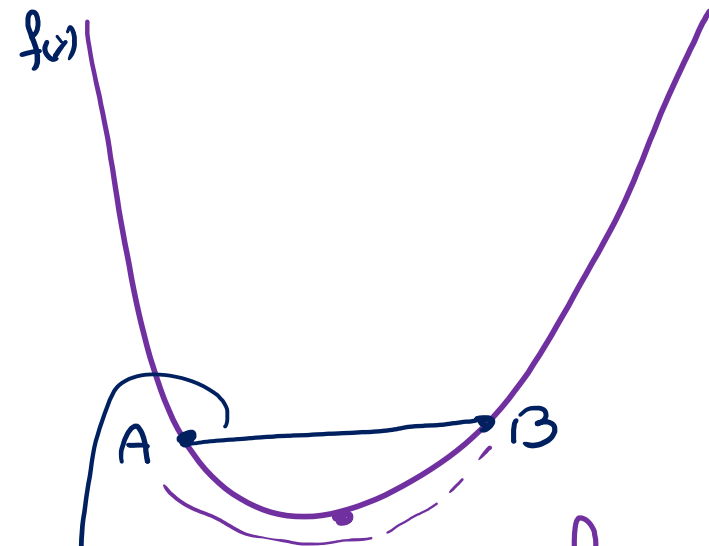


$$E[f(x)] \leq f(E[x])$$

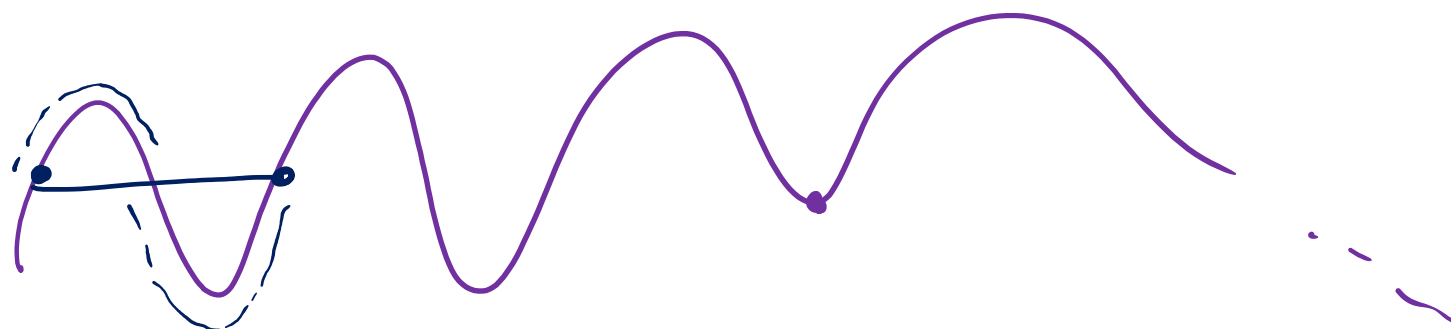
convex

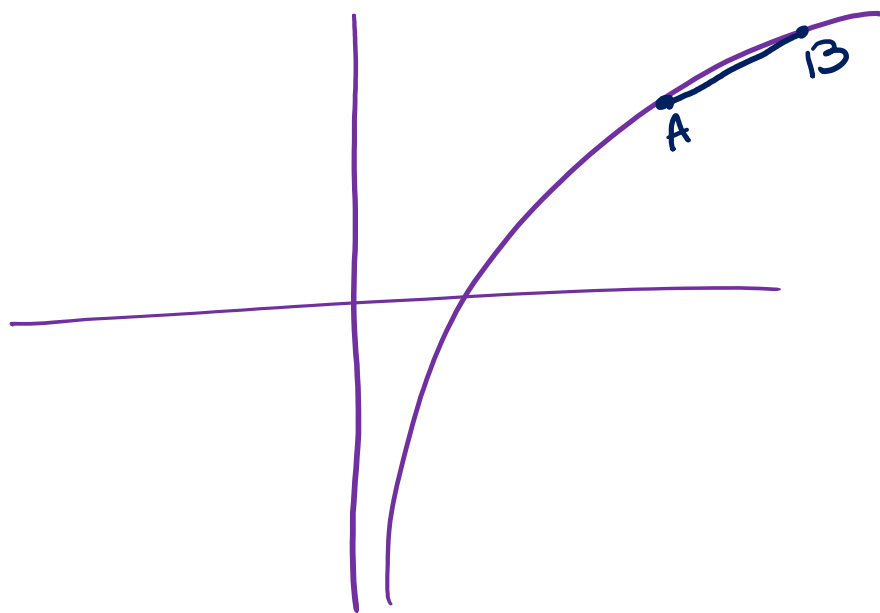
$$f(x) = x^2$$

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



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





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

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

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

$$\frac{Q(s)}{P(s)} = g(s)$$

---


$$-KL[P][Q] = \sum P(s) \log \frac{Q(s)}{P(s)} = \sum \underbrace{P(s)} \underbrace{\log g(s)}$$

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

$$= \leq \log \sum P(s) g(s)$$

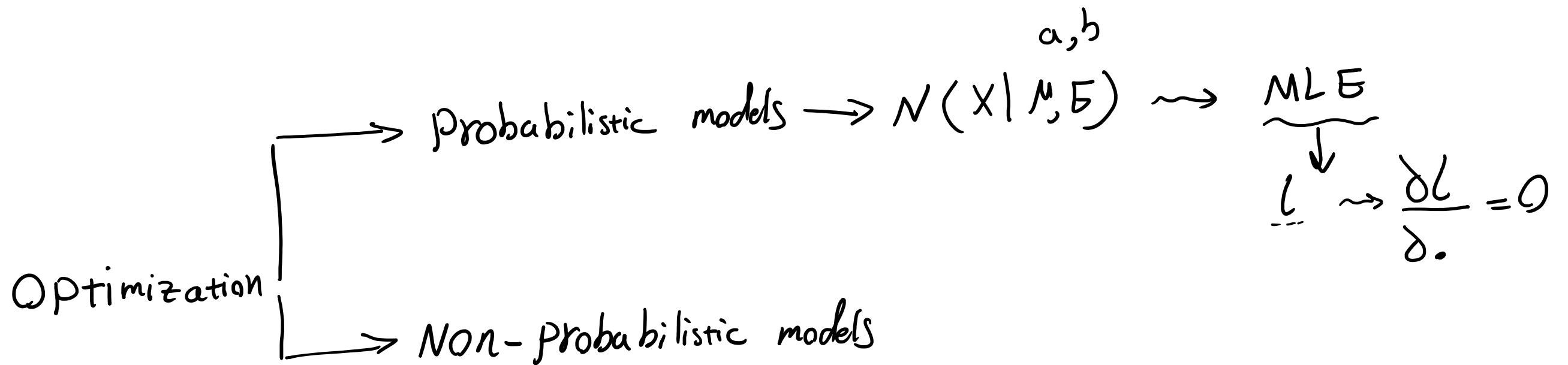
$$= \leq \log \sum \cancel{P(s)} \frac{Q(s)}{\cancel{P(s)}}$$

$$= \leq \log (\sum Q(s))$$

$$-KL[P][Q] \leq \log 1$$

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





No - constraint

$$f(x, y) = x^2 + y^2$$

$$\frac{\partial f(x, y)}{\partial x} = 0$$

$$\frac{\partial f(x, y)}{\partial y} = 0$$

Objective function

Equality constraint

$$f(x, y) = x^2 + y^2$$

constraint function

$$\text{s.t. } x - y = 8$$

$$g(x, y) = x - y - 8$$

↓  
Lagrange function

$$\frac{\partial \mathcal{L}}{\partial \theta} = 0$$

Inequality constraint

$$f(x, y) = x^2 + y^2$$

$$\text{s.t. } x - y \leq 8$$

$$\frac{\partial \mathcal{L}}{\partial \theta} = 0$$

& satisfy KKT conditions

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

$$f'' = \begin{bmatrix} \frac{f''}{\delta^2 M} & \frac{f''}{\delta M \delta S} \\ \frac{f''}{\delta S \delta M} & \frac{f''}{\delta^2 S} \end{bmatrix}$$

$M$ : # hours study ML  
 $S$ : # " sleep

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

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

$f(M, S) = 6M^2 + 3S^2$  → objective function  
 $M + S = 24$   
 $M - S = 10$   
 $\Rightarrow$   $g(M, S) = M + S - 24 = 0$  → constraint function  
 $h(M, S) = M - S - 10$   
 $L(M, S, \lambda_1, \lambda_2) = f(M, S) - \lambda_1 g(M, S) - \lambda_2 h(M, S)$

$\lambda, \alpha, \beta$

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

s.t.  $M + S = 24$

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

$$\frac{\partial L}{\partial M} = 0 \Rightarrow 12M - \lambda = 0 \Rightarrow M = \frac{\lambda}{12} = \frac{96}{12} = 8$$

$$\frac{\partial L}{\partial S} = 0 \Rightarrow 6S - \lambda = 0 \Rightarrow S = \frac{\lambda}{6} = \frac{96}{6} = 16$$

$$M + S = 24$$

$$8 + 16 = 24$$

$$\frac{\partial L}{\partial \lambda} = 0 \Rightarrow -(M + S - 24) = 0 \Rightarrow M + S = 24 \Rightarrow \frac{\lambda}{12} + \frac{\lambda}{6} = 24 \Rightarrow \lambda = 96$$

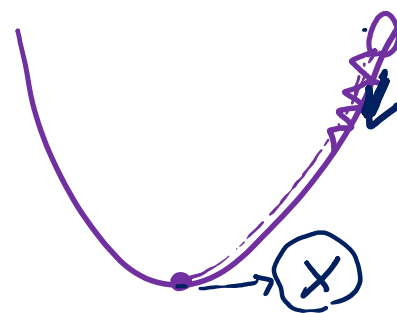
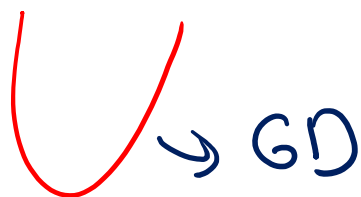
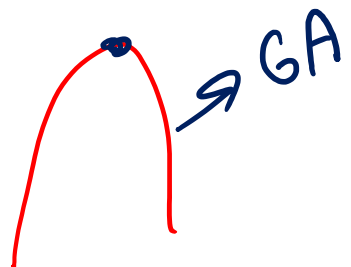
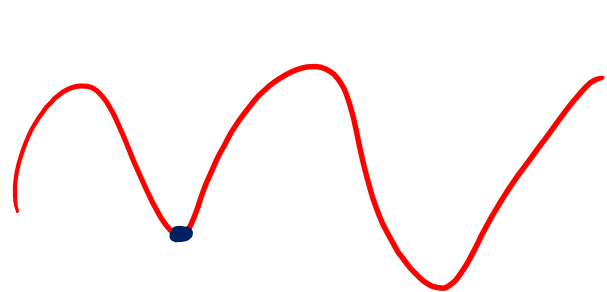
$$L(\mu, \lambda, \gamma) = f(\mu, \lambda) - \gamma g(\mu, \lambda)$$

$$\begin{aligned} \nabla f(\mu, \lambda) &\approx \nabla g(\mu, \lambda) \\ \nabla f(\mu, \lambda) &= \gamma \nabla g(\mu, \lambda) \end{aligned}$$

$$\nabla_{\mu, \lambda, \gamma} L(\mu, \lambda, \gamma) = 0$$

$$\nabla (f(\mu, \lambda) - \gamma g(\mu, \lambda)) = 0$$

$$\nabla f(\mu, \lambda) - \cancel{\gamma} \nabla g(\mu, \lambda) = 0$$



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

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

$$x^{t+1} \leftarrow x^t - \alpha \frac{\partial f(x)}{\partial x} \rightarrow \text{GD}$$

learning step = 0.002

$$x^{t+1} \leftarrow x^t - \alpha 2x^t$$

$$L(\mu, \sigma, \delta) = f(\mu, \sigma) - \delta g(\mu, \sigma) \quad \mu, \sigma, \delta$$

$$\mu^{t+1} \leftarrow \mu^t - \alpha \frac{\partial L}{\partial \mu}$$

$$\sigma^{t+1} \leftarrow \sigma^t - \alpha \frac{\partial L}{\partial \sigma}$$

$$\delta^{t+1} \leftarrow \delta^t - \alpha \frac{\partial L}{\partial \delta}$$

$$\begin{aligned} \min f(M, S) &= 6M^2 + 3S^2 \\ \text{s.t. } M+S &\leq 24 \Rightarrow \underbrace{M+S-24}_{g(M, S)} \leq 0 \end{aligned}$$

KKT conditions

$$\textcircled{1} \text{ Stationary condition } \Rightarrow L = \underbrace{6M^2 + 3S^2}_{\min} + \lambda (M+S-24)$$

$$\textcircled{2} \text{ Primal feasibility } g(M, S) \leq 0$$

$$\textcircled{3} \text{ Dual feasibility } \lambda \geq 0$$

$$\textcircled{4} \text{ Complementary slackness } g(M, S) \lambda = 0$$

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

$\swarrow$   
 $\searrow$

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

$\lambda = 0$   $\leftarrow$  Inactive solution

# Using Gradient Descent **as an alternative** to solve the equality constraint optimization example

```
import numpy as np

def minimize_gd():
    LEARNING_RATE = 0.01
    TOLERANCE = 1e-6
    # Initialize M, S, and lambda (make sure M + S = 24)
    M = 12.0
    S = 12.0
    lm = 0.0 # Lagrange Multiplier

    while True:
        # Calculate the gradients of the Lagrangian w.r.t. M, S, and lambda
        gradientM = 12 * M - lm
        gradientS = 6 * S - lm
        gradientLambda = - (M + S - 24)

        # Update M, S, and lambda
        newM = M - LEARNING_RATE * gradientM
        newS = S - LEARNING_RATE * gradientS
        newLambda = lm - LEARNING_RATE * gradientLambda

        # If the changes in M, S, and lambda are smaller than the tolerance, we break the loop
        if np.abs(newM - M) < TOLERANCE and np.abs(newS - S) < TOLERANCE and np.abs(newLambda - lm) < TOLERANCE:
            break

        M = newM
        S = newS
        lm = newLambda

    print("Minimum occurs at M = ", M , ", S = " , S , ", with lambda = ", lm)
    print("Minimum value of z = " , (6*M * M + 3*S * S))

minimize_gd()
```

This won't converge; do you know why?

Example 1:

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

Example 2:

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



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

$$\text{Sum} \left( \begin{bmatrix} 0 & 0 & -2 \\ 2 & 0 & -3 \\ 1 & 0 & -1 \end{bmatrix} \right) = |-3|$$

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						

$$\begin{bmatrix} 1 & 0 & -1 & 1 & 0 & -1 \\ 0 & 1 & 2 & 2 & 1 & 3 \end{bmatrix} \quad \begin{bmatrix} 1 & 0 & -1 \\ 1 & 2 & 1 \end{bmatrix}$$