



Machine Learning CS 4641

Optimization

Nimisha Roy

Lecturer, SCI

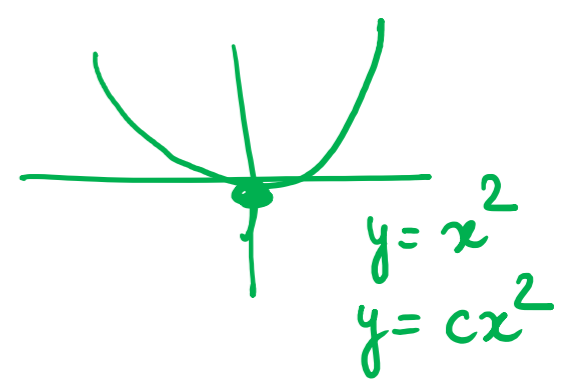
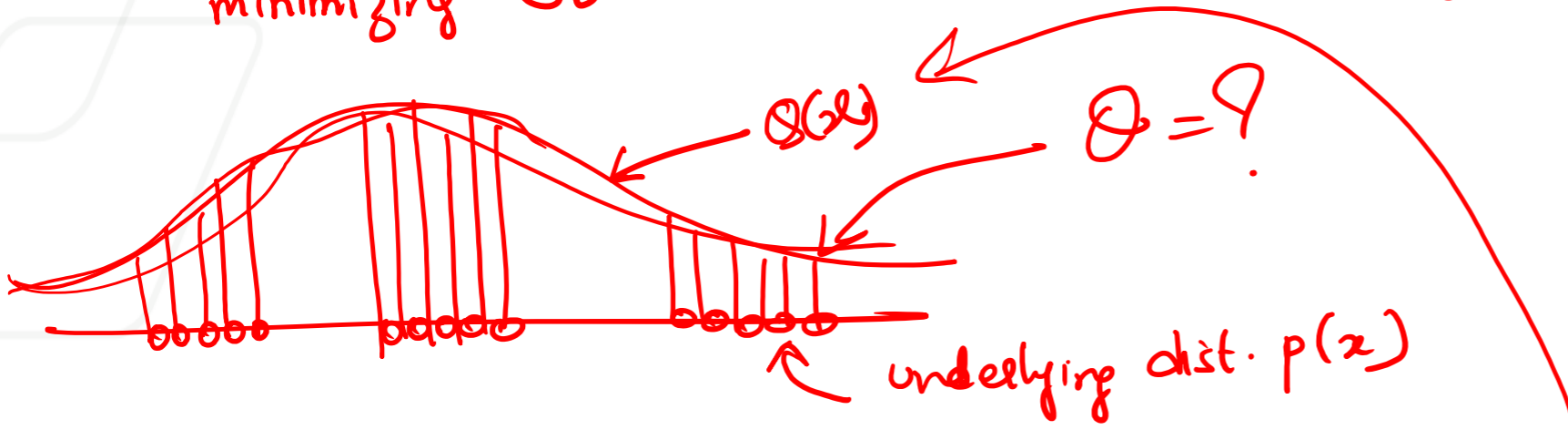
Director of Online Undergraduate Initiatives

College of Computing



$$CE = \sum_x p(x) \log \frac{1}{q(x)}$$

minimizing CE is the same as minimizing -ve log likelihood.



$$L(\theta | x) = \prod_i f(x^{i3} | \theta)$$

$$\log L(\theta | x) = \ell(\theta | x) = \sum_i \log f(x^{i3} | \theta)$$

$$= \min - \sum_i \log f(x^{i3} | \theta)$$

$$= \min \left(-\frac{1}{n} \sum_i \log f(x^{i3} | \theta) \right)$$

$$= \min -E[\log f(x | \theta)]$$

$$= \min -E[\log Q(x)]$$

minimizing -ve log likelihood = $\min -E[\log Q(x)]$

$$\min -E[\log P(x) - \log P(x) + \log Q(x)]$$

$$\min -E[\log P(x)] - E\left[\log \frac{Q(x)}{P(x)}\right]$$

$$\min \left(-\sum_x p(x) \log p(x) \right) \left(-\sum_x p(x) \log \frac{q(x)}{p(x)} \right)$$

$$\min. \left(H(P) + KL[P](Q) \right)$$

$\min CE$

$$\min. -\log \text{likelihood} = \min. CE$$

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

$$E[g(x)] \\ = \sum_x p(x)g(x)$$

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:

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

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

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

Labeling target values: Label encoding (ordinal) and One-hot encoding

• **Input Matrix (X):**

• **Target Labels (Y):**

• **Label Encoding (Ordinal):**

cat \rightarrow 1, fish \rightarrow 2, dog \rightarrow 3

\hat{Y}_p \leftarrow predicted labels

$$X = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} \quad n \times d$$

$$Y_a = \begin{bmatrix} \text{cat} \\ \text{dog} \\ \text{fish} \\ \vdots \end{bmatrix}$$

$$Y_a = \begin{bmatrix} 1 \\ 2 \\ 3 \\ \vdots \end{bmatrix}$$

model \rightarrow

multiclass classification

$$\hat{Y}_p = \begin{bmatrix} 0.8 \\ 1.5 \\ 2.75 \\ \vdots \end{bmatrix} \times$$

• **One-Hot Encoding:**

Cat \rightarrow [1 0 0], fish \rightarrow [0 1 0], dog \rightarrow [0 0 1]

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

sums to 1
model \rightarrow

$$\hat{Y}_p = \begin{bmatrix} 100 & 80 & 20 \\ 30 & 70 & 210 \\ 40 & 500 & 60 \\ \vdots & \vdots & \vdots \end{bmatrix}$$

softmax \rightarrow

$$\begin{bmatrix} 0.8 & 0.1 & 0.1 \\ 0.2 & 0.3 & 0.5 \\ 0.3 & 0.65 & 0.05 \\ \vdots & \vdots & \vdots \end{bmatrix}$$

Label encoding (One-hot encoding) Loss Computation

Minimize Objective Function:

$$\sum_i \min \left\| y_a^{\varepsilon i3} - \hat{y}_p^{\varepsilon i3} \right\|_2^2$$

MSE



$$y_a = [1 \ 0 \ 0]$$

$$\hat{y}_p = [0.8 \ 0.1 \ 0.1]$$

$$\frac{(1 - 0.8)^2 + (0 - 0.1)^2 + (0 - 0.1)^2}{}$$

In an ideal world,

$$\hat{y}_p = [1 \ 0 \ 0]$$

$$\text{error} = 0 + 0 + 0 \dots = 0$$

Maximize Objective Function:

$$\min - \sum_i y_a^{\varepsilon i3} \cdot \hat{y}_p^{\varepsilon i3 T}$$

max

$$\sum_i y_a^{\varepsilon i3} \cdot \hat{y}_p^{\varepsilon i3 T}$$



$$1 \times 0.8 + 0 \times 0.1 + 0 \times 0.1 = 0.8$$

In an ideal world

$$1 \times 1 + 0 \times 0 + 0 \times 0 = 1$$

$$1 + 1 + \dots = n$$

Interpretation:

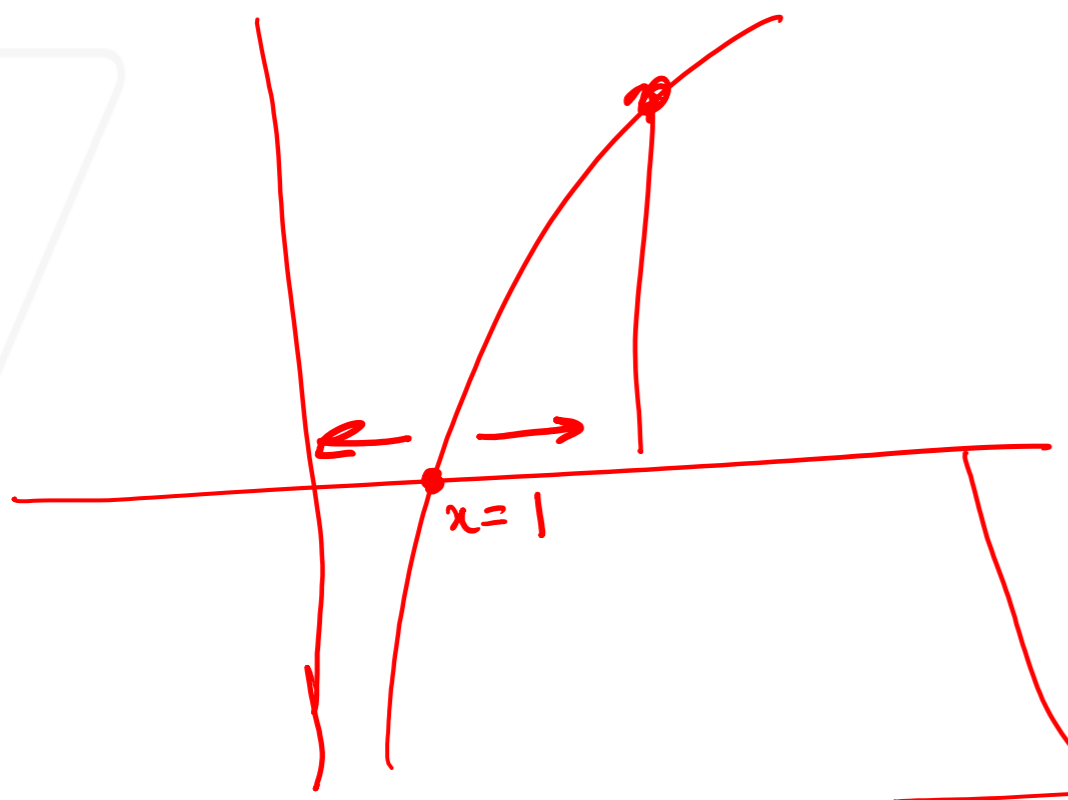
- Minimum squared Error works but isn't ideal for classification.
- Minimize negative sum of dot product is also good but Cross-Entropy / NLL is better.

Why Cross Entropy is better than dot product

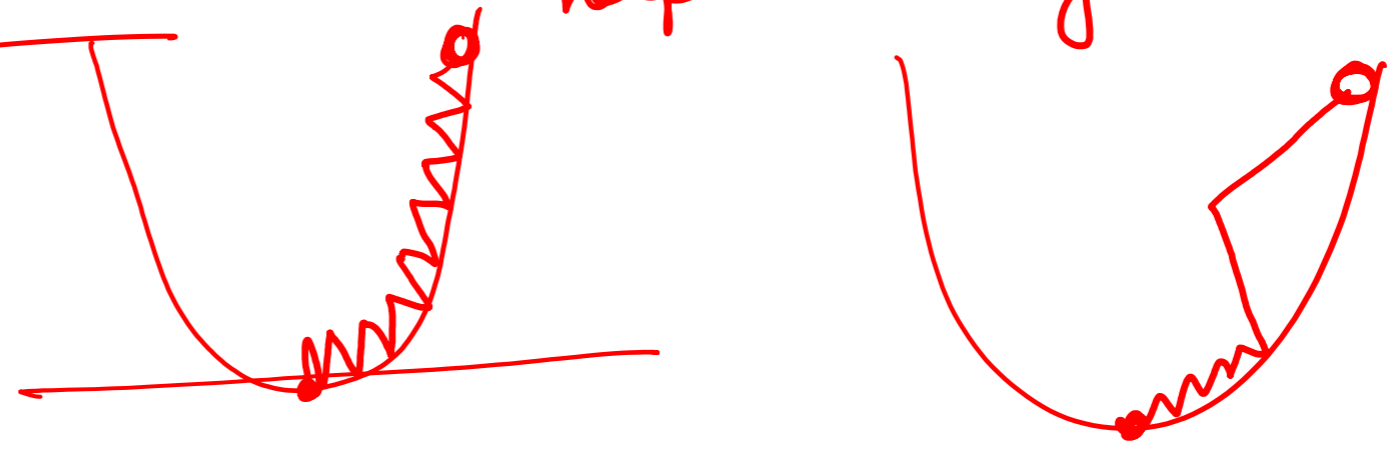
$$H(p, q) = - \sum_{x \in \mathcal{X}} p(x) \log q(x)$$

$$= \min_i \sum_i y_a^{i3} \cdot \log \hat{y}_p^{i3T} \quad \text{Case 1}$$

$$\text{Dot product} = \min_i \sum_i y_a^{i3} \cdot \hat{y}_p^{i3T} \quad \text{Case 2}$$



Log curve imposes a lot of penalty during optimizes which helps in easy convergence.



Logarithmic penalty encourages big update leading to faster convergence

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

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

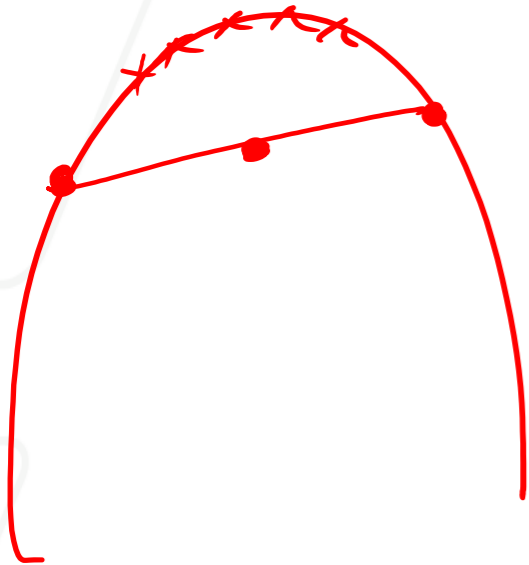
$$\begin{aligned}\sum_s P(s) \log \frac{Q(s)}{P(s)} &\leq \log \sum_s P(s) \frac{Q(s)}{P(s)} && \text{By Jensen Inequality} \\ &= \log \sum_s Q(s) = \log 1 = 0\end{aligned}$$

log function is concave or convex?

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

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

Concave Function: Jensen Inequality



$$\underline{\mathbb{E}[f(X)] \leq f(\mathbb{E}[X])}$$

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

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

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

$< 0 \rightarrow$ Concave

Unique

Strictly concave functions have \nearrow global maximum if it exists

Convex Function: Jensen Inequality

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$$

$$f''(x) > 0$$

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

$$x = [1 \quad 2]$$

$$p(x) = [\frac{1}{2} \quad \frac{1}{2}]$$

$$\mathbb{E}[X] = \sum_x x \cdot p(x) = 1 \times \frac{1}{2} + 2 \times \frac{1}{2} = 1.5$$

$$\mathbb{E}[f(x)] = \sum_x f(x) p(x) = 1^2 \times \frac{1}{2} + 2^2 \times \frac{1}{2} = 2.5$$

$$f(\mathbb{E}[x]) = 1.5^2 = 2.25$$

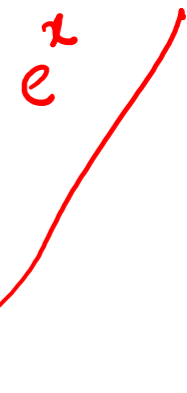
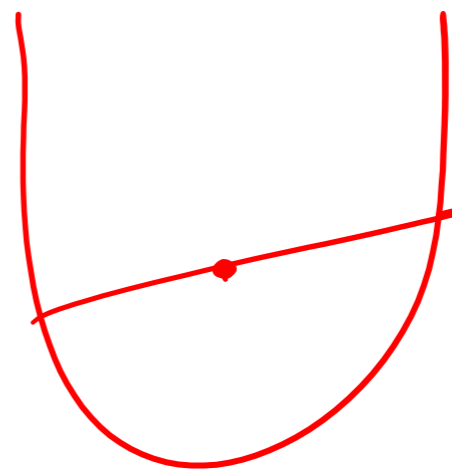
unique

Strictly convex functions have ₁ global minimum if it exists

$$f(x) = 3x + 4y$$

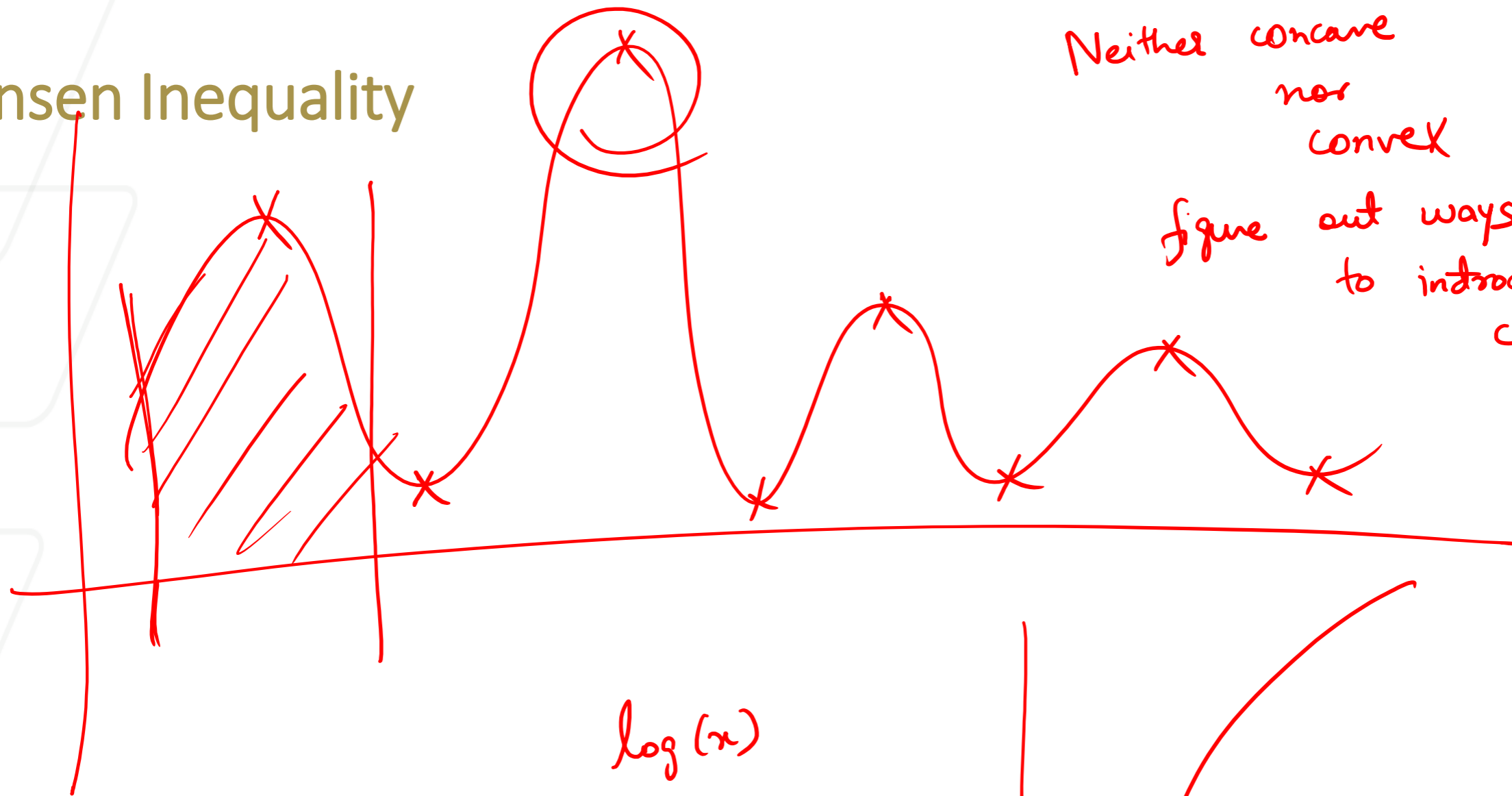
Hessian Matrix =

$$\begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial x \partial y} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix}$$



det(Hessian)
 > 0
 Convex
 < 0
 Concave

Jensen Inequality



Neither concave
nor
convex

figure out ways
to introduce
constraints

$\log(x)$

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

KL Divergence is always non negative

$$-KL[P][Q] = \sum_x p(x) \log \frac{q(x)}{p(x)}$$

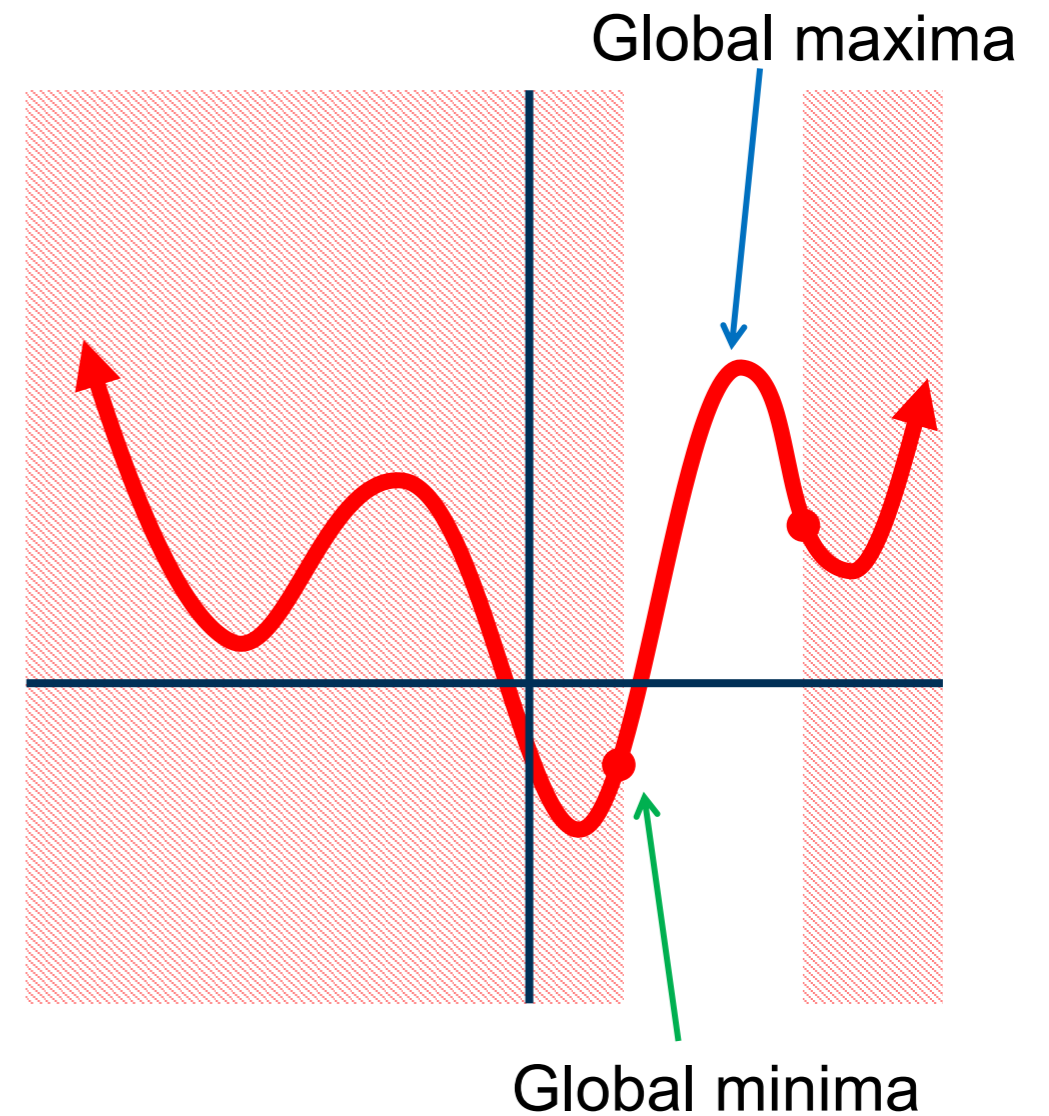
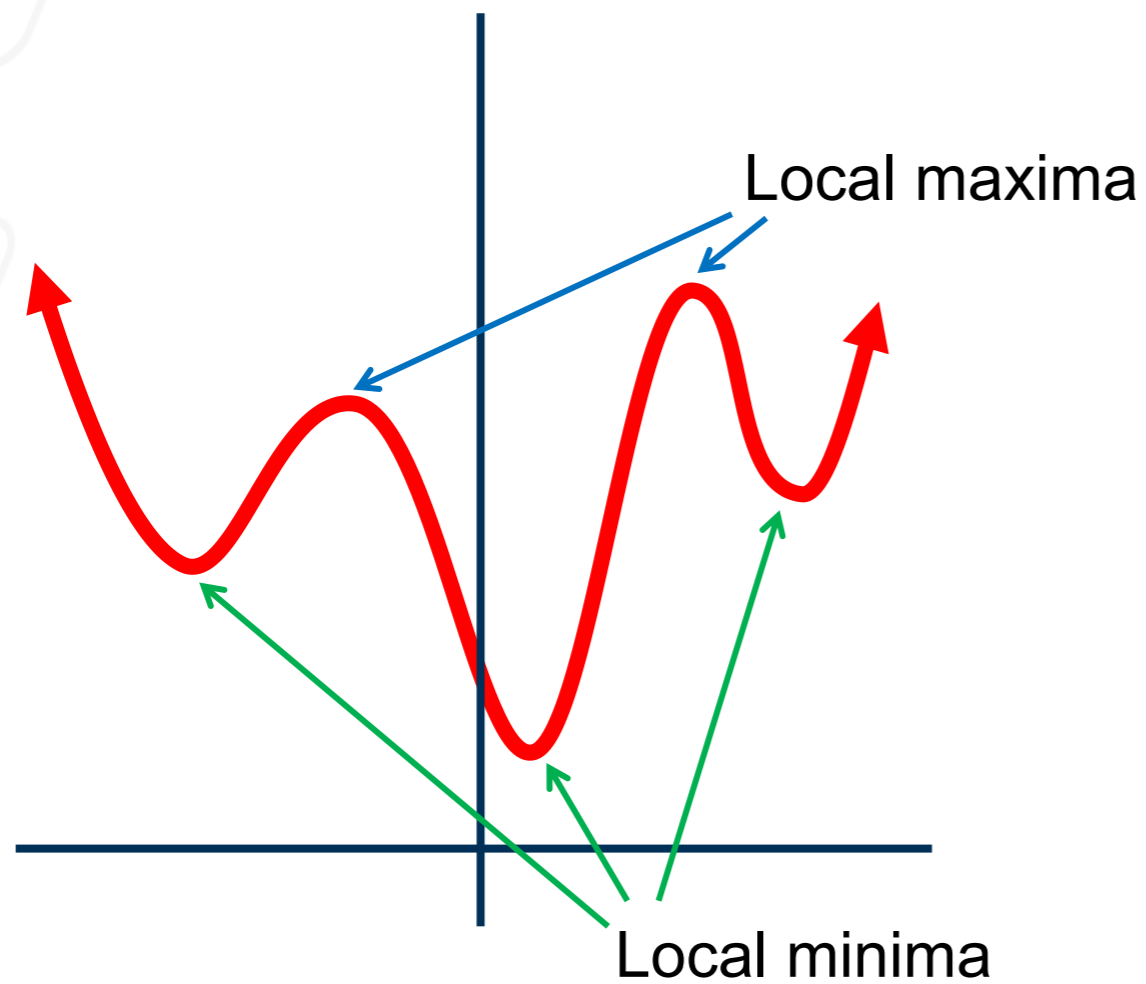
$q(x)$

$$\begin{aligned} &= \sum_x p(x) \log q(x) &= E[\log q(x)] \\ & &\leq \log(E[q(x)]) \\ & &\leq \log\left(E\left[\frac{q(x)}{p(x)}\right]\right) \\ & &\leq \log\left(\sum_x p(x) \cdot \frac{q(x)}{p(x)}\right) \\ & &\leq \log \sum_x q(x) \\ & &\leq \log 1 \end{aligned}$$

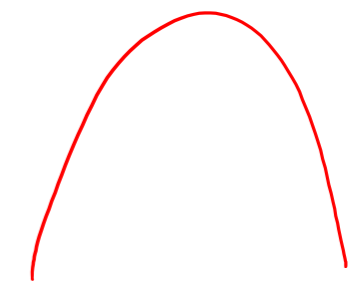
$$-KL[P][Q] \leq 0$$

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

Unconstrained and constrained optimization

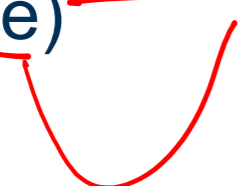


Basic unconstrained optimization problem



objective
function

concave

- Objective $f(\mathbf{x})$ the quantity we are trying to optimize (maximize or minimize)

- The variables x_1, x_2, \dots, x_m which can be represented in vector form as \mathbf{x}
(Note: the subscript m indicates a feature NOT a datapoint in our dataset)
- Finding the stationary points where the function is no longer increasing or decreasing

Equality constrained optimization

- An equality constrained optimization problem is expressed as:

$$\min_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})$$

subject to $g(\mathbf{x}, \mathbf{y}) = 0$

constraint is $x+y=5$
 $g(x,y) = x+y-5$

Inequality constrained optimization

- An inequality constrained optimization problem is expressed as:

$$\min_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})$$

subject to $g(\mathbf{x}, \mathbf{y}) \leq 0$

KKT conditions.

Linear, Quadratic, and Nonlinear Programming

Linear Programming (LP)

Objective: Linear function of x .

Constraints: Linear.

Example: $\min 3x_1 + 5x_2$ s.t. $x_1 + 2x_2 \leq 6$

Quadratic Programming (QP)

Objective: Quadratic

Constraints: Linear.

Example: $\min x^2 + y^2 + 3x + 5y$ s.t. $x + y \leq 4$

Nonlinear Programming (NLP)

Objective: Nonlinear.

Constraints: Can be nonlinear.

Example: $\min x^2 + y^2 + 3x + 5y$ s.t. $x^2 + y \leq 4$

WILL DEAL
WITH THESE 2
ONLY.

Recap

Amount of info needed to encode r.v. x with true pdf $p(x)$ but assuming to follow predicted pdf $q(x)$

$$CE = - \sum_x p(x) \log q(x) = H[P] + KL[P][Q]$$

• Cross Entropy

good objective or loss function in ML?

MLE
minimize -ve log likelihood

• Loss Functions

it leads to quicker convergence

$f(x) \rightarrow$ objective

MSE ↓

Dot product ↑

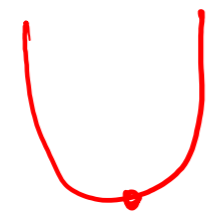
CE ↓

• KL Divergence

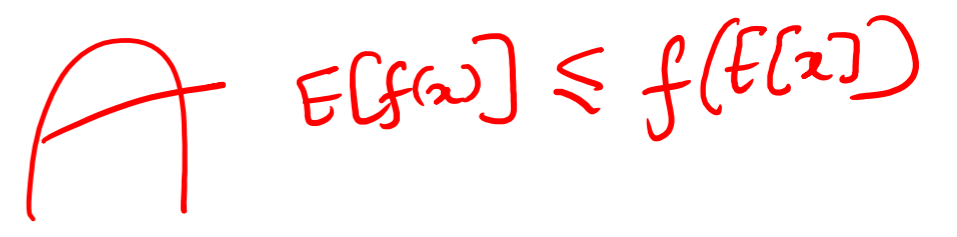
a kind of distance metric between 2 distributions

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

• Jensen Inequality

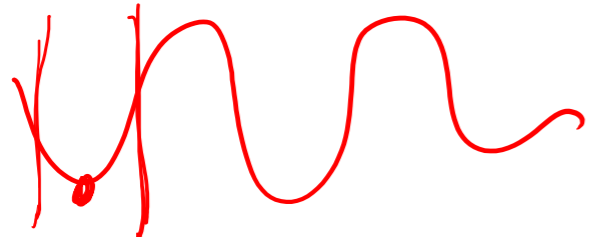


$$E[f(x)] \geq f(E[x])$$



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

• Unconstrained and Constrained Optimization

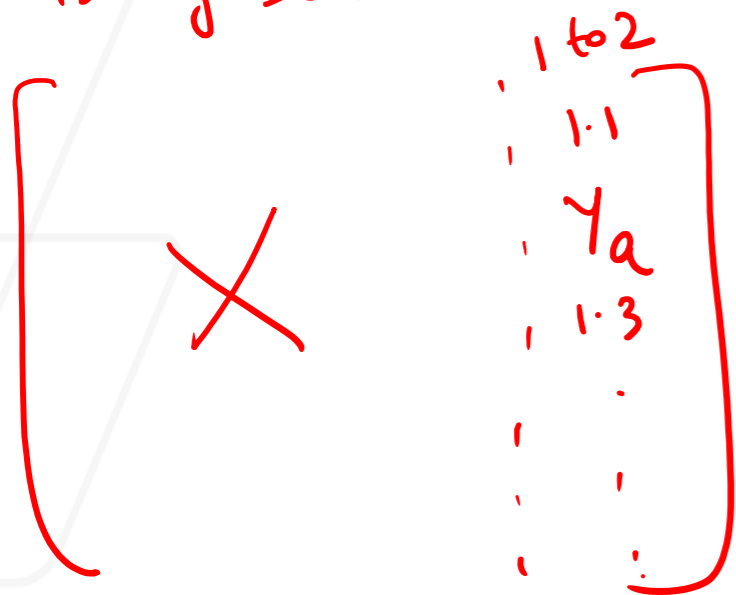


Linear: $f(x)$ is linear
 $g(x)$ is linear

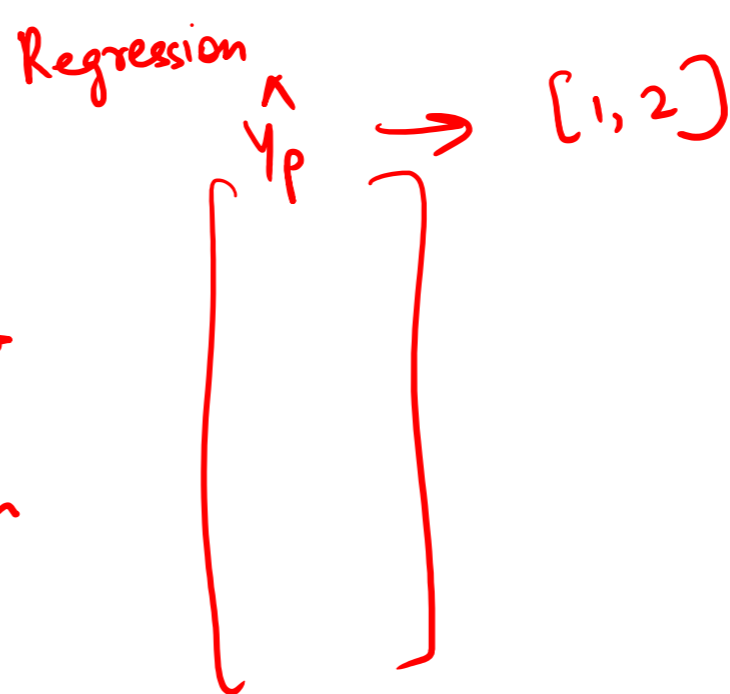
Quadratic: $f(x)$ is quadratic
 $g(x)$ is linear

• Linear, Quadratic and Non-Linear Programming

Training Dataset



Select Model
 $\theta_{ini} = \text{random}$

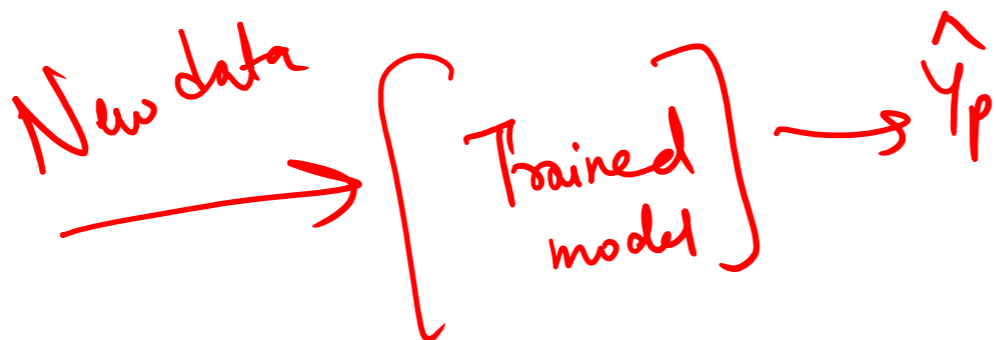


Construct an objective function

$$CE = \sum_i y_a^{(i)} \cdot \log \hat{y}_p^{(i)}$$

CE = very high value

CE will be very small
↓
Then θ is optimal
↓
model



Equality constrained optimization problem

$$\min f(x) \quad \text{s.t. } g(x)=0$$

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

$$g(x,y) = x + y - 3000$$

There is a notion of a Lagrangian expression where we combine both the objective function and our constraints

Lagrangian expression \rightarrow $L(x, \lambda) = f(x) - \lambda g(x)$

objective function $f(x)$ Lagrange multiplier λ constraint $g(x)$

s.t. $\lambda \neq 0$

$f(x)$ \leftarrow objective function needs to be optimized

$$\frac{\partial L}{\partial x} = 0, \quad \frac{\partial L}{\partial y} = 0$$

$$\frac{\partial f(x)}{\partial x} = 0 \Rightarrow \text{most optimal } \theta$$

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

amt. spent on food $\rightarrow x^2$ amt. spent on rent $\rightarrow y^2$

s.t. $x + y = 3000$

$g(x) = 0$

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

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

$$\frac{\partial x}{\partial x} = 0$$

$$\frac{\partial y}{\partial y} = 0$$

$$x=0, y=0$$

Equality constrained optimization problem

$$f(x) = x^2 + y^2 \quad \text{s.t.} \quad \begin{aligned} x + y &= 3000 \\ x + 5y &= 1 \end{aligned}$$

$$L(x, y, \lambda) = x^2 + y^2 - \lambda_1(x + y - 3000) - \lambda_2(x + 5y - 1)$$

$\frac{\partial L}{\partial x} = 0 \quad \frac{\partial L}{\partial y} = 0$

- With multiple equality constraints the Lagrangian is expressed as follows:

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \sum_{i=1}^k \lambda_i g_i(\mathbf{x})$$

$$\text{s.t.} \quad \lambda_i \neq 0$$

Equality constrained optimization problem

$$L(x, y, \lambda) = f(x, y) - \lambda g(x, y)$$
$$\nabla L(x, y, \lambda) = \nabla f(x, y) - \lambda \nabla g(x, y) = 0$$
$$\nabla f(x, y) = \lambda \nabla g(x, y)$$

- To solve this optimization problem, we only need to satisfy the stationarity condition, which states:

$$\nabla L(\mathbf{x}, \lambda) = 0$$

- We also need to satisfy

$$\lambda_i \neq 0$$

- This appeals to the intuition that at the optimal point, the gradient of the equality constraint is proportional to the gradient of the objective function

Equality constrained optimization: Example 1

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

$$g(x, y) = x + 5y - 1$$

$$L(x, y, \lambda) = x^2 + y^2 - \lambda(x + 5y - 1)$$

$$\frac{\partial L}{\partial x} = 0$$

$$2x - \lambda = 0 \Rightarrow \lambda = 2x \Rightarrow x = \lambda/2$$

$$\frac{\partial L}{\partial y} = 0$$

$$2y - 5\lambda = 0 \Rightarrow \lambda = \frac{2y}{5} \Rightarrow y = 5\lambda/2$$

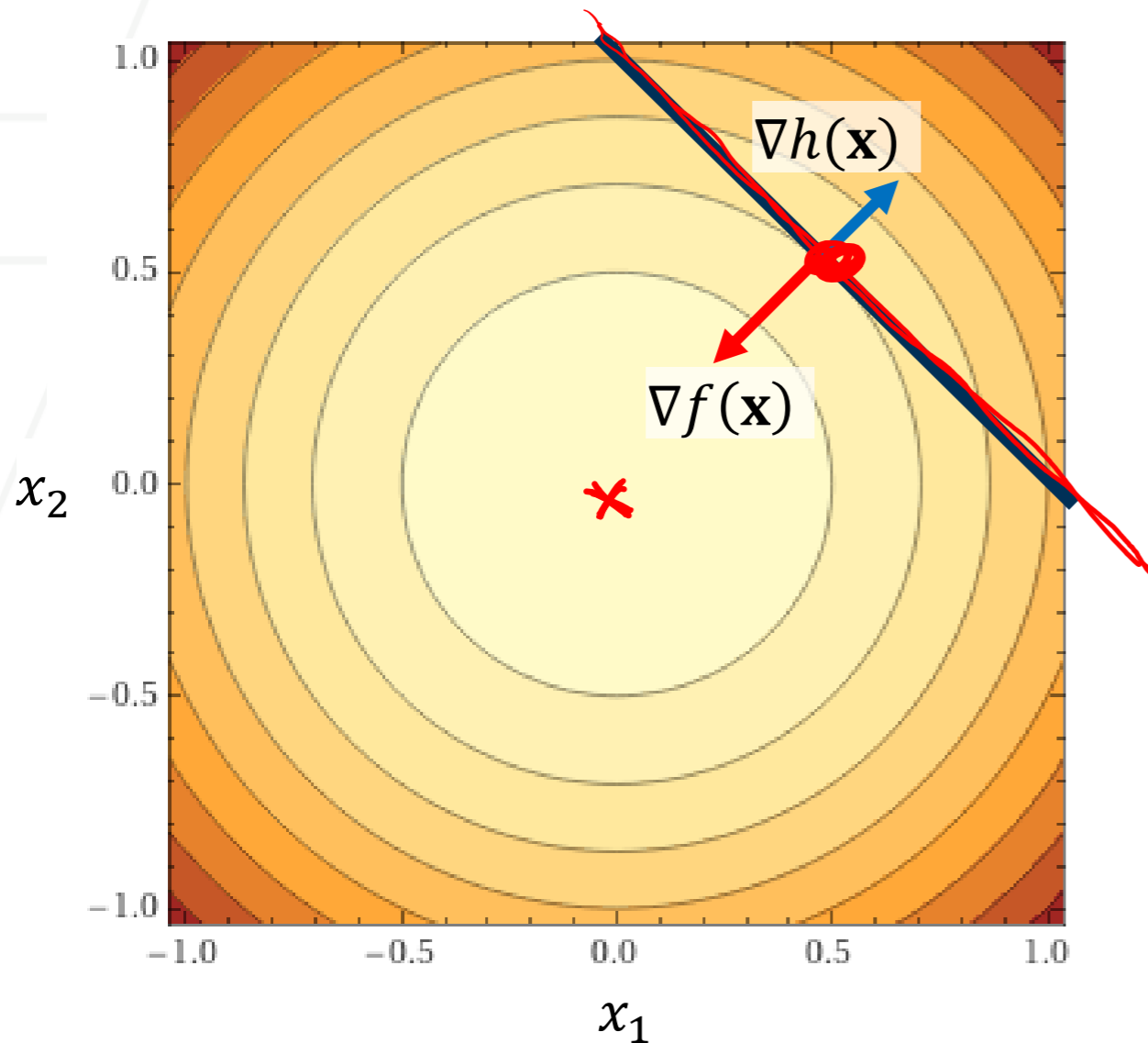
$$g(x) = 0$$

$$x + 5y = 1$$

$$\frac{\lambda}{2} + 5 \times \frac{5\lambda}{2} = 1 \Rightarrow \frac{26\lambda}{2} = 1 \Rightarrow \lambda = 1/13$$

$$\boxed{\begin{array}{l} x = 1/26 \\ y = 5/26 \end{array}}$$

Equality constrained optimization: *Example 2*



$$\begin{aligned} \max_{\mathbf{x}} \quad & 1 - x_1^2 - x_2^2 \\ \text{s. t.} \quad & x_1 + x_2 = 1 \end{aligned}$$

$$\begin{aligned} f(x_1, x_2) &= -1 + x_1^2 + x_2^2 \\ g(x_1, x_2) &= x_1 + x_2 - 1 \\ L(x_1, x_2, \lambda) &= -1 + x_1^2 + x_2^2 - \lambda(x_1 + x_2 - 1) \end{aligned}$$

Equality constrained optimization: *Example*

$$\begin{array}{ll} \max_{\mathbf{x}} & 1 - x_1^2 - x_2^2 \\ \text{s.t.} & x_1 + x_2 = 1 \end{array}$$

1. Write the Lagrangian for this problem

$$L(\mathbf{x}, \lambda) = 1 - x_1^2 + x_2^2 + \lambda(x_1 + x_2 - 1)$$

2. Calculate the gradient with respect to \mathbf{x}

$$\frac{\partial L}{\partial x_1} = -2x_1 + \lambda = 0 \rightarrow \lambda = 2x_1$$

$$\frac{\partial L}{\partial x_2} = -2x_2 + \lambda = 0 \rightarrow \lambda = 2x_2$$

Equality constrained optimization: *Example*

3. Solve for \mathbf{x}^* and λ

Using our previous two derivatives and the equality constraint, we obtain:

$$x_1 + x_2 = 1 \rightarrow x_1 + x_1 = 1 \rightarrow x_1^* = x_2^* = \frac{1}{2}$$

From $\frac{\partial L}{\partial x_1} = 0$ we find the value of λ

$$\lambda = 2x_1 \rightarrow \lambda = 1$$

Inequality constrained optimization problem

- Constraints that limit how small or big variables can be. These are inequality constraints noted as $g_j(\mathbf{x})$
- An optimization problem is usually expressed as:

$$\begin{array}{l} \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{subject to } g_j(\mathbf{x}) \leq 0, j = 1, 2, \dots, m \end{array}$$

objective

constraint

Inequality constrained optimization problem

↳ we have to add $\lambda g(x)$

- The optimization problem can be rewritten as

$$\min_{\mathbf{x}} f(\mathbf{x}) + \sum_{j=1}^m \lambda_j g_j(\mathbf{x})$$

subject to $\lambda \geq 0$

Handwritten annotations: A red circle around $f(\mathbf{x})$, a red circle around the summation term, and arrows pointing from the summation term to ≥ 0 and ≤ 0 .

- We need to satisfy the **Karush-Kuhn-Tucker (KKT)** conditions on \mathbf{x} , λ which is **necessary and sufficient** for optimality

Karush-Kuhn-Tucker Conditions

- Stationarity condition:

$$\min L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \sum_{j=1}^m \lambda_j g_j(\mathbf{x})$$

$\lambda \geq 0$
 $g_j(\mathbf{x}) \leq 0$
 $\nabla L(\mathbf{x}, \lambda) = 0$

- Primal feasibility

$$g_j(\mathbf{x}) \leq 0, j = 1, 2, \dots, m$$

- Dual Feasibility

$$\lambda_j \geq 0$$

- Complementary slackness

$$\lambda_j g_j(\mathbf{x}) = 0, j = 1, \dots, m$$

If $g(\mathbf{x}) < 0 \rightarrow$ **inactive** solution since $\lambda = 0$

If $g(\mathbf{x}) = 0 \rightarrow$ **active** solution since λ is at play and > 0 , which means constraint is active

$\lambda > 0, g(\mathbf{x}) = 0$
 equality constrained optimization

inactive Sol.
 unconstrained optimization

Active Sol.

Karush-Kuhn-Tucker Interpretations

- The complementary slackness condition applies only to inequality constraints
- It states that for a given inequality constraint...
...either $g_j(x) = 0$ or $\lambda_j = 0$
- Whenever the dual variable $\lambda_j > 0$ and therefore $g_j(x) = 0$, we say that the **constraint is active** or that the **constraint is tight at \mathbf{x}**
- The reverse $\lambda_j = 0$ and therefore $g_j(x) < 0$ means that the constraint is inactive, meaning that it is not effectively impacting the solution

Rule: For each inequality constraint, **either** it is active and has a positive multiplier, or it is inactive and its multiplier is zero.

Primal vs. Dual Form

$$f(x, y) = \frac{x^2}{2} + \frac{y^2}{2}$$

$$g(x, y) = x + y - 24$$

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

$$\frac{\partial L}{\partial x} = 0$$

$$\frac{\partial L}{\partial y} = 0$$

$$\begin{aligned} x - \lambda &= 0 & x &= \lambda \\ y - \lambda &= 0 & y &= \lambda \end{aligned}$$

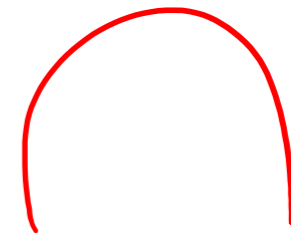
Quadratic Prog. is popular because it leads to

Primal Form

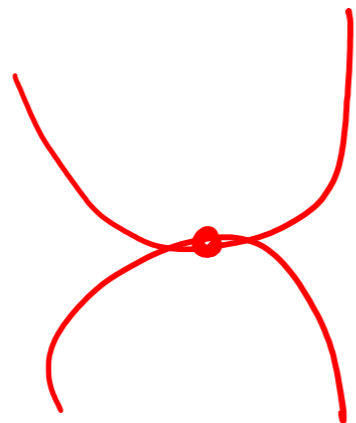
strong duality

$$\begin{aligned} L(\lambda) &= \frac{\lambda^2}{2} + \frac{\lambda^2}{2} - \lambda(\lambda + \lambda - 24) \\ &= \lambda^2 - \lambda^2 - \lambda^2 + 24\lambda = \underline{\underline{-\lambda^2 + 24\lambda}} \end{aligned}$$

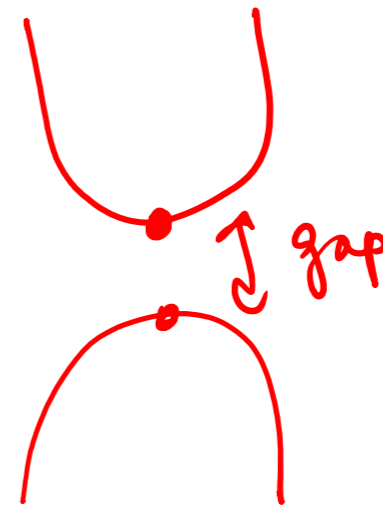
Dual Form



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



Strong Duality



Weak duality

Primal vs Dual Form

Primal Problem

Original optimization problem

Dual Problem

Derived problem as a **function of λ** only.

Analogy:

"Instead of directly solving for x and y together, we solve for what x and y *should look like* given a certain λ (like a candidate penalty price). Then we search over λ to find the best price that satisfies the constraint. This flips a 2D search into a 1D search."

Motivation:

In big problems (many x 's), solving for x in terms of λ can be done analytically or numerically, but then the dual problem might be much lower-dimensional \rightarrow computationally cheaper.

Primal vs Dual Form

Weak Duality

- Dual optimum \leq primal optimum (for minimization problems).
- Always true.

Strong Duality

- Dual optimum = primal optimum.
- Holds for **convex problems** with feasible solution.

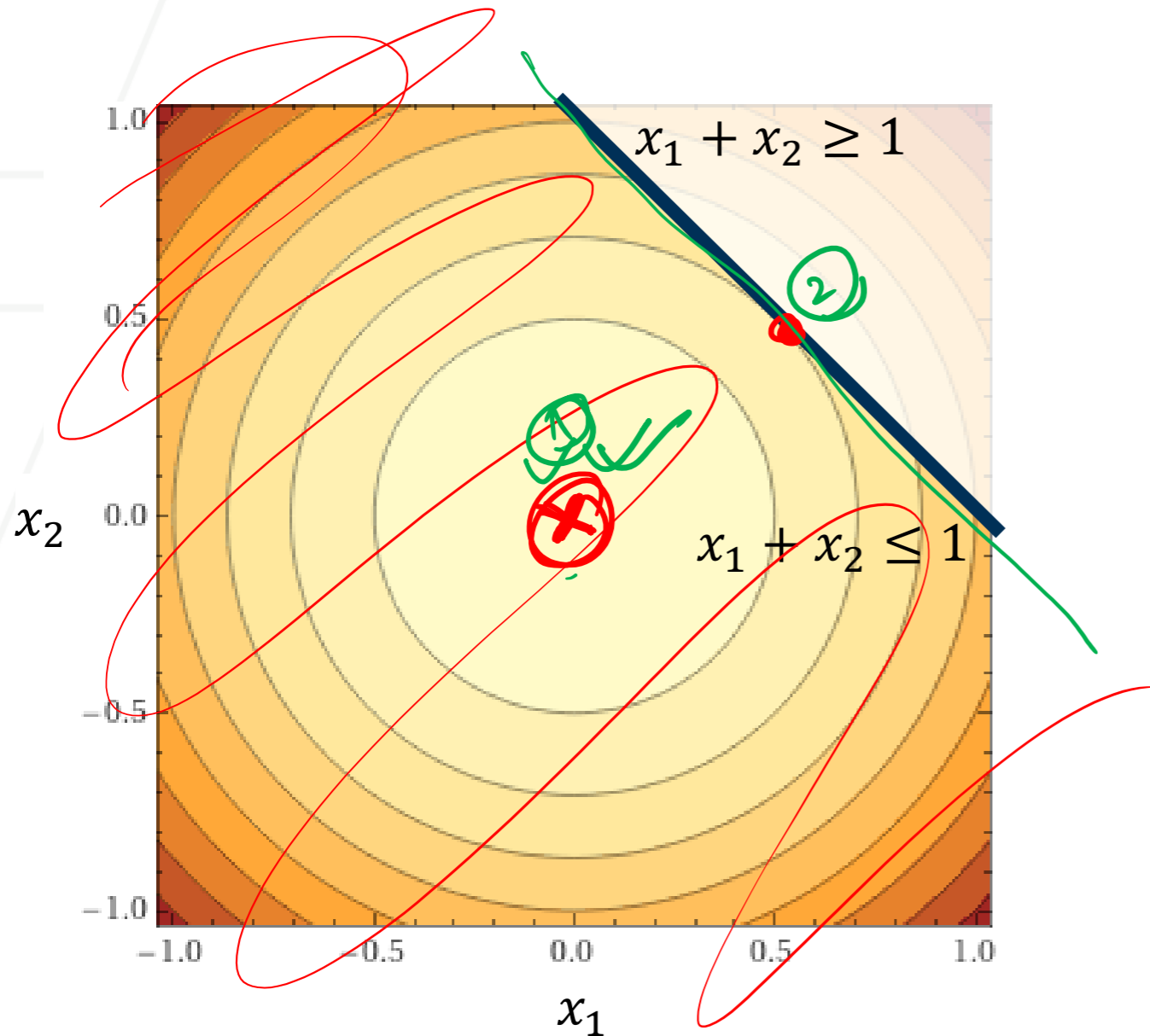
This is why **Quadratic Programming** is popular: QP is convex and has efficient solvers that work on the dual.

Why care?

Many ML models (e.g., SVMs) are solved via the dual.

Inequality constrained optimization: Example 1

$$f(x) = x^2 + y^2 \quad \text{s.t.} \quad x + y \geq 1$$
$$g(x, y) = -x - y + 1 \leq 0$$



$$\max_x \quad 1 - x_1^2 - x_2^2$$
$$\text{s.t.} \quad x_1 + x_2 \leq 1$$

$$g(x_1, x_2) = x_1 + x_2 - 1 \leq 0$$

$$g(x) \leq 0$$
$$\lambda \geq 0$$

Inequality constrained optimization: *Example 1*

1. Write the Lagrangian for this problem

$$L(\mathbf{x}, \lambda) = \underbrace{-1 + x_1^2 + x_2^2}_{f(x)} + \lambda \underbrace{(x_1 + x_2 - 1)}$$

2. Calculate the gradient with respect to \mathbf{x} (stationarity condition)

$$\frac{\partial L}{\partial x_1} = 2x_1 + \lambda = 0 \rightarrow \lambda = -2x_1$$

$$\frac{\partial L}{\partial x_2} = 2x_2 + \lambda = 0 \rightarrow \lambda = -2x_2$$

$$x_1 = x_2$$

Inequality constrained optimization: *Example 1*

3. Use the complementary slackness

~~Case 1:~~ $\lambda > 0$ and $g(x) = 0$

Active Sol.

$$x_1 + x_2 - 1 = 0 \rightarrow 2x_1 = 1$$

$$x_1 = x_2 = \frac{1}{2}$$

4a. Check if the dual feasibility condition is satisfied

$$\lambda = -2x_1 = -1$$

Dual feasibility X

This solution **does NOT** satisfy the dual feasibility condition and it is thus **unfeasible**

Inequality constrained optimization: *Example 1*

3. Use the complementary slackness

Case 2: $\lambda = 0$ and $\underline{g(x) < 0}$

Inactive sol.

$$\lambda = -2x_1$$

$$x_1 = x_2 = 0$$

Dual feasibility
 $\lambda \geq 0$

4b. Check if the primal feasibility condition is satisfied:

$$g(x_1, x_2) = 0 + 0 - 1 = -1 < 0$$

Primal feasibility

$$g(x) \leq 0$$

The solution is **feasible** and therefore the optimal!

Inequality constrained optimization: *Example 1*

Summary of the KKT conditions for the problem:

$$\left. \begin{aligned} \lambda &= -2x_1 \\ \lambda &= -2x_2 \end{aligned} \right\} \text{Stationarity}$$
$$x_1 + x_2 - 1 \leq 0 \quad \text{Primal feasibility}$$
$$\lambda \geq 0 \quad \text{Dual feasibility}$$
$$\lambda(x_1 + x_2 - 1) = 0 \quad \text{Complementary slackness}$$

Gradient descent

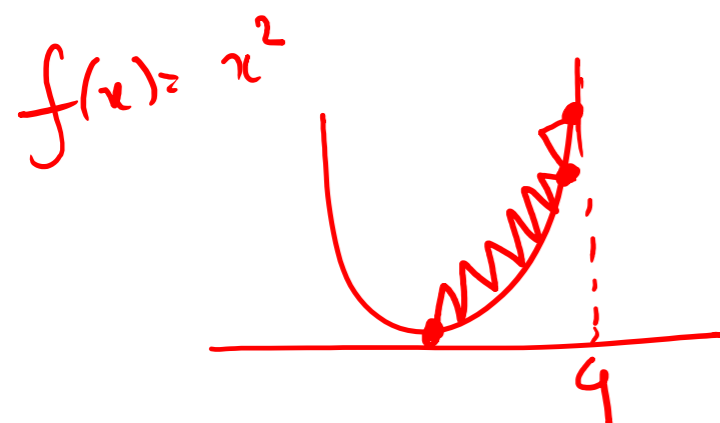
$$f(x) = e^x + 2x$$

Numerical Approaches



- In machine learning, some of the optimization problems suffer from the following challenges:
 - A closed form solution may be computationally intractable
 - Functions are complex: no easy derivative = 0 solution.
 - High dimensions: solving $\nabla f(x) = 0$ analytically is impossible.
 - All the data may not be available at the time of training (e.g. streaming data)
- Gradient descent is an iterative minimization technique for differentiable functions on a domain
- **Intuition:** the function $f(\mathbf{x})$ decreases the fastest by going from x_n to x_{n+1} in the opposite direction of the gradient of $f(\mathbf{x})$

$$x_{n+1} = x_n + \alpha \nabla f(x_n)$$

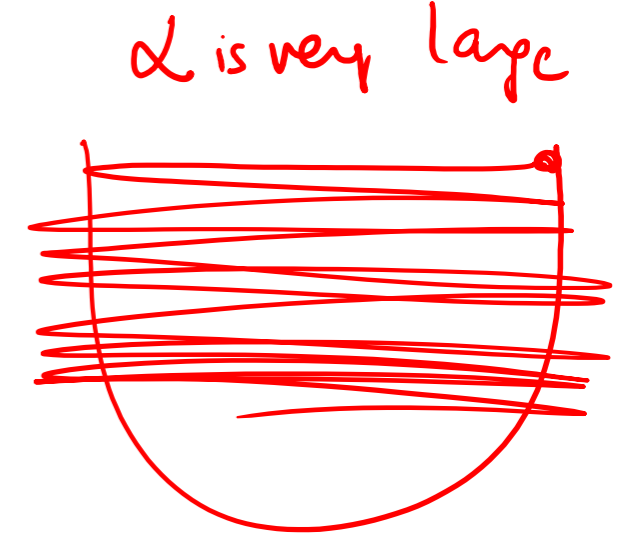


$$\mathbf{x}_{n+1} = \mathbf{x}_n - \alpha_n \nabla f(\mathbf{x}_n)$$

$$x_0 = 4 \quad x_1 = x_0 - \alpha_0 \nabla f(x_0)$$

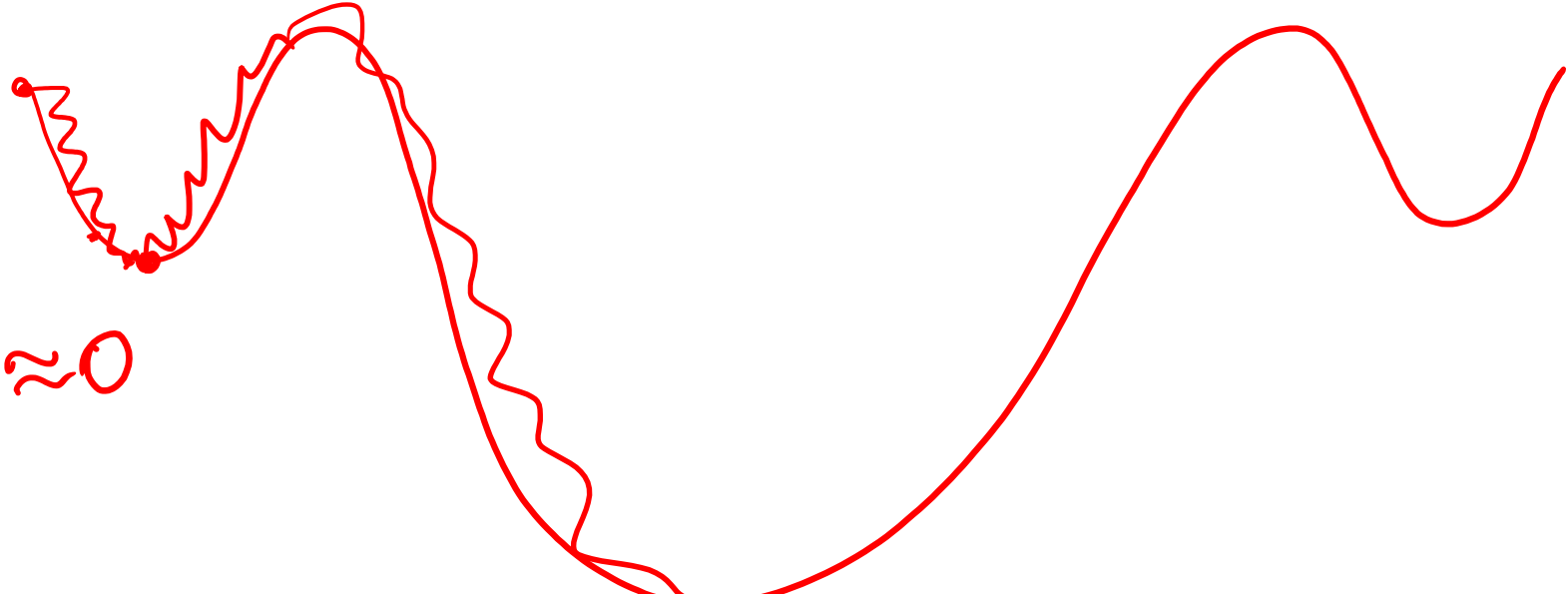
learning rate
HYPERPARAMETER

Learning rate



- The choice of the learning rate α_n plays an important role in the success of the minimization:
 - If α_n is **too small** then the number of iterations necessary to reach the minimum will be too large
 - If α_n is **too large** then we may overshoot the minimum
- It can be set to a constant or varied iteratively to improve the performance of the algorithm

Momentum



$\nabla f(x) \approx 0$

- The gradient descent technique is very susceptible to local minima, therefore, it can be modified by adding a momentum term

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \alpha_n g(\mathbf{x}_n)$$

Where:

$$g(\mathbf{x}_n) = \mu g(\mathbf{x}_{n-1}) + \nabla f(\mathbf{x}_n)$$

μ = momentum constant

- The update direction is a running average of past gradients.
- The momentum constant also needs to be fine tuned along with the learning rate to each problem, typically picked as a value between 0 and 1

Quick Knowledge Check

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

- Which property of a function is required to apply Jensen's inequality? A. Differentiability B. Convexity or concavity C. Existence of a global minimum D. Finite domain
- When is an inequality constraint considered **inactive** at the solution? A. When the constraint holds with equality B. When $\lambda > 0$ C. When the unconstrained optimum satisfies the constraint D. When the gradient is zero
- Which KKT condition ensures that either the constraint is tight or the multiplier is zero? A. Primal feasibility B. Dual feasibility C. Stationarity D. Complementary slackness
- What type of optimization landscape most strongly motivates the use of **momentum**? A. Linear objectives B. Perfectly symmetric bowls C. Narrow valleys and shallow local minima