

EE769 Introduction to ML

Basic Mathematics for ML

Amit Sethi, EE, IITB

asethi, 7483

Objective

- Revise and gain comfort with the following:
 - Vectors and matrices
 - Calculus and convex optimization
 - Probability and statistics

Scalar-vector operations

$$\alpha \in \mathbb{R}$$

$$\bar{x} \in \mathbb{R}^{n \times 1}$$

$$\alpha \bar{x} = \begin{bmatrix} \alpha x_1 \\ \alpha x_2 \\ \alpha x_3 \end{bmatrix}$$

.

$$\bar{x} = \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \in \mathbb{R}^{3 \times 1}$$

$$\bar{x} + a = \begin{bmatrix} x_1 + a \\ x_2 + a \\ x_3 + a \end{bmatrix}$$

$$\bar{x} + a \mathbf{1} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + a \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$\sqrt{n \times 1} \quad \sqrt{3 \times 1}$

Scalar-vector operations

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

x_1 scalar

$\mathbf{x} \in \mathbb{R}^3$

$a \in \mathbb{R}$

$\mathbf{x} + a$

Not a math function

$$\mathbf{x} + a = \begin{bmatrix} x_1 + a \\ x_2 + a \\ x_3 + a \end{bmatrix} = \mathbf{x} + a \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$
$$= \mathbf{x} + a \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

\downarrow \uparrow \uparrow \uparrow
 \mathbf{x} a 1 s

Vector-vector operations

$$\bar{x} + \bar{y} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \end{bmatrix}$$

$\bar{x} \in \mathbb{R}^n$ $\bar{y} \in \mathbb{R}^n$

$$\bar{x} \cdot \bar{y} = \bar{x}^\top \bar{y} = \langle x, y \rangle = \frac{x_1 y_1 + x_2 y_2}{\text{Scalar}}$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_{N \times 1} \odot \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}_{N \times 1} = \begin{bmatrix} x_1 y_1 \\ x_2 y_2 \end{bmatrix}_{N \times 1}$$

$\begin{bmatrix} x_1 & x_2 \end{bmatrix} \in \mathbb{R}^{1 \times 2}$ $\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \in \mathbb{R}^{2 \times 1}$

Vector-vector operations

$$\mathbf{x} + \mathbf{y} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ x_3 + y_3 \end{bmatrix}$$

\mathbf{x} & \mathbf{y}
have the same
dimensions

$$\begin{aligned}\mathbf{x} \cdot \mathbf{y} &= x_1 y_1 + x_2 y_2 + x_3 y_3 = \mathbf{x}^T \mathbf{y} \\ &= [x_1 \ x_2 \ x_3] \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \langle \mathbf{x}, \mathbf{y} \rangle\end{aligned}$$

$$\mathbf{x} \odot \mathbf{y} = \begin{bmatrix} x_1 y_1 \\ x_2 y_2 \\ x_3 y_3 \end{bmatrix} \in \mathbb{R}^3$$

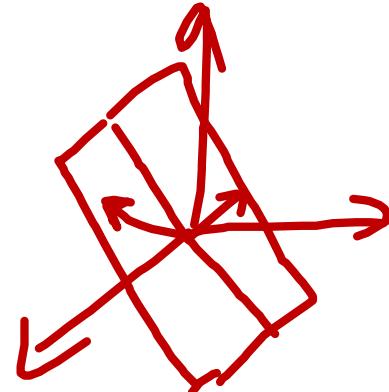
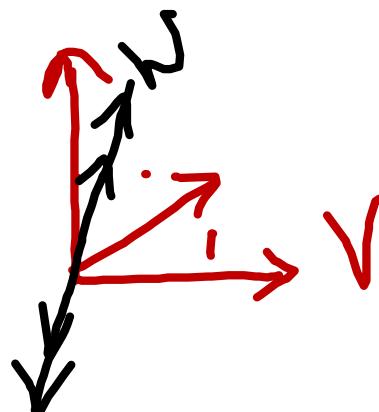
Sub-spaces spanned by vectors

$$\forall \bar{v}_1, \bar{v}_2 \in V$$

$$\begin{array}{l} (1) \quad \bar{v}_1 + \bar{v}_2 \in V \\ (2) \quad c\bar{v}_1 \in V \end{array}$$

$$\forall \bar{w}_1, \bar{w}_2 \in W \subseteq V$$

$$\begin{array}{l} (1) \quad \bar{w}_1 + \bar{w}_2 \in W \\ (2) \quad c\bar{w}_1 \in W \end{array}$$



$$\alpha_1 \bar{w}_1 + \alpha_2 \bar{w}_2$$

Sub-spaces

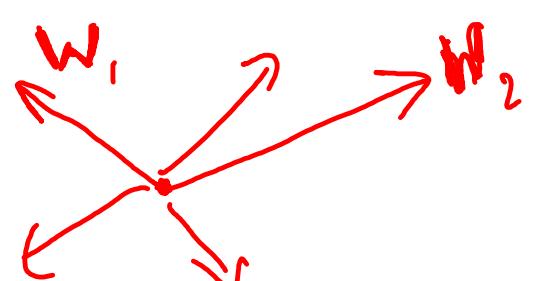
$$\forall \mathbf{v}_1, \mathbf{v}_2 \in V$$

Vector Space

$$\begin{aligned} & \textcircled{1} \quad \mathbf{v}_1 + \mathbf{v}_2 \in V \\ & \textcircled{2} \quad c \mathbf{v}_1 \in V \end{aligned}$$

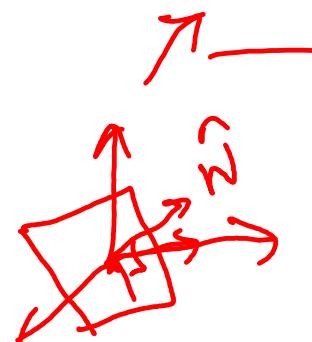
$$\rightarrow w_1, w_2 \in W \subset V \rightarrow w_1 + w_2 \in W \quad \text{Sub-Space}$$

$$c w_1 \in W$$



$$\alpha_1 w_1 + \alpha_2 w_2$$

$$\alpha_1, \alpha_2 \in \mathbb{R}$$



$$w \cdot \hat{w} = 0$$

any fixed

Matrix-matrix operations

Matrix-matrix operations

$$X \in \mathbb{R}^{M \times N}, Y \in \mathbb{R}^{M \times N}, x \in \mathbb{R}^N, x \in \mathbb{R}^{N \times 1}$$

Row-wise

$$x + y = \begin{bmatrix} (x_{11} + y_{11}) & \dots & (x_{1N} + y_{1N}) \\ \vdots & & \vdots \\ (x_{m1} + y_{m1}) & \dots & (x_{mN} + y_{mN}) \end{bmatrix}$$

$$xz = \begin{bmatrix} (x_{11} z_{11} + \dots + x_{1N} z_{1P}) & \dots & (x_{m1} z_{1P} + \dots + x_{mN} z_{1P}) \\ \vdots & \ddots & \vdots \\ \vdots & \ddots & (x_{m1} z_{1P} + \dots + x_{mN} z_{NP}) \end{bmatrix}$$

$$x \odot y = \begin{bmatrix} (x_{11} y_{11}) & & \\ & \ddots & \\ & & (x_{mN} y_{mN}) \end{bmatrix}$$

Transpose, determinant, and inverse of a matrix

$$\mathbf{X}^T = \begin{bmatrix} x_{11} & x_{21} & \cdots & x_{M1} \\ \vdots & \ddots & x_{ji} & \ddots \\ x_{1N} & \ddots & \ddots & x_{MN} \end{bmatrix} \in \mathbb{R}^{N \times M}$$

$$\mathbf{X} \in \mathbb{R}^{2 \times 2} \quad \mathbf{X} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad \det(\mathbf{X}) = ad - bc$$

$$\mathbf{X} \mathbf{X}^{-1} = \mathbf{I} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$3 \times 3^{-1} = 1$$
$$\begin{bmatrix} 1 & -2 \\ 3 & 4 \end{bmatrix} \times \begin{bmatrix} 0.4 & 0.2 \\ -0.3 & 0.1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$
$$\mathbf{X}^{-1} = \frac{1}{|\mathbf{X}|} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

Rank of a matrix

$$\bar{X} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \in \mathbb{R}^{2 \times 3} \quad \text{Rank}(\bar{X}) \leq \min(m, n)$$

$$\bar{Y} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 8 & 12 \end{bmatrix} \quad \text{Rank}(\bar{Y}) = 1$$

↑ ↑ ↑ Rank deficient

$$\bar{X} = \begin{bmatrix} 1 & 2 \\ 4 & 5 \end{bmatrix}$$

full rank matrix

$$\bar{Y} = \begin{bmatrix} 1 & 2 \\ 4 & 8 \end{bmatrix}$$

Rank deficient

Rank of a matrix

$$X = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \in \mathbb{R}^{2 \times 3}$$

Not invertible

$$\text{Rank}(X) \leq \min(M, N)$$

$$Y = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 8 & 12 \end{bmatrix}$$

Rank(Y) = 1
Rank deficient

If $X \in \mathbb{R}^{n \times m}$, Rank = \underline{m}
full Rank

Pseudo-inverse of a matrix

$$\bar{X}_{M \times N} \quad \bar{X}^+_{N \times M} = (\bar{X}^H \bar{X})^{-1} \bar{X}^H$$

$$\bar{X}^+ \bar{X} = [(\bar{X}^H \bar{X})^{-1} \bar{X}^H] X = I_{N \times N}$$

Conjugate transpose

Pseudo-inverse of a matrix

$$\cancel{X} \quad \cancel{X}^+ \quad (\cancel{X}\cancel{X}^+) \cancel{X} = \cancel{X}$$

$$\cancel{X}^+ = (\cancel{X}^H \cancel{X})^{-1} \cancel{X}^H$$

Conjugate transpose

$M \times N$
 $N \times M$
 $N \times N$

Eigen Decomposition

$$\bar{A} \bar{v}_i = \lambda_i \bar{v}_i \quad \bar{v}_i \rightarrow \text{eigen vector of } \bar{A}$$

$N \times N \quad N \times 1 \quad i \in \{1, \dots, n\} \quad \lambda_i \rightarrow \text{eigen value of } \bar{A}$

$$\bar{A} = \bar{Q} \bar{\Lambda} \bar{Q}^{-1} \quad \|\bar{v}_i\|_2^2 = 1 \leftarrow$$

$$v_i^T v_j = \delta_{ij}$$

$$\delta_{ij} = 0, \quad i \neq j$$

$$= 1, \quad i = j$$

Eigen Decomposition

$$A \mathbf{v}_i = \lambda_i \mathbf{v}_i$$

A : $N \times N$, \mathbf{v}_i : $N \times 1$, λ_i : Scalar

\mathbf{v}_i is an eigen vector of A

λ_i is the "value" of A

\mathbf{v}_i is a unit vector

$$A = Q \Lambda Q^{-1}$$

Q : $N \times N$ matrix of eigenvectors, Λ : Diagonal matrix of eigenvalues

$Q = [\mathbf{v}_1 \ \mathbf{v}_2 \ \mathbf{v}_3 \ \dots \ \mathbf{v}_N]$

$\Lambda = \begin{bmatrix} \lambda_1 & & & & \\ & \lambda_2 & & & \\ & & \ddots & & \\ & & & \lambda_N & \end{bmatrix}$

Tensors

$\alpha \in \mathbb{R}$; $\bar{x} \in \mathbb{R}^N$; $\bar{y} \in \mathbb{R}^{M \times N}$
 $\underline{T} \in \mathbb{R}^{M \times N \times P \times \dots \times Q}$.

transpose($T, [3, 1, 2, 0]$)
D, 1, 2, 3

0, 1
1, 0

Tensors

$$\alpha \in \mathbb{R}$$

$$x \in \mathbb{R}^N$$

$\overset{0}{M} \times \overset{1}{N} \times \overset{2}{P} \times \overset{3}{Q}$

$$T \in \mathbb{R}$$

$\nwarrow \uparrow \uparrow$

$$R = \text{transpose}(T, \underline{[3, 1, 2, 0]})$$

order of
permutation
of dimensions

Functions

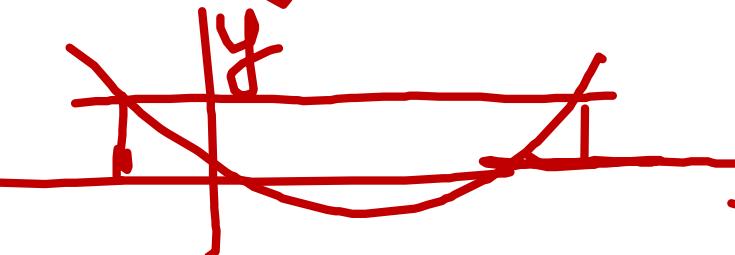
$$f: X \rightarrow Y$$

$$x \in X$$

$$f(x) \in Y$$

$$y = ax^2 + bx + c$$

many \leftarrow \rightarrow one



$$x^2 + y^2 = 1 \Rightarrow y = \pm \sqrt{1 - x^2}$$

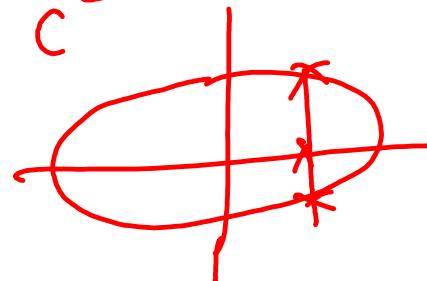
X NOT
a function

Functions

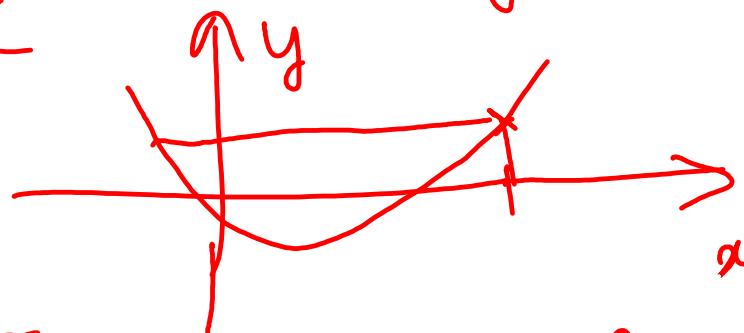
$$f: X \rightarrow Y \quad x \in X \quad f(x) \rightarrow y$$

$$f(x) = ax^2 + bx + c$$

$$ax^2 + by^2 = c^2$$



many to one



Not a function

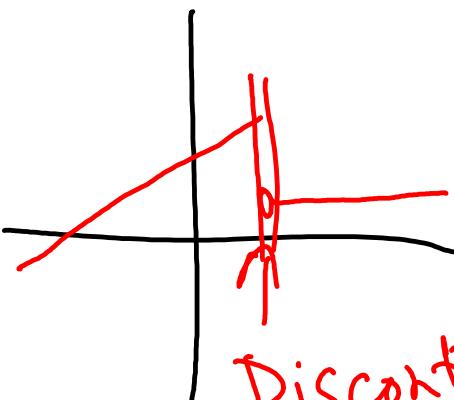
Continuity

x is continuous

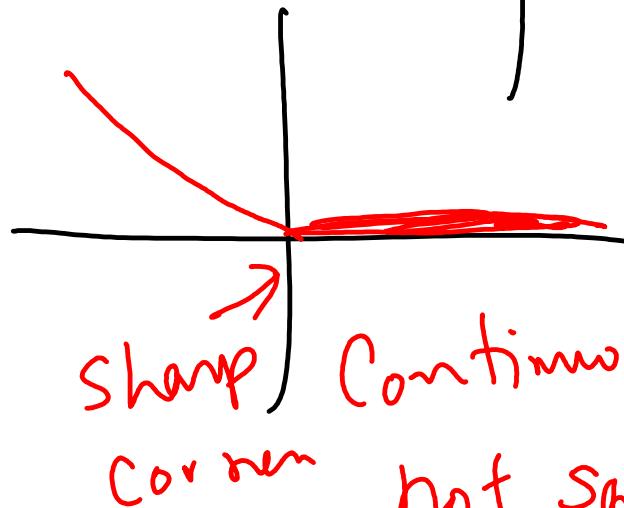
$$\text{if } f(x + \Delta x) = f(x) \text{ for } \Delta x \rightarrow 0$$

Lipschitz continuity

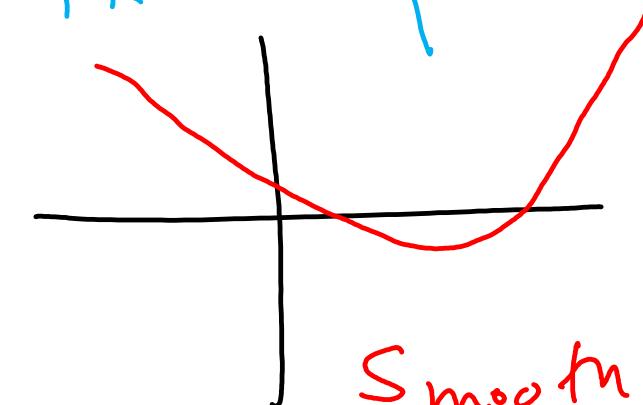
$$+k|f(x_1) - f(x_2)| \leq k|x_1 - x_2|$$



Discontinuous



Sharp corner
Continuous

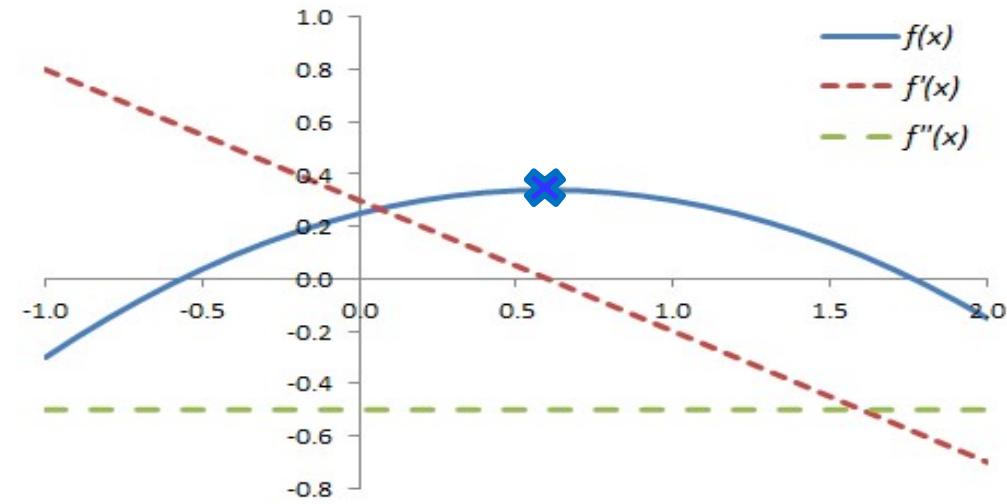
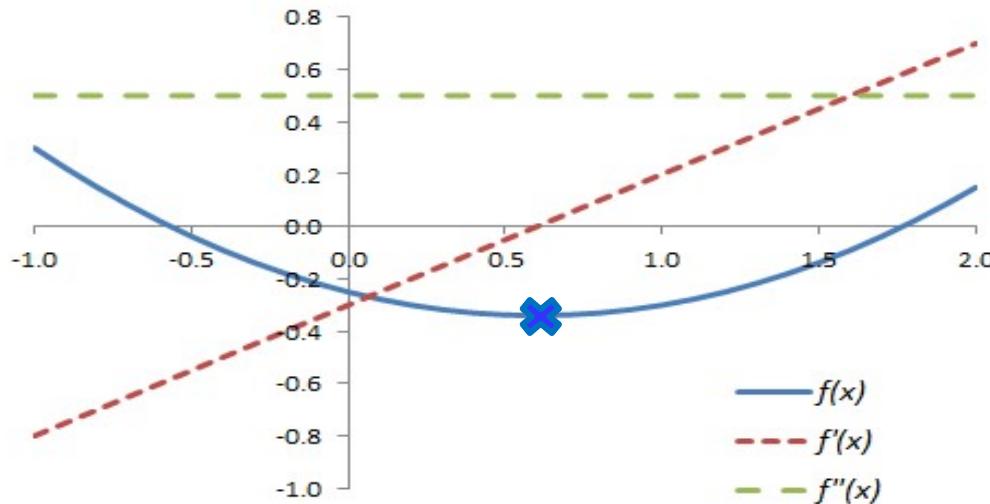


Smooth
continuous

$$y = \sqrt{x} \quad y = x^2$$

Lipschitz continuity

Derivative of a function of a scalar



E.g. $f(x) = ax^2 + bx + c, \quad f'(x) = 2ax + b, \quad f''(x) = 2a$

- Derivative $f'(x) = \frac{d f(x)}{d x}$ is the rate of change of $f(x)$ with x
- It is zero when the function is flat (horizontal), such as at the minimum or maximum of $f(x)$
- It is positive when $f(x)$ is sloping up, and negative when $f(x)$ is sloping down
- To move towards the maxima, taking a small step in a direction of the derivative

Derivative of a function

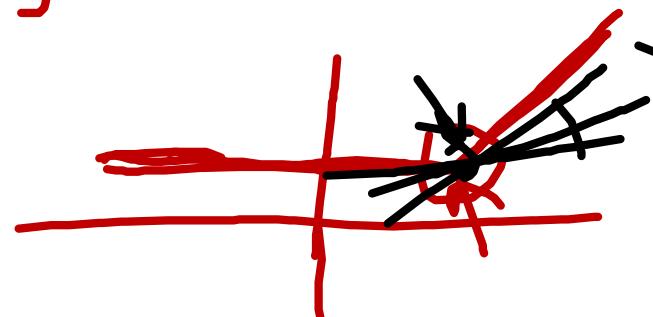
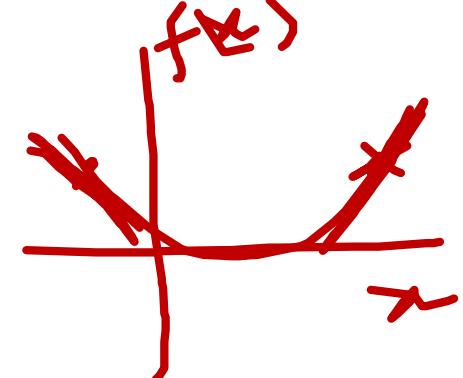
$$\text{Let } f(x+\Delta x) = f(x) + x$$

$\Delta x \rightarrow 0$

$$f'(x) = \lim_{\Delta x \rightarrow 0}$$

$$\frac{f(x+\Delta x) - f(x)}{\Delta x}$$

continuous



Sub-tangent

Sub-gradient

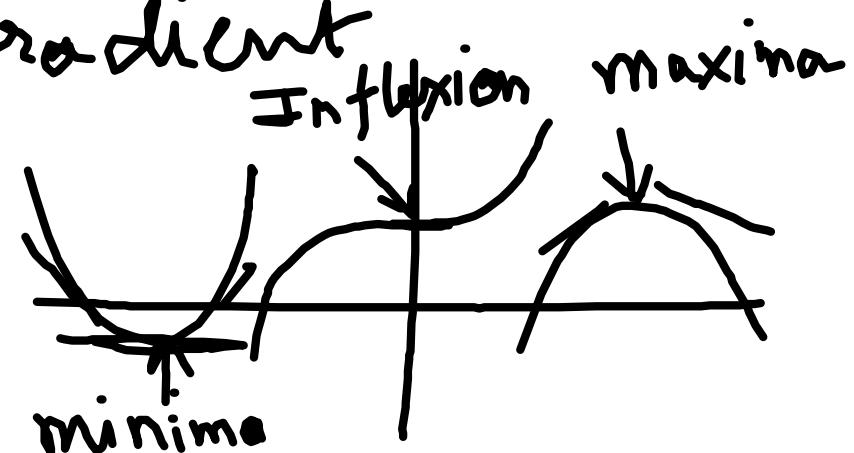
$f'(x) = 0 \rightarrow$ critical

point

$f''(x) > 0 \rightarrow$ min

$< 0 \rightarrow$ max

$= 0 \rightarrow$ inflexion

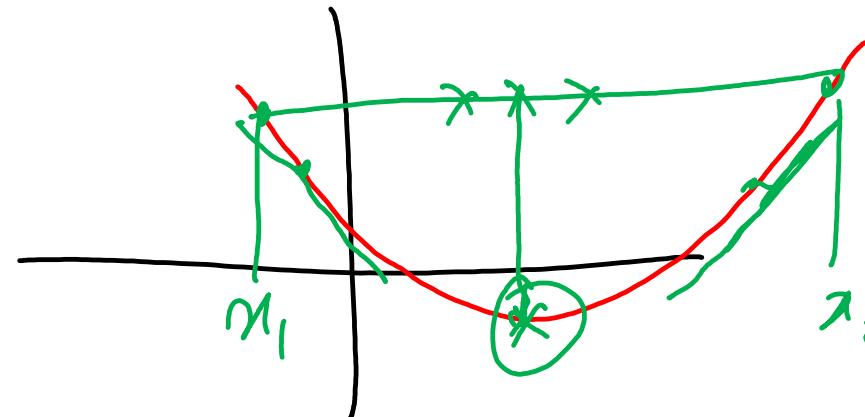


Derivative of a function

Gradient

$$f(x)$$

$$f'(x) = 2ax + b$$



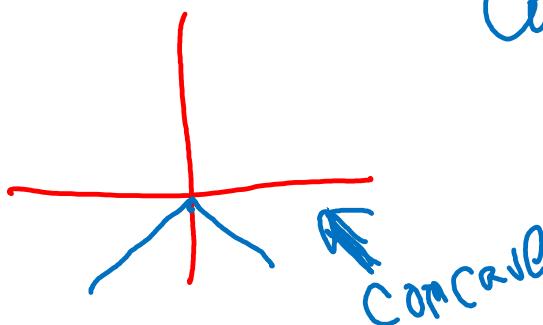
$$\lambda f(x_1) + (1-\lambda) f(x_2)$$

$$\geq f(\lambda x_1 + (1-\lambda)x_2)$$

$$an^2 + bn + c$$

Convex function

Convex



Subgradient

Sub-derivatives

are slopes of
subtangents

$$f(x) = |x|$$

$$f'(x) = \begin{cases} -1, & \text{if } x < 0 \\ +1, & \text{if } x > 0 \end{cases}$$

Critical points of a function

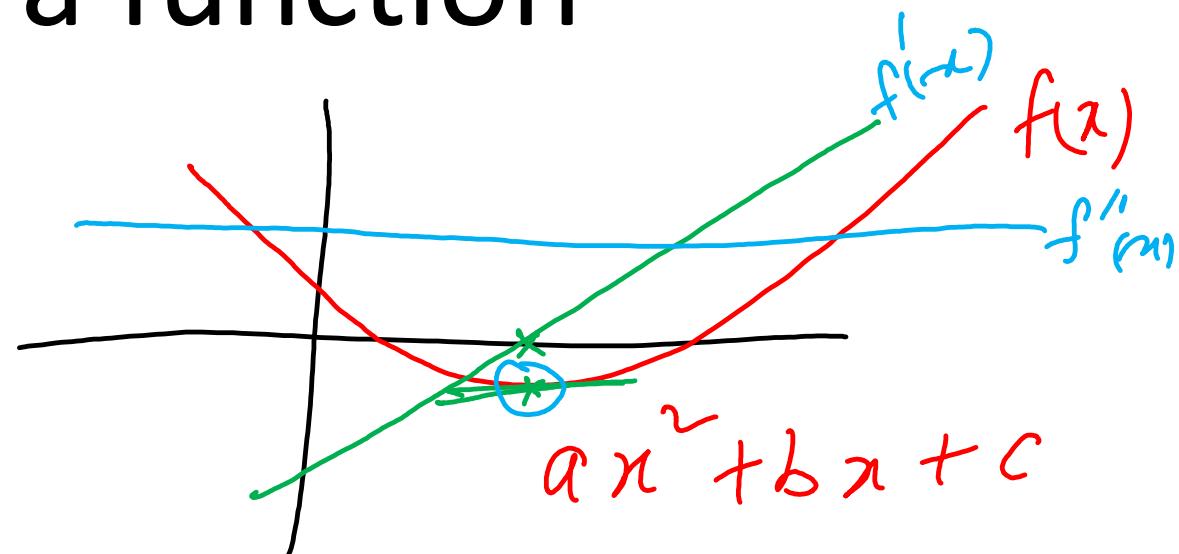
$$f'(x) = 0$$

$$f'(x) = 2ax + b$$

$$x = \frac{-b}{2a}$$

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

$\geq 0 \leftarrow$ Convex
Concave fn. $f''(x) \leq 0$



Maxima \rightarrow if $f''(x) < 0$

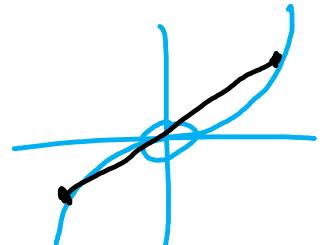
Minimas \rightarrow if $f''(x) > 0$

Inflection

$$f(x) = x^3$$

$$f'(x) = 3x^2$$

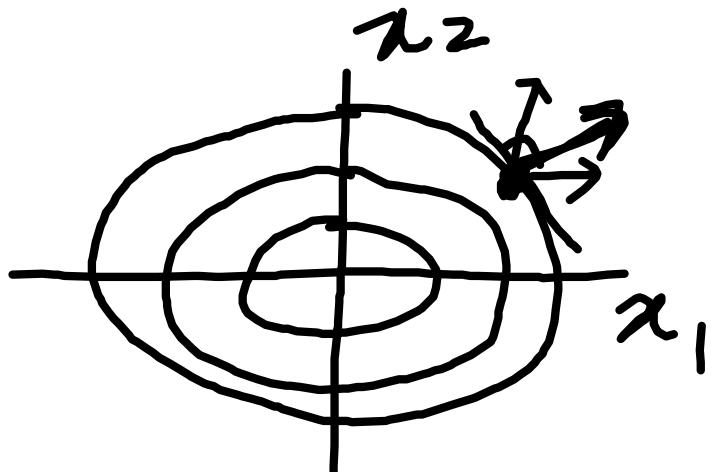
$$f''(x) = 6x$$



Multivariate functions

$$y = f(x_1, x_2) = ax_1^2 + bx_2^2 \leftarrow$$

$$\nabla y = \begin{bmatrix} \partial y / \partial x_1 \\ \partial y / \partial x_2 \end{bmatrix} = \begin{bmatrix} 2ax_1 \\ 2bx_2 \end{bmatrix} \leftarrow$$

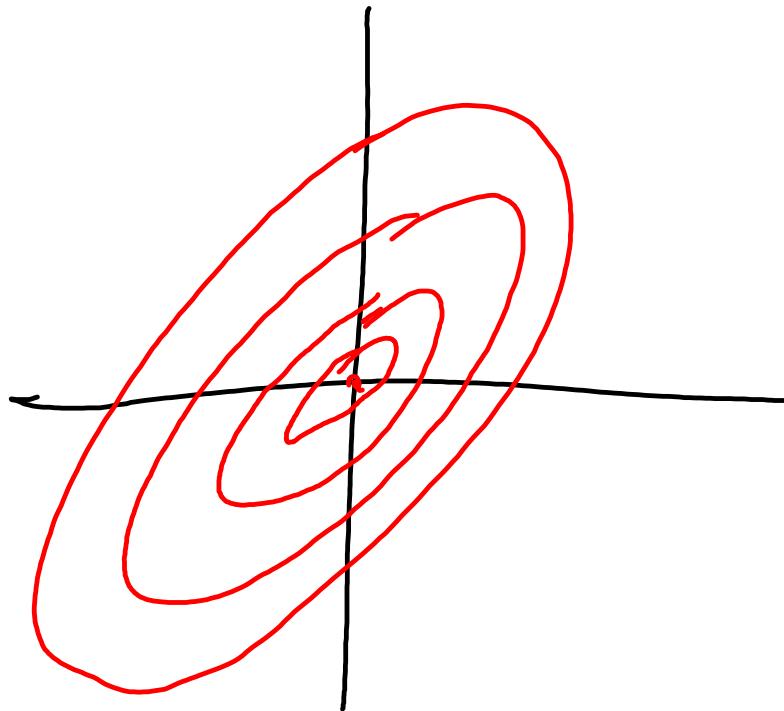


G. D. $\overset{\equiv}{\text{in}}$ $z - d$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \underset{\text{---}}{\leftarrow} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - n \nabla y$$

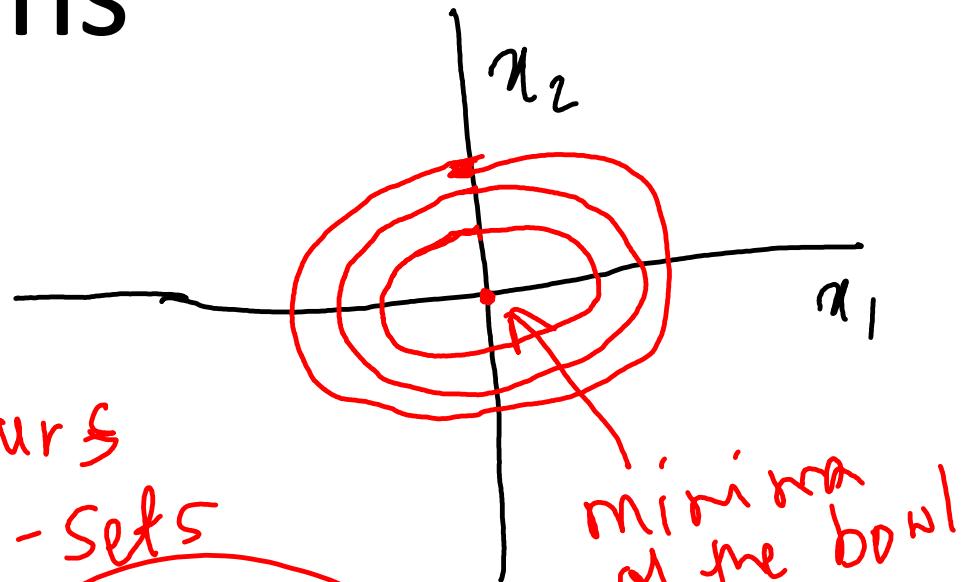
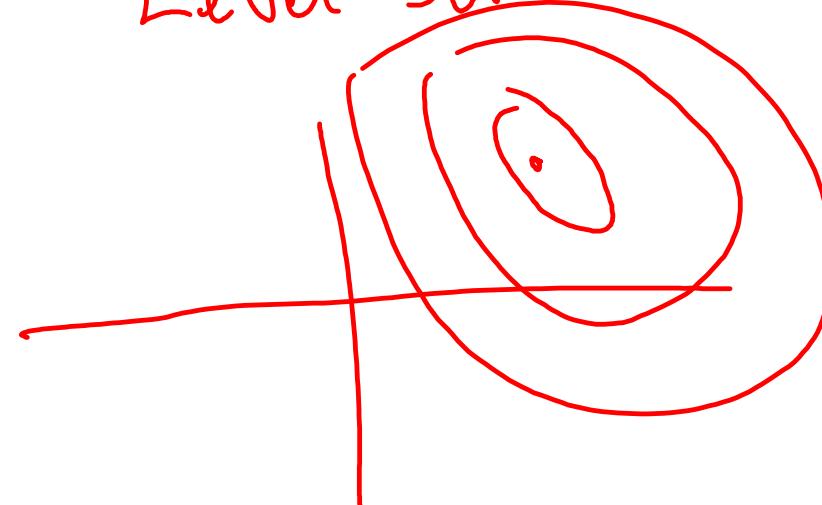
Multivariate functions

$$y = a x_1^2 + b x_2^2$$

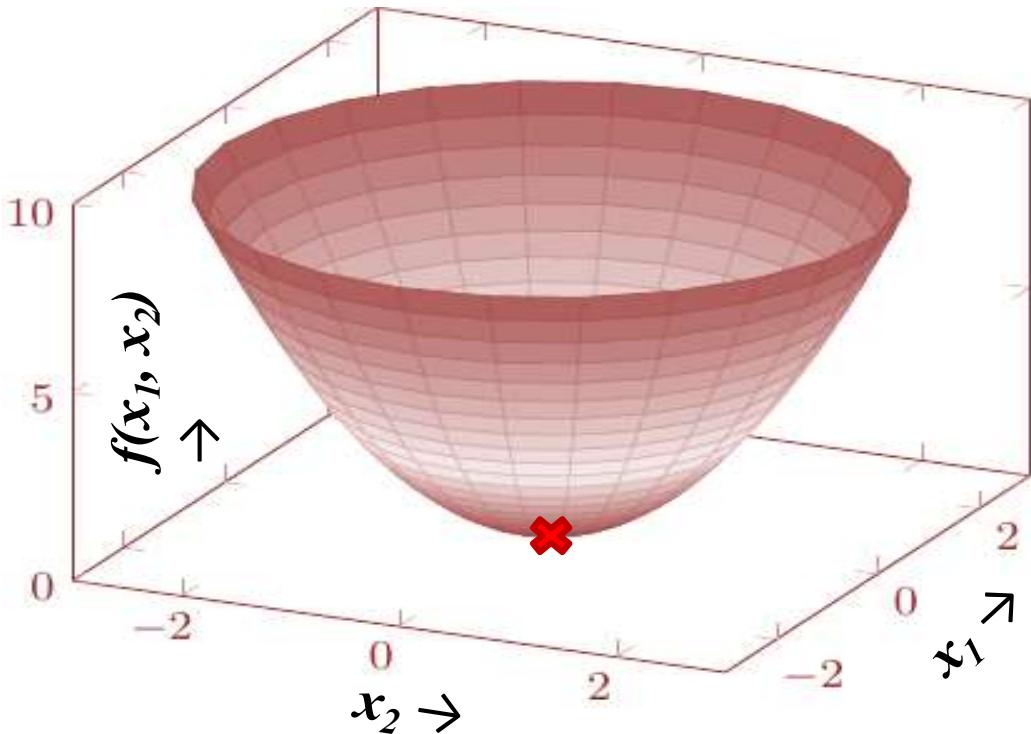


Contours

Level-Sets

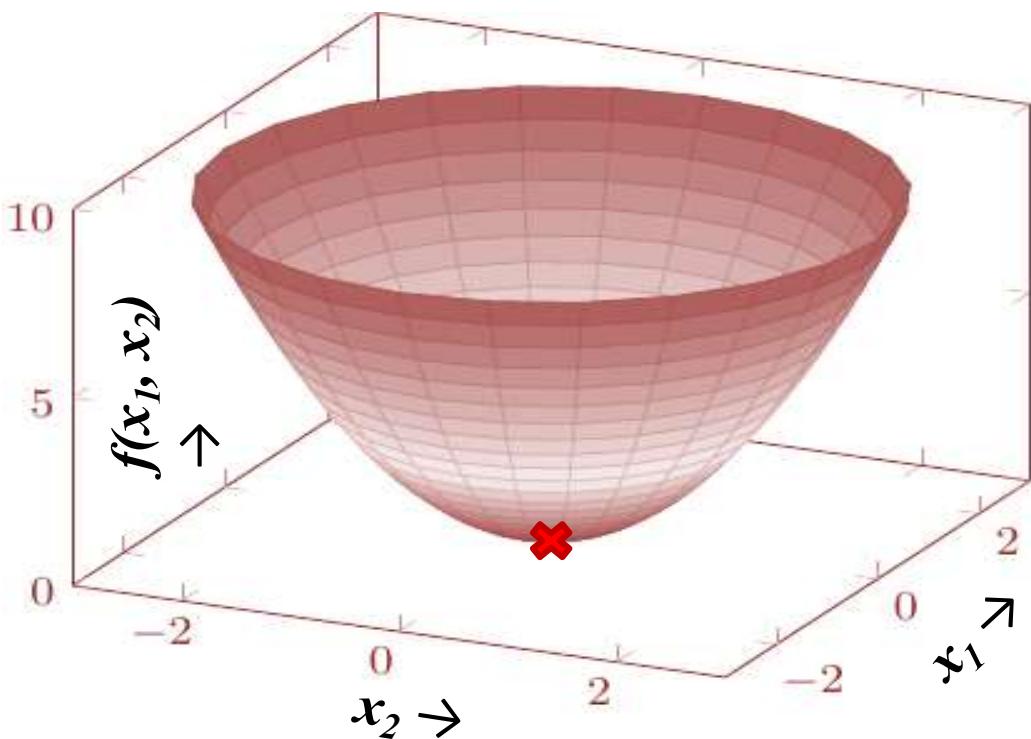


Gradient of a function of a vector

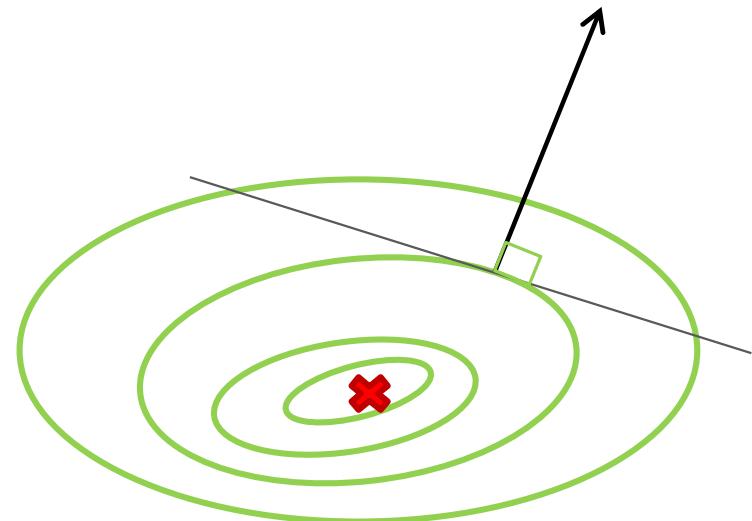


- Derivative with respect to each dimension, holding other dimensions constant
- $\nabla f(\mathbf{x}) = \nabla f(x_1, x_2) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix}$
- At a minima or a maxima the gradient is a zero vector
The function is flat in every direction
- At a minima or a maxima the gradient is a zero vector

Gradient of a function of a vector



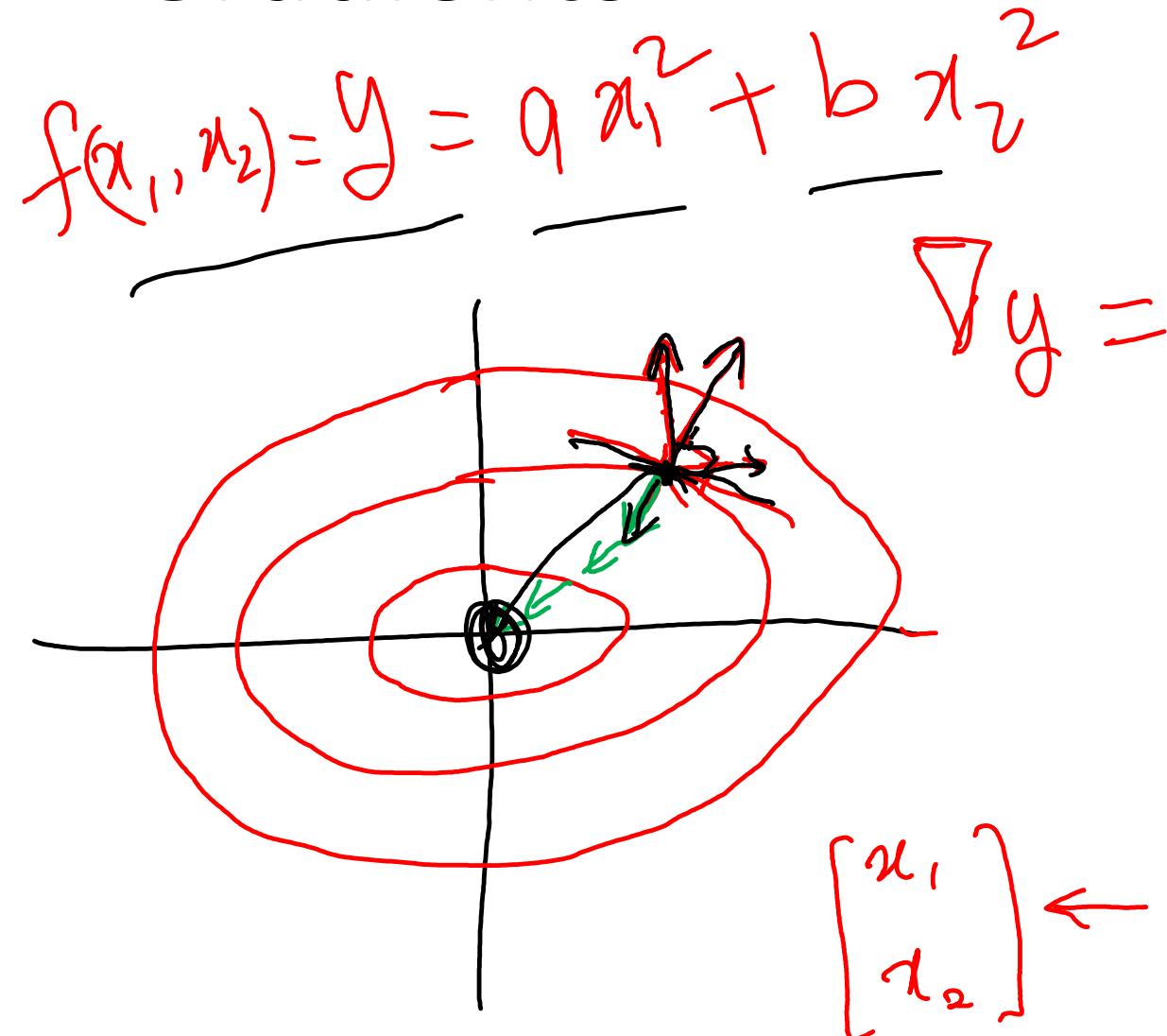
- Gradient gives a direction for moving towards the minima
- Take a small step towards negative of the gradient



Example of gradient

- Let $f(\mathbf{x}) = f(x_1, x_2) = 5x_1^2 + 3x_2^2$
- Then $\nabla f(\mathbf{x}) = \nabla f(x_1, x_2) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 10x_1 \\ 6x_2 \end{bmatrix}$
- At a location $(2,1)$ a step in $\begin{bmatrix} 20 \\ 6 \end{bmatrix}$ or $\begin{bmatrix} 0.958 \\ 0.287 \end{bmatrix}$ direction will lead to maximal increase in the function

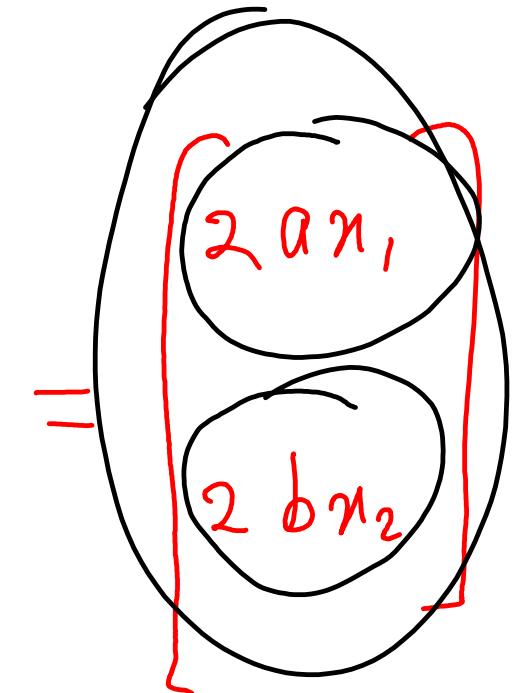
Gradients



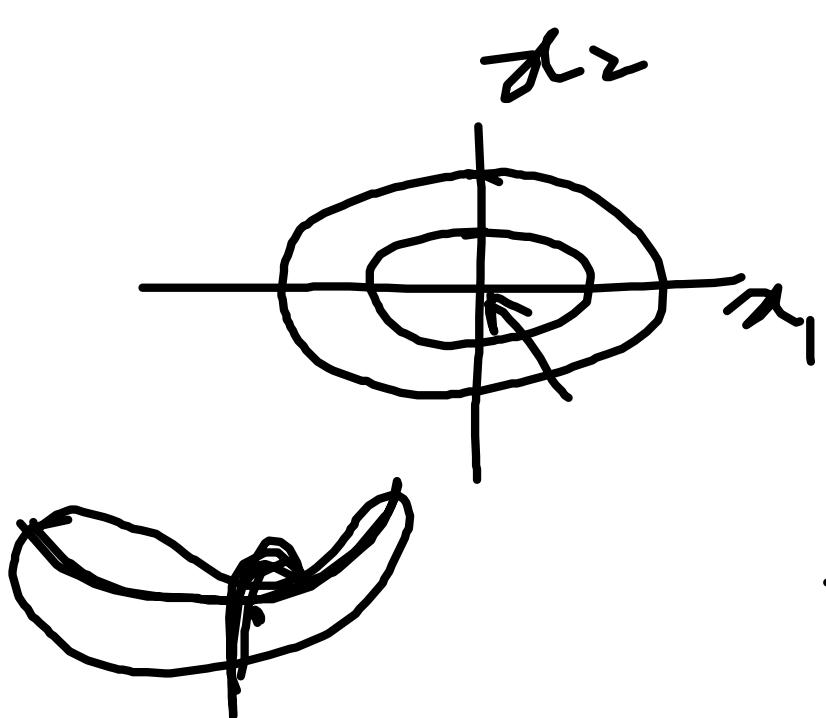
Direction of steepest

ascend

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \leftarrow \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - n \begin{bmatrix} \frac{\partial y / \partial x_1}{\partial y / \partial x_2} \end{bmatrix}$$



Maxima and minima for multivariate functions



$$y = ax_1^2 + bx_2^2$$

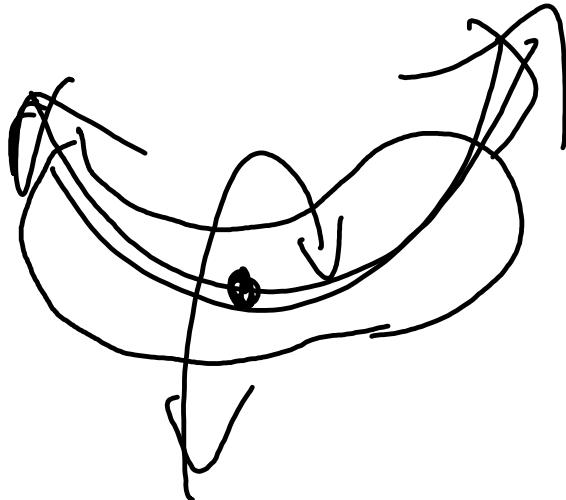
$$\nabla y = 0 = \begin{bmatrix} 2ax_1 \\ 2bx_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

=====

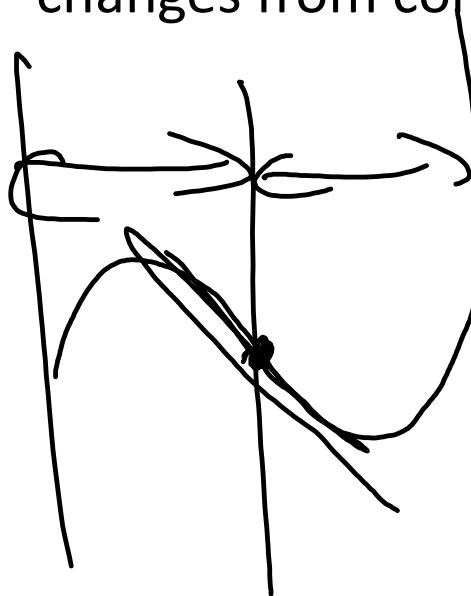
If $a > 0 ; b < 0$
 $\Rightarrow (0, 0)$ is max for x_2
min for x_1
Saddle point

When functions are neither convex nor concave

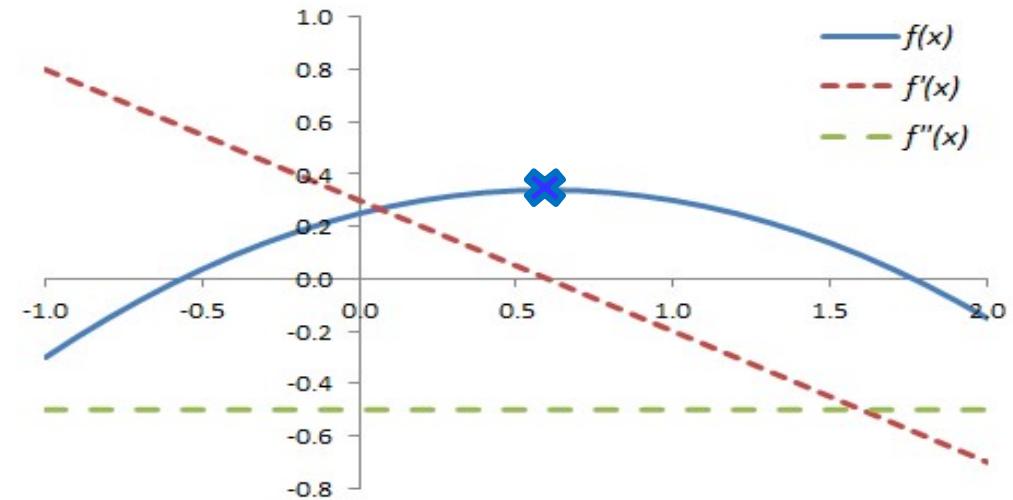
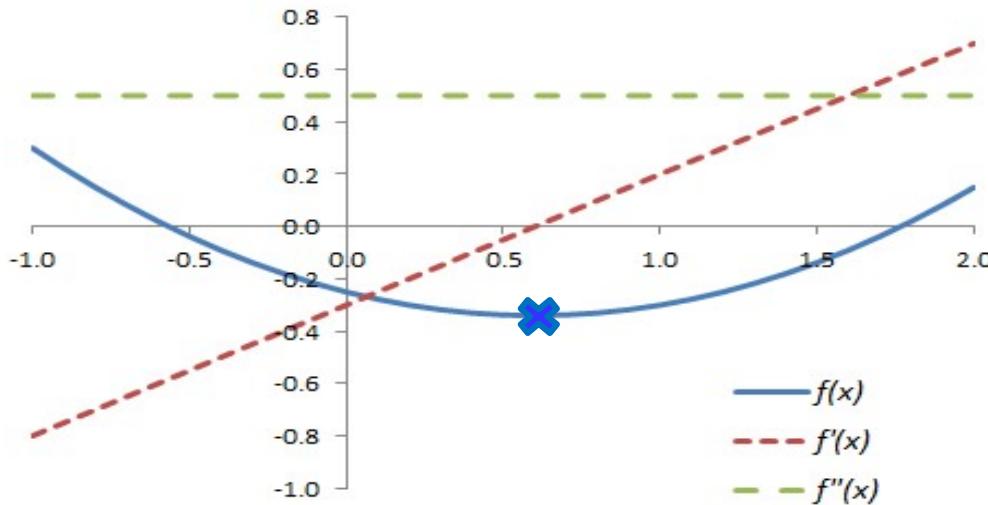
Saddle point
where derivative is zero
and one direction is convex,
and the other in concave



Inflection point
where second derivative
is zero and the function
changes from convex to concave



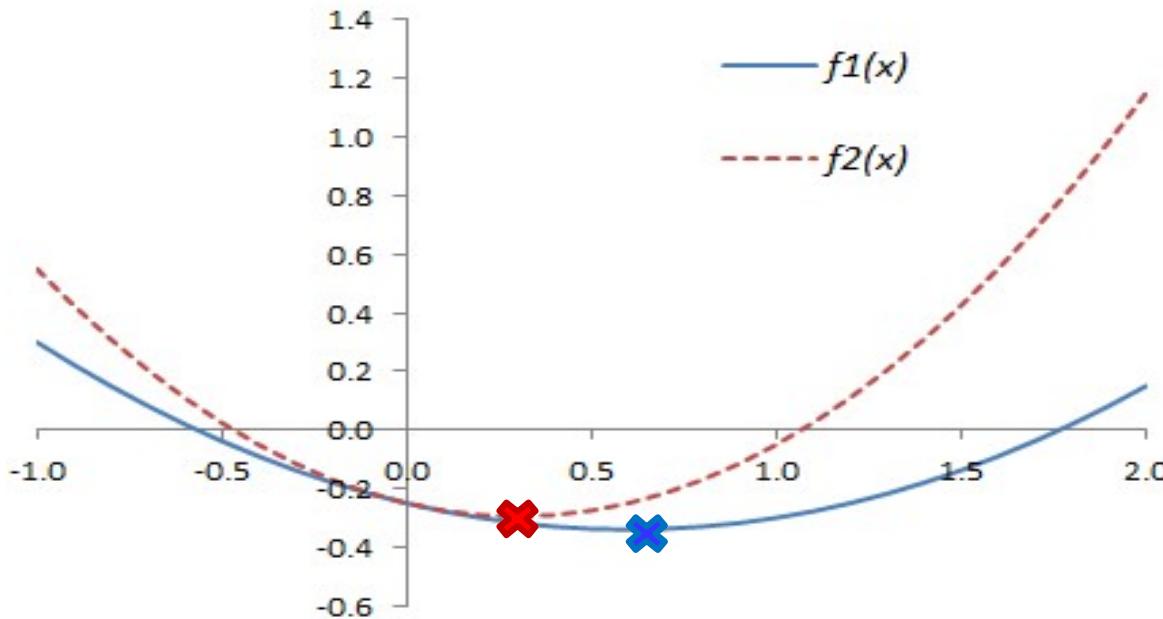
Double derivative



E.g. $f(x) = ax^2 + bx + c$, $f'(x) = 2ax + b$, $f''(x) = 2a$

- Double derivative $f''(x) = \frac{d^2 f(x)}{d x^2}$ is the derivative of derivative of $f(x)$
- Double derivative is positive for convex functions (have a single minima), and negative for concave functions (have a single maxima)

Double derivative



$$\begin{aligned}f(x) \\= ax^2 + bx + c, \\f'(x) = 2ax + b, \\f''(x) = 2a\end{aligned}$$

- Double derivative tells how far the minima might be from a given point.
- From $x = 0$ the minima is closer for the red dashed curve than for the blue solid curve, because the former has a larger second derivative (its slope reverses faster)

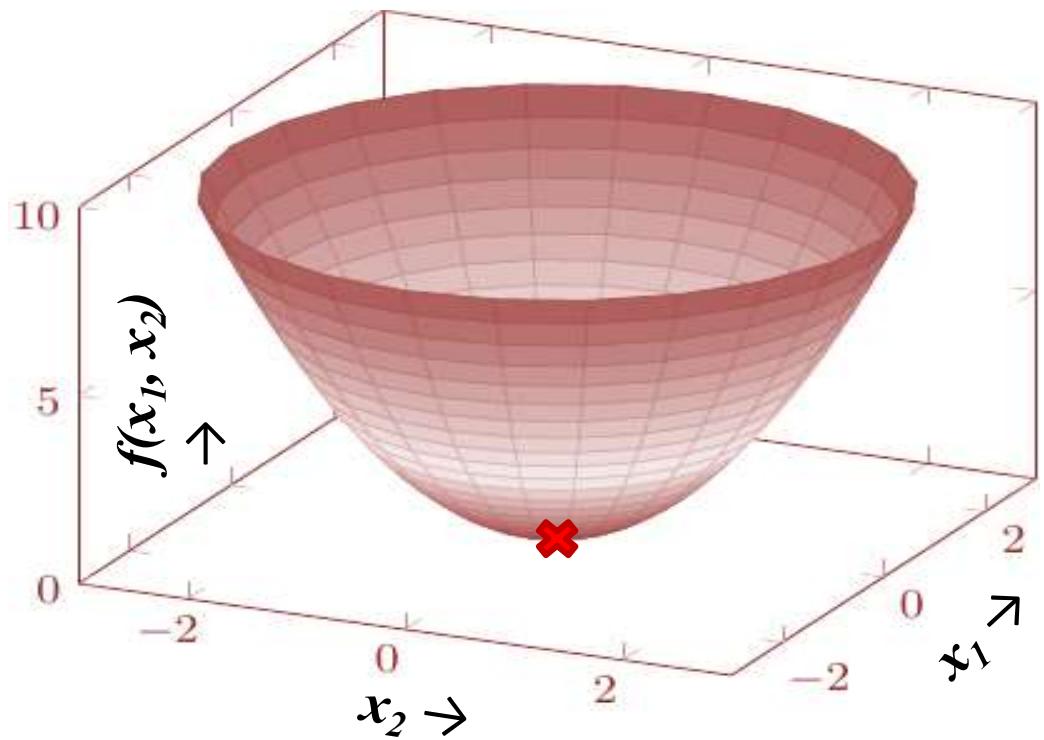
Perfect step size for a paraboloid

- Let $f(x) = ax^2 + bx + c$
- Assuming $a < 0$
- Minima is at: $x^* = -\frac{b}{2a}$
- For any x the perfect step would be:

$$-\frac{b}{2a} - x = -\frac{2ax+b}{2a} = -\frac{f'(x)}{f''(x)}$$

- So, the perfect learning rate is: $\eta^* = \frac{1}{f''(x)}$
- In multiple dimensions, $x \leftarrow x - H(f(x))^{-1} \nabla(f(x))$
- Practically, we do not want to compute the inverse of a Hessian matrix, so we approximate Hessian inverse

Hessian of a function of a vector



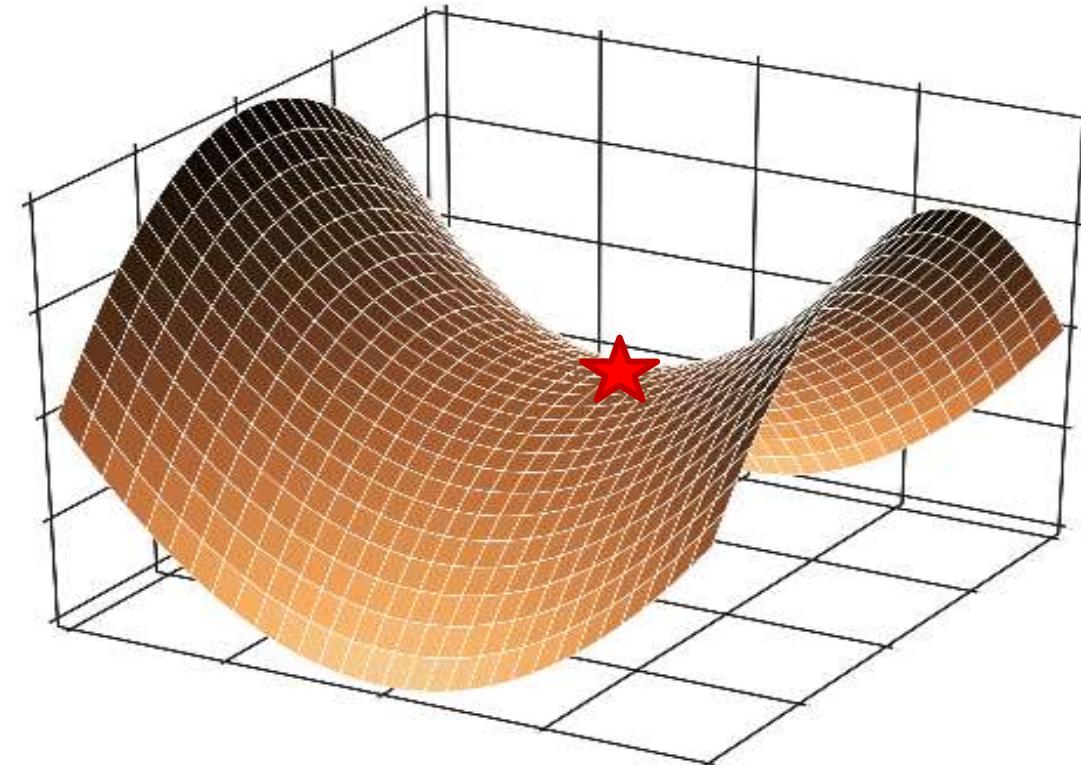
- Double derivative with respect to a pair of dimensions forms the Hessian matrix:
- If all eigenvalues of a Hessian matrix are positive, then the function is convex

Example of Hessian

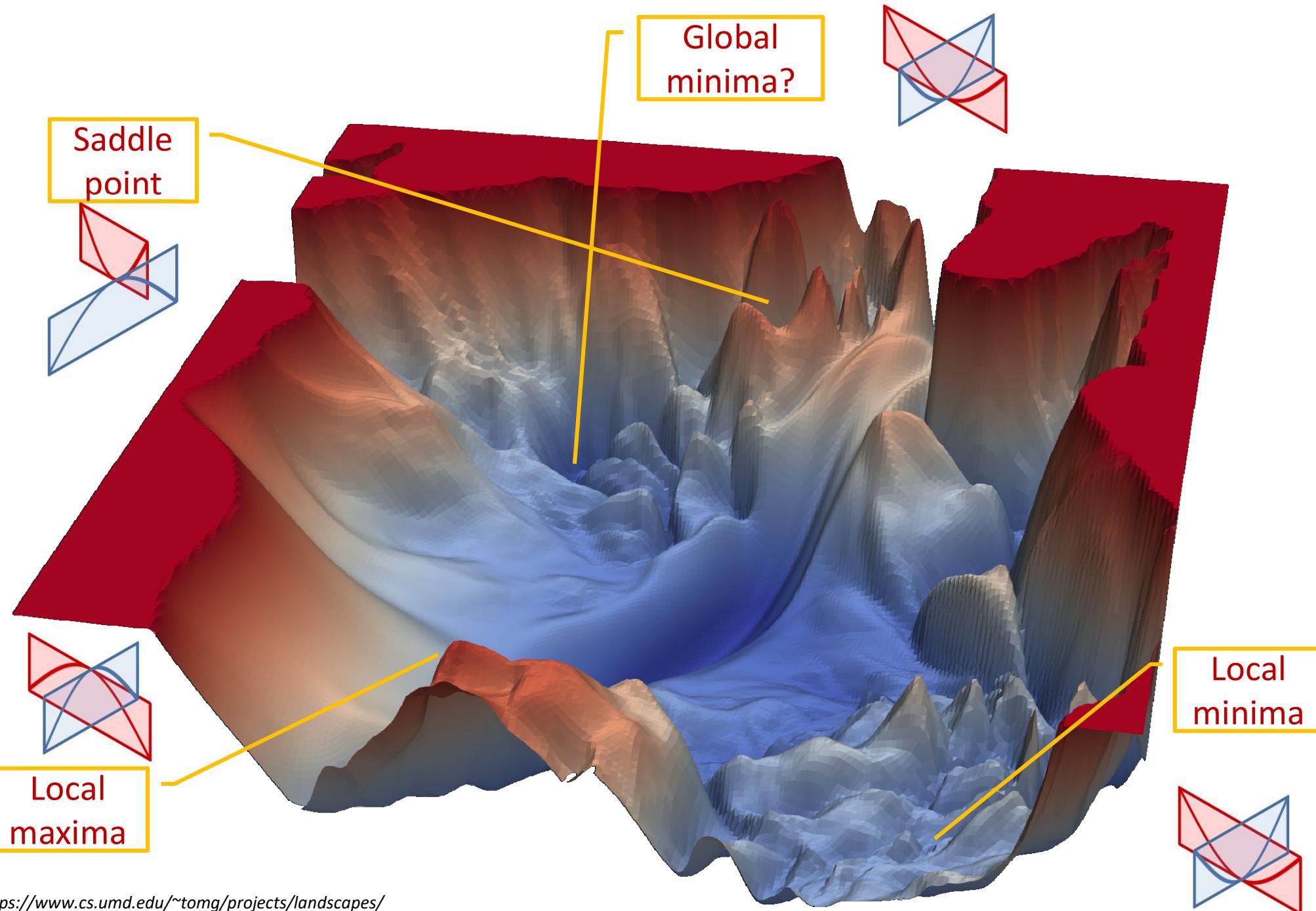
- Let $f(\mathbf{x}) = f(x_1, x_2) = 5x_1^2 + 3x_2^2 + 4x_1x_2$
- Then $\nabla f(\mathbf{x}) = \nabla f(x_1, x_2) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 10x_1 + 4x_2 \\ 6x_2 + 4x_1 \end{bmatrix}$
- And, $H(f(\mathbf{x})) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 10 & 4 \\ 4 & 6 \end{bmatrix}$

Saddle points, Hessian and long local furrows

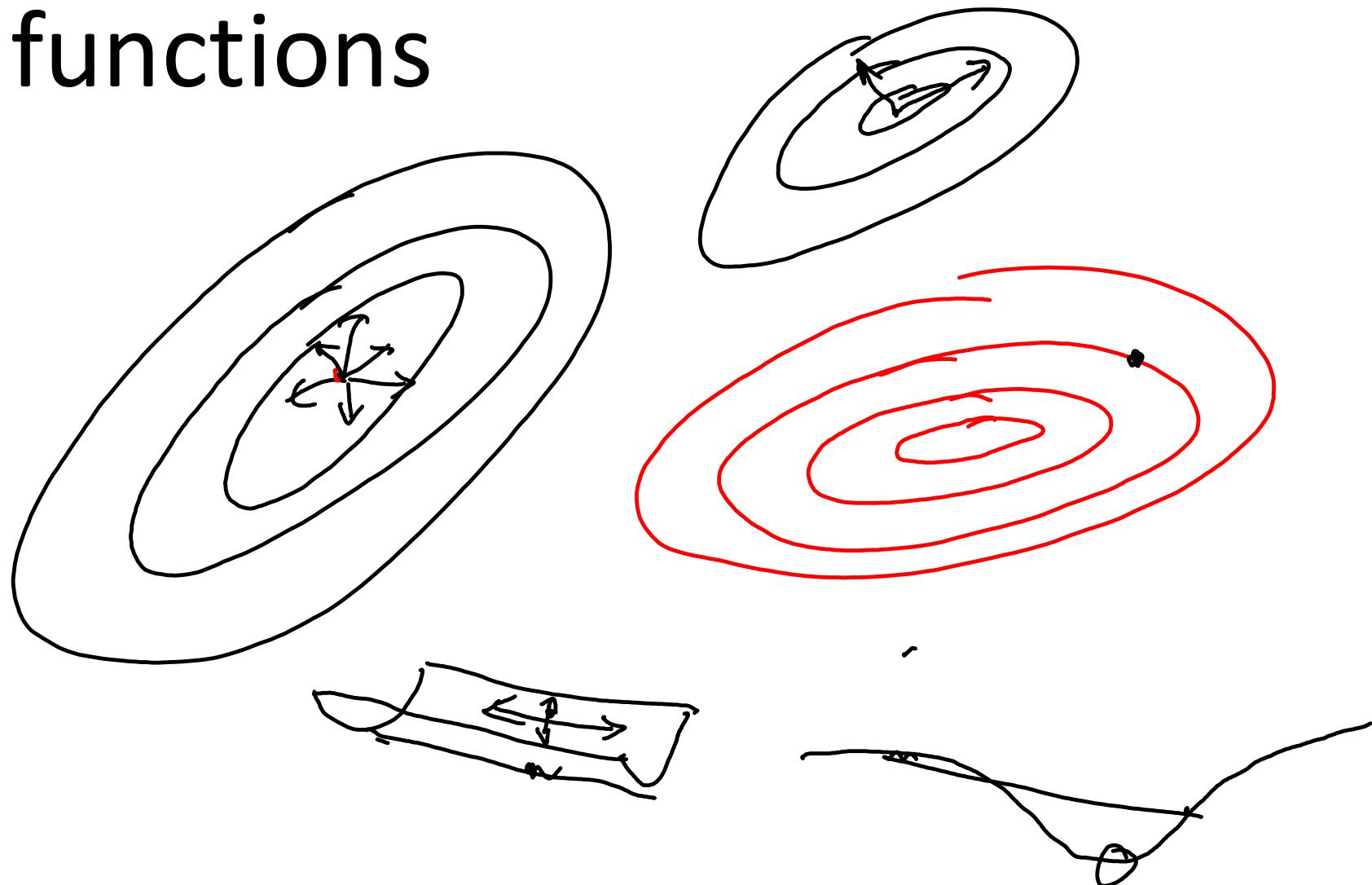
- Some variables may have reached a local minima while others have not
- Some weights may have almost zero gradient
- At least some eigenvalues may not be negative



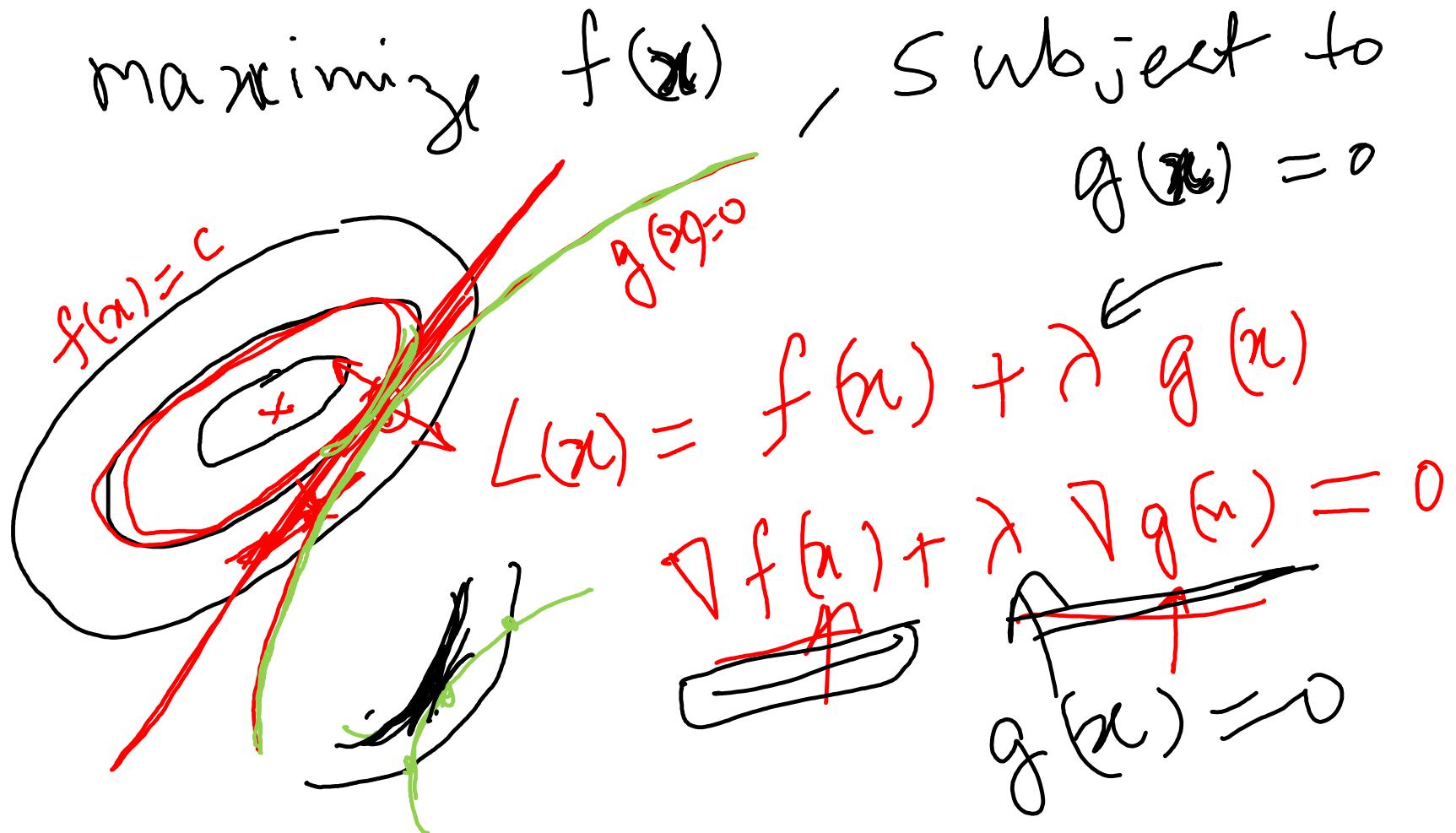
A realistic picture



Some functions

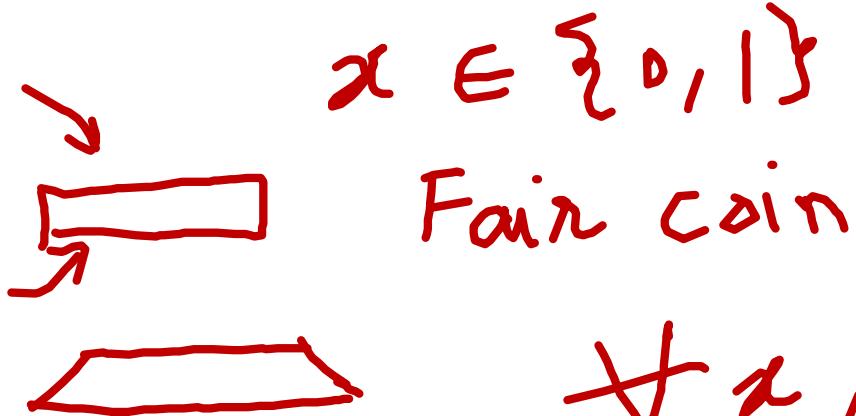


Constrained optimization using Lagrange multiplier



Random variable

$$X = x \quad \{H, T\} \quad M.E., C.E$$



$$x \in \{0, 1\}$$

Fair coin

$$\begin{aligned} P(X=0) &= 1 - P(X=1) \\ &= 0.5 \end{aligned}$$

$$\forall x, P(X=x) \geq 0$$

$$\sum_x P(X=x) = 1$$

Random variable

$X = x$ $\{H, T\}$ Event space

$X \in \{0, 1\}$ M.E.C.E.

Fair coin



0	1
---	---

$$P(X=0) = 1 - P(X=1) = 0.5$$

Biased coin

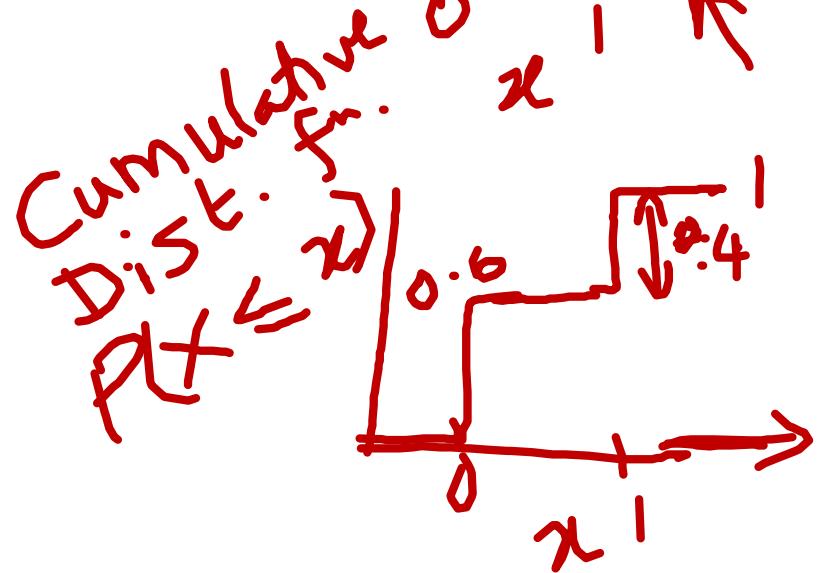
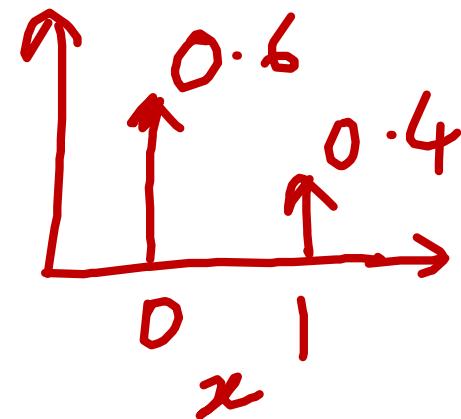
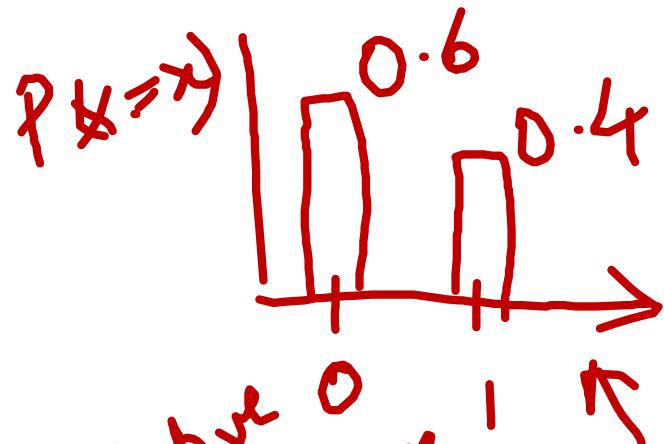


$$P(X=0) = 0.6$$

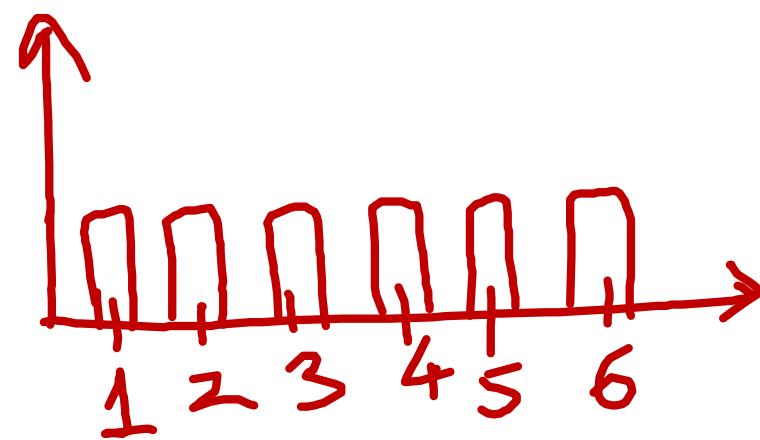
$$x_n, P(X=x) \geq 0 \quad P(X=1) = 0.4$$

$$\sum_x P(X=x) = 1$$

Probability mass function



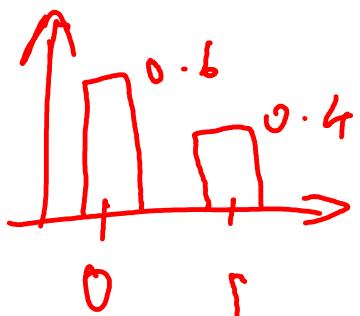
6-faced dice



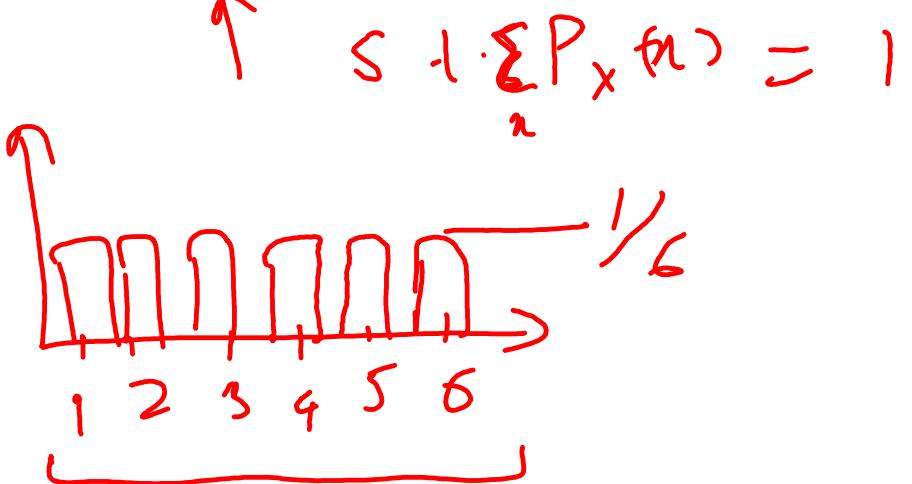
Probability mass function

$$P(X=x)$$

Discrete R.V.



$$P_x(x) : x \rightarrow [0, 1]$$



$$\sum_{x \in S} P_x(x) = 1$$

Some common PMFs

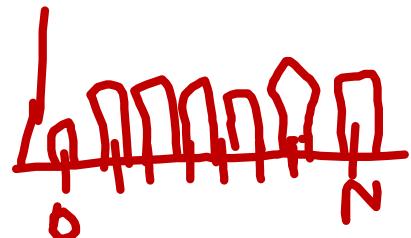
Bernoulli Dist.

$$x \in \{0, 1\}$$

$$\mu = P(X=1)$$

$$\text{Bern}(x|\mu) = \mu^x (1-\mu)^{1-x} = \begin{cases} \mu, & \text{if } x=1 \\ 1-\mu, & \text{if } x=0 \end{cases}$$

Binomial Dist.



$$\mathbb{E}(x) = \sum_x x P(x=x)$$

$$\mathbb{E}(x - \bar{x})^2 = \mu(1-\mu)$$

$$\text{Bin}(y|\mu, N) = \binom{N}{y} \mu^y (1-\mu)^{N-y} \mathbb{E}[f(x)]$$

q

$$= \sum_x f(x)(x=x)$$

Some common PMFs

①

Bernoulli

$$x \in \{0, 1\}$$

$$\mu = P(x=1) \quad \mu \in [0, 1]$$

$$\text{Bern}(x|\mu) = \mu^x (1-\mu)^{1-x}$$

given $\begin{cases} \mu, & \text{if } x=1 \\ 1-\mu, & \text{if } x=0 \end{cases}$

Variance of x

$$= \mathbb{E}(x - \mathbb{E}(x))^2$$

$$= \mathbb{E}(x - \mu)^2$$

$$= \mathbb{E}(x^2 - 2x\mu + \mu^2)$$

$$= \mathbb{E}(x^2) - 2\mathbb{E}(x)\mu + \mu^2$$

$$= \mu(1-\mu)$$

$$\mathbb{E}(x) = \sum_x x P(x) = \underline{\underline{0 \cdot (1-\mu) + 1 \cdot \mu}} = \mu$$

$$\mathbb{E}(f(x)) = \sum_x f(x) P(x) = f(0)(1-\mu) + f(1)\mu$$

$$\mathbb{E}(x^2) = 0(1-\mu) + 1(\mu) = \mu$$

Binomial distribution

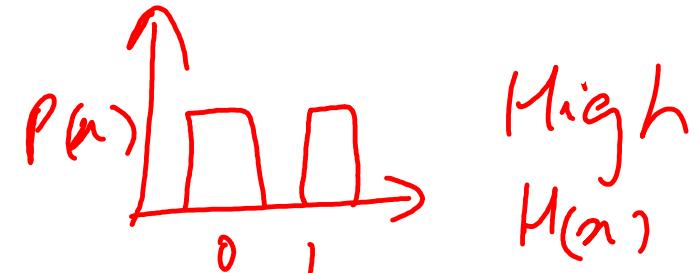
$$\begin{aligned}
 & \mu, N \quad \text{I.I.P.} \\
 & \text{Bin}(y | \mu, N) \\
 & = \binom{N}{y} \mu^y (1-\mu)^{N-y} \\
 & \quad \overbrace{\hspace{10em}}^{N \times y} \quad \text{Diagram showing a binomial distribution curve with bars at } y=0, 1, 2, 3 \\
 & E(y) = N\mu
 \end{aligned}$$

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

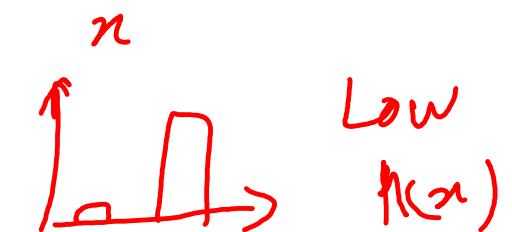
Entropy

Entropy of a random variable

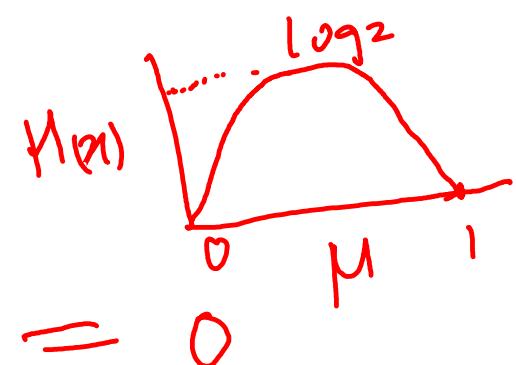
$$H(x) = -E[\log(P(x))]$$



$$= - \sum_x \log[P(x)] P(x)$$



$$= - [(1-\mu)\log(1-\mu) + \mu \log \mu]$$



$$\mu = 0.5 \Rightarrow H(x) = -\log\left(\frac{1}{2}\right) = \log 2$$

Let $\mu \rightarrow 0$

$$H(x) = -[(1-\mu)\log(1-\mu) + \mu \log \mu] = 0$$

Joint, conditional, and marginal probability

Dice : $x \in \{1, 2, 3, 4, 5, 6\}$

~~Joint~~ $y \in \{0, 1\}$, $0 \Rightarrow \text{even}$

~~Joint~~ $z \in \{0, 1\}$, $x < 3$

Conditional

$$P(y|z) = \frac{P(y,z)}{P(z)}$$

$P(y,z) \neq 0$

		$z=0$	$z=1$	
		$y=0$	$y=1$	
y	$z=0$	$\frac{1}{6}$	$\frac{2}{6} = \frac{1}{3}$	$P(z)$
	$z=1$	$\frac{1}{6}$	$\frac{2}{6} = \frac{1}{3}$	

$$\sum_y \sum_z P(y,z) = 1$$

Marginal Dist

$$P(B) = \sum_y P(y,z)$$

$$P(y) = \sum_z P(y,z)$$

$$P(y|z=0)$$

$$P(y=0|z=0)$$

$$= \frac{1/6}{1/3} = \frac{1}{2}$$

Joint, conditional, and marginal probability

Dice $x \in \{1, 2, 3, 4, 5, 6\}$

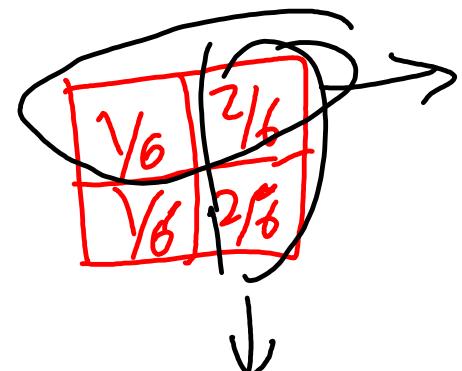
$y \in \{0, 1\}$ $0 \Rightarrow$ Even, 1 means odd

$z \in \{0, 1\}$ $0 \Rightarrow x < 3$, $1 \Rightarrow x \geq 3$

$$P(y, z)$$

$$\sum_y \sum_z P(y, z) = 1 \quad P(y, z) \in [0, 1]$$

		$z=0$	$z=1$
$y=0$	$x \in \{1, 2\}$	$\frac{1}{6}$	$\{4, 6\} \frac{2}{6}$
	$x \in \{1\}$	$\frac{1}{6}$	$\{3, 5\} \frac{2}{6}$



$P(y, z)$	$y=0$	$z=0$	$z=1$
	$\cancel{1/6}$	$\cancel{2/6}$	$\cancel{2/6}$
	$\cancel{1/6}$		
$y=1$			

$$P(y|z)$$

↑
given

Conditional

$$\text{Marginal } P(z) = \sum_y P(y, z) = \frac{2}{6} + \frac{2}{6} = \frac{4}{6}$$

$$P(x_{\text{Even}} | x \geq 3) = \frac{2/6}{2/6 + 2/8} = \frac{1}{2}$$

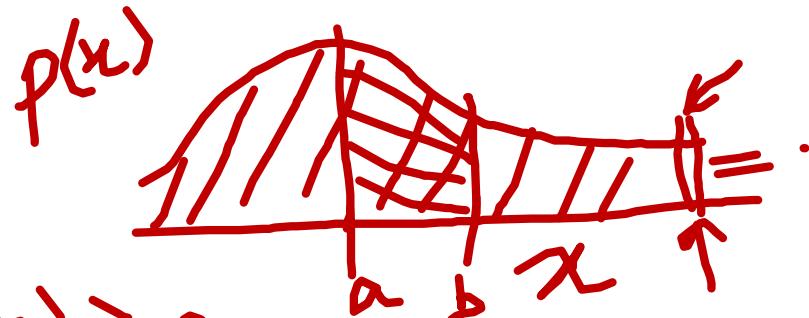
$$P(y|z) = \frac{P(y, z)}{P(z)} = \frac{2/6}{4/6} = \frac{1}{2}$$

$$P(y, z) = P(y|z) \underbrace{P(z)}$$

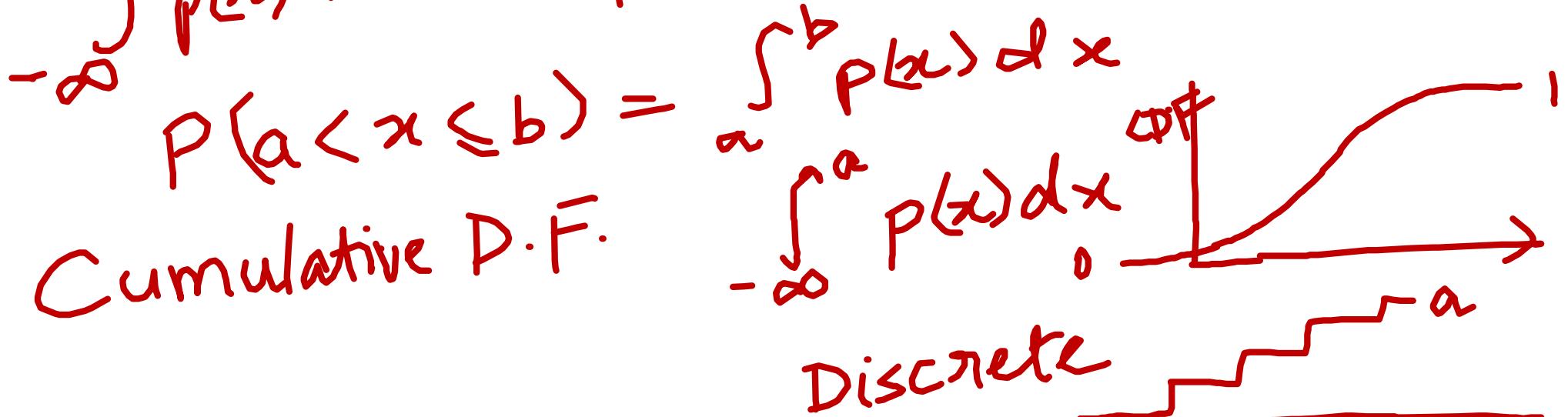
Continuous random variable and PDF

$$P(X=x) = 0 \quad x = 170 \text{ cm}$$

$$p(x)$$



$$\int_{-\infty}^{\infty} p(x) dx = 1 \quad p(x) > 0$$



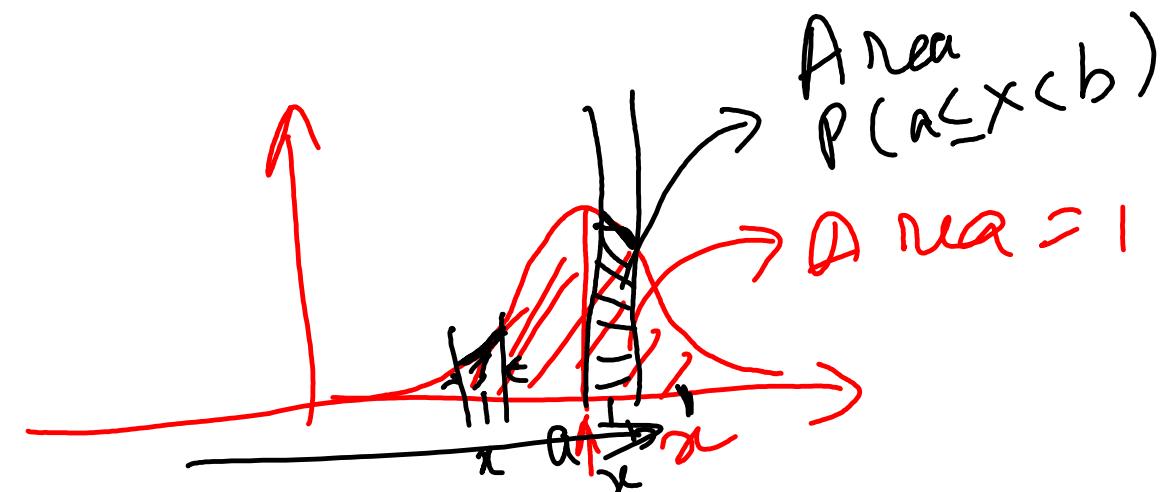
Cumulative D.F.

Continuous random variable and PDF

$$P(x=x) = 0$$

Prob. density fn.

$$p(x) \geq 0 \quad \int_{-\infty}^{\infty} p(x) = 1$$



$$170 \text{ cm} \quad 169.99 \\ 170.01$$

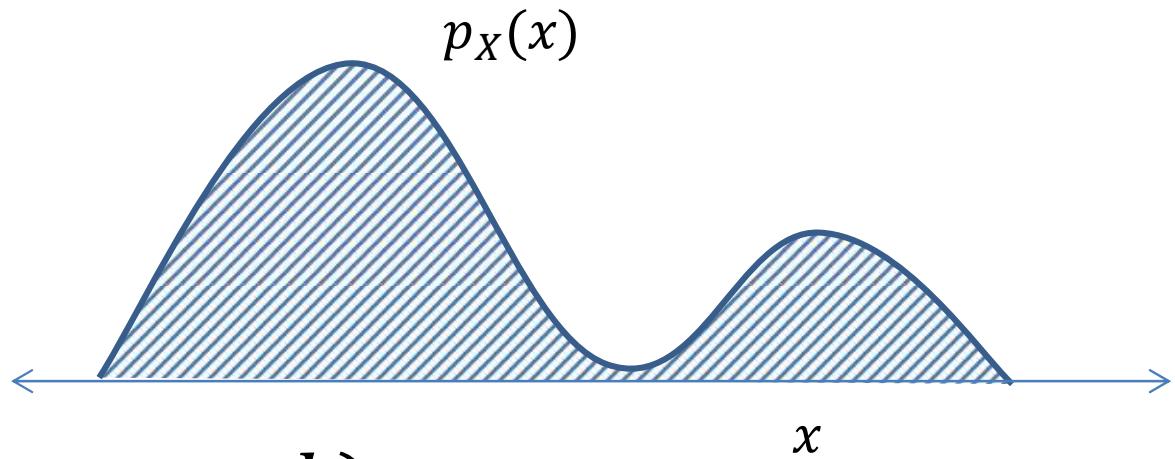
Cont. R.V. Cumulative DF = $\int_{-\infty}^x p_x(y) dy$

Discrete "



Probability density function (PDF)

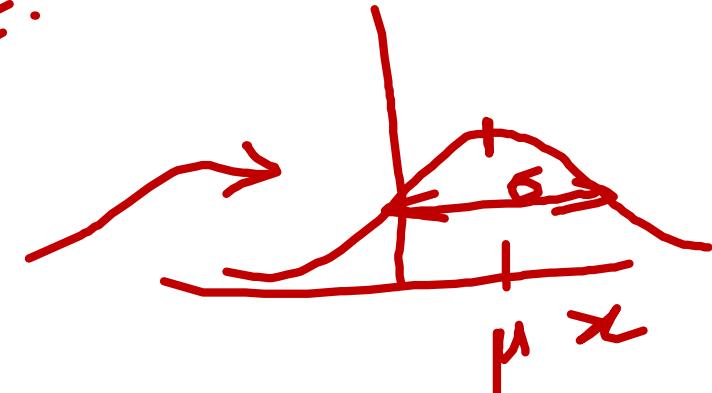
- PDF: $p_X(x) \geq 0$
- $\int_a^b p_X(\hat{x})d\hat{x} = Pr(a \leq x \leq b)$
- $\int_{-\infty}^{+\infty} p_X(\hat{x})d\hat{x} = 1$
- CDF: $P_X(x) = \int_{-\infty}^x p_X(\hat{x})d\hat{x}$
- $p_X(x) = \frac{dP_X(x)}{dx}$



Some common PDFs

Gaussian / Normal Dist.

$$N(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



$$E(x) = \int_{-\infty}^{\infty} x p(x) dx = \mu$$

$$E[(x-\mu)^2] = \sigma^2$$

Beta dist. $x \in [0, 1]$

Beta($x ; a, b$)

$\text{Unif}(x ; a, b)$

$$= \frac{(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$$

Some common PDFs

Gaussian / Normal

$$N(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

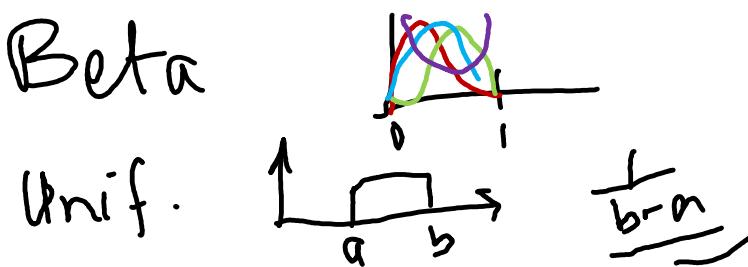
$$H(x) = \int_{-\infty}^{\infty} (\log p(x)) p(x) dx$$



$$\mathbb{E}[f(x)] = \int_{-\infty}^{\infty} f(x) p(x) dx$$

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} x p(x) dx$$

Beta



Unif.

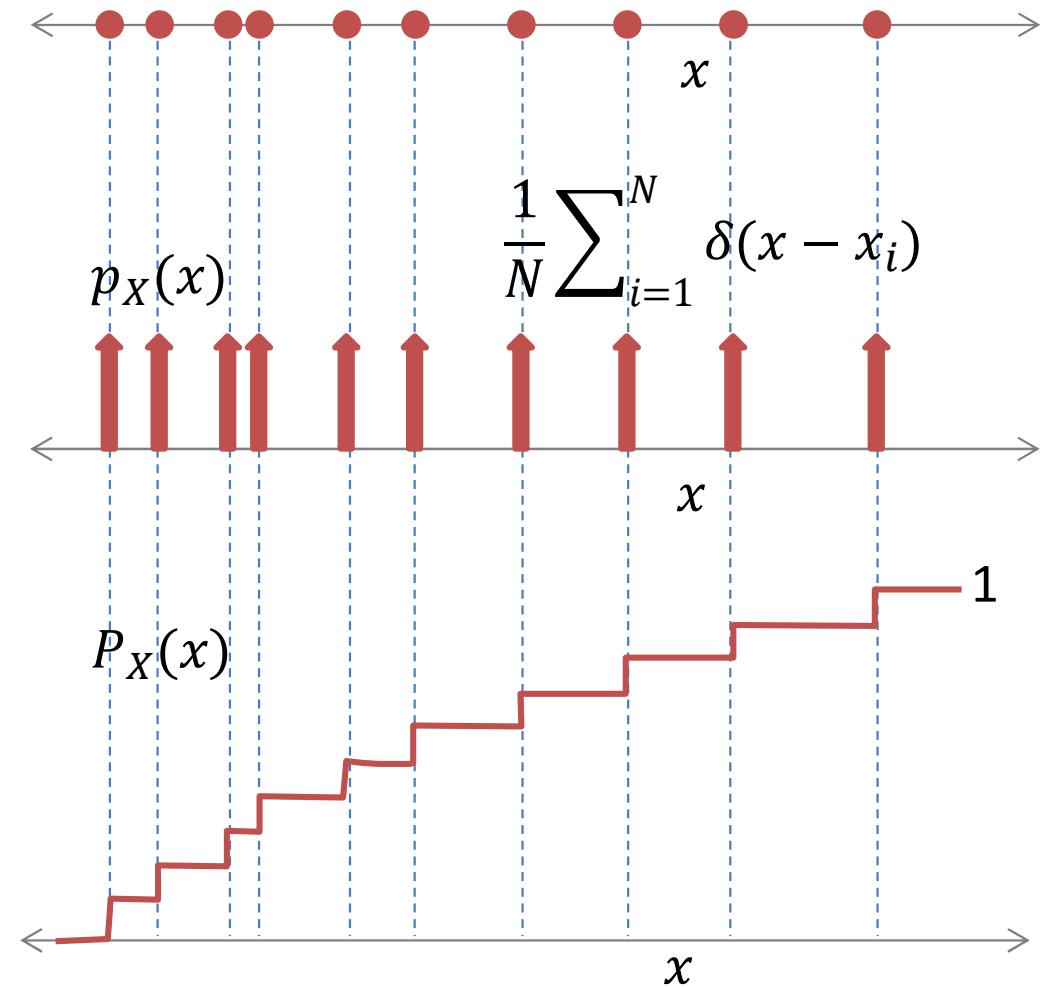
a, b

$$\text{Beta}(x|a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$$

for $0 \leq x \leq 1$
0 otherwise

Empirical distribution

- Given samples
- What is the “best” explanation of the data?
- Empirical PDF (train of Dirac delta functions scaled $1/N$)
- CDF has N steps



Problem with empirical distribution

- Rote learning
- Does not generalize
- If we randomly split the sample into training and validation
- How likely is the validation subset based on the distribution learnt from the training subset?

Between two distributions, which is better?

- The one under which the data has higher likelihood
- If $p_A(x_i) > p_B(x_i)$ then p_A is a better distribution to explain the observation x_i
- This can be maximized by having a Dirac delta function $\delta(x - x_i)$
- What about for the entire data sample X ?

The IID Assumption

- Independence of x_i and x_j for $i \neq j$.
- $p(x_i, x_j) = p(x_i)p(x_j)$
- All samples drawn from the same (identical) distributed
- $x_i, x_j \sim p_X$

Implication of the IID assumption

- Likelihood of the whole data factorizes:
- $$\begin{aligned} p(X) &= p(x_1, \dots, x_N) = p_X(x_1) \times \dots \times p_X(x_N) \\ &= \prod_{i=1}^N p_X(x_i) \end{aligned}$$
- $L(X) = \log p(X) = \sum_i \log p_X(x_i)$

MLE of parameterized distribution

- Between two distributions p_A is a better explanation than p_B of the entire data X if $\prod_i p_A(x_i) > \prod_i p_B(x_i)$
- By extension, if a family of distributions is parameterized by θ , then we are interested in
$$\begin{aligned}\arg \max_{\theta} \prod_i p_{\theta}(x_i) &= \arg \max_{\theta} \sum_i \log p_{\theta}(x_i) \\ &= \arg \max_{\theta} L_{\theta}(X)\end{aligned}$$

Example 1: Exponential distribution

- $p_\lambda(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$
- Log likelihood $L_\lambda(X) = \sum_i \log p_\lambda(x_i)$
- Differentiating $\frac{\partial L_\lambda(X)}{\partial \lambda} = \sum_{x_i \geq 0} \left(\frac{\partial \log \lambda}{\partial \lambda} - \frac{\partial \lambda x_i}{\partial \lambda} \right) = 0$
- (assuming all samples are non-negative)
- Implies $\lambda = \frac{N}{\sum_i x_i} = \text{inverse of the sample mean}$

Example 2: Uniform distribution

- $p_{a,b}(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{otherwise} \end{cases}$
- Log likelihood $L_{a,b}(X) = \sum_i \log p_{a,b}(x_i)$
- $= K \log(0) - (N - K) \log(b - a)$
- Reduce the K (terms outside $[a, b]$) to zero, and minimize $(b - a)$ by differentiating wrt a, b
- So, $a = \min x_i, b = \max x_i$

Example 3: Gaussian distribution

- $p_{\mu,\sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$
- MLE by differentiating log likelihood wrt μ, σ
- gives $\mu = \bar{x} = \frac{\sum_i x_i}{N}$; i.e., sample mean
- and, $\sigma = \sqrt{\frac{\sum_i (x_i - \bar{x})^2}{N}}$

Sufficient statistics

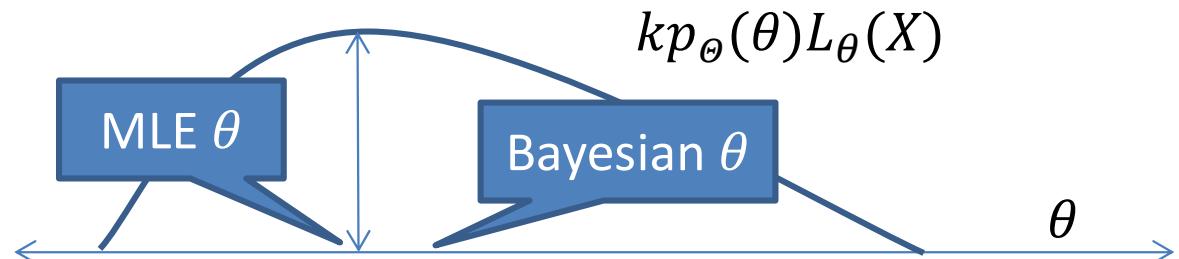
- Statistic is a function of the sample
$$T(X) = T(x_1, \dots, x_N)$$
- For some distributions, computing a few statistics is sufficient for MLE estimate
- Gives complete information about the distribution
- Examples:
 - Sample mean and variance for Gaussian distribution
 - Sample mean for exponential distribution
 - Max and min for uniform distribution

Comparing two parametric distributions

- Let there be two candidate families of distributions $p_\theta(X)$ and $p_\phi(X)$ to explain the data
- Can we compare $\max L_\theta(X)$ and $\max L_\phi(X)$?
- Yes, we can, but we might overfit
- Narrow down the family of distributions based on domain knowledge (e.g. physical phenomenon)
- E.g. “Can the random variable take negative values?”

MLE vs. Bayesian estimate

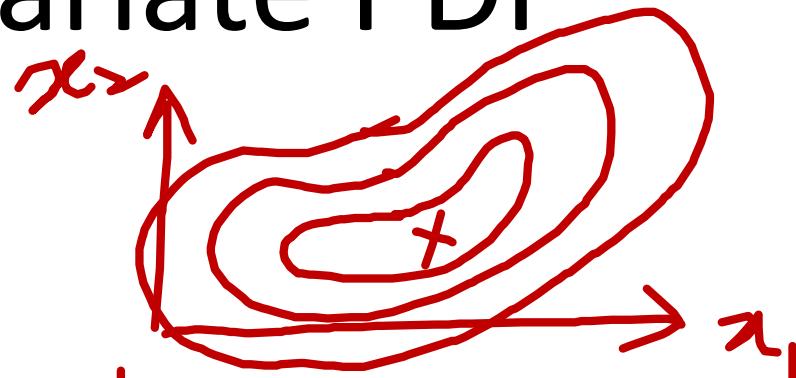
- MLE finds θ that maximizes $L_\theta(X)$
- Bayesian estimate takes the expected value of θ w.r.t. $L_\theta(X)$
- Bayesian estimate: $\int \theta L_\theta(X) d\theta / \int L_\theta(X) d\theta$
- We can also incorporate a prior over θ
- $\int \theta p_\theta(\theta) L_\theta(X) d\theta / \int p_\theta(\theta) L_\theta(X) d\theta$



Multivariate PDF

$$p(x_1, x_2) \geq 0$$

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p(x_1, x_2) dx_1 dx_2 = 1$$



$$\int_{-\infty}^{\infty} p(x_1, x_2) dx_1 = p(x_2)$$

$$p(x_2 | x_1 = a) = \frac{p(x_1 = a, x_2)}{p(x_1 = a)}$$

Multivariate PDF

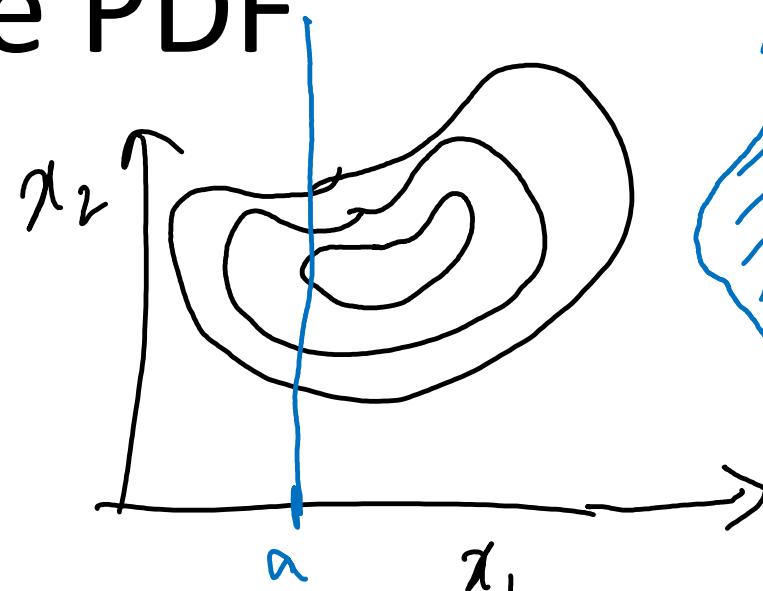
$$P(x_1, x_2)$$

$$P(x_1, x_2) \geq 0$$

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p(x_1, x_2) dx_1 dx_2 = 1$$

$$p(x_1) = \int_{-\infty}^{\infty} p(x_1, x_2) dx_2$$

marginal



$$P(\mathbf{x})$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$P(x_2 | x_1 = a)$$

$$= \frac{P(x_1 = a, x_2)}{P(x_1 = a)}$$

Condition

Multivariate Gaussian

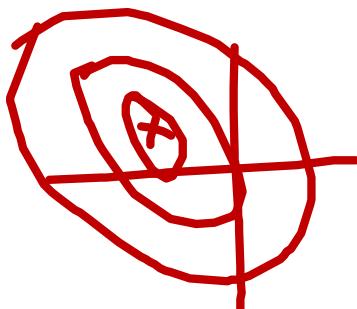
$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$p(x_1, x_2)$$

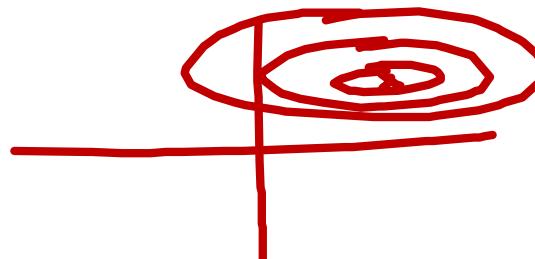
$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \sigma_{11}^2 & \sigma_{12} \\ \sigma_{21} & \sigma_{22}^2 \end{bmatrix}$$

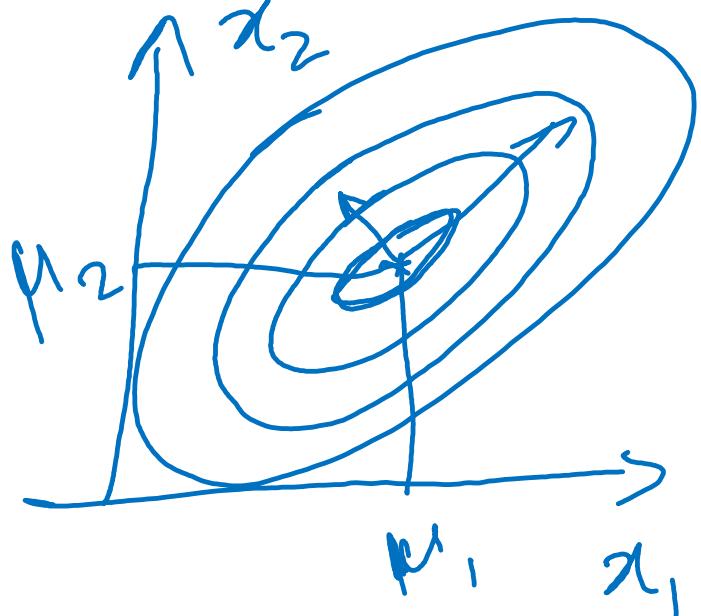
$$p(x) = \frac{1}{(2\pi)^K \text{Det}(\Sigma)}$$



$$p(x) = \frac{1}{(2\pi)^K \text{Det}(\Sigma)} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$



Multivariate Gaussian PDF



$$p(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^k \det(\Sigma)}} \times e^{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})}$$

where $\boldsymbol{\mu} = [\mu_1, \mu_2]^T$ and $\Sigma = \begin{bmatrix} \sigma_{11}^2 & \sigma_{12} \\ \sigma_{21} & \sigma_{22}^2 \end{bmatrix}$

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \sigma_{11}^2 & \sigma_{12} \\ \sigma_{21} & \sigma_{22}^2 \end{bmatrix}$$

