

# REVIEW

## MATRIX CALCULUS

FUNCTIONS

$$f: \mathbb{R} \rightarrow \mathbb{R}$$

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$f: \mathbb{R}^d \rightarrow \mathbb{R}^p$$

EXAMPLE

$$x^2$$

LOSS

PROJECTION/  
NEURALNET  
- LAYER

VALUE

$$\mathbb{R}$$

$$\mathbb{R}$$

$$\mathbb{R}^p$$

FIRST DER

$$\mathbb{R}$$

$$\mathbb{R}^d$$

(GRADIENT)

$$\mathbb{R}^{d \times p}$$

(JACOBIAN)

SECOND DER

$$\mathbb{R}$$

$$\mathbb{R}^{d \times d} \quad (\mathbb{S}^d)$$

(HESSIAN)

$$\mathbb{R}^{d \times p \times p}$$

(HIGHER DIM TENSOR)

GRADIENT

$$\nabla_x f(\bar{x}) = \begin{bmatrix} \frac{\partial f(\bar{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_d} \end{bmatrix}$$

$$\bar{x} = (x_1, \dots, x_d)$$

DIRECTION TO INCREASE  
THE VALUE OF  $f(x)$  THE  
MOST

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial a_{11}}(A) & \dots & \\ & \ddots & \\ & & \frac{\partial f(A)}{\partial a_{mn}} \end{bmatrix}$$

GENERALIZATION

## HESSIAN

$$\nabla_x^2 f(x) \quad f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\nabla_x^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1 \partial x_1} & \dots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_d \partial x_1} & \dots & \frac{\partial^2 f(x)}{\partial x_d \partial x_d} \end{bmatrix}$$

Ex

$$\nabla_x b^T x = \begin{bmatrix} \vdots \\ \frac{\partial b^T x}{\partial x_i} \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ \frac{\partial}{\partial x_i} (b_1 x_1 + b_2 x_2 + \dots + b_d x_d) \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ b_i \\ \vdots \end{bmatrix} = b$$

## PRODUCT RULE (MULTIVARIABLE)

$$\nabla_x x^T A x = \nabla_x \cancel{x}^T A x + \nabla_x x^T A \cancel{x} = A x + A^T x = 2 A x \text{ (if } A = A^T \text{)}$$

$$\nabla_A \log \overbrace{|A|}^{\det A} = A^{-1}$$

## PROBABILITY THEORY

### ELEMENTS

→ SAMPLE SPACE  $\Omega$  ex  $\{HH, HT, TH, TT\}$

→ EVENT  $A \subseteq \Omega$  (SAMPLE SET)

→ EVENT SPACE SET OF ALL SUBSET  
 $F = \{A_1, \dots, A_n\}$

→ PROB MEASURE  $P: F \rightarrow \mathbb{R} \ (0, 1)$

$$P(A) \geq 0 \quad \forall A \in F$$

$$P(\Omega) = 1$$

IF  $A_1, A_2, \dots$  DISJOINT SET OF EVENTS ( $A_i \cap A_j = \emptyset$  WHEN  $i \neq j$ )  
THEN

$$P(\cup A_i) = \sum P(A_i)$$

### CONDITIONAL

LET  $B$  ANY EVENT SUCH THAT  $P(B) \neq 0$

$$P(A|B) := \frac{P(A \cap B)}{P(B)} \rightarrow \text{INTERSECTION}$$

- $A \perp B$  IF AND ONLY IF  $P(A \cap B) = P(A)P(B)$  ( $\perp$  INDEPENDENT)
- $A \perp B$  "  $P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A)P(B)}{P(B)} = P(A)$   
THE  $P(A)$  DOESN'T CHANGE WHETHER  $P(B)$  OCCURS OR NOT

### RANDOM VARS

$$\omega_0 = HHTHTTTHTT$$

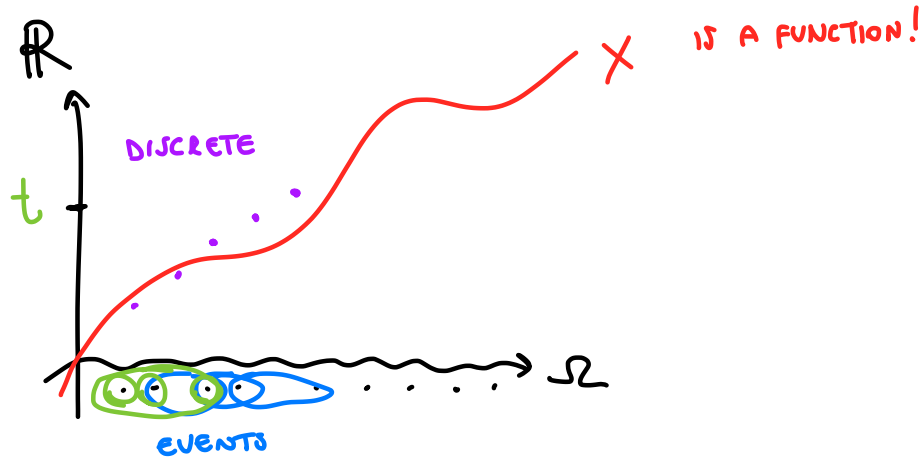
A RV IS  $X: \Omega \rightarrow \mathbb{R}$  MAPPING FROM OUTCOME TO  $\mathbb{R}$

$$\# \text{ OF HEADS } X(\omega_0) = 5$$

# OF TOSSES UNTIL TAILS:  $X(\omega_0) = 4$

$$VAL(X) := X(\Omega)$$

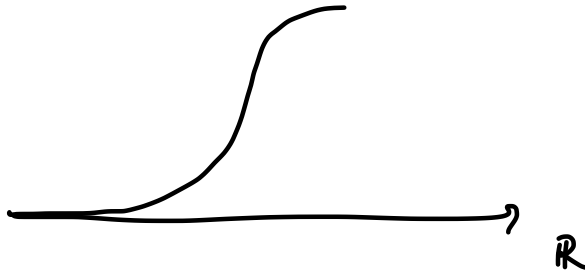
Ex



# CUMULATIVE DISTRIBUTION FUNCTION (CDF)

$$F_x(x) = P(X \leq x)$$

$$P[\{\omega: x(\omega) \leq t\}] \quad \text{MEASURE PROB. ON } \Omega$$



DISCRETE RV:  $\text{VAL}(X)$  COUNTABLE

$$P(X=k) := P(\{\omega | X(\omega) = k\})$$

PROBABILITY MASS FUNCTION

$$p_X: \text{VAL}(X) \rightarrow [0, 1]$$

$$p_X(x) := P(X=x)$$

$$\sum_{x \in \text{VAL}(X)} p_X(x) = 1$$

CONTINUOUS RV:  $\text{VAL}(X)$  UNCOUNTABLE

$$P(a \leq X \leq b) := P(\{\omega | a \leq X(\omega) \leq b\})$$

PROBABILITY DENSITY FUNCTION

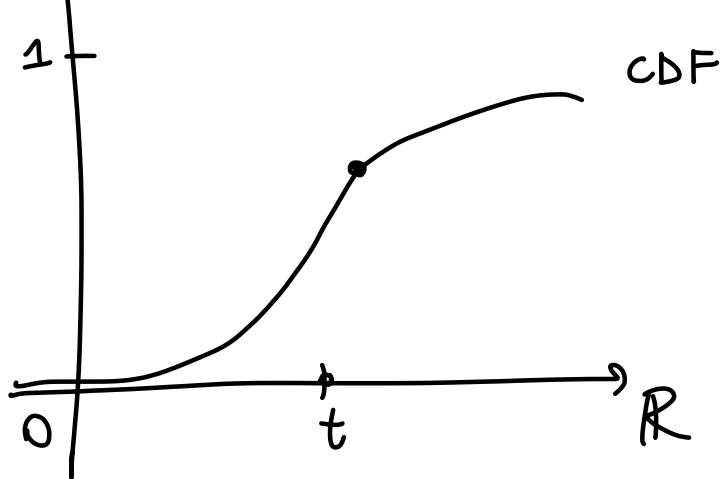
$$f_X: \mathbb{R} \rightarrow \mathbb{R}$$

$$f_X(x) := \frac{d}{dx} F_X(x)$$

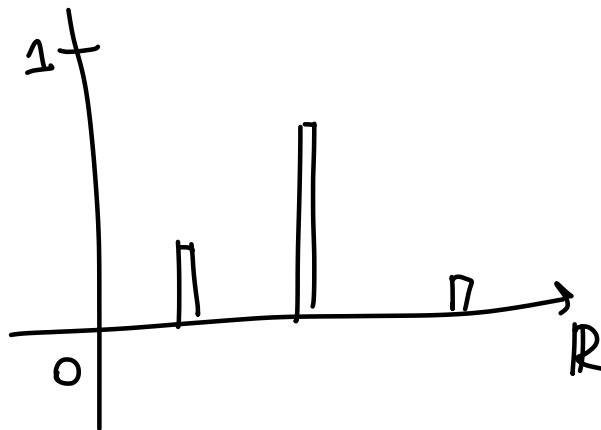
$$f_X(x) \neq P(X=x)$$

$$\int_{-\infty}^{+\infty} \underbrace{f_X(x)}_{\text{blue}} dx = 1 \quad \text{blue } P(x \in X \in x+dx)$$

PDF

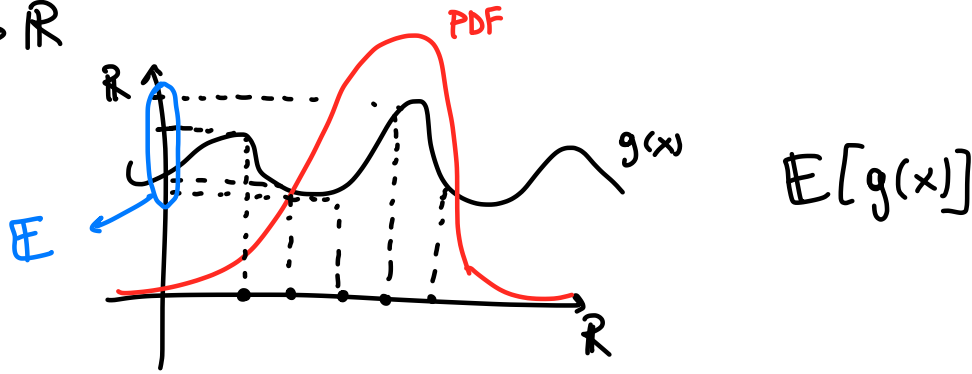


PMF



## EXPECTED VALUE AND VARIANCE

$$g: \mathbb{R} \rightarrow \mathbb{R}$$



LET  $X$  BE A DISCRETE RV WITH PMF  $p_x$

$$E[g(x)] := \sum_{x \in \text{Val}(x)} g(x) p_x(x)$$

IF CONTINUOUS (PDF)

$$E[g(x)] := \int_{-\infty}^{+\infty} g(x) f_x(x) dx$$

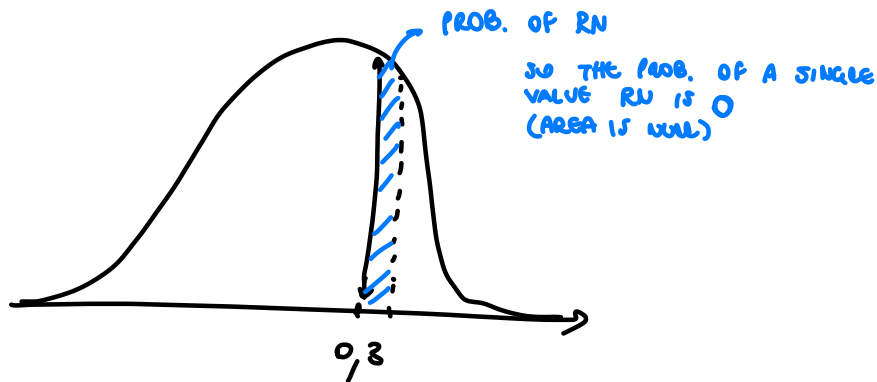
SAMPLES



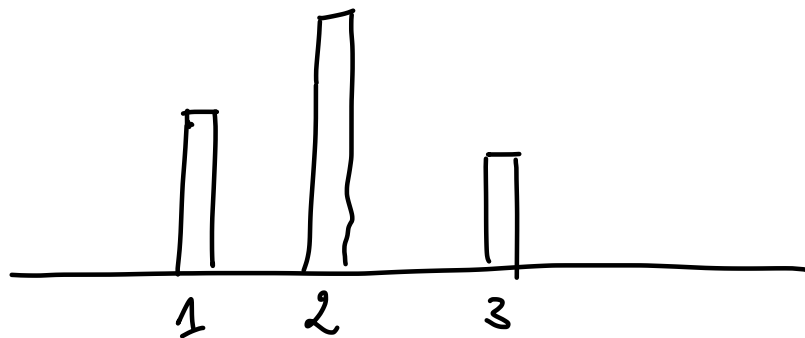
$$\lim_{N \rightarrow +\infty}$$

$$\underbrace{\frac{1}{N} \sum_{i=1}^N g(x^{(i)})}_{\text{MONTE CARLO ESTIMATE}} \rightarrow \int_{-\infty}^{\infty} \underbrace{g(x)}_{\substack{\text{SOME} \\ \text{FUNCTION} \\ \text{OF } x}} \underbrace{p(x)}_{\substack{\text{DENSITY}}} dx$$

LAW OF LARGE NUMBERS  
(TRUE EXPECTATIONS)

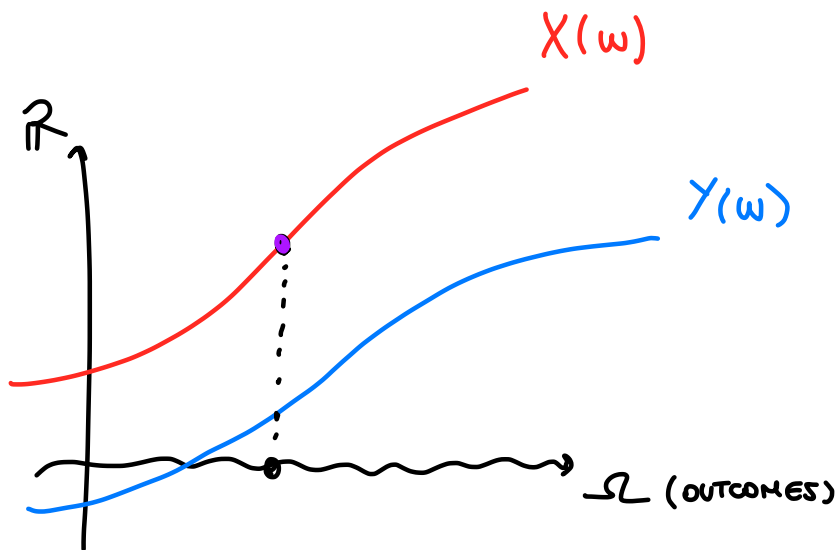


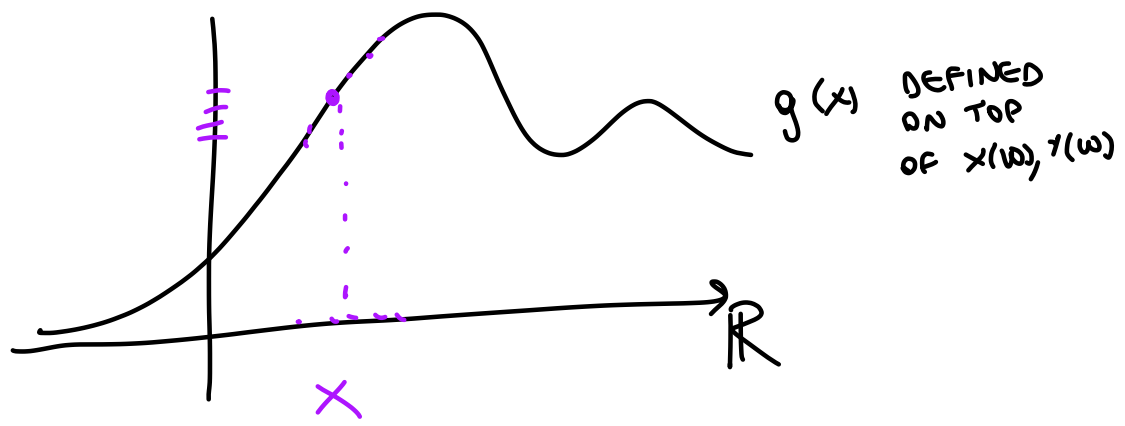
$$P(x = 0.3) \neq$$



OUTCOMES

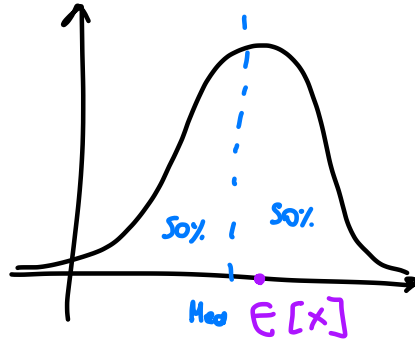
EVENTS  $\leftarrow$  PROB





### VARIANCE

$$\text{Var}(X) := \mathbb{E}[(X - \mathbb{E}[X])^2]$$



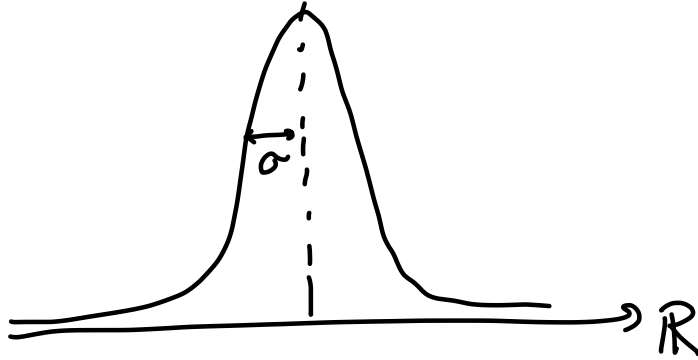
IF SYMMETRICAL (THE DISTRIBUTION)  
THEN  $\text{Med} \equiv \mathbb{E}[X]$

# DISTRIBUTION EX

DISTRIBUTION	PDF OR PMF	MEAN	VARIANCE
BERNOULLI ( $p$ )	$\begin{cases} p & \text{if } x=1 \\ 1-p & \text{if } x=0 \end{cases}$	$p$	$p(1-p)$
BINOMIAL ( $n, p$ )	$\binom{n}{k} p^k (1-p)^{n-k}$ for $k = 0, 1, \dots, n$	$np$	$np(1-p)$
GEOMETRIC ( $p$ )	$p(1-p)^{k-1}$ for $k=1, \dots$	$\frac{1}{p}$	$\frac{1-p}{p^2}$
POISSON ( $\lambda$ )	$\frac{e^{-\lambda} \lambda^k}{k!}$ for $k=0, 1, \dots$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
GAUSSIAN ( $\mu, \sigma^2$ )	$\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ for all $x \in (-\infty, \infty)$	$\mu$	$\sigma^2$
EXPONENTIAL ( $\lambda$ )	$\lambda e^{-\lambda x}$ for all $x \geq 0, \lambda > 0$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$

SPACE OVER THE DEFINED DISTRIBUTION :  $\mathbb{R}, \mathbb{R}^2, \dots$

PARAMETER, SHAPE OF THE DISTRIBUTION



WE CAN HAVE 2 RV

CDF  $F_{X,Y}(x,y) = P(X \leq x, Y \leq y)$

PMF  $p_{X,Y}(x,y) = P(X=x, Y=y)$

MARGINAL PMF  $p_X(x) = \sum_y p_{X,Y}(x,y)$

BIVARIATE PDF

$$f_{X,Y}(x,y) = \frac{\partial^2 F_{X,Y}(x,y)}{\partial x \partial y}$$

MARGINAL PDF

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x,y) dy$$

FOR INSTANCE, THE JOINT PDF, PMF

$$p(x, y)$$

DISC

CONT

$$p(x) = \sum_y p(x, y) = \int p(x, y) dy$$

$$p(y) = \sum_x p(x, y) = \int p(x, y) dx$$

### BAYES' THEOREM

GIVEN THE CONDITIONAL PROBABILITY OF AN EVENT  $P(x|y)$

WANT TO FIND THE "REVERSE" CONDITIONAL PROBABILITY,  $P(y|x)$

$$P(y|x) = \frac{P(x|y) P(y)}{P(x)} \quad \left( P(x|y) = \frac{P(x \cap y)}{P(y)} \right)$$

$$\text{WHERE } p(x) = \sum_{y' \in \text{value } y} P(x|y') P(y')$$

X AND Y ARE CONTINUOUS

$$f(y|x) = \frac{f(x|y)f(y)}{f(x)}$$

WHERE  $f(x) = \int_{y' \in \text{value}_y} f(x|y') f(y') dy'$

$$P(x, y) = P(x) P(y|x)$$

JOINT

MARGINED

CONDITIONAL

CHAIN RULE

$$= P(y) P(x|y)$$

BAYES

$$P(y|x) = \frac{P(y) P(x|y)}{P(x)} = \frac{P(y) P(x|y)}{\sum_{y'} P(y') P(x|y')}$$

JOINT  $\rightarrow \sum P(x, y')$

## INDEPENDENCE (OF 2 RN)

$X, Y$  ARE INDIA.

$$P_{XY}(x, y) = P_X(x) P_Y(y)$$

$$P_{Y|X}(x, y) = P_Y(y)$$

FOR CONT. RN

$$P_{XY}(x, y) \rightarrow f_{XY}(x, y)$$

EX

TOSS A COIN AND SPIN 1-7 NUMBERS. PROB OF GETTING ODD NUMBER AND A TAIL?

$$P_{XY}(x, y) = P_X(x) P_Y(y) = \frac{1}{2} \cdot \frac{4}{7} = \frac{2}{7}$$



Ex

$X, Y$  : 2 CONT. RV

$g, \mathbb{R}^2 \rightarrow \mathbb{R}$  : A FUNC. OF  $Y$  AND  $X$

$$\mathbb{E}(g(x, y)) = \int_{x \in \text{VAL}(x)} \int_{y \in \text{VAL}(y)} g(x, y) f_{XY}(x, y) dx dy$$

Ex

$$g(x, y) = 3x$$

$$f_{X,Y} = 4xy \quad 0 < x < 1, 0 < y < 1$$

$$\mathbb{E}(g(x, y)) = \int_0^1 \int_0^1 12x^2y \, dx dy$$

## COVARIANCE

$$\text{VAR}(x) = \mathbb{E}[x^2] - \mathbb{E}[x]^2$$

$$\text{Cov}(x, y) = \mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y]$$

$$\text{Cov}(x, x) = \text{VAR}[x]$$

# MULTIVARIATE GAUSSIAN

$x \in \mathbb{R}^n$  OBSERVATION

MEAN  $\mu$

COVARIANCE MAT  $\Sigma$

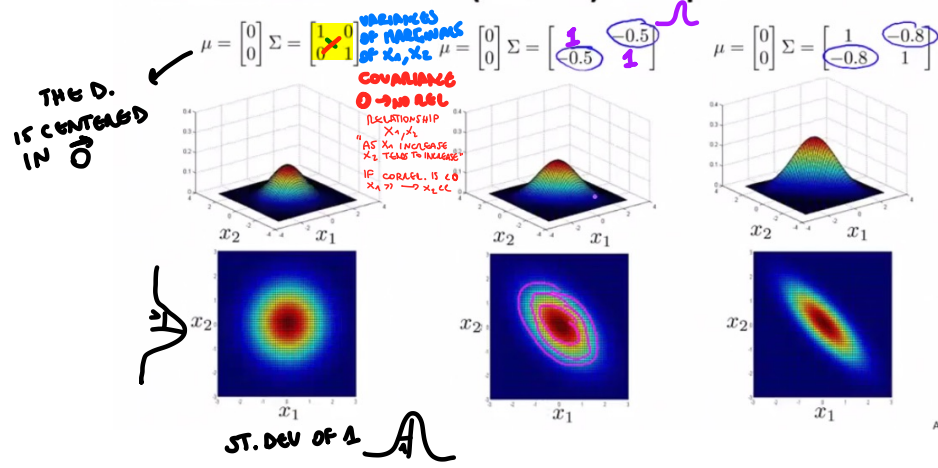
- INVERTIBLE

- POSITIVE SEMIDEFINITE

$$P(x; \underbrace{\mu, \Sigma}_{\text{PARAMS}}) = \underbrace{\frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}}}_{\text{NORMALIZING CONSTANT}} \exp\left[-\frac{1}{2} \underbrace{(x-\mu)^T \Sigma^{-1} (x-\mu)}_{\text{QUADRATIC FORM}}\right]$$

$$\iiint P(x; \mu, \Sigma) dx_1 dx_2$$

## Multivariate Gaussian (Normal) examples

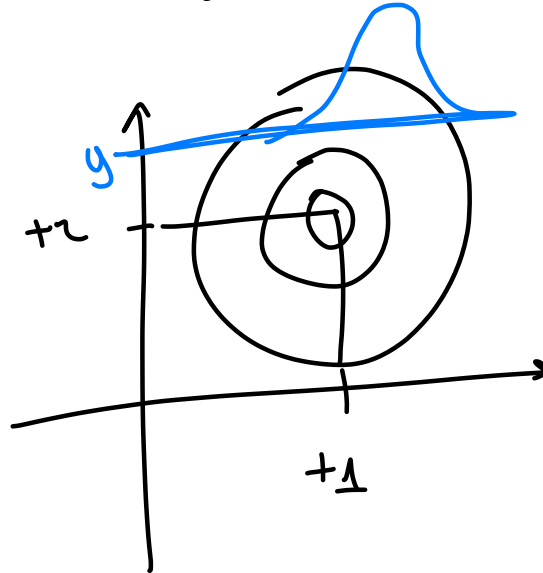


## CONDITIONAL EXPECTATION

$E[X]$  — CONSTANT

$E[X|Y]$  — RV

$E[X|Y=y]$  — FUNCTION OF  $y$



$E[X|Y=y]$

LAW OF TOTAL  $E$

$$E[X] = E[E[X|Y]] \quad \text{TRUE FOR ANY } X \text{ AND ANY } Y$$

$$P(a|b) = \frac{P(b|a)P(a)}{P(b)}$$

$$P(a|b, c) = \frac{P(b|a, c)P(a|c)}{P(b|c)}$$

# STATISTICS

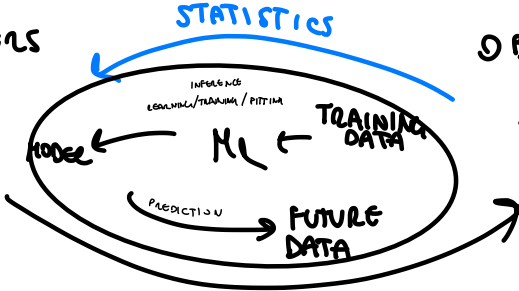
PARAMETERS

$$\mu, \Sigma$$

STATISTICS

OBSERVATION (DATA)

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



- METHOD OF  
MOMENTS

- MAXIMUM LIKELIHOOD  
ESTIMATION

PROBABILITY

TRAINING DATA

$$(x, y)$$

MAX LIKELIHOOD ESTIMATION (MLE)

EX GAUSSIAN DATA

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

$X = x^{(1)}, \dots, x^{(n)}$ ,  $x^{(i)} := i^{\text{th}} \text{ example} \in \mathbb{R}^d$ , SAMPLE INDEPENDENTLY  
AND IDENTICALLY DISTRIBUTED  
I.I.D.

VAL OF FUNCTION  
IF DATA

PARAMS

$$P(x; \mu, \Sigma) = \frac{1}{\sqrt{2\pi} |\Sigma|^{\frac{1}{2}}} \exp\left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right]$$

PROB. DENSITY

$$P(x^{(1)}, \dots, x^{(n)}) = \prod_{i=1}^n P(x^{(i)})$$

JOINT PROBABILITY

LET'S GENERALIZE

$$P(x^{(1)}, \dots, x^{(n)}; \mu, \Sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi} |\Sigma|^{\frac{1}{2}}} \exp\left[-\frac{1}{2} (x^{(i)} - \mu)^T \Sigma^{-1} (x^{(i)} - \mu)\right]$$

DEFINING LIKELIHOOD FUNCTION

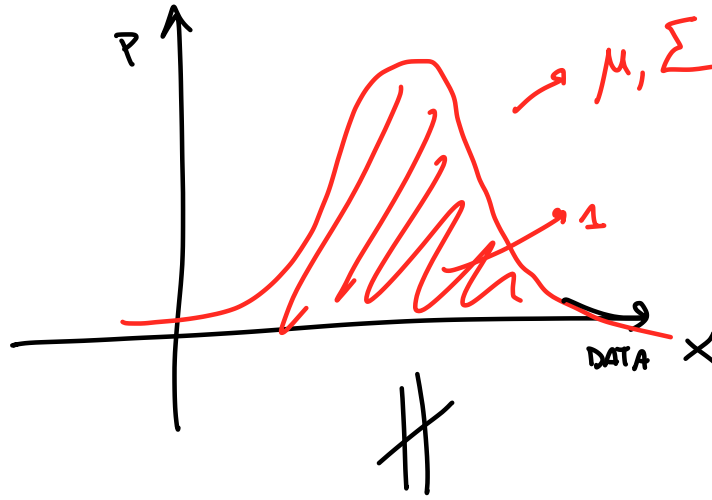
$$L(\underbrace{\mu, \Sigma}_{\text{NOW THE PARAMS ARE THE VARS}}; \underbrace{x^{(1)}, \dots, x^{(n)}}_{\text{PARAMETERS}}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi} |\Sigma|^{\frac{1}{2}}} \exp\left[-\frac{1}{2} (x^{(i)} - \mu)^T \Sigma^{-1} (x^{(i)} - \mu)\right]$$

SO IT IS A FUNCTION OF  $\mu, \Sigma$

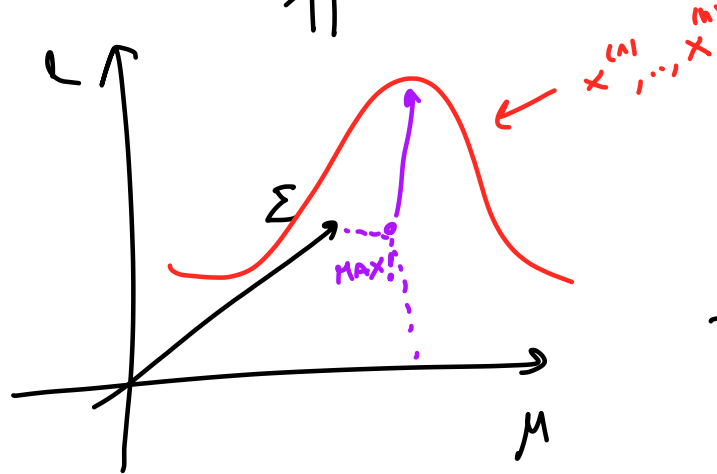
PROB. OF THE DATA GIVEN PARAMS  
LIKEL. OF THE PARAMS GIVEN DATA

$$\int p(x) dx = 1$$

$$\int L(\mu, \Sigma) d\mu d\Sigma = ? \quad \text{MAY NOT EVEN EXIST!}$$



$$P(x)$$



$$L(\mu, \Sigma)$$

THEN  $\hat{\mu}$  AND  $\hat{\Sigma}$  USED  
FOR PREDICTIONS



$$L(\theta; x) = \prod_{i=1}^n L(\theta; x^{(i)})$$

ESTIMATED (OUTPUT OF ESTIM. PROCESS)

$$\begin{aligned} \hat{\theta}_{MLE} &= \arg \max_{\theta} \prod_{i=1}^n L(\theta; x^{(i)}) \\ &= \arg \max_{\theta} \log \prod_{i=1}^n L(\theta; x^{(i)}) \end{aligned}$$

MAXIMIZE LOG. FUNC. (SAME  $\theta$ )

STD RECIPE  $\Rightarrow$

$$\arg \max_{\theta} \sum_{i=1}^n l(\theta; x^{(i)}) \quad l(\theta) = \log L(\theta)$$

NOW LET'S ASSUME  $M$ -GAUSSIAN  $l$   $\textcircled{=}$

$$\begin{aligned} \hat{\mu}, \hat{\Sigma} &= \arg \max_{\mu, \Sigma} \sum_{i=1}^n \log \left[ \frac{1}{(2\pi)^{\frac{M}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left[ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right] \right] \\ &= \arg \max_{\mu, \Sigma} \left[ \sum_{i=1}^n \left\{ K - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\} \right] \end{aligned}$$

LET'S MAX TO BOTH VARS

$$\mu$$

$$\nabla_{\mu} \sum_{i=1}^n K - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x^{(i)} - \mu)^T \Sigma^{-1} (x^{(i)} - \mu)$$

$$\nabla_x x^T A x = 2Ax$$

$\Sigma$  IS SYMM  $\Sigma = \Sigma^T$

$$\sum_{i=1}^n -\frac{1}{2} \left[ x^{(i)T} \cancel{\Sigma^{-1}} x^{(i)} - x^{(i)T} \Sigma^{-1} \mu - \mu^T \Sigma^{-1} x^{(i)} + \mu^T \Sigma^{-1} \mu \right]$$

$$= \sum_{i=1}^n -\frac{1}{2} \left[ -x^{(i)T} \Sigma^{-1} - \Sigma^{-1} x^{(i)T} + 2 \Sigma^{-1} \mu \right]$$

$$= \sum_i \left[ x^{(i)T} \Sigma^{-1} - \Sigma^{-1} \mu \right] = 0$$

SUMMATION

$$n \Sigma^{-1} \mu = \sum_{i=1}^n \Sigma^{-1} x^{(i)}$$

$$\cancel{\sum}^{-1} \mu = \cancel{\sum}^{-1} \left[ \frac{1}{n} \sum_{i=1}^n x^{(i)} \right]$$

$$\Rightarrow \boxed{\mu = \frac{1}{n} \sum_{i=1}^n x^{(i)}}$$

$$\left( \sum \right)$$

$$\frac{d}{d\Sigma} \Sigma^{-1} ?$$

$$S = \Sigma^{-1} \quad \nabla_S \mathcal{L} = 0 \quad (\Leftrightarrow) \quad \nabla_{\Sigma^{-1}} \mathcal{L} = 0 \quad \begin{array}{l} \text{SOLVE FOR } S, \\ \Sigma \end{array} \text{ INVERT AND FIND}$$

$$\nabla_S \sum_{i=1}^n \frac{1}{2} \log |S| - \frac{1}{2} (x^{(i)} - \mu)^T S (x^{(i)} - \mu)$$

$$= -\frac{1}{2} \left[ n S^{-1} - \sum_{i=1}^n [(x^{(i)} - \mu)(x^{(i)} - \mu)^T] \right] = 0 \quad \dots$$

$$\nabla_A x^T A x = ? \quad [x x^T]$$

$$\nabla_A A = \left[ \frac{d}{da_{11}} A \quad \dots \quad \frac{d}{da_{nm}} A \right] \quad x^T A = (x_1 a_{11} + x_2 a_{21} + \dots + x_n a_{n1}, \dots, x_1 a_{1m} + \dots + x_n a_{nm})^T$$

$$X^T A X = (x_1^T x_1 a_{11} + \dots + x_n^T x_n a_{nn}, \dots, x_1^T x_n a_{1n} + \dots + x_n^T x_n a_{nn})^T$$

$$\nabla_A X^T A X = \nabla_A (x_1^T x_1 a_{11} + \dots + x_n^T x_n a_{nn}, \dots, x_1^T x_n a_{1n} + \dots + x_n^T x_n a_{nn})^T$$

$$\nabla_A X^T A X = X^T X$$

$$S^{-1} = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu)(x^{(i)} - \mu)^T = \sum$$

$$\approx \mathbb{E}[(x - \mathbb{E}[x])^2]$$

DONE! 😊