

Johns Hopkins University
EN. 601.475 Machine Learning

Probability & Linear Algebra Review

Xuan Zhang
Sep. 2, 2022

I. Probability

- Random Variables
- Bayes' Rule
- Bernoulli distributions & Sigmoid function
- Categorical distributions & Softmax function
- Gaussian distribution

II. Bayes Optimal Rule

III. Linear Algebra

- Vector norms
- matrix multiplication
- vector derivatives

I. Probability

Probability

two interpretations:

- frequentist

probabilities represent long run frequencies of events

- Bayesian (this course)

probability is used to quantify our uncertainty about sth.

♥ can be used to model uncertainty about one-off events

Probability

$\overset{\text{event}}{\uparrow}$
 $\Pr(A)$: probability that A is true
 $\in [0, 1]$

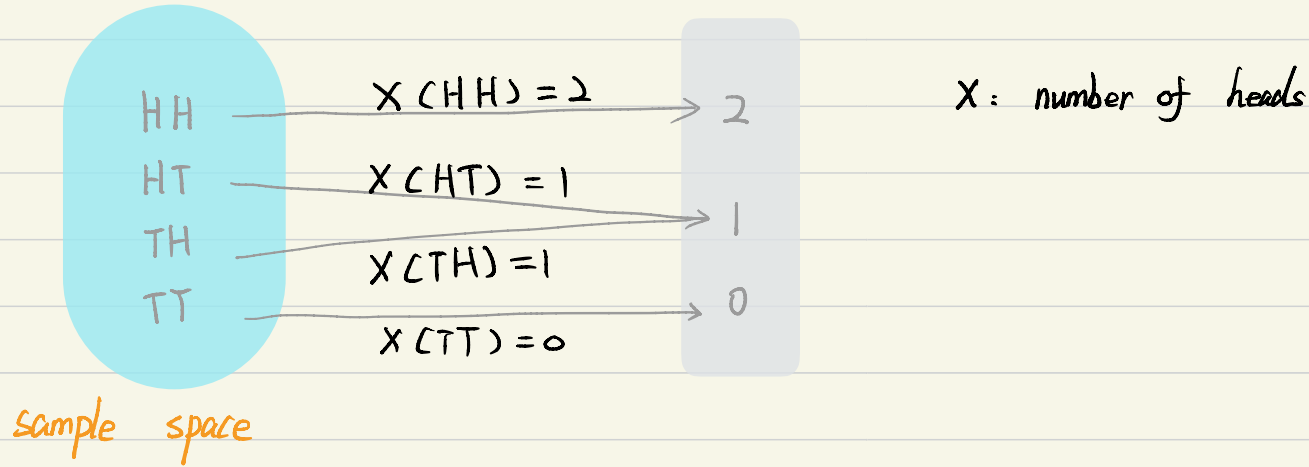
joint probability: $\Pr(A, B)$

if A & B are independent: $\Pr(A, B) = \Pr(A) \Pr(B)$

conditional probability: $\Pr(B|A) \triangleq \frac{\Pr(A, B)}{\Pr(A)}$

conditionally independent $A \perp B | C$
 $\Pr(A, B|C) = \Pr(A|C) \Pr(B|C)$

Random Variable (RV)



RVs are functions. RV is a numeric function of the outcome.

Random Variable (RV)

discrete RV:

sample space is finite or countably infinite

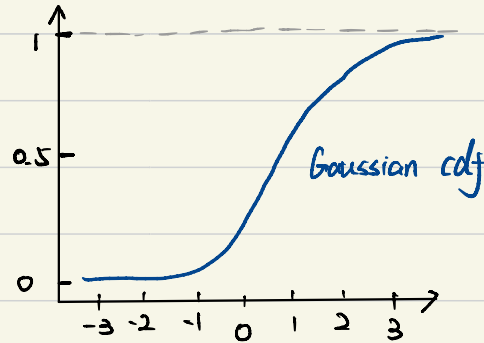
continuous RV:

infinite number of values between two values

Cumulative distribution function (cdf)

$$P(x) \triangleq \Pr(X \leq x)$$

e.g. $\Pr(a < X \leq b) = P(b) - P(a)$



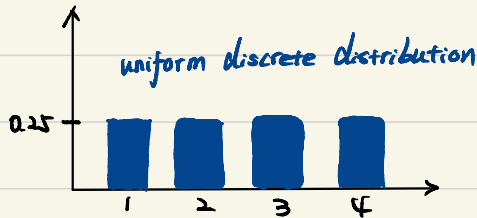
cdf's are monotonically non-decreasing functions.

Random Variable (RV)

probability mass function (pmf)
(discrete RV)

$$p(x) \triangleq \Pr(X=x)$$

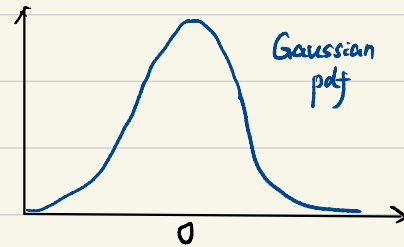
$$\sum_{x \in X} p(x) = 1$$



probability density function (pdf)
(continuous RV)

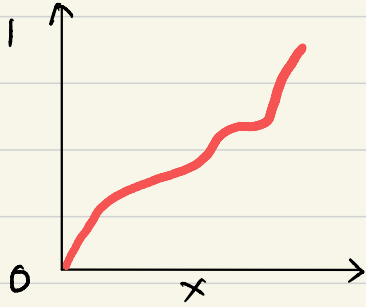
$$p(x) \triangleq \frac{d}{dx} P(x)$$

$$\int_{-b}^{\infty} p(x) dx = 1$$

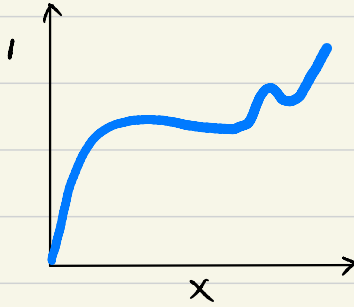


Which of the following are valid cdf?

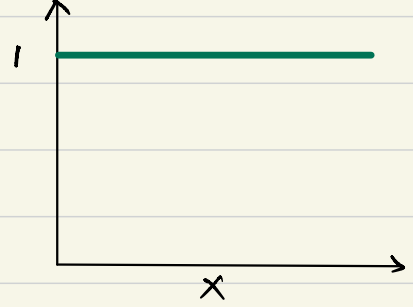
a.



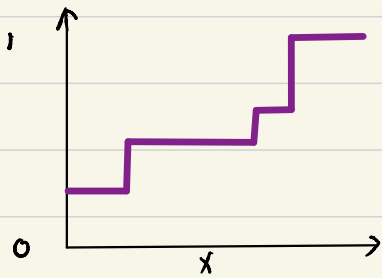
b.



c.



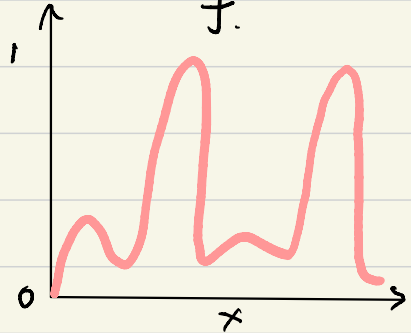
d.



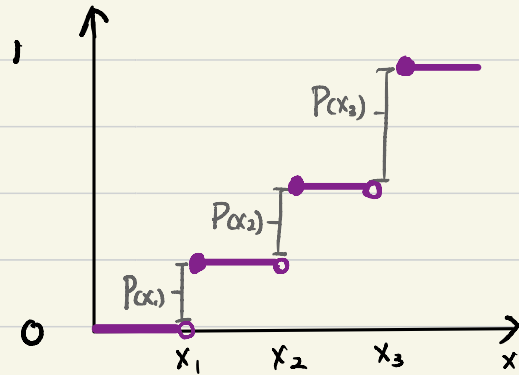
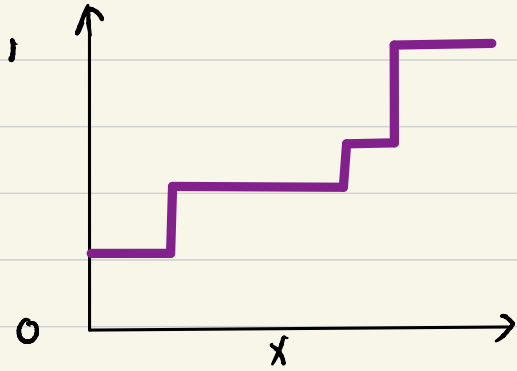
e.



f.



Is it a cdf for a discrete random variable?



Terminology

X

a random variable

$x \sim X$

a sample value of a RV X

$$\Pr(x) = \Pr(X=x)$$

probability of event X has value x

$$\Pr(\bar{x}) = 1 - \Pr(x)$$

probability of x not happening

$P(x)$

cumulative distribution function (cdf)

$p(x)$

probability mass function (pmf)

probability density function (pdf)

Related Random Variables

joint distribution:

$$p(x, y) = p(X=x, Y=y)$$

marginal distribution:

$$p(X=x) = \sum_y p(X=x, Y=y) \rightarrow \text{sum rule}$$

conditional distribution:

$$p(Y=y | X=x) = \frac{p(X=x, Y=y)}{p(X=x)}$$

$$\Downarrow$$
$$p(x, y) = p(x) p(y|x) \rightarrow \text{product rule}$$

chain rule of probability.

$$p(x_1, \dots, x_p) = p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2) \dots p(x_p | x_1, \dots, x_{p-1})$$

Moments of a distribution

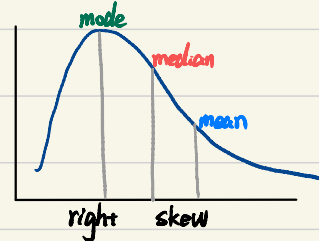
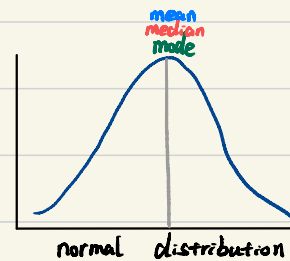
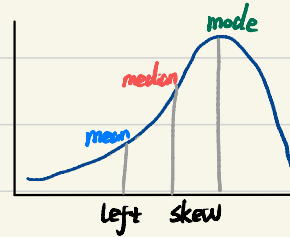
- mean (expected value):

$$E[X] \triangleq \sum_{x \in \mathcal{X}} x p(x) \quad E[X] \triangleq \int_{\mathcal{X}} x p(x) dx$$

(discrete RV) (continuous RV)

linearity of expectation:

$$E\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n E[X_i]$$



- Variance

$$V[X] \triangleq E[(X - \mu)^2] = \int (x - \mu)^2 p(x) dx$$
$$= E[X^2] - \mu^2$$

$$E[X^2] = \sigma^2 + \mu^2$$

$$V\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n V[X_i]$$

$$V[aX + b] = a^2 V[X]$$

- mode

$$x^* = \operatorname{argmax}_x p(x)$$

Bayes' rule

H : unknown (or hidden) quantity

$Y=y$: observed data

prior: what we know about H
before seeing data

likelihood: distribution of data
we expect to see if $H=h$

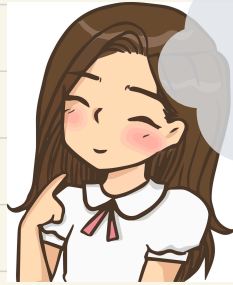
$$p(H=h | Y=y) = \frac{p(H=h) p(Y=y | H=h)}{p(Y=y)}$$

posterior: new belief state about H

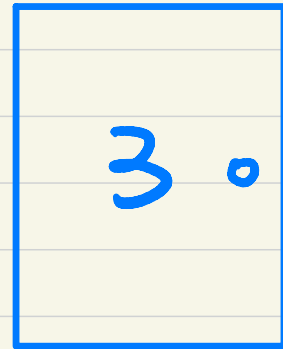
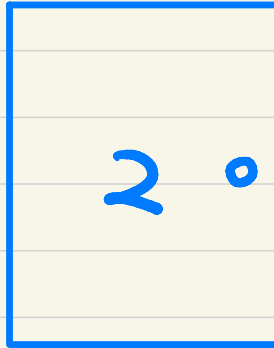
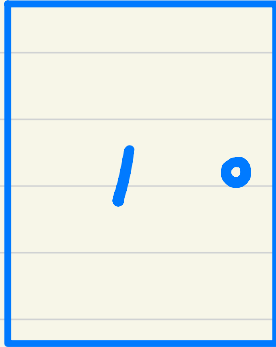
marginal likelihood

$$\text{posterior} \propto \text{prior} \times \text{likelihood}$$

Example: The Monty Hall Problem



\$1,000,000

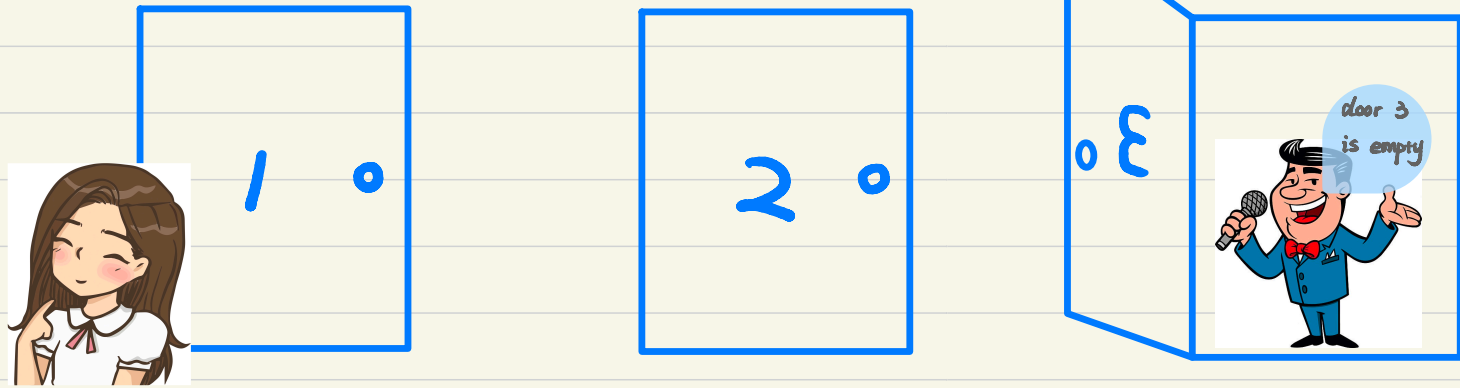


Example: The Monty Hall Problem

What's your choice?

A. stick with door 1

B. switch to door 2



Example: The Monty Hall Problem

H : (hidden quantity) the prize is behind a door $\{1, 2, 3\}$

Y : (observation) a door is opened $\{1, 2, 3\}$

We want to compare $P(H=1|Y=3)$ vs. $P(H=2|Y=3)$

Example: The Monty Hall Problem

H: (hidden quantity) the prize is behind a door $\{1, 2, 3\}$

Y: (observation) a door is opened $\{1, 2, 3\}$

We want to compare $P(H=1|Y=3)$ vs. $P(H=2|Y=3)$

$$P(H=1|Y=3) = \frac{\overset{\substack{\frac{1}{2} \text{ why?} \\ \uparrow}}{P(Y=3|H=1)} \overset{\substack{\frac{1}{3} \\ \uparrow}}{P(H=1)}}{P(Y=3)} = \frac{1}{3}$$

(Note: A green arrow points from the $P(Y=3)$ denominator to the $P(Y=3)$ term in the equation below.)

$$\begin{aligned} P(Y=3) &= P(Y=3|H=1)P(H=1) + P(Y=3|H=2)P(H=2) + P(Y=3|H=3)P(H=3) \\ &= \frac{1}{2} \times \frac{1}{3} + 1 \times \frac{1}{3} + 0 \times \frac{1}{3} = \frac{1}{2} \end{aligned}$$

$$P(H=2|Y=3) = 1 - P(H=1|Y=3) = \frac{2}{3}$$

choose door 2!

Bernoulli and Binomial Distribution

Bernoulli distribution

toss a coin

θ : probability that it lands head

$$p(Y=1) = \theta \quad p(Y=0) = 1-\theta$$

↑ head ↓ tail

$$Y \sim \text{Ber}(\theta)$$

pmf: $\text{Ber}(y|\theta) = \begin{cases} 1-\theta & \text{if } y=0 \\ \theta & \text{if } y=1 \end{cases}$

↕

$$\text{Ber}(y|\theta) \triangleq \theta^y (1-\theta)^{1-y}$$

Binomial distribution

toss a coin N times

number of heads: $S \triangleq \sum_{n=1}^N \mathbb{I}(y_n=1)$

$$\text{Bin}(s|N, \theta) \triangleq \binom{N}{s} \theta^s (1-\theta)^{N-s}$$

Bernoulli Distribution and Sigmoid function

Bernoulli distribution

$$\text{Ber}(y|\theta) \triangleq \theta^y (1-\theta)^{1-y}$$

Sigmoid (logistic) function

We want to predict a binary variable $y \in \{0, 1\}$ given inputs $x \in \mathcal{X}$

$$p(y|x, \theta) = \text{Ber}(y|f(x; \theta)) \quad \underline{0 \leq f(x; \theta) \leq 1}$$

can be relaxed

$$p(y|x, \theta) = \text{Ber}(y|\sigma(f(x; \theta))) \quad \sigma(): \text{sigmoid function} \quad 0 \leq \sigma(a) \leq 1$$

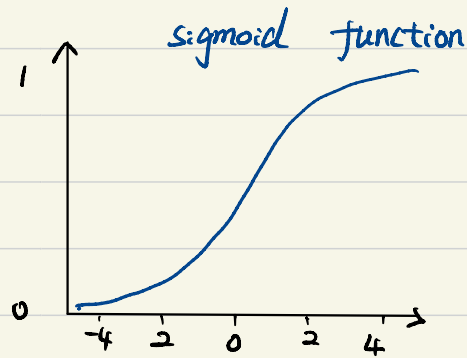
$$\sigma(a) \triangleq \frac{1}{1+e^{-a}}, \text{ where } a = f(x; \theta)$$

Sigmoid function

$$\sigma(a) \triangleq \frac{1}{1+e^{-a}}, \text{ where } a = f(x; \theta)$$

$$p(y=1 | x, \theta) = \frac{1}{1+e^{-a}} = \sigma(a)$$

$$p(y=0 | x, \theta) = 1 - \frac{1}{1+e^{-a}} = \sigma(-a)$$



log odds: $\log\left(\frac{p}{1-p}\right) = \log\left(\frac{e^a}{1+e^a} \frac{1+e^a}{1}\right) = \log(e^a) = \boxed{a}$, where $p = p(y=1 | x, \theta)$

↓
logit

(Spoiler) Binary logistic regression

$$p(y | x; \theta) = \text{Ber}(y | \sigma(w^T x + b)) \quad \Leftrightarrow \quad p(y=1 | x; \theta) = \sigma(w^T x + b) = \frac{1}{1+e^{-(w^T x + b)}}$$

decision boundary: $x^* \rightarrow p(y=1 | x=x^*, \theta) = 0.5$

Categorical and multinomial distributions

Categorical Distribution

roll a C -sided dice

$$C > 2$$

$$P(y=c | \theta) \triangleq \prod_{c=1}^C \theta_c^{\mathbb{I}(y=c)}$$

$$P(y=c | \theta) = \theta_c$$

Multinomial Distribution

roll a C -sided dice N times

$$P(s | N, \theta) \triangleq \binom{N}{s_1, \dots, s_C} \prod_{c=1}^C \theta_c^{s_c},$$

$$\text{where } s_c \triangleq \sum_{n=1}^N \mathbb{I}(y_n = c)$$

$\binom{N}{s_1, \dots, s_C}$: multinomial coefficient
number of ways to divide a set of size $N = \sum_{c=1}^C s_c$ into subsets with sizes s_1 to s_C

Softmax function

In conditional case:

$$p(y|x; \theta) = \text{Cat}(y | f(x; \theta)) \quad 0 \leq f(x; \theta) \leq 1$$

↓
can be relaxed

$$p(y|x; \theta) = \text{Cat}(y | S(a)_c)$$

$S(\cdot)$: softmax function

$$0 \leq S(a)_c \leq 1 \\ \sum_{c=1}^C S(a)_c = 1$$

$$S(a) \triangleq \left[\frac{e^{a_1}}{\sum_{c'=1}^C e^{a_{c'}}}, \dots, \frac{e^{a_C}}{\sum_{c'=1}^C e^{a_{c'}}} \right]$$

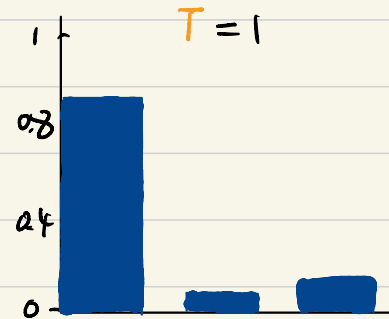
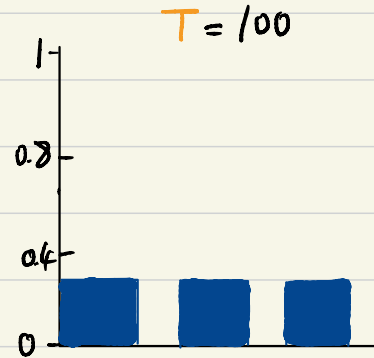
$$a = f(x; \theta) \quad \text{logits}$$

Softmax function

$$S(a) \triangleq \left[\frac{e^{a_1}}{\sum_{c=1}^C e^{a_{c'}}}, \dots, \frac{e^{a_C}}{\sum_{c=1}^C e^{a_{c'}}} \right]$$

temperature T : change the output distribution

$$\text{as } T \rightarrow 0, S(a/T)_c = \begin{cases} 1.0 & \text{if } c = \operatorname{argmax}_{c'} a_{c'} \\ 0.0 & \text{otherwise} \end{cases}$$



Gaussian (normal) distribution

$$Y \sim N(\mu, \sigma^2)$$

pdf:

$$N(y | \mu, \sigma^2) \triangleq \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y-\mu)^2}$$

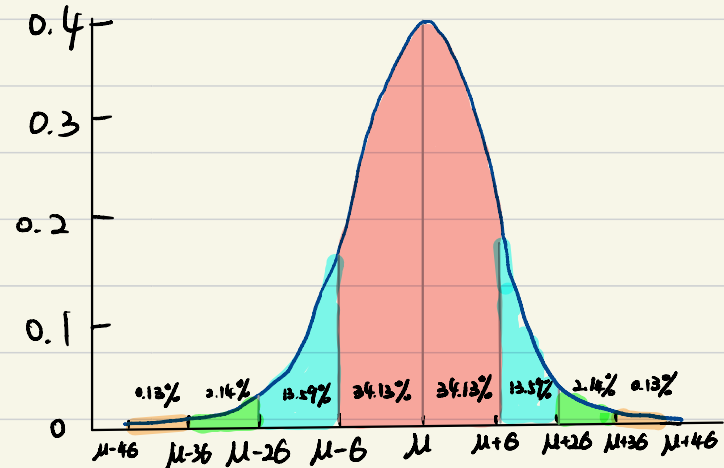
mean:

$$\mu = E[N(\cdot | \mu, \sigma^2)]$$

standard deviation:

$$\sigma = \text{std}[N(\cdot | \mu, \sigma^2)]$$

Normal Distribution



II. Bayes Optimal Rule

Bayes Optimal Rule

We are searching for a "f" that minimizes the expected loss.

measures distance between true label and prediction.

$$f^* = \underset{f}{\operatorname{argmin}} \mathbb{E}_{X,Y} [\operatorname{Loss}(Y, f(X))]]$$

$f^*: X \rightarrow Y$

The expectation on $P_{X,Y}$,
i.e. the data distribution.

prediction

the function that achieves
the minimal possible population risk

true label

population risk

can't compute because $P_{X,Y}$ is unknown.

III. Linear Algebra

Notation

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

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

one-hot vector (unit vector):

$$e_i = (0, \dots, 0, 1, 0, \dots, 0)$$

matrix $A \in \mathbb{R}^{m \times n}$

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}$$

A_{ij} : the entry of A in the i -th row and j -th column

$A_{i,:}$: the i -th row

$A_{:,j}$: the j -th column

Tensor

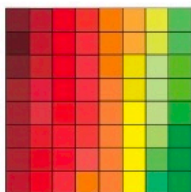
tensor = multidimensional array

vector



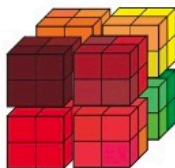
$$\mathbf{v} \in \mathbb{R}^{64}$$

matrix



$$\mathbf{X} \in \mathbb{R}^{8 \times 8}$$

tensor



$$\mathbf{X} \in \mathbb{R}^{4 \times 4 \times 4}$$

tensor:

a generalization of a 2d array to more than 2 dimensions.

order or rank of the tensor:
the number of dimensions

Vector Norms

norm of a vector $\|x\|$:

a measure of the "length" of the vector

properties:

1. non-negativity $\|x\| \geq 0$

2. definiteness $\|x\| = 0$ if $x = 0$

3. absolute value homogeneity $\|tx\| = |t| \|x\|$ for all $x \in \mathbb{R}^n$, $t \in \mathbb{R}$

4. triangle inequality $\|x+y\| \leq \|x\| + \|y\|$ for all $x, y \in \mathbb{R}^n$

Vector Norms

p-norm

$$\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}, \text{ for } p \geq 1$$

2-norm (Euclidean norm)

$$\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

note that $\|x\|_2^2 = x^T x$

1-norm

$$\|x\|_1 = \sum_{i=1}^n |x_i|$$

max-norm

$$\|x\|_\infty = \max_i |x_i|$$

Matrix Multiplication

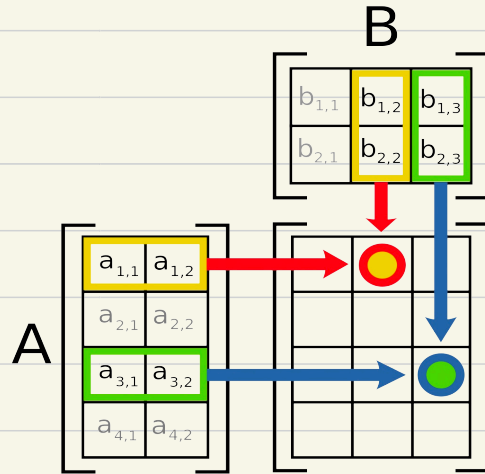
can only occur when $A \in \mathbb{R}^{N \times M}$ $B \in \mathbb{R}^{M \times D}$

$$A \times B = C \in \mathbb{R}^{N \times D}$$

$$C_{i,j} = \sum_{m=1}^M a_{i,m} \times b_{m,j}$$



$$C_{i,j} = A[i, :] \cdot B[:, j]$$



Matrix Derivatives and Gradients

consider a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$

gradient of f : $\nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$

consider a function $f: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$

$$\frac{\partial f}{\partial X} = \begin{pmatrix} \frac{\partial f}{\partial x_{11}} & \dots & \frac{\partial f}{\partial x_{1n}} \\ \vdots & & \vdots \\ \frac{\partial f}{\partial x_{m1}} & \dots & \frac{\partial f}{\partial x_{mn}} \end{pmatrix}$$

consider a function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$

Jacobian $m \times n$

$$J_f(x) = \frac{\partial f}{\partial x^T} \triangleq \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{pmatrix} = \begin{pmatrix} \nabla f_1(x)^T \\ \vdots \\ \nabla f_m(x)^T \end{pmatrix}$$

Consider a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$
that is twice differentiable

Hessian $n \times n$

$$H_f = \frac{\partial^2 f}{\partial x^2} = \nabla^2 f = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

Matrix Derivatives Idioms

$$\frac{\partial (a^T x)}{\partial x} = a$$

$$\frac{\partial}{\partial X} (a^T X b) = a b^T$$

$$\frac{\partial (b^T A x)}{\partial x} = A^T b$$

$$\frac{\partial}{\partial X} (a^T X^T b) = b a^T$$

$$\frac{\partial (x^T A x)}{\partial x} = (A + A^T) x$$

There are more!

Refer to « The Matrix Cookbook », Petersen et al., Chapter 2.



Questions?