

①

Marginal Probability:

Probability of specific event occurring independent of other events

$$P(A) = \sum_b P(A, B=b) \text{ or } \sum_b P(A \cap B=b)$$

Conditional probability

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \text{ or } \frac{P(A, B)}{P(B)}$$

Given B already happened what is probability of A

Joint probability :— two or more events occurring at same time
 $P(A \cap B) \text{ or } P(A, B)$

Sum Rule

for disjoint events

$$P(X=x_i) = \sum_{j=1}^L P(X=x_i, Y=y_j)$$

Product Rule

$$P(X=x_i, Y=y_j) = P(Y=y_j | X=x_i) \cdot P(X=x_i)$$

or

$$P(X=x_i, Y=y_j) = P(X=x_i | Y=y_j) \cdot P(Y=y_j)$$

Bayes theorem

$$P(X=x_j | Y=y_i) = \frac{P(Y=y_i | X=x_j) \cdot P(X=x_j)}{P(Y=y_i)}$$

By sum Rule $P(Y=y_i) = \sum_{j=1}^L P(Y=y_i | X=x_j) \cdot P(X=x_j)$

(3)

Probability distribution function

function that describes the likelihood of a continuous random variable taking a particular value

PDF → Continuous random variable

PMF → discrete random variable.

key characteristics

① non negative.

$$f(x) \geq 0 \text{ for all } x \text{ (random variable)}$$

② Total area under the curve.

total area under the PDF curve across the entire range of possible values equals 1

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

③ probability of an interval :- probability of that the random variable falls within a specific interval $[a, b]$ can be

found by $P(a \leq X \leq b) = \int_a^b f(x) dx$

Cumulative distribution function(CDF)
~ is a function that describes the probability that a random variable takes on a value less than or equal to a specific value.

For random variable X

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

small x represent specific value

→ As x approach negative infinity CDF approach 0

$$\lim_{x \rightarrow -\infty} F(x) = 0$$

→ As x approach positive infinity CDF approach 1

$$\lim_{x \rightarrow \infty} F(x) = 1$$

Application

CDFs allow for the calculation of probability of a random variable falling within a range

$$P(a \leq X \leq b) = F(b) - F(a)$$

Expectation :-

(3)

~ provides a measure of "center" or "average" of a random variable's distribution. It gives an idea of what value you can expect to obtain from random variable after many trials

⇒ Discrete Random variable :-

If X is a discrete random variable with possible values n_1, n_2, \dots, n_n and corresponding probabilities

$P(X=n_i) = p_i$ then Expectation $E(x)$

$$E(x) = \sum_{i=1}^n n_i * p_i$$

⇒ Continuous Random variable :-

If X is continuous random variable with a PDF $f(n)$ then

$$E(x) = \int_{-\infty}^{\infty} n \cdot f(n) dn.$$

(6)

Properties of Expectation.

① Linearity

$$E(ax+by) = aE(x) + bE(y)$$

a and b are constants, x is any random variable

② non negative :-

if X is non negative random variable

then $E(X) \geq 0$

③ Expectation of Constant : if c is a constant

then $E(c) = c$

Variance

⑦

it measures how much the values of random variable X differ from the expected value (mean)

$$\begin{aligned}\text{Var}(X) &= E[(X - E(X))^2] \\ &= E(X^2) - (E(X))^2\end{aligned}$$

Properties of Variance

① Non negativity

$$\text{Var}(X) \geq 0 \quad \text{for all } X$$

② Units :- Variance is in squared unit of random variable

③ Linearity:- For any random variable X and Y and constant a and b

$$\text{Var}(ax + by) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y)$$

if X and Y are independent then $\text{Cov}(X, Y) = 0$

④ Variance of constant is zero

Covariance

measure of relationship b/w two random variables
it indicates extent to which variable change together.

positive Covariance :- both increase together.

negative Covariance :- one increase other decrease
zero Covariance :- no linear relationship

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

if X takes value $x_1, x_2, x_3, \dots, x_n$

with probability $p_1, p_2, p_3, \dots, p_n$

Y takes value $y_1, y_2, y_3, \dots, y_n$ with same probability

optionally

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

$$\text{Cov}(x,y) = \sum_{j=1}^n (x_i - E(x))(y_i - E(y)) \cdot p_i$$

or

$$\boxed{\text{Cov}(x,y) = E(xy) - E(x)E(y)}$$

For sample size n

using \star

$$\boxed{\text{Cov}(x,y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}$$

(9)

Gaussian Distribution or

Normal distribution

bell shaped curve.

for single real value x

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

μ is mean

σ is standard deviation

Variance σ^2

Multivariate Gaussian Distribution

For single real value variable x

$$N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

μ = mean

D = dimension

Σ Covariance matrix $D \times D$

Likelihood function for Gaussian distribution

~ expresses the probability of observing the data given specific parameters (mean and variance)

For a set of independent observations

$\mathbf{x} = [x_1, x_2, \dots, x_n]$ from a gaussian distribution with unknown mean μ and variance σ^2

probability of dataset given by

$$P(\mathbf{x} | \mu, \sigma^2) = \prod_{n=1}^N N(x_n | \mu, \sigma^2)$$

product of gaussian distribution
Stand for product

Log likelihood function is

$$\ln(P(\mathbf{x} | \mu, \sigma^2)) = -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 - \frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2)$$

To estimate parameters μ and σ^2
we maximize the log likelihood function

(18)

$$\mu_{ML} = \frac{1}{n} \sum_{i=1}^n n_i$$

Differentiate w.r.t μ

$$\sigma^2_{ML} = \frac{1}{n} \sum_{i=1}^n (n - \mu_{ML})^2$$

which is mean

variable

Independence :-

Two events A and B are independent if the occurrence of one does not affect the occurrence of other.

Ex Flipping a coin, getting head
Rolling a die and getting 4

$$P(A \cap B) = P(A) \times P(B)$$

Mutual Exclusivity

Two events A and B are mutually exclusive if they cannot occur at the same time.

Ex Rolling a die and getting 2
Rrolling a die and getting 5

$$P(A \cap B) = 0$$

key differences

independence:- Events can occur together

mutual exclusivity:- Events cannot occur together.

Entropy :-

Entropy is a measure of uncertainty or randomness associated with a random variable.

Given a discrete random variable X with possible outcomes x_1, x_2, \dots, x_n and corresponding probabilities $P(x_1), P(x_2), \dots, P(x_n)$ then entropy $H(X)$ is defined as

$$H(X) = - \sum_{i=1}^n P(x_i) \log_2 P(x_i)$$

Higher Entropy :- indicates more uncertainty.
if all outcomes are equally likely, entropy is maximum.

lower Entropy :- indicates less uncertainty.
if outcome is certain (probability of 1) then entropy is zero.

Properties

- ① non negative, always greater than or equal to zero
- ② additive

$$H(X, Y) = H(X) + H(Y)$$

Joint Entropy

(15)

Quantifies measure of uncertainty associated with two random variables considered together, it quantifies the total amount of information required to describe the outcome of both variables simultaneously.

For two discrete R.V X and Y

$$H(X,Y) = - \sum_{x \in X} \sum_{y \in Y} P(x,y) \log_2 P(x,y)$$

$P(x,y)$ is joint probability distribution of X and Y

Higher joint entropy:- indicates greater uncertainty or complexity in relationship b/w the variables

Lower joint entropy:- suggests a more predictable or constrained relationship.

Properties:-

① Relation to individual Entropies:- The joint entropy is related to the individual entropies

$$H(X,Y) = H(X) + H(Y|X)$$

where $H(Y|X)$ is conditional entropy of Y given X

② if x and y are independent

(1)

$$H(X, Y) = H(X) + H(Y)$$

Conditional Entropy

(1)

Conditional entropy is a measure of the amount of uncertainty remaining about a random variable given that the value of another random variable is known. It quantifies how much additional information is needed to describe the outcome of one variable when the outcome of another variable is already known of y given x .

For two discrete R.V X and Y , conditional Entropy

$$H(Y|X) = - \sum_{n \in X} p(n) \sum_{y \in Y} p(y|n) \log_2 p(y|n)$$

where

$p(y|n)$ is conditional probability of y given n
 $p(n)$ is probability of X

High conditional entropy → indicate that knowing x does not provide much information about y (y remains uncertain)

low conditional Entropy → suggest that knowing x significantly reduces uncertainty about y .

Properties

Relationship to Joint Entropy

$$H(Y|X) = H(X, Y) - H(X)$$

Conditional entropy of X and Y can be derived from the total uncertainty in X and Y minus uncertainty in X alone.

If Y is independent of X :-

$$H(Y|X) = H(Y)$$

In this case knowing X provides no additional information about Y .

Chain Rule of Entropy

For a sequence of random variable $X_1, X_2, X_3 \dots X_n$ the joint entropy $H(X_1, X_2, \dots, X_n)$ can be expressed using the chain rule.

$$H(X_1, X_2, \dots, X_n) = H(X_1) + H(X_2 | X_1) + H(X_3 | X_1, X_2) + \dots + H(X_n | X_1, X_2, \dots, X_{n-1})$$

Cross Entropy

~ measures the difference b/w two random probability distributions. it quantifies how well one probability distribution approximates another, often used in context such as classification task.

For two discrete probability distribution P and Q the cross entropy $H(P, Q)$ is defined as

$$H(P, Q) = - \sum_x P(x) \log_2 Q(x)$$

where, P is true distribution (often ground truth)

Q is the estimated distribution (often model prediction)

The sum is taken over all possible outcomes x .

Higher cross entropy :- indicates a larger difference b/w two distributions meaning Q does not accurately represent P

Lower cross entropy :- indicate the Q as a good approximation of P

Mutual Information

mutual information btw two R.V measure the amount of information that one conveys about the other. it measure the average reduction in uncertainty about X that result from learning about Y

for two discrete R.V X and Y , mutual information is defined as

$$I(X;Y) = H(X) - H(X|Y)$$

or it can be expressed as joint entropy

$$I(X;Y) = H(X) + H(Y) - H(X,Y)$$

zero mutual information :- if $I(X;Y) = 0$ the variables are independent, knowing one does not provide any information about the other

positive mutual information :- indicates a dependency btw the variables, knowing one variable reduces uncertainty about the other.

Properties

$$\text{Chain rule } I(X;Y,Z) = I(X;Y) + I(X;Z|Y)$$

① Symmetry $I(X;Y) = I(Y;X)$

~~$I(X,Y,Z) = I(X) + I(Y|X) + I(Z|X,Y)$~~

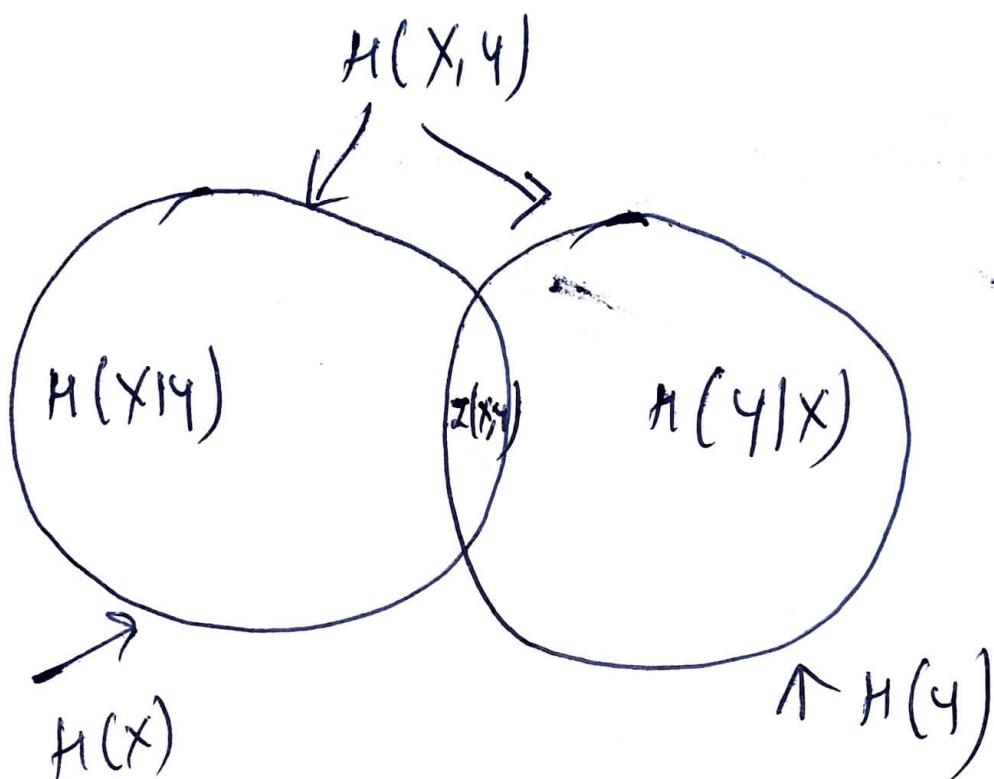
Venn diagram summary of concepts and relationship

Entropy $H(X)$, $H(Y)$

Joint Entropy $H(X,Y) = H(X) + H(Y)$

Conditional Entropy $H(X|Y)$, $H(Y|X)$

Mutual Information $I(X;Y) = H(X) - H(X|Y)$



Distance $D(x,y)$ Btw two random variable
 x and y

The amount by which the joint entropy of two random variables exceeds their mutual information is a measure of the "distance" b/w them

$$D(x,y) = H(x,y) - I(x;y)$$

↓ ↓

Joint mutual.

kullback-leibler distance, or kl divergence

quantifies the difference btwn two probability distributions. Specifically, it measures how one probability distribution diverges from second, expected distribution.

For two discrete probability P (true distribution) and Q (the approximati distribution) the kl divergence from Q to P is defined as

$$D_{KL}(P||Q) = \sum_n P(x) \log \frac{P(x)}{Q(x)}$$

$P(x)$ is the probability of event x under distribution P

$Q(x)$ is the probability of event x under distribution Q .

My notes
on KL Divergence

Covariance matrix

is a matrix that summarizes the covariances between multiple random variables. each element in matrix represents the covariance b/w two variables and the diagonal elements represent the variance of each variable

For n random variables X_1, X_2, \dots, X_n , the Covariance matrix Σ is defined as

$$\Sigma = \begin{bmatrix} \text{Cov}(X_1, X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_n) \\ \text{Cov}(X_2, X_1) & \text{Cov}(X_2, X_2) & \dots & \text{Cov}(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_n, X_1) & \text{Cov}(X_n, X_2) & \dots & \text{Cov}(X_n, X_n) \end{bmatrix}$$

→ if Covariance matrix have only diagonal elements then all random variable X_1, X_2, \dots, X_p are independent to each other or uncorrelated

AutoCovariance Covariance with same R.V

$$C = E[(\underline{x} - \mu_n)(\underline{x} - \mu_n)^H]$$

Covariance matrix is always Conjugate Symmetric (Hermitian)

(25)

Q Calculate SVD, eigen vectors
and eigenvalues

$$\begin{bmatrix} 4 & 0 \\ 3 & -5 \end{bmatrix}$$

Soln

$$A = U \Sigma V^T$$

Step 1 Find V^T

$$A^T \cdot A = \begin{bmatrix} 4 & 3 \\ 0 & -5 \end{bmatrix} \begin{bmatrix} 4 & 0 \\ 3 & -5 \end{bmatrix} = \begin{bmatrix} 16+9 & 0-15 \\ 0-15 & 0+25 \end{bmatrix}$$

$$= \begin{bmatrix} 25 & -15 \\ -15 & 25 \end{bmatrix}$$

To Find Eigen values

$$|A^T A - \lambda I| = 0 \rightarrow \text{characteristic eqn.}$$

$$\begin{vmatrix} 25-\lambda & -15 \\ -15 & 25-\lambda \end{vmatrix} = 0$$

$$\boxed{\lambda^2 - 50\lambda + 500 = 0}$$

$$S_1 = \text{Trace}(A^T A) = 25 + 25 = 50 \quad (20)$$

$$S_2 = \begin{vmatrix} 25 & -15 \\ -15 & 25 \end{vmatrix} = 625 - 225 = 400$$

$$\lambda^2 - 50\lambda + 400 = 0$$

$$\lambda^2 - 40\lambda - 10\lambda + 400 = 0$$

$$\boxed{\lambda = 40, \lambda = 10} \quad \text{eigen values}$$

Eigen Vector at 40

$$A^T A \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 40 \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$25v_1 - 15v_2 = 40v_1 \quad \text{---(1)}$$

$$-15v_1 + 25v_2 = 40v_2 \quad \text{---(2)}$$

consider (1)

$$-15v_1 - 15v_2 = 0$$

$$v_1 = -v_2$$

choosing $v_1 = 1$ we get

$$v_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\text{L}_2 \text{ Norm} = \sqrt{(1)^2 + (-1)^2} = \sqrt{2}$$

$$\text{Eigenvector} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix}$$

eigen vector at 10

$$A^T A \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 10 \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$25v_1 - 15v_2 = 10v_1 \quad \text{---(1)}$$

$$-15v_1 + 25v_2 = 10v_2 \quad \text{---(2)}$$

consider (1)

$$15v_1 - 15v_2 = 0$$

$$v_1 = v_2$$

choosing $v_1 = 1$ we get

$$v_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\text{L}_2 \text{ Norm} = \sqrt{(1)^2 + (1)^2}$$

$$\text{Eigenvector} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

Eigen Vectors matrix

27

$$V = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

Step 2 Calculate Σ which is singular values at diagonal in descending order.

$$\sigma^1 = \sqrt{40} = 2\sqrt{10} \quad \sigma^2 = \sqrt{10}$$

$$\Sigma = \begin{bmatrix} 2\sqrt{10} & 0 \\ 0 & \sqrt{10} \end{bmatrix}$$

Step 3 Calculate U

$$A \cdot A^T = \begin{bmatrix} 4 & 3 \\ 3 & -5 \end{bmatrix} \begin{bmatrix} 4 & 3 \\ 0 & -5 \end{bmatrix} = \begin{bmatrix} 16+0 & 12+0 \\ 12+0 & 9+25 \end{bmatrix}$$

$$= \begin{bmatrix} 16 & 12 \\ 12 & 34 \end{bmatrix}$$

$$|A \cdot A^T - \lambda I| = 0$$

Same as Step ①

$$16x_1 + 12x_2 = 0 \quad \lambda^2 - S_1\lambda + S_2 = 0$$

$$\left\{ \begin{array}{l} S_1 = 16+34 = 50 \\ S_2 = 544-144 = 400 \end{array} \right.$$

$$\lambda^2 - 50\lambda - 10\lambda + 400 = 0$$

$$\lambda = 40, \lambda = 10$$

Eigenvalues 40, 10

Eigen vector at $\lambda = 40$

$$A A^T \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = 40 \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

$$\begin{bmatrix} 16 & 12 \\ 12 & 34 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} 40u_1 \\ 40u_2 \end{bmatrix}$$

$$16u_1 + 12u_2 = 40u_1 \quad \text{---(1)}$$

$$12u_1 + 34u_2 = 40u_2 \quad \text{---(2)}$$

Consider (1)

$$12u_2 = 24u_1$$

$$u_2 = 2u_1$$

$u_1 = 1$ we get

$$u_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

Normalize $u_1 =$

$$\|u_1\| = \sqrt{1^2 + 2^2} = \sqrt{5}$$

$$U = \begin{bmatrix} \frac{1}{\sqrt{5}} & -\frac{2}{\sqrt{5}} \\ \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{5}} \end{bmatrix} \quad \Sigma = \begin{bmatrix} 2\sqrt{10} & 0 \\ 0 & \sqrt{10} \end{bmatrix} \quad V^T = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

$$A = U \Sigma V^T = \begin{bmatrix} \frac{1}{\sqrt{5}} & -\frac{2}{\sqrt{5}} \\ \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{5}} \end{bmatrix} \begin{bmatrix} 2\sqrt{10} & 0 \\ 0 & \sqrt{10} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

2

Eigen vector at $\lambda = 10$

$$\begin{bmatrix} 16 & 12 \\ 12 & 34 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} 10u_1 \\ 10u_2 \end{bmatrix}$$

$$16u_1 + 12u_2 = 10u_1 \quad \text{---(1)}$$

$$12u_1 + 34u_2 = 10u_2 \quad \text{---(2)}$$

consider (1)

$$-12u_2 = 6u_1$$

$$u_1 = -2u_2$$

choose $u_2 = 1$

$$u_1 = -2$$

$$u_2 = \begin{bmatrix} -2 \\ 1 \end{bmatrix}$$

$$\|u_2\| = \sqrt{(-2)^2 + 1^2} = \sqrt{5}$$

Correlation

~ measure the strength and direction of a linear relationship between two different variables

For two R.V X and y Correlation.

Coefficient (often denoted ρ or r) is

$$\rho_{x,y} = \frac{\text{Cov}(x,y)}{s_x s_y}$$

$\text{Cov}(x,y)$ is Covariance between X and y

s_x and s_y are standard deviations of X and y

The Correlation Coefficient range from -1 to +1

+1 : perfect positive linear relationship

-1 : perfect negative linear relationship

0 : No linear relationship

Correlation is widely used in statistics
and machine learning to explore how one variable
might predict or effect another.

Auto Correlation

~ measure correlation of a signal with a delayed copy of itself. it assesses how a signal or time series correlates with itself over various time lags. For a time series x_t , the autocorrelation at lag k (denoted $\rho(k)$) is defined as

$$\rho(k) = \frac{\text{Cov}(x_t, x_{t+k})}{\sigma_x^2}$$

where

$\text{Cov}(x_t, x_{t+k})$ is covariance btw the time series and itself at a time lag of k

σ_x^2 is variance of time series

key differences

- Correlation is btw two different variables; autocorrelation is within a single variable or time series at different time steps
- Correlation tells you how two variable move in relation to each other. while autocorrelation tells you how a variable relate to itself over time.

Stationarity

~ is a key concept in time series analysis referring to a time series whose statistical properties do not change over time, in other words, a time series is stationary if its mean, variance and auto correlation structure remain constant over time. Stationarity is important because many statistical models especially timeseries forecasting, require data to be stationary for accurate modeling and prediction.

Types

① Strict stationarity

A time series is said to be strictly stationary if the joint distribution of any subset of observations is the same, no matter when the observations are taken.

Mathematically for a time series X_t , it means that for any time points t_1, t_2, \dots, t_n and any

any lag k , the joint distribution of
 $(X_{t_1}, X_{t_2}, \dots, X_{t_n})$ is identical to that of
 $(X_{t_1+k}, X_{t_2+k}, \dots, X_{t_n+k})$

② weak stationarity (second order stationarity)

~ is a less strict condition and is more commonly used in practice. A time series is weakly stationary if it satisfies the following conditions

i) Constant mean :- mean of series $E[X_t]$ is constant and does not depend on time t

$$E[X_t] = \mu \quad \text{for all } t$$

ii) Constant variance :- variance of series $\text{Var}(X_t)$ is constant and does not depend on time t .

$$\text{Var}(X_t) = \sigma^2, \quad \text{for all } t$$

iii) Constant autocovariance structure :- autocovariance btw values at two points t and $t+k$ depends only on the time lag k , not on specific time t

$$\text{Cov}(X_t, X_{t+k}) = r(k)$$