

Introduction ¶

A major task of probability theory is applying known probability densities to derive number characteristics such as expectation values, variances etc. Statistical inference, however, focuses on estimating parameters of unknown probability densities first, and then apply the estimated distributions.

The following is a summary of some important probability distributions. The relations among these distributions are also studied in order to highlight the specific conditions for each specific distribution. For a more complete list of probability distributions, see the following link:

https://en.wikipedia.org/wiki/List_of_probability_distributions (https://en.wikipedia.org/wiki/List_of_probability_distributions) For their number characteristics, click each individual link.

Probability distribution for discrete time and space

The starting point: Bernoulli distribution

The PMF is:

$$f(x, p) = p^x (1 - p)^{1-x} \quad x \in \{0, 1\}$$

Bernoulli distribution is also called 0-1 distribution. Its expectation and variance are p and $p(1 - p)$ respectively.

Probability of x occurrences after n Bernoulli trials: Binomial distribution

$$f(x, n, p) = C_n^x p^x (1 - p)^{n-x} \quad x = 0, 1, 2, \dots, n$$

Flipping n times and assume x up and $n - x$ down. There are C_n^x such cases, and therefore the above formula.

Conditions are required for the Binomial distribution:

- Two possible outcomes for each experiment. The probability for each outcome is constant.
- Each experiment is **independent** and identical.

These conditions are also required in other five distributions below that are strongly related binomial distribution. The expectation and variance are np and $np(1 - p)$, which are n times of those of Bernoulli distribution.

Probability of first occurrence after n Bernoulli trials: Geometric distribution

$$f(n, p) = p(1 - p)^{n-1}$$

Note we cannot obtain this by setting $x = 1$ in the binomial distribution. This is because we only want the $x = 1$ case for the first time, but not all possible cases where we flip n times and have only one occurrence.

Probability of first x th occurrences after n Bernoulli trials: Negative Binomial distribution

Also called Pascal distribution. If the x th occurrence occurs on the n th trial, that implies $(x - 1)$ th occurrences must have happened in the previous $(n - 1)$ trials. Because these $(x - 1)$ occurrences can occur in any order, the probability of this happening is just Binomial PMF:

$$P(X = x - 1; n - 1) = C_{n-1}^{x-1} p^{x-1} (1 - p)^{n-x}$$

Furthermore, the probability that we will have an occurrence on the n th trial is p , therefore

$$f(x; p) \equiv P(X = x) = C_{n-1}^{x-1} p^x (1 - p)^{n-x}$$

This is true for $n = x, x + 1, \dots$. For $n < x$, the probability is zero.

There are alternative formula in the following link https://en.wikipedia.org/wiki/Negative_binomial_distribution (https://en.wikipedia.org/wiki/Negative_binomial_distribution).

Probability distribution for continuous time and space

Sometimes we want to know the probability of an event occurring within a continuous time period or space. If we separate the continuous time or space into many tiny time or space intervals, then it is possible to analyze these distributions with the same methods introduced in previous chapter.

The value of PDF at value x is not the probability of $P(X = x)$. So PDF $f(x)$ can take on values larger than one (but the integral of $f(x)$ over any subset of R will be at most one).

From binomial to Poisson distribution

We now calculate the probability that there would be x events in t time, given some average number of occurrences ($\lambda = \nu t$) in time t . This can be handled with Bernoulli trials.

- Divide time t period into n equal intervals.
- Probability of occurrence in any interval is $p = \lambda/n$.

If the conditions for binomial distribution (see previous chapter) are satisfied, then the probability that the event will occur x times within n trials is

$$P(N = x, n) = C_n^x \left(\frac{\lambda}{n} \right)^x \left(1 - \frac{\lambda}{n} \right)^{n-x} = \frac{\lambda^x}{x!} e^{-\lambda} \quad (n \rightarrow \infty) = \frac{(vt)^x}{x!} e^{-vt}$$

Conditions required to apply Poisson distribution

- From its relation to binomial distribution, we know Poisson distribution must also follow the conditions similar to binomial: (1) Each event is independent and identical. (2) The mean occurrence rate vt or λ (corresponding to np in binomial) is constant. **Note definition of λ .** There might be different definitions for different distributions. (3). The probability of two or more occurrences at same time or location is negligible.
- From the derivation above, the formula is true only for very big n . So Poisson distribution is only true for rare event, or small $p = \frac{\lambda}{n}$.

From these conditions we know that when $n \rightarrow \infty$ and when p is small, then binomial distribution can be replaced by Poisson distribution, which is easier to calculate. However, we must also be careful of other conditions such as iid, constant rate to apply Poisson formula. For example, the number of student who arrive at the student union per minute is not a constant (low rate during class time for example), and the arrivals of individual students are not independent (they tends to come in groups).

From geometrical to exponential distribution

When doing experiments in each tiny time interval during a long time period. What is the probability it takes time duration of t to have the first occurrence? This is similar to the geometric distribution for e.g. binomial coin flipping. What is the probability of first occurrence after n trials? Here taking time duration t is similar to taking how n trials.

Assuming T_1 is the time for first event. Then the probability that first event has not occurred is

$$P(T_1 > t) \equiv P(X = 0; t, v) = \frac{(vt)^0}{0!} e^{-vt} = e^{-vt}$$

The CDF for the occurrence of first event is therefore $P(T \leq t) = 1 - e^{-vt}$, and PDF is derivative $P(T_1 = t) = ve^{-vt}$

From negative binomial to Gamma distribution

Similar to geometric distribution, we can use Poisson distribution to derive the Gamma distribution.

Assuming T_k is the time until k th event. Then the probability of k or more occurrences in time t is

$$P(T_k \leq t) = \sum_{x=k}^{\infty} P(X_t = x) = 1 - \sum_{x=0}^{k-1} \frac{(vt)^x}{x!} e^{-vt}$$

Thus the PDF is:

$$P(T_k = t) = \frac{d}{dt}P(T_k \leq t) = \frac{v(vt)^{k-1}}{(k-1)!} e^{-vt}$$

This special form of Gamma distribution give the probability of k th event after a time period of t . This is similar to the k th events after, e.g., flipping n times of coins, as described by negative binomial distribution. **Note the v here has the same definition of λ used in several places discussed later. But this definition is different from that in the Poisson distribution before.** General Gamma is given below.

Gamma distribution and its derivative distributions

General

https://en.wikipedia.org/wiki/Gamma_distribution (https://en.wikipedia.org/wiki/Gamma_distribution)

A shape parameter is a kind of numerical parameter of a parametric family of probability distributions. Such a parameter must affect the shape of a distribution rather than simply shifting it (as a location parameter does) or stretching/shrinking it (as a scale parameter does).

The gamma distribution is a two-parameter family of continuous probability distributions. It is a generalization of the exponential distribution. The exponential distribution, Erlang distribution, and chi-squared distribution are special cases of the gamma distribution. There are **three different parametrizations** (all positive real parameters) in common use:

- With a shape parameter κ and a scale parameter θ .
- With a shape parameter $\alpha = \kappa$ and an inverse scale parameter $\beta = \frac{1}{\theta}$, called a rate parameter.
- With a shape parameter κ and a mean parameter $\mu = \kappa\theta = \frac{\alpha}{\beta}$

The parameterization with κ and θ appears to be more common in econometrics and certain other applied fields, where for example the gamma distribution is frequently used to model waiting times. The parameterization with α and β is more common in Bayesian statistics.

Normally κ or α can be understood as the number of events k in the discrete case. When κ or α is an integer, $\Gamma(\kappa + 1) = \kappa!$. The other scale parameter θ or inverse scale parameter β is related to the rate of the events. In fact the β here is just the rate parameter λ , v , or r in many other distributions discussed later. However, note the λ in Poisson distribution is often defined as vt , so the λ in Poisson case is normally different from other distributions.

Gamma function and gamma distribution

Gamma function is defined as $\Gamma(\kappa) = \int_0^\infty x^{\kappa-1} e^{-x} dx$ $\kappa > 0$. Making the substitution $y = x/\lambda$ where λ is a positive real constant, the gamma function becomes $\Gamma(\kappa) = \int_0^\infty (\lambda y)^{\kappa-1} e^{-\lambda y} \lambda dy$ $\kappa > 0$. Divide both sides by $\Gamma(\kappa)$ yielding

$$\int_0^{\infty} \frac{\lambda^{\kappa} y^{\kappa-1} e^{-\lambda y}}{\Gamma(\kappa)} dy = 1,$$

which is suitable for a PDF of a random variable.

A continuous random variable X with pdf $f(x) = \frac{\lambda^{\kappa} x^{\kappa-1} e^{-\lambda x}}{\Gamma(\kappa)}$ $x > 0$ for some real constant $\lambda > 0$ and $\kappa > 0$ is a gamma(λ, κ) random variable. Check the $f(x)$ for $\kappa > 1, \kappa = 1, \kappa < 1$.

Properties of gamma function:

$$\Gamma(n+1) = n! \quad n = 0, 1, \dots$$

$$\Gamma(\alpha) = (\alpha-1)\Gamma(\alpha-1) \quad \alpha > 0$$

$$\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$$

Special cases of gamma distribution

Exponential distribution

When $\kappa = 1$

$$f(x) = \lambda e^{-\lambda x} \quad x > 0$$

Comparing to $P(T_1 = t) = \nu e^{-\nu t}$ obtained earlier, the λ here is just the rate ν , or sometimes denoted as r . Both ν and r are rate and have a dimension of 1 over time. However, the λ in the Poisson distribution before is defined as $\lambda = \nu t = rt$. So the λ there is different from the λ here. **Be careful of this point.** In many places, λ, ν, r etc., have the same meaning for rate. However, in Poisson case, it is often defined as νt or rt .

Erlang distribution

When κ is a positive integer k

$$f(x) = \frac{\lambda^k x^{k-1} e^{-\lambda x}}{(k-1)!} \quad x > 0$$

If X_1, X_2, \dots, X_k are iid exponential (λ) random variable, then $X_1 + X_2 + \dots + X_k \sim \text{Erlang}(\lambda, k)$. **Note the definition of λ here is same as an example later, but different from that of Poisson.**

χ^2 distribution

When κ is $k/2$, where k is a parameter known as the degrees of freedom, and $\lambda = 1/2$

$$f(x) = \frac{x^{k/2-1} e^{-x/2}}{\Gamma(k/2) 2^{k/2}} \quad x > 0$$

The χ^2 distribution is related to the standard normal distribution. If a random variable Z has the standard normal distribution, then Z^2 has the χ^2 distribution with one degree of freedom. If Z_1, Z_2, \dots, Z_k are independent standard normal variables, then $Z_1^2 + Z_2^2 + \dots + Z_k^2$ has a χ^2 distribution with k degrees of freedom, as described by the PDF above.

Gamma-distribution related application

If the events are occurring according to the Poisson distribution the time till the occurrence of the first event is described by the exponential distribution. The time till the occurrence of the second event is described by the Gamma distribution with $k = 2$ (summation of the time intervals between two consecutive events). Comparing this to the discrete counterparts of geometric and negative binomial distributions discussed earlier.

For example it is known that under free flow conditions vehicular arrival pattern follows Poisson distribution. If the mean arrival rate is 5 vehicles/minute, the headways between consecutive vehicular arrivals follow exponential distribution with the following PDF $f(t) = \frac{1}{5} e^{-5t}$.

The time gap between every two vehicles follows the Gamma distribution with PDF $f(t) = \frac{1}{5} (\frac{1}{5} t) e^{-5t}$. The time at which k th vehicle

arrives at the measurement location follows the PDF $f(t) = \frac{1}{5} \frac{(\frac{1}{5} t)^{k-1}}{(k-1)!} e^{-\frac{1}{5} t}$. Generally this can be written as

$$f(t) = \lambda \frac{(\lambda t)^{k-1}}{\Gamma(k)} e^{-\lambda t}$$

Note the definition of λ is different from that in Poisson distribution. From above, we can understand that gamma distribution is a generalization of exponential distribution. To have an intuition on gamma distribution, imagine the random variable as the waiting time of some event. Or in the discrete case, the number of total flips to have the k th event.

Conditional probability and Bayes rule

Conditional probability

see machine learning notes

Bayes rule

see machine learning notes

Likelihood function vs probability

https://en.wikipedia.org/wiki/Likelihood_function (https://en.wikipedia.org/wiki/Likelihood_function)

- In statistics, a likelihood function (or just the likelihood) is a **particular function of the parameter** of a statistical model given data. In informal contexts, "likelihood" is often used as a synonym for "probability". In statistics, the two terms have different meanings. **Probability is function of x given parameter θ , while likelihood is function of θ given x . where x is the outcome of random variable X .** Likelihood is used with each of the four main foundations of statistics: frequentism, Bayesianism, likelihoodism, and AIC-based.
- Let X be a discrete random variable (discrete) with probability mass function p depending on a parameter θ , then the function $L(\theta | x) = p_\theta(x) = P_\theta(X = x)$ considered as a function of θ , is the likelihood function of θ , given the outcome x of the random variable X . For continuous random variable, we have similar definitions.
- An intuitive understanding: Because likelihood is function of parameters for a given x , we can think whether x is LIKELY (hence likelihood) to be in the probability distribution of given by θ_1 , or θ_2A more specific example about biased coin tossing in explaining expectation maximization algorithm (by Do and Batzoglou, 2008). After a tossing experiment, i.e. the x , we may calculate whether this experiment is LIKELY from coin A or coin B, which are modeled by different parameters.
- **Confusion of different notations**
 - When judging whether it is a likelihood or a probability, we need make sure which is the function argument, but not the specific written form. The $p_\theta(x)$ above is more like a function of x given θ . But it describes a function of θ given x . So it is a likelihood function.
 - When we describe the probability (not likelihood) of "the value x of X given the parameter value θ , we often write it as $P(X = x | \theta)$. Although formally it is like a conditional probability, it is not. So instead, we write it as $P(X = x; \theta)$ to emphasize that it is not a conditional probability.
 - The likelihood is sometimes written as $L(\theta | x)$ and sometimes as $L(x | \theta)$. In other words, the order of the appearance of θ does not matter. The key is whether θ is taken as a function argument (variable). Anyway, the $|$ sign does not indicate conditional probability.
- In the Maximum Likelihood Estimation (MLE), we are just maximizing the likelihood (**use this fact to remember the concept of likelihood**). In deriving many supervised learning algorithms such as those in generalized linear models, we usually maximize the likelihood with the form $P(y | x; \theta)$. Here y is conditioned on x but parameterized on θ . However, because the main purpose is treating $P(y | x; \theta)$ as a function of θ and find its optimal value, we are maximizing likelihood, rather conditional probability (**Likelihood seems contain the meaning of 'might be', 'It is not for sure yet' as we are just trying to estimate the value of parameter. This is another way to remember the concept of likelihood in statistics.**
- In deriving unsupervised learning algorithms with expectation maximization (EM), we are still maximizing likelihood but with an iterative approach. The likelihood function takes a form of joint probability, but essentially we should take it a likelihood as it is a function of parameters.

Expectation

Mean/average (or just expectation), variance, covariance, correlation etc., are obtained by expectation operator acting on the corresponding quantity. In other words, they are all 'expectation values'. However, the 'expectation value' is often used to refer to only the mean or average of a random variable. Other related concepts include sample mean, sample variance, sample covariance, covariance matrix, correlation matrix.

Expectation value / mean / average

- Finite case: $E(X) = \sum_{i=1}^k x_i p_i$. The expression of $E(x) = \frac{1}{k} \sum_{i=1}^k x_i$ is a special case of $p_i = \frac{1}{k}$.
- Countably infinite case $E(X) = \sum_{i=1}^{\infty} x_i p_i$, where $\sum_{i=1}^{\infty} |x_i| p_i$ converges.
- Absolutely continuous case $E(X) = \int_{\mathbb{R}} x f(x) dx$.

Variance

The variance of a random variable is the **expectation value** of the squared deviation from the mean.

- General definition $Var(X) = E[(X - E[X])(X - E[X])] = E((x - \mu)^2)$, or $Var(X) = E(X^2) - (E(X))^2$.
- Discrete random variable $Var(X) = \sum_{i=1}^n p_i (x_i - \mu)^2$, or $Var(X) = \sum_{i=1}^n p_i x_i^2 - \mu^2$. For a set of n equally likely values, p_i can be replaced by $\frac{1}{n}$. In this case, it can also be written as other forms without referring to mean (see Wikipedia).
- Continuous random variable $Var(X) = \int (x - \mu)^2 f(x) dx = \int x^2 f(x) dx - \mu^2$.

Covariance

The covariance is the expectation value of quantity $(X - E[X])(Y - E[Y])$. Variance is the special case of covariance where X and Y are same. So variance is also called auto-variance.

- General definition $cov(X, Y) = E[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y]$.
- Discrete variables If the random variable pair (X, Y) can take on values (x_i, y_i) for $i = 1, 2, \dots, n$ with equal probabilities $\frac{1}{n}$, the $cov(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - E(X))(y_i - E(Y))$. Like variance, it can also be expressed without directly referring to the means. If each pair of value is not with equal probability, then $\frac{1}{n}$ need be replaced with probability p_i for each pair of data.
- When $X = Y$, then covariance becomes variance.

Correlation

It is just the scaled form of covariance. $\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$.

Correlation (covariance) and dependence

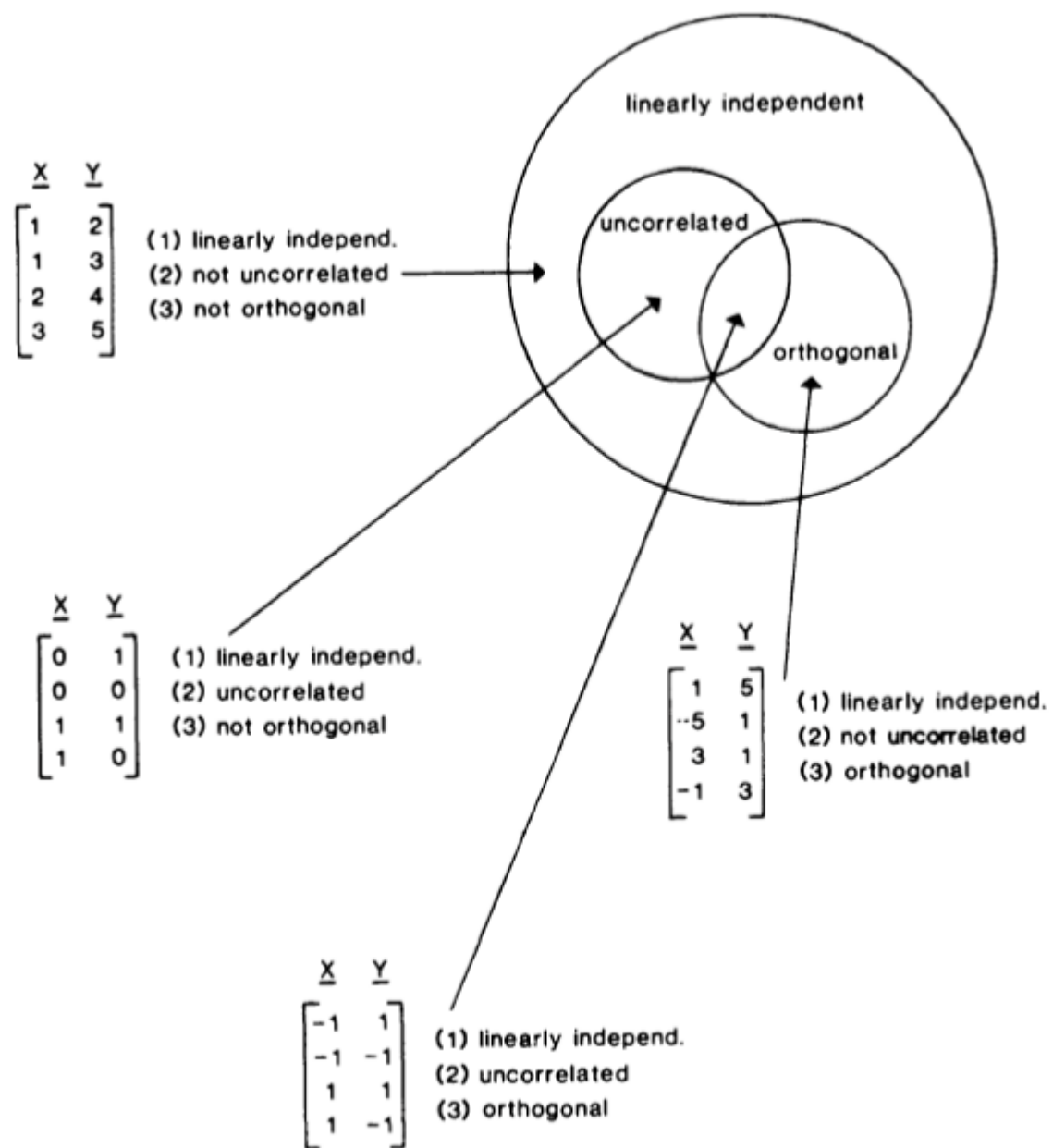
- Non-zero off-diagonal elements of the covariance matrix $X^T X$ or XX^T (data centered) indicate that the relevant variables must be linearly dependent.
- The non-zero covariance/correlation elements, however, give no clue about the nonlinear dependence among variables. The nonlinearity requires defining nonlinear correlation / covariance, which should include more than two factors in the product of normal covariance definition.

Correlation and basis change

- Two correlated variables will become linearly uncorrelated in the eigen-basis of covariance matrix. This does not mean that the correlation between the two variables are gone, but just means that two new variables, which are linear combination of the original variables, are no longer linear correlated in the new basis.

Uncorrelated, orthogonal and independent

- If we have only diagonal nonzero elements in covariance matrix, then it is safe to say there is no linear correlation, or linear dependence among relevant variables. However, it is in general not true to say that the variables are independent. They are linearly independent, but not necessarily nonlinearly independent.
- Assuming X is uniformly distributed in $(-1, 1)$, then it can be shown $\text{cov}(X, X^2) = 0$. The relationship between X and X^2 is non-linear, while correlation and covariance are measures of linear dependence between two variables. This shows two uncorrelated variables **does not in general imply that they are independent**.
- Independence implies uncorrelation but uncorrelation DOES NOT imply independence. However, when two variables are Gaussian, then uncorrelation and independence are equivalent.
- Relation of LINEAR independence, orthogonality and uncorrelation
 - Linearly independent vectors are those vectors that do not fall along the same line; that is, there is no multiplicative constant that will expand, contract, or reflect one vector onto the other!
 - Orthogonal vectors are a special case of linearly independent variables!
 - Uncorrelated vectors imply that once each variable is centered then the vectors are perpendicular.



Random vector

- The number characteristics such as variance, covariance etc. introduced earlier are all define on the **scalar random variable**, meaning the outcome of a random variable can only be a scalar. Now we extend the usual random variable to random vector.
- A random vector is a random variable with multiple dimensions. Each element of the vector is a scalar random variable. Each element has either a finite number of observed empirical values or a finite or infinite number of potential values. The potential values are specified by a theoretical joint probability distribution. **Comments: a random vector is also a random variable, except that this variable is a high-dimensional variable. Unlike the scalar outcome from a usual random variable, the outcome of a random vector is a vector.**

Covariance matrix

https://en.wikipedia.org/wiki/Covariance_matrix (https://en.wikipedia.org/wiki/Covariance_matrix) A covariance matrix is a matrix whose element in the i, j position is the covariance between the $i - th$ and $j - th$ elements of a **random vector**.

Defining a \mathbf{X} , \mathbf{Y} as random vectors (note they are not the common scalar random variables X, Y), X_i and Y_i are scalar random variables.

If the entries in the column vector $\mathbf{X} = (X_1, X_2, \dots, X_n)^T$ are random variables (scalar), each with finite variance and expected values, then the covariance matrix $K_{\mathbf{X}\mathbf{X}}$ is the matrix whose (i, j) entry is the covariance (in the sense of scalar random variables):

$$K_{X_i X_j} = \text{cov}[X_i, X_j] = E[(X_i - E[X_i])(X_j - E[X_j])]$$

The definition above is equivalent to the matrix equality

$$K_{\mathbf{X}\mathbf{X}} = \text{cov}[\mathbf{X}, \mathbf{X}] = E[(\mathbf{X} - \mu_{\mathbf{X}})(\mathbf{X} - \mu_{\mathbf{X}})^T] = E[\mathbf{X}\mathbf{X}^T] - \mu_{\mathbf{X}}\mu_{\mathbf{X}}^T$$

where $\mu_{\mathbf{X}} = E[\mathbf{X}]$.

Comments:

- Because \mathbf{X} is a vector, so $\mathbf{X}\mathbf{X}^T$ is a matrix with elements such as $X_1X_1, X_1X_2, \dots, X_nX_n$, where X_i is scalar random variables but not a single outcome of the random variable. Then taking expectation value of each elements with $E(X_iX_j)$, we obtain the (i, j) element of the covariance matrix.
- Be sensitive to the transpose sign appeared in the definition of $K_{\mathbf{X}\mathbf{X}}$. This indicates \mathbf{X} is a random vector, rather a scalar random variable.

Conflicting nomenclatures and notations

Recall the variance of scalar random variable $\text{Var}(X) = E[(X - E[X])(X - E[X])]$. We find that the **covariance matrix** $K_{\mathbf{X}\mathbf{X}}$ is just the extension of the variance of scalar random variables. In other words, $K_{\mathbf{X}\mathbf{X}}$ is also called **variance of the random vector** \mathbf{X} .

Therefore, depending whether you focus on the \mathbf{X} as a whole or as a collection of individual scalar random variables, $K_{\mathbf{X}\mathbf{X}}$ is sometimes

called $\text{var}(\mathbf{X})$ or $\text{cov}(\mathbf{X})$. Also, sometimes $K_{\mathbf{X}\mathbf{X}}$ is also called **variance-covariance matrix** since the diagonal terms are in fact variances (in the sense of scalar random variables). Finally $K_{\mathbf{X}\mathbf{X}}$ is also called auto-covariance matrix.

$R_{\mathbf{X}\mathbf{X}}$ vs $\text{corr}(\mathbf{X})$

https://en.wikipedia.org/wiki/Covariance_matrix (https://en.wikipedia.org/wiki/Covariance_matrix)

$$R_{\mathbf{X}\mathbf{X}} = E[\mathbf{X}\mathbf{X}^T]$$

is called **correlation matrix**, and

$$\text{corr}(\mathbf{X}) = (\text{diag}(K_{\mathbf{X}\mathbf{X}}))^{-\frac{1}{2}} K_{\mathbf{X}\mathbf{X}} (\text{diag}(K_{\mathbf{X}\mathbf{X}}))^{-\frac{1}{2}}$$

is called **the matrix of correlation coefficients**. This basically normalized each entry of the covariance matrix with a specific dominator.

Sample mean, variance, covariance matrix

https://en.wikipedia.org/wiki/Sample_mean_and_covariance (https://en.wikipedia.org/wiki/Sample_mean_and_covariance)

We introduced the random vector \mathbf{X} in earlier sections. This can be understood as a high-dimension random variable. Each outcome of this random variable is thus also a vector denoted as \mathbf{x} , each entry of \mathbf{x} thus corresponds to the corresponding scalar random variable, which is an entry in the random vector \mathbf{X} .

Assuming the $\mathbf{x} \in \mathbb{R}^d$ and we have n outcomes (samples), i.e. $\mathbf{x}_i, i = 1, 2, \dots, n$. Then the **sample mean vector** has the form

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_i^n \mathbf{x}_i = \begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_i \\ \vdots \\ \bar{x}_d \end{pmatrix}$$

Comments: Here we have sample mean vector, instead of the scalar sample mean, because the population is described by a random vector \mathbf{X} rather than a usual scalar random variable X .

Following the same notation, the previously defined covariance matrix $K_{\mathbf{X}\mathbf{X}}$ can be written as

$$\begin{aligned} K_{\mathbf{X}\mathbf{X}} &= \text{cov}[\mathbf{X}, \mathbf{X}] = E[(\mathbf{X} - \mu_{\mathbf{X}})(\mathbf{X} - \mu_{\mathbf{X}})^T] \\ &= \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \equiv A A^T \end{aligned}$$

where each term in the sum has rank 1. The total rank of covariance matrix is at most $n - 1$.

https://en.wikipedia.org/wiki/Sample_mean_and_covariance (https://en.wikipedia.org/wiki/Sample_mean_and_covariance). If $n = d$, then the matrix is singular and not good for numerical calculations. From here, we know that in machine learning, the bigger the number of samples, the better chance we have stable numerical calculations.

In the above equations we write the covariance matrix in the form of AA^T , where A is formed by stacking n columns of \mathbf{x}_i . Note this is because the multiplication of two matrices can always be expressed as outer product. That is,

$$C = AB \iff C = \sum_{i=1}^n a_i b_i^T,$$

where a_i is column vector of A , and b_i^T is row vector of B . Note however, all the equations here are using column vector to store a collection of random variables. However, in machine learning, we often use rows to store different samples of data. In other words, the random vector we introduced above is a row vector rather than a column vector. In this case, the covariance matrix has the form of $A^T A$.