

Properties of Expected Values (Expectations), III

3. $E(b_1 + X_1) = b_1 + E(X_1)$.

4. $E(\sum_{i=1}^n X_i) = \sum_{i=1}^n E(X_i)$, n is finite.

5. $E(\sum_{i=1}^n b_i X_i) = \sum_{i=1}^n b_i E(X_i)$.

$$E\left(\sum_{i=1}^n \frac{1}{n} X_i\right) = \sum_{i=1}^n \frac{1}{n} E(X_i) = \frac{1}{n} \sum_{i=1}^n E(X_i).$$

Covariance

Assume that $E(|X_1|)$, $E(|X_2|)$ and $E(|X_1 X_2|)$ are all finite. Then the covariance between X_1 and X_2 is defined as

$$\text{Cov}(X_1, X_2) = E(X_1 - E(X_1))(X_2 - E(X_2)).$$

An alternative equivalent definition is (by linearity of expectations).

$$\text{Cov}(X_1, X_2) = E(X_1 X_2) - E(X_1)E(X_2).$$

bilinearity of covariances: let a_1 and a_2 be const's

$$\text{Cov}(a_1 X_1, a_2 X_2) = a_1 \cdot a_2 \cdot \text{Cov}(X_1, X_2)$$

Symmetry: notice $\text{Cov}(X_1, X_2) = \text{Cov}(X_2, X_1)$.

Main property: bilinearity of covariances

Proposition: Let X_1, \dots, X_n and Y_1, \dots, Y_m be rvs with well-defined covariances $Cov(X_i, Y_j)$ for every i and j . Let a_1, \dots, a_n and b_1, \dots, b_m be any constants. Then

$$Cov\left(\sum_{i=1}^n a_i X_i, \sum_{j=1}^m b_j Y_j\right) = \sum_{i=1}^n \sum_{j=1}^m a_i b_j Cov(X_i, Y_j). \quad (\star)$$

Let C be a matrix such that $C_{ij} = Cov(X_i, Y_j)$.

let $a = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix}$; $b = \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix}$. $a^T = (a_1, \dots, a_n)$.

Then $(\star) = a^T \cdot C \cdot b$
 $= b^T \cdot C^T \cdot a$

Variance is a special case of covariance

Variance of a rv. If $X_1 = X_2$, then

$$\begin{aligned} \text{Cov}(X_1, X_2) &= \text{Cov}(X_1, X_1) = E(X_1^2) - (E(X_1))^2 = \text{Var}(X). \\ &= E([X_1 - E(X_1)]^2) \geq 0 \end{aligned}$$

Q: Is there a difference between $E(X^2)$ and $(E(X))^2$?

A: yes. Since $\text{Var}(X_1) \geq 0$, $E(X^2) \geq (E(X))^2$

Variance of a linear combination of rvs

$$\text{Var}\left(\sum_{j=1}^n a_j X_j\right) = \sum_{i=1}^n \sum_{j=1}^n a_i \cdot a_j \cdot \text{Cov}(X_i, X_j)$$

using matrix algebra, $= a^T \cdot C \cdot a$

Random variables X_1 and X_2 are said to be uncorrelated if

$$\text{Cov}(X_1, X_2) = 0.$$

Correlation

Correlation coefficient (loosely, correlation) between rvs X_1 and X_2 that have a finite second moment is defined as

$$\text{Corr}(X_1, X_2) = \frac{\text{Cov}(X_1, X_2)}{\sqrt{\text{Var}(X_1) \cdot \text{Var}(X_2)}}.$$

Let $\rho = \text{Corr}(X_1, X_2)$. If $\rho = 0$, the rvs are uncorrelated.

Q: Why is $|\rho| \leq 1$?

Independence and correlation

Let X_1 and X_2 be rv's with the joint cdf $F_{1,2}$, joint pdf/pmf $f_{1,2}$, marginal cdf's F_1 and F_2 and marginal pdf's/pmf f_1 , f_2 .

Recall that X_1 and X_2 are independent if and only if

$F_{1,2}(x_1, x_2) = F_1(x_1) \cdot F_2(x_2)$ if and only if

$f_{1,2}(x_1, x_2) = f_1(x_1) \cdot f_2(x_2)$ for every x_1, x_2 .

Independence of X_1 and X_2 does not guarantee existence of moments.

However if $E(X_1^2) < \infty$ and $E(X_2^2) < \infty$ and X_1 and X_2 are independent, we have

$$\textcircled{\star} \quad E(X_1 \cdot X_2) = E(X_1) E(X_2) \Rightarrow \text{Cov}(X_1, X_2) = 0, \quad = E(X_1) \cdot E(X_2)$$

i.e., X_1 and X_2 are uncorrelated.

$$E(X_1 X_2) = \int \int x_1 \cdot x_2 \cdot \underbrace{f_{1,2}(x_1, x_2)}_{f_1(x_1) \cdot f_2(x_2) \text{ by indep.}} dx_1 dx_2 = \left(\int x_1 f_1(x_1) dx_1 \right) \cdot \left(\int x_2 f_2(x_2) dx_2 \right).$$

$\text{Cov}(X_1, X_2) = 0$ does not imply that X_1 and X_2 are independent.

Parametric family of distributions

Notation: $\mathcal{F} = \{F_\theta : \theta \in \Theta\}$ is the parametric family of distributions indexed by θ .

Parameterization is the correspondence between θ and F_θ .

Common setup: $\underbrace{X_1, X_2, \dots, X_n}_{\text{iid sample from } F_\theta} \stackrel{\text{iid}}{\sim} F_\theta$.

WMS: X_1, X_2, \dots, X_n is a “random sample” from F_θ .

Warning: in general, “random” does not mean “independent”.

Goals of probability and statistics

Goal of probability: determine $\Pr(T(X_1, \dots, X_n) \in A)$, where T is some function.

To motivate goals of statistics, consider a game:

1. Mother Nature chooses $\theta \in \Theta$ and generates

$$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} F_\theta.$$

2. Goal of statistics/statistician: use info in the sample X_1, \dots, X_n to make inference/learn/guess the value of θ that the nature has chosen.

What is meant by statistical inference?

Inference = estimation + hypothesis testing

Estimation:

1. Point estimation: use of data X_1, X_2, \dots, X_n to “guess” θ using a rv $T(X_1, \dots, X_n)$ that is close to θ in some prob. sense.
2. Set/interval estimation: find a random set/interval $S(X_1, \dots, X_n)$ such that $\Pr(\theta \in S(X_1, \dots, X_n))$ is high.

Hypothesis testing: use the sample to determine if a hypothesis $\theta = \theta_0$ is likely to be true. “Do the data support the hypothesis that $\theta = \theta_0$?”

Other goals of statistics (besides inference): modeling, prediction.

Check out the video about Ritz Casino scam for an amazing illustration of how modeling and prediction are useful in real life

<https://www.youtube.com/watch?v=GnaOM4W-hDE>

Identifiability of a parameterization

Recall the goal of statistics and the game:

1. Mother Nature chooses $\theta \in \Theta$ and generates $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} F_\theta$.
2. Goal of statistics/statistician: use info in the sample X_1, \dots, X_n to make inference/learn/guess the value of θ that the nature has chosen, or the distribution F_θ .

A parameterization $\theta \mapsto F_\theta$ is called identifiable if $\theta \neq \theta'$ implies $F_\theta \neq F_{\theta'}$.

Ex. Non-identifiable parameterization

$$F_\theta = \text{Normal}(\mu_1 + \mu_2, \sigma^2); \quad \theta = (\mu_1, \mu_2, \sigma^2).$$

\uparrow weight of owner \uparrow weight of a pet

$$\theta' = (\mu_1 + c, \mu_2 - c, \sigma^2)$$

corresponds to the same distribution

In linear regression, this is known as "multicollinearity".

Estimators and estimates

Def. An *estimator* of $g(\theta)$ is a statistic $T(X_1, X_2, \dots, X_n)$ that is used to “guess” the value of $g(\theta)$.

Def. An *estimate* of $g(\theta)$ is the value of the estimator $T(X_1, \dots, X_n)$ when $(X_1, \dots, X_n) = (x_1, \dots, x_n)$.

Notice: $T(X_1, \dots, X_n)$ is a rv (an estimator), while $T(x_1, \dots, x_n)$ is a fixed number (an estimate).

Our interest is in finding good estimators.

Two principal methods are the method of moments and the maximum likelihood estimation (discussed later).

To compare the “goodness” estimators, it is necessary to define several criteria, below. It is assumed here that θ is a scalar; this can be generalized to the case when θ is a vector.

Bias

Def. *Bias* of an estimator $\hat{\theta}$ of θ is

$$Bias(\theta|\hat{\theta}) = E(\hat{\theta}) - \theta.$$

Notice that, typically, the bias is a function of θ (and, possibly, of other parameters).

If $Bias(\theta|\hat{\theta}) = 0$ for every value of θ , the estimator $\hat{\theta}$ is called *unbiased*.

Examples:

MSE and its bias-variance decomposition

Def. *Mean squared error* of an estimator $\hat{\theta}$ is

$$MSE(\theta|\hat{\theta}) = E\{(\theta - \hat{\theta})^2\}.$$

Bias-variance decomposition of the MSE:

$$MSE(\theta|\hat{\theta}) = \{Bias(\theta|\hat{\theta})\}^2 + Var(\hat{\theta}).$$

Consequence: bias-variance tradeoff.