

Mathematical Statistics I

Chapter 4: Expected Values

Jesse Wheeler

1. Discrete random variables
2. Continuous random variables
3. Expectation of functions of random variables
4. Variance and Standard Deviation
 Bias-Variance Tradeoff
5. Covariance and Correlation

Discrete random variables

Introduction

- This material comes primarily from Rice (2007, Chapter 4).
- We will cover the ideas of expected value, variance, as well as higher-order moments.
- This includes topics such as conditional expectation, which is one of the fundamental ideas behind many branches of statistics and machine learning.
- For instance, most regression / prediction algorithms are built with the idea of minimizing some conditional expectation.

Expectation: Discrete random variables

Definition: Expectation of discrete random variables

Let X be a discrete random variable with pmf $p(x)$, which takes values in the space \mathcal{X} . The **expected value** of X is

$$E(X) = \sum_{x \in \mathcal{X}} x p(x),$$

provided that $\sum_{x \in \mathcal{X}} |x| p(x) < \infty$; otherwise, the expectation is not defined.

- This is not the most mathematically precise definition of expectation, but a more complete treatment of the topic is outside the scope of this course (See Resnick, 2019).

Expectation: Discrete random variables II

- The concept of the expected value parallels the notion of a *weighted average*.
- That is, we weight each possibility $x \in \mathcal{X}$ by their corresponding probability: $\sum_x x p(x)$.
- $E(X)$ is also referred to as the **mean** of X , and is typically denoted μ or μ_X .
- If the function p is thought of as a weight, then $E(X)$ is the center; that is, if we place the mass $p(x_i)$ at the points x_i , then the balancing point is $E(X)$.
- Like with the pmf and cdf, we often use subscripts to denote which probability law we are using for the expectation, if it is not clear: $E_X(X)$.

Expectation: Discrete random variables III

Roulette

A roulette wheel has the numbers 1 through 36, as well as 0 and 00. If you bet \$1 that an odd number comes up, you win or lose \$1 according to whether that event occurs. If X denotes your net gain, $X = 1$ with probability $18/38$ and $X = -1$ with probability $20/38$. The expected value of X is

$$E(X) = 1 \times \frac{18}{38} + (-1) \times \frac{20}{38} = -\frac{1}{19}.$$

- As you might imagine, the expected value coincides in the limit with the actual average loss per game, if you play many games (Chapter 5).

Expectation: Discrete random variables IV

- Most casino games have a negative expected value by design; you may win some money, but if a large number of games are played, the house will come out on top.

Expectation: Discrete random variables V

Geometric Random Variable

Suppose that items are produced in a plant are independently defective with probability p . If items are inspected one by one until a defective item is found, then how many items must be inspected on average?

Solution:

Expectation: Discrete random variables VI

Poisson Distribution

The $\text{Poisson}(\lambda)$ distribution has pmf $p(k) = \frac{\lambda^k}{k!} e^{-\lambda}$, for all $k \geq 0$. Thus, if $X \sim \text{Pois}(\lambda)$, then what is $E[X]$?

Solution:

Continuous random variables

Expectation: Continuous random variables

Definition: Expectation of continuous random variables

Let X be a continuous random variable with pdf $f(x)$, which takes values in the space \mathcal{X} . The **expected value** of X is

$$E(X) = \int_{x \in \mathcal{X}} x f(x) dx.$$

provided that $\int_{x \in \mathcal{X}} |x| f(x) dx < \infty$, otherwise the expectation is undefined.

- As before, this is not the most mathematically precise definition of expectation, but a more complete treatment of the topic is outside the scope of this course (See Resnick, 2019).

Expectation: Continuous random variables II

- We can still think of $E(X)$ as the center of mass of the density.

Expectation: Continuous random variables III

Exponential(λ) expectation

Let X have an Exponential(λ) density, with $\lambda > 0$. Thus, the pdf of X is given by

$$f_X(x) = \frac{1}{\lambda} e^{-x/\lambda}, \quad 0 \leq x < \infty$$

Find $E[X]$.

Expectation: Continuous random variables IV

Solution.

Expectation: Continuous random variables V

Gamma Density

If X follows a gamma density with parameters α and λ , then the pdf of X is

$$f(x) = \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x}, \quad x \geq 0.$$

Find $E(X)$.

Expectation: Continuous random variables VI

Solution.

Expectation of functions of random variables

Functions of random variables

- We are often interested in functions of random variables:
 $Y = g(X)$.
- Ideas that we have already covered enable us to calculate $E(Y)$.
- For instance, you could use the change-of-variables theorem to get the density of Y , then use the definition to calculate $E[Y]$.
- Fortunately, we don't have to do this. We can instead calculate $E[Y]$ by integrating (or summing) with respect to X :

$$E[g(X)] = \int_{x \in \mathcal{X}} g(x) f(x) dx.$$

- We will justify this for the discrete case.

Functions of random variables II

Theorem 4.1: Expectation of transformed random variables

Suppose that X is a random variable and that $Y = g(X)$ for some function g . Then,

- If X is discrete with pmf $p(x)$:

$$E(Y) = \sum_x g(x) p(x),$$

provided that $\sum_x |g(x)|p(x) < \infty$.

- If X is continuous with pdf $f(x)$:

$$E(Y) = \int_{-\infty}^{\infty} g(x) f(x) dx,$$

provided that $\int |g(x)|f(x) dx < \infty$.

Functions of random variables: proof

Proof:

Functions of random variables: proof II

- The proof for the continuous case is similar, but does require a measure-theoretic approach to integration.
- One important thing to note is that $g(E(X))$ is not usually equal to $E(g(x))$.
- For example, let Z be a standard normal. We know that $E[Z] = 0$, because it's symmetric. However, $P(|Z| > 0) = 1$, thus we can readily deduce that $E[|Z|] \geq 0 = |E[Z]|$.
- This idea can be extended to show that if for all non-negative random variables X that have finite expectation, if $g(x) \leq x$ for some function g , then $E[g(X)] \leq E[X]$.

Expected value of indicator functions

- Another important example of expectations is **indicator** random variables.
- For example, suppose that X is a random variable. Then $Y = 1[X \in A]$ for some $A \subset \mathcal{X}$ is a random variable.

Indicator Random Variable

Let X follow a standard normal distribution, and $A = [-1, 1]$. Then $Y = 1[X \in A]$ is defined as the random variables such that $Y(\omega) = 1$ if $X(\omega) \in A$, and $Y(\omega) = 0$ otherwise.

Expected value of indicator functions II

- Expectations of indicator variables are **probabilities**. Let $Y = 1[X \in A]$.

$$\begin{aligned} E(Y) &= E(1[X \in A]) \\ &= \int_{x \in \mathcal{X}} 1[X \in A] f(x) dx \\ &= \int_{x \in A} f(x) dx = P(X \in A). \end{aligned}$$

- This fact is useful for deriving some important inequalities.
- First, we will show that the expectations of interest actually exist.

Expected value of indicator functions III

- Let X be a continuous random variable with expectation $E(X)$. From our definition, this implies that $\int |x| f(x) dx < \infty$.
- Now suppose that for some random variable $Y = g(X)$ such that $|Y| \leq |X|$. Then we can deduce that $\int |y| f(x) dx < \infty$, and therefore $E[Y]$ exists.
- Now suppose that φ is a non-decreasing, non-negative function, and that for some $a \in \mathbb{R}$, $\varphi(a) > 0$. Then, for all $x \geq a$, $\varphi(x)/\varphi(a) \geq 1$.

Expected value of indicator functions IV

- Define $Y = 1[X \geq a]$. Note that for all possible outcomes $\omega \in \Omega$,

$$Y = 1[X \geq a] \leq \varphi(X)/\varphi(a)1[X \geq a] \leq \varphi(X)/\varphi(a).$$

- Taking expectations of everything (which we argued preserves inequalities),

$$E(1[X \geq a]) = P(X \geq a) \leq \frac{E[\varphi(X)]}{\varphi(a)} = E[\varphi(X)/\varphi(a)].$$

- This inequality is known as **Markov's (general) inequality**, and is very useful for bounding the probability of particular events.

Expected value of indicator functions V

- Specifically, if $\varphi(x) = |x|^p$, with $p > 0$, then because $|X|$ is always positive, φ is non-negative, non-decreasing, and therefore

$$P(|X| \geq a) \leq \frac{E[|X|^p]}{a^p},$$

- If we restrict ourselves to the case where X is non-negative, we get the most standard version of the inequality:

$$P(X \geq a) \leq E(X)/a.$$

Expected value of indicator functions VI

Markov's Inequality in Action

Suppose that an individual is taken randomly from a population that has an average salary of \$50,000. If we assume that salary from the population is approximately independently and identically distributed, we can provide an upper-bound for the probability that the individual is wealthy.

Let X_i be the salary of individual i , randomly drawn from said population. Even though all we know is the average salary, Markov's inequality tells us that:

$$P(X \geq 200,000) \leq \frac{50,000}{200,000} = \frac{1}{4}.$$

Expected value of indicator functions VII

- Returning to expectations of functions of random variables, we can extend to the multi-variate case

Expected value of indicator functions VIII

Theorem 4.2: functions of multiple variables

Suppose that X_1, \dots, X_n are jointly distributed RVs and $Y = g(X_1, \dots, X_n)$. Then

- IF X_i are discrete with pmf $p(x_1, \dots, x_n)$, then

$$E(Y) = \sum_{x_1, \dots, x_n} g(x_1, \dots, x_n) p(x_1, \dots, x_n).$$

- If X_i are continuous with pdf $f(x_1, \dots, x_n)$, then

$$E(Y) = \int_{\mathcal{X}_1, \dots, \mathcal{X}_n} g(x_1, \dots, x_n) f(x_1, \dots, x_n) dx_1 \dots dx_n.$$

In both cases, we need the sum (or integral) of $|g|$ to converge.

Expected value of indicator functions IX

- The proof for the discrete case of Theorem 4.2 follows directly that of Theorem 4.1
- An immediate consequence of Theorem 4.2 is the following

Corollary 4.2.1

If X and Y are independent random variables, and g and h are fixed functions, then

$$E[g(X)h(Y)] = \left(E[g(X)] E[h(Y)] \right),$$

provided that the expectations on the right-hand side exist.

Expected value of indicator functions X

Example: Breaking sticks

A stick of unit-length is broken randomly (uniformly) in two places. What is the average length of the middle piece?

We will interpret this problem to mean that the locations of the two break-points are independent uniform random variables, U_1 and U_2 , and we need to computing $E|U_1 - U_2|$.

Solution:

Linear Combinations of Random Variables

- A useful property of expectation is that it is a **linear operator**.

Theorem 4.3: Linear combinations

If X_1, \dots, X_n are jointly distributed random variables with expectations $E(X_i)$, respectively, and $Y = a + \sum_{i=1}^n b_i X_i$, then,

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

Linear Combinations of Random Variables II

Proof.

Linear Combinations of Random Variables III

- The previous theorem is extremely useful for calculating expected values.
- An obvious example is **sums** of random variables, such as the arithmetic average.
- It's also useful because some distributions can be expressed as the sum of other distributions.
- For instance, we saw in a previous example that the sum of two exponential random variables has a Gamma distribution. Thus, if we know the mean of an exponential, we can readily calculate the mean of a Gamma distribution.

Linear Combinations of Random Variables IV

Expectation of a binomial distribution

Let Y follow a Binomial(p, q) distribution. Find the expected value of Y .

Solution:

Linear Combinations of Random Variables V

Example: Baseball Card Collection

Suppose that you collect baseball cards, that there are n distinct cards, and that on each trial you are equally likely to get a card of any of the types. How many trials would you expect to go through until you had a complete set of cards?

Linear Combinations of Random Variables VI

Example: Group Testing

Suppose that a large number, n of blood samples are screened for a rare disease. If each sample is taken individually, n tests will be required. An alternative approach is group individuals into m groups of size k , pool the blood samples for each group together and perform a test on the pooled sample. If the pooled test is negative, we know all individuals in the group do not have the rare disease; however, if the test is positive, we can then do tests on each individual in the smaller group. What is the expected number of tests that will be conducted using this approach?

Linear Combinations of Random Variables VII

Example: Counting DNA “words”

Within DNA patterns, we might be interested in finding the number of times a particular combination of letters (or “word”) occurs in a DNA sequence. This can be useful for determining if a region of DNA has unusually large occurrences of specific sequences. Assume each sequence is randomly composed of letters A, C, G, T , and that for each location in the sequence, each letter has probability $1/4$. For example, consider occurrence of the “word” $TATA$.

ACTATATAGATATA

In the above sequence, we would count $TATA$ 3 times (counting overlaps). In a sequence of length N , what is the expected number of times a word of length q occurs?

Expected value as a predictor

- One useful property of the expectation is that it serves as a good predictor for the value of a random variable.
- Suppose X is a random variable with well-defined expectation, and that we want to make a prediction for the value of X .
- Denote our predicted value of X as b .
- One common way to measure accuracy using the Mean-Squared Error (MSE), which is defined as:

$$\text{MSE}(b) = E[(X - b)^2].$$

- Here, the closer b is to X , the smaller $(X - b)^2$ is. We take the expectation because X is random.
- By this measure, the best predictor would minimize this error.

Expected value as a predictor II

Theorem: Expectation and MSE

If X is a random variable, then the value b that minimizes $E[(X - b)^2]$ is $b = E[X]$:

$$\operatorname{argmin}_b E[(X - b)^2] = E[X].$$

Proof:

Some comments on expected values

- An important thing to notice about the theorem for linear combinations is that we do not require independence.
- The last example demonstrates this principle. Though I_n is Bernoulli distributed, $\sum_n I_n$ is **NOT** binomial distributed, because the I_n are not independent.
- As an example, if our word is $TATA$, then $I_1 = 1$ implies that $I_2 = 0$, since a $TATA$ at position 1 implies that the second letter starts with A , and thus $TATA$ cannot occur at position 2.
- Despite this, we can still calculate the expected value of a sum by taking the sum of expected values.

Some comments on expected values II

- The expected value can be used as an indication of the central value of the density or frequency function.
- Because of this, the expected value is sometimes referred to as a **location parameter**.
- The expected value is not the only type of location parameter. For instance, the *median* is also a type of location parameter.
- We have seen a lot of parallel between the expected value of a discrete random variable and that of a continuous random variable. This is not a coincidence.
- Specifically, we generally just “swap” and integration with summation, and pdf with pmfs.

Some comments on expected values III

- With a more rigorous definition of expectation, we could define expectation as a **Lebesgue-Stieltjes** integral, with respect to some measure P .
- That is, $E(X) = \int_{\Omega} X dP$, where P is a probability measure. If the probability measure is a counting measure, then the integral *is* a sum.
- Note that this definition does not require the existence of a pdf; in fact, there distributions where the expectation is well-defined, but the pdf is not. These types of distributions do not come up often in standard examples.

Variance and Standard Deviation

Variance

- The expected value is useful for summarizing the average or expected behavior of a random variable.
- We are also often interested in the “spread” of a random variable.
- That is, if the expected value is the center (or location) of a distribution, we want an indication of how dispersed a distribution is around this center.
- The two most common ways to express this idea is the **variance** and **standard deviation** of a random variable.

Variance II

Definition: Variance

If X is a random variable with expected value $E(X)$, then the **variance** of X is

$$\text{Var}(X) = E\left[(X - E(X))^2\right],$$

provided the expectation exists.

Variance III

- Letting $\mu = E[X]$, we can use the identity $g(x) = (X - \mu)^2$, and our expression for $E[g(X)]$ to get a way of calculating the variance.
- If X is a discrete random variable, then by Theorem 4.1,

$$\text{Var}(X) = \sum_i (x_i - \mu)^2 p(x_i),$$

- If X is a continuous random variable, then

$$\text{Var}(X) = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$$

Definition: Standard deviation

If X is a random variable, then the standard deviation of X is the square-root of the variance, provided it exists.

- The variance is often denoted by σ^2 , and the standard deviation σ .
- Because $(X - E(X))^2 \geq 0$, $\text{Var}(X) \geq 0$.
- Formally, the variance is the mean of the squared distance between X and $E[X]$. If most values of X are close to the mean, this value is small; and vice-versa if most values of X are far away from $E[X]$.
- By this definition, the units for the variance are squared units.

Variance V

- That is, if X is measured in meters, then the variance is measured in square-meters, and the standard deviation is measured in meters.

Theorem 4.4: linear transformation of a single variable

Let X be a random variable, and assume that $\text{Var}(X)$ exists. Then if $Y = a + bX$, then $\text{Var}(Y) = b^2\text{Var}(X)$.

Proof.

Variance VII

- This result makes a lot of sense: adding a constant only “shifts” a distribution, it does not affect the spread.
- The multiplier does change the spread, and because we’re squaring the difference, the multiplier is also squared.
- From this result, we can also see that the standard deviation also changes in a natural way.
- Specifically, if σ_Y, σ_X denote the standard deviations of X and Y , respectively, then

$$\sigma_Y = |b|\sigma_X.$$

- We take the absolute value, because variance and standard deviation are always positive, though the multiplier b might be negative.

Variance VIII

Example: Bernoulli distribution

Let X be a Bernoulli(p) distributed random variable. What is the variance of X ?

Example: Normal distribution

Let $X \sim N(\mu, \sigma^2)$. What is $\text{Var}(X)$?

- Using the definition of variance, we will derive a very famous inequality.

Theorem 4.5: Chebyshev's Inequality

Let X be a random variable with $E[X] = \mu$, and $\text{Var}(X) = \sigma^2$. Then for any $t > 0$,

$$P(|X - \mu| > t) \leq \frac{\sigma^2}{t^2}.$$

Variance XI

- This theorem bounds the probability that the difference between X and $E[X]$ is larger than t .
- If σ^2 is small, then the probability that X deviates far away from the mean is also small.
- By letting $t = k\sigma$, we get a bound on the probability that a variable will be k -standard deviations away from the mean:

$$P(|X - \mu| > k\sigma) \leq \frac{1}{k^2},$$

- For instance, the probability that any arbitrary random variable X will be more than 4σ away from $E[X]$ is less than $1/16$.

Variance XII

- While applicable to all random variables with well-defined variances, it is not the most optimal bound we can achieve.
- For instance, if $X \sim N(\mu, \sigma^2)$, then
$$P(|X - \mu| > 1.96 \times \sigma) = 0.05 < 1/4$$

Corollary: zero variance

Let X be a random variable with $\text{Var}(X) = 0$. Then
$$P(X = \mu) = 1.$$

Theorem 4.6: Variance Calculation

Let X be a random variable such that $\text{Var}(X)$ exists. Then

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

where $\mu = E(X)$.

Variance XIV

- Theorem 4.6 is sometimes useful to help us calculate the variance of a random variable.
- Other times, the variance is known, and the theorem helps us calculate $E(X^2)$.

Example: Uniform distribution

Let $X \sim U(0, 1)$. Use Theorem 4.6 to find $\text{Var}(X)$.

Measurement Error

- Often, values of interest cannot be known precisely, but instead must be determined by experimental procedures.
- For instance: measurements of weight, length, voltage, or intervals of time can be complex, and generally involve potential sources of error.
- The National Institute of Standards and Technology (NIST) in the US are charged with developing and maintaining measurement standards.
- Statisticians have historically been employed by these organizations to help with this endeavor.

Measurement Error II

- Typically, there are two main types of measurement error: **random** vs **systematic**.
- For instance, a sequence of repeated independent measurements made from the same instrument or experimental procedure may not give the same value each time. These uncontrollable differences are modeled as **random** error.
- However, there may be a **systematic** error that affects all measurements, such as poorly calibrated instruments, or errors that are associated with the method of measurement.

Measurement Error III

- Suppose that the true value of a quantity being measured is x_0 . We have a random measurement X , which is modeled as

$$X = x_0 + \beta + \epsilon.$$

- Here, β is the systematic error, and ϵ is the random component of the error.

Measurement Error IV

Definition: Bias

Let x_0 be the true value of a measurement, modeled as a random variable X such that

$$X = x_0 + \beta + \epsilon,$$

where $E(\epsilon) = 0$, $\text{Var}(\epsilon) = \sigma^2$. Then, we have

$$E[X] = x_0 + \beta.$$

The value $\beta = E(X - x_0)$ is called the **bias** of the random variable, and we say that X is an unbiased estimate of x_0 if $\beta = 0$.

Measurement Error V

- The two factors that impact the quality of our estimator is the bias β and the variance σ^2 .
- If both $\beta = 0$ and $\sigma^2 = 0$, then we get a perfect measurement.
- Ideally, we want an estimator that minimizes the bias and the variance, though as we will see (Math 4451) there is a principle known as the **bias-variance** trade-off, which suggests that efforts to minimize bias often result in larger variance (and vice-versa).
- Many approaches in statistics we will cover next semester aim at finding estimators that are unbiased ($\beta = 0$), while having minimum variance as possible (that is, the minimum-variance unbiased estimator (MVUE)).

Theorem 4.7: Mean Squared Error

Let X be a random variable representing a random estimate for value x_0 . The mean-squared error of the estimator X is defined as $\text{MSE}(X) = E[(X - x_0)^2]$. If β is the bias of the estimator and σ^2 the variance, then

$$\text{MSE}(X) = \beta^2 + \sigma^2.$$

Covariance and Correlation

Covariance

- The variance of a random variable is a measure of its variability.
- The *covariance* of two random variables is a measure of their joint-variability.
- It's also used to measure how closely associated two random variables are.

Definition: Covariance

If X and Y are jointly distributed random variables with expectations μ_X and μ_Y , the covariance of X, Y is:

$$\text{Cov}(X, Y) = E[(X - \mu_x)(Y - \mu_y)]$$

Covariance II

- The covariance is the average value of the product of the deviation of X from its mean, and Y from its mean.
- If X and Y are positively associated, we expect that if a value of X is larger than its mean, then the value of Y is also larger than its mean.
- In this case, the covariance is positive.
- Example: Suppose X is a random variable representing height of an adult male, and Y is the weight. In this case, we expect heights larger than average will also have weights larger than average, so the covariance is positive.

Covariance III

Calculating Covariance

Let X and Y be random variables. Then

$$\text{Cov}(X, Y) = E[XY] - E[X]E[Y].$$

Proof.

Covariance IV

- One important example is when X and Y are independent:
- In this case, we have shown that $E[XY] = E[X]E[Y]$.
- Therefore, $\text{Cov}(X, Y) = E[X]E[Y] - E[X]E[Y] = 0$.
- **however**, the inverse is not true: Just because $\text{Cov}(X, Y) = 0$ does *not* imply X and Y are independent.

Example: Calculating Covariance

Let (X, Y) be jointly defined random variables is joint pdf $f(x, y) = 2x + 2y - 4xy$, for all $0 \leq x, y \leq 1$. Calculate the covariance $\text{Cov}(X, Y)$.

Solution:

Solution cont...

Covariance Properties

- Covariance has several useful properties that can help with calculations.
- One of them is that the covariance is **bilinear** operator.

Theorem: Bilinear Covariance

Let $X_i, i = 1, 2, \dots, n$ and $Y_j, j = 1, 2, \dots, m$ be a collection of random variables, and a, c, b_i, d_j be real numbers for all i and j .

Then:

$$\text{Cov}\left(a + \sum_{i=1}^n X_i, c + \sum_{j=1}^m Y_j\right) = \sum_{i=1}^n \sum_{j=1}^m b_i d_j \text{Cov}(X_i, Y_j).$$

In particular,

$$\text{Cov}(aX + bW, cY + dZ) = ac\text{Cov}(W, Y) + bc\text{Cov}(X, Y) + ad\text{Cov}(X, Z) + bd\text{Cov}(W, Z)$$

References and Acknowledgements

Resnick S (2019). *A probability path*. Springer.

Rice JA (2007). *Mathematical statistics and data analysis*, volume 371. 3 edition. Thomson/Brooks/Cole Belmont, CA.

- Compiled on October 17, 2025 using R version 4.5.1.
- Licensed under the [Creative Commons Attribution-NonCommercial](#) license. Please share and remix non-commercially, mentioning its origin.
- We acknowledge [students and instructors for previous versions of this course / slides](#).

