

Chapter 12

Linear Algebra in Probability & Statistics

12.1 Mean, Variance, and Probability

We are starting with the three fundamental words of this chapter: *mean, variance, and probability*. Let me give a rough explanation of their meaning before I write any formulas:

The **mean** is the *average value* or expected value

The **variance** σ^2 measures the average *squared distance* from the mean m

The **probabilities** of n different outcomes are positive numbers p_1, \dots, p_n adding to 1.

Certainly the mean is easy to understand. We will start there. But right away we have two different situations that you have to keep straight. On the one hand, we may have the results (*sample values*) from a completed trial. On the other hand, we may have the expected results (*expected values*) from future trials. Let me give examples:

Sample values Five random freshmen have ages **18, 17, 18, 19, 17**

Sample mean $\frac{1}{5}(18 + 17 + 18 + 19 + 17) = \mathbf{17.8}$

Probabilities The ages in a freshmen class are 17 (**20%**), 18 (**50%**), 19 (**30%**)

A random freshman has **expected age** $\mathbf{E}[x] = (0.2)17 + (0.5)18 + (0.3)19 = \mathbf{18.1}$

Both numbers 17.8 and 18.1 are correct averages. The sample mean starts with N samples x_1, \dots, x_N from a completed trial. Their mean is the *average* of the N observed samples:

$$\text{Sample mean} \quad m = \mu = \frac{1}{N}(x_1 + x_2 + \dots + x_N) \quad (1)$$

The **expected value of x** starts with the probabilities p_1, \dots, p_n of the ages x_1, \dots, x_n :

$$\text{Expected value } m = \mathbf{E}[x] = p_1x_1 + p_2x_2 + \dots + p_nx_n. \quad (2)$$

This is $p \cdot x$. Notice that $m = \mathbf{E}[x]$ tells us what to expect, $m = \mu$ tells us what we got.

By taking many samples (large N), the sample results will come close to the probabilities. The “Law of Large Numbers” says that with probability 1, the sample mean will converge to its expected value $\mathbf{E}[x]$ as the sample size N increases. A fair coin has probability $p_0 = \frac{1}{2}$ of tails and $p_1 = \frac{1}{2}$ of heads. Then $\mathbf{E}[x] = (\frac{1}{2})0 + \frac{1}{2}(1)$. The fraction of heads in N flips of the coin is the sample mean, expected to approach $\mathbf{E}[x] = \frac{1}{2}$.

This does *not* mean that if we have seen more tails than heads, the next sample is likely to be heads. The odds remain 50-50. The first 100 or 1000 flips do affect the sample mean. *But 1000 flips will not affect its limit*—because you are dividing by $N \rightarrow \infty$.

Variance (around the mean)

The **variance σ^2** measures expected distance (squared) from the expected mean $\mathbf{E}[x]$. The **sample variance S^2** measures actual distance (squared) from the sample mean. The square root is the **standard deviation σ or S** . After an exam, I email μ and S to the class. I don’t know the expected mean and variance because I don’t know the probabilities p_1 to p_{100} for each score. (After teaching for 50 years, I still have no idea what to expect.)

The deviation is always deviation *from the mean*—sample or expected. We are looking for the size of the “spread” around the mean value $x = m$. Start with N samples.

$$\text{Sample variance } S^2 = \frac{1}{N-1} \left[(x_1 - m)^2 + \dots + (x_N - m)^2 \right] \quad (3)$$

The sample ages $x = 18, 17, 18, 19, 17$ have mean $m = 17.8$. That sample has variance 0.7:

$$S^2 = \frac{1}{4} \left[(.2)^2 + (-.8)^2 + (.2)^2 + (1.2)^2 + (-.8)^2 \right] = \frac{1}{4}(2.8) = \mathbf{0.7}$$

The minus signs disappear when we compute squares. Please notice! Statisticians divide by $N - 1 = 4$ (and not $N = 5$) so that S^2 is an unbiased estimate of σ^2 . One degree of freedom is already accounted for in the sample mean.

An important identity comes from splitting each $(x - m)^2$ into $x^2 - 2mx + m^2$:

$$\begin{aligned} \text{sum of } (x_i - m)^2 &= (\text{sum of } x_i^2) - 2m(\text{sum of } x_i) + (\text{sum of } m^2) \\ &= (\text{sum of } x_i^2) - 2m(Nm) + Nm^2 \\ \text{sum of } (x_i - m)^2 &= (\text{sum of } x_i^2) - Nm^2. \end{aligned} \quad (4)$$

This is an equivalent way to find $(x_1 - m)^2 + \dots + (x_N - m)^2$ by adding $x_1^2 + \dots + x_N^2$.

Now start with probabilities p_i (never negative!) instead of samples. We find expected values instead of sample values. The variance σ^2 is the crucial number in statistics.

$$\text{Variance } \sigma^2 = E[(x - m)^2] = p_1(x_1 - m)^2 + \cdots + p_n(x_n - m)^2. \quad (5)$$

We are squaring the distance from the expected value $m = E[x]$. We don't have samples, only expectations. We know probabilities but we don't know experimental outcomes.

Example 1 Find the variance σ^2 of the ages of college freshmen.

Solution The probabilities of ages $x_i = 17, 18, 19$ were $p_i = 0.2$ and 0.5 and 0.3 . The expected value was $m = \sum p_i x_i = 18.1$. The variance uses those same probabilities:

$$\begin{aligned} \sigma^2 &= (0.2)(17 - 18.1)^2 + (0.5)(18 - 18.1)^2 + (0.3)(19 - 18.1)^2 \\ &= (0.2)(1.21) + (0.5)(0.01) + (0.3)(0.81) = 0.49. \end{aligned}$$

The **standard deviation** is the square root $\sigma = 0.7$.

This measures the spread of 17, 18, 19 around $E[x]$, weighted by probabilities .2, .5, .3.

Continuous Probability Distributions

Up to now we have allowed for n possible outcomes x_1, \dots, x_n . With ages 17, 18, 19, we only had $n = 3$. If we measure age in days instead of years, there will be a thousand possible ages (too many). Better to allow *every number between 17 and 20*—a continuum of possible ages. Then the probabilities p_1, p_2, p_3 for ages x_1, x_2, x_3 have to move to a **probability distribution** $p(x)$ for a whole continuous range of ages $17 \leq x \leq 20$.

The best way to explain probability distributions is to give you two examples. They will be the **uniform distribution** and the **normal distribution**. The first (uniform) is easy. The normal distribution is all-important.

Uniform distribution Suppose ages are uniformly distributed between 17.0 and 20.0. All ages between those numbers are “equally likely”. Of course any one exact age has no chance at all. There is zero probability that you will hit the exact number $x = 17.1$ or $x = 17 + \sqrt{2}$. What you can truthfully provide (assuming our uniform distribution) is **the chance** $F(x)$ **that a random freshman has age less than** x :

$$\begin{aligned} \text{The chance of age less than } x = 17 \text{ is } F(17) &= 0 && x \leq 17 \text{ won't happen} \\ \text{The chance of age less than } x = 20 \text{ is } F(20) &= 1 && x \leq 20 \text{ will happen} \\ \text{The chance of age less than } x \text{ is } F(x) &= \frac{1}{3}(x - 17) && \mathbf{F \text{ goes from 0 to 1}} \end{aligned}$$

That formula $F(x) = \frac{1}{3}(x - 17)$ gives $F = 0$ at $x = 17$; then $x \leq 17$ won't happen. It gives $F(x) = 1$ at $x = 20$; then $x \leq 20$ is sure. Between 17 and 20, the graph of the **cumulative distribution** $F(x)$ increases linearly for this uniform model.

Let me draw the graphs of $F(x)$ and its derivative $p(x)$ = “probability density function”.

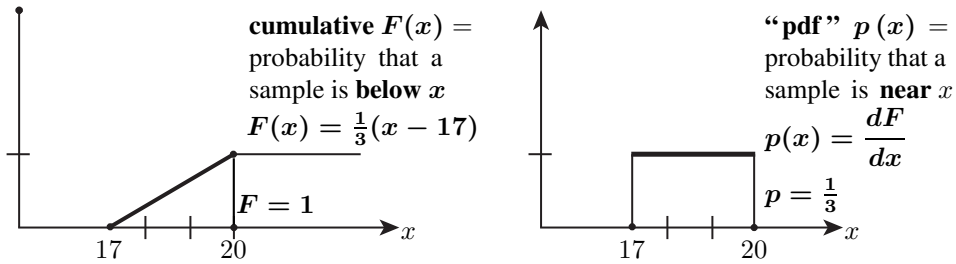


Figure 12.1: $F(x)$ is the cumulative distribution and its derivative $p(x) = dF/dx$ is the **probability density function (pdf)**. For this uniform distribution, $p(x)$ is constant between 17 and 20. The total area under the graph of $p(x)$ is the total probability $F = 1$.

You could say that $p(x) dx$ is the probability of a sample falling in between x and $x + dx$. This is “infinitesimally true”: $p(x) dx$ is $F(x + dx) - F(x)$. Here is the full truth:

$$F = \text{integral of } p \quad \text{Probability of } a \leq x \leq b = \int_a^b p(x) dx = F(b) - F(a) \quad (6)$$

$F(b)$ is the probability of $x \leq b$. I subtract $F(a)$ to keep $x \geq a$. That leaves $a \leq x \leq b$.

Mean and Variance of $p(x)$

What are the mean m and variance σ^2 for a probability distribution? Previously we added $p_i x_i$ to get the mean (expected value). With a continuous distribution we **integrate $x p(x)$** :

Mean $m = E[x] = \int_{x=17}^{20} x p(x) dx = \int_{x=17}^{20} (x) \left(\frac{1}{3}\right) dx = 18.5$

For this uniform distribution, the mean m is halfway between 17 and 20. Then the probability of a random value x below this halfway point $m = 18.5$ is $F(m) = \frac{1}{2}$.

In MATLAB, $x = \text{rand}(1)$ chooses a random number uniformly between 0 and 1. Then the expected mean is $m = \frac{1}{2}$. The interval from 0 to x has probability $F(x) = x$. The interval below the mean m always has probability $F(m) = \frac{1}{2}$.

The variance is the average squared distance to the mean. With N outcomes, σ^2 is the sum of $p_i(x_i - m)^2$. For a continuous random variable x , the sum changes to an **integral**.

Variance $\sigma^2 = E[(x - m)^2] = \int p(x) (x - m)^2 dx \quad (7)$

When ages are uniform between $17 \leq x \leq 20$, the integral can shift to $0 \leq x \leq 3$:

$$\sigma^2 = \int_{17}^{20} \frac{1}{3}(x - 18.5)^2 dx = \int_0^3 \frac{1}{3}(x - 1.5)^2 dx = \frac{1}{9}(x - 1.5)^3 \Big|_{x=0}^{x=3} = \frac{2}{9}(1.5)^3 = \frac{3}{4}.$$

That is a typical example, and here is the complete picture for a uniform $p(x)$, 0 to a .

Uniform distribution for $0 \leq x \leq a$ **Density** $p(x) = \frac{1}{a}$ **Cumulative** $F(x) = \frac{x}{a}$

Mean $m = \frac{a}{2}$ halfway **Variance** $\sigma^2 = \int_0^a \frac{1}{a} \left(x - \frac{a}{2}\right)^2 dx = \frac{a^2}{12}$ (8)

The mean is a multiple of a , the variance is a multiple of a^2 . For $a = 3$, $\sigma^2 = \frac{9}{12} = \frac{3}{4}$. For one random number between 0 and 1 (mean $\frac{1}{2}$) the variance is $\sigma^2 = \frac{1}{12}$.

Normal Distribution : Bell-shaped Curve

The normal distribution is also called the ‘‘Gaussian’’ distribution. It is the most important of all probability density functions $p(x)$. The reason for its overwhelming importance comes from repeating an experiment and averaging the outcomes. The experiments have their own distribution (like heads and tails). *The average approaches a normal distribution.*

Central Limit Theorem (informal) The average of N samples of ‘‘any’’ probability distribution approaches a normal distribution as $N \rightarrow \infty$.

Start with the ‘‘standard normal distribution’’. It is symmetric around $x = 0$, so its mean value is $m = 0$. It is chosen to have a standard variance $\sigma^2 = 1$. It is called $\mathbf{N}(0, 1)$.

$$\boxed{\text{Standard normal distribution} \quad p(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}.} \quad (9)$$

The graph of $p(x)$ is the **bell-shaped curve** in Figure 12.2. The standard facts are

$$\begin{aligned} \text{Total probability} &= 1 & \int_{-\infty}^{\infty} p(x) dx &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-x^2/2} dx = 1 \\ \text{Mean } \mathbf{E}[x] &= 0 & m &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x e^{-x^2/2} dx = 0 \\ \text{Variance } \mathbf{E}[x^2] &= 1 & \sigma^2 &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (x - 0)^2 e^{-x^2/2} dx = 1 \end{aligned}$$

The zero mean was easy because we are integrating an odd function. Changing x to $-x$ shows that “integral = $-$ integral”. So that integral must be $m = 0$.

The other two integrals apply the idea in Problem 12 to reach 1. Figure 12.2 shows a graph of $p(x)$ for the normal distribution $\mathbf{N}(0, \sigma)$ and also its cumulative distribution $F(x) = \text{integral of } p(x)$. From the symmetry of $p(x)$ you see *mean = zero*. From $F(x)$ you see a very important practical approximation for opinion polling :

The probability that a random sample falls between $-\sigma$ and σ is $F(\sigma) - F(-\sigma) \approx \frac{2}{3}$.

This is because $\int_{-\sigma}^{\sigma} p(x) dx$ equals $\int_{-\infty}^{\sigma} p(x) dx - \int_{-\infty}^{-\sigma} p(x) dx = F(\sigma) - F(-\sigma)$.

Similarly, the probability that a random x lies between -2σ and 2σ (“less than two standard deviations from the mean”) is $F(2\sigma) - F(-2\sigma) \approx 0.95$. If you have an experimental result further than 2σ from the mean, it is fairly sure to be not accidental: chance = 0.05. Drug tests may look for a tighter confirmation, like probability 0.001. Searching for the Higgs boson used a hyper-strict test of 5σ deviation from pure accident.

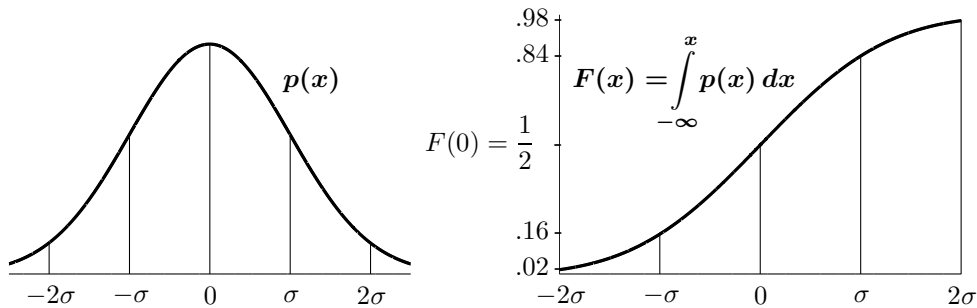


Figure 12.2: The standard normal distribution $p(x)$ has mean $m = 0$ and $\sigma = 1$.

The normal distribution with any mean m and standard deviation σ comes by shifting and stretching the standard $\mathbf{N}(0, 1)$. **Shift x to $x - m$.** **Stretch $x - m$ to $(x - m)/\sigma$.**

Gaussian density $p(x)$

Normal distribution $\mathbf{N}(m, \sigma)$

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-m)^2/2\sigma^2} \quad (10)$$

The integral of $p(x)$ is $F(x)$ —the probability that a random sample will fall below x . The differential $p(x) dx = F(x + dx) - F(x)$ is the probability that a random sample will fall between x and $x + dx$. There is no simple formula to integrate $e^{-x^2/2}$, so this cumulative distribution $F(x)$ is computed and tabulated very carefully.

N Coin Flips and $N \rightarrow \infty$

Example 2 Suppose x is 1 or -1 with equal probabilities $p_1 = p_{-1} = \frac{1}{2}$.

The mean value is $m = \frac{1}{2}(1) + \frac{1}{2}(-1) = 0$. The variance is $\sigma^2 = \frac{1}{2}(1)^2 + \frac{1}{2}(-1)^2 = 1$.

The key question is the *average* $A_N = (x_1 + \cdots + x_N)/N$. The independent x_i are ± 1 and we are dividing their sum by N . The expected mean of A_N is still zero. The law of large numbers says that this sample average approaches zero with probability 1. How fast does A_N approach zero? **What is its variance σ_N^2 ?**

$$\text{By linearity } \sigma_N^2 = \frac{\sigma^2}{N^2} + \frac{\sigma^2}{N^2} + \cdots + \frac{\sigma^2}{N^2} = N \frac{\sigma^2}{N^2} = \frac{1}{N} \text{ since } \sigma^2 = 1. \quad (11)$$

Example 3 Change outputs from 1 or -1 to $x = 1$ or $x = 0$. Keep $p_1 = p_0 = \frac{1}{2}$.

The new mean value $m = \frac{1}{2}$ falls halfway between 0 and 1. The variance moves to $\sigma^2 = \frac{1}{4}$:

$$m = \frac{1}{2}(1) + \frac{1}{2}(0) = \frac{1}{2} \quad \text{and} \quad \sigma^2 = \frac{1}{2} \left(1 - \frac{1}{2}\right)^2 + \frac{1}{2} \left(0 - \frac{1}{2}\right)^2 = \frac{1}{4}.$$

The average A_N now has mean $\frac{1}{2}$ and variance $\frac{1}{4N^2} + \cdots + \frac{1}{4N^2} = \frac{1}{4N} = \sigma_N^2$. (12)

This σ_N is half the size of σ_N in Example 2. This must be correct because the new range 0 to 1 is half as long as -1 to 1. Examples 2-3 are showing a law of linearity.

The new 0 – 1 variable x_{new} is $\frac{1}{2} x_{\text{old}} + \frac{1}{2}$. So the mean m is increased to $\frac{1}{2}$ and the variance is *multiplied* by $(\frac{1}{2})^2$. A shift changes m and the rescaling changes σ^2 .

$$\text{Linearity } x_{\text{new}} = ax_{\text{old}} + b \text{ has } m_{\text{new}} = am_{\text{old}} + b \text{ and } \sigma_{\text{new}}^2 = a^2 \sigma_{\text{old}}^2. \quad (13)$$

Here are the results from three numerical tests: random 0 or 1 averaged over N trials.

[48 1's from $N = 100$] [5035 1's from $N = 10000$] [19967 1's from $N = 40000$].

The standardized $X = (x - m)/\sigma = (A_N - \frac{1}{2}) / 2\sqrt{N}$ was [-.40] [.70] [-.33].

The Central Limit Theorem says that the average of many coin flips will approach a normal distribution. Let us begin to see how that happens: **binomial approaches normal**.

For each flip, the probability of heads is $\frac{1}{2}$. For $N = 3$ flips, the probability of heads all three times is $(\frac{1}{2})^3 = \frac{1}{8}$. The probability of heads twice and tails once is $\frac{3}{8}$, from three sequences HHT and HTH and THH. These numbers $\frac{1}{8}$ and $\frac{3}{8}$ are pieces of $(\frac{1}{2} + \frac{1}{2})^3 = \frac{1}{8} + \frac{3}{8} + \frac{3}{8} + \frac{1}{8} = 1$. *The average number of heads in 3 flips is 1.5.*

$$\text{Mean } m = (3 \text{ heads})\frac{1}{8} + (2 \text{ heads})\frac{3}{8} + (1 \text{ head})\frac{3}{8} + 0 = \frac{3}{8} + \frac{6}{8} + \frac{3}{8} = 1.5 \text{ heads}$$

With N flips, Example 3 (or common sense) gives a mean of $m = \sum x_i p_i = \frac{1}{2}N$ heads.

The variance σ^2 is based on the *squared distance* from this mean $N/2$. With $N = 3$ the variance is $\sigma^2 = \frac{3}{4}$ (which is $N/4$). To find σ^2 we add $(x_i - m)^2 p_i$ with $m = 1.5$:

$$\sigma^2 = (3 - 1.5)^2 \frac{1}{8} + (2 - 1.5)^2 \frac{3}{8} + (1 - 1.5)^2 \frac{3}{8} + (0 - 1.5)^2 \frac{1}{8} = \frac{9 + 3 + 3 + 9}{32} = \frac{3}{4}.$$

For any N , the variance is $\sigma_N^2 = N/4$. Then $\sigma_N = \sqrt{N}/2$.

Figure 12.3 shows how the probabilities of 0, 1, 2, 3, 4 heads in $N = 4$ flips come close to a bell-shaped Gaussian. That Gaussian is centered at the mean value $N/2 = 2$. To reach the standard Gaussian (mean 0 and variance 1) we shift and rescale that graph. If x is the number of heads in N flips—the average of N zero-one outcomes—then x is shifted by its mean $m = N/2$ and rescaled by $\sigma = \sqrt{N}/2$ to produce the standard X :

Shifted and scaled

$$X = \frac{x - m}{\sigma} = \frac{x - \frac{1}{2}N}{\sqrt{N}/2} \quad (N = 4 \text{ has } X = x - 2)$$

Subtracting m is “centering” or “detrrending”. The mean of X is zero.

Dividing by σ is “normalizing” or “standardizing”. The variance of X is 1.

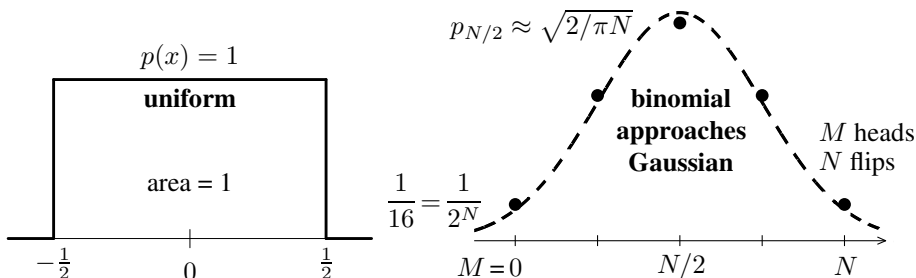


Figure 12.3: The probabilities $p = (1, 4, 6, 4, 1)/16$ for the number of heads in 4 flips. These p_i approach a Gaussian distribution with variance $\sigma^2 = N/4$ centered at $m = N/2$. For X , the Central Limit Theorem gives convergence to the normal distribution $\mathbf{N}(0, 1)$.

It is fun to see the Central Limit Theorem giving the right answer at the center point $X = 0$. At that point, the factor $e^{-X^2/2}$ equals 1. We know that the variance for N coin flips is $\sigma^2 = N/4$. The center of the bell-shaped curve has height $1/\sqrt{2\pi\sigma^2} = \sqrt{2/N\pi}$.

What is the height at the center of the coin-flip distribution p_0 to p_N (the binomial distribution)? For $N = 4$, the probabilities for 0, 1, 2, 3, 4 heads come from $(\frac{1}{2} + \frac{1}{2})^4$.

$$\text{Center probability } \frac{6}{16} \quad \left(\frac{1}{2} + \frac{1}{2}\right)^4 = \frac{1}{16} + \frac{4}{16} + \frac{6}{16} + \frac{4}{16} + \frac{1}{16} = 1.$$

The binomial theorem in Problem 8 tells us the center probability $p_{N/2}$ for any even N :

$$\text{The center probability } \left(\frac{N}{2} \text{ heads, } \frac{N}{2} \text{ tails} \right) \text{ is } \frac{1}{2^N} \frac{N!}{(N/2)!(N/2)!}$$

For $N = 4$, those factorials produce $4!/2!2! = 24/4 = 6$. For large N , Stirling's formula $\sqrt{2\pi N}(N/e)^N$ is a close approximation to $N!$. Use Stirling for N and twice for $N/2$:

$$\begin{array}{l} \text{Limit of coin-flip} \\ \text{Center probability} \end{array} \quad p_{N/2} \approx \frac{1}{2^N} \frac{\sqrt{2\pi N}(N/e)^N}{\pi N(N/2e)^N} = \frac{\sqrt{2}}{\sqrt{\pi N}} = \frac{1}{\sqrt{2\pi\sigma}}. \quad (14)$$

At that last step we used the variance $\sigma^2 = N/4$ for the coin-tossing problem. The result $1/\sqrt{2\pi\sigma}$ matches the center value (above) for the Gaussian. The Central Limit Theorem is true: The “binomial distribution” approaches the normal distribution as $N \rightarrow \infty$.

Monte Carlo Estimation Methods

Scientific computing has to work with errors in the data. Financial computing has to work with unsure numbers and uncertain predictions. All of applied mathematics has moved to **accepting uncertainty in the inputs and estimating the variance in the outputs**.

How to estimate that variance? Often probability distributions $p(x)$ are not known. What we can do is to try different inputs b and compute the outputs x and take an average. This is the simplest form of a **Monte Carlo method** (named after the gambling palace on the Riviera, where I once saw a fight about whether the bet was placed in time). Monte Carlo approximates an expected value $E[x]$ by a sample average $(x_1 + \dots + x_N)/N$.

Please understand that every x_k can be expensive to compute. We are not just flipping coins. Each sample comes from a set of data b_k . *Monte Carlo randomly chooses this data b_k , it computes the outputs x_k , and then it averages those x 's.* Decent accuracy for $E[x]$ often requires many samples b and huge computing cost. The error in approximating $E[x]$ by $(x_1 + \dots + x_N)/N$ is normally of order $1/\sqrt{N}$. *Slow improvement as N increases.*

That $1/\sqrt{N}$ estimate came for coin flips in equation (11). Averaging N independent samples x_k of variance σ^2 reduces the variance to σ^2/N .

“Quasi-Monte Carlo” can sometimes reduce this variance to σ^2/N^2 : a big difference! The inputs b_k are selected very carefully—not just randomly. This QMC approach is surveyed in the journal *Acta Numerica* 2013. The newer idea of “Multilevel Monte Carlo” is outlined by Michael Giles in *Acta Numerica* 2015. Here is how it works.

Suppose it is much simpler to simulate another variable $y(b)$ close to $x(b)$. Then use N computations of $y(b_k)$ and only $N^* < N$ computations of $x(b_k)$ to estimate $E[x]$.

2-level Monte Carlo

$$E[x] \approx \frac{1}{N} \sum_1^N y(b_k) + \frac{1}{N^*} \sum_1^{N^*} [x(b_k) - y(b_k)].$$

The idea is that $x - y$ has a smaller variance σ^* than the original x . Therefore N^* can be smaller than N , with the same accuracy for $E[x]$. We do N cheap simulations to find the y 's. Those cost C each. We only do N^* expensive simulations involving x 's. Those cost C^* each. The total computing cost is $NC + N^*C^*$.

Calculus minimizes the overall variance for a fixed total cost. The optimal ratio N^*/N is $\sqrt{C/C^*} \sigma^*/\sigma$. Three-level Monte Carlo would simulate x , y , and z :

$$E[x] \approx \frac{1}{N} \sum_1^N z(\mathbf{b}_k) + \frac{1}{N^*} \sum_1^{N^*} [\mathbf{y}(\mathbf{b}_k) - z(\mathbf{b}_k)] + \frac{1}{N^{**}} \sum_1^{N^{**}} [x(\mathbf{b}_k) - \mathbf{y}(\mathbf{b}_k)].$$

Giles optimizes N, N^*, N^{**}, \dots to keep $E[x] \leq \text{fixed } E_0$, and provides a MATLAB code.

Review : Three Formulas for the Mean and the Variance

The formulas for m and σ^2 are the starting point for all of probability and statistics. There are three different cases to keep straight: **sample** values X_i , **expected** values (discrete p_i), and a range of **expected** values (continuous $p(x)$). Here are the mean and the variance:

Samples X_1 to X_N	$m = \frac{X_1 + \dots + X_N}{N}$	$S^2 = \frac{(X_1 - m)^2 + \dots + (X_N - m)^2}{N - 1}$
n possible outputs with probabilities p_i	$m = \sum_1^n p_i x_i$	$\sigma^2 = \sum_1^n p_i (x_i - m)^2$
Range of outputs with probability density	$m = \int x p(x) dx$	$\sigma^2 = \int (x - m)^2 p(x) dx$

A natural question: Why are there no probabilities p on the first line? How can these formulas be parallel? Answer: *We expect a fraction p_i of the samples to be $X = x_i$.* If this is exactly true, $X = x_i$ is repeated $p_i N$ times. Then lines 1 and 2 give the same m .

When we work with samples, we don't know the p_i . We just include each output X as often as it comes. We get the "empirical" mean instead of the expected mean.

Problem Set 12.1

- 1 Add 7 to every output x . What happens to the mean and the variance? What are the new sample mean, the new expected mean, and the new variance?
- 2 We know: $\frac{1}{3}$ of all integers are divisible by 3 and $\frac{1}{7}$ of integers are divisible by 7. What fraction of integers will be divisible by 3 or 7 or both?
- 3 Suppose you sample from the numbers 1 to 1000 with equal probabilities $1/1000$. What are the probabilities p_0 to p_9 that the last digit of your sample is $0, \dots, 9$? What is the expected mean m of that last digit? What is its variance σ^2 ?
- 4 Sample again from 1 to 1000 but look at the last digit of the sample *squared*. That square could end with $x = 0, 1, 4, 5, 6, \text{ or } 9$. What are the probabilities $p_0, p_1, p_4, p_5, p_6, p_9$? What are the (expected) mean m and variance σ^2 of that number x ?

- 5 (a little tricky) Sample again from 1 to 1000 with equal probabilities and let x be the first digit ($x = 1$ if the number is 15). What are the probabilities p_1 to p_9 (adding to 1) of $x = 1, \dots, 9$? What are the mean and variance of x ?
- 6 Suppose you have $N = 4$ samples 157, 312, 696, 602 in Problem 5. What are the first digits x_1 to x_4 of the squares? What is the sample mean μ ? What is the sample variance S^2 ? Remember to divide by $N - 1 = 3$ and not $N = 4$.

- 7 Equation (4) gave a second equivalent form for S^2 (the variance using samples):

$$S^2 = \frac{1}{N-1} \text{sum of } (x_i - m)^2 = \frac{1}{N-1} [(\text{sum of } x_i^2) - Nm^2].$$

Verify the matching identity for the expected variance σ^2 (using $m = \sum p_i x_i$):

$$\sigma^2 = \text{sum of } p_i (x_i - m)^2 = (\text{sum of } p_i x_i^2) - m^2.$$

- 8 If all 24 samples from a population produce the same age $x = 20$, what are the sample mean μ and the sample variance S^2 ? What if $x = 20$ or 21, 12 times each?
- 9 Computer experiment as on page 541: Find the average $A_{1000000}$ of a million random 0-1 samples! What is $X = (A_N - \frac{1}{2}) / 2\sqrt{N}$?

- 10 The probability p_i to get i heads in N coin flips is the binomial number $b_i = \binom{N}{i}$ divided by 2^N . The b_i add to $(1+1)^N = 2^N$ so the probabilities p_i add to 1.

$$p_0 + \dots + p_N = \left(\frac{1}{2} + \frac{1}{2}\right)^N = \frac{1}{2^N}(b_0 + \dots + b_N) \text{ with } b_i = \frac{N!}{i!(N-i)!}$$

$$N=4 \text{ leads to } b_0 = \frac{24}{24}, b_1 = \frac{24}{(1)(6)} = 4, b_2 = \frac{24}{(2)(2)} = 6, p_i = \frac{1}{16}(1, 4, 6, 4, 1).$$

Notice $b_i = b_{N-i}$. *Problem:* Confirm that the mean $m = 0p_0 + \dots + Np_N$ equals $\frac{N}{2}$.

- 11 For any function $f(x)$ the expected value is $E[f] = \sum p_i f(x_i)$ or $\int p(x) f(x) dx$ (discrete probability or continuous probability). Suppose the mean is $E[x] = m$ and the variance is $E[(x - m)^2] = \sigma^2$. **What is $E[x^2]$?**

- 12 Show that the standard normal distribution $p(x)$ has total probability $\int p(x) dx = 1$ as required. A famous trick multiplies $\int p(x) dx$ by $\int p(y) dy$ and computes the integral over all x and all y ($-\infty$ to ∞). The trick is to replace $dx dy$ in that double integral by $r dr d\theta$ (polar coordinates with $x^2 + y^2 = r^2$). Explain each step:

$$2\pi \int_{-\infty}^{\infty} p(x) dx \int_{-\infty}^{\infty} p(y) dy = \iiint_{-\infty}^{\infty} e^{-(x^2+y^2)/2} dx dy = \int_{\theta=0}^{2\pi} \int_{r=0}^{\infty} e^{-r^2/2} r dr d\theta = 2\pi.$$