

Averages of Random Variables (Chapters 8&9) ①

(as long as they are identically and independently distributed), or, Why Statistics Work

Suppose that X_1, X_2, X_3, \dots of random variables which are independent and have exactly the same distribution. In particular, they all have the same expected value $E(X_i) = \mu$ and the same variance $V(X) = \sigma^2$.

The average of the first n of these

$$\begin{aligned} \bar{X}_n &= \frac{X_1 + X_2 + \dots + X_n}{n} \\ &= \frac{1}{n}X_1 + \frac{1}{n}X_2 + \dots + \frac{1}{n}X_n. \end{aligned}$$

$$\begin{aligned} E(\bar{X}_n) &= E\left(\frac{1}{n}X_1 + \frac{1}{n}X_2 + \dots + \frac{1}{n}X_n\right) \\ &= \frac{1}{n}E(X_1) + \frac{1}{n}E(X_2) + \dots + \frac{1}{n}E(X_n) \\ &= \frac{1}{n}\mu + \frac{1}{n}\mu + \dots + \frac{1}{n}\mu \\ &= \frac{n}{n}\mu = \mu \end{aligned}$$

②

$$\begin{aligned} V(\bar{X}_n) &= V\left(\frac{1}{n}X_1 + \frac{1}{n}X_2 + \dots + \frac{1}{n}X_n\right) && \text{since the } X_i \text{ are independent} \\ &= V\left(\frac{1}{n}X_1\right) + V\left(\frac{1}{n}X_2\right) + \dots + V\left(\frac{1}{n}X_n\right) \\ &= \frac{1}{n^2}V(X_1) + \frac{1}{n^2}V(X_2) + \dots + \frac{1}{n^2}V(X_n) \\ &= \frac{1}{n^2}\sigma^2 + \frac{1}{n^2}\sigma^2 + \dots + \frac{1}{n^2}\sigma^2 \\ &= \frac{n}{n^2}\sigma^2 = \frac{\sigma^2}{n} \end{aligned}$$

Notice that as n increases, the variance of the average decreases, so we expect the average to be closer to the expected value μ .

This idea is theoretically summarized in the Laws of Large Numbers.

The text only gives what others call the Weak Law of Large Numbers.

(Weak) Law of Large Numbers (Ch. 8) ③

Suppose that X_1, X_2, X_3, \dots is a sequence of independent and identically distributed random variables with $E(X_i) = \mu$. Then, for all $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{\overbrace{X_1 + \dots + X_n}^{\bar{X}_n}}{n} - \mu\right| \geq \varepsilon\right) = 0.$$

(Strong) Law of Large Numbers (Not in textbook.)

Suppose X_1, X_2, X_3, \dots is a sequence of independent and identically distributed random variables with $E(X_i) = \mu$. Then

$$P\left(\lim_{n \rightarrow \infty} \frac{\overbrace{X_1 + \dots + X_n}^{\bar{X}_n}}{n} = \mu\right) = 1.$$

This is why statistics should work: the probability that the average of your data points approaches the true expected value as the number of data points increases is 1.

More computationally useful: relate things to the standard normal distribution. (4)

We'll modify our random variable to have the same expected value and variance as the standard normal distribution.

If X is a random variable with expected value μ and variance σ^2 , the normalization of X is $Z = \frac{X - \mu}{\sigma}$. Then

$$\begin{aligned} E(Z) &= E\left(\frac{X - \mu}{\sigma}\right) = \frac{1}{\sigma} E(X - \mu) \\ &= \frac{1}{\sigma} (E(X) - E(\mu)) = \frac{1}{\sigma} (E(X) - E(X)) = 0 \end{aligned}$$

$$\begin{aligned} \& V\left(\frac{X - \mu}{\sigma}\right) &= \frac{1}{\sigma^2} V(X - \mu) = \frac{1}{\sigma^2} V(X) + \frac{1}{\sigma^2} V(-\mu) \\ &= \frac{1}{\sigma^2} V(X) = \frac{1}{\sigma^2} \sigma^2 = 1 \end{aligned}$$

constants have variance 0

just like for the standard normal.

5

Suppose that X_1, X_2, X_3, \dots is a sequence of random variables that are independent and identically distributed, with $E(X_i) = \mu$ and variance $V(X_i) = \sigma^2$,

$$\begin{aligned} \text{Let } S_n^* &= \frac{Z_1 + Z_2 + \dots + Z_n}{\sqrt{n}} \\ &= \frac{\frac{1}{\sigma}(X_1 - \mu) + \frac{1}{\sigma}(X_2 - \mu) + \dots + \frac{1}{\sigma}(X_n - \mu)}{\sqrt{n}} \\ &= \frac{(X_1 - \mu) + (X_2 - \mu) + \dots + (X_n - \mu)}{\sigma\sqrt{n}} \\ &= \frac{(X_1 + X_2 + \dots + X_n) - n\mu}{\sigma\sqrt{n}} \end{aligned}$$

$$\begin{aligned} E(S_n^*) &= E\left(\frac{X_1 + \dots + X_n - n\mu}{\sigma\sqrt{n}}\right) \\ &= \frac{1}{\sigma\sqrt{n}} E(X_1 + \dots + X_n - n\mu) \\ &= \frac{1}{\sigma\sqrt{n}} \left(\underbrace{E(X_1) + \dots + E(X_n)}_{n\mu} - n\mu \right) \\ &= \frac{1}{\sigma\sqrt{n}} \cdot 0 = 0 \end{aligned}$$

6

$$\begin{aligned}
V(S_n^*) &= V\left(\frac{X_1 + \dots + X_n - n\mu}{\sigma\sqrt{n}}\right) \\
&= \frac{1}{(\sigma\sqrt{n})^2} V(X_1 + \dots + X_n - n\mu) \\
&= \frac{1}{\sigma^2 n} \left[V(X_1) + \dots + V(X_n) + \cancel{V(-n\mu)} \right] \quad \begin{array}{l} \text{since} \\ \text{constants} \\ \text{have no} \\ \text{variance} \end{array} \\
&= \frac{1}{\sigma^2 n} \cdot [n\sigma^2] = \frac{\cancel{n}\sigma^2}{\sigma^2\cancel{n}} = 1
\end{aligned}$$

This brings us to the Central Limit Theorem

If X_1, X_2, \dots is a sequence of independent and identically distributed random variables with $E(X_i) = \mu$ and $V(X_i) = \sigma^2$, then

$$\lim_{n \rightarrow \infty} P(a \leq S_n^* \leq b) = \lim_{n \rightarrow \infty} P\left(a \leq \frac{X_1 + \dots + X_n - n\mu}{\sigma\sqrt{n}} \leq b\right)$$

$$= \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx \quad \text{ie the distribution of } S_n^* \text{ approaches that of the standard normal distribution as } n \text{ increases.}$$

Q: How big does n have to be for the standard normal to give useful approximations?

A: Empirically, about $n=30$, provided there isn't some systematic bias. More is better!

Example: We toss a ^{fair} coin 30 times. (7)
 (X counts heads.)
 Each toss has expected value $\frac{1}{2}$
 & variance $\frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$ $\sigma = \sqrt{\frac{1}{4}} = \frac{1}{2}$

$$P(13 \leq \overset{X_1 + \dots + X_{30}}{\cancel{S_{30}}} \leq 15) \quad S_{30}^* = \frac{X_1 + \dots + X_{30} - 30 \cdot \frac{1}{2}}{\frac{1}{2} \sqrt{30}}$$

$$= P\left(\frac{13-15}{\frac{1}{2}\sqrt{30}} \leq S_{30}^* \leq \frac{15-15}{\frac{1}{2}\sqrt{30}}\right)$$

$$= P\left(\frac{-2}{\frac{1}{2}\sqrt{30}} \leq S_{30}^* \leq 0\right)$$

$$= P\left(\frac{-4}{\sqrt{30}} \leq S_{30}^* \leq 0\right)$$

$$\approx P\left(\frac{-4}{5.5} \leq S_{30}^* \leq 0\right) \approx P(-0.7303 \leq S_{30}^* \leq 0)$$

↓
standard normal

$$\approx P(-0.7303 \leq Z \leq 0)$$

$$= P(Z \leq 0) - P(Z < -0.7303)$$

$$\approx 0.5 - 0.2327 \approx 0.2673$$