

Central Limit Theorem

Lecturer: Jie Fu

- Functions of random variables.
- Expectation, moment, variance of a random variable:
 - Can be defined for both continuous and discrete RV.
- Important property of variance.
- Central Limit Theorem
- Applications of C.L.T.

Expectation

- The **expectation** of a continuous random variable X is defined by

$$\mathbf{E}[X] = \int_{-\infty}^{\infty} x f_X(x) dx$$

- The **expectation** of a discrete random variable X is defined by

$$\mathbf{E}[X] = \sum_k x_k P(X = x_k)$$

- As for discrete random variables, the expectation can be interpreted as
 - "**center of gravity**" of the PDF
 - anticipated **average** value of X in a large number of independent repetitions of the experiment.

Example

1. Let X be the outcome of rolling a **fair 6-sided die**. The probability mass function (PMF) is:

The expectation is calculated as:

2. Let X be an exponential distribution with pdf.

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

The expectation is calculated as:

Function of random variable

- For any real-valued function $g(\cdot)$, $Y = g(X)$ is also a random variable.

- The expectation of $g(X)$ is

$$\mathbf{E}[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx .$$

Example

- A company manufactures LED lightbulbs, and the lifespan (in years) of each bulb follows an Exponential distribution with average life space of 5 years. The company offers a warranty where if a lightbulb fails within 3 years, it is replaced for free. The replacement cost per bulb is \$10.
- What is the expected cost the company need to pay per light bulb under the warranty replacement?

Moments and variance

- The n th **moment** of X is defined by $\mathbf{E}[X^n]$.
- The **variance** of X is defined by

$$\begin{aligned}\mathbf{Var}[X] &= \mathbf{E}[(X - \mathbf{E}[X])^2] \\ &= \int_{-\infty}^{\infty} (x - \mathbf{E}[X])^2 f_X(x) dx\end{aligned}$$

- Please verify the equality.

Property of variance

- $0 \leq \mathbf{Var}[X] = \mathbf{E}[X^2] - (\mathbf{E}[X])^2$
- If $Y = aX$, then

$$\mathbf{E}[Y] = a\mathbf{E}[X], \quad \mathbf{Var}[Y] = a^2\mathbf{Var}[X].$$

- If $Y = X + b$, then

$$\mathbf{E}[Y] = \mathbf{E}[X] + b, \quad \mathbf{Var}[Y] = \mathbf{Var}[X].$$

- if $Y = aX + b$, then

$$\begin{aligned} \mathbf{E}[Y] &= \\ \mathbf{Var}[Y] &= \end{aligned}$$

-
- If $Y = X_1 + X_2$ for two independent RV, then
$$\mathbf{E}[Y] =$$
$$\mathbf{Var}[Y] =$$

- Expectation, moment, variance of a random variable:
 - Can be defined for both continuous and discrete RV.
- Important property of variance.
- Next: Central limit theorem.

Motivating problem:

- A machine process parts, one at a time, in a time independently and uniformly distributed in $[1,5]$.
- What is the probability the machine processes at least 100 parts in 320 time units?

-
- Let X_1, \dots, X_n be a sequence of **independent identically** distributed random variable with mean μ and variance σ^2
 - Let $S_n = X_1 + X_2 + \dots + X_n$, what is the mean of S_n ? What is the variance of S_n ?

Background

The distribution of S_n spreads out as n increases.

But the situation is different if we consider the **sample mean**

$$M_n = \frac{X_1 + \cdots + X_n}{n} = \frac{S_n}{n}.$$

The sample mean is itself a RV (why?) so we can compute its mean and variance

Background

- Given our calculation:

$$\mathbf{E}[M_n] = \mu, \quad \text{var}(M_n) = \frac{\sigma^2}{n}.$$

- The **variance** of M_n **decreases to zero as n increases**.
- Thus, the bulk of the distribution of M_n must be very close to the mean μ as **n increases**.

Background

- We will also consider a quantity which is intermediate between S_n and M_n .
- Z_n is defined as follows.
 1. subtract $n\mu$ from S_n , to obtain the zero-mean random variable $S_n - n\mu$
 2. then divide by $\sigma\sqrt{n}$, to form the random variable.

$$Z_n = \frac{S_n - n\mu}{\sigma\sqrt{n}}$$

$$S_n \sim N(n\mu, n\sigma^2)$$

$$S_n \sim N(1000 \times 0.1, \underbrace{1000 \times 0.09}_{90})$$

$$\frac{S_n - n\mu}{\sqrt{n\sigma^2}} = \frac{S_n - n\mu}{\sigma\sqrt{n}} = \frac{S_n - 100}{\sqrt{90}} = \sqrt{0.09} \times \sqrt{1000}$$

Formally

- Let X_1, \dots, X_n be a sequence of independent identically distributed random variable with mean μ and variance σ^2

- Define

$$Z_n = \frac{S_n - n\mu}{\sigma\sqrt{n}} = \frac{X_1 + \dots + X_n - n\mu}{\sigma\sqrt{n}}.$$

- Let's use jupyter notebook

The Central Limit Theorem

- *Theorem (The Central Limit Theorem)* The CDF of $Z_n = \frac{X_1 + \dots + X_n - n\mu}{\sigma\sqrt{n}}$ converges to standard normal CDF

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-x^2/2} dx$$

in the sense that

$$\lim_{n \rightarrow \infty} P(Z_n \leq z) = \Phi(z)$$

- The central limit theorem is surprisingly general.
- Besides **independence**, and the implicit assumption that the **mean and variance are finite**, it places **no other requirement** on the distribution of the X_i ,
 - which could be discrete, continuous, or mixed.

Going back to our example

Example 2

- A call center receives customer calls according to an exponential distribution with a mean wait time of 4 minutes.

Questions:

- 1. If a single customer calls, what is the probability that they wait more than 5 minutes?
- 2. If we take a random sample of 40 customers, what is the probability that their average wait time is more than 5 minutes?



Advanced thinking: Polling

- p : fraction of population that will vote "yes" in a referendum
- i -th random people polled: 1 : yes, 0: no
- Let X_i be the random variable.
- $M_n = \frac{X_1 + X_2 + \dots + X_n}{n}$ the fraction of "yes" in our sample.
- We would like small error:

$$|M_n - p| \leq 0.01$$

- How many samples to generate so that the probability of error greater than 0.01 is smaller than 0.05?