

Concept of Random Variables

Goals of Messi in a club Football Matches

| Match No | Goal Number |
|----------|-------------|
| 01 | 1 |
| 02 | 3 |
| 03 | 2 |
| 04 | 2 |
| 05 | 0 |
| 06 | 1 |

Here for the same variable Messi Goals number are different for different matches. Also the goal values don't have a relation with each other. This concept is random variable concept.

There are two types of Random variables.

- 1) Discrete Random Variable
- 2) Continuous Random Variable.

1) Discrete Random Variable: (Examples)

- 1) Number of Heads after tossing a coin 10 times.
- 2) Goals of Messi in World cup
- 3) Number of children in each section of class 6.

Continuous Random Variable: (Examples)

- Heights of the students in a class
- Time remaining for the next bus to come
- Family members consumption of coke from a 2L bottle.

Probability Distribution on Discrete Random Variables:

Let X be a discrete random variable which denotes the values $x_0, x_1, x_2, \dots, x_n$

$$P(X=x_i) = f(x_i) = f_i$$

$$x_i = \{x_0, x_1, x_2, \dots, x_n\}$$

The value of f_i depends on x_i i.e. $i \rightarrow \{x_0, x_1, x_2, \dots, x_n\}$

This Function, f_i called the Probability Mass Function (PMF).

The set of ordered pair (x_i, f_i) is called the discrete probability distribution of X .

Ordered pair $(a, b) = (u, v)$ iff, $a=u$ and $b=v$

This distribution represent as below:

$$X: x_0, x_1, x_2, \dots, x_n$$

$$f_i: f_0, f_1, f_2, \dots, f_n$$

Let's take an example. Suppose we toss two coin. at a time.

$$\text{Sample Space } (S) = \{ HH, HT, TH, TT \}$$

X is a random variable which says,

X : No. of head ~~appear~~ appears, for every coin toss. 2 times.

$$X(HH) = 2, \quad X(HT) = 1, \quad X(TH) = 1, \quad X(TT) = 0$$

Range of X is $\{0, 1, 2\}$

$$P(X=0) \rightarrow \text{Means No Head} = \frac{1}{4} \quad \{ P(X=0) = f_0 \}$$

$$P(X=1) \rightarrow \text{Means 1 Head} = \frac{2}{4} = \frac{1}{2} \quad \{ P(X=1) = f_1 \}$$

$$P(X=2) \rightarrow \text{Means 2 Head} = \frac{1}{4} \quad \{ P(X=2) = f_2 \}$$

Let's represent the distribution \rightarrow

| | | | |
|---------------|---------------|---------------|---------------|
| $X :$ | 0 | 1 | 2 |
| (PMF) $f_i :$ | $\frac{1}{4}$ | $\frac{1}{2}$ | $\frac{1}{4}$ |

Properties :

$$1) f_i \geq 0 \rightarrow \left(\frac{1}{4} \geq 0, \frac{1}{2} \geq 0, \frac{1}{4} \geq 0 \right)$$

$$2) \sum f_i = 1 \rightarrow \left(\frac{1}{4} + \frac{1}{2} + \frac{1}{4} = 1 \right)$$

Binomial Distribution:

I have noted down the concept well in PW-skills part.

Question: You tossed a coin 3 times. What is the probability of getting 2 Heads?

In Binomial Distribution, PMF = $\frac{n!}{x!(n-x)!}$

$$P(x) = {}^n C_x p^x (1-p)^{n-x}$$

$$\therefore P(x) = {}^3 C_x \left(\frac{1}{2}\right)^x \left(1 - \frac{1}{2}\right)^{3-x}$$

$$= {}^3 C_x \left(\frac{1}{2}\right)^x \left(\frac{1}{2}\right)^{3-x}$$

$$\therefore P(2) = {}^3 C_2 \left(\frac{1}{2}\right)^2 \left(\frac{1}{2}\right)^{3-2}$$

$$= {}^3 C_2 \times \frac{1}{4} \times \frac{1}{2}$$

$$= \frac{3}{8}$$

n = number of trials

P = Probability of success

$$n = 3$$

$$P = \frac{1}{2}$$

↳ The probability of getting Head for each toss

x = Number of success

(4 means number of heads)
 $= 2$

Question 2: In a factory that produces light bulbs. It is known that 10% of the bulbs are defective. A inspector randomly selects 15 light bulbs from a recent batch. What is the probability that 3 of the 15 selected bulbs are defective.

→ Here the probability of success (finding a defective bulb) = 10% = 0.10

PMF Function, $P(X=K) = {}^nC_K P^K (1-P)^{n-K}$

Here, $n = 15$ (number of trials)

$K = 3$ (number of success) → Getting a defective bulb

$P = 0.10$ (probability of defective)

$$P(X=K) = {}^{15}C_3 (0.10)^3 (1-0.10)^{15-3}$$

$$= 0.332$$

Ans

Question 3:

A website administrator monitors the traffic to their site and observe that the avg click rate for a particular ad is 20%. If the add is shown to 100 people, what is the probability that it would be clicked by 15 people?

Here, $n = 100$

$p = 0.2$

$K = 15$

$$\therefore \text{PMF} \rightarrow P(X=K) = {}^nC_K p^K (1-p)^{n-K}$$

$$= {}^{100}C_{15} (0.2)^{15} (1-0.2)^{100-15}$$

$$= 0.201$$

Ans

Bernouli Distribution:

This concept also noted down well in the PW-sikls part.

Question 1: In a series of basketball free throws, a player makes a successful shot (denotes by 1) with a probability of 0.75. What is the probability that the player makes ^a 8 successful shots in a sequence of ^a 8 free throws.

Here, $n = 8$

$K = 8$

$p = 0.75$

$$P(X=K) = {}^nC_K p^K (1-p)^{n-K}$$

$$= {}^8C_8 p^8 (1-p)^{8-8}$$

$$= p$$

$$= 0.75 \quad \underline{\text{Ans}}$$

Question 2:

In a deck of 20 playing cards, 6 cards are red and the rest are black.

If a card is drawn at random from the deck, where red is considered a success. What is the probability of drawing a red card.

$$p = \frac{6}{20} \\ = \frac{3}{10} = 0.3$$

For Bernoulli Distribution

$$P(X=K) = {}^nC_K \cdot p^K \cdot (1-p)^{n-K}$$

Now, $n=1$
 $K=1$ } For Bernoulli

$$\therefore P(X=K) = 1 \times p^1 \times (1-p)^{1-1} \\ = p \rightarrow \text{For success}$$

$q = (1-p) \rightarrow \text{For not success}$

Example of Probability Distribution: (How we represent it)

Suppose, you have two dice. You need to show the probability distribution of the sum of the two dice after throwing.

Sum of the dices after throwing can be \rightarrow

$\{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$

Let's make the dice table \rightarrow

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|----|----|----|--------------|--------------|--------------|---------------|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | | | | |
| 2 | 3 | 4 | 5 | 6 | 7 | 8 | | | | |
| 3 | 4 | 5 | 6 | 7 | 8 | 9 | | | | |
| 4 | 5 | 6 | 7 | 8 | 9 | 10 | | | | |
| 5 | 6 | 7 | 8 | 9 | 10 | 11 | | | | |
| 6 | 7 | 8 | 9 | 10 | 11 | 12 | | | | |

$$P(2) = \frac{1}{36}$$

$$P(3) = \frac{2}{36} = \frac{1}{18}$$

$$P(4) = \frac{3}{36} = \frac{1}{12}$$

$$P(5) = \frac{4}{36} = \frac{1}{9}$$

$$P(6) = \frac{5}{36}$$

$$P(7) = \frac{6}{36}$$

$$P(8) = \frac{5}{36}$$

$$P(9) = \frac{4}{36}$$

$$P(10) = \frac{3}{36}$$

$$P(11) = \frac{2}{36}$$

$$P(12) = \frac{1}{36}$$

So, Representation of probability distribution

Dice (sum) \rightarrow possible outcomes

| | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|-------------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| P(sum) | | | | | | | | | | | |
| P(sum) | $\frac{1}{36}$ | $\frac{2}{36}$ | $\frac{3}{36}$ | $\frac{4}{36}$ | $\frac{5}{36}$ | $\frac{6}{36}$ | $\frac{5}{36}$ | $\frac{4}{36}$ | $\frac{3}{36}$ | $\frac{2}{36}$ | $\frac{1}{36}$ |

But the sample space was really less for the previous exam for which we ~~can~~ could draw the tables and calculate the probability.

But that can't be done if we wanted to get the probability of each sum for 10000, 1M Dice. right?

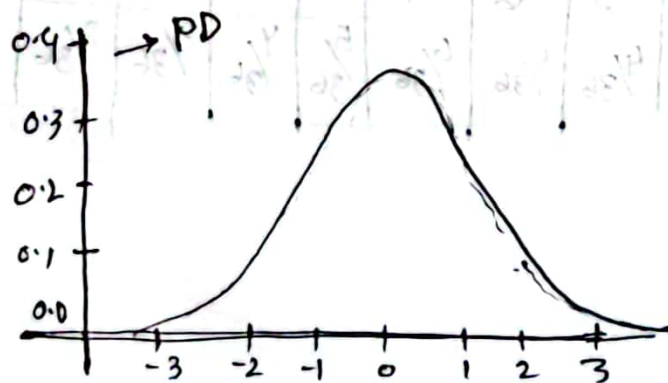
So, what we can do is we can create Probability distribution functions (PMF, PDF, CDF) using which we can plot graph and watch the probability distribution.

(I have already noted down PMF, PDF, CDF concepts in 'PW-Skills Section'.)

For getting the probability of each sum from rolling 2 dice,

$$\text{PMF} = \begin{cases} 1/36 & \text{when } x \in \{2, 12\} \\ 2/36 & \text{when } x \in \{3, 11\} \\ \vdots & \\ 0 & \text{otherwise} \end{cases}$$

As PDF is little bit complex, I would note it again here briefly.

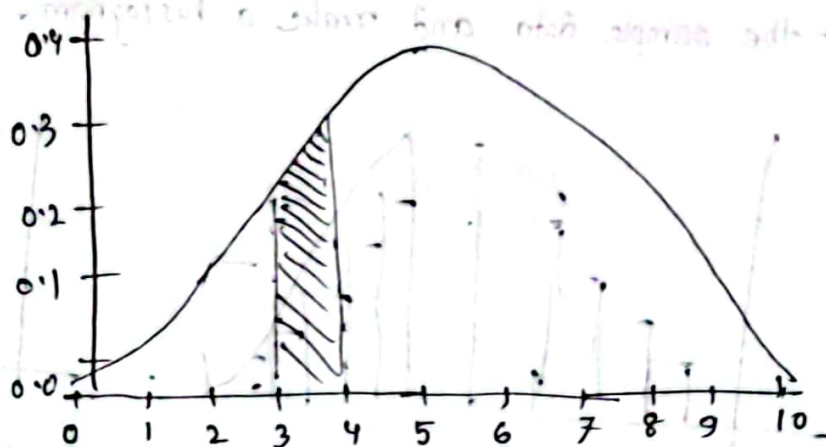


in PMF graph, y-axis we have **probability**

But in PDF, in y-axis we have **probability Density**.

What is probability Density? What is the reason to use it?

→ In CDF, we use continuous Random variables. Suppose, we created a PDF where in x axis the cgpa are given. Now suppose, I want to find whose cgpa is 3.874. Between cgpa 3 to 4 there can be infinite numbers. The count is so big the probability of finding only one value in that range is 0. So that's why we can't use probability in Y axis for PDF. We have to use Probability density.



using PDF we can't say the probability if 3.874 instead we can find the probability of having a value between 3 to 4 by calculating the area.

$$\int_3^4 f(x) dx$$

→ This integration will provide the value for probability for 3-4.

Now, the rectangle area that we got, we can make it more thinner

to get the single value probabilities like 3.874

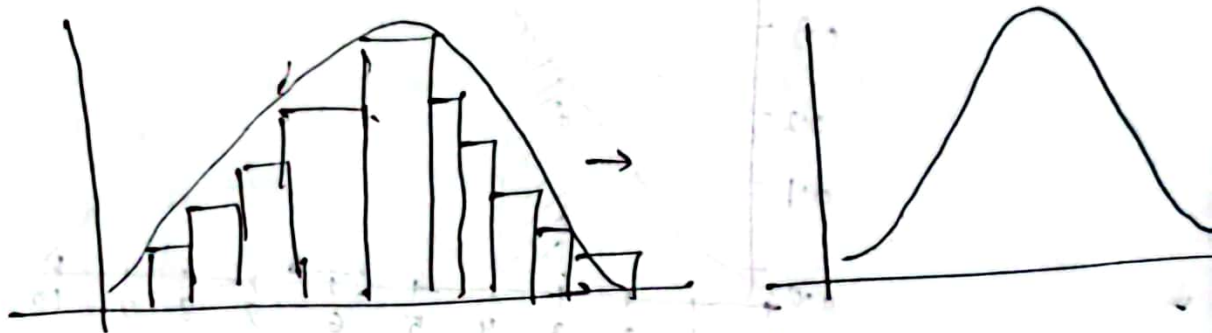
How to estimate Density:

There are two types of Density estimation:

- 1) Parametric
- 2) Non Parametric

Parametric Distribution:

First we plot the sample data and make a histogram.



Suppose this is the histogram. Then we will check the histogram looks like which distribution. Here, we can see that the histogram looks like Normal Distribution. Now we do everything according to Normal Distribution. We will get the μ (Mean) and σ (std) from the available data. Now from the sample data, we have to estimate the population Mean and population std (σ).

Then we have to use the PDF equation to calculate probability ^{density}.

$$PDF = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

for each sample data value.

Non Parametric Technique:

When after plotting the histogram, the graph is not similar to any probability Distribution. Then we apply Non parametric Technique.

Plot CDF for PDF:

If you Integrate PDF, you can get CDF (Integration)

~~If~~ If you Differentiate CDF, you get PDF (Differentiation).

CDF for PDF is noted down on PW-skills section.

Uniform Distribution:

The concept is noted down in PW-skills section well.

Some Examples:

- ① The height of a person randomly selected from a group of individuals whose height range from 5'6" to 5'10" would follow a continuous uniform distribution.
- ② The time taken for a machine to produce a product, where the production time ranges from 5 to 10 minutes, would follow a continuous uniform distribution.

- ③ The distance that a randomly selected car travels on a tank of gas, where the distance ranges from 300 to 400 miles, would follow a continuous uniform distribution.

Problem 1: Buses in a city are scheduled to arrive at a particular bus stop ~~at~~ every 15 minutes. With the first bus arriving at exactly 8 AM. If a passenger arrives at the bus stop at a random time between 8 AM to 9 AM, what is the probability that they will wait less than 5 minutes for the next bus?

Here, $a = 0 \text{ min}$
 $b = 60 \text{ min}$ [Because Range $9 \text{ AM} - 8 \text{ AM} = 60 \text{ min}$]
 $x_2 = 5 \text{ min}$
 $x_1 = 0 \text{ min}$

$$\text{Pr}(0 \leq x \leq 5) = (x_2 - x_1) \frac{1}{b - a}$$

$$= \frac{5}{60} = \frac{1}{12} \quad \underline{\text{Ans}}$$

Problem 2: In a parking lot, the available parking spaces are uniformly distributed throughout the lot. The lot has a total of 200 spaces. On average 40% of the spaces are occupied. What is the probability that the next car, arriving, finds an available parking space?

Here, Unoccupied spaces = $(100 - 40) = 60\%$

$$\text{Number of unoccupied spaces} = 0.6 \times 200 \\ = 120$$

$$\therefore \text{Probability of finding space} = \frac{120}{200} \\ = 0.6 = 60\% \quad \underline{\text{Ans}}$$

Problem 3: An internet router's download speed, uniformly distributed between 50 mbps and 100 mbps. Calculate the probability that a random user experiences a download speed at least 70 mbps?

$$\text{Here, } b = 100 \\ a = 50$$

$$x_2 = 100 \quad \left[\begin{array}{l} \text{because at least 70 mbps means} \\ \text{minimum must be 70 mbps} \end{array} \right.$$

$$x_1 = 70$$

So max can be 100 mbps

$$\therefore \text{Pr}(\text{speed at least 70 mbps}) = \frac{100 - 70}{100 - 50} \times \frac{1}{100 - 50} \\ = \frac{30}{50} = 0.6 = 60\% .$$

Ans

Normal Distribution: Normal distribution is noted in PW-skills section.

But the Standard Normal Distribution part is not there.

Here is the explanation of that,

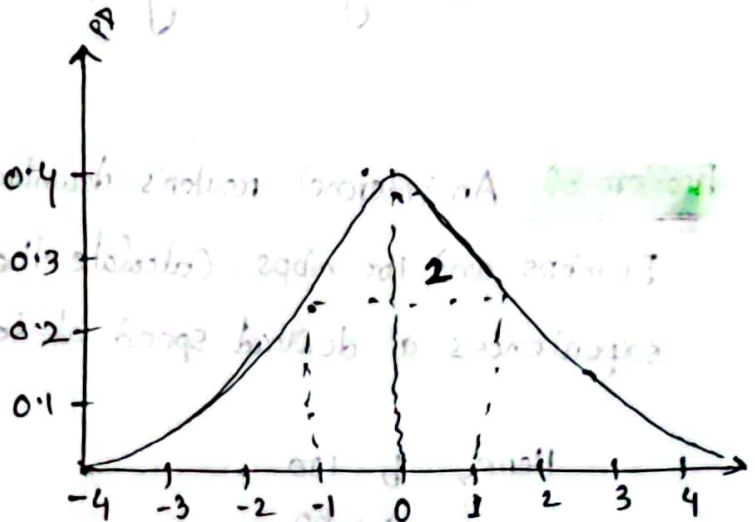
Standard Normal Distribution:

Parameters: $\mu=0, \sigma=1$

Notation: $X \sim N(0, 1^2)$

$$f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-0}{1^2} \right)^2}$$

$$= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} x^2}$$



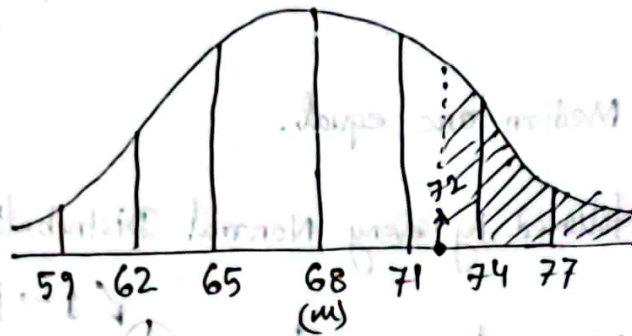
Standardization: To convert any normal distribution to the standard form

$$Z = \frac{\sum_{i=1}^n X_i - n\mu}{\sigma \sqrt{n}}$$

Standardization is very crucial to compare variables of different magnitudes.

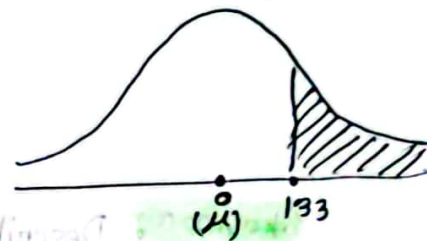
Problem: Suppose, the heights of adult males in a certain population follow a normal distribution with a mean of 68 inch, std of 3 inches. What is the probability that a randomly selected adult male from this population taller than 72 inches?

$$x \sim N(68, 3)$$



Standardized value of 72, $z = \frac{72 - 68}{3} \left[\frac{x - \mu}{\sigma} \right]$

$$= 1.33$$



Using the Z score in Z score table, we found ~~0.9824~~ 0.90824, this probability (~~98.7~~ 90.8%)

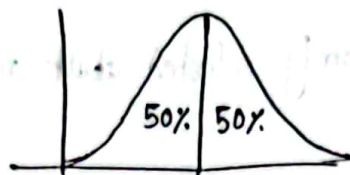
So, in the population ~~98.7~~ 90.8% people are shorter or equal to 72 inches.

$$\begin{aligned} \therefore \text{Taller people} &= (1 - 0.90824) \\ &= 0.09176 \\ &= 9.2\% \end{aligned}$$

The **advantage** of converting normal distribution to standard normal distribution is because then we can find the Z score and using the Z score and Z table we can find the desired outcome.

Problem: Properties of Normal Distribution:

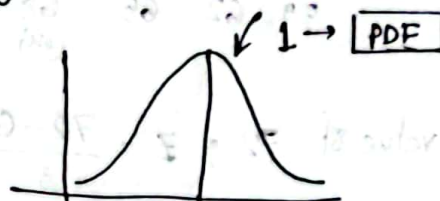
1) Symmetry



2) Mean, Mode and Median are equal.

3) Empirical rules followed by every Normal Distribution

4) Area under the curve = 1



Skewness: Describes the degree to which a dataset deviates from the normal distribution. (More noted on PW-skills section)

Moments in Statistics:

1st Moment → Mean

2nd Moment → Variance

3rd Moment → Skewness

4th Moment → Kurtosis.

Formula of sample skew: →

$$\frac{n}{(n-1)(n-2)} \sum \left(\frac{x - \bar{x}}{s} \right)^3 \rightarrow \text{(3rd moment)}$$