

Disclaimer: Please note that any topics that are not covered in today's lecture will be covered in the next lecture.

Content

- Empirical vs Theoretical Probability
- Expectations
- Binomial Distribution
- Bernoulli Distribution

▼ Distribution Functions

Probability Density Function (PDF):

- The PDF is a function that describes the probability density of a continuous random variable over its range.
- The term "density" here is similar to how tightly data is packed around a specific point, like cars on a road.

Probability Mass Function (PMF):

- The PMF is a function that describes the probability of a discrete random variable taking on a specific value.

Cumulative Distribution Function (CDF):

- The CDF is a function that gives the probability that a random variable is less than or equal to a specified value.

Let's implement this using a height dataset

We will going to work on the height dataframe that we saw above for now

```
df_height = df_hw["Height"]
df_height.head()
```

```
0    73.847017
1    68.781904
2    74.110105
3    71.730978
4    69.881796
Name: Height, dtype: float64
```

```
# minimum height
min_height = df_height.min()
min_height
```

```
54.2631333250971
```

```
# maximum height
max_height = df_height.max()
max_height
```

```
78.9987423463896
```

```
total = len(df_height)
total
```

```
10000
```

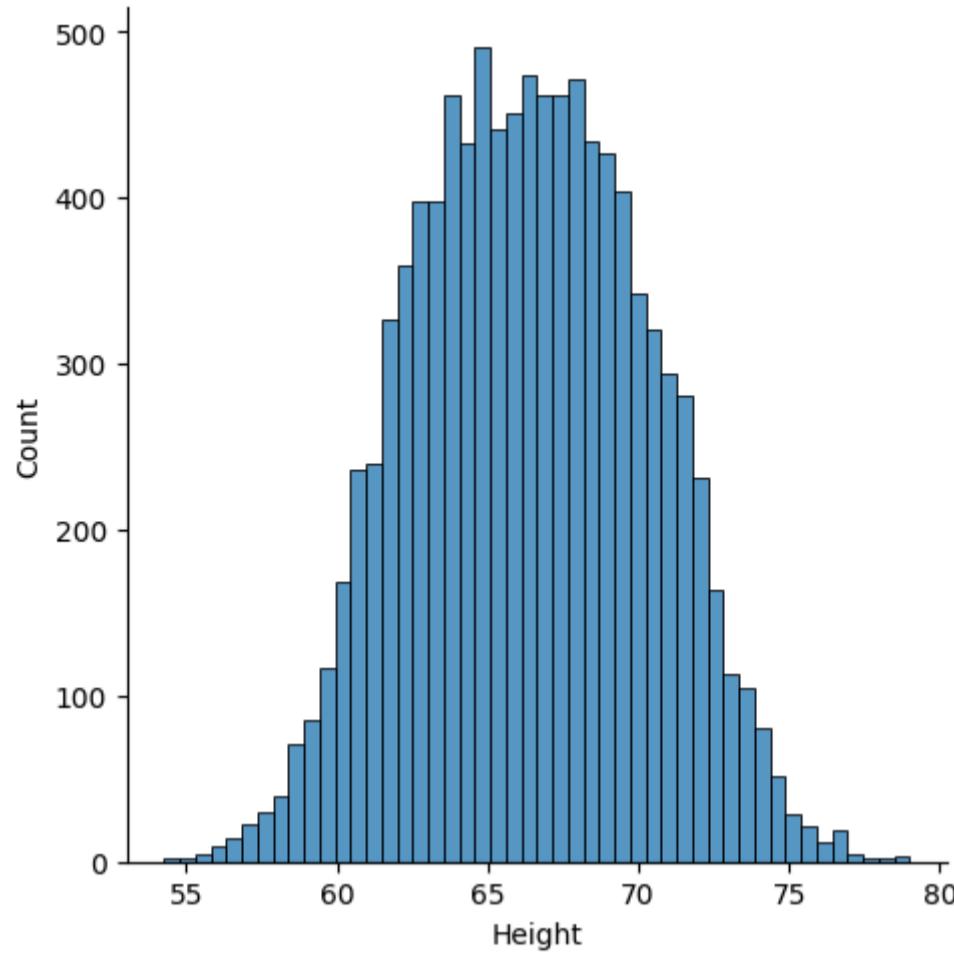
To plot this type of distribution we generally use Histograms or Distribution plots

▼ Histogram

It is a graphical representation of a dataset's distribution, showing the frequency or probability of different values within the data.

```
sns.displot(df_height)
```

<seaborn.axisgrid.FacetGrid at 0x7ed0633c3670>



Q.What we can understand from this distribution?

- Each bar in the histogram represents one of the intervals or ranges,
- The height of the bar indicates the frequency or number of data points falling within that interval.

Count:

- It indicates the "**frequency**", which means in the particular bar or range of height, how many values are there.
 - We can assert this like, **around 500 people have their height in the range of 63 - 65 (that on bar)**

This is what histograms or distribution plots tell about the data

Now let's have a look into some distribution functions

▼ Probability Mass Function (PMF)

The PMF is a function that describes the probability of a discrete random variable taking on a specific value.

It associates each possible value of the random variable with its probability of occurrence.

Example: Rolling a Fair Six-Sided Die

- For example, if we have a discrete random variable X representing the outcome of rolling a fair six-sided die, the PMF might look like
 - $P(X = 1) = \frac{1}{6}$, $P(X = 2) = \frac{1}{6}$, and so on.

Probability Mass Function



▼ Probability Density Function (PDF)

PDF is used for **continuous random variables**, as opposed to PMF, which is for discrete variables.

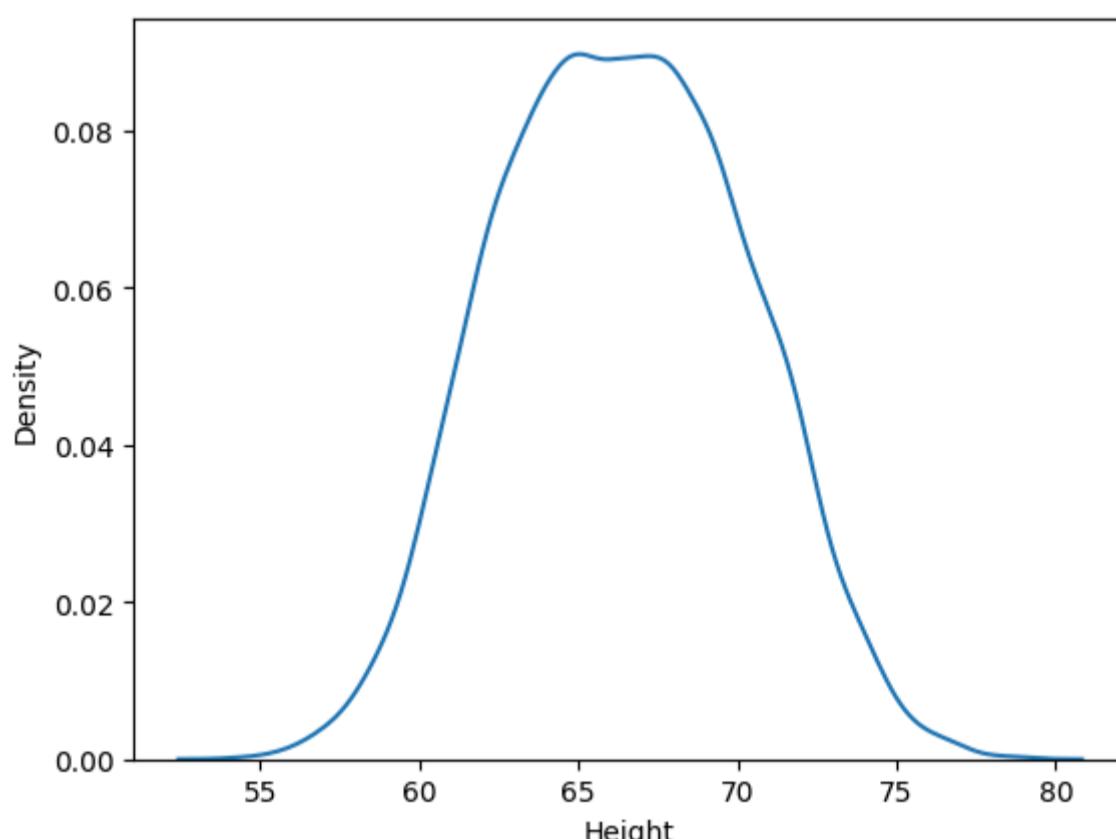
PDF does not provide the probability of a specific value but gives the **probability of the random variable falling within a certain interval**

- For instance, it answers questions like "What are the chances that the next height chosen will fall between 62 and 65 inches?"

We can visualize a PDF by using distribution plots like histograms or KDE (Kernel Density Estimation) plots.

```
sns.kdeplot(df_height)
```

```
<Axes: xlabel='Height', ylabel='Density'>
```



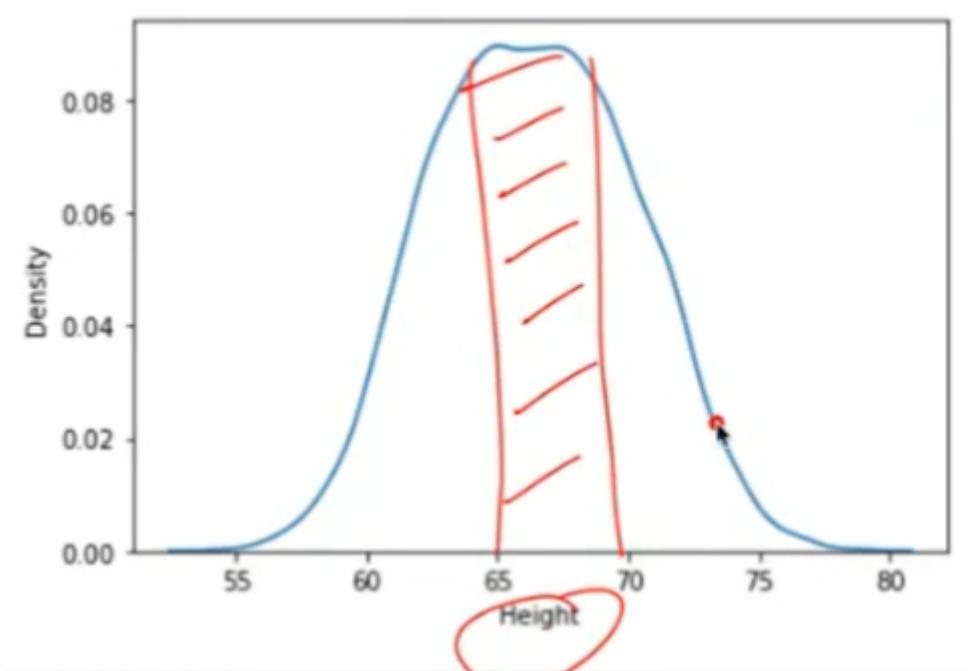
Example:

If we have a continuous random variable Y representing the height of people in a population,

The **PDF might represent the probability that a randomly chosen person has a height within a certain range, such as between 65 and 70.**

- We will find out the area under that interval to find the probability

```
<AxesSubplot:xlabel='Height', ylabel='Density'>
```



Next up we have:

▼ Cumulative Distribution Function (CDF)

The Cumulative Distribution Function (CDF) describes the probability that a random variable takes on a value less than or equal to a given value.

In the context of this dataset, in CDF, we talk about fractions of people who are less than the given height

- Let's say you take 60 inches, then what fraction of the people have less than or equal to this value? This fraction is calculated using CDF
- It gives you the cumulative probability up to a certain point.

Example:

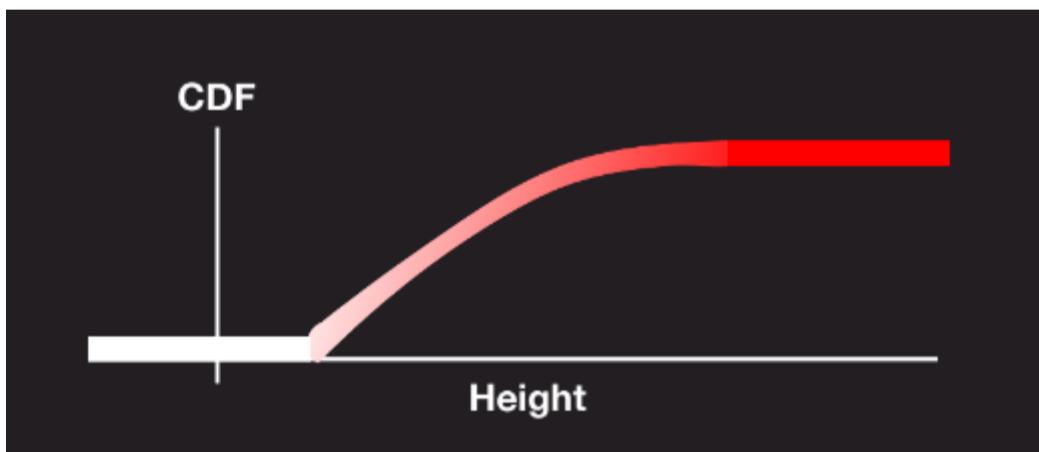
If you have a random variable Z representing the number of heads in three coin tosses,

The CDF would tell you the probability that Z is less than or equal to a certain number, like $P(Z \leq 2)$.

How to calculate CDF?

The CDF is calculated by accumulating the probabilities for each height value.

- As you move along the X-axis (height values) on the CDF graph, you're essentially adding up the probabilities
- It shows how likely it is to find someone with a height less than or equal to that value.



- The CDF graph typically starts at 0% on the Y-axis (probability) when height is at its minimum (in our dataset)
- It ends at 100% when height is at its maximum.
- The curve starts at the left and gradually climbs towards the right.
- The steepness of the curve at a particular point represents how quickly the probability is accumulating

Conclusion

So, the PDF shows you the probability of a specific height, while the CDF shows you the probability of heights up to a certain value in your dataset.

Percentile $25 \rightarrow 63.5$

CDF: $63.5 \rightarrow 0.25$

Conclusion :

In summary, the relationships are as follows:

- The PMF is used for discrete random variables.
- The PDF is used for continuous random variables.
- The CDF is used for both discrete and continuous random variables to provide cumulative probabilities.

These functions are essential tools in probability and statistics for understanding the behaviour of different probability distributions.

▼ Case study on Empirical vs Theoretical Probability

▼ Casino Case Study

A bag has **3 Red** and **2 Blue** balls.

You pick a ball, write its colour, and **put it back** in the bag. This is done **4 times** in total.

If all 4 times, the **Red balls** was drawn, you **win Rs 150**.

Otherwise you **lose Rs 10**.

Question : Would engaging in this game result in a profit or loss for you?

Discuss:

- Problem mentions that once you've noted the color, you put it back in the bag
 - What does this mean in the probability language?
 - It means that the balls are drawn **with replacement**
- The step of taking out the ball is repeated 4 times

Whether you end up gaining or losing will depend on how many red balls are drawn.

Therefore, let's define a random variable X to denote the number of red balls drawn.

- Hence, X will be a discrete random variable.
- Possible outcomes of X : 0, 1, 2, 3 or 4

Empirical Approach

▼ Motivation for Empirical Approach

We know that it is possible to get 7 heads on 10 coin tosses, when using a fair coin.

So how would one go about proving that $P(\text{Heads}) = 0.5$ for a fair coin?

How many of you have heard about the scientist who wanted to prove this?

In order to do so, he tossed a fair coin 10,000 times repeatedly, and noted down his observation on each toss.

The idea was

- Though 7 heads is probable for 10 tosses, when a coin is tossed for 10,000 times, the number of heads should be approximately 5,000

This process of simulating the experiment, and repeating it multiple times, is done in an effort to calculate probability value (of getting heads in this example).

This value is known as **Empirical Probability**.

- The idea is make estimates using real-world data/observations

Let's try to estimate whether we will have a profit or loss after playing this game, using probabilities calculated by empirical approach.

For this, we will simulate this situation in Python code.

```
import math as m
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

Let's simulate the given Casino problem using `np.random.choice()`

Since we have,

- 3 red balls and
- 2 blue balls

We can represent the possible outcomes as a list `["R", "R", "R", "B", "B"]`

Since the ball is being drawn 4 times, we set `size = 4`

Notice that every time we execute the following code, we will get a different result, and they're being chosen randomly each time.

```
#Code to be shared to learners
rolls=np.random.choice(["R","R","R","B","B"],size=4)
rolls

array(['R', 'B', 'B', 'R'], dtype='<U1')
```

Recall that X represents the number of red balls drawn in a simulation.

▼ How to evaluate X from a simulation?

1. Create a Boolean mask having all observed "R" as True and "B" as False
2. Use `np.count_nonzero()` to count number of True in the mask.

```
# rolls=="R" creates a Boolean mask
# count_nonzero() counts the number of non-zero OR True elements in the passed list.

np.count_nonzero(rolls=="R")
```

2

Let's take inspiration from the scientist and perform this simulation 10,000 times using a code, and note our observations.

Discuss:

- We already know how to simulate a ball draw
- Let's store the no of reds observed in a variable `num_red`
- And store this value for all 10,000 simulations into a list `red_values`

```
red_values=[]

for person in range(10000):
    rolls=np.random.choice(["R","R","R","B","B"],size=4)
    num_red=np.count_nonzero(rolls=="R")
    red_values.append(num_red)
pd.value_counts(red_values)
```

```
3    3552
2    3394
1    1496
4    1281
0     277
dtype: int64
```

```
# red_values
```

Let's do a `.value_counts()` to see the frequency of values it contains.

```
pd.value_counts(red_values, normalize=True)
```

```
3    0.3552
2    0.3394
1    0.1496
4    0.1281
0    0.0277
dtype: float64
```

We are aware that passing `normalize=True` in `value_counts()` gives us the result in percentage of their occurrence.

We can see that the probability of drawing 3 red balls is 0.3552, 2 red balls is 0.3394 and so on..

Based on this data, how many red balls we will get on an average based on simulations we have done 10,000 times?

```
# This is empirical value
np.mean(red_values)
```

```
2.4064
```

▼ Expectation using Empirical Approach

How do you think, this mean was calculated from these frequency values?

As you learnt in the last class, this was calculated as a result of **Weighted Average**

So, for the given frequency count, we can see that this is calculated as:

$$\text{Mean} = \frac{4(1281) + 3(3552) + 2(3394) + 1(1496) + 0(277)}{1281 + 3552 + 3394 + 1496 + 277} = \frac{4(1281) + 3(3552) + 2(3394) + 1(1496) + 0(277)}{10000}$$

$$(4*(1281) + 3*(3552) + 2*(3394) + 1*(1496) + 0*(277)) / (10000)$$

```
2.4064
```

Now that we've verified this,

Let's represent the same equation in a slightly different format.

$$\text{Mean} = 4\frac{1281}{10000} + 3\frac{3552}{10000} + 2\frac{3394}{10000} + 1\frac{1496}{10000} + 0\frac{277}{10000}$$

If you closely look at the value counts table, you will see that this can be represented as the following formula:

$$E(X) = \sum_i X_i * P(X = X_i)$$

where

- X was our random variable that denotes the no of red balls drawn.
- $P(X = X_i)$ represents the probability of X getting a value of X_i
- $E(X)$ is known as the **Expected value** of the random variable X

Let's define it formally:

Expectation of a random variable X , is the weighted average of the values that X takes, with the weights being the probabilities.

$$E(X) = \sum X^* P[X]$$

$$E(X) = 1 \times P[X=1] + 2 \times P[X=2] + 3 \times P[X=3] + 4 \times P[X=4]$$

Until now, we simulated the event 10,000 times, and found an expected value of random variable X using the data observed.

This is known as the **Empirical Approach** of solving the problem.

▼ Theoretical Approach

Now, let's solve this case study using theoretical approach and observe the difference in the result

Now let's discuss the theoretical approach to this Casino case study

Let's look at the problem statement once more.

A bag has **3 Red** and **2 Blue** balls.

You pick a ball, write its colour, and **put it back** in the bag. This is done **4 Times** in total.

If all 4 times, the **Red balls** was drawn, you **win Rs 150**.

Otherwise you **lose Rs 10**.

Question: Would engaging in this game result in a profit or loss for you?

Casino case study

A bag has 3 red and 2 blue balls.



You pick a ball, write its colour, and put it back in the bag. This is done 4 times in total.

If all 4 times, the red ball was drawn, you win Rs 150. In any other case, you lose Rs 10.

Would you play this game?

Let "X" denote the number of red balls when you draw 4 balls with replacement
Here, X is an example of what is called a "Random Variable"

Theoretical approach: Compute probability using rules

What is the probability of 1 red ball in 1 pick?

$$P[\bullet] = 3/5$$

What is the probability of 1 blue ball in 1 pick?

$$P[\bullet] = 2/5$$

What is the probability of 2 red balls in 2 picks?

$$P[\bullet\bullet] = (3/5)(3/5)$$

What is the probability of 1 red ball in first pick and 1 blue ball in second?

$$P[\bullet\bullet] = (3/5)(2/5)$$

What is the probability of 1 blue ball in first pick and 1 red ball in second?

$$P[\bullet\bullet] = (2/5)(3/5)$$

$$P[\bullet\bullet\bullet] = (3/5)(3/5)(3/5)(2/5)$$

$$P[\bullet\bullet\bullet] = (3/5)(3/5)(3/5)(2/5)$$

$$P[\bullet\bullet\bullet] = (2/5)(2/5)(2/5)(2/5)$$

Let's define 2 events:

- *R*: Drawing a red ball
- *B*: Drawing a blue ball

What would be the probability of obtaining a red ball once?

$$P(R) = \frac{3}{5}$$

Similarly, we know that $P(B) = \frac{2}{5}$

What is the probability of drawing a red ball twice?

$$P(RR) = \frac{3}{5} * \frac{3}{5}$$

What is the probability of drawing a red ball followed by a blue ball?

$$P(RB) = \frac{3}{5} * \frac{2}{5}$$

These values are easy to evaluate when we are drawing the balls just twice.

In our case study, we are **drawing it 4 times**. Let's consider that case.

Like before, we define X as a random variable that denotes the no of red balls drawn.

What would be the probability of obtaining 1 red ball?

For $X = 1$, we can have 4 possible cases as drawn below:

- BBBR
- BBRB
- BRBB
- RBBB

Let's look at the probability value of each of these individual cases:

- Case 1: $\frac{2}{5} * \frac{2}{5} * \frac{2}{5} * \frac{3}{5}$
- Case 2: $\frac{2}{5} * \frac{2}{5} * \frac{3}{5} * \frac{2}{5}$

and so on

So we can see that for all these 4 cases, we can write their probability as: $(\frac{2}{5})^3 * (\frac{3}{5})^1$

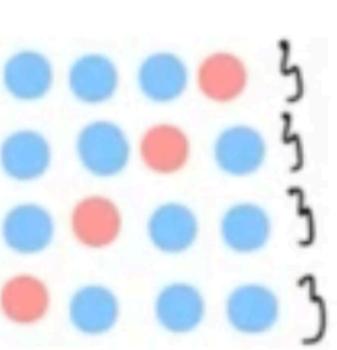
Since there are 4 such cases, we write the total probability of $X = 1$ as:

$$P(X = 1) = \text{case 1 OR case 2 OR case 3 OR case 4}$$

$$P(X = 1) = 4 * (\frac{2}{5})^3 * (\frac{3}{5})^1$$

Theoretical Probability

$P[R] = \frac{3}{5}$ $P[B] = \frac{2}{5}$
 $P[RR] = \left(\frac{3}{5}\right)\left(\frac{3}{5}\right)$ $P[RB] = \left(\frac{3}{5}\right)\left(\frac{2}{5}\right)$

$X=1$ 
The probability calculation shows four terms:
 $\left(\frac{2}{5}\right)\left(\frac{2}{5}\right)\left(\frac{2}{5}\right)\left(\frac{3}{5}\right)$ $\left(\frac{2}{5}\right)\left(\frac{2}{5}\right)\left(\frac{3}{5}\right)\left(\frac{2}{5}\right)$
These are grouped by a brace under the term $\left(\frac{2}{5}\right)^3\left(\frac{3}{5}\right)^1$.

$P[X=1] = 4 \left(\frac{2}{5}\right)^3\left(\frac{3}{5}\right)^1$

What would be the probability of getting 2 red balls out of the 4 balls drawn?

Let's look at the different orientations possible for $X = 2$.

- We have 6 possibilities.

Let's look at the probability of each of these orientations:

- Case 1: $\frac{2}{5} * \frac{2}{5} * \frac{3}{5} * \frac{3}{5}$

... and so on

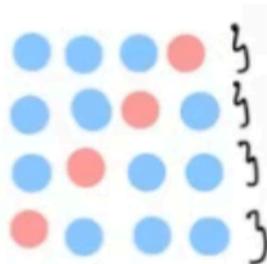
So, at the end of the day, we know that probability for each of these individual cases would be: $(\frac{2}{5})^2 * (\frac{3}{5})^2$

Since either of these 6 cases are possible, the total probability becomes:

$$P(X = 2) = 6 * (\frac{2}{5})^2 * (\frac{3}{5})^2$$

Theoretical Probability

$$X=1$$



$$\left(\frac{2}{5} \right) \left(\frac{2}{5} \right) \left(\frac{2}{5} \right) \left(\frac{3}{5} \right)$$

$$\left(\frac{2}{5} \right)^3 \left(\frac{3}{5} \right)^1$$

$$P[R] = \frac{3}{5} \quad P[B] = \frac{2}{5}$$

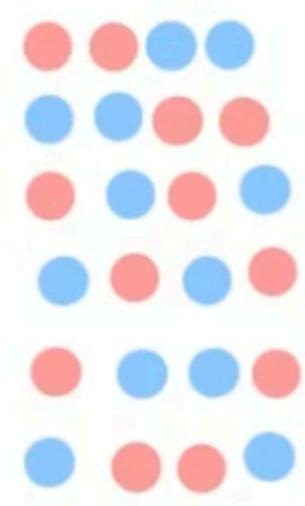
$$P[RR] = \left(\frac{3}{5} \right) \left(\frac{3}{5} \right) \quad P[RB] = \left(\frac{3}{5} \right) \left(\frac{2}{5} \right)$$

$$X=2$$

$$P[X=2] = 6 \left(\frac{3}{5} \right)^2 \left(\frac{2}{5} \right)^2$$

$$\left(\frac{2}{5} \right) \left(\frac{2}{5} \right) \left(\frac{3}{5} \right) \left(\frac{3}{5} \right)$$

$$\left(\frac{3}{5} \right) \left(\frac{3}{5} \right) \left(\frac{2}{5} \right) \left(\frac{2}{5} \right)$$



Conclusion

Can we write this 4 and 6 in a different format?

Recall the combinatorics lecture.

We know that

- $4 = {}^4C_1$
- $6 = {}^4C_2$

With this in mind, when we take a look at the results of $P(X = 1)$ and $P(X = 2)$, can we derive some general expression?

$$P(X = k) = {}^4C_k \left(\frac{3}{5} \right)^k \left(\frac{2}{5} \right)^{4-k}$$

Notice that here, 4 is nothing but the no of times a ball was drawn from the bag, i.e. **no of trials**

We can use this derived equation to find probability for all valid values of the random variable X :

- $P(X = 0) = {}^4C_0 \left(\frac{3}{5} \right)^0 \left(\frac{2}{5} \right)^4$
- $P(X = 1) = {}^4C_1 \left(\frac{3}{5} \right)^1 \left(\frac{2}{5} \right)^3$
- $P(X = 2) = {}^4C_2 \left(\frac{3}{5} \right)^2 \left(\frac{2}{5} \right)^2$
- $P(X = 3) = {}^4C_3 \left(\frac{3}{5} \right)^3 \left(\frac{2}{5} \right)^1$
- $P(X = 4) = {}^4C_4 \left(\frac{3}{5} \right)^4 \left(\frac{2}{5} \right)^0$

$${}^4C_k \left(\frac{3}{5}\right)^k \left(\frac{2}{5}\right)^{4-k}$$

$$\begin{aligned} X = 0 & \quad {}^4C_0 \left(\frac{3}{5}\right)^0 \left(\frac{2}{5}\right)^4 \\ X = 1 & \quad {}^4C_1 \left(\frac{3}{5}\right)^1 \left(\frac{2}{5}\right)^3 \\ X = 2 & \quad {}^4C_2 \left(\frac{3}{5}\right)^2 \left(\frac{2}{5}\right)^2 \\ X = 3 & \quad {}^4C_3 \left(\frac{3}{5}\right)^3 \left(\frac{2}{5}\right)^1 \\ X = 4 & \quad {}^4C_4 \left(\frac{3}{5}\right)^4 \left(\frac{2}{5}\right)^0 \end{aligned}$$

0 red	1 red	2 red	3 red	4 red
••••	••••	••••	••••	••••
••••	••••	••••	••••	••••
••••	••••	••••	••••	••••
••••	••••	••••	••••	••••
2 2 2 2	2 2 2 3	2 2 3 3	2 3 3 3	3 3 3 3
5 5 5 5	5 5 5 5	5 5 5 5	5 5 5 5	5 5 5 5
4C_0	4C_1	4C_2	4C_3	4C_4

Now we've understood this in theory, but

How can we compute this in code?

We will use built-in functions of the `math.comb()` library.

```
import math
```

How will we find the value of 4C_0 ?

```
math.comb(4, 0)
```

```
1
```

As you can see this gave us the result of $\frac{4!}{0!*(4-0)!}$

Similarly, we can find 4C_1 as:

```
math.comb(4, 1)
```

```
4
```

Let's evaluate the probability values $P(X)$ for all possible values of $X = \{0, 1, 2, 3, 4\}$

```
# P(X=0)
math.comb(4,0)* (3/5)**0 * (2/5)**4
0.02560000000000005
```

```
# P(X=1)
math.comb(4,1)* (3/5)**1 * (2/5)**3
0.153600000000004
```

```
# P(X=2)
math.comb(4,2)* (3/5)**2 * (2/5)**2
```

```
0.3456000000000001
```

```
# P(X=3)
```

```
math.comb(4,3)* (3/5)**3 * (2/5)**1
```

```
0.3455999999999996
```

```
# P(X=4)
```

```
math.comb(4,4)* (3/5)**4 * (2/5)**0
```

```
0.1296
```

Let's compare these probability results to what we evaluated through the Empirical approach

Notice that these values are very close.

As discussed earlier, **if we increase the no of simulations, the observed result would be more and more closer to these theoretical values.**

Hence, proved.

```
[ ] pd.value_counts(red_values,normalize=True)
```

```
3    0.3552
2    0.3394
1    0.1496
4    0.1281
0    0.0277
dtype: float64
```

▼ Binomial Distribution

You might not be aware, but while solving this Casino case study, we've also been deriving the equation for **Binomial Distribution**.

Let's summarize our findings, and look at this distribution formally.

Binomial distribution is a **discrete probability distribution** of the number of successes in n **independent** experiments sequence.

A Binomial trial will always have **two possible outcomes**:

- Success / Win
- Failure / Loss

We defined a **discrete random variable** X that denoted number of red balls drawn.

- Note that the event of drawing a ball is independent.
- X will be called a **Binomial RV**

Also, we were given some parameters in our problem, let's define them:

- n : No of independent trials
 - In our example, we draw balls 4 times, hence $n = 4$
- p : Probability of success in one trial
 - In our example, this denotes the probability of drawing a red ball, hence $p = \frac{3}{5}$
 - Therefore, $(1 - p)$ becomes the probability of failure in each trial (i.e. drawing a blue ball, in this example)

Using these parameters, we can re-write the equation we derived in general form: $P(X = k) = {}^nC_k (p)^k (1 - p)^{n-k}$

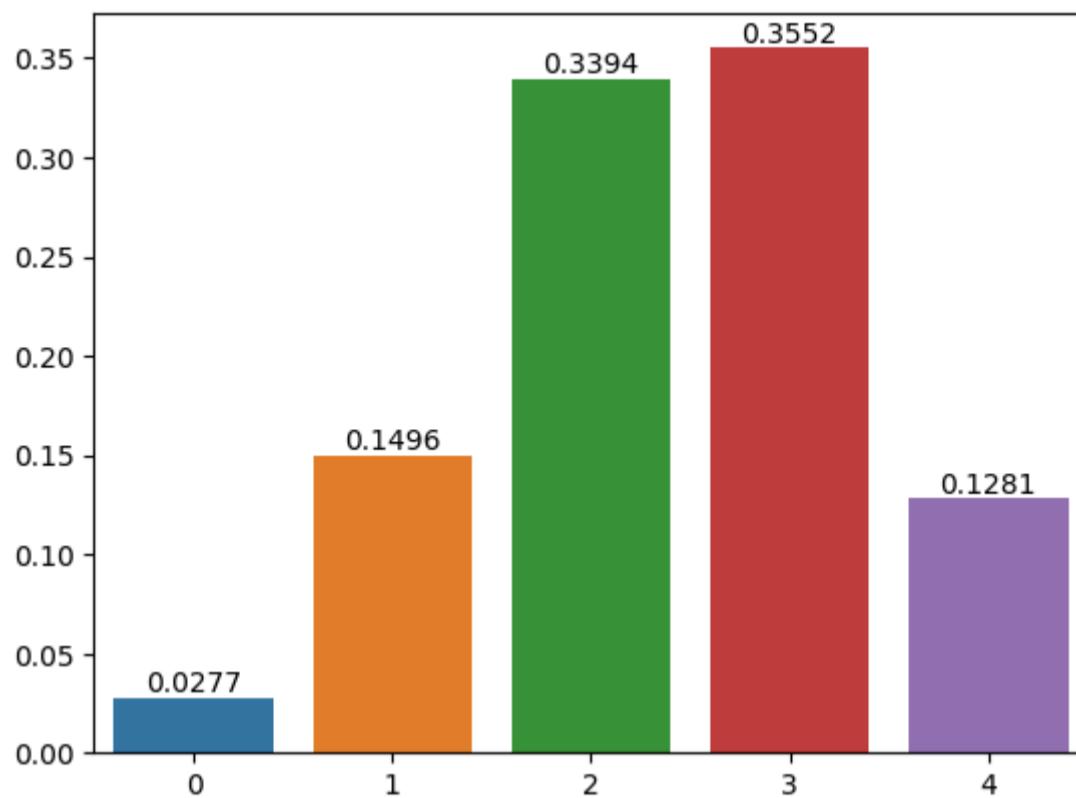
Let's plot our calculated values to see what Binomial distribution looks like.

```
x = pd.value_counts(red_values, normalize=True)
x
```

```
3    0.3552
2    0.3394
1    0.1496
4    0.1281
```

```
0    0.0277  
dtype: float64
```

```
ax = sns.barplot(x = x.index, y = x.values)  
  
for i in ax.containers:  
    ax.bar_label(i,)
```



This is the **Probability Mass Function (PMF)** of our given Binomial experiment, which is called as Binomial Probability Distribution

- The graph shows the probability of obtaining each possible number of successes (k) in n trials.
- The height of each bar represents the probability of that particular outcome.
- The sum of all the probabilities equals 1.

The `scipy.stats.binom` library gives us a built in function that eases the calculation of PMF values (i.e. value of $P(X)$ for specific values of X).

Instead of using the formula $P(X = k) = {}^nC_k (p)^k (1 - p)^{n-k}$, we can directly use this function to calculate the PMF value.

We just need to specify the 3 parameters:

- n
- k
- p

```
from scipy.stats import binom
```

```
prob_0_red = binom.pmf(n=4,p=3/5,k=0)  
prob_0_red
```

```
0.02559999999999994
```

```
prob_1_red = binom.pmf(n=4,p=3/5,k=1)  
prob_1_red
```

```
0.1535999999999996
```

```
prob_2_red = binom.pmf(n=4,p=3/5,k=2)  
prob_2_red
```

```
0.3456
```

```
prob_3_red = binom.pmf(n=4,p=3/5,k=3)  
prob_3_red
```

```
0.3456000000000001
```

```
prob_4_red = binom.pmf(n=4,p=3/5,k=4)  
prob_4_red
```

Notice that these values are the same as what we calculated using `math.comb`

❖ Expectation using theoretical approach

How will we calculate the theoretical expectation value?

We know the formula: $E(X) = \sum_i X_i P(X = X_i)$

- Here, we saw that we can calculate the probability values using `scipy.stats.binom`
- And that random variable $X = \{0, 1, 2, 3, 4\}$

```
expectation_theoretical= (0*prob_0_red) + (1*prob_1_red) + (2*prob_2_red) + (3*prob_3_red) + (4*prob_4_red)
expectation_theoretical
```

2.4000000000000004

Note that this is close to the **Empirical Expected value** we calculated.

Alternately, there is a built-in function to find this expected value in `stats.binom`

Here, we need to pass the following arguments to `args`:

- n , and
- p

```
binom.expect(args=(4,3/5))
```

2.4000000000000004

Variance in Binomial Distribution

Recall that variance tells you how much the actual results might vary from the expected average (mean), helping you understand whether your observations are likely due to chance or if there's something more going on, like bias in the coin.

Formula for Variance in Binomial Distribution:

The formula for variance in a binomial distribution is:

$$\sigma^2 = n * p * (1 - p)$$

Where:

- n is the number of trials (or coin flips in our example).
- p is the probability of success on each trial (the probability of getting heads in our coin flip example).
- $(1 - p)$ represents the probability of failure on each trial (the probability of getting tails).
 - $(1 - p)$ is also denoted as q also, so

$$\sigma^2 = n * p * (1 - p) \text{ or}$$

$$\sigma^2 = npq$$

We learnt about the concept of Binomial distribution.

But we still haven't answered our question :

❖ Would engaging in this game result in a profit or loss for you?

Let's define another random variable Y that denotes the amount of money won/lost through gambling.

- Therefore, possible values of $Y : \{150, -10\}$

Let's create a table for this random variable Y , with its possible values and corresponding probabilities.

- Case of winning Rs 150 ($Y = 150$)
 - $P(Y = 150)$ would be the same as $P(X = 4)$
- Case of loosing Rs 10 ($Y = -10$)
 - $P(Y = -10) = P(X = 0) + P(X = 1) + P(X = 2) + P(X = 3) = 1 - P(X = 4)$

```
# P(Y=150)
prob_4_red
```

0.1296

```
#P(Y = -10)
1 - prob_4_red
```

0.8704000000000001

What would be expected value of Y ?

$$E(Y) = \sum_i Y_i P(Y = Y_i) = (150 * 0.1296) + (-10 * 0.8704000000000001)$$

```
expected_y = (150*0.1296) + (-10*0.8704000000000001)
expected_y
```

10.73599999999997

Conclusion of the case study:

This value means that if we play many many times, at the end of the day, we are expected to have profit of Rs 10.736

Casino case study



A bag has 3 red and 2 blue balls.

You pick a ball, write its colour, and put it back in the bag. This is done 4 times in total.

If all 4 times, the red ball was drawn, you win Rs 150. In any other case, you lose Rs 10.

Would you play this game?

What are all the outcomes?

0 red	1 red	2 red	3 red	4 red
● ● ● ●	● ● ● ○	● ● ○ ○	● ○ ○ ○	○ ○ ○ ○
● ● ● ○	● ● ○ ○	● ○ ○ ○	○ ○ ○ ○	
● ● ○ ○	● ○ ○ ○	○ ○ ○ ○		
● ○ ○ ○	○ ○ ○ ○			
○ ○ ○ ○				

2 2 2 2	2 2 2 3	2 2 3 3	2 3 3 3	3 3 3 3
5 5 5 5	5 5 5 5	5 5 5 5	5 5 5 5	5 5 5 5
4C_0	4C_1	4C_2	4C_3	4C_4

Let “ X ” denote the number of red balls when you draw 4 balls with replacement
Here, X is an example of what is called a “Random Variable”

Let “ Y ” be the amount won. This is also another example of a random variable

What are all the outcomes for “ Y ”?

“ $Y = 150$ ” If we get 4 red balls

“ $Y = -10$ ” Otherwise

Y	$P[Y]$	
150	${}^4C_4 \left(\frac{3}{5}\right)^4$	0.1296
-10	${}^4C_0 \left(\frac{2}{5}\right)^4 + {}^4C_1 \left(\frac{2}{5}\right)^3 \left(\frac{3}{5}\right)^1 + {}^4C_2 \left(\frac{2}{5}\right)^2 \left(\frac{3}{5}\right)^2 + {}^4C_3 \left(\frac{2}{5}\right)^1 \left(\frac{3}{5}\right)^3$	0.8704

$$E[Y] = (150)(0.1296) + (-10) * (0.8704) = 10.736$$

Conditions of Binomial Experiment

1. The experiment must consist of a **fixed number of trials (n)**, with only 2 possible outcomes: Success or Failure

- Here, we had fixed $n = 4$
- So, we cannot increase or decrease it in between
- We defined success as the event of drawing a red ball: $P(\text{Success}) = \frac{3}{5}$

2. Individual trials are **identical and independent**.

- This needs to hold true, otherwise, the probability values might change for different trials.

- In this example also, the trials of drawing balls were identical and independent, as we were replacing the balls after each draw.
- Hence, it contained exactly same number of balls of each color.

3. The random variable denotes the number of success in n trials.

→ Binomial Random Variable

$n \rightarrow$ no. of trials

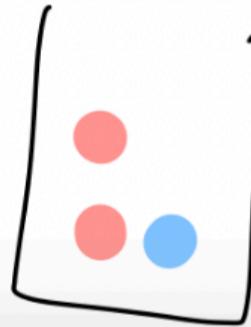
$p \rightarrow$ probability of success in one trial -

$$P[X=k] = {}^n C_k p^k (1-p)^{n-k}$$

① Find no. of trials → Success (3/5)
 (4) → Failure (2/5)

↳ Bernoulli

② Trials are identical & independant,
 p remains same from trial to trial

[] ③ Random Variable denotes the no. of success in n trials.

▼ Bernoulli Trials

In the above example, it was conveyed clearly that we will draw a ball 4 times from the bag.

But let's consider the case when a ball is drawn from the bag, only one, i.e. **one trial**.

Like before, we still define X as the event of getting a red ball.

Therefore, on drawing the ball, we get 2 possibilities:

- Getting a red ball (Success)
 - We know that probability for this will be: $P(\text{Success}) = p = \frac{3}{5}$
- Getting a blue ball (Failure)
 - Probability: $P(\text{Failure}) = 1 - p = \frac{2}{5}$

This is known as a **Bernoulli Trial**.

Essentially, it is the **special case** of Binomial trial, where $n = 1$

Hence, it must also follow the condition that there must be only 2 possible outcomes:

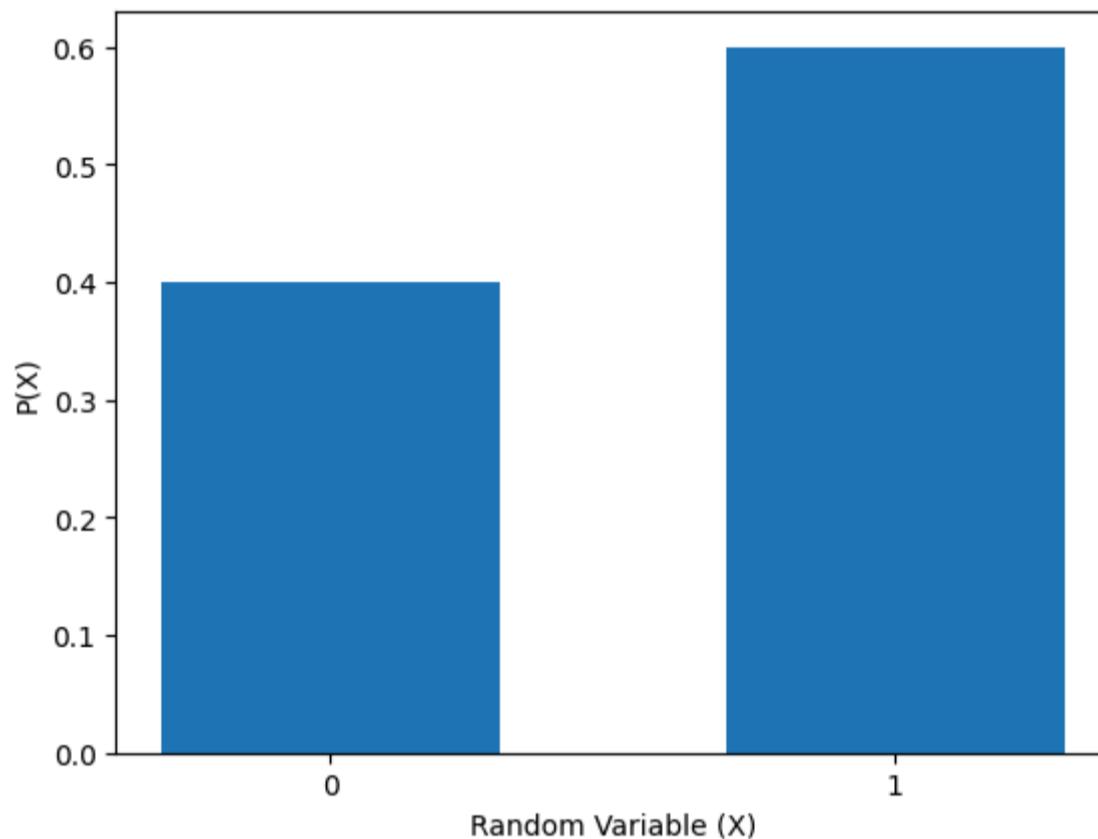
- Success, or
- Failure

Let's plot this, to see what **Bernoulli distribution** looks like.

```
x = [0, 1]
y = [2/5, 3/5]

plt.bar(x, y, width=0.6, tick_label=["0", "1"])
plt.xlabel("Random Variable (X)")
plt.ylabel("P(X)")

Text(0, 0.5, 'P(X)')
```



Another example of Bernoulli distribution

Consider the situation of passing or failing an exam.

- Let's assume the probability to pass the exam is 95%,
- Therefore the probability to fail will be 5%.

In this case, if the event to pass the exam is considered, then the Bernoulli event will contain the probability of passing the exam.

Similarly, it goes for failing the exam.

To summarize, what is the difference between Binomial and Bernoulli distribution?

- Bernoulli deals with the outcome of the single trial of the event, whereas Binomial deals with the outcome of the multiple trials of the single event.
- Hence, we can define **Binomial distribution** in another way:

It is the collection of Bernoulli trials for the same event, i.e., it contains more than 1 Bernoulli event for the same scenario for which the Bernoulli trial is calculated.

Let's take a look at another example problem

Dice Example

You toss 2 dice. If both dice are 6, you get Rs 2. Else, if one dice is 6, you get Rs 1. Otherwise, you do not get anything.

Let's define a random variable X that represents the amount of money won.

- Hence, it can take the values: $X = \{0, 1, 2\}$

Answer the following questions.

What is the probability of getting the following?

- Rs 0
- Rs 1
- Rs 2

Let's visualize the possible outcomes through a table.

- Possible values of Dice 1 along the row
- Possible values of Dice 2 along the column
- Value corresponding to a row and col, represents the money won (X)

		D_2	1	2	3	4	5	6		X	$P(X)$
		# of 6	1	2	3	4	5	6		0	${}^2C_0 \left(\frac{5}{6}\right)^2$
	1	1	0	0	0	0	0	1	1	0	${}^2C_0 \left(\frac{5}{6}\right)^2$
	2	2	0	0	0	0	0	1	1	1	${}^2C_1 \left(\frac{1}{6}\right) \left(\frac{5}{6}\right)$
	3	3	0	0	0	0	0	1	1	2	${}^2C_2 \left(\frac{1}{6}\right)^2$
	4	4	0	0	0	0	0	1	1		
	5	5	0	0	0	0	0	1	1		
	6	6	1	1	1	1	1	2	2		

$\frac{5 * 5}{36}$

$\frac{5 * 1 + 1 * 5}{36}$

$\frac{1 * 1}{36}$

Finding $P(X = 0)$

From the table we can see that we will get 0 Rs for 25 outcomes, hence $P(X = 0) = \frac{25}{36}$

Finding $P(X = 1)$

From the table, $P(X = 1) = \frac{10}{36}$

Finding $P(X = 2)$

From the table, $P(X = 2) = \frac{1}{36}$

Now, let's see if we can obtain the same answers using the Binomial formula

Before we get to solving, let's define the parameters:

- **What will be the value of n?**
 - Since we are throwing 2 dice, $n = 2$
- **What will be the value of p?**
 - p is defined as the probability of success in one trial
 - So how do we define success here?
 - Obtaining a 6
 - Therefore, $p = \text{probability of getting a 6 in a single dice roll, i.e. } p = \frac{1}{6}$

We know the Binomial formula is: $P(X = k) = {}^nC_k (p)^k (1 - p)^{n-k}$

Therefore,

- $P(X = 0) = {}^2C_0 \left(\frac{1}{6}\right)^0 \left(\frac{5}{6}\right)^2 = 1 * 1 * \frac{25}{36} = \frac{25}{36}$
- $P(X = 1) = {}^2C_1 \left(\frac{1}{6}\right)^1 \left(\frac{5}{6}\right)^1 = 2 * \frac{1}{6} * \frac{5}{6} = \frac{10}{36}$
- $P(X = 2) = {}^2C_2 \left(\frac{1}{6}\right)^2 \left(\frac{5}{6}\right)^0 = 1 * \frac{1}{36} * 1 = \frac{1}{36}$

These are the exact answers we got using the table above!!

Alternately, we could've evaluated the binomial formula using code as:

```
binom.pmf(n=2, p=1/6, k=0)  
0.6944444444444443
```

Now answer the second question.

What is the expected value of money won?

We can find this using the formula: $E(X) = \sum_i X_i P(X = X_i)$

$$\begin{aligned} &= (0 * \frac{25}{36}) + (1 * \frac{10}{36}) + (2 * \frac{1}{36}) \\ &= \frac{1}{3} \end{aligned}$$

Alternately, we can use the `stats.binom.expect()` function

```
binom.expect(args=(2,1/6))  
0.3333333333333326
```