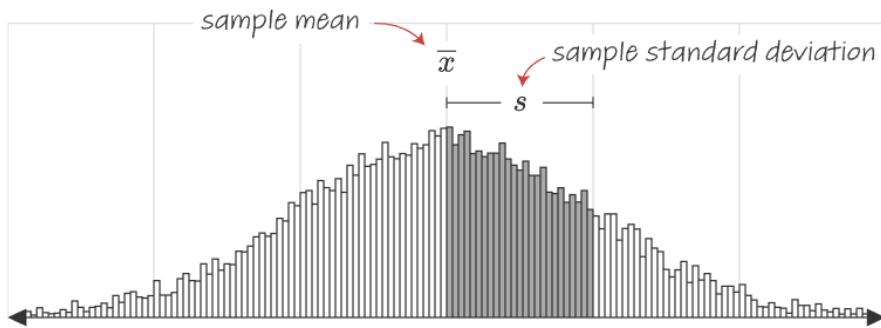


Sample Standard Deviation



$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

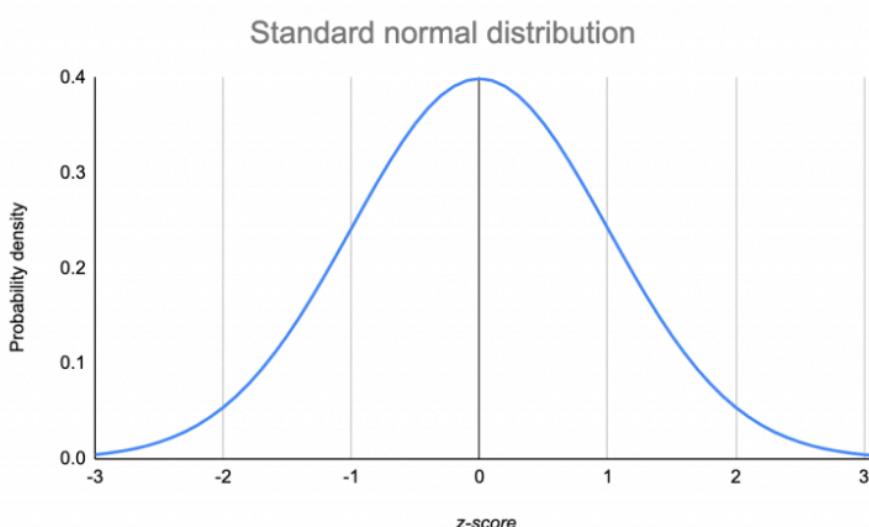
- The value for the standard deviation describes how closely the data set is to the mean. One standard deviation away from the mean on either side contains approximately 68.3% of the samples, two standard deviations contains approximately 95.4% of the samples, and so on. Because this formula calculates an approximate value for the true standard deviation there will be some discrepancy between the actual normal distribution and the one modeled by the sample standard deviation.

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$$

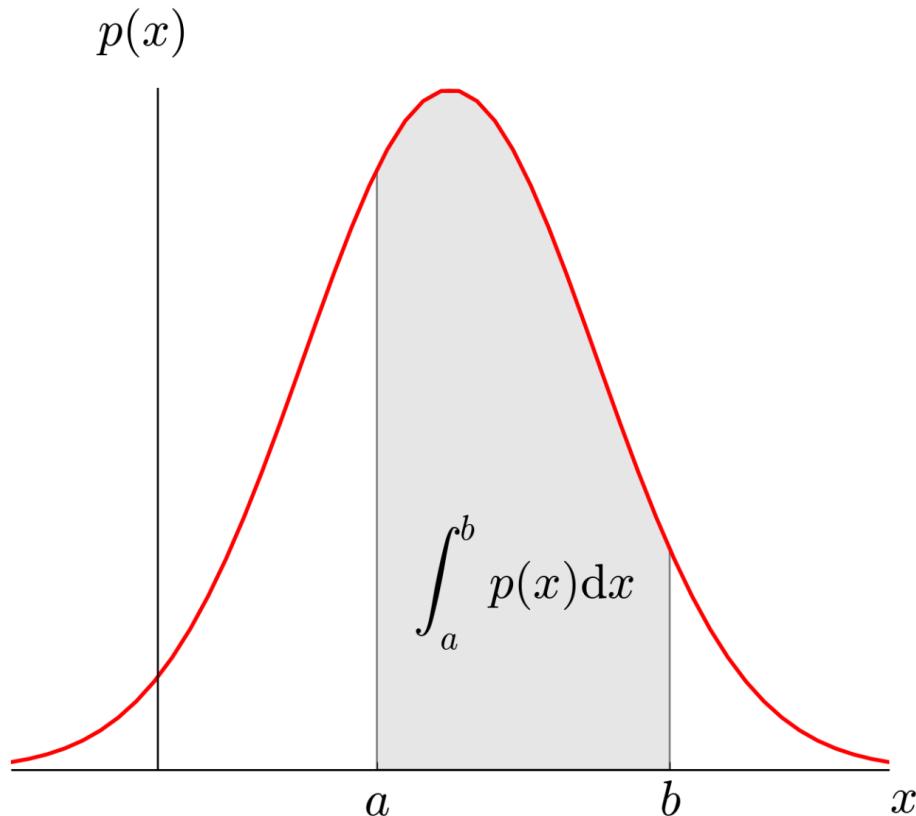
The Standard Normal Distribution

- A standard normal random variable is a normally distributed random variable with mean $\mu = 0$ and standard deviation $\sigma = 1$. It will always be denoted by the letter Z.



The Probability Density Function(PDF)

- defines the probability function representing the density of a continuous random variable lying between a specific range of values. In other words, the probability density function produces the likelihood of values of the continuous random variable. Sometimes it is also called a probability distribution function or just a probability function.



Probability Density Function Properties

Let x be the continuous random variable with density function $f(x)$, and the probability density function should satisfy the following conditions:

- For a continuous random variable that takes some value between certain limits, say a and b , the PDF is calculated by finding the area under its curve and the X-axis within the lower limit (a) and upper limit (b). Thus, the PDF is given by

$$P(x) = \int_a^b f(x) dx$$

- The probability density function is non-negative for all the possible values, i.e. $f(x) \geq 0$, for all x .
- The area between the density curve and horizontal X-axis is equal to 1, i.e.
$$\int_{-\infty}^{\infty} f(x) dx = 1$$
- Due to the property of continuous random variables, the density function curve is **continued** for all over the given range. Also, this defines itself over a range of continuous values or the domain of the variable.

In []:

1

What is a Cumulative Distribution Function?

- The Cumulative Distribution Function (CDF), of a real-valued random variable X , evaluated at x , is the probability function that X will take a value less than or equal to x . It is used to describe the probability distribution of random variables in a table. And with the help of these data, we can easily create a CDF plot in an excel sheet.
- In other words, CDF finds the cumulative probability for the given value. To determine the probability of a random variable, it is used and also to compare the probability between values under certain conditions. For discrete distribution functions, CDF gives the probability values till what we specify and for continuous distribution functions, it gives the area under the probability density function up to the given value specified.

Cumulative Distribution Function Formula

The CDF defined for a discrete random variable and is given as

$$F_X(x) = P(X \leq x)$$

Where X is the probability that takes a value less than or equal to x and that lies in the semi-closed interval $(a,b]$, where $a < b$.

Therefore the probability within the interval is written as

$$P(a < X \leq b) = F_X(b) - F_X(a)$$

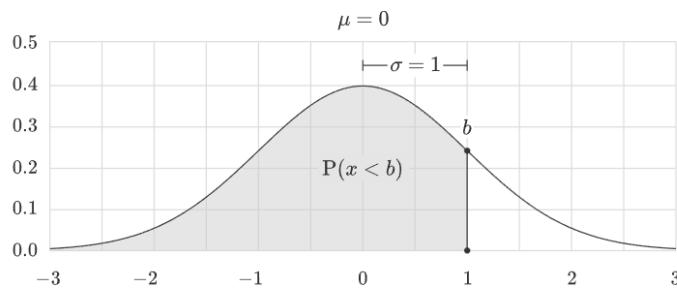
The CDF defined for a continuous random variable is given as;

$$F_X(x) = \int_{-\infty}^x f_X(t)dt$$

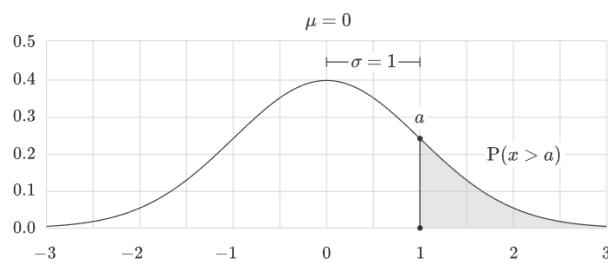
Here, X is expressed in terms of integration of its probability density function f_X .

In case, if the distribution of the random variable X has the discrete component at value b ,

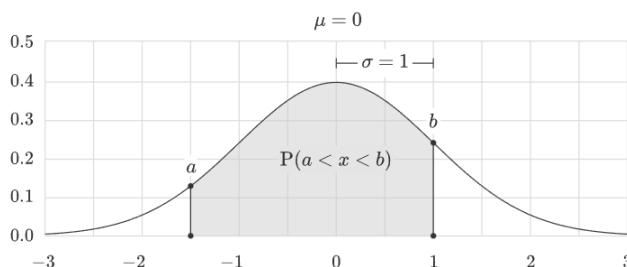
$$P(X = b) = F_X(b) - \lim_{x \rightarrow b^-} F_X(x)$$



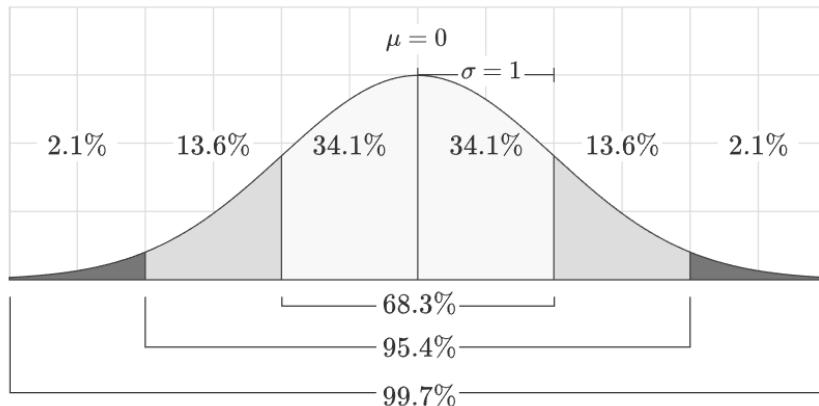
$$P(x > a) = 1 - P(x < a)$$



$$P(a < x < b) = P(x < b) - P(x < a)$$



Probability	Integral	Description
$P(x < b)$	$\int_{-\infty}^b f(x)dx$	The probability of an event occurring below a threshold b .
$P(x > a)$	$\int_a^{\infty} f(x)dx$	The probability of an event occurring above a threshold a .
$P(a < x < b)$	$\int_a^b f(x)dx$	The probability of an event occurring between a and b .



In [2]: 1 from scipy.stats import norm

In [3]: 1 norm.cdf(1)

Out[3]: 0.8413447460685429

In [4]: 1 norm.cdf(-1)

Out[4]: 0.15865525393145707

In [6]: 1 (norm.cdf(1) - norm.cdf(-1))*100 # $P[x = -1 \text{ to } x = 1] = 68.26\%$ mean +/- 1 σ

Out[6]: 68.26894921370858

In [7]: 1 (norm.cdf(2) - norm.cdf(-2))*100 # $P[x = -2 \text{ to } x = 2] = 95.44\%$ mean +/- 2 σ

Out[7]: 95.44997361036415

In [8]: 1 (norm.cdf(3) - norm.cdf(-3))*100 # $P[x = -3 \text{ to } x = 3] = 99.73\%$ mean +/- 3 σ

Out[8]: 99.73002039367398

Z - Scores (Standardisation)

$$z = \frac{x - \mu}{\sigma}$$

A Z-score is a numerical measurement that says , how far data point x is away from mean value.

It's a standardized measure

If you're asked to find standardized values, use this formula to make your calculations:
Standardized Values:

calculate a standardized value (a z-score), using the above formula. The symbols are:

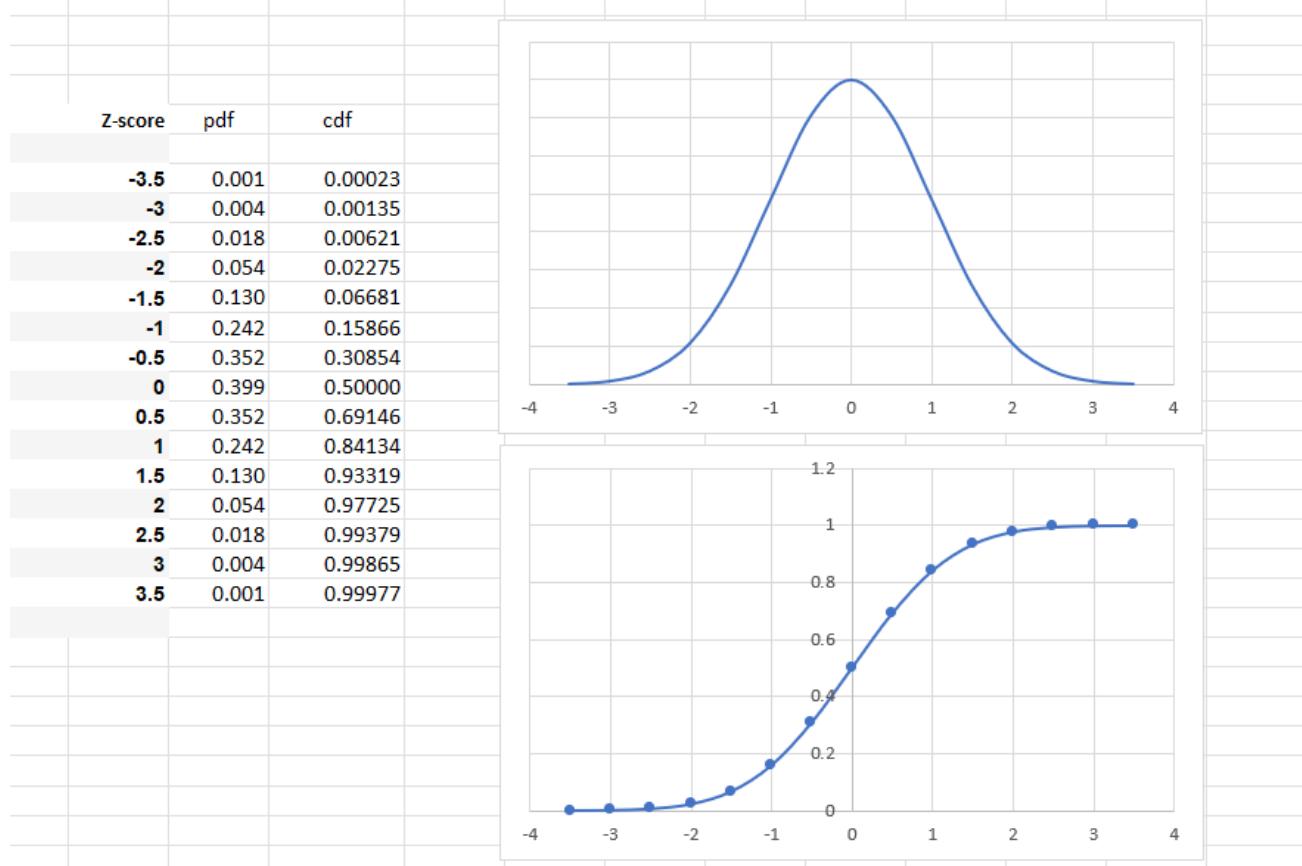
X: the observation (a specific value that you are calculating the z-score for).

Mu(μ): the mean.

Sigma(σ): the standard deviation.

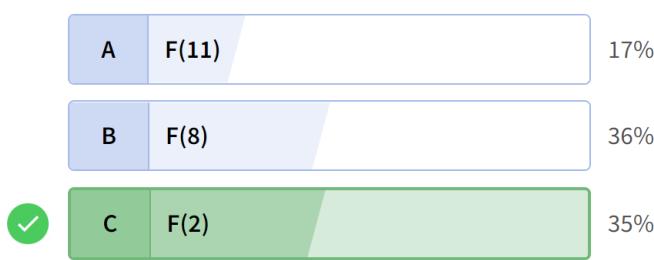
In []:

1



Let X be a Gaussian random variable with mean 3 and variance 16, and let F denote the CDF of a "standard" Gaussian. What is P(X < 11)

72 users have participated



```
1 μ = 3
2 σ = squar_root(16) = 4
3
4 P[X<=11] = (11-3)/4
5      = 8 / 4
6      = 2
7
```

```
In [16]: 1 norm.cdf(2) # cummulative probability of z = 2
```

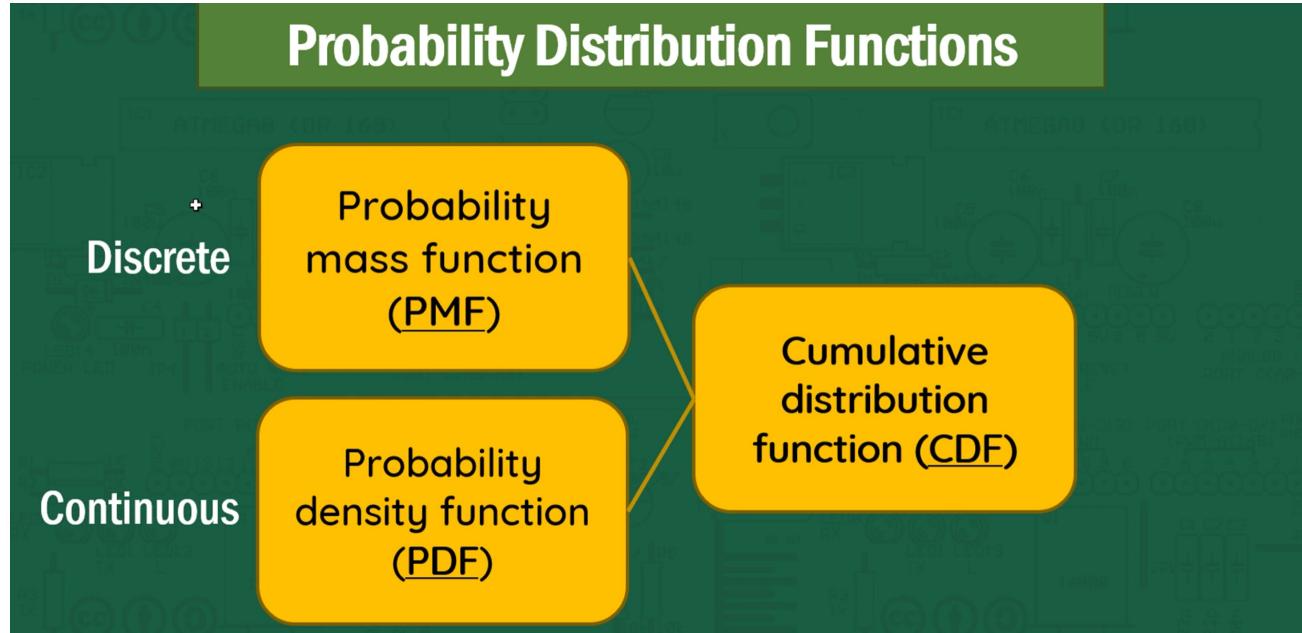
```
Out[16]: 0.9772498680518208
```

```
In [ ]: 1 # probability of 11 when, mean = 3 , varinace = 16
```

```
In [15]: 1 norm.cdf(11 , loc = 3 , scale = 4) # P[x<=11]
2 # = P[(x-3)/4 <= (11-3)/4]
3 # = P[(x-3)/4 <= (8)/4]
4 # = P[Z <= 2]
```

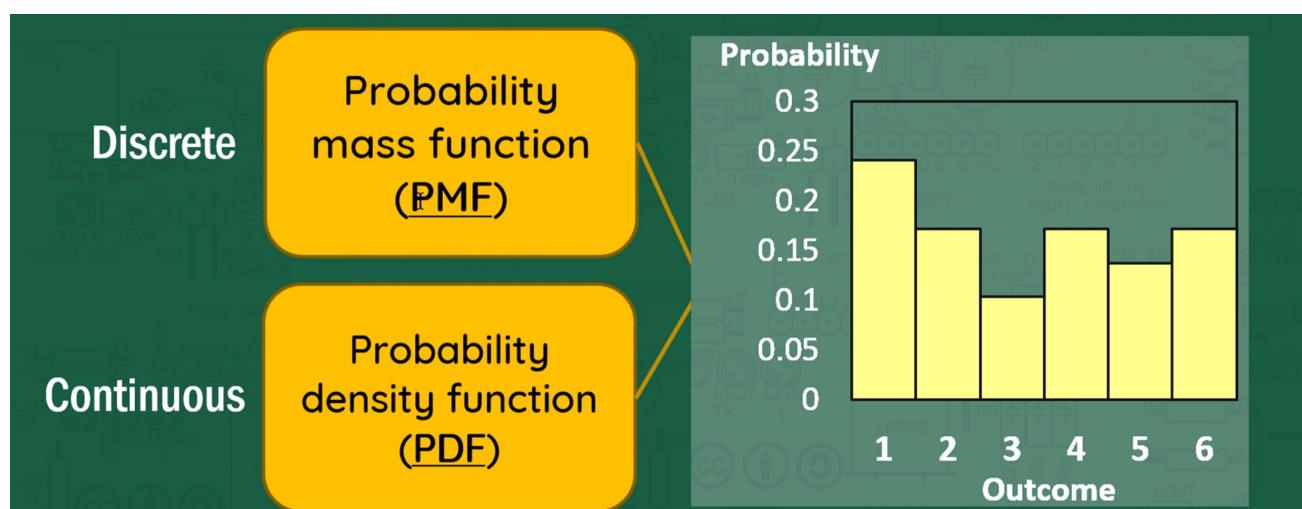
```
Out[15]: 0.9772498680518208
```

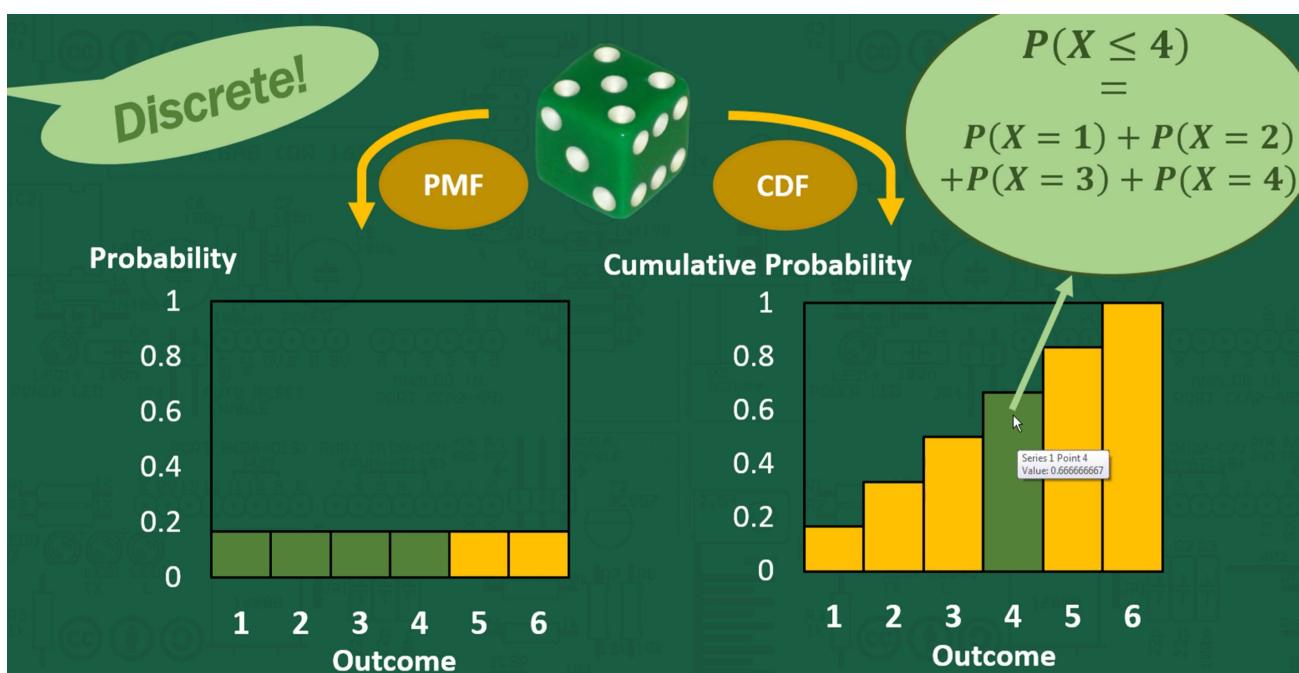
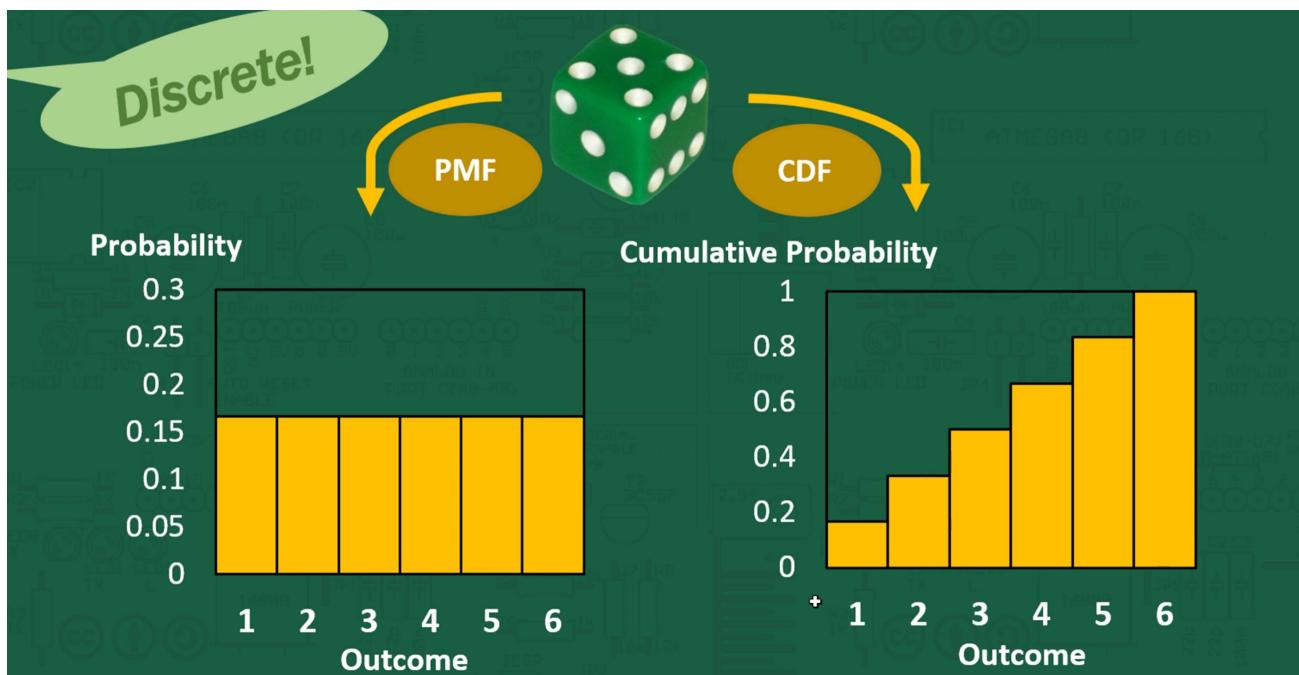
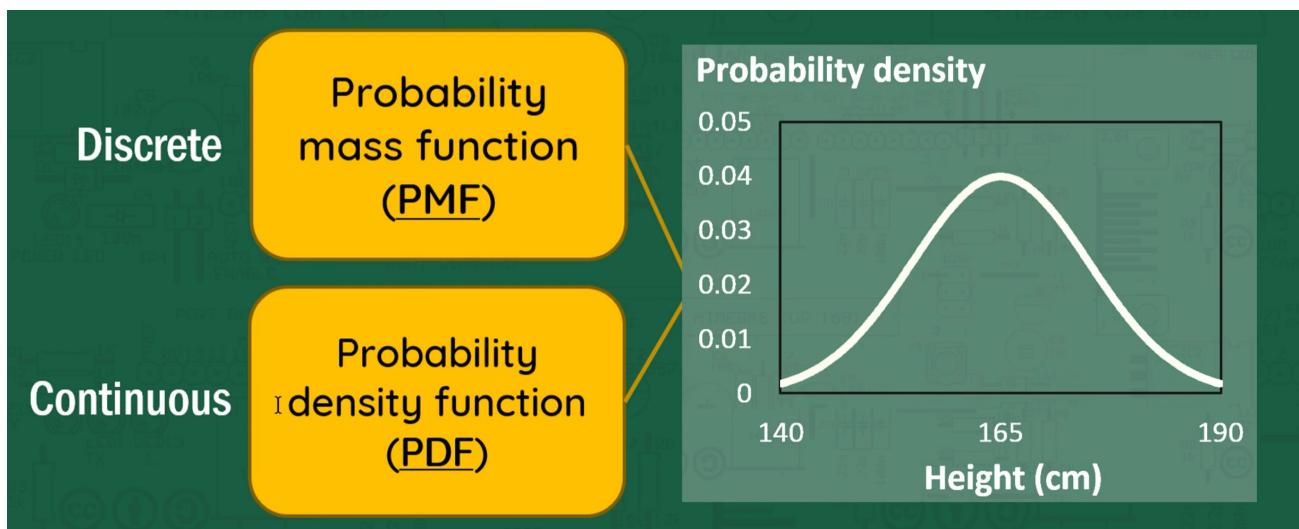
```
In [ ]: 1 # snapshots from : https://www.youtube.com/playlist?list=PLTNMv857s9WVzutwxaMb0YZKW7hoveGLS
```



PMF : probability of each outcomes (discrete random variable)

PDF : probability of continuous random variable

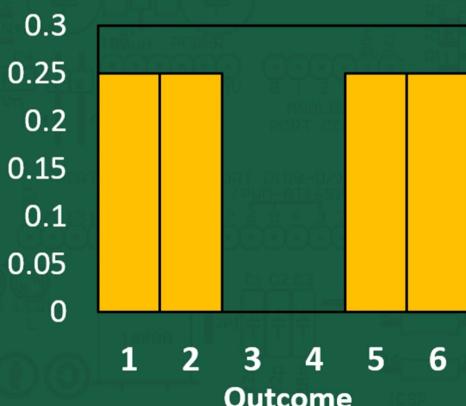




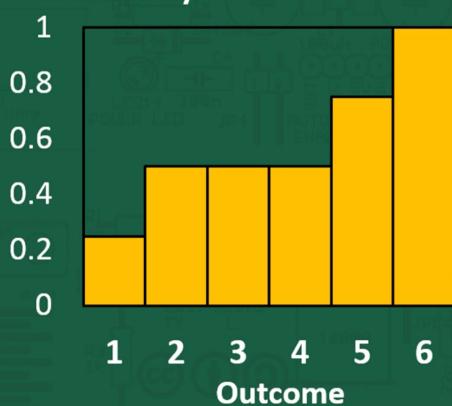
Discrete!



Probability



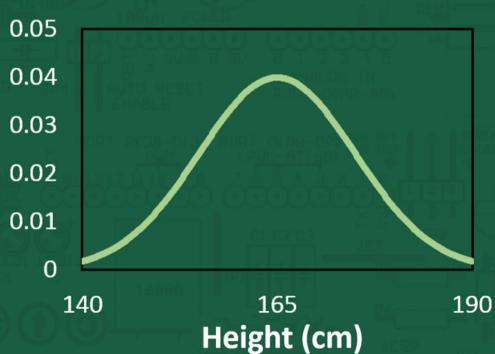
Cumulative Probability



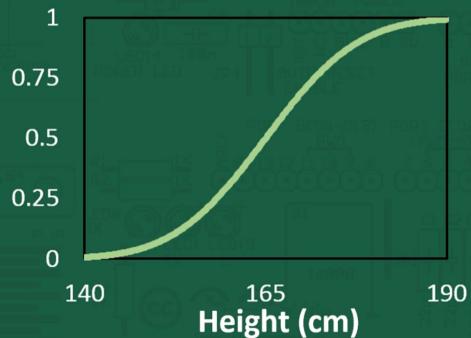
Continuous!



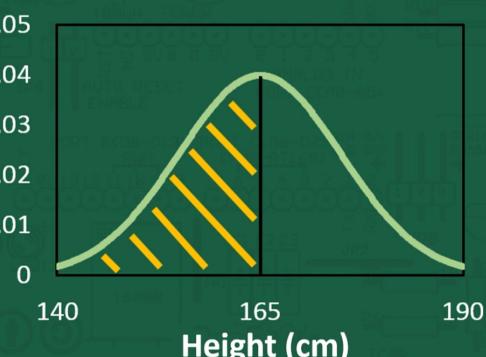
Probability density



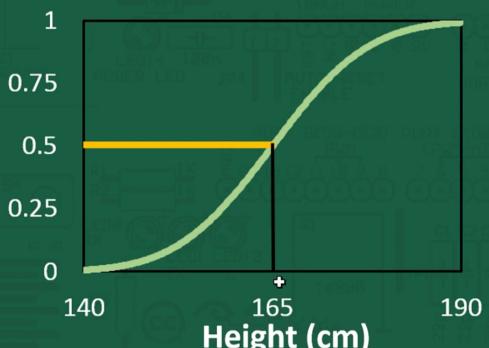
Cumulative probability



Probability density



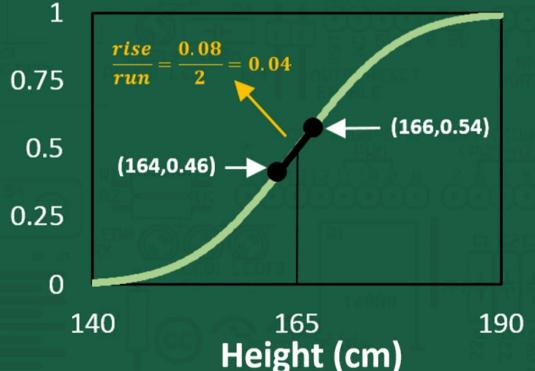
Cumulative probability



Probability density

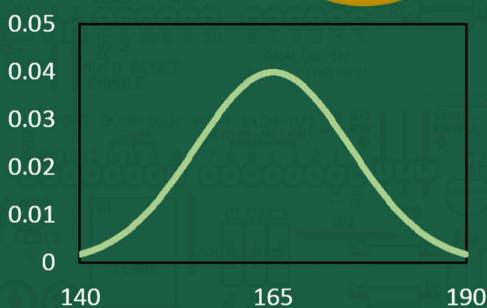


Cumulative probability



Probability density

$f(x)$



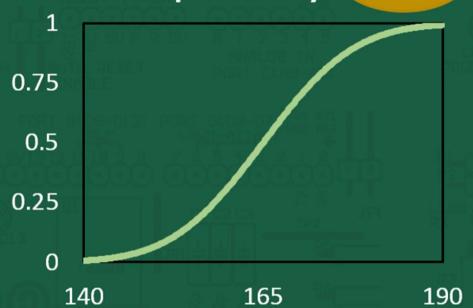
Cumulative probability

$F(x)$

Gradient
Area to the left

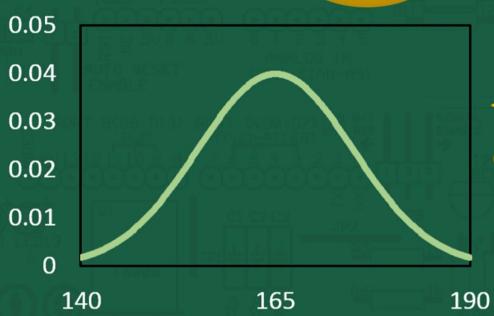
Cumulative probability

$F(x)$



Probability density

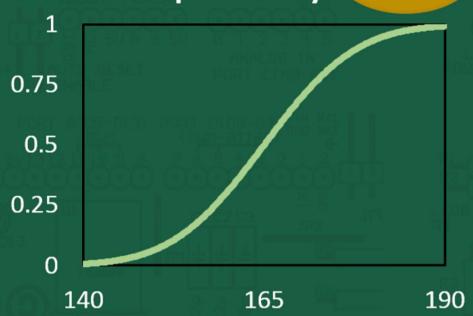
$f(x)$



Cumulative probability

$F(x)$

$$\int_{-\infty}^x f(x) dx = F(x)$$



In []:

1

In []:

1

In []:

1

Variance and its sq-root, the standard deviation are both measures of the spread or variability in data. Two or more data set could have the same mean , but very different variances .

variance is based on squared difference between each observation and mean .

the average squared deviation : How far is each data point from the center (mean) on average

Variance

- Variance is the expected value of the squared variation of a random variable from its mean value, in probability and statistics. Informally, variance estimates how far a set of numbers (random) are spread out from their mean value.
- The value of variance is equal to the square of standard deviation, which is another central tool.
- Variance is symbolically represented by σ^2 , s^2 , or $Var(X)$.

In []:

1

1 Expected value of variance :
 2 $\text{Var}(X) = E[(X - \mu)^2]$

1 a variance is a measure of how far a set of data (numbers) are spread out from their mean (average) value.
 2
 3 Variance means to find the expected difference of deviation from actual value. Therefore, variance depends on the standard deviation of the given data set.
 4
 5 The more the value of variance, the data is more scattered from its mean and if the value of variance is low or minimum, then it is less scattered from mean. Therefore, it is called a measure of spread of data from mean.
 6
 7 For the purpose of solving questions, the formula for variance is given by:
 8
 9 $\text{Var}(X) = E[(X - \mu)^2]$
 10
 11
 12 Put into words; this means that variance is the expectation of the squared deviation of a random set of data from its mean value. Here,
 13
 14 X = Random variable
 15
 16 “ μ ” is equal to $E(X)$ so the above equation may also be expressed as,
 17
 18 $\text{Var}(X) = E[(X - E(X))^2]$
 19
 20 $\text{Var}(X) = E[X^2 - 2X E(X) + (E(X))^2]$
 21
 22 $\text{Var}(X) = E(X^2) - 2 E(X) E(X) + (E(X))^2$
 23
 24 $\text{Var}(X) = E(X^2) - (E(X))^2$
 25
 26 Sometimes the covariance of the random variable itself is treated as the variance of that variable. Symbolically,
 27
 28 $\text{Var}(X) = \text{Cov}(X, X)$

As we know already, the variance is the square of standard deviation, i.e.,

$$\text{Variance} = (\text{Standard deviation})^2 = \sigma^2$$

The corresponding formulas are hence,

$$\text{Population standard deviation } \sigma =$$

$$\sqrt{\frac{\sum(X-\mu)^2}{N}}$$

and Sample standard deviation s =

$$\sqrt{\frac{\sum(x-x)^2}{n-1}}$$

Where X (or x) = Value of Observations

μ = Population mean of all Values

n = Number of observations in the sample set

\bar{x}

= Sample mean

N = Total number of values in the population

Calculate the variance of a random variable X which is the outcome of a fair dice: {1, 2, 3, 4, 5, 6}

57 users have participated



```

1 Var(X) = E[ ( X - E(X) )^2 ]
2   = E(X^2) - (E(X))^2
3
4   X = {1,2,3,4,5,6}
5
6   E(X) = 3.5 = 7/2
7
8   X^2 = {1,4,9,16,25,36}
9   E(X^2) = (1*(1/6)) + (4*(1/6)) + (9*(1/6)) + (16*(1/6)) + (25*(1/6)) + (36*(1/6))
10
11  E(x^2) = 15.166 = 91/6
12
13  Var(X) = E[ ( X - E(X) )^2 ]
14   = E(X^2) - (E(X))^2
15   = 91/6 - (7/2)^2
16   = 2.91
17   = 35/12

```

In [17]: 1 $(1*(1/6)) + (4*(1/6)) + (9*(1/6)) + (16*(1/6)) + (25*(1/6)) + (36*(1/6))$

Out[17]: 15.166666666666666

In [19]: 1 $(91/6) - ((7/2)^2)$

Out[19]: 2.9166666666666666

In [20]: 1 $35/12$

Out[20]: 2.9166666666666665

In []: 1

In []: 1 Expected values of variance : $\Sigma ((x - E(x))^2 * P(x))$

In []: 1

Expected Value and Variance of a Discrete Random Variable

Expected Value (or mean) of a Discrete Random Variable

For a discrete random variable, the expected value, usually denoted as μ or $E(x)$, is calculated using:

$$\mu = E(X) = \sum x_i f(x_i)$$

The formula means that we multiply each value, x , in the support by its respective probability, $f(x)$, and then add them all together. It can be seen as an average value but weighted by the likelihood of the value.

```
1 | we were given the following discrete probability distribution:  
2 |  
3 | 0   1   2   3   4  
4 | 1/5 1/5 1/5 1/5 1/5  
5 | What is the expected value?  
6 |  
7 |
```

$$\begin{aligned}\mu = E(X) &= \sum x f(x) = 0 \left(\frac{1}{5}\right) + 1 \left(\frac{1}{5}\right) + 2 \left(\frac{1}{5}\right) + 3 \left(\frac{1}{5}\right) + 4 \left(\frac{1}{5}\right) \\ &= 2\end{aligned}$$

$$\text{Var}(X) = \left[0^2 \left(\frac{1}{5}\right) + 1^2 \left(\frac{1}{5}\right) + 2^2 \left(\frac{1}{5}\right) + 3^2 \left(\frac{1}{5}\right) + 4^2 \left(\frac{1}{5}\right) \right] - 2^2 = 6 - 4 = 2$$

$$\text{SD}(X) = \sqrt{2} \approx 1.4142$$

Variance of a Discrete Random Variable

$$\sigma^2 = \text{Var}(X) = \sum (x_i - \mu)^2 f(x_i)$$

The formula means that we take each value of x , subtract the expected value, square that value and multiply that value by its probability. Then sum all of those values.

There is an easier form of this formula we can use.

$$\sigma^2 = \text{Var}(X) = \sum x_i^2 f(x_i) - E(X)^2 = \sum x_i^2 f(x_i) - \mu^2$$

The formula means that first, we sum the square of each value times its probability then subtract the square of the mean. We will use this form of the formula in all of our examples.

Type *Markdown* and *LaTeX*: α^2

In []:

1

Covariance

```
1 | Covariance is a measure of the relationship between two random variables and to what extent,  
they change together. Or we can say, in other words, it defines the changes between the two  
variables, such that change in one variable is equal to change in another variable.
```

cov(x,y) = covariance between x and y variable

xi = data value of x yi = data value of x

x-bar = mean of x y-bar = mean of y N = number of data values

$$cov_{x,y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N - 1}$$

Type *Markdown* and *LaTeX*: α^2

In []:

1	
---	--

In []:

1	
---	--

In []:

1	
---	--

In []:

1	
---	--

In []:

1	
---	--

Joint Distribution :

```
1 coin toss :  
2   dice :  
3  
4       h1 h2 h3 h4 h5 h6  
5       t1 t2 t3 t4 t5 t6  
6   X coin : {h,t} {h=1,t=0}  
7   Y dice : {1,2,3,4,5,6}  
8  
9       P[x = 0 , y = 3] = 1/12  
10  
11      P[y = 3] = ({y=3} and {x=0}) U ({y=3} and {x=1})  
12      = P({y=3} and {x=0}) + P({y=3} and {x=1})  
13
```

```
1 # generalised :  
2 x,y have a joint pmf:  
3  
4     P[x = i , y = j]  
5  
6 P[x=i] =  
7 = (P[x=i] and P[y=1]) U (P[x=i] and P[y=2]) U .....U (P[x=i] and P[y=j])  
8 = sum( P[x=i and y=k] )  
9  
10 here, P[x = i] is marginal probability  
11
```

In []:

1	
---	--

In []:

1	
---	--

```
1 # 3 types of batteries  
2  
3 NEW    WORKING    DEFECTIVE  
4      3          4          5  
5  
6 pick randomly 3 batteries .  
7
```

```

8 X : no of new batteries
9 Y : no of working batteries
10
11 Q:
12
13 Y
14
15 3
16 2
17 1
18 0 1 2 3 X
19
20 X=0,y=0 means all batteries got defective
21 P[X=0 and y=0] = 5c3 / 12c3
22
23 P[x=0 and y=1] = 4c1 * 5c2 / 12c3
24 P[x=1 and y=1] = 3c1 * 4c1 5c1 / 12c3
25 P[x=1 and y=2] = 3c1 * 4c2 / 12c3
26 ....
27 ..
28 for entire table.
29
30
31 if asked P(y=1) = P(x=0,y=1) + P(x=1,y=1) + P(x=2,y=1) + P(x=3,y=1)
32

```

In []: 1

In []: 1

In []: 1
2
3

In []: 1

In []: 1

In []: 1

In []: 1

In [90]: 1 import scipy.stats
2

```
In [88]: 1 import scipy.stats  
2 scipy.stats.norm(loc=100, scale=12)  
3 #where loc is the mean and scale is the std dev  
4 #if you wish to pull out a random number from your distribution  
5 scipy.stats.norm.rvs(loc=100, scale=12)  
6  
7 #To find the probability that the variable has a value LESS than or equal  
8 #let's say 113, you'd use CDF cumulative Density Function  
9 scipy.stats.norm.cdf(113,100,12)  
10 # Output: 0.86066975255037792  
11 #or 86.07% probability  
12  
13 #To find the probability that the variable has a value GREATER than or  
14 #equal to Let's say 125, you'd use SF Survival Function  
15 scipy.stats.norm.sf(125,100,12)  
16 # Output: 0.018610425189886332  
17 #or 1.86%  
18  
19 #To find the variate for which the probability is given, let's say the  
20 #value which needed to provide a 98% probability, you'd use the  
21 #PPF Percent Point Function  
22 scipy.stats.norm.ppf(.98,100,12)  
23 # Output: 124.64498692758187
```

Out[88]: 124.64498692758187

```
In [91]: 1 scipy.stats.norm(loc=100, scale=12)  
2
```

Out[91]: <scipy.stats._distn_infrastructure.rv_frozen at 0x13b076a7d30>

```
In [96]: 1 scipy.stats.norm.rvs(loc=100, scale=12)  
2
```

Out[96]: 86.58028716991166

```
In [98]: 1 scipy.stats.norm.cdf(113,100,12)  
2
```

Out[98]: 0.8606697525503779

```
In [99]: 1 scipy.stats.norm.sf(125,100,12)  
2
```

Out[99]: 0.018610425189886332

```
In [100]: 1 scipy.stats.norm.cdf(125,100,12)  
2
```

Out[100]: 0.9813895748101137

```
In [101]: 1 1-0.9813895748101137  
2
```

Out[101]: 0.018610425189886315

In []:

1

In []:

1 0.5,0.2,0.7

In [105]:

1 0.5*0.2*0.7

Out[105]: 0.06999999999999999

In []:

1

```
In [174]: 1 arr = [float(x) for x in "1.764052345967664 0.4001572083672233 0.9787379841057392 2.240893199201"]
           ◀ ▶
```

```
In [ ]: 1
```

```
In [175]: 1 arr,len(arr)
```

```
Out[175]: ([1.764052345967664,
 0.4001572083672233,
 0.9787379841057392,
 2.240893199201458,
 1.8675579901499675,
 -0.977277879876411,
 0.9500884175255894,
 -0.1513572082976979,
 -0.10321885179355784,
 0.41059850193837233,
 0.144043571160878,
 1.454273506962975,
 0.7610377251469934,
 0.12167501649282841,
 0.44386323274542566,
 0.33367432737426683,
 1.4940790731576061,
 -0.20515826376580087,
 0.31306770165090136,
 0.35100577202017218]
```

```
In [176]: 1 arr = np.array(arr)
```

```
In [177]: 1 std = arr.std()
2 std
```

```
Out[177]: 1.1133678186681364
```

```
In [178]: 1 mean = arr.mean()
2 mean
```

```
Out[178]: 0.12524032797040804
```

```
In [179]: 1 upper = mean + (0.025774048*std)
2 upper
```

```
Out[179]: 0.1539363235704159
```

```
In [180]: 1 lower = mean - (0.025774048*std)
```

```
In [181]: 1 lower,upper
```

```
Out[181]: (0.0965443323704002, 0.1539363235704159)
```

```
In [182]: 1 for i in arr:
2     if i < upper and i > lower:
3         print("yess")
```

```
yess
yess
```

```
In [ ]: 1
```

```
In [ ]: 1
```

Central Limit Theorem:

If simple random samples of size n is selected from a distribution with mean μ and standard deviation $= \sigma$, then the distribution of sample means approach a normal distribution, as sample size increases, **irrespective of parent distribution of population**. The greater the sample size, the better is the approximation (typically $n \geq 30$).

In this condition, Sample mean follows approximately normal distribution

with Mean of sampling means:

$$\mu_{\bar{x}} = \mu$$

and

Standard deviation of Sample mean or Standard error:

$$SE = \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

~~~ *Thank you* ~~~

---

[socratic.org/users/VarunStatsHelp](https://socratic.org/users/VarunStatsHelp)

In [ ]:

1