

▼ cdf and ppf

- cdf: Cumulative Distribution Function
- ppf: Percent Point Function
- cdf & ppf are inverse of each other

▼ scipy.stats.norm

```
1 import scipy
2 from scipy import stats
3 from scipy.stats import norm
4 from scipy.stats import t
```

▼ norm.cdf(Z)

- cdf: Cumulative Distribution Function
- takes Z-Statistics, value on x-axis
- returns area under the curve towards the left side of the mean
- returned area under the curve represents probability / P-Value
- total area under all the curve or total probability = 1

```
1 norm.cdf(-1)
2 # cdf takes values from the x-axis, which can be +ve/-ve
3 # cdf returns area under the curve to the left side of the value on x-axis

0.15865525393145707
```

▼ norm.ppf(area/prob)

- ppf: Percent Point Function
- takes area under the curve or the probability
- area under the curve as argument can't be -ve
- returns Z-Statistics or Z-Score on X-axis on left side of the mean

```
1 norm.ppf(0.15865525393145707)
2 # ppf: Percent Point Function
3 # Inverse of CDF: Cumulative distribution function
```

-1.0

```
1 norm.ppf(0.025)
2 # takes 0.025 or 2.5% of area under the curve on the left side of the mean
3 # returns Z-Statistics on X-axis
```

-1.9599639845400545

```
1 norm.ppf(0.975)
2 # takes 0.025 or 2.5% of area under the curve on the left side of the mean
3 # returns Z-Statistics on X-axis
```

1.959963984540054

```
1 norm.ppf(0.5)
2 # takes 0.5 or 50% area under the curve on the left side of the curve
3 # total area under the curve or total probability = 1
4 # so Z-Statistics = 0
```

0.0

▼ Hypothesis Testing

- Null Hypothesis H_0 : No effect exists in the population.
- Alternative Hypothesis H_A / H_1 / H_2 : The effect exists in the population.

▼ Steps involved in Hypothesis analysis

Step 1: Establish Null Hypothesis H_0 & Alternative Hypothesis (H_A / H_1 / H_2)

- If Null Hypothesis, H_0 is found to be true, then
 - NO action , NO decision , NO changes
- 1. Right Tail Test : when error is on right side of Null Hypothesis on number line
 - e.g. delivery time ≤ 40 min, speed limit ≤ 80 kmph
- 2. Left Tail Test : when error is on left side of Null Hypothesis on number line

- e.g. voltage ≤ 220

3. Two Tail Test : when error is both sides of an interval

- e.g. Blood Pressure level ≤ 80 and Blood Pressure level ≥ 120
- e.g. Voltage ≤ 180 & voltage ≥ 240

- Null Hypothesis, H_0 will always have an equal to sign
- e.g. salary promised $\geq 30,000$ is Null Hypothesis

Step 2: Establish Confidence Level(C.L.) and the Level of Significance (LoS) for the analysis

- Confidence level (C.L.)
 - default value is ≥ 0.95 or 95%
- Level of Significance (LoS):
 - $LoS = 1 - C.L.$
 - default value of LoS is ≤ 0.05 or 5%

Step 3: Choose the appropriate statistical test for performing the analysis

Step 4: use the appropriate functions in python to calculate the P-Value

Step 5: Compare the P-value (which you got from step 4) against the Level of Significance LoS (which you got from Step 2)

- NOTE:
 - If P-value $< LoS$, then reject H_0
 - If P-value $> LoS$, then don't reject H_0
- Data used in hypothesis is prone to error as it is not population data, it is sample data
- Statistical Sample = Population parameter + Actual Difference + Chance Variation [not a mathematical equation]
- Number of possible Hypothesis = Number of variables + 1
 - for $y = mx + c$, number of variables = 1
 - so, number of possible hypothesis = 2

▼ Types of Error

1. Type 1 Error
2. Type 2 Error

Type 1 Error

- rejecting a Null Hypothesis H_0 which is actually true is Type 1 Error
- rejecting due to limited information
- e.g.: rejecting a good girl / important phone call / some good movie due to limited information

▼ Type 2 Error

- Not rejecting the Null Hypothesis H_0 , which is actually False
- e.g.: not rejecting a bad movie
- LoS : max allowed probability of committing Type 1 error
- P-Value : actual probability of committing Type 1 Error
- assume maxprob(getting affecting) / LoS = 7%
- if p-value / actual probability is 23%, then we reject H_0
- if p-value / actual probability is 3%, then we accept H_0
- if the Sample Statistics and Plain Value(P-Value)/Population Parameter do not meet the criteria, then P-value gives the probability that the difference is just a matter of chance, there is no actual difference and the action can be avoided.

General Note (Hypothesis Testing)

- H_0 : X does not have an impact on Y
 - X : predictor / variables
 - Y : response / result
- if $P\text{-Value} < 0.05$ ($P < |Z|$)
 - H_0 is rejected
 - X has impact on Y . $y = mx + c$
 - influence of X on Y cannot be ignored
 - It is statistically significant
 - it is not a matter of chance variation

▼ Tools for continuous data

1. Sample t Tests / Sample Z tests

- one-sample t test / one-sample Z test
- two-sample t test / two-sample Z test

2. ANOVA tests

- one-way ANOVA
- two-way ANOVA

3. one-Sample Variance Test

▼ Sample t test / Sample Z test

1. one-sample t test / one-sample Z test
2. two-sample t test / two-sample Z test

- uses ztest
- used to compare mean against a claimed value

▼ one-sample t test / one-sample Z test

- one-sample t test is used to compare the mean of a sample against a claimed value and establish whether the difference between them is statistically significant or not
- z-test
 - population S.D. should be known
 - we work with sample S.D.
 - Sample size, $n > 30$
- t-test
 - Sample S.D should be known
 - Sample size, $n \leq 30$
- data should follow normal distribution

```
1 x = [4, 7, 6, 8, 9, 7, 2, 3, 5, 6, 7, 8]
2 # create a list of sample data
```

```

1 import numpy as np
2 import pandas as pd
3
4 import statsmodels
5 from statsmodels import stats
6 from statsmodels.stats import weightstats
7 # for ztest
8
9
10 from scipy import stats
11 from scipy.stats import norm
12 # for P-value / probability

```

```

1 np.mean(x)
2 # actual sample Mean to be used to compare with claimed mean
3 # and then establish whether difference is statistically significant or not

```

6.0

assuming that, H_0 : Mean(x) is atleast 7.2

Step1: establishing H_0 & H_A

=> H_0 : $\mu \geq 7.2$ and H_A : $\mu < 7.2$

```

- SS = PP + ActualDifference + ChanceVariation
- 6    7.2                      could be a matter of chance variation

```

Step2: CL = 95% , LoS = 5% = 0.05

Step3: Choose one-sample t test / one-sample Z test to compare a sample mean with the claim value

Step4: use appropriate python function to calculate P-value

▼ statsmodels.stats.weightstats.ztest(list/ndarray, value=Claimed Value, alternative=Sign Of H_A)

```

1 statsmodels.stats.weightstats.ztest(x, value=7.2, alternative='smaller')
2 # statsmodels.stats.weightstats.ztest(list/ndarray, value=Claimed Value, alternative=Sign Of  $H_A$ )
3 # returns (TestStatistics, P-Value)

(-1.9497692171126306, 0.025601816055026008)

```

▼ norm.cdf(testStatistics)

- to cross verify the P-Value with the P-Value we got from ztest() using the TestStatistics value

```
1 norm.cdf(-1.9497692171126306)
2 # norm.cdf(testStatistics) returns P-value / Probability

0.025601816055026008
```

Step5: compare the P-Value with LoS

- $P\text{-Value}(0.025601816055026008) < LoS(0.05)$, so we reject the H_0 .
- the difference between 6 and 7.2 is statistically significant, cannot be ignored, so action is required
- Next steps to take an action and repeat (here repeating with different claimed value)

In next repetition, assuming that, H_0 : Mean(x) is atleast 6.5

Step1: establishing H_0 & H_A

=> H_0 : $\mu \geq 6.5$ and H_A : $\mu < 6.5$

```
1 statsmodels.stats.weightstats.ztest(x, value=6.5, alternative='smaller')
2 # statsmodels.stats.weightstats.ztest(list/ndarray, value=ClaimedValue, alternative=SignOfAltHypothesis)
3 # returns (TestStatistics, P-Value)

(-0.812403840463596, 0.2082799716383955)
```

- Now, $P\text{-Value}(0.2082799716383955) > LoS(0.05)$, so we do not reject the H_0 .
- But, the difference between 6 and 6.5 is statistically significant which cannot be ignored, so action is required, and we repeat Step1 to Step5 with different claimed value

▼ scipy.stats.t

```
1 import scipy
2 from scipy import stats
3 from scipy.stats import t
```

▼ scipy.stats.t.interval(C.L., dof, \bar{x} , StdErrMean)

- calculating interval, and then calculating the probability / area under the curve for the interval as Z-statistics on x-axis using CDF method

```
1 x = stats.t.interval(0.95, 2000-1, 27000, 1000/2000**0.5)
2 # interval(C.L., dof,  $\bar{x}$ , StdErrMean)
3 # Standard Error of the Mean, StdErrMean = S.D. / sqrt(SampleSize)
4 # returns (lower_value, higher_value)
5 x
```

```
(26956.147321103283, 27043.852678896717)
```

▼ scipy.stats.norm

```
1 import scipy
2 from scipy import stats
3 from scipy.stats import norm
4 norm.cdf(1) - norm.cdf(-1) #  $1-\sigma$ 
```

```
0.6826894921370859
```

▼ norm.cdf(Z)

```
1 norm.cdf(1)
2 # area under curve/probability for z-statistics = 1
```

```
0.8413447460685429
```

```
1 norm.cdf(-1)
2 # area under curve/probability for z-statistics = -1
```

```
0.15865525393145707
```

```
1 norm.cdf(26956.147321103283)
2 # z-statistics is lower end of interval
3 # area under curve/probability for z-statistics = 26956.147321103283
```

```
1.0
```

```
1 norm.cdf(27043.852678896717)
2 # z-statistics is upper end of interval
3 # area under curve/probability for z-statistics = 27043.852678896717
```

```
1.0
```



```
1 norm.cdf(26956.147321103283) - norm.cdf(27043.852678896717)
2 # area under the curve / probability between the intervals
```

```
0.0
```

Note

- Non-parametric test can be used without considering the distribution