

▼ Chi-Square test

Chi square test is a statistical method to determine if two categorical variables have a significant correlation between them.

Both those variables should be from same population and they should be categorical like – Yes/No, Male/Female, Red/Green etc.

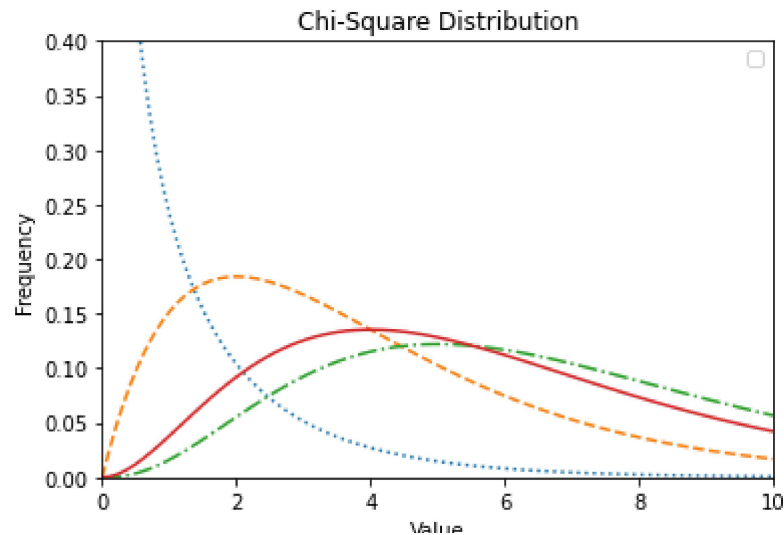
For example, we can build a data set with observations on people's ice-cream buying pattern and try to correlate the gender of a person with the flavour of the ice-cream they prefer.

If a correlation is found we can plan for appropriate stock of flavours by knowing the number of gender of people visiting.

We use various functions in **numpy and scipy library** to carry out the chi-square test.

```
1 from scipy import stats
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 x = np.linspace(0, 10, 100)
6 fig,ax = plt.subplots(1,1)
7
8 linestyle = [':', '--', '-.', '-']
9 deg_of_freedom = [1, 4, 7, 6]
10 for df, ls in zip(deg_of_freedom, linestyle):
11     ax.plot(x, stats.chi2.pdf(x, df), linestyle=ls)
12
13 plt.xlim(0, 10)
14 plt.ylim(0, 0.4)
15
16 plt.xlabel('Value')
17 plt.ylabel('Frequency')
18 plt.title('Chi-Square Distribution')
19
20 plt.legend()
21 plt.show()
```

No handles with labels found to put in legend.



```

1 # chi-squared test with similar proportions
2 from scipy.stats import chi2_contingency
3 from scipy.stats import chi2
4 # contingency table
5 table = [ [10, 20, 30],
6           [6, 9, 17]]
7 print(table)
8 stat, p, dof, expected = chi2_contingency(table)
9 print('dof=%d' % dof)
10 print(expected)
11 print('chi2_stat=%.2f' % stat)
12 # interpret test-statistic
13 prob = 0.95
14 critical = chi2.ppf(prob, dof)
15 print('probability=%.3f, critical=%.3f, stat=%.3f' % (prob, critical, stat))
16 if abs(stat) >= critical:
17     print('Dependent (reject H0)')
18 else:
19     print('Independent (fail to reject H0)')
20 # interpret p-value
21 alpha = 1.0 - prob
22 print('significance=%.3f, p=%.3f' % (alpha, p))
23 if p <= alpha:

```

```
24 print('Dependent (reject H0)')
25 else:
26 print('Independent (fail to reject H0)')

[[10, 20, 30], [6, 9, 17]]
dof=2
[[10.43478261 18.91304348 30.65217391]
 [ 5.56521739 10.08695652 16.34782609]]
chi2_stat=0.27
probability=0.950, critical=5.991, stat=0.272
Independent (fail to reject H0)
significance=0.050, p=0.873
Independent (fail to reject H0)

1 # Load libraries
2 from sklearn.datasets import load_iris
3 from sklearn.feature_selection import SelectKBest
4 from sklearn.feature_selection import chi2
5
6 # Load iris data
7 iris_dataset = load_iris()
8
9 # Create features and target
10 X = iris_dataset.data
11 y = iris_dataset.target
12
13 # Convert to categorical data by converting data to integers
14 X = X.astype(int)
15
16 # Two features with highest chi-squared statistics are selected
17 chi2_features = SelectKBest(chi2, k = 2)
18 X_kbest_features = chi2_features.fit_transform(X, y)
19
20 # Reduced features
21 print('Original feature number:', X.shape[1])
22 print('Reduced feature number:', X_kbest_features.shape[1])
23
```

Original feature number: 4

Reduced feature number: 2

```
1 from scipy.stats import chi2_contingency
2
3 # defining the table
4 data = [[207, 282, 241], [234, 242, 232]]
5 stat, p, dof, expected = chi2_contingency(data)
6
7 # interpret p-value
8 alpha = 0.05
9 print("p value is " + str(p))
10 if p <= alpha:
11     print('Dependent (reject H0)')
12 else:
13     print('Independent (H0 holds true)')
14
```

p value is 0.1031971404730939
Independent (H0 holds true)

1

✓ 0s completed at 10:03 AM



Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.