

Chi-Squared Test of Independence Chi-Squared Test of Independence is a key concept in probability that describes a situation where knowing the value of one variable tells you nothing about the value of another. For instance, the month you were born probably doesn't tell you anything about which web browser you use, so we'd expect birth month and browser preference to be independent. On the other hand, your month of birth might be related to whether you excelled at sports in school, so month of birth and sports performance might not be independent.

The chi-squared test of independence tests whether two categorical variables are independent. The test of independence is commonly used to determine whether variables like education, political views and other preferences vary based on demographic factors like gender, race and religion. Let's generate some fake voter polling data and perform a test of independence.

Importing Libraries

```
1 import numpy as np
2 import pandas as pd
3 import scipy.stats as stats
```

```
""" METHODOLOGY 01: MANUAL CALCULATION """
```

▼ STEP 1: GENERATE A RANDOM DATASET

Generate under a random factor

<https://docs.scipy.org/doc/numpy/reference/generated/numpy.random.seed.html>

```
1 np.random.seed(10)

1 # Sample data randomly at fixed probabilities
2 voter_race = np.random.choice(a=["asian", "black", "hispanic", "other", "white"],
3                                p=[0.05, 0.15, 0.25, 0.05, 0.5],
4                                size=1000)

1 # Sample data randomly at fixed probabilities
2 voter_party = np.random.choice(a=["democrat", "independent", "republican"],
3                                 p=[0.4, 0.2, 0.4],
4                                 size=1000)
```

```

1 # Binding 2 arrays (voter_race and voter_party) to make a DataFrame
2 voters = pd.DataFrame({"race":voter_race,
3                         "party":voter_party})
4 # You can check the data of DataFrame by calling it
5 voters

```



	race	party
0	white	democrat
1	asian	republican
2	white	independent
3	white	republican
4	other	democrat
...
995	white	republican
996	hispanic	independent
997	black	independent
998	white	republican
999	black	democrat

1000 rows × 2 columns

```

1 # Create a CrossTab from DataFrame, Assign the column names and row names
2 voter_tab = pd.crosstab(voters.race, voters.party, margins=True)
3 voter_tab.columns = ["democrat", "independent", "republican", "row_totals"]
4 voter_tab.index = ["asian", "black", "hispanic", "other", "white", "col_totals"]
5 # You can check the data of CrossTab by calling it
6 voter_tab
7

```




	democrat	independent	republican	row_totals
asian	21	7	32	60
black	65	25	64	154
hispanic	107	50	94	251
other	15	8	15	38
white	189	96	212	497
col_totals	397	186	417	1000

STEP 2: GET THE "OBSERVED" TABLE AND "EXPECTED" TABLE

Calculate the "observed" table: "Observed" table can be extracted from our CrossTab by exclude the row_totals and col_totals You can see row_totals is in the index of 4 (in column) and col_totals is in the index of 6 (in row). [0:5, 0:3] means "we will take the rows from 0 index to 5 index and columns from 0 index to 3 index and assign to new CrossTab that named [observed]

```
1 observed = voter_tab.iloc[0:5, 0:3]
2 # You can check the data of observed table by calling it
3 observed
```



	democrat	independent	republican
asian	21	7	32
black	65	25	64
hispanic	107	50	94
other	15	8	15
white	189	96	212

Calculate the "expected" table: "Expected" table can be calculated using below formula: $\text{total_rows} \times \text{total_columns} / \text{total_observations}$ And these factors can be get by: - total_rows = voter_tab["row_totals"] - total_columns = voter_tab["col_totals"] - total_observations = 1000 Please note that the "loc" function in below code is used to switch the index base on column name to row name

```
1 expected = np.outer(voter_tab["row_totals"][0:5],
2                     voter_tab.loc["col_totals"][0:3]) / 1000
3 # Now convert into a DataFrame, Assign the column names and row names
4 expected = pd.DataFrame(expected)
5 expected.columns = ["democrat", "independent", "republican"]
6 expected.index = ["asian", "black", "hispanic", "other", "white"]
7 # You can check the data of expected table by calling it
8 expected
9
```



STEP 3: CALCULATE THE CHI SQUARE VALUE and CRITICAL VALUE

Chi square formula: $\text{chi square} = \text{total of } [(\text{observed} - \text{expected})^2 / \text{expected}]$

Note: We call `.sum()` twice: once to get the column sums and a second time to add the column sums together, returning the sum of the entire 2D table.

```
1 chi_squared_stat = (((observed-expected)**2)/expected).sum().sum()
2 print(chi_squared_stat)
```

```
7.169321280162059
```

Find the critical value for confidence of 95% and degree of freedom (df) of 8 Why $df = 8$? Degree of freedom formula: $df = (\text{total rows} - 1) \times (\text{total columns} - 1) = (5 - 1) \times (3 - 1) = 4 \times 2 = 8$

```
1 crit = stats.chi2.ppf(q = 0.95, # Find the critical value for 95% confidence*
2                             df = 8) # *
3
4 print("Critical value")
5 print(crit)
6
7 p_value = 1 - stats.chi2.cdf(x=chi_squared_stat, # Find the p-value
8                             df=8)
9 print("P value")
10 print(p_value)
```

```
Critical value
15.50731305586545
P value
0.518479392948842
```

STEP 4: MAKE THE CONCLUSION

Because $\text{chi_squared_stat} < \text{crit}$

When your p-value is less than or equal to your significance level, you reject the null hypothesis. The data favors the

alternative hypothesis.

With some code to accept the alternative hypothesis if we fail to reject the null hypothesis:

```
1 if chi_squared_stat < crit:
2     print("""At 0.95 level of significance, we reject the null hypotheses and accept H1.
3 They are not independent.""")
4 else:
5     print("""At 0.95 level of significance, we accept the null hypotheses.
6 They are independent.""")
```

➞ At 0.95 level of significance, we reject the null hypotheses and accept H1.
They are not independent.

METHODOLOGY 02: CALCULATE USING SCIPY.STATS LIBRARY

```
1 """ METHODOLOGY 02: CALCULATE USING SCIPY.STATS LIBRARY"""
2 stats = stats.chi2_contingency(observed=observed)
3 # You can check the returned data by calling it
4 # The returned data includes: chi_squared_stat, p_value, df, expected_crosstab
5 print(stats)
```

➞ (7.169321280162059, 0.518479392948842, 8, array([[23.82 , 11.16 , 25.02],
[61.138, 28.644, 64.218],
[99.647, 46.686, 104.667],
[15.086, 7.068, 15.846],
[197.309, 92.442, 207.249]]))

1

