Never miss a tutorial:

# Machine Learning Mastery

Making Developers Awesome at Machine Learning

Picked for you:

Statistics for Machine Learning (7-Day Mini-Course)

**Click to Take the FREE Statistics Crash-Course**

Search...

A Gentle Introduction to k-fold Cross-Validation

# A Gentle Introduction to the Chi-Squared Test for Machine Learning

**Brownlee** on June 15, 2018 in **Statistics**

How to Calculate Bootstrap Confidence Intervals For Machine Learning Results in Python

Tweet    Share    Share

Last Updated on October 31, 2019

A Gentle Introduction to Normality Tests in Python

A common problem in applied machine learning is determining whether input features are relevant to the outcome to be predicted.

Statistical Significance Tests for Comparing Machine Learning Algorithms

This is the problem of feature selection.

In the case of classification problems where input variables are also categorical, we can use statistical tests to determine whether the output variable is dependent or independent of the input variables. If independent, then the input variable is a candidate for a feature that may be irrelevant to the problem and removed from the dataset.

**Loving the Tutorials?**

The Statistics for Machine Learning EBook is where you'll find the *Really Good* stuff.

The Pearson's chi-squared statistical hypothesis is an example of a test for independence between categorical   >> SEE WHAT'S INSIDE

In this tutorial, you will discover the chi-squared statistical hypothesis test for quantifying the independence of pairs of categorical variables.

After completing this tutorial, you will know:

- Pairs of categorical variables can be summarized using a contingency table.
- The chi-squared test can compare an observed contingency table to an expected table and determine if the categorical variables are independent.
- How to calculate and interpret the chi-squared test for categorical variables in Python.

**Kick-start your project** with my new book Statistics for Machine Learning, including *step-by-step tutorials* and the *Python source code* files for all examples.

Let's get started.

Start Machine Learning

**Never miss a tutorial:**

- **Update Jun/2018**: Minor typo fix in the interpretation of the critical values from the test (thanks Andrew).
- **Update Oct/2019**: Fixed language around factor/levels (thanks Marc)

**Picked for you:**

[Statistics for Machine Learning (7-Day Mini-Course)](#)

[A Gentle Introduction to k-fold Cross-Validation](#)

[How to Calculate Bootstrap Confidence Intervals For Machine Learning Results in Python](#)

[A Gentle Introduction to Normality Tests in Python](#)

[Statistical Significance Tests for Comparing Machine Learning Algorithms](#)

**Loving the Tutorials?**

The [Statistics for Machine Learning](#) EBook is where you'll find the *Really Good* stuff.

>> SEE WHAT'S INSIDE

## Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

A Gentle Introduction to the Chi-Squared Test for Machine Learning
Photo by [NC Wetlands](#), some rights reserved

# Tutorial Overview

This tutorial is divided into 3 parts, they are:

1. Contingency Table
2. Pearson's Chi-Squared Test
3. Example Chi-Squared Test

---

## Need help with Statistics for Machine Learning?

Take my free 7-day email crash course now (with sample code).

Click to sign-up and also get a free PDF Ebook version of the course.

Download Your FREE Mini-Course

Start Machine Learning

## Contingency Table

A categorical variable is a variable that may take on one of a set of labels.

An example might be sex, which may be summarized as male or female. The variable or factor is '*sex*' and the labels or levels of the variable are '*male*' and '*female*' in this case.

We may wish to look at a summary of a categorical variable as it pertains to another categorical variable. For example, sex and interest, where interest may have the labels '*science*', '*math*', or '*art*'. We can collect observations from people collected with regard to these two categorical variables; for example:

```
1 Sex,     Interest
2 Male,    Art
3 Female, Math
4 Male,    Science
5 Male,    Math
6 ...
```

We can summarize the collected observations in a [table] with one variable corresponding to columns and another variable corresponding to rows. Each cell in the table corresponds to the count or frequency of observations that correspond to the row and column categories.

Historically, a table summarization of two categorical variables in this form is called a contingency table.

For example, the *Sex=rows* and *Interest=columns* table with contrived counts might look as follows:

```
1          Science,    Math,    Art
2 Male        20,        30,     15
3 Female      20,        15,     30
```

The table was called a contingency table, by Karl Pearson, because the intent is to help determine whether one variable is contingent upon or depends upon the other variable. For example, does an interest in math or science depend on gender, or are they independent?

This is challenging to determine from the table alone; instead, we can use a statistical method called the Pearson's Chi-Squared test.

# Pearson's Chi-Squared Test

The Pearson's Chi-Squared test, or just Chi-Squared test for short, is named for Karl Pearson, although there are variations on the test.

The Chi-Squared test is a statistical hypothesis test that assumes (the null hypothesis) that the observed frequencies for a categorical variable match the expected frequencies for the categorical variable. The test calculates a statistic that has a chi-squared distribution, named for the Greek capital letter Chi (X) pronounced "ki" as in kite.

Given the Sex/Interest example above, the number of observations for a category (such as male and female) may or may not the same. Nevertheless, we can calculate the expected frequency of observations in each Interest group and see whether the partitioning of interests by Sex results in similar or different frequencies.

The Chi-Squared test does this for a contingency table, first calculating the expected frequencies for the groups, then determining whether the division of the groups, called the observed frequencies, matches the expected frequencies.
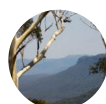
The result of the test is a test statistic that has a chi-squared distribution and can be interpreted to reject or fail to reject the assumption or null hypothesis that the observed and expected frequencies are the same.

> *When observed frequency is far from the expected frequency, the corresponding term in the sum is large; when the two are close, this term is small. Large values of X^2 indicate that observed and expected frequencies are far apart. Small values of X^2 mean the opposite: observeds are close to expecteds. So X^2 does give a measure of the distance between observed and expected frequencies.*

— Page 525, *Statistics*, Fourth Edition, 2007.

The variables are considered independent if the observed and expected frequencies are similar, that the levels of the variables do not interact, are not dependent.

> *The chi-square test of independence works by comparing categorically the data that you have collected (known as the observed frequencies) with the frequencies that you would expect to get in each cell of a table by chance alone ...*

— Page 162, [Statistics in Plain English](#), Third Edition, 2010.

We can interpret the test statistic in the context of the chi-squared distribution with the requisite number of degress of freedom as follows:

- **If Statistic >= Critical Value**: significant result, reject null hypothesis (H0), dependent.
- **If Statistic < Critical Value**: not significant result, fail to reject null hypothesis (H0), independent.

The degrees of freedom for the chi-squared distribution is calculated based on the size of the contingency table as:

```
1  degrees of freedom: (rows - 1) * (cols - 1)
```

In terms of a p-value and a chosen significance level (alpha), the test can be interpreted as follows:

- **If p-value <= alpha**: significant result, reject null hypothesis (H0), dependent.
- **If p-value > alpha**: not significant result, fail to reject null hypothesis (H0), independent.

For the test to be effective, at least five observations are required in each cell of the contingency table.

Next, let's look at how we can calculate the chi-squared test.

# Example Chi-Squared Test

The Pearson's chi-squared test for independence can be calculated in Python using the chi2_contingency() SciPy function.
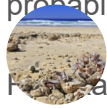
The function takes an array as input representing the contingency table for the two categorical variables. It returns the calculated statistic and p-value for interpretation as well as the calculated degrees of freedom and table of expected frequencies.

```
1  stat, p, dof, expected = chi2_contingency(table)
```

We can interpret the statistic by retrieving the critical value from the chi-squared distribution for the probability and number of degrees of freedom.

For example, a probability of 95% can be used, suggesting that the finding of the test is quite likely given the assumption of the test that the variable is independent. If the statistic is less than or equal to the critical value, we can fail to reject this assumption, otherwise it can be rejected.

```
1  # interpret test-statistic
2  prob = 0.95
3  critical = chi2.ppf(prob, dof)
4  if abs(stat) >= critical:
5      print('Dependent (reject H0)')
6  else:
7      print('Independent (fail to reject H0)')
```

We can also interpret the p-value by comparing it to a chosen significance level, which would be 5%, calculated by inverting the 95% probability used in the critical value test.

```
1  # interpret p-value
2  alpha = 1.0 - prob
3  if p <= alpha:
4      print('Dependent (reject H0)')
5  else:
6      print('Independent (fail to reject H0)')
```

We can tie all of this together and demonstrate the chi-squared significance test using a contrived contingency table.

A contingency table is defined below that has a different number of observations for each population (row), but a similar proportion across each group (column). Given the similar proportions, we would expect the test to find that the groups are similar and that the variables are independent (fail to reject the null hypothesis, or H0).

```
1  table = [   [10, 20, 30],
2              [6,  9,  17]]
```

The complete example is listed below.

```
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  # interpret test-statistic
12  prob = 0.95
13  critical = chi2.ppf(prob, dof)
14  print('probability=%.3f, critical=%.3f, stat=%.3f' % (prob, critical, stat))
15  if abs(stat) >= critical:
16      print('Dependent (reject H0)')
```

```
17  else:
18      print('Independent (fail to reject H0)')
19  # interpret p-value
20  alpha = 1.0 - prob
21  print('significance=%.3f, p=%.3f' % (alpha, p))
22  if p <= alpha:
23      print('Dependent (reject H0)')
24  else:
25      print('Independent (fail to reject H0)')
```

Running the example first prints the contingency table. The test is calculated and the degrees of freedom (dof) is reported as 2, which makes sense given:

```
1  degrees of freedom: (rows - 1) * (cols - 1)
2  degrees of freedom: (2 - 1) * (3 - 1)
3  degrees of freedom: 1 * 2
4  degrees of freedom: 2
```

Next, the calculated expected frequency table is printed and we can see that indeed the observed contingency table does appear to match via an eye-ball check.

The critical value is calculated and interpreted, finding that indeed the variables are dependent (fail to reject H0). The interpretation of the p-value makes the same finding.

```
1  [[10, 20, 30], [6, 9, 17]]
2
3  dof=2
4
5  [[10.43478261 18.91304348 30.65217391]
6   [ 5.56521739 10.08695652 16.34782609]]
7
8  probability=0.950, critical=5.991, stat=0.272
9  Independent (fail to reject H0)
10
11  significance=0.050, p=0.873
12  Independent (fail to reject H0)
```

# Extensions

This section lists some ideas for extending the tutorial that you may wish to explore.

- Update the chi-squared test to use your own contingency table.
- Write a function to report on the independence given observations from two categorical variables
- Load a standard machine learning dataset containing categorical variables and report on the independence of each.

If you explore any of these extensions, I'd love to know.

# Further Reading

This section provides more resources on the topic if you are looking to go deeper.

## Books

- Chapter 14, The Chi-Square Test of Independence, Statistics in Plain English, Third Edition, 2010.
- Chapter 28, The Chi-Square Test, Statistics, Fourth Edition. 2007.

The Statistics for Machine Learning EBook is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

Start Machine Learning

## API

**Never miss a tutorial:**

- scipy.stats.chisquare() API
- scipy.stats.chi2_contingency() API
- sklearn.feature_selection.chi2() API
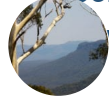
**Picked for you:**

# Articles

- Statistics for Machine Learning (7-Day Mini-Course)
- Chi-squared test on Wikipedia
- Pearson's chi-squared test on Wikipedia
- Contingency table on Wikipedia
- A Gentle Introduction to k-fold Cross-Validation
- How is chi test used for feature selection in machine learning? on Quora

# Summary

How to Calculate Bootstrap Confidence Intervals For Machine Learning Results in Python

In this tutorial, you discovered the chi-squared statistical hypothesis test for independence of pairs of categorical variables.

Specifically, you learned:

A Gentle Introduction to Normality Tests in Python

- Pairs of categorical variables can be summarized
- The chi-squared test can compare an observe
- determine if the categorical variables are inde
- Statistical Significance Tests for Comparing Machine Learning Algorithms
- How to calculate and interpret the chi-squared test for categorical variables in Python.

Do you have any questions?
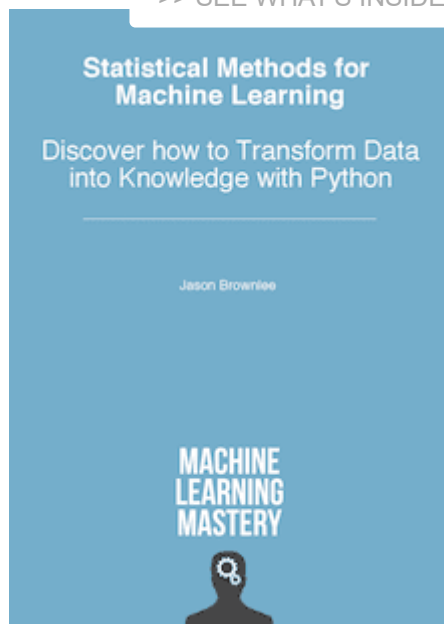Ask your questions in the comments below and I will do my best to answer.

## Loving the Tutorials?

The Statistics for Machine Learning EBook is where you'll find the *Really Good* stuff.

# Get a Handle on Statistics for Machine Learning!

>> SEE WHAT'S INSIDE

**Develop a working understanding of statistics**

...by writing lines of code in python

Discover how in my new Ebook:
Statistical Methods for Machine Learning

It provides **self-study tutorials** on topics like:
*Hypothesis Tests, Correlation, Nonparametric Stats, Resampling*, and much more...

**Discover how to Transform Data into Knowledge**

Skip the Academics. Just Results.

SEE WHAT'S INSIDE

Start Machine Learning

---

**Start Machine Learning** ✕

You can master applied Machine Learning
**without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

**START MY EMAIL COURSE**

**Never miss a tutorial:**

Tweet       Share            Share

**Picked for you:**

About Jason Brownlee

Jason Brownlee, PhD is a machine learning specialist who teaches developers how to get results with modern machine learning methods via hands-on tutorials.

View all posts by Jason Brownlee →

Statistics for Machine Learning (7-Day Mini-Course)

‹ How to Calculate the 5-Number Summary for Your Data in Python

Statistical Significance Tests for Comparing Machine Learning Algorithms ›

A Gentle Introduction to k-fold Cross-Validation

How to Calculate Bootstrap Confidence Intervals For Machine Learning Results in Python

Responses to *A Gentle Introduction to the Chi-Squared Test for Machine Learning*

A Gentle Introduction to Normality Tests in Python

**Elie Kawerk** June 19, 2018 at 5:27 am #

Hi Jason,

Statistical Significance Tests for Comparing Machine Learning Algorithms

Thanks for this nice post.

What statistical test should be used to test the dependence of a continuous variable on a categorical variable (ex: weight and gender).

Best,
Elie

---

**Start Machine Learning**                                    ✕

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

**START MY EMAIL COURSE**

---

**Loving the Tutorials?**

The Statistics for Machine Learning EBook is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

June 19, 2018 at 6:38 am #                                    REPLY ↩

Good question. I have not seen a test that can do this directly.

Often, the continuous variable is made discrete/ordinal and the chi-squared test is used. It will give a results, but I'm not sure how statistically valid this would be.

---

**DearML** July 2, 2019 at 8:37 pm #                                    REPLY ↩

Is there any way to get the correlation between all the input features only but with binary values which is 0 and 1 (converted from true and false)?

---

**Jason Brownlee** July 3, 2019 at 8:33 am #                                    REPLY ↩

Perhaps if you convert input fe[...]              Start Machine Learning

**Never miss a tutorial:**

DearML July 3, 2019 at 3:00 pm #

its a discrete variables like for example

df = pd.DataFrame({
    'y1': [1,1,1,1,1,1,1,1,0,1,0,0],
    'y2': [1,1,1,1,1,1,1,1,0,1,1,0],
    y3: [0,1,0,0,0,1,0,0,0,1,1,0],
    y4: [0,1,1,1,0,0,1,1,0,0,1,0],
})

Here it should be a strong correlation this all are features ( independent/i

is there any methods I can use to f Please help.

Chi squared might be a g

**DearML** July 5, 2019 at 2:24 pm #

Chi squared is about input and output. isn't it? What about cosine similarity? i think it will work.

rownlee July 6, 2019 at 8:22 am #

Chi squared is only concerned with two categorical variables. They may or may not be inputs or outputs to a model.

What about cosine similarity exactly?

**DearML** July 8, 2019 at 1:55 pm #

cosine similarity can give me the similarity of two different vectors. here in my example above, it will say that y1 and y2 are related with some more than ~95%

**Jason Brownlee** July 9, 2019 at 8:04 am #

**Never miss a tutorial:**

Details here:
https://en.wikipedia.org/wiki/Cosine_similarity

in   twitter   f   email   rss

**Picked for you:**

Judith Vazquez June 26, 2018 at 1:09 am #          REPLY ↰

Statistics for Machine Learning (7-Day Mini-Course)

Hi Elie,

You might try using the binning technique. Please see below

http://www.saedsayad.com/binning.htm

A Gentle Introduction to k-fold Cross-Validation

Hope it helps 🙂

How to Calculate Bootstrap Confidence Intervals For Machine Learning Results in Python

**Olutobi Adeyemi** October 15, 2018

Analysis of variance will work for t

A Gentle Introduction to Normality Tests in Python

Adi June 16, 2019 at 6:24 pm #

Statistical Significance Tests for Comparing Machine Learning Algorithms

A 2-sample KS Test (https://docs.scipy
0.14.0/reference/generated/scipy.stats.ks_2samp.html)

**Loving the Tutorials?**          Jason Brownlee June 17, 2019 at 8:19 am #          REPLY ↰

The Statistics for Machine Learning EBook is
where you'll find the **Really Good** stuff.

Nice!

―――― >> SEE WHAT'S INSIDE ――――

Rishabh March 31, 2020 at 5:15 am #          REPLY ↰

Independent two sample t test

ana July 1, 2020 at 2:59 am #          REPLY ↰

https://en.wikipedia.org/wiki/Correlation_ratio

Andrew V. June 21, 2018 at 3:36 am #          REPLY ↰

Hi Jason, great article! One quick thing: shouldn't the above read: "If statistic > critical value
then significant result" and "If statistic <= critical value then non-significant result"? The statistic's value
and p-value should be inversely related.

**Never miss a tutorial:**

REPLY ↩

on    wn June 21, 2018 at 6:22 am #

Yes, thanks. That was a typo in the explanation. Fixed.

**Picked for you:**

Statistics for Machine Learning (7-Day
Mini-Course)
**Hani** December 25, 2018 at 10:51 pm #

REPLY ↩
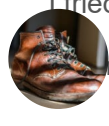
Hi ,

A Gentle Introduction to k-fold Cross-
validation can I loop the chisq to check the Target vs. all other variables in one step.
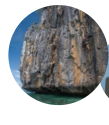and will let me know the p-value dof etc… of any c

I tried and it didnt work out….

How to Calculate Bootstrap Confidence
Intervals For Machine Learning Results in
Python
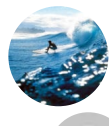


A Gentle Introduction to Normality Tests in
Python **Jason Brownlee** December 26, 2018 at 6

Perhaps write a for-loop to check all v

Statistical Significance Tests for
Comparing Machine Learning Algorithms

**SK** January 24, 2019 at 10:12 pm #

REPLY ↩

How do we perform chi-squared test for finding terms that are the most correlated with each
class ?

**Loving the Tutorials?**

The Statistics for Machine Learning EBook is
where you'll find the *Really Good* stuff.

>> SEE WHAT'S INSIDE            uary 25, 2019 at 8:44 am #

REPLY ↩

Perhaps calculate the test for each term, then rank order the results?

**Cody** March 8, 2019 at 3:43 am #

REPLY ↩

Very helpful and easy to understand. Thank you very much.

**Jason Brownlee** March 8, 2019 at 7:55 am #

REPLY ↩

Thanks, I'm glad it helped.

Start Machine Learning

---

**Start Machine Learning**                                    ✕

You can master applied Machine Learning
**without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

**START MY EMAIL COURSE**

**Never miss a tutorial:**

Sachin April 20, 2019 at 5:23 am #                                    REPLY ↩

Such a nice written article! Thanks for your time!

**Picked for you:**

Statistics for Machine Learning (7-Day Mini-Course)

Jason Brownlee April 20, 2019 at 7:43 am #                           REPLY ↩

Thanks, I'm glad it helped.

A Gentle Introduction to k-fold Cross-Validation

Vaishali Bhadwaj June 11, 2019 at 7:28 pm #

Hi

How to Calculate Bootstrap Confidence Intervals For Machine Learning Results in Python

I have 3 categorical variables in my data set.

Happiness, Income and Degree

A Gentle Introduction to Normality Tests in Python

Need to find below.

A survey was conducted among 2800 customers on status, sex, age, age-group, race, happiness, no. of income group etc. had been captured for that purpose.

a. Is there any relationship in between labour force Statistical Significance Tests for Comparing Machine Learning Algorithms educational qualification is somehow controlling the marital status? c. Is happiness is driven by earnings or marital status?

---

**Loving the Tutorials?**

The Statistics for Machine Learning EBook is Jason Brownlee June 12, 2019 at 7:55 am #          REPLY ↩

where you'll find the **Really Good** stuff.

Perhaps try using the chi squared test?

>> SEE WHAT'S INSIDE

---

BM August 18, 2019 at 3:53 am #                                      REPLY ↩

How to creat the cotingency table in python

---

Jason Brownlee August 18, 2019 at 6:49 am #                         REPLY ↩

See this:

https://www.statsmodels.org/stable/contingency_tables.html

---

Patrick C. September 12, 2019 at 8:44 am #                          REPLY ↩

Start Machine Learning

---

**Start Machine Learning**                                          ✕

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

**START MY EMAIL COURSE**

You could also have a look at the pandas crosstab functions ~ https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.crosstab.html

**Never miss a tutorial:**

**Jason Brownlee** September 12, 2019 at 1:48 pm #   REPLY

Thanks for the note Patrick.

**Bruno Ambrozio** September 30, 2019 at 3:15 am #   REPLY

Great content! Thanks!

Doubt:

If I understood well, with this chi-squared test you might group up categories, but if so, how to find out the ones that actually represent significant results?

Eg.: Let's say that you managed to reject the null hypothesis. Which categories would account for the significant differences?

Do we have to apply a Fisher exact test in each cat in contingency tables?

Thanks!

**Jason Brownlee** September 30, 2019 at 6:17 am #   REPLY

What do you mean group? Do you mean the categories for a given variable?

**Bruno Ambrozio** September 30, 2019 at 7:58 pm #   REPLY

Yes. In your example you have 3 categories been tested: Science, Math and Art. Let's say your result concludes you have evidence enough to reject the null hypothesis (the variables are dependent). But, how do you know which one (or whether all of them) account for such result?

Let's consider another example, where you have 34 categories (Degrees of Freedom = 33). You also manage to reject the null hypothesis. So, how do you know which of those 34 categories were responsible for the final result ($p \le alpha$)?

**Jason Brownlee** October 1, 2019 at 6:50 am #   REPLY

Yes, that is one discrete random variable that has 3 states or events.

The test comments on the random variable, not the states.

Does that help?

**Start Machine Learning**

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Start Machine Learning

**Never miss a tutorial:**

**Yogesh Naicker** October 9, 2019 at 8:24 pm #

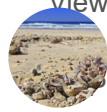Can need a lea provide python code for the below 4 categorical variables???

The table shows the contingency table of marital status by education. Use Chi-Square test for testing Homogenity contingency table of marital status by education.

**Picked for you:**

View the table by executing the following command python

prettytable import PrettyTable

Statistics for Machine Learning (7-Day Mini-Course)

PrettyTable(['Marital Status','Middle school', 'High School','Bachelor','Masters','PhD'])

t.add_row(['Single',18,36,21,9,6])

t.add_row(['Married',12,36,45,36,21])

A Gentle Introduction to k-fold Cross-Validation

row(['Divorced',6,9,9,3,3])

d_row(['Widowed',3,9,9,6,3])

print (t)

t()

How to Calculate Bootstrap Confidence Intervals For Machine Learning Results in Python

othesis

Null Hypothesis: There is no difference in distributio

marital status.

A Gentle Introduction to Normality Tests in Python

nate Hypothesis: There is a Difference

Coding

import chi2_contingency and chi2 from scipy.stat

Statistical Significance Tests for Comparing Machine Learning Algorithms

clare a 2D array with the values mentioned in the contingency table of marital status by education.

3.Calculate and print the values of

– Chi-Square Statistic

**Loving the Tutorials?**

– Degree of Freedom

– P-value

The Statistics for Machine Learning EBook is where you'll find the *Really Good* stuff.

– Hint: Use chi2_contingency() function

4.Assume the alpha value to be 0.05

>> SEE WHAT'S INSIDE

5.Compare the P-value with alpha and decide whether or not to reject the null hypothesis.

– If Rejected print "Reject the Null Hypothesis"

– Else print "Failed to reject the Null Hypothesis"

Sample output 2.33 4.5 8.9 Reject the Null Hypothesis

REPLY ↩

**Jason Brownlee** October 10, 2019 at 6:57 am #

Looks like homework. Perhaps try posting to stackoverflow?

REPLY ↩

**Sachin Ladhad** October 17, 2019 at 1:22 pm #

from scipy.stats import chi2_contin

from scipy.stats import chi2

**Never miss a tutorial:**

table= [ [18,31,21,9,6],[12,36,45,36,21], [6,9,9,3,3],[3,9,9,6,3] ]

stat,p,dof,expected = chi2_contingency(table)

prob 95

critical = chi2.ppf(prob, dof)

**Picked for you:**

if abs(p) <= 0.05:

    print(stat, dof ,p ,'Reject the Null Hypothesis')

Statistics for Machine Learning (7-Day
Mini-Course) else print(stat, dof ,p ,'Failed to reject the Null Hypothesis')

output

A Gentle Introduction to k-fold Cross 21.0528564352978821 12 0.0499013559023993 Reject the Null Hypothesis
Validation

Help needed : Please let me know why the

How to Calculate Bootstrap Confidence
Intervals For Machine Learning Results in
Python

**Jason Brownlee** October 17, 2

I have some suggestions here

A Gentle Introduction to Normality Tests in https://machinelearningmastery.com/fa
Python work-for-me

Statistical Significance Tests for
Comparing Machine Learning Algorithms

**Mary** November 29, 2019 at 7:48 am #

Your table indicates that for the "Single" road the values are 18,36,21,9,6

**Loving the Tutorials?** and for your code you have 18,*31*,21…

That 31 should be 36

The Statistics for Machine Learning EBook is

where you

```
1  from scipy.stats import chi2_contingency
2  from scipy.stats import chi2
3  table = [[18,36,21,9,6],[12,36,45,36,21],[6,9,9,3,3],[3,9,9,6,3]]
4  stat,p,dof,expected = chi2_contingency(table)
5  prob = 0.95
6  critical = chi2.ppf(prob, dof)
7  if abs(stat) >= 0.05:
8      print(stat, dof, p, 'Reject the Null Hypothesis')
9  else:
10     print(stat, dof, p, 'Failed to reject the Null Hypothesis')
```

**Jason Brownlee** November 29, 2019 at 1:40 pm #

Thanks for sharing.

**Marc Hansen** October 31, 2019 at 7:12 am #

REPLY

Thank you for the clear explanations.

Start Machine Learning

---

**Start Machine Learning** ✕

You can master applied Machine Learning
**without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

REPLY

In the text you say: "The variable is 'sex' and the labels or factors of the variable are 'male' and 'female' in this case."

on me "The variable or factor is 'sex' and the labels or levels of the variable are 'male' and 'female' in this case."

ref https://stattrek.com/statistics/dictionary.aspx?definition=factor

Statistics for Machine Learning (7-Day Mini-Course)

**Jason Brownlee** October 31, 2019 at 7:32 am #    REPLY

A Gentle Introduction to k-fold Cross-Validation

Yes, you're right. Fixed, thanks.

How to Calculate Bootstrap Confidence Intervals For Machine Learning Results in Python

**Sandeep** December 29, 2019 at 9:18 am #

Thanks Jason. Good read it is.

A Gentle Introduction to Normality Tests in Python

**Jason Brownlee** December 30, 2019 at 5

You're welcome

Statistical Significance Tests for Comparing Machine Learning Algorithms

**James Tizard** February 12, 2020 at 2:23 pm #    REPLY

**Loving the Tutorials?**
Great tutorial, thanks!

The Statistics for Machine Learning EBook is
I'm wondering how to do chi2 test where survey respondents could select multiple answers.
where you'll find the *Really Good* stuff.
For example: which OS do you use? A) Windows B) Linux C)Mac

Results >> SEE WHAT'S INSIDE ey, 500 say windows, 400 say Mac and 200 say linux. Total is greater than the number of respondents.

Can I compare windows and mac by creating the following contingency table and running the test?

OS, Not OS
Mac 400, 600
Windows 500, 500

**Jason Brownlee** February 13, 2020 at 5:36 am #    REPLY

Good question, I'm not sure off the cuff when it comes to multiple answers. It messes up the contingency table.

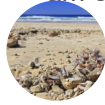You might have to hit the books or ping a statistician / post on crossvalidated.

**Never miss a tutorial:**

Eric Ren June 6, 2020 at 2:30 am #

REPLY ↩

Hi Jason,

Very nice article, clearly explained the Chi2 test. I have one question to ask. When reading the sklearn feature selection using the Chi2 test here: https://scikit-learn.org/stable/modules/feature_selection.html, I am confused by the example, in which the Iris data is used to demo the Chi 2 test on non categorical which is not even frequency or count. Is it wrong?

**Picked for you:**

Statistics for Machine Learning (7-Day Mini-Course)

Jason Brownlee June 6, 2020 at 7:58 am #

REPLY ↩

A Gentle Introduction to k-fold Cross-Validation

Probably.

How to Calculate Bootstrap Confidence Intervals For Machine Learning Results in Python

Saurabh Agarwal August 12, 2020 at 5:28 pm #

A Gentle Introduction to Normality Tests in Python

A very clearly written Nartidelty

Statistical Significance Tests for Comparing Machine Learning Algorithms

Jason Brownlee August 13, 2020 at 6:08

Thank you!

---

**Loving the Tutorials?**

The Statistics for Machine Learning EBook is where you'll find the *Really Good* stuff.

>> SEE WHAT'S INSIDE

Dhruv Modi August 20, 2020 at 8:01 pm #

REPLY ↩

Hi Jason,

Does c ndependent variables of an imbalanced data having bad rate just 1%?

---

Jason Brownlee August 21, 2020 at 6:27 am #

REPLY ↩

The test requires at least 20 examples in each cell of the contingency table I believe.

---

Hridaya Saboo June 11, 2021 at 4:04 pm #

REPLY ↩

Thank you. It is a really important and very practical question. Can you please provide more insights into this?

---

**Start Machine Learning**

✕

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

**START MY EMAIL COURSE**

Start Machine Learning

**Never miss a tutorial:**

**Jason Brownlee** June 12, 2021 at 5:23 am #     REPLY ↩

...n't have any more insight to give, perhaps check some of the references in the further reading section.

**Picked for you:**

Statistics for Machine Learning (7-Day Mini-Course)

**Fidan** September 23, 2020 at 10:26 am #     REPLY ↩

Hi Jason.

A Gentle Introduction to k-fold Cross-Validation

...at article. I have one question: If we done the chi-square test on a sample dataset and the result ...een two categorical variables are dependent. ...

population?

How to Calculate Bootstrap Confidence Intervals For Machine Learning Results in Python

**Jason Brownlee** September 23, 2020 at ...

A Gentle Introduction to Normality Tests in Python

...Yes, and than it is a sample estimate... confidence.

Statistical Significance Tests for Comparing Machine Learning Algorithms

**Kenny** October 20, 2020 at 11:14 pm #     REPLY ↩

Thanks Jason for the Good and informative article.
Sometimes I get mixed up between chi-square Goodness of fit and chi-square Tests of Independence.
Can we use the terms interchangeably or are they different to each other?

**Loving the Tutorials?**

The Statistics for Machine Learning EBook is where you'll find the *Really Good* stuff.

>> SEE WHAT'S INSIDE     :ober 21, 2020 at 6:40 am #     REPLY ↩

Same thing I believe, different use case.

**Kenny** October 22, 2020 at 8:43 pm #     REPLY ↩

Thanks Jason for the clarification.
In scipy there are 2 different function for chi-square-
1)scipy.stats.chisquare
2)scipy.stats.chi2_contingency
Do you mind telling which one to use for which use-case, please?

**Jason Brownlee** October 23, 2020 at 6:09 am #     REPLY ↩

Perhaps check the API docun...

Start Machine Learning

**Start Machine Learning**

You can master applied Machine Learning
**without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

**Never miss a tutorial:**

[in] [Twitter] **Cu** **sly** **ding** **xah** November 5, 2020 at 7:31 am #

REPLY ↩

**Picked for you:**

Hey JB,

Have you ever explored the reason why sklearn's chi2 gives different values for the test statistic and p-value compared to performing the test by hand or using chi2_contingence from scipy?

Statistics for Machine Learning (7-Day Mini-Course)

I can't seem to find a satisfactory answer, and I'm hoping the good doctor (you) might have some insight.

Cheers

A Gentle Introduction to k-fold Cross-Validation

**Jason Brownlee** November 5, 2020 at 7:

How to Calculate Bootstrap Confidence Intervals For Machine Learning Results in Python

I have not.

A Gentle Introduction to Normality Tests in Python **Hamed** November 19, 2020 at 9:57 pm #

Hey,

Statistical Significance Tests for Comparing Machine Learning Algorithms
very new in this area.

ve x_1 and x_2 and y. How can I see the dependence of y to x_1 and x_2?

x=[[1,0],[1,0],[0,1],[1,1],[1,0],[1,0],[1,0],[1,1],[0,1],[0,1]]
y=[0,0,1,1,0,0,0,1,1,1]

**Loving the Tutorials?**

The Statistics for Machine Learning EBook is where you'll find the **Really Good** stuff. **Jason Brownlee** November 20, 2020 at 6:45 am #

REPLY ↩

>> SEE WHAT'S INSIDE    e tutorial to calculate the dependence?

**Rara** July 10, 2021 at 3:01 am #

REPLY ↩

Hi! I'm new to data science. Would like to understand like how do you decide which test to use if chi-square or one-sample t-test, independent sample t-test, paired sample t test in A/B testing?

**Jason Brownlee** July 10, 2021 at 6:12 am #

REPLY ↩

Good question, see this:

https://machinelearningmastery.com/statistical-hypothesis-tests-in-python-cheat-sheet/

**Start Machine Learning**

## Leave a Reply

**Never miss a tutorial:**

(in) (twitter) (f) (email) (rss)

**Picked for you:**

Statistics for Machine Learning (7-Day Mini-Course)

A Gentle Introduction to k-fold Cross-Validation

Name (required)

Email (will not be published)

How to Calculate Bootstrap Confidence Intervals For Machine Learning Results in Python

Website

A Gentle Introduction to Normality Tests in Python

SUBMIT COMMENT

Statistical Significance Tests for Comparing Machine Learning Algorithms

Welcome!
I'm *Jason Brownlee* PhD
and I **help developers** get results with **machine learning**.
Read more

### Loving the Tutorials?

The Statistics for Machine Learning EBook is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

### Start Machine Learning ✕

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

**START MY EMAIL COURSE**

Start Machine Learning

**Never miss a tutorial:**

## Picked for you:

Statistics for Machine Learning (7-Day Mini-Course)

A Gentle Introduction to k-fold Cross-Validation

How to Calculate Bootstrap Confidence Intervals For Machine Learning Results in Python

LinkedIn | Twitter | Facebook | Newsletter | RSS

A Gentle Introduction to Normality Tests in Python

Disclaimer | Terms | Contact | Sitemap | Search

Statistical Significance Tests for Comparing Machine Learning Algorithms

## Loving the Tutorials?

The Statistics for Machine Learning EBook is where you'll find the *Really Good* stuff.

>> SEE WHAT'S INSIDE

### Start Machine Learning

✕

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Start Machine Learning