

Data Science and Big Data Analytics

Experiment 4: Chi-Square Test, ANOVA, Correlation

Aim: To perform Chi-square test, ANOVA, Pearson Correlation, Correlation Heatmap.

Description: An Analysis of Variance (ANOVA) is a statistical test employed to compare two or more means together, which are determined through the analysis of variance. One-way ANOVA tests are utilized to analyze differences between groups and determine if the differences are statistically significant.

One-way ANOVAs compare two or more independent group means, though in practice they are most often used when there are at least three independent groups.

In order to carry out an ANOVA on the Gapminder dataset, we'll need to transform some of the features, as these values in the dataset are continuous but ANOVA analyses are appropriate for situations where one variable is categorical and one variable is quantitative.

The Chi-Square test of independence is utilized when both explanatory and response variables are categorical. You likely also want to use the Chi-Square test when the explanatory variable is quantitative and the response variable is categorical, which you can do by dividing the explanatory variable into categories.

The Chi-Square test of independence is a statistical test used to analyze how significant a relationship between two categorical variables is. When a Chi-Square test is run, every category in one variable has its frequency compared against the second variable's categories. This means that the data can be displayed as a frequency table, where the rows represent the independent variables and the columns represent the dependent variables.

The Pearson Correlation test is used to analyze the strength of a relationship between two provided variables, both quantitative in nature. The value, or strength of the Pearson correlation, will be between +1 and -1.

A correlation of 1 indicates a perfect association between the variables, and the correlation is either positive or negative. Correlation coefficients near 0 indicate very weak, almost non-existent, correlations. While there are other ways of measuring correlations between two variables, such as Spearman Correlation or Kendall Rank Correlation, Pearson correlation is probably the most commonly used correlational test.

STEP – 1:

Import Modules

```
import statsmodels.formula.api as smf
```

```
import statsmodels.stats.multicomp as multi  
import scipy
```

```
from scipy.stats import pearsonr
import pandas as pd
from seaborn import regplot
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
```

HeatMap

CODE:

```
import matplotlib.pyplot as mp
import pandas as pd
import seaborn as sb

# import file with data
data = pd.read_csv("countrylifgdp.csv")
print(data.info())

print(data.corr())

# plotting correlation heatmap
dataplot = sb.heatmap(data.corr(), cmap="YlGnBu", annot=True)

# displaying heatmap
mp.show()
```

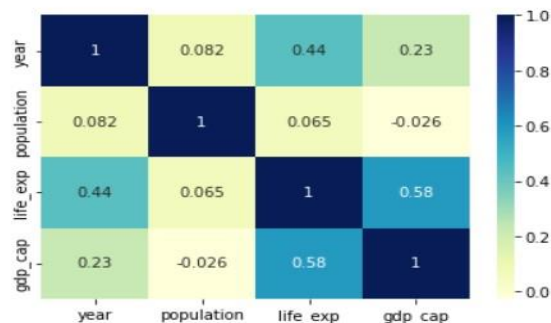
OUTPUT:

```
In [4]: # prints data that will be plotted
# columns shown here are selected by corr() since
# they are ideal for the plot
print(data.corr())

# plotting correlation heatmap
dataplot = sb.heatmap(data.corr(), cmap="YlGnBu", annot=True)

# displaying heatmap
mp.show()
```

	year	population	life_exp	gdp_cap
year	1.000000	0.082308	0.435611	0.227318
population	0.082308	1.000000	0.064955	-0.025600
life_exp	0.435611	0.064955	1.000000	0.583706
gdp_cap	0.227318	-0.025600	0.583706	1.000000



Output Analysis:

Life_exp and gdp_cap have the highest positive correlation, population and gdp_cap has negative correlation.

2) Pearson Correlation

CODE:

```
df_clean = data.dropna()
```

```
def plt_regression(x, y, data, label_1, label_2):
```

```
    reg_plot = regplot(x=x, y=y, fit_reg=True,
```

```
    data=data)plt.xlabel(label_1)
```

```
    plt.ylabel(label_2)
```

```
    plt.show()
```

```
plt_regression('year', 'population', df_clean, 'Year', 'Population')
```

```
plt_regression('year', 'life_exp', df_clean, 'Year', 'Life expectancy')
```

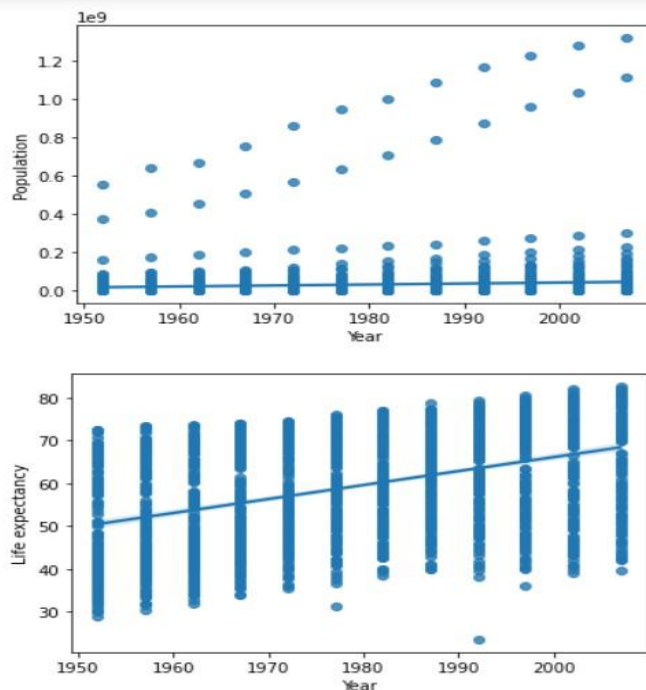
```
print('Assoc. - life expectancy and internet use rate')
```

```
print(pearsonr(df_clean['year'], df_clean['population']))
```

```
print('Assoc. - between employment rate and internet use rate')
```

```
print(pearsonr(df_clean['year'], df_clean['life_exp']))
```

OUTPUT:



```
Assoc. - life expectancy and internet use rate  
(0.08230807783830922, 0.0006716006496141012)  
Assoc. - between employment rate and internet use rate  
(0.4356112240540735, 7.54679462559958e-80)
```

Output Analysis:

As we saw in the heatmap year, population and year, life_Exp have positive correlation in the same way pearson graph also has a positive slope.

3) ANOVA

CODE:

```
def bin(dataframe, cols):  
    # Create new columns that store the binned data  
    for col in cols:  
        new_col_name = "{}_bins".format(col)  
        dataframe[new_col_name] = pd.qcut(dataframe[col], 10)  
  
df3 = data.copy()  
cols = ['year']  
cols1 = ['population']  
norm_cols = ['population', 'gdp_cap', 'life_exp']  
# This creates new columns filled with the binned column data  
bin(df3, cols)  
bin(df3, norm_cols)  
bin(df3, cols1)  
  
anova_df = df3[['year_bins', 'population']].dropna()  
  
relate_df = df3[['year_bins', 'population']]  
  
anova = smf.ols(formula='population ~ C(year_bins)', data=anova_df).fit()  
  
print(anova.summary())  
  
# We may also want to check the mean and standard deviation for the groups  
mean = relate_df.groupby("year_bins").mean()  
sd = relate_df.groupby("year_bins").std()  
print(mean)  
print(sd)
```

OUTPUT:

```
st8888/notebooks/Untitled8.ipynb
```

jupyter Untitled8 Last Checkpoint: 05/31/2022 (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 C

```
print(mean)
print(sd)
```

OLS Regression Results

Dep. Variable:	population	R-squared:	0.007
Model:	OLS	Adj. R-squared:	0.001
Method:	Least Squares	F-statistic:	1.282
Date:	Thu, 15 Sep 2022	Prob (F-statistic):	0.242
Time:	11:13:57	Log-Likelihood:	-33902.
No. Observations:	1704	AIC:	6.782e+04
Df Residuals:	1694	BIC:	6.788e+04
Df Model:	9		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
Intercept	1.786e+07	6.29e+06	2.837	0.005	5.51e+06	3.02e+07
C(year_bins)[T.Interval(1957.0, 1962.0, closed='right')]	2.564e+06	1.09e+07	0.235	0.814	-1.88e+07	2.39e+07
C(year_bins)[T.Interval(1962.0, 1967.0, closed='right')]	4.801e+06	1.09e+07	0.440	0.660	-1.66e+07	2.62e+07
C(year_bins)[T.Interval(1967.0, 1972.0, closed='right')]	7.333e+06	1.09e+07	0.673	0.501	-1.41e+07	2.87e+07
C(year_bins)[T.Interval(1972.0, 1979.5, closed='right')]	9.819e+06	1.09e+07	0.901	0.368	-1.16e+07	3.12e+07
C(year_bins)[T.Interval(1979.5, 1987.0, closed='right')]	1.377e+07	8.9e+06	1.546	0.122	-3.69e+06	3.12e+07
C(year_bins)[T.Interval(1987.0, 1992.0, closed='right')]	1.813e+07	1.09e+07	1.663	0.096	-3.25e+06	3.95e+07
C(year_bins)[T.Interval(1992.0, 1997.0, closed='right')]	2.098e+07	1.09e+07	1.925	0.054	-4.01e+05	4.24e+07
C(year_bins)[T.Interval(1997.0, 2002.0, closed='right')]	2.36e+07	1.09e+07	2.165	0.031	2.22e+06	4.5e+07
C(year_bins)[T.Interval(2002.0, 2007.0, closed='right')]	2.616e+07	1.09e+07	2.400	0.017	4.78e+06	4.75e+07

Omnibus: 2403.699 Durbin-Watson: 0.187
Prob(Omnibus): 0.000 Jarque-Bera (JB): 438135.300
Skew: 8.286 Prob(JB): 0.00
Kurtosis: 79.787 Cond. No. 8.74

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

year_bins
population

Activate
Go to Settings

earch 23°C Mostly cloudy

```

In [29]: count_table = pd.crosstab(df3['year_bins'], df3['population_bins'])
print(count_table)

def chi_sq_test(table):

    print("Results for:")
    print(str(table))

    # Get column percentages

```

OUTPUT Analysis:

P-Value is $0.242 > 0.05$ [Level of Significance] which implies that we don't have a significant relationship between population and year.

4) Chi Square Test

CODE:

```
count_table = pd.crosstab(df3['year_bins'], df3['population_bins'])
```

```
print(count_table)
```

```
def chi_sq_test(table):
```

```
    print("Results for:")
```

```
    print(str(table))
```

```
    # Get column percentages
```

```
col_sum = table.sum(axis=0)
col_percents = table/col_sum
print(col_percents)

chi_square = scipy.stats.chi2_contingency(table)
print("Chi-square value, p-value, expected_counts")
print(chi_square)

print()

print("Initial Chi-
square:")
chi_sq_test(count_table)
print(" ")

chi_sq_test(count_table_3)
chi_sq_test(count_table_4)
chi_sq_test(count_table_5)
chi_sq_test(count_table_6)
```

OUTPUT:

