



Python Programming

Statistics

Dr. Chun-Hsiang Chan
Department of Geography
National Taiwan Normal University

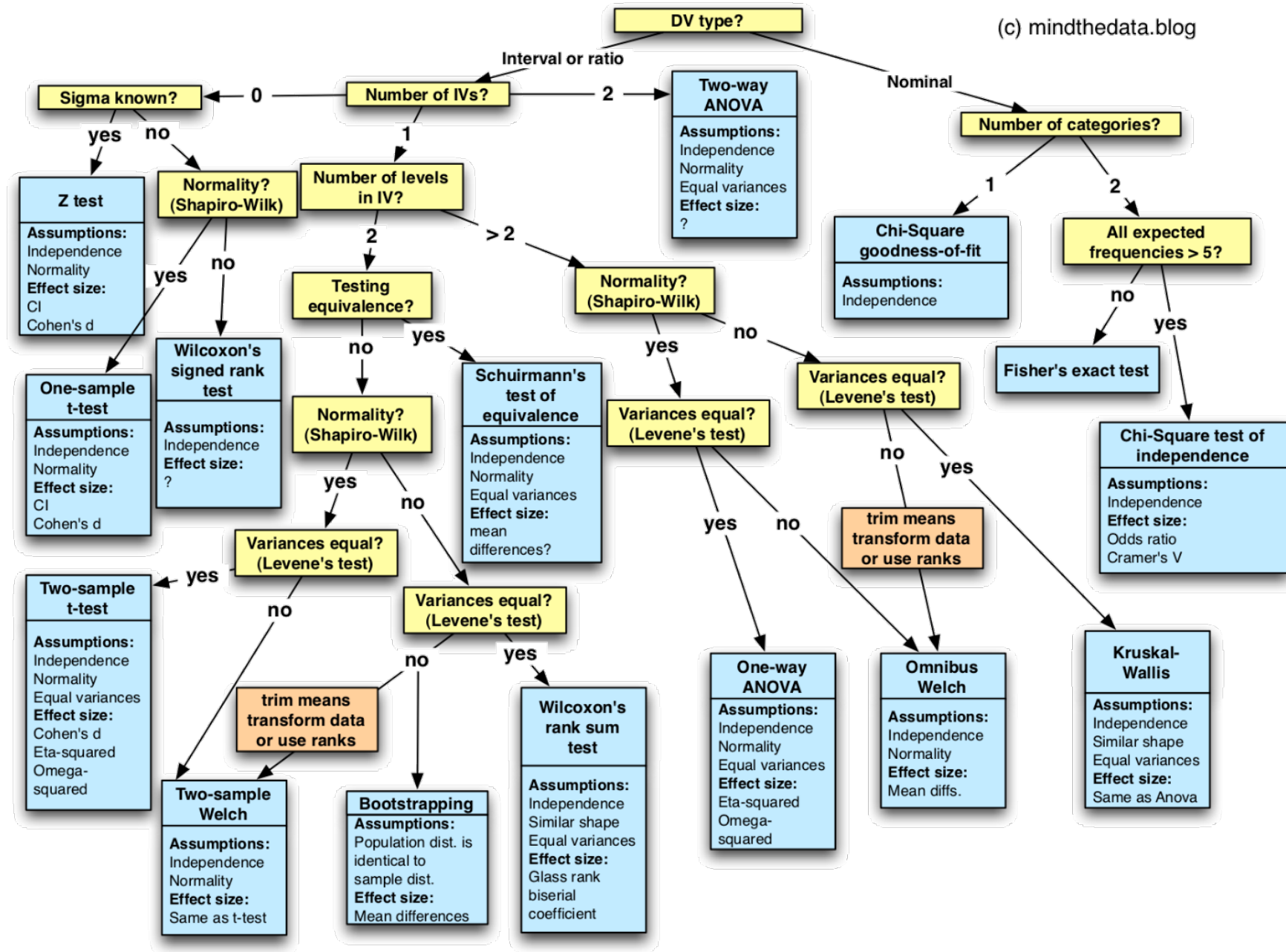


Outlines

- Statistical Analysis Road Map
- Normality Test
- F Test
- Levene's
- T Test
- Correlation Analysis
- Assignment

Statistical Analysis Road Map

(c) mindthedata.blog



Statistical Analysis Road Map

Table 1 Choice of statistical test from paired or matched observation

Variable	Test
Nominal	McNemar's Test
Ordinal (ordered categories)	Wilcoxon
Quantitative (discrete or non-normal)	Wilcoxon
Quantitative (normal)	Paired t-test

Table 2 Parametric and nonparametric tests for comparing two or more groups

Parametric Test	Situation	Nonparametric Test
t-test	Two independent population	Wilcoxon rank sum test
t-test		Mann-Whitney U test
One way analysis of variance	Three or more populations	Kruskal Wallis test
Paired t-test	Paired population	Sign test
		Wilcoxon rank sign test
Pearson correlation	Correlation	Spearman correlation

Source:

<https://www.healthknowledge.org.uk/public-health-textbook/research-methods/1b-statistical-methods/parametric-nonparametric-tests>

Statistical Analysis Road Map

Table 3 Choice of statistical test for independent observations

Input variable	Outcome variable						
		Nominal	Categorical (>2)	Ordinal	Quantitative Discrete	Quantitative Non-normal	Quantitative Normal
Nominal		χ^2 or Fisher's	χ^2	χ^2 -trend or Mann-Whitney	Mann-Whitney	Mann-Whitney or log-rank ^a	T-test
Categorical (>2)		χ^2	χ^2	Kruskal-Wallis ^b	Kruskal-Wallis ^b	Kruskal-Wallis ^b	ANOVA ^c
Ordinal		χ^2 -trend or Mann-Whitney	^e	Spearman rank	Spearman rank	Spearman rank	Spearman rank or Linear regression ^d
Quantitative Discrete		Logistic regression	^e	^e	Spearman rank	Spearman rank	Spearman rank or Linear regression ^d
Quantitative Non-normal		Logistic regression	^e	^e	^e	Plot data and Pearson or Spearman rank	Plot data and Pearson or Spearman rank and Linear regression
Quantitative Normal		Logistic regression	^e	^e	^e	Linear regression ^d	Pearson or Linear regression

^a If data are censored. ^b The Kruskal-Wallis test is used for comparing ordinal or non-Normal variables for more than two groups, and is a generalisation of the Mann-Whitney U test. ^c Analysis of variance is a general technique, and one version (one way analysis of variance) is used to compare Normally distributed variables for more than two groups, and is the parametric equivalent of the Kruskal-Wallis test. ^d If the outcome variable is the dependent variable, then provided the residuals (the differences between the observed values and the predicted responses from regression) are plausibly Normally distributed, then the distribution of the independent variable is not important. ^e There are a number of more advanced techniques, such as Poisson regression, for dealing with these situations. However, they require certain assumptions and it is often easier to either dichotomise the outcome variable or treat it as continuous.

Source: <https://www.healthknowledge.org.uk/public-health-textbook/research-methods/1b-statistical-methods/parametric-nonparametric-tests>

Normality Test :: Shapiro–Wilk Test

- The Shapiro–Wilk test is a test of normality in frequentist statistics. **The Shapiro–Wilk test tests the null hypothesis that a sample x_1, \dots, x_n came from a normally distributed population.** The test statistic is

$$W = \frac{\left(\sum_{i=1}^n a_i x_{(i)}\right)^2}{\sum_{i=1}^n (x_i - \bar{x})^2}, \text{ where } x_i$$

(with parentheses enclosing the subscript index i ; not to be confused with x_i) is the i th order statistic, i.e., the i th-smallest number in the sample; $\bar{x} = \frac{1}{n}(x_1 + x_2 + \dots + x_n)$, $8 \leq n \leq 50$ is the sample mean.

Shapiro–Wilk Test

$$W = \frac{\left(\sum_{i=1}^n a_i x_{(i)}\right)^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$\bar{x} = \frac{1}{n}(x_1 + x_2 + \dots + x_n)$ is the sample mean.

- The coefficient a_i are given by: $(a_1, \dots, a_n) = \frac{m^T V^{-1}}{C}$, where C is a vector norm: $C = \|V^T m\| = \sqrt{m^T V^{-1} V^{-1} m}$ and the vector m , $m = (m_1, \dots, m_n)^T$ is made of the expected values of the order statistics of independent and identically distributed random variables sampled from the standard normal distribution; finally, V is the covariance matrix of those normal order statistics.

Shapiro–Wilk Test

$\alpha = 0.05$				$n = 15$				Data Pairs					
Raw Data		Sorted Data		a value		Upper value		Lower value		Difference	a*Difference		
1	20	1	18	a1	0.5150	#15	22	#01	18	4	2.0600		
2	19	2	18	a2	0.3306	#14	21	#02	18	3	0.9918		
3	18	3	18	a3	0.2495	#13	21	#03	18	3	0.7485		
4	19	4	18	a4	0.1878	#12	21	#04	18	3	0.5634		
5	22	5	19	a5	0.1353	#11	20	#05	19	1	0.1353		
6	18	6	19	a6	0.0880	#10	20	#06	19	1	0.0880		
7	21	7	19	a7	0.0433	#09	19	#07	19	0	0.0000		
8	19	8	19										
9	21	9	19										
10	18	10	20										
11	18	11	20										
12	19	12	21										
13	20	13	21										
14	21	14	21										
15	19	15	22										

$\sum_{i=1}^n a_i x_{(i)}$	4.59	$\left(\sum_{i=1}^n a_i x_{(i)}\right)^2$	21.0681
$\sum_{i=1}^n (x_i - \bar{x})^2$	23.733	W	0.886541
		W critical	0.881

Chun-Hsiang Chan (2026)

Kolmogorov-Smirnov Test

- The Kolmogorov–Smirnov test (K-S test or KS test) is a **nonparametric test** of the equality of continuous (or discontinuous), one-dimensional probability distributions that can be used to compare a sample with a reference probability distribution (one-sample K–S test), or to compare two samples (two-sample K–S test), where n is larger than 50.
- **The null distribution of this statistic is calculated under the null hypothesis that the sample is drawn from the reference distribution (in the one-sample case) or that the samples are drawn from the same distribution (in the two-sample case).**

Kolmogorov-Smirnov Test

- The two-sample K–S test is one of the most useful and general **nonparametric** methods for comparing two samples, as it is **sensitive to differences in both locations and shape of the empirical cumulative distribution functions of the two samples.**

$$F_n = \frac{\text{number of (elements in the sample } \leq x)}{n} = \frac{1}{n} \sum_{i=1}^n 1_{(-\infty, x]}(X_i),$$

- The Kolmogorov-Smirnov statistic for a given cumulative distribution function $F(x)$ is

$$D_n = \sup_x |F_n(x) - F(x)|,$$

where \sup_x is the supremum of the set of distances.

- Intuitively, the statistic takes the largest absolute difference between the two distribution functions across all x values.

Normality Test

```
# import packages  
import os  
import numpy as np  
import pandas as pd  
from scipy.stats import shapiro  
from statsmodels.stats.diagnostic import kstest_normal  
from scipy.stats import normaltest
```

Normality Test

```
# load data
```

```
# source: https://en.wikipedia.org/wiki/2024\_United\_States\_presidential\_election#Results
```

```
df = pd.read_excel('US Presidential Election.xlsx')
```

```
df.head()
```

	States	Trump_Votes	Trump_Percent	Trump_Electoral_Votes	Harris_Votes	Harris_Percent	Harris_Electoral_Votes	Total_Votes
0	Alabama	1462616	0.646	9	772412	0.341	–	2265090
1	Alaska	184458	0.545	3	140026	0.414	–	338177
2	Arizona	1770242	0.522	11	1582860	0.467	–	3390161
3	Arkansas	759241	0.642	6	396905	0.336	–	1182676
4	California	6081697	0.383	–	9276179	0.585	54	15865475

Normality Test

```
# detect NA
```

```
print('number of NA:', np.sum(pd.isna(df).values, axis=0))
```

```
number of NA: [0 0 0 0 0 0 0 0]
```

```
# detect dash
```

```
print('number of -: ', np.sum(df=='-', axis=0).values)
```

```
number of -: [ 0  0  0 20  0  0 31  0]
```

```
# fill dash with 0
```

```
df[df=='-'] = 0
```

```
print('number of -: ', np.sum(df=='-', axis=0).values)
```

```
number of -: [0 0 0 0 0 0 0 0]
```

Normality Test

```
# Shapiro-Wilk normality test
w, p = shapiro(df['Trump_Votes'].astype(np.float64))
print('w:',w, 'p:', p)
w: 0.7904571294784546 p: 4.3635461111080076e-07

# Kolmogorov-Smirnov normality test
ks, pval = kstest_normal(df['Trump_Votes'].astype(np.float64),
dist='norm', pvalmethod='table')
print('Kolmogorov-Smirnov Test:', ks, pval)
Kolmogorov-Smirnov Test: 0.17761281609974033 0.000999999999999998899

# D'Agostino and Pearson's: test that combines skew and kurtosis to
produce an omnibus test of normality
nor, pval = normaltest(df['Trump_Votes'].astype(np.float64))
print('D'Agostino and Pearson's:',nor, pval)
D'Agostino and Pearson's: 28.85096438328366 5.433662821590572e-07
```

F Test

- The definitional equation of sample variance is

$$s^2 = \frac{1}{n-1} \sum_i (y_i - \bar{y})^2$$

- The fundamental technique is a partitioning of the total sum of squares SS into components related to the effects used in the model.

$$SS_{Total} = SS_{treatments} + SS_{Error}$$
$$\sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y} \dots)^2 = \sum_{i=1}^k n_i (\bar{y}_{i.} - \bar{y} \dots)^2 + \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i.})^2$$

F Test

- The F -test is used for comparing the factors of the total deviation. For example, in one-way, or single-factor ANOVA, statistical significance is tested for by comparing the F test statistic.

$$F = \frac{\text{variance between treatments}}{\text{variance within treatments}}$$
$$F = \frac{MS_{Treatments}}{MS_{Error}} = \frac{\frac{SS_{Treatments}}{I - 1}}{\frac{SS_{Error}}{n_T - 1}}$$

where MS is mean square, I is the number of treatments and n_T is the total number of cases.

Levene's Test

- Levene's test is an inferential statistic used to assess the equality of variances for a variable calculated for two or more groups.
- Levene's test has been used in the past before a comparison of means to inform the decision on whether to use a pooled t-test or the Welch's t-test for two sample tests or analysis of variance or Welch's modified oneway ANOVA for multi-level tests.
- Levene's test is equivalent to a 1-way between-groups analysis of variance (ANOVA) with the dependent variable being the absolute value of the difference between a score and the mean of the group to which the score belongs (shown below as $Z_{ij} = |Y_{ij} - \bar{Y}_{i.}|$).

Levene's Test

- The test statistic, W , is equivalent to the F statistic that would be produced by such an ANOVA, and is defined as follows:

$$W = \frac{N - k}{k - 1} \cdot \frac{\sum_{i=1}^k N_i (Z_{i.} - Z_{..})^2}{\sum_{i=1}^k \sum_{j=1}^k (Z_{ij} - Z_{i.})^2}$$

where

- k is the number of different groups to which the sampled cases belong
- N_i is the number of cases in the i th group
- Y_{ij} is the value of the measured variable for the j th case from the i th group
- $Z_{ij} = \begin{cases} |Y_{ij} - \bar{Y}_{i.}|, & \bar{Y}_{i.} \text{ is a mean of the } i\text{-th group} \\ |Y_{ij} - \tilde{Y}_{i.}|, & \tilde{Y}_{i.} \text{ is a median of the } i\text{-th group} \end{cases}$

Variance Test

```
from scipy.stats import levene
from scipy.stats import f_oneway

# F test
F = f_oneway(df['Trump_Votes'], df['Harris_Votes'])
print('F test: statistics=', F.statistic, 'p-value=', F.pvalue)
F test: statistics= 0.020419207963903285 p-value= 0.8866600125553097

# Levene's test
print("Levene's test:", levene(df['Trump_Votes'],
                               df['Harris_Votes']))
Levene's test: LeveneResult(statistic=0.06853717697239912, pvalue=0.7940171597643334)
```

T Test

- In addition to one-sample t-test, there are three types of t-test, including paired t-test, and two-sample independent t-test (assume that the variance of two samples or populations are [not] equal).
- In the following slides, we will give some examples to show their differences.

Question X

How do we determine whether the variances between two samples or populations are equal?

Paired T Test

- If the two samples or populations are from matched or paired sources or a replicated measurement, you must select a paired t-test.

$$t = \frac{\overline{X}_D - \mu_0}{\frac{s_D}{\sqrt{n}}}$$

\overline{X}_D and s_D are the average and standard deviation of the differences between all pairs, the constant μ_0 is zero if we want to test whether the average of the difference is significantly different, and n is the number of pairs.

Source: https://en.wikipedia.org/wiki/Student%27s_t-test

Source: <https://www.statisticssolutions.com/free-resources/directory-of-statistical-analyses/paired-sample-t-test/>

Source: <https://www.omnicalculator.com/statistics/t-test#p-value-from-t-test>

Two-sample Independent T-test

- If the variance of two samples or populations are equal (or very similar).

$$t = \frac{\bar{x}_1 - \bar{x}_2 - \mu_0}{s_p \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

s_p is the pooled standard deviation, defined by

$$s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

Source: https://en.wikipedia.org/wiki/Student%27s_t-test

Source: <https://www.statisticssolutions.com/free-resources/directory-of-statistical-analyses/paired-sample-t-test/>

Source: <https://www.omnicalculator.com/statistics/t-test#p-value-from-t-test>

Two-sample Independent T-test

- If the variance of two samples or populations is **unequal** (or very dissimilar), refer to Welch's t-test.

$$t = \frac{\bar{x}_1 - \bar{x}_2 - \mu_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Source: https://en.wikipedia.org/wiki/Student%27s_t-test

Source: <https://www.statisticssolutions.com/free-resources/directory-of-statistical-analyses/paired-sample-t-test/>

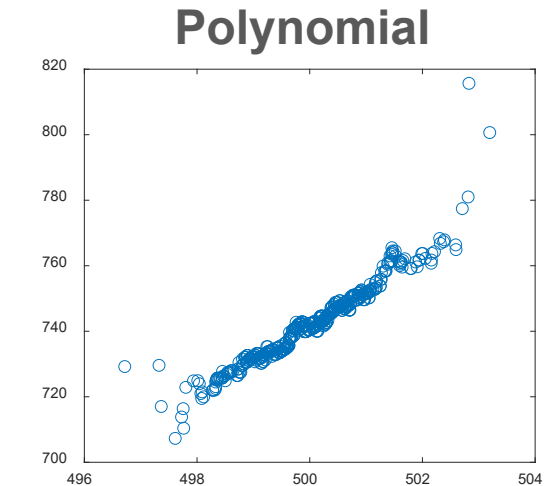
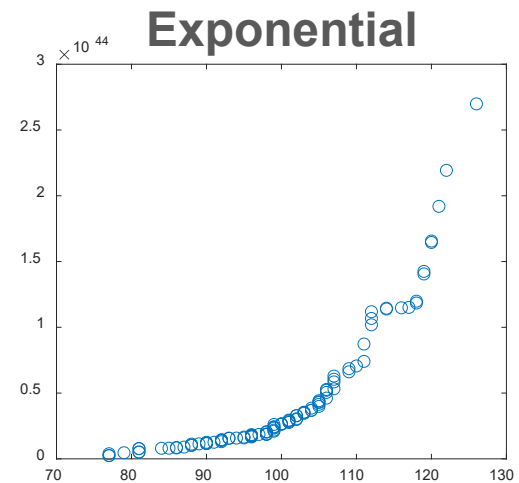
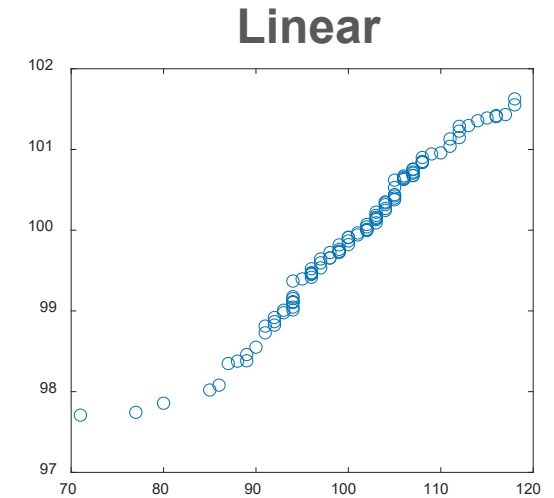
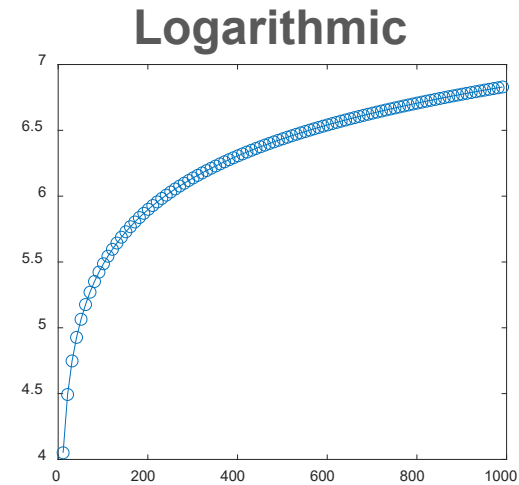
Source: <https://www.omnicalculator.com/statistics/t-test#p-value-from-t-test>

Code :: Two-sample Independent T-test

```
from scipy.stats import ttest_ind
# two independent variables - t test
print('Votes:', ttest_ind(df['Trump_Votes'], df['Harris_Votes'], axis=0,
                          equal_var=True))
Votes: Ttest_indResult(statistic=0.14289579407352498, pvalue=0.886660012555303)
```

Correlation Analysis

- Correlation analysis is an inferential statistic to describe the relationship or association between one variable and another.
- Most formulae for correlation analyses are developed for linear relationships; therefore, other relationships (e.g., logistic, exponential, and cubic) are unsuitable. A nonlinear relationship could adopt the performance of curve-fitting results.



Correlation Analysis

Variable Y/X	Quantitative X	Ordinal X	Nominal X
Quantitative Y	Pearson r	Biserial r_b	Point Biserial r_{pb}
Ordinal Y	Biserial r_b	Spearman ρ / Tetrachoric r_{tet}	Rank Biserial r_{rb}
Nominal Y	Point Biserial r_{pb}	Rank Biserial r_{rb}	Phi, L, C, V, Lambda

- There are two important outcomes from correlation analyses: significance and coefficient.
- **Significant correlation:** the consistency of the association between one variable and the other.
- **Coefficient of correlation:** the direction (i.e., positive or negative) and magnitude (i.e., value) of correlation between variables.

Pearson Correlation Coefficient r

- **Pearson correlation coefficient**, also known as Pearson product-moment correlation coefficient (PPMCC), is to measure the linear correlation between two variables or data.
- The definition of Pearson correlation coefficient is calculated by the covariance of the two variables divided by the product of their standard deviations. Its value ranges from -1 to +1.

$$\rho = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}, \text{ when it is applied for population}$$

$$r = \frac{\text{cov}(X, Y)}{S_X S_Y}, \text{ when it is applied for sample}$$

Pearson Correlation Coefficient r

$$\rho = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}, \text{ where } \text{cov}(X, Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)],$$

$$\text{then } \rho = \frac{\mathbb{E}[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}, \text{ and ...}$$

$$\mu_X = \mathbb{E}[X]; \mu_Y = \mathbb{E}[Y];$$

$$\sigma_X^2 = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2;$$

$$\sigma_Y^2 = \mathbb{E}[(Y - \mathbb{E}[Y])^2] = \mathbb{E}[Y^2] - (\mathbb{E}[Y])^2;$$

$$\mathbb{E}[(X - \mu_X)(Y - \mu_Y)] = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$

$$\rho_{XY} = \frac{\mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]}{\sqrt{\mathbb{E}[X^2] - (\mathbb{E}[X])^2} \sqrt{\mathbb{E}[Y^2] - (\mathbb{E}[Y])^2}}$$

Pearson Correlation Coefficient r

- Testing using t -distribution with degrees of freedom $n - 2$, where standard error is denoted as,

$$\sigma_r = \sqrt{\frac{1 - r^2}{n - 2}}$$

- Therefore, t value is ...

$$t = \frac{r}{\sigma_r} = r \sqrt{\frac{n - 2}{1 - r^2}}$$

The inverse function for determining the critical values for r is ...

$$r = \frac{t}{\sqrt{n - 2 + t^2}}$$

Pearson Correlation Coefficient r

Sleeping/Day, X_i	Relax/Day, Y_i
7.5	1
8	12
9.1	2
6	10
10	5
8.4	6.1
9.1	7
2.4	8.2
6.7	7
6.8	6
9	4.5

$$\rho = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{-2.2370}{2.0011 \times 3.0509} = -0.3664$$

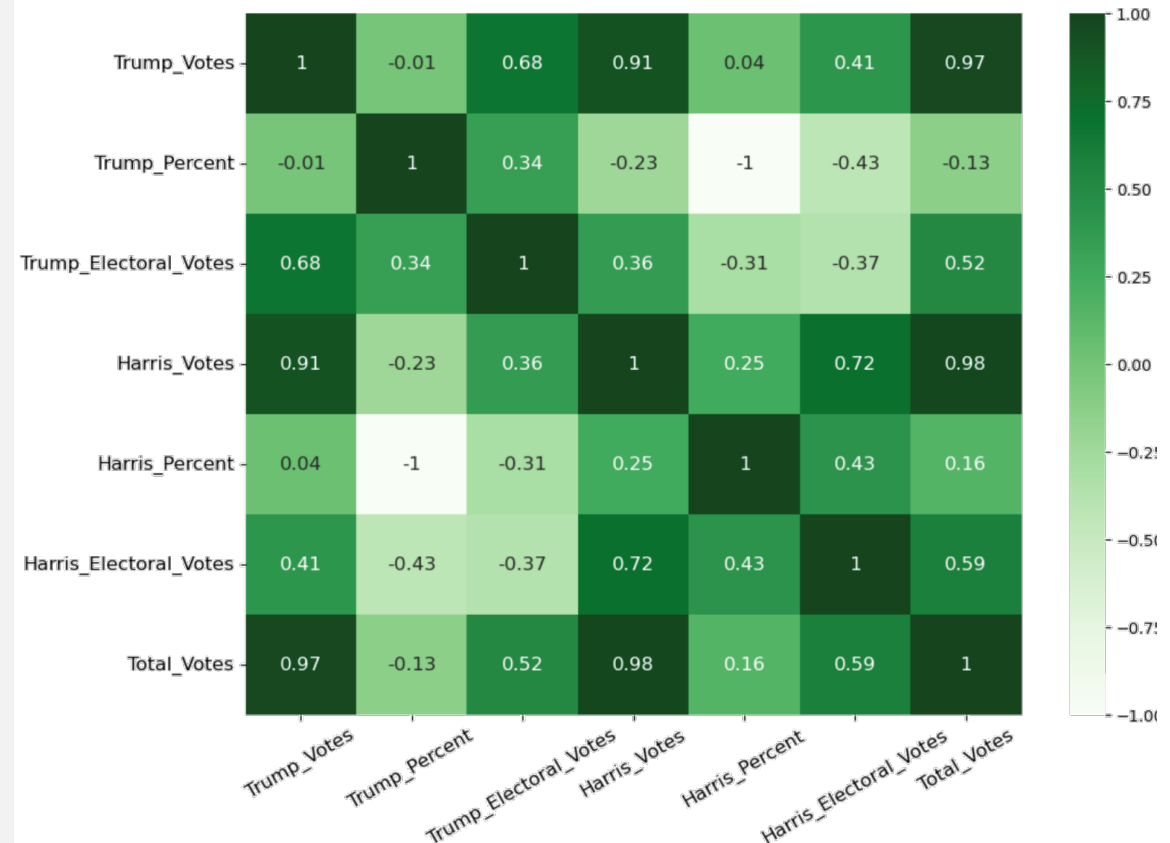
$$t = \frac{r}{\sigma_r} = r \sqrt{\frac{n-2}{1-r^2}}$$

$$t = -0.3664 \times \sqrt{\frac{11-2}{1-(-0.3664)^2}} = -1.18143$$

$$t_{-1.18143, 9} = 0.133853$$

Pearson Correlation Coefficient r

```
import matplotlib.pyplot as plt
import seaborn as sns
# calculate Pearson
# correlation matrix
corr = df1.iloc[:,1:].corr()
corr = np.round(corr,2)
# visualize correlation results
plt.figure(figsize=(11,8))
sns.heatmap(corr,
            cmap="Greens",annot=True,
            annot_kws = {'size': 12})
plt.xticks(fontsize=12, rotation=30)
plt.yticks(fontsize=12)
plt.show()
```



Spearman Rank Correlation ρ

- **The Spearman correlation coefficient** is defined from the Pearson correlation coefficient between the rank variables.
- For a sample of size n , the n raw scores X_i, Y_i are converted to ranks $R(X_i), R(Y_i)$, and r_s is computed as

$$r_s = \rho_{R(X)R(Y)} = \frac{\text{cov}(R(X), R(Y))}{\sigma_{R(X)}\sigma_{R(Y)}}$$

where ρ denotes the Pearson correlation coefficient with rank variables, $\text{cov}(R(X), R(Y))$ is the covariance of the rank variables, $\sigma_{R(X)}$ and $\sigma_{R(Y)}$ are the standard deviations of the rank variables.

Spearman Rank Correlation ρ

- Only if all n ranks are distinct integers, it can be computed using the popular formula.

$$r_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

where $d_i = R(X_i) - R(Y_i)$ is the difference between the two ranks of each observation, n is the number of observations.

- Significance measurement could be obtained from t distribution, where degree of freedom is $n - 2$.

$$t = r_s \sqrt{\frac{n - 2}{1 - r^2}}$$

Spearman Rank Correlation ρ

PR, X_i	Reading/Day, Y_i	x_i rank	y_i rank	d_i	d_i^2
60	1	1	2	-1	1
65	1	2	2	0	0
71	2	3	4	-1	1
75	1	4	2	2	4
78	5	5	6	-1	1
81	6	6	7.5	-1.5	2.25
85	7	7	9.5	-2.5	6.25
89	8	8	11	-3	9
91	7	9	9.5	-0.5	0.25
95	6	10	7.5	2.5	6.25
99	4	11	5	6	36

- **Spearman Rank Corr.**

$$r_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

$$r_s = 1 - \frac{6 \times 67}{11(11^2 - 1)}$$

$$r_s = 0.695455$$

$$t = r_s \sqrt{\frac{n - 2}{1 - r^2}}$$

$$= 2.870262$$

$$t_{2.87, df=9, two} = 0.018469$$

Pearson Correlation Coefficient r

Spearman Rank Correlation ρ

```
# correlation analysis
# Pearson Correlation: linear parametric
from scipy.stats import pearsonr
rho_p, pval_p = pearsonr(df['Trump_Votes'], df['Harris_Votes'])
print('Pearson Correlation:', rho_p, pval_p)
Pearson Correlation: 0.9138721222845388 8.127625915117507e-21

# Spearman Correlation: nonlinear or nonparametric
from scipy.stats import spearmanr
rho_s, pval_s = spearmanr(df['Trump_Votes'], df['Harris_Votes'])
print('Spearman Correlation:', rho_s, pval_s)
Spearman Correlation: 0.9351131221719456 1.0108735645658283e-23
```

Assignment #01

Format

```
def myPearsonR(xData, yData):  
    # annotation  
    ...  
    ...  
    ...  
    return r, pval
```

- Define a function that produces a Pearson correlation coefficient and *p-value* with x and y data.
- You cannot directly call `scipy.stats.pearsonr` or any built-in correlation functions.

The End

Thank you for your attention!

Email: chchan@ntnu.edu.tw

Website: <https://toodou.github.io/>

