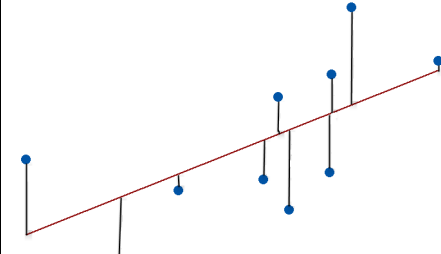# Harold's Statistics
## ANOVA and Chi-Squared
## Cheat Sheet
24 June 2022

## Analysis of Variance (ANOVA) for Linear Regression

| Term | Formula | Description |
|------|---------|-------------|
| **Problem Statement** | *How do we measure both the explained and unexplained variances?* | |
| **Residual Sum of Squares (SSE)** | $SSE = \sum e_i^2 = \sum \left(Y_i - \hat{Y}_i\right)^2$ | Estimator errors |
| **Residual Degrees of Freedom (df)** | $df = n - p$ | Number of regression parameters |
| **Residual Mean Square (MSE)** | $MSE = \dfrac{SSE}{n - p}$ | Measures the amount of error in statistical models. 0 = no error.  |
| **Residual Standard Error (s)** | $s = \sqrt{MSE}$ | Estimates the standard deviation of the residuals |
| **Python** | ```python
import pandas as pd
import statsmodels.api as sm
from statsmodels.formula.api import ols
scores = pd.read_csv('ExamScores.csv')

# Creates a linear regression model
results = ols('Exam4 ~ Exam1', data=scores).fit()
print(results.summary())

# The explained and unexplained variance can be obtained
# from the analysis of variance table
aov_table = sm.stats.anova_lm(results, typ=2)
print(aov_table)
```<br>`          sum_sq    df      F      PR(>F)`<br>`Exam1    217.166351  1.0   3.517655  0.066808`<br>`Residual 2963.333649 48.0        NaN       NaN` | |

# One-Way Analysis of Variance (One-Way ANOVA)

| Test | Hypotheses | Research Question |
|---|---|---|
| **Problem Statement** | *How do we determine whether a statistically significant difference exists among the means of three or more groups or populations?* | |
| **One-Way ANOVA** | Tests for an association between a single <u>categorical</u> predictor variable and a response variable. | |
| **Factor** | A categorical predictor variable. | |
| **Level** | A possible value of a factor. | |
| **F-test** | 1. Hypotheses<br>  $H_0 : \mu_1 = \mu_2 = \ldots = \mu_k$<br>  $H_a : \mu_i \neq \mu_j$, for some $i \neq j$<br>2. F-test (F statistic)<br>  $$F = \dfrac{between\ group\ variance}{within\ group\ variance}$$<br>3. p-value (Prob (F-statistic))<br>4. Significance level $\alpha$ = 0.05 | Assumes the k population means ($\mu$) are independent.<br>If p-value < $\alpha$, then reject $H_0$.<br>If p-value > $\alpha$, then favor $H_0$. |
| **Python 1** | <pre>`import pandas as pd`<br>`import scipy.stats as st`<br>`scores = pd.read_csv('ExamScores.csv')`<br>`# Statistics of each exam`<br>`exam1_score = scores[['Exam1']]`<br>`exam2_score = scores[['Exam2']]`<br>`exam3_score = scores[['Exam3']]`<br>`exam4_score = scores[['Exam4']]`<br>`print(st.f_oneway(exam1_score, exam2_score, exam3_score, exam4_score))`</pre><br>`F_onewayResult(`<br>`statistic=array([ 3.85696089]),`<br>`pvalue=array([ 0.01034867]))` | |
| **Python 2** | <pre>`import statsmodels.api as sm`<br>`import pandas as pd`<br>`from statsmodels.formula.api import ols`<br>`df = pd.read_csv('ExamScoresGrouped.csv')`<br>`mod = ols('Scores ~ Exam',data=df).fit()`<br>`aov_table = sm.stats.anova_lm(mod, typ=2)`<br>`print(aov_table)`</pre><br>`          sum_sq     df        F    PR(>F)`<br>`Exam      2400.735    3.0  3.856961  0.010349`<br>`Residual  40666.220  196.0     NaN       NaN` | |

## Post-Hoc Tests

| Test | Hypotheses | Research Question |
|------|------------|-------------------|
| **Problem Statement** | *If the ANOVA null hypothesis is rejected, further analysis is required because the F-test does not determine which groups have different means.*<br>*Helps us identify groups that are significantly different than others.* | |
| **Post-Hoc Analysis** | Determines which groups have different means, which group has the highest or lowest mean, and other relationships between the groups. | |
| **Tukey Honestly Significant Difference (HSD)** | A procedure that gives the 95% confidence intervals for the mean difference between pairwise groups and determines which mean difference is statistically significant. | |
| **Python** | <br>```python\nimport pandas as pd\nfrom statsmodels.stats.multicomp import (pairwise_tukeyhsd,\nMultiComparison)\ndf = pd.read_csv(ExamScoresGrouped.csv')\nmod = MultiComparison(df['Scores'], df['Exam'])\nprint(mod.tukeyhsd())\n```<br><br>```\nMultiple Comparison of Means - Tukey HSD,FWER=0.05\n=================================================\ngroup1 group2 meandiff  lower     upper   reject\n-------------------------------------------------\nExam1  Exam2    -3.3   -10.7652  4.1652  False\nExam1  Exam3    -9.36  -16.8252 -1.8948  True\n``` | |

# Chi-Square Tests for Comparing Categorical Variables

| Goodness-of-Fit Test – Chi-Square | | |
|---|---|---|
| **Problem Statement** | *The chi-square distribution is used to test how close the distribution of a population is to a theoretical distribution.*<br>*The chi-squared test statistic measures how different the observed <u>counts</u> are compared to the expected counts, assuming the null hypothesis is true.* | |
| $\chi^2$**-test** | 1. Hypotheses (two-sided)<br>   $H_0$ : The random variable follows the expected distribution<br>   $H_a$ : The random variable does not follow the expected distribution<br>2. $\chi^2$-test statistic<br>3. p-value (Prob ($\chi^2$-statistic))<br>4. Significance level α = 0.05 | For a chi-squared distribution<br>If p-value < α, then reject $H_0$.<br>If p-value > α, then favor $H_0$.<br><br>If $H_0$ is rejected, insufficient evidence exists to conclude that the distribution does not follow the expected distribution. |
| **Expected Frequencies for a Chi-Square** | $$E = np$$ | $p = proportion$<br>$n = sample\ size$ |
| **Chi-Square Test Statistic** | $$\chi^2 = \frac{(n-1)s^2}{\sigma^2}$$<br><br>$$\chi^2 = \sum \frac{(observed - expected)^2}{expected}$$ | Large $\chi^2$ values are evidence against the null hypothesis, which states that the percentages of observed and expected match (as in, any differences are attributed to chance). |
| **Degrees of Freedom** | $$df = k - 1$$ | $k = number\ of\ possible\ values$<br>$(categories)\ for\ the\ variable$<br>$under\ consideration$ |
| **Python** | ```python
from scipy.stats import chisquare
statistic, pvalue = chisquare([61, 17, 11, 15, 6], f_exp=[55, 25.3, 13.2, 11, 5.5])
print(statistic)
print(pvalue)
``` | |
| | ```
5.244137022397893
0.26315206062015767
``` | |

## Independence Test – Chi-Square (2 Variables)

| | |
|---|---|
| **Problem Statement** | *Determine whether two or more variables from a single population are independent by comparing the <u>distributions</u> of the variables over two or more categories.* |

| | | |
|---|---|---|
| $\chi^2$**-test** | 1. Hypotheses (two-sided)<br>   $H_0$: The two variables are independent<br>   $H_a$: The two variables are not independent<br>2. $\chi^2$-test statistic<br>3. p-value (Prob ($\chi^2$-statistic))<br>4. Significance level $\alpha = 0.05$ | A rule of thumb is that all expected individual cell counts should be at least five (5). Combine cells if less than 5.<br><br>For a chi-squared distribution<br>If p-value < $\alpha$, then reject $H_0$, insufficient evidence exists to conclude that the two variables are not independent.<br>If p-value > $\alpha$, then favor $H_0$. |
| **Expected Frequencies for a Chi-Square** | $$E = \frac{rc}{n}$$ | Contingency tables:<br>$r = \#\ of\ rows$<br>$c = \#\ of\ columns$<br>$n = sample\ size$ |
| **Chi-Square Test Statistic** | $$\chi^2 = \frac{(O-E)^2}{E}$$ $$\chi^2 = \sum \frac{(observed\ count - expected\ count)^2}{expected\ count}$$ | |
| **Degrees of Freedom** | $$df = (r-1)(c-1)$$ | $rc = number\ of\ possible\ values\ for$ <br> $the\ two\ variables\ under$ <br> $consideration$ |
| **Contingency Tables** | A contingency table is constructed from the values of the variables and categories along the rows and columns. An expected cell count is calculated by multiplying the row total by the column total and dividing by the overall total. | |

**Observed**

| Category A | Group B1 | Group B2 | Total by Category A |
|---|---|---|---|
| | **Category B** | | |
| Group A1 | 5 | 8 | 13 |
| Group A2 | 6 | 9 | 15 |
| Group A3 | 7 | 10 | 17 |
| Total by Category B | 18 | 27 | 45 |

*df = (3-1)(2-1) = 2*

**Expected**

| Category A | Group B1 | Group B2 | Total by Category A |
|---|---|---|---|
| | **Category B** | | |
| Group A1 | $\frac{13 \cdot 18}{45} = 5.2$ | $\frac{13 \cdot 27}{45} = 7.8$ | 13 |
| Group A2 | $\frac{15 \cdot 18}{45} = 6$ | $\frac{15 \cdot 27}{45} = 9$ | 15 |
| Group A3 | $\frac{17 \cdot 18}{45} = 6.8$ | $\frac{17 \cdot 27}{45} = 10.2$ | 17 |
| Total by Category B | 18 | 27 | 45 |

**Python**

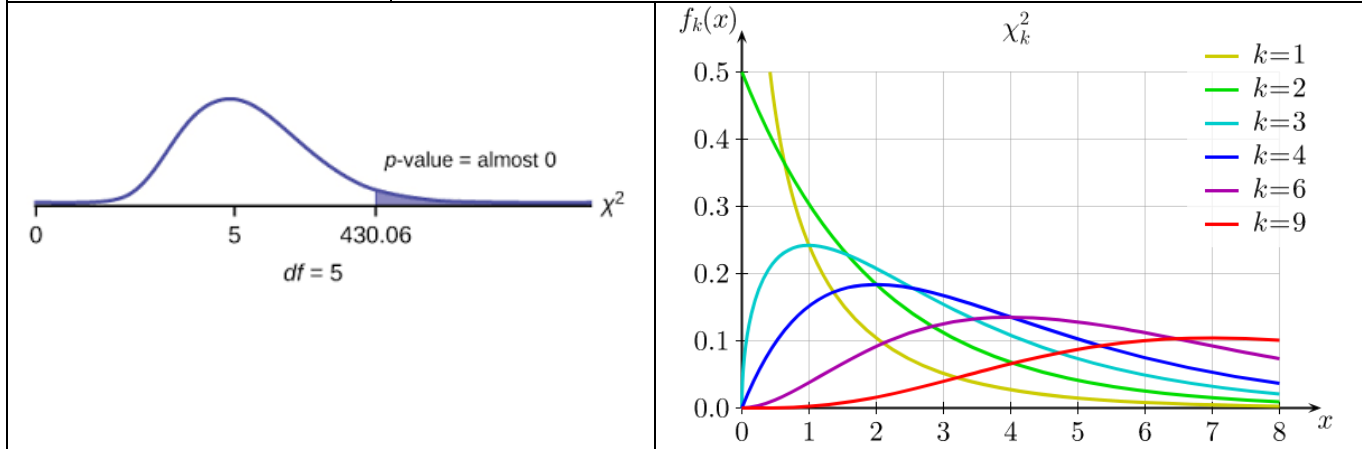```python
import numpy as np
from scipy.stats import chi2_contingency

# Construct a contingency table
parole = np.array([[405,1422], [240,470], [151,275]])

# Calculate the test statistic, p-value, df, and expected counts
chi2, p, df, ex = chi2_contingency(parole)
print(chi2, p, df, ex)
```
```
X2-test: 53.87860692066112
p-value: 1.9971429926442894e-12
df: 2
[[ 490.81741478 1336.18258522]
 [ 190.73911576  519.26088424]
 [ 114.44346946  311.55653054]]
```

## Homogeneity Test – Chi-Square (Multiple Variables)

| Problem Statement | *Determines if two or more populations (or subgroups of a population) have the same unknown distribution of a single categorical variable.* |
|---|---|

| $\chi^2$-test | 1. Hypotheses (two-sided)<br>   $H_0$ : The distribution of the frequency of one variable is the same across all sampled populations<br><br>   $p_{1,1} = p_{1,2} = ... = p_{1,J}$<br>   $p_{2,1} = p_{2,2} = ... = p_{2,J}$<br>   ...<br>   $p_{I,.1} = p_{I,2} = ... = p_{I,J}$<br><br>   $p_{1,1} = p_{2,1} = ... = p_{I,1}$<br>   $p_{1,2} = p_{2,2} = ... = p_{I,2}$<br>   ...<br>   $p_{,1,J} = p_{2,J} = ... = p_{I,J}$<br>   $H_a$ : At least one of the probability statements is false.<br>2. $\chi^2$-test statistic<br>3. p-value (Prob ($\chi^2$-statistic))<br>4. Significance level $\alpha = 0.05$ | i = 1, 2,…, I represent the categories of the first variable.<br>j = 1,2,…, J represent the categories of the second variable.<br><br>For a chi-squared distribution<br>If p-value < $\alpha$, then reject $H_0$.<br>If p-value > $\alpha$, then favor $H_0$.<br><br>If $H_0$ is rejected, conclude that the distribution of one variable is not the same across categories of the other variable. |

| Chi-Square Test Statistic | The test for homogeneity is performed in the same way as the test for independence. |
|---|---|



| Python | ```python
import numpy as np
from scipy.stats import chi2_contingency

z = np.array([[551, 580], [244, 289], [387, 503], [452, 618], [443,742]])

chi2, p = chi2_contingency(z)
print(chi2)
print(p)
```<br>X2-test: 32.2443<br>p-value: 1.70526e-6 |
|---|---|