**Python Inferential Statistics**

Lucas Hasting

University of North Alabama

DA 380 - Programming for Data Analytics

Dr. Michael Floren

2/24/2026

**Python Inferential Statistics**

## Contents

## List of Tables

## Excel Inferential Statistics Practice Using BRFSS

First, an independent-samples $t$-test was conducted to compare average general health between people who exercise and people who don't. People who did not exercise ($\bar{x} \approx 3.11$, $s \approx 1.08$), on average, had a higher score than people who did exercise ($\bar{x} \approx 2.46$, $s \approx 0.96$). This difference was statistically significant, $t(587) \approx 10.47$, $p < 0.01$. These results suggest that the general health differs between people who exercise and people who do not.

Next, a one-way analysis of variance (ANOVA) was conducted to examine differences in mean general health across 11 groups defined by annual household income. Mean and standard deviation of general health values for each group are approximately shown in Table 1. The overall effect was statistically significant, $F(10, 1170) \approx 20.00$, $p < 0.01$. These results indicate that average general health differs among annual household income groups.

Then, a chi-squared test of independence was conducted to examine the association between exercise participation and median annual household income. The test indicated that the association was statistically significant, $\chi^2(1, N = 1180) = 79.06$, $p < 0.01$. Examination of the contingency table, shown in Table 2, showed that people who did not exercise tend to be below (or equal to) the median annual household income; however, more people tend to participate in exercise in general. This pattern suggests that people who tend to not participate in exercise are more likely to have a lower annual household income than those who do, and vice versa.

Finally, a multiple linear regression analysis was conducted to examine whether exercise participation and disability status were associated with BMI. The overall regression model was statistically significant, $F(2, 1344) \approx 17.97$, $p < 0.01$, and explained 2.60% of the variance in BMI ($R^2 \approx 0.026$). Exercise participation was a significant predictor of BMI ($\beta = 2.03$, $p < 0.01$), and disability status was significant ($\beta \approx 1.14$,

$p < 0.01$) as well. These results indicate that although there is not a strong linear relationship of exercise participation and disability status with BMI, the variables do have some relationship to BMI, and thus are associated with BMI.

| Group | $\bar{x}$ | $s$ |
|:---:|:---:|:---:|
| 1 | 3.58 | 1.24 |
| 2 | 3.24 | 1.10 |
| 3 | 3.12 | 1.18 |
| 4 | 3.31 | 1.14 |
| 5 | 3.01 | 1.03 |
| 6 | 2.69 | 1.01 |
| 7 | 2.61 | 1.00 |
| 8 | 2.44 | 0.91 |
| 9 | 2.18 | 0.86 |
| 10 | 2.17 | 0.84 |
| 11 | 2.00 | 0.75 |

**Table 1**

*Mean and Standard Deviation of General Health by Annual Household Income*

| Median Annual Household Income | Participated in Exercise | Did not Participate in Exercise | Total |
|---|---|---|---|
| Less than or equal median household income | 420 | 235 | 655 |
| More than median household income | 457 | 68 | 525 |
| Total | 877 | 303 | 1180 |

**Table 2**

*Contingency Table for Median Annual Household Income and Exercise Participation*

**Appendix A**

**Python File**

```
 1  '''
 2  Name: Lucas Hasting
 3  Course: DA 380
 4  Instructor: Dr. Michael Floren
 5  Date: 2/24/2026
 6  Description: Get inferential stats from brfss1_cleaned.cvs
 7  '''
 8
 9  #import needed libraries
10  import numpy as np
11  import pandas as pd
12  from scipy import stats
13  import statsmodels.api as sm
14  import statsmodels.formula.api as smf
15
16  #load the data
17  df = pd.read_csv("brfss1_cleaned.csv")
18
19  #is general health different between people who exercise vs. not? - T-TEST
20  print("GENHLTH vs. EXERANY2:",end="\n\n")
21  gen_hlth_exer = df.loc[df["EXERANY2"] == 1, "GENHLTH"].dropna()
22  gen_hlth_not_exer = df.loc[df["EXERANY2"] == 2, "GENHLTH"].dropna()
23
24  print(stats.ttest_ind(gen_hlth_not_exer, gen_hlth_exer, equal_var=False),
        end="\n\n")
25  print("Exercise:\n",gen_hlth_exer.agg(["mean", "std"]),end="\n\n",sep="")
26  print("No Exercise:\n",gen_hlth_not_exer.agg(["mean", "std"]),end="\n\n",
        sep="")
27
28  #general health across income groups - are they different? - ANOVA
```

```python
29 print("INCOME3 vs. GENHLTH:",end="\n\n")
30 df["INCOME3_no_na"] = df["INCOME3"].dropna()
31 df["GENHLTH_no_na"] = df["GENHLTH"].dropna()
32 model_genhlth = smf.ols("GENHLTH_no_na ~ C(INCOME3_no_na)", data=df).fit()
33 print(sm.stats.anova_lm(model_genhlth))
34 print(df.groupby("INCOME3_no_na")["GENHLTH_no_na"].agg(["mean", "std"]),
      end="\n\n")
35
36 #association between exercise participation and median annual household
      income? - Ch^2
37
38 print("INCOME3 vs. GENHLTH:",end="\n\n")
39 df["EXERANY2_no_na"] = df["EXERANY2"].dropna()
40 df["income_bin_no_na"] = df["income_bin"].dropna()
41 print(df.groupby("EXERANY2_no_na")["income_bin_no_na"].value_counts())
42 print(stats.chi2_contingency(pd.crosstab(df["income_bin_no_na"], df["
      EXERANY2_no_na"])),end="\n\n")
43
44 #BMI linearily correlated with disabled and exercise participation? - OLS
      Regression
45
46 print("BMI vs. EXERANY2 and disabled:",end="\n\n")
47 df["BMI_no_na"] = df["bmi"].dropna()
48 df["disabled_no_na"] = df["disabled"].dropna()
49 model_bmi = smf.ols("BMI_no_na ~ C(EXERANY2_no_na) + C(disabled_no_na)",
      data=df).fit()
50 print(sm.stats.anova_lm(model_bmi))
51 print(model_bmi.summary())
```

## Appendix B

## Python File Output

```
 1 GENHLTH vs. EXERANY2:
 2
 3 TtestResult(statistic=np.float64(10.469372106336646), pvalue=np.float64
     (1.2428555292489605e-23), df=np.float64(586.6549314911825))
 4
 5 Exercise:
 6 mean     2.458296
 7 std      0.955383
 8 Name: GENHLTH, dtype: float64
 9
10 No Exercise:
11 mean     3.114058
12 std      1.081881
13 Name: GENHLTH, dtype: float64
14
15 INCOME3 vs. GENHLTH:
16
17                      df        sum_sq     mean_sq          F         PR(>F)
18 C(INCOME3_no_na)    10.0    187.317149   18.731715   20.000458   1.892123e-34
19 Residual          1170.0   1095.780226    0.936564        NaN            NaN
20                   mean         std
21 INCOME3_no_na
22 1.0             3.576923   1.238485
23 2.0             3.236842   1.101208
24 3.0             3.116279   1.179372
25 4.0             3.313725   1.140003
26 5.0             3.006579   1.026127
27 6.0             2.692308   1.010381
28 7.0             2.613953   1.002236
29 8.0             2.439024   0.914740
```

```
30 9.0              2.179348  0.859157

31 10.0             2.173469  0.837590

32 11.0             2.000000  0.746299

33

34 INCOME3 vs. GENHLTH:

35

36 EXERANY2_no_na   income_bin_no_na

37 1.0              1.0                    457

38                  0.0                    420

39 2.0              0.0                    235

40                  1.0                     68

41 Name: count , dtype: int64

42 Chi2ContingencyResult(statistic=np.float64(79.05919206447463), pvalue=np.
     float64(6.027662302195648e-19), dof=1, expected_freq=array
     ([[486.80932203, 168.19067797],

43     [390.19067797, 134.80932203]]))

44

45 BMI vs. EXERANY2 and disabled:

46

47                         df        sum_sq        mean_sq           F          PR
     (>F)

48 C(EXERANY2_no_na)      1.0    1289.123160   1289.123160   28.073209   1.363661e
     -07

49 C(disabled_no_na)      1.0     361.244477    361.244477    7.866814   5.107537e
     -03

50 Residual            1344.0   61716.547295     45.920050         NaN
     NaN

51                             OLS Regression Results

52 ===============================================================================

53 Dep. Variable:             BMI_no_na    R-squared:
     0.026

54 Model:                           OLS    Adj. R-squared:
```

```
     0.025
55 Method:                  Least Squares    F-statistic:
     17.97
56 Date:                 Sun, 22 Feb 2026    Prob (F-statistic):            1.99
     e-08
57 Time:                        17:43:55    Log-Likelihood:
     -4487.2
58 No. Observations:                1347    AIC:
     8980.
59 Df Residuals:                    1344    BIC:
     8996.
60 Df Model:                           2
61 Covariance Type:            nonrobust
62 ================================================================================

63                             coef     std err          t      P>|t|
     [0.025      0.975]
64 --------------------------------------------------------------------------------

65 Intercept                27.8740      0.238    117.343      0.000
     27.408      28.340
66 C(EXERANY2_no_na)[T.2.0]  2.0290      0.439      4.625      0.000
      1.168       2.890
67 C(disabled_no_na)[T.1]    1.1377      0.406      2.805      0.005
      0.342       1.933
68 ===========================================================================

69 Omnibus:                      740.497    Durbin-Watson:
     1.990
70 Prob(Omnibus):                  0.000    Jarque-Bera (JB):
     14851.974
71 Skew:                           2.113    Prob(JB):
     0.00
```

```
72 Kurtosis:                              18.709    Cond. No.
       2.76
73 ==============================================================================

74
75 Notes:
76 [1] Standard Errors assume that the covariance matrix of the errors is
       correctly specified.
```