

```
In [1]: import numpy as np
import pandas as pd
import statsmodels.api as sm
from scipy import stats
import matplotlib.pyplot as plt
import seaborn as sns
import statsmodels.formula.api as smf
import statsmodels.api as sm
from statsmodels.formula.api import ols
from statsmodels.stats.anova import anova_lm
from statsmodels.graphics.factorplots import interaction_plot
```

Problem 1 Statement:

1a) A popular Two-Wheelers company claims that its best-selling model averages 18 km per litre of petrol. But recently a government agency used a sample of 80 two-wheelers of this model and finds the sample mean to be 16.8 km/litre. We also know from previous studies that the population standard deviation is 3 km/litre. Can we expect (within 2 standard errors) that we could select such a sample if indeed the population mean is 18 km/litre?

Solution

Given The Sample A has mean $\mu_1 = 16.8$ km/Lt, sample Size = $n_1=80$.
The Population mean = $\mu = 18$ km/Lt with Standard Deviation = $\sigma = 3$.

Can we expect a sample whose population mean is 18km/Lt & within 2 standard error?

Formulation of Hypothesis

Null Hypothesis H_0 : $\mu_1 = 18$ km/Lt.

Alternative Hypothesis H_A : $\mu_1 \neq 18$ km/Lt.

```
In [2]: mu=18
sigma = 3
n1=80
mu1 =16.8

In [3]: # For A sample
zstat=(mu1-mu)/(sigma/np.sqrt(n1))
zstat

Out[3]: -3.577708763999661

In [4]: zcrit=stats.norm.ppf(.01)
zcrit

Out[4]: -2.3263478740408408

In [5]: zstat < zcrit

Out[5]: True

In [6]: std_err=sigma/np.sqrt(n1)
std_err

Out[6]: 0.33541019662496846

In [7]: # For B Sample
#condition given is 2 std err is there for the new sample is std_err1
#so first we find sample size for the same

In [8]: std_err1=2*std_err

In [9]: n2=(sigma/std_err1)**2
round(n2,2)

Out[9]: 20.0

In [10]: # now we need to find the mean (mu2) of this sample size(n2) & it should be in the range of 2 standard error

In [11]: zcrit2=stats.norm.ppf(.05/2)
zcrit2

Out[11]: -1.9599639845400545

In [12]: mu2=((zcrit2*sigma)/np.sqrt(n2))*mu
mu2

Out[12]: 16.685216189135126

In [13]: zcrit3=stats.norm.ppf(1-.05/2)
zcrit3

Out[13]: 1.959963984540054

In [14]: mu3=((zcrit3*sigma)/np.sqrt(n2))*mu
mu3

Out[14]: 19.31478381086487
```

so the range is 16.685 to 19.314. it is possible to have it with sample size 20

Problem 2 Statement:

```
In [15]: #2a) Marriott International is an American multinational diversified hospitality company that manages and franchises hotels, resorts and other accommodations. The company is a member of the S&P 500.
#Top Level Middle Level Bottom Level
# 8 8 5
# 7 7 6
# 6 6 7
# 7 9 6
# 9 10 7
# 9 9 8
# 10

#Assuming level of significance as 0.05, formulate the null and alternative hypotheses and determine which test to use
```

Solution

The scores given by Top Level, Middle Level & bottom Level are in ordinal type.

So Anova one way test can't be performed as the data is ordinal type but there is a test which has below conditions are followed also in the given question.

- 1.It's applicable on 3 independent samples
- 2.Must have minimum 5 or more observations
- 3.The data samples size can be different

The test that can be conducted in this case is Kruskal -Wallis H Test.

Formulation of Null Ho & Alternative Ha Hypothesis

Null Hypothesis H_0 : Evaluations are same at all levels.

Alternative Hypothesis H_A : One or more evaluations are different.

```
In [16]: from Kruskals import Kruskals

In [17]: Top=[8,7,6,7,9]
Top

Out[17]: [8, 7, 6, 7, 9]

In [18]: Middle=[8,7,6,9,10,9]
Middle

Out[18]: [8, 7, 6, 9, 10, 9]

In [19]: Bottom=[5,6,7,6,7,8,10]
Bottom

Out[19]: [5, 6, 7, 6, 7, 8, 10]

In [45]: H,pval=stats.kruskal(Top,Middle,Bottom)

In [46]: print ("Kruskal -Wallis H Test p-value=", pval)

alpha_level = 0.05

if pval < alpha_level:
    print('We have enough evidence to reject the null hypothesis in favour of alternative hypothesis')
    print('We conclude that One or more evaluations are different.')
else:
    print('We do not have enough evidence to reject the null hypothesis in favour of alternative hypothesis')
    print('We conclude that Evaluations are same at all levels.')

Kruskal -Wallis H Test p-value= 0.35270982873167683
We do not have enough evidence to reject the null hypothesis in favour of alternative hypothesis
We conclude that Evaluations are same at all levels.
```

Decision Rule-

Since pvalue is > 0.05. so at 95% confidence level we fail to reject Null Hypothesis and conclude that evaluation are same at all levels.

```
In [ ]:
In [ ]:
```

Problem 3 Statement:

3b) The table shows the quantity of soaps sold for different brands at different locations, collected over 20 days. Conduct a two-way ANOVA with interaction at $\alpha = 5\%$ to test the effects of brands, locations and interaction on sales.

The file (soaps.csv) is there.

```
In [22]: data = pd.read_csv ('SOAPS- Q3.csv')

In [23]: data.head()

Out[23]:
   Loc  Brand  Qty
0    1     X   20
1    2     X   20
2    1     X   16
3    2     X   21
4    1     X   24

In [24]: data.shape

Out[24]: (120, 3)

In [25]: data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 120 entries, 0 to 119
Data columns (total 3 columns):
#   Column  Non-Null Count  Dtype
---  --
0    Loc    120 non-null         int64
1    Brand   120 non-null         object
2    Qty     120 non-null         int64
dtypes: int64(2), object(1)
memory usage: 2.9+ KB

In [26]: data.describe()

Out[26]:
          Loc      Qty
count  120.000000  120.000000
mean     1.500000   24.908333
std       0.502096   5.130230
min       1.000000   10.000000
25%       1.000000   21.750000
50%       1.500000   25.000000
75%       2.000000   28.000000
max       2.000000   39.000000

In [27]: data.Loc = pd.Categorical(data.Loc)

In [28]: data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 120 entries, 0 to 119
Data columns (total 3 columns):
#   Column  Non-Null Count  Dtype
---  --
0    Loc    120 non-null         category
1    Brand   120 non-null         object
2    Qty     120 non-null         int64
dtypes: category(1), int64(1), object(1)
memory usage: 2.2+ KB

In [29]: data['Loc'].value_counts()

Out[29]:
2    60
1    60
Name: Loc, dtype: int64

In [30]: data['Brand'].value_counts()

Out[30]:
Y    40
X    40
Z    40
Name: Brand, dtype: int64
```

Formulation of Hypothesis

Null Hypothesis H_0 : The Qty sold variable with respect to Location equal.

Alternative Hypothesis H_A : At least one of the Qty sold variable with respect to Location unequal.

Null Hypothesis H_0 : The Qty sold variable with respect to each Brand equal.

Alternative Hypothesis H_A : At least one of the Qty sold variable with respect to Brand is unequal.

Null Hypothesis H_0 : The Qty sold variable with respect to Location and Brand and their interaction is equal.

Alternative Hypothesis H_A : At least one of the Qty sold variable with respect to Location and Brand and their interaction is unequal.

Assumptions for ANOVA

- 1. All populations under consideration have normal distribution
- 2. All populations under consideration have equal variances.
- 3. The sample is a random sample, i.e. the observations are collected independently of each other.

```
In [31]: sns.distplot(data['Qty'])
plt.show()

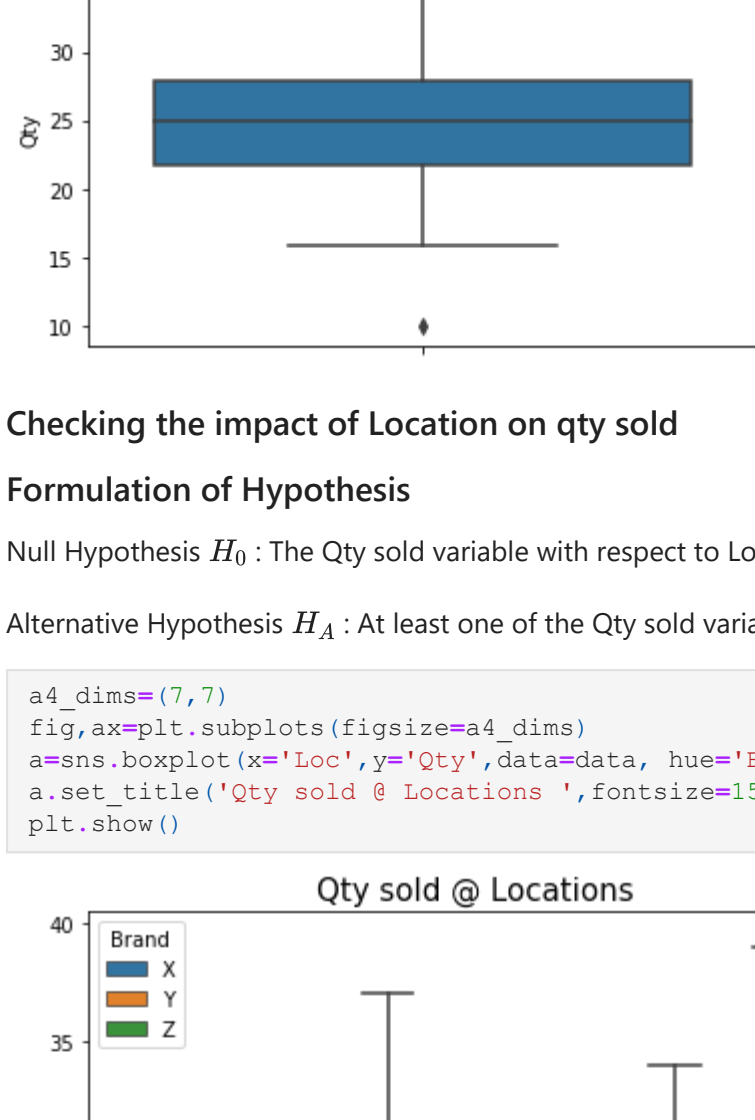
In [32]: #Assumption 1: Normality
w,p_value=stats.shapiro(data['Qty'])
print('w={}'.format(w), 'p_value={}'.format(p_value))

w=0.9890562891960144 p_value=0.4538714587688446

Since p-value of the test is very large, so the response follows the normal distribution.

In [33]: sns.boxplot(x='Qty',data=data,orient='v')

Out[33]: <matplotlib.axes._subplots.AxesSubplot at 0x10ee28927f0>
```



Checking the impact of Location on qty sold

Formulation of Hypothesis

Null Hypothesis H_0 : The Qty sold variable with respect to Location equal.

Alternative Hypothesis H_A : At least one of the Qty sold variable with respect to Location unequal.

```
In [34]: a4_dims=(7,7)
fig,ax=plt.subplots(figsize=a4_dims)
a=sns.boxplot(x='Loc',y='Qty',data=data, hue='Brand')
a.set_title('Qty sold @ Locations',fontsize=15)
plt.show()

In [35]: formula='Qty~ Loc'
models=formula,data).fit()
aov_table=anova_lm(mod)
print(aov_table)
```

	df	sum_sq	mean_sq	F	PR(>F)
Loc	1.0	7.008333	7.008333	38.170340	1.829940e-13
Residual	118.0	3124.983333	26.482910	NaN	NaN

Since F value is small and p is more than alpha we fail to reject null hypothesis. so qty sold is same at all locations

Checking the impact of brand on the qty sold

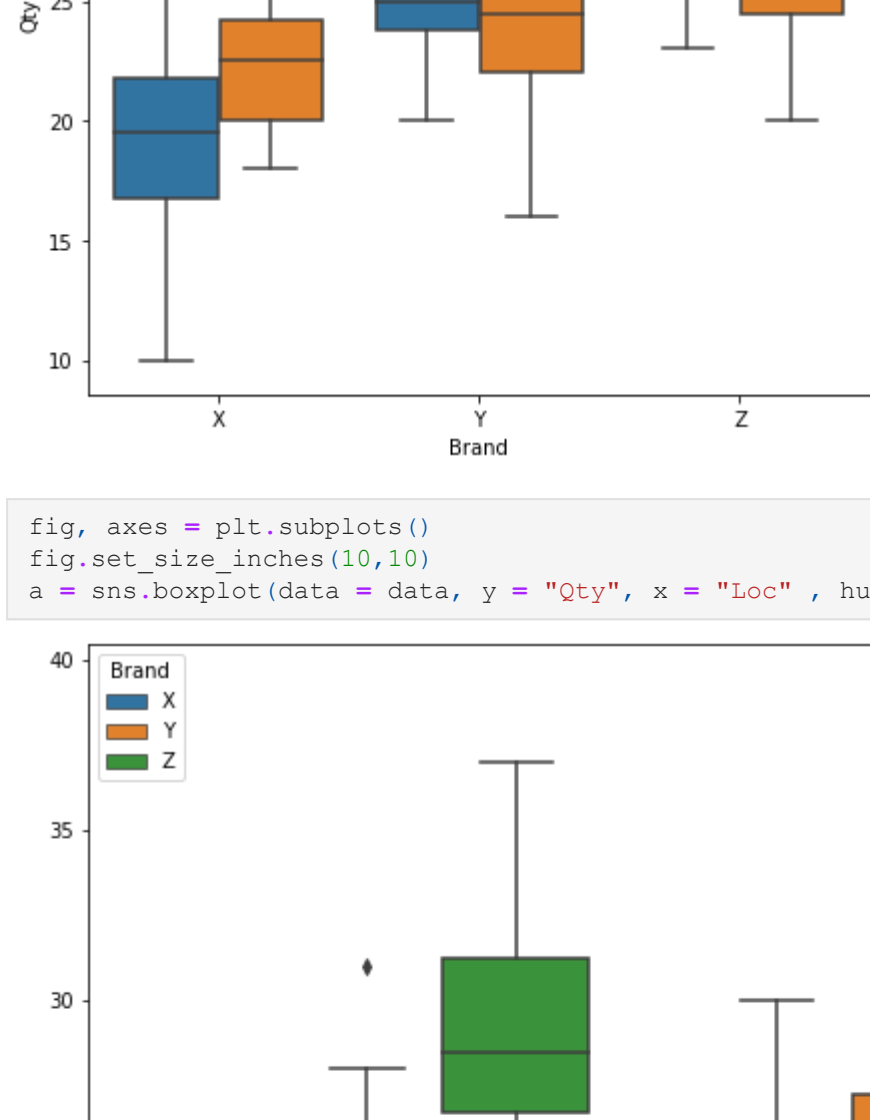
Formulation of Hypothesis

Null Hypothesis H_0 : The Qty sold variable with respect to each Brand equal.

Alternative Hypothesis H_A : At least one of the Qty sold variable with respect to Brand is unequal.

```
In [36]: a4_dims=(7,7)
fig,ax=plt.subplots(figsize=a4_dims)
a=sns.boxplot(x='Brand',y='Qty',data=data, hue='Loc')
a.set_title('Qty sold by the Brands',fontsize=15)
plt.show()

In [37]: formula='Qty~ Brand'
mod1=ols(formula,data).fit()
aov_table1=anova_lm(mod1)
print(aov_table1)
```



```
In [38]: formula='Qty~ Brand'
mod1=ols(formula,data).fit()
aov_table1=anova_lm(mod1)
print(aov_table1)
```

	df	sum_sq	mean_sq	F	PR(>F)
Brand	2.0	1240.316667	620.158333	39.279968	1.055160e-13
Residual	117.0	1891.675000	16.168162	NaN	NaN

Since F value is large and p is small and less than alpha we reject null hypothesis. so qty sold is impacted by the brands.

Checking the impact on qty sold by the brand and location and their interaction

Formulation of Hypothesis

Null Hypothesis H_0 : The Qty sold variable with respect to Location and Brand and their interaction is equal.

Alternative Hypothesis H_A : At least one of the Qty sold variable with respect to Location and Brand and their interaction is unequal.

```
In [39]: formula2='Qty~ C(Brand)+C(Loc)'
model2=ols(formula2, data).fit()
aov_table2=anova_lm(model2)
print(aov_table2)

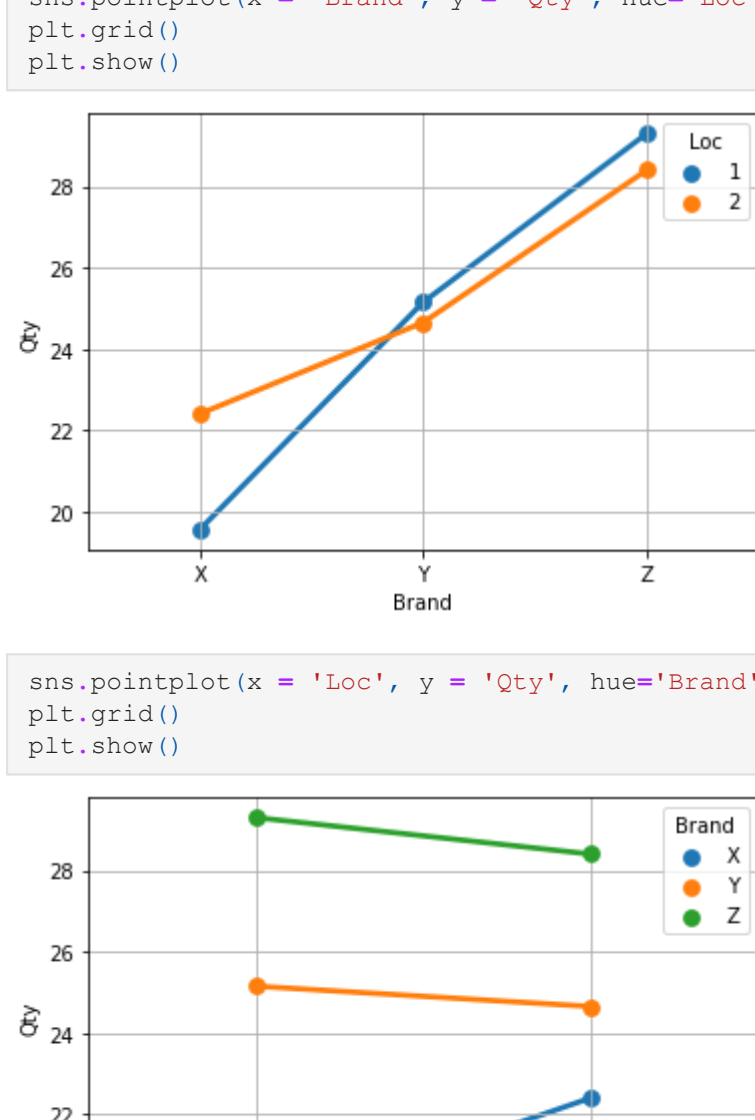
In [40]: formula3='Qty~ C(Brand)+C(Loc)+ C(Brand):C(Loc)'
model3=ols(formula3, data).fit()
aov_table3=anova_lm(model3)
print(aov_table3)
```

	df	sum_sq	mean_sq	F	PR(>F)
C(Brand)	2.0	1240.316667	620.158333	39.279968	1.055160e-13
C(Loc)	1.0	7.008333	7.008333	0.443898	5.065930e-01
C(Brand):C(Loc)	2.0	84.816667	42.408333	2.686085	7.246036e-02
Residual	114.0	1799.850000	15.788158	NaN	NaN

Brand impact is visible in the below plot too.

```
In [43]: sns.pointplot(x='Brand', y='Qty', hue='Brand',data=data,ci=None)
plt.grid()
plt.show()

In [44]: sns.pointplot(x='Loc', y='Qty', hue='Brand',data=data,ci=None)
plt.grid()
plt.show()
```



Since Fstat value (2.68) is more than 1 but is less than F crit value(3.07)and p is 0.07 and more than alpha (0.05) we fail to reject null hypothesis for the interaction. so qty sold with respect to location and their interaction is equal. Individually the location has no impact & the brand has great impact on the qty sold..

