

ONE WAY ANOVA python.

```
import numpy as np
import pandas as pd
import statsmodels.api as sm
from statsmodels.formula.api import ols. # for the model part.
```

data part :

```
data = pd.read_csv(' ')
data
```

model part :

```
model = ols('dependent ~ indep', variable var, data=data).fit()
```

```
anova_result = sm.stats.anova_lm(model, type=2)
```

print(anova_result)

gives ANOVA table. p-value will be there
 $SSB \rightarrow \text{indep-var}$ rej/accept H_0 .
 $SSW \rightarrow \text{residual}$

Tukey Test :

if we accept H_0 , it means atleast one pair has significant difference. To determine where that difference is, we use Tukey's Test

Tukey's Test

```
from statsmodels.stats.multicomp import pairwise_tukeyhsd
```

```
tukey = pairwise_tukeyhsd(endog=dependent variable,  

                           groups=indep var ↳ categorical var)
```

tukey.summary() as we

OR tukey.results_table used a dataframe,

it gives pairwise reject H_0 | endog = data['*dep var name*']
 - if true, (it means reject H_0) | groups = data['']

i.e. means are significantly different)

- if false (accept H_0 ; i.e. means are nearer)

TWO WAY ANOVA (PYTHON)

```

import pandas as pd
import numpy as np
import statsmodels.api as sm
from statsmodels.formula.api import ols
import statsmodels.stats.multicomp as mc
from statsmodels.stats.multicomp import MultiComparison
from statsmodels.stats.multicomp import pairwise_tukeyhsd
from scipy import stats

```

without replication.

repeat each input
4 times

create dat aframe

```

df = pd.DataFrame ({'Lab': np.repeat(['Lab1','Lab2','Lab3','Lab4'], 4),
                    'Sample': np.tile (np.repeat(['Type1','Type2','Type3','Type4'], 1), 4),
                    'Recovery': [ 16 observations ] })

```

df

2 way anova w/o replication

→ if we do with replicates
+ c(Lab): c(Sample)

```

model = ols('Recovery ~ c(Lab) + c(Sample)', data=df).fit()

```

```

anova_result = sm.stats.anova_lm(model, type=2)

```

multi comparison

```

mc = MultiComparison(df['Recovery'], df['Lab'])

```

```

mc_result = mc.tukeyhsd(0.05)

```

```

mc_result.summary()

```

do comparison like this for each Indep. var.

2way ANOVA with replication.

```

model = ols('depvar ~ c(indvar1) + c(indvar2) + c(indvar1):c(indvar2)')

```

data=df).fit

interaction effect

MANOVA python

```

import pandas as pd
import numpy as np
from statsmodels.formula.api import ols
from statsmodels.multivariate.manova import MANOVA
from statsmodels.stats.multicomp import pairwise_tukeyhsd
from scipy import stats

# data
df =

```

if 2 way manova:
 ('depvar1 + depvar2
 ~ indepvar1 + indepvar2)
 for univar part:
 ('depvar1 ~ indepvar1 + ind 2')

test manova-

```

maov = MANOVA.from_formula ('depvar1 + depvar2 + .. ~ indep.var',
print (maov.mv_test())
data = df )

```

gives 2 tables: intercept, independent variable.

see Wilks lambda - if $p < \alpha$, rej H₀.

as H_a is accepted i.e there is significant differences in manova, we saw combined effect,

we have to now see for each dependent variable so we do anova.

* intercept is the estimate of the dependent variable when all the indep vars are 0

anova.

```
fit1 = ols('depvar1 ~ indep var', data = df).fit()
```

```
anova1 = sm.stats.anova_lm(fit1, type=2)
```

anova1

do similarly for the dependent variables

multiple comparison.

```
mcl = pairwise_tukeyhsd(df['depvar1'], df['indepvar1'])
```

mcl - result 3-table.

do for all.

T-Tests in python.

A pair t-test (`stats.ttest_rel`) is the same as an independent samples t-test on the difference in 2 groups (`stats.ttest_1samp`)

Eg:-
`before = np.array([1, 2, ...])`
`after = np.array([1, 2, 3])`

→ `stats.ttest_rel(before, after, alternative='>')`
 OR

`difference = before - after`

→ `stats.ttest_1samp(difference, popmean=0, alternative='>')`
 it gives p-value & t-statistic.
 (2 sided p-value if alternative isn't mentioned)

$$SST = SSW + SSB$$

$$\sum_{i=1}^c \sum_{j=1}^m (x_{ij} - \bar{x})^2 = \sum_{j=1}^c n_j (\bar{x}_j - \bar{x})^2 + \sum_{i=1}^c \sum_{j=1}^m (x_{ij} - \bar{x}_j)^2$$

M	T	W	T	F	S	S
Page No.:	1					
Date:						

YOUVA

ANOVA ONE WAY

(Solving)

given: Wheat Production of 3 wheat varieties in tonnes.

fields	A	B	C
F1	6	5	5
F2	7	5	4
F3	3	3	3
F4	8	7	4

factors affecting the result = 1
which is variety of wheat

Step 1: H_0 and H_a .

$\rightarrow H_0: \mu_A = \mu_B = \mu_C$ (production of type A, B, C is same, it implies that the variety of wheat doesn't affect the production of wheat significantly)

$\rightarrow H_a$: atleast one of μ_A, μ_B, μ_C are significantly different (it implies that variety of wheat affects the production of wheat significantly).

Step 2: SS Between (SSB). $SSB = n_1(\bar{x}_1 - \bar{x})^2 + n_2(\bar{x}_2 - \bar{x})^2 + \dots$

$$A: n_1 = 4 \quad \bar{x}_1 = (6+7+3+8)/4 = 6 \quad \bar{x}_1 = 6$$

$$B: n_2 = 4 \quad \bar{x}_2 = (5+5+3+7)/4 = 5 \quad \bar{x}_2 = 5$$

$$C: n_3 = 4 \quad \bar{x}_3 = (5+4+3+4)/4 = 4 \quad \bar{x}_3 = 4$$

$$\text{Grand mean: } \bar{x} = (6+5+4)/3 = 5 \quad \bar{x} = 5$$

$$SSB = n_1(\bar{x}_1 - \bar{x})^2 + n_2(\bar{x}_2 - \bar{x})^2 + n_3(\bar{x}_3 - \bar{x})^2$$

$$SSB = 4(6-5)^2 + 4(5-5)^2 + 4(4-5)^2 = 8 \quad \bullet SSB = 8$$

Step 3: SS Within (SSW).

$$SSW = \sum (x_{1i} - \bar{x}_1)^2 + \sum (x_{2i} - \bar{x}_2)^2 + \sum (x_{3i} - \bar{x}_3)^2$$

- find diff b/w each data point & its column mean
- square that deviation & add all up.



$$A: (6-6)^2 + (7-6)^2 + (3-6)^2 + (8-6)^2 = 14$$

$$B: (5-5)^2 + (5-5)^2 + (3-5)^2 + (7-5)^2 = 8$$

$$C: (5-4)^2 + (4-4)^2 + (3-4)^2 + (4-4)^2 = 2$$

$$SSW = 14 + 8 + 2 = 24$$

Step 4: SS Total

$$SST = \sum (x_{ij} - \bar{x})^2 \rightarrow \text{diff b/w each}$$

$$\text{OR } SST = SSB + SSW$$

data point & grand mean b/w square & add.

$$SST = SSW + SSB = 24 + 8 = 32$$

Step 5: ANOVA TABLE $m = 3 \quad n = 12$

source of variation	SS	deg. of freedom	MS	Ratio = MSB/MSW	Fcrit
between sample	8	$3-1 = 2$	$8/2 = 4$	$4/2.67$	$= F(2, 9)$
within sample	24	$12-3 = 9$	$24/9 = 2.67$	$= 1.5$	$= 4.26$
Total	32	$12-1 = 11$	$\frac{SSW}{df \text{ within}}$		

- SSB degree of freedom = $m-1$: $m = \text{no. of samples}$
- SSW degree of freedom = $n-m$: $n = \text{no. of elements in all samples}$
- SST degree of freedom = $n-1$
- Fcritical: $F(\alpha, df_{\text{between}}, df_{\text{within}})$

Step 6: Conclusion

Ratio $>$ Fcrit : Reject $H_0 \quad p < \alpha$, reject H_0 .

Here, we accept H_0 .

i.e. production of wheat variety A, B, C is same

i.e. the variety does not affect production significantly

TWO WAY ANOVA: used to estimate how the mean of a quantitative variable changes due to the levels of 2 categorical variables. How 2 independent variables in combination affect a dependent variable.

2 way anova with replication Eg:-

	genotype AA	Aa	aa	yes interaction
Temp 10°C	_____	_____	_____	..
Temp 20°C	_____	_____	_____	..
Temp 30°C	_____	_____	_____	..

without replication :

	AA	Aa	aa	no interaction
Temp 10°	-	-	-	
Temp 20°	-	-	-	
Temp 30°	-	-	-	

— X —

gender score age grp.

factors \rightarrow age
 \rightarrow gender.

B		H ₀ : gender will have no significant effect on score
B		
G	for factor 1:	
G	H ₀ : no sig. diff.	
B	b/w means of	\rightarrow age will have no significant effect on score
G	row factor	
:	H _a : sig. diff b/w	
:	means of row fact.	
	for factor 2:	
	column factor	\rightarrow gender & age interaction will have no effect on score.

* factor is an independent variable whose values are controlled and varied by the experimenter.

* Tukey test allows us to interpret the statistical significance of our ANOVA test & find out which specific groups mean compared w. each other are different.

it is used to check if there is and find out where the statistical significance is occurring.
used for pairwise comparison.

* MANOVA: used to examine a dataset with multiple dependent variables at a time.
we can find the effect of one or more independent variables on two or more dependent variables.

one way ~~ANOVA~~: 1 indep var.

2 way : 2 indep var