

March 2011						
Wk	Mo	Tu	We	Th	Fr	Sa
9		1	2	3	4	5
10	7	8	9	10	11	12
11	14	15	16	17	18	19
12	21	22	23	24	25	26
13	28	29	30	31		

Tuesday

1

March

2011

Week 9 • 60-305

```
from scipy.stats import mannwhitneyu
stats, p = mannwhitneyu(df1.Design1, df1.Design2)
print(stats, p)
```

↙

0.9641

p > 0.05 → accept H₀

(Both design1 & design2 have same sales).

## # Kruskal-Wallis Test

```
df2 = pd.read_excel('kruskal_wallis.xlsx', sheet_name=0)
```

df2

```
from scipy.stats import kruskal
stats, p = kruskal(df2.Design1, df2.Design2, df2.Design3)
```

print(stats, p)

↙

2.734

0.284

p > 0.05

∴ H₀ → accepted

(design1/2/3 are same, there is no diff. in sales.)

## # Chi-square

```
df3 = pd.read_excel('chisquaretest.xlsx', sheet_name=0)
```

df4 = df3.dropna()

↳ removes all Null value & store in df4.

2

# Wednesday

March 2011

March						2011	
Wk	Mo	Tu	We	Th	Fr	Sa	Su
9		1	2	3	4	5	6
10	7	8	9	10	11	12	13
11	14	15	16	17	18	19	20
12	21	22	23	24	25	26	27
13	28	29	30	31			

Week 9 ■ 61-304

(When Both var. are categorical).

```
from scipy.stats import chi2_contingency
```

```
chitable = pd.crosstab(af4.Gender, af4.Smoking).
```

charitable

stats, p, dof, expected = chi2\_contingency(chitable)

```
print(stats, p)
```

 ↓

3.17.12

0.2048

$p > 0.05$  (we accept  $H_0$ )

No dependency on gender & smoking.

## # Practical Implementation of Parametric Test

## \* One-Sample Test

```
df5 = pd.read_excel('One Sample.xlsx', sheet_name=0)  
df5.head()
```

from scipy.stats import ttest\_1samp → population mean.

stats, p = ttest\_1samp(af5\_Height, 65)

```
print(stats, p)
```

11.498

1.087.e-26

-26  
exponential

$P < 0.05$

reject  $H_0$   
accept  $H_1$

	April 2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
13					1	2	3
14	4	5	6	7	8	9	10
15	11	12	13	14	15	16	17
16	18	19	20	21	22	23	24
17	25	26	27	28	29	30	

Thursday

March

2011

3

TWO

Week 9 • 62-303

## # Sample paired t-test

```
df6 = pd.read_excel('pairedsample.xlsx', sheet_name=0)
df6.head()
```

```
from scipy.stats import ttest_rel
stats, p = ttest_rel(df6.English, df6.Math)
print(stats, p)
```

36.3156      3.0717 e-128      P < 0.05 (accept H1)

exponential

## # TWO sample independent t-test

```
df7 = pd.read_excel('IndependentSample.xlsx', sheet_name=3)
df7.head()
```

```
from scipy.stats import ttest_ind
stats, p = ttest_ind(df7.Nmathlete, df7.athlete)
print(stats, p)
```

13.368

4.116 e-33

P < 0.55 (accept H1)

(there is difference

two time of athlete &

## # One-sample f-test (One Way ANOVA)

- ANOVA is used to compare 2 or more sample mean.

4

Friday

March

2011

	March 2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
9		1	2	3	4	5	6
10	7	8	9	10	11	12	13
11	14	15	16	17	18	19	20
12	21	22	23	24	25	26	27
13	28	29	30	31			

Week 9 • 63-302

OWT

Statistical Technique	Dependent variable	independn variable	Purpose of Technique.
1. ANOVA (analysis of variance).	continuous	categorical (also called factors)	It is used to compare the mean of two or more sample mean
2. 10 ANCOVA (analysis of covariance)	continuous	continuous & categorical	It is used to compare of covariance of sample mean.
3			
4	class A v/s class B	(Between comparision).	
5	within your class it's self (within class A).	(within comparision).	
6			
#	ANOVA is used in Between as well as within comparision.		
	{ depn var → continuous		
	{ indepn categorical var → factors. also called.		
◦	One way ANOVA (1 depn var & 1 independn var.)		
◦	Two way ANOVA (1 depn var & 2 independn var.)		
◦	Multi way ANOVA (1 depn var & more than 2 indepn. var.)		

April 2011						
Wk	Mo	Tu	We	Th	Fr	Sa Su
13					1	2 3
14	4	5	6	7	8	9 10
15	11	12	13	14	15	16 17
16	18	19	20	21	22	23 24
17	25	26	27	28	29	30

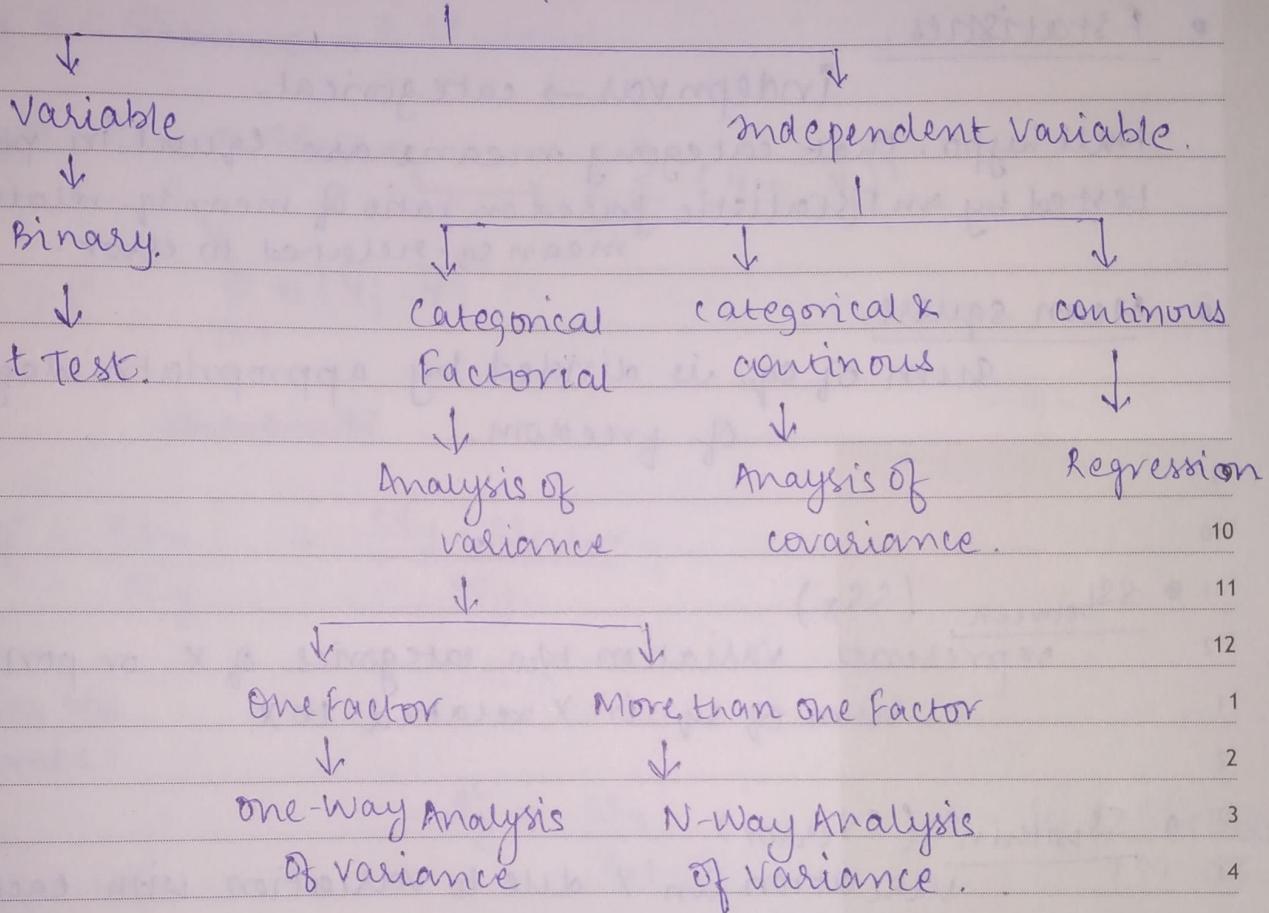
Saturday

5

March 2011

Week 9 • 64-301

Metric Dependent Variable.



• When Both variable are continuous → Regression

• Indepn var → categorical → ANOVA

• Indepn var → categorical + continuous → ANCOVA

Sunday

65-300

6

# One way Analysis of Variance.

•  $\eta^2$  → deciding factor.

→ how much indepn. var is going to affect dependent. var.

$$0 < \eta^2 < 1$$

$\eta^2 = 0.725 \rightarrow 72.5\% \text{ indepn. var. is affecting depn. var.}$

7

Monday

March 2011

March						
Wk	Mo	Tu	We	Th	Fr	Sa
9		1	2	3	4	5
10	7	8	9	10	11	12
11	14	15	16	17	18	19
12	21	22	23	24	25	26
13	28	29	30	31		27

Week 10 • 66-299

## • F statistics.

indepn var → categorical.

Null hypo. that category means are equal in popul'n is tested by an F statistic based on ratio of mean sq. related to x & mean sq. related to error.

## • Mean square.

Sum of sq. is divided by appropriate degree of freedom.

9

10

## • SS<sub>between</sub> (SS<sub>x</sub>)

represents variation b/w categories of x or portion of sum of sq. in y related to x

## • SS<sub>within</sub> (SS<sub>error</sub>)

variation in y due to variation within each of categories of x.

## • SS<sub>y</sub> → total variation in y

continuous  
categorical

# ① Identify the Dependent & Independent vars.

steps involved in ANOVA.

② Decompose Total Variation

③ Measure effects

④ Test the significance  
interpret result.

	Mo	Tu	We	Th	Fr	Sa	Su
						1	2
13						3	
14	4	5	6	7	8	9	10
15	11	12	13	14	15	16	17
16	18	19	20	21	22	23	24
17	25	26	27	28	29	30	

Tuesday

March 2011

8

Week 10 • 67-298

$$SS_y = SS_{\text{between}} + SS_{\text{within}}$$

$$SS_y = SS_x + S_{\text{error}} \rightarrow \sum_{j=1}^c \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_j)^2$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 \quad \sum_{j=1}^c n_j (\bar{y}_j - \bar{y})^2$$

where  $\bar{y}_j$

$$n^2 = \frac{SS_x}{SS_y} = \frac{SS_y - S_{\text{error}}}{SS_y}$$

Varies b/w

0 and 1.

$$S_y^2 = \frac{SS_x}{c-1}$$

$$or \quad S_y^2 = \frac{S_{\text{error}}^2}{(N-c)}$$

$$H_0: \mu_0 = \mu_1 = \mu_2 = \dots = \mu_c$$

= Mean sq. due to x

$$H_1: \neq$$

$$= MS_x$$

= Mean sq.  
due to  
error

$$= MS_{\text{error}}$$

$$F = \frac{SS_x/c-1}{S_{\text{error}}/(N-c)} \xrightarrow{\text{NO. of categories}}$$

$$F = \frac{MS_x}{MS_{\text{error}}}$$

compare calculated F stats &  
F value by table.

9

Wednesday

March 2011

March 2011						
Wk	Mo	Tu	We	Th	Fr	Sa
9		1	2	3	4	5
10	7	8	9	10	11	12
11	14	15	16	17	18	19
12	21	22	23	24	25	26
13	28	29	30	31		

Week 10 • 68-297

Day 13.

Output variable (target var.) → Dependent var.  
 Input var. → Independent var.

- One Way → One Dependent + 1 Indepen.

Ex Sales (Depen<sup>n</sup>) v/s Promotion of item (Indepn<sup>n</sup>)

On Python,

9      F ← ~~F < 0.05~~      \* Two way. → 1 depend<sup>n</sup> + 2 Independ<sup>n</sup>.  
 10     F      P > 0.05.

Ex

sales v/s Promotion, coupon  
dependn                  Indepn.

### # ANCOVA.

Analysis of covariance

Used when depn → continuous

indepn var → continuous + categorical

Ex DV → Sales

IDV → Promotion, coupon, Client Rating.

### • Python Implementation.

```
df = pd.read_excel('ANCOVA.xlsx', sheet_name=0)
df.head()
```

### ② one way ANOVA.

	2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
13					1	2	3
14	4	5	6	7	8	9	10
15	11	12	13	14	15	16	17
16	18	19	20	21	22	23	24
17	25	26	27	28	29	30	

Thursday

March

2011

10

Week 10 • 69-296

import statsmodels.api as sm  
from statsmodels.formula.api import ols → ordinary least square

model = ols('Sales ~ c(Promotion)', df).fit()

↓      ↓      ↓  
DV      categorical      IDV

DV & IDV  
are differentiated  
by ~

- typ=3 (discriminant analysis)
- typ=1 (Regression)

standard rule  
for anova &  
ancova

one-wayanova = sm.stats.anova\_lm(model, typ = 2)      10

print(onewayanova).

(SS <sub>xi</sub> )	# D/P:	sum_sq	df	F	Prob.value.	
B/W variation ↪ c(Prom <sup>n</sup> )	106.0667	2.0	17.94	0.000011		2
Residual	79.800	27.0	NAN	NAN		3
within variation ↪ (SSerror)		↓	calculated			4
		N-C	F value.			5
						6

BCz of promotion, there is significant impact on sales.

## ② TwoWay ANOVA.

model = ols('Sales ~ c(Promotion) + c(coupon)', df).fit()

twoWay = sm.stats.anova\_lm(model, typ = 2).

twoWay.

#		sum_sq	df	F	PR(>F)	Behavior
c(Promotion)	106.0667	2.0	52.098	8.032 e-10	2 ↗ affec	
c(coupon)	53.333	1.0	52.392	1.0950 e-07	sale	
Residual	26.4667	26.0	Nan	Nan		

11

Friday

March 2011

March 2011						
Wk	Mo	Tu	We	Th	Fr	Sa
9		1	2	3	4	5
10	7	8	9	10	11	12
11	14	15	16	17	18	19
12	21	22	23	24	25	26
13	28	29	30	31		

Week 10 • 70-295

① Multinay.ANCova.

model = ols('Sales ~ C(Promotion) + C(Coupon) + ClientRating', df).fit()

ancova = sm.stats.anova\_lm(model, typ=2)

	PR(>F)	
Promotion	1.301 e-09	{ has impact on sales.
Coupon	1.4715 e-07	
ClientRating	3.745 e-01	→ NO impact on sales P>0.05

9

10 # Finding eta( $\eta^2$ ).

11 def anova\_table(onedayanova):

1 onedayanova['eta\_sq'] = oneday[:-1]['sum\_sq']/  
2 sum(onedayanova['sum\_sq'])

3 cols = ['sum\_sq', 'df', 'F', 'PR(&gt;F)', 'eta\_sq']

4 oneday = onedayanova[cols]

5 return onedayanova;

6 anova\_table(onedayanova)

anova\_table(twoway).

anova\_table(ancova)

①