

# Unit3\_HW

April 21, 2020

## 1 HW-U3

```
[1]: import numpy as np
import pandas as pd
import scipy.stats as stats
import statistics as sta
import seaborn as sns
import matplotlib.pyplot as plt
import math
import pingouin as pg
import statsmodels.stats.anova as anova
import statsmodels.api as sm
from statsmodels.formula.api import ols
import statsmodels.stats.multicomp
from IPython.core.interactiveshell import InteractiveShell
InteractiveShell.ast_node_interactivity = 'all'
sns.set_style("darkgrid")
```

### 1.1 HW-U3-1 : CI and NHST:

对于随机样本 $x_1, x_2, x_3$  (用如下python代码产生)

```
n1=25
np.random.seed(100)
x1=stats.norm.rvs(3,3,n1)+stats.uniform.rvs(-1,1,n1)
x2=stats.f.rvs(2,30,0,1,n1)**2+stats.uniform.rvs(-1,1,n1)
x3=stats.uniform.rvs(-1,1,n1)**2+x1
```

```
[2]: n1=25
np.random.seed(100)
x1=stats.norm.rvs(3,3,n1)+stats.uniform.rvs(-1,1,n1)
x2=stats.f.rvs(2,30,0,1,n1)**2+stats.uniform.rvs(-1,1,n1)
x3=stats.uniform.rvs(-1,1,n1)**2+x1
```

(1) 请检验 $x_1, x_2, x_3$ 的正态性; 然后根据正态性, 完成下面两个计算:

```
[3]: print(stats.kurtosis(x1),stats.skew(x1))
      print(stats.kurtosis(x2),stats.skew(x2))
      print(stats.kurtosis(x3),stats.skew(x3))
```

```
-0.7593896261044124 -0.08719026216018647
3.633915083065232 2.206928065140578
-0.8634789564853813 -0.19772153292015818
```

可以看出, x1和x3的正态性较好, x2的正态性较差

## (2) 计算x1,x2,x3对应总体均值的99% CI

```
[4]: def t_ci(data,alpha):
      mean,std,length=np.mean(data),np.std(data,ddof=1),len(data)
      ci_len=stats.t.isf(alpha,length-1)*std/np.sqrt(length)
      return (mean-ci_len,mean+ci_len)

      def bootstrap_ci(data,alpha,n_boots=200):
          means = []
          for i in range(n_boots):
              random_sample=np.random.choice(data,len(data),replace=True)
              means.append(np.array(random_sample).mean())
          # Compute the percentiles of choice for the bootstrapped means
          ci_l,ci_h = np.percentile(means, [alpha*100,(1-alpha)*100])
          return ci_l,ci_h
```

```
[5]: t_ci(x1,0.005)
      bootstrap_ci(x2,0.005)
      t_ci(x3,0.005)
```

```
[5]: (1.293685723789949, 4.1922367478347855)
```

```
[5]: (0.40974113071592927, 6.773569286898744)
```

```
[5]: (1.5525518558575795, 4.4572687891730665)
```

## (3) 计算x3, x1总体均值差值的95% CI

```
[6]: data=x3-x1
      t_ci(data,0.025)
```

```
[6]: (0.14955544577583174, 0.3743427276300801)
```

## (4) 利用置信区间和NHST两种方法推断x1来自的总体均值是否大于2.0

置信区间法: 由 (2) 知, x1的总体均值最有可能在 (1.29,4.19) 之间, 因此总体均值有可能大于2.0.

NHST: 设 $H_0: \bar{x}_1 \leq 2.0$ ,  $H_1: \bar{x}_1 > 2.0$ .

```
[7]: length=len(x1)
per=stats.t.isf(0.005,length-1)
t,p=stats.ttest_1samp(x1,2)
p
t,per
```

[7]: 0.16452133248027612

[7]: (1.433832017252969, 2.796939504772805)

由于 $t < per$ , 因此接受 $H_0$ , 即总体均值不大于2.0

#### (5) 利用置信区间和NHST两种方法推断 $x_1, x_3$ 来自的总体均值是否相等; 并计算effect size (Cohen' s d)

```
[8]: t_ci(x3-x1,0.005)
t,p=stats.ttest_1samp(x3-x1,2)
t,p
cohen_d=np.mean(x3-x1)/np.std(x3-x1,ddof=1)
cohen_d
```

[8]: (0.10963626462978451, 0.41426190877612734)

[8]: (-31.916047477418633, 3.605187302482976e-21)

[8]: 0.9620408037436069

$\mu_3 - \mu_1$ 的置信区间为(0.11,0.41), 因此推断两者总体均值不相等。

根据NHST结果,  $p = 3.61 \times 10^{-21}$ , 因此也可以推断出总体均值不相等。

effect size(Cohen' s d)=0.96>0.8, 说明两者有比较明显的差异。

## 1.2 HW-U3-2 : ANOVA 睡眠治疗实验

(1) 表单SleepExp\_1.csv 是招募60名被试, 随机分成三种不同剂量组 (10mg, 50mg, 100mg)进行试验, 表单Scores是治疗后被试的评分, 请推断不同剂量组间是否有治疗效果差异?

```
[9]: df=pd.read_csv('SleepExp_1.csv')

amount=['10mg','50mg','100mg']
data=[]
# scipy.stats
for i in range(3):
    data.append(df[df['Dosage']==amount[i]]['Scores'])
f,p=stats.f_oneway(data[0],data[1],data[2])
f,p
# pg
df.anova(dv='Scores',between='Dosage')
```

[9]: (10.480888179350163, 0.00013298547134746072)

```
[9]:      Source  ddof1  ddof2      F      p-unc    np2
0 Dosage      2      57  10.481  0.000133  0.269
```

$p=0.000133 < 0.05$ , 说明不同剂量组间存在显著治疗效果差异。

**(2) 表单SleepExp\_2.csv是招募20名被试，每个被试连续进行了三种剂量治疗的（10mg, 50mg, 100mg）实验，表单Scores是每个剂量治疗后被试的评分，请推断不同剂量组间是否有治疗效果差异？**

```
[10]: df=pd.read_csv('SleepExp_2.csv')
      res=anova.AnovaRM(df, 'Scores', 'Subjects', within=['Dosage']).fit()
      print(res)
```

```

              Anova
=====
      F Value Num DF  Den DF Pr > F
-----
Dosage  2.4209  2.0000 38.0000 0.1024
=====
```

$Pr=0.1024 > 0.05$  故推断不同剂量组间没有明显的治疗效果差异。

**(3) 表单SleepExp\_3.csv是招募了30名被试，每个被试连续进行了三种剂量治疗的（10mg, 50mg, 100mg）实验，表单Scores是每个剂量治疗后被试的评分，请推断剂量、性别、及剂量与性别相互作用的效应分别对治疗评分的影响是否显著？**

```
[11]: df=pd.read_csv('SleepExp_3.csv')
      df['Gender'].value_counts()
      pg.mixed_anova(df, 'Scores', within='Dosage', subject='Subjects', between='Gender')
```

```
[11]: Female      57
      Male       33
      Name: Gender, dtype: int64
```

```
[11]:      Source      SS  DF1  DF2      MS      F      p-unc    np2    eps
0      Gender  84.627    1   28  84.627  2.178  1.511432e-01  0.072    -
1      Dosage 931.005    2   56  465.502 19.115  4.699411e-07  0.406  0.957
2  Interaction  47.620    2   56   23.810  0.978  3.825007e-01  0.034    -
```

根据p值可知，性别和剂量与性别相互作用的效应对治疗评分影响显著，而剂量的影响并不显著。

**(4) 表单SleepExp\_4.csv是招募了15名被试，每个被试分别在春季，秋季都连续进行了三种剂量治疗的（10mg, 50mg, 100mg）实验，表单Scores是每个剂量治疗后被试的评分，请推断剂量、季节、及剂量与季节相互作用的效应分别对治疗评分的影响是否显著？**

```
[12]: df=pd.read_csv('SleepExp_4.csv')
res=anova.AnovaRM(df, 'Scores', 'Subjects', within=['Dosage', 'Season']).fit()
print(res)
pg.rm_anova(df, 'Scores', within=['Dosage', 'Season'], subject='Subjects')
```

```

              Anova
=====
              F Value Num DF  Den DF Pr > F
-----
Dosage          36.5660  2.0000 28.0000 0.0000
Season           6.7376  1.0000 14.0000 0.0212
Dosage:Season    0.8290  2.0000 28.0000 0.4469
=====
```

```
[12]:
```

	Source	SS	ddof1	ddof2	MS	F	p-unc	\
0	Dosage	2400.691	2	28	1200.346	36.566	1.555154e-08	
1	Season	368.954	1	14	368.954	6.738	2.115423e-02	
2	Dosage * Season	78.718	2	28	39.359	0.829	4.468987e-01	

  

	p-GG-corr	np2	eps
0	3.795746e-07	0.723	0.788
1	2.115423e-02	0.325	1.000
2	4.273312e-01	0.056	0.818

由于 $0.0000, 0.0212 < 0.05$ , 而 $0.4469 > 0.05$ , 因此认为剂量和季节对治疗评分影响显著而两者相互作用的效应并不显著。

**(5) 表单SleepExp\_5.csv从上海、北京招募了90名被试, 随机分成三种剂量治疗组 (10mg, 50mg, 100mg) 进行睡眠实验, 表单Scores是每个被试治疗后的评分, 请推断剂量、城市、及剂量与城市相互作用的效应分别对治疗评分的影响是否显著?**

```
[13]: df=pd.read_csv('SleepExp_5.csv')
# pg
df.anova(dv='Scores', between=['Dosage', 'City'])
# statsmodels
model = ols('Scores ~ Dosage*City', df).fit()
anova_table = sm.stats.anova_lm(model, typ=2)
print(anova_table)
```

```
[13]:
```

	Source	SS	DF	MS	F	p-unc	np2
0	Dosage	3820.314	2	1910.157	90.684590	1.042341e-21	0.683460
1	City	16.548	1	16.548	0.785615	3.779598e-01	0.009266
2	Dosage * City	292.849	2	146.425	6.951497	1.608049e-03	0.142008
3	Residual	1769.355	84	21.064	NaN	NaN	NaN

  

	sum_sq	df	F	PR(>F)
Dosage	3820.313727	2.0	90.684583	1.042343e-21

City	16.547897	1.0	0.785610	3.779613e-01
Dosage:City	292.849482	2.0	6.951506	1.608036e-03
Residual	1769.354510	84.0	NaN	NaN

由p值可知，剂量对治疗评分影响显著，城市对治疗评分的影响不显著，但剂量和城市相互作用对治疗评分的影响比较显著。