Sardar Patel Institute of Technology,Mumbai
Department of Electronics and Telecommunication Engineering
B.E. Sem-VII (2022-2023) Data Analytics

## Experiment: Exploratory Data Analysis (EDA)

**Name: Rahul Alshi**     **UID:** 2019110001     **BE ETRX**     **DA LAB 3**

**Aim:** Perform Statistical Analysis (tests) on SAS

**Dataset Overview**

The dataset 'SASHELP.BASEBALL' contains 24 columns with 322 rows containing the information of each baseball player from different countries of the world in his career. :

Code :

```
1  data first;
2  set SASHELP.BASEBALL;
3  run;
4  proc means data= SASHELP.BASEBALL mean median mode std var min max;
5  run;
6
7  proc means data= SASHELP.BASEBALL nmiss;
8  run;
9
10 proc print =SASHELP.BASEBALL;
11 where Div = "NW";
12 run;
13
14 proc sql;
15 select count(distinct Div) as Div,
16        count(distinct Team) as Team,
17        count(distinct CrHits) as CrHits
18   from SASHELP.BASEBALL;
19 quit;
20
21 proc freq data=SASHELP.BASEBALL;
22    tables Team;    /* _ALL_ is the defaul */
23 run;
24
25 data nHits;
26 set SASHELP.BASEBALL(keep=_NUMERIC_            /* all numeric variables */
27                  ); /* two character variables */
28 run;
29
30 proc print data= nHits(obs=5);
31 run;
32
```

```sas
18    from SASHELP.BASEBALL;
19  quit;
20
21  proc freq data=SASHELP.BASEBALL;
22    tables Team;    /* _ALL_ is the defaul */
23  run;
24
25  data nHits;
26  set SASHELP.BASEBALL(keep=_NUMERIC_              /* all numeric variables */
27                    ); /* two character variables */
28  run;
29
30  proc print data= nHits(obs=5);
31  run;
32
33  proc means data=nHits nmiss;
34  run;
35
36  ods graphics / reset width=6.4in height=4.8in imagemap;
37  proc sgplot data=SASHELP.BASEBALL;
38  vbox  nAtBat / category=nHome;
39  yaxis grid;
40  run;
41  ods graphics / reset;
42
43  proc ttest data = SASHELP.BASEBALL SIDES=L;
44  class nAtBat;
45  var nHits;
46  run;
47
48
49
50
```
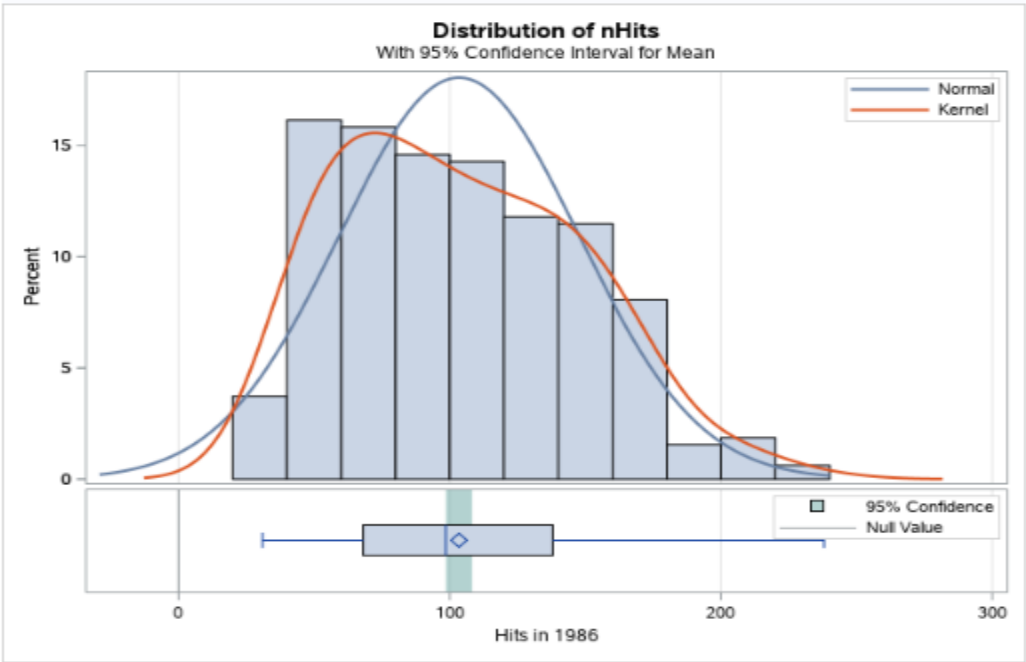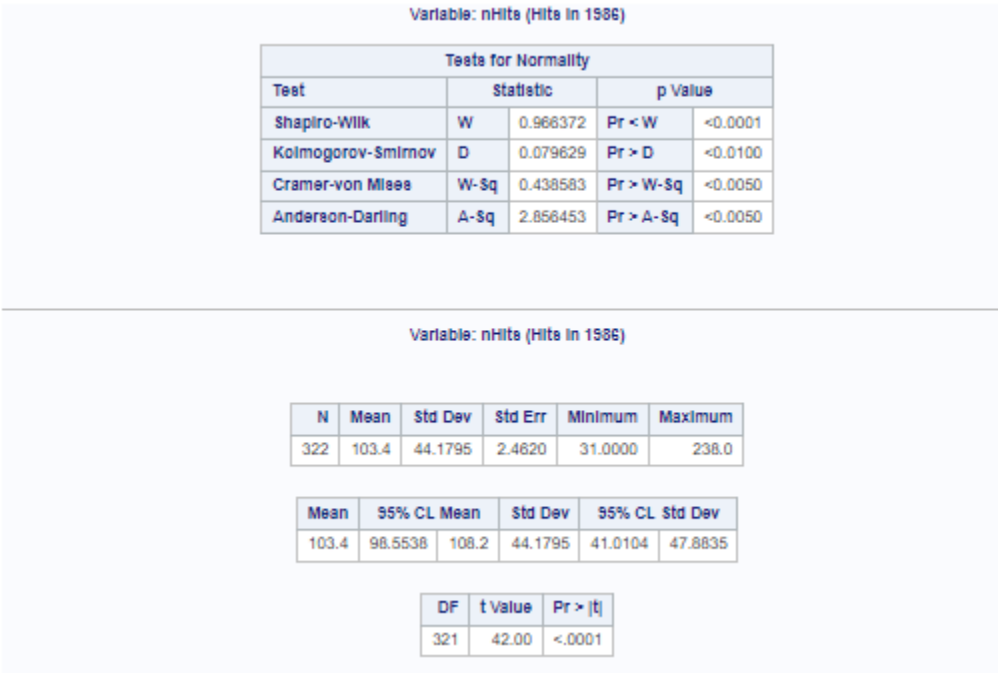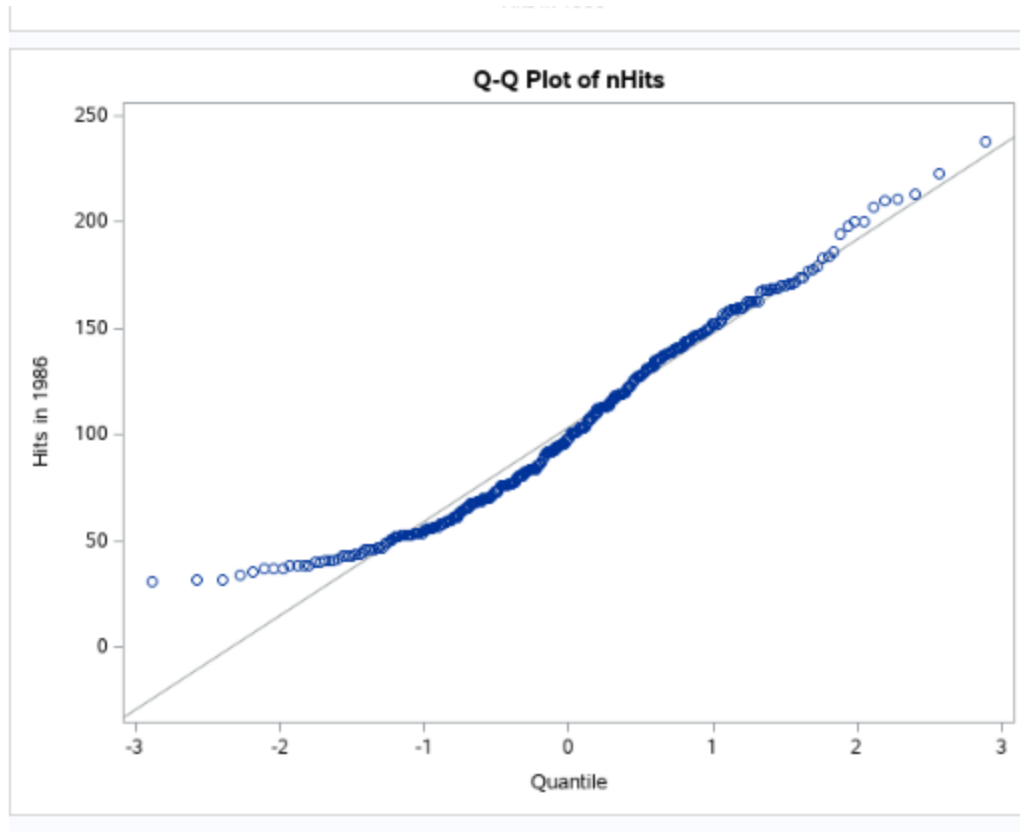
```sas
 1  /*
 2   *
 3   * Task code generated by SAS Studio 3.8
 4   *
 5   * Generated on '11/16/22, 2:53 PM'
 6   * Generated by 'u62384181'
 7   * Generated on server 'ODAWS01-APSE1.ODA.SAS.COM'
 8   * Generated on SAS platform 'Linux LIN X64 3.10.0-1062.9.1.el7.x86_64'
 9   * Generated on SAS version '9.04.01M6P11072018'
10   * Generated on browser 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) Apple
11   * Generated on web client 'https://odamid-apse1.oda.sas.com/SASStudio/m
12   *
13   */
14
15  ods noproctitle;
16  ods graphics / imagemap=on;
17
18  /* Test for normality */
19  proc univariate data=SASHELP.BASEBALL normal mu0=0;
20      ods select TestsForNormality;
21      var nAtBat;
22  run;
23
24  /* t test */
25  proc ttest data=SASHELP.BASEBALL sides=2 h0=0 plots(showh0);
26      var nAtBat;
27  run;
```

Output :

**Tests for Normality**

| Test | Statistic | | p Value | |
|---|---|---|---|---|
| Shapiro-Wilk | W | 0.966372 | Pr < W | <0.0001 |
| Kolmogorov-Smirnov | D | 0.079629 | Pr > D | <0.0100 |
| Cramer-von Mises | W-Sq | 0.438583 | Pr > W-Sq | <0.0050 |
| Anderson-Darling | A-Sq | 2.856453 | Pr > A-Sq | <0.0050 |

Variable: nHits (Hits in 1986)

| N | Mean | Std Dev | Std Err | Minimum | Maximum |
|---|---|---|---|---|---|
| 322 | 103.4 | 44.1795 | 2.4620 | 31.0000 | 238.0 |

| Mean | 95% CL Mean | | Std Dev | 95% CL Std Dev | |
|---|---|---|---|---|---|
| 103.4 | 98.5538 | 108.2 | 44.1795 | 41.0104 | 47.8835 |

| DF | t Value | Pr > |t| |
|---|---|---|
| 321 | 42.00 | <.0001 |



Distribution of nHits — With 95% Confidence Interval for Mean

**Q-Q Plot of nHits**

**Conclusion:**

1. Performed statistical analysis (Hypothesis testing) on SAS Studio for baseball dataset.

2. A t test is a statistical test that is used to compare the means of two groups. It is often used in hypothesis testing to determine whether a process or treatment actually has an effect on the population of interest, or whether two groups are different from one another.

3. The t value for the one sample t test for the no of hits in the year 1986 is calculated to be 42

4. Few insights we found from the dataset:
   - With 95% confidence interval for mean on distribution of n Hits , the distribution is highest for 100 Hits in 1986 with above 15%
   - The no of hits goes on increasing per quantile upto 250 from the Q-Q plot.
   - Similarly we can perform t- tests for 2 sampled and 1 sampled means wwith different attributes in the dataset.