Name: Amit Birajdar

Class: BTech CS-B

Roll No.: B014

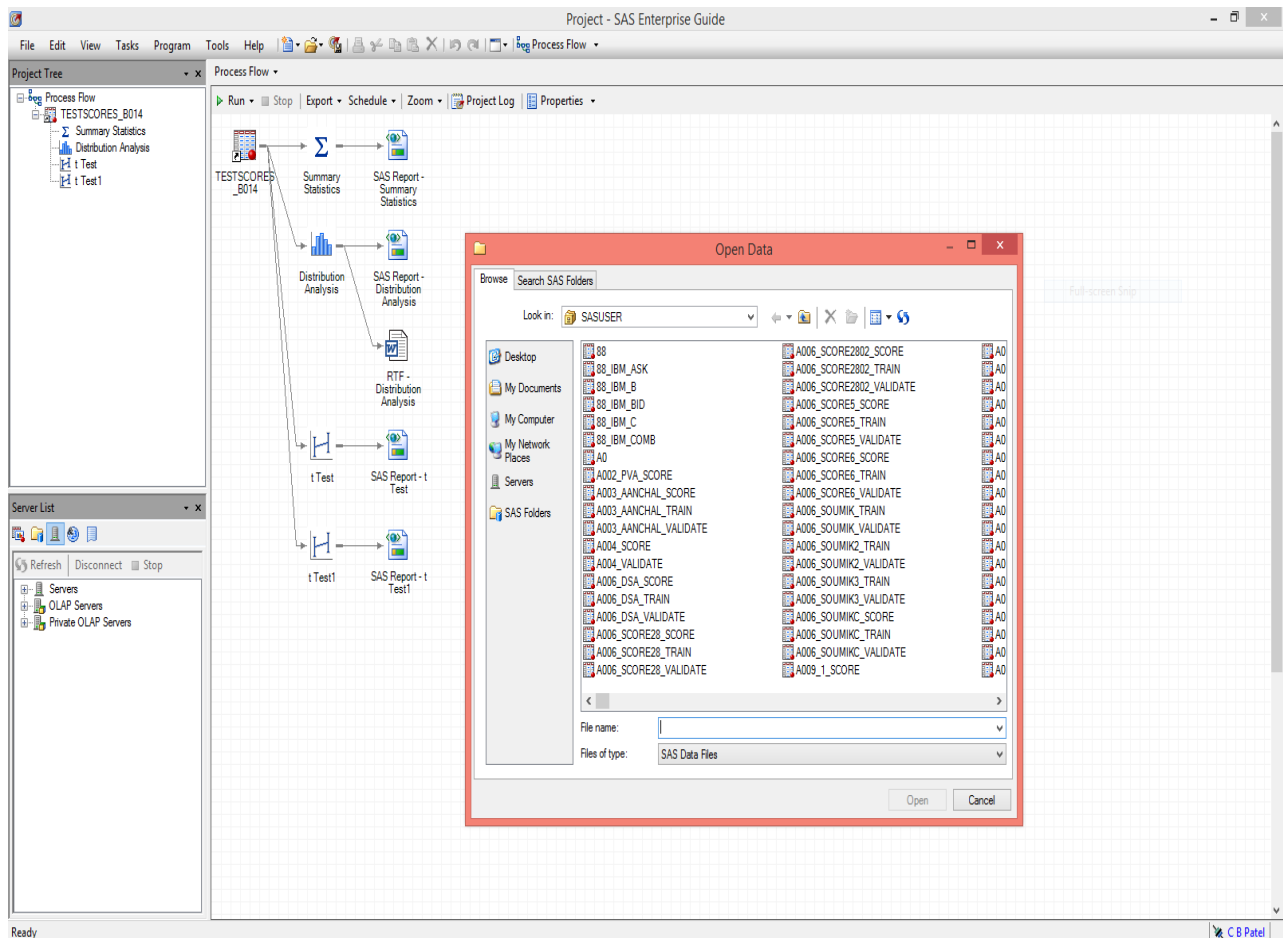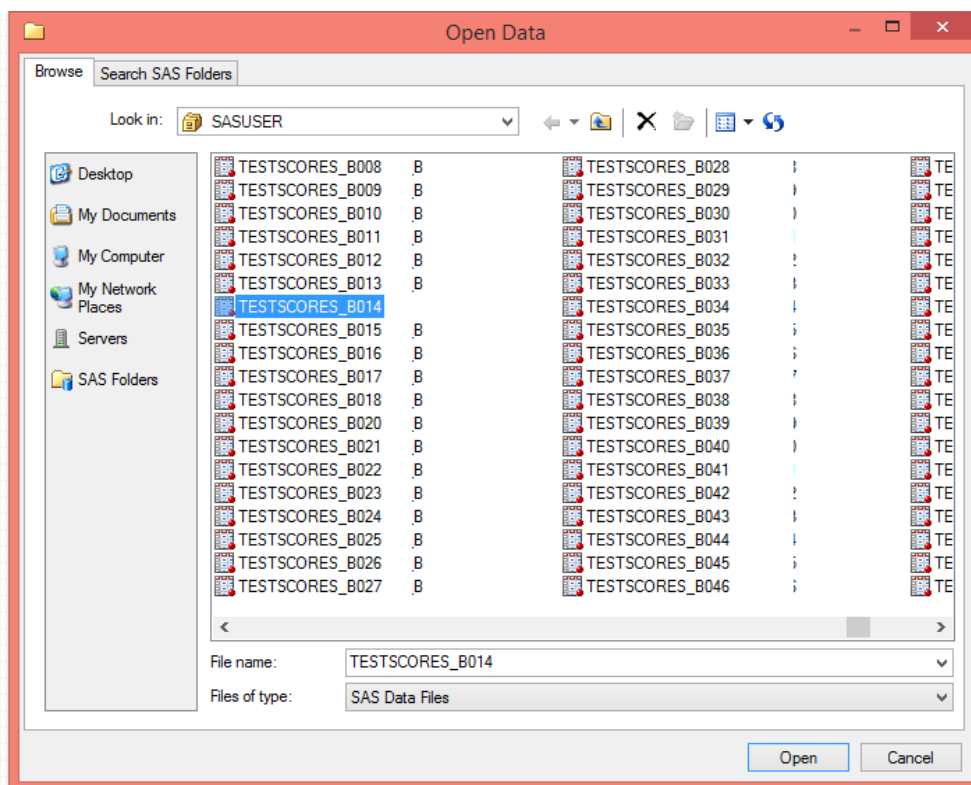# PREDICTIVE MODELLING - ASSIGNMENT 1

**Aim**: To implement:

1. Filter and Sort
2. Statistical Analysis
3. Distribution Analysis
4. Confidence Intervals
5. t-test – One sample and Two sample

**Data-set used**: SAT Test Score

**Importing a data source:**

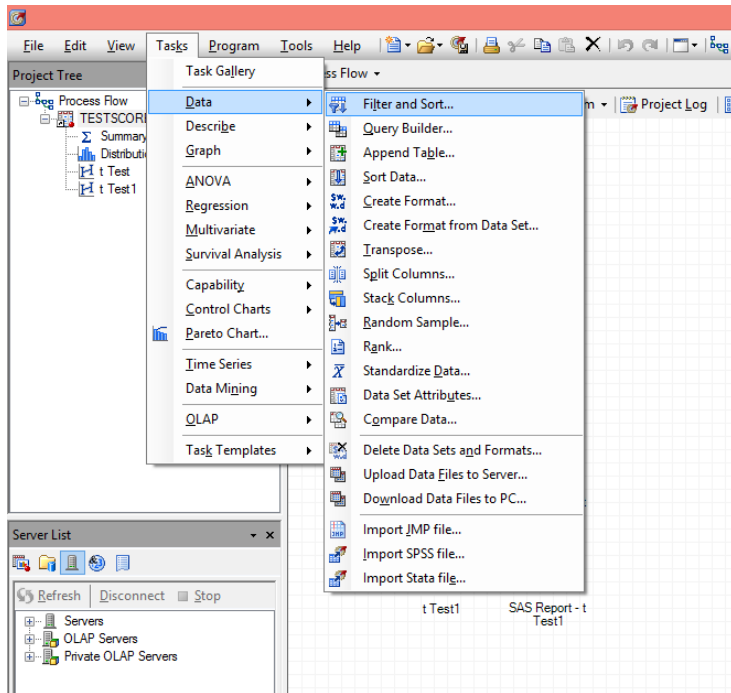File -> Open -> Data -> Servers -> Local -> SASUSER library -> TESTSCORES_B014

**Open Data**

Browse | Search SAS Folders

Look in: SASUSER

- Desktop
- My Documents
- My Computer
- My Network Places
- Servers
- SAS Folders

| | | |
|---|---|---|
| TESTSCORES_B008 | TESTSCORES_B028 | TE |
| TESTSCORES_B009 | TESTSCORES_B029 | TE |
| TESTSCORES_B010 | TESTSCORES_B030 | TE |
| TESTSCORES_B011 | TESTSCORES_B031 | TE |
| TESTSCORES_B012 | TESTSCORES_B032 | TE |
| TESTSCORES_B013 | TESTSCORES_B033 | TE |
| TESTSCORES_B014 | TESTSCORES_B034 | TE |
| TESTSCORES_B015 | TESTSCORES_B035 | TE |
| TESTSCORES_B016 | TESTSCORES_B036 | TE |
| TESTSCORES_B017 | TESTSCORES_B037 | TE |
| TESTSCORES_B018 | TESTSCORES_B038 | TE |
| TESTSCORES_B020 | TESTSCORES_B039 | TE |
| TESTSCORES_B021 | TESTSCORES_B040 | TE |
| TESTSCORES_B022 | TESTSCORES_B041 | TE |
| TESTSCORES_B023 | TESTSCORES_B042 | TE |
| TESTSCORES_B024 | TESTSCORES_B043 | TE |
| TESTSCORES_B025 | TESTSCORES_B044 | TE |
| TESTSCORES_B026 | TESTSCORES_B045 | TE |
| TESTSCORES_B027 | TESTSCORES_B046 | TE |

File name: TESTSCORES_B014

Files of type: SAS Data Files

Open | Cancel

Double click on the data set in the process flow window to view its contents.

TESTSCORES_B014

Filter and Sort | Query Builder | Data ▾ | Describe ▾ | Graph ▾ | Analyze ▾ | Export ▾ | Send To ▾ |

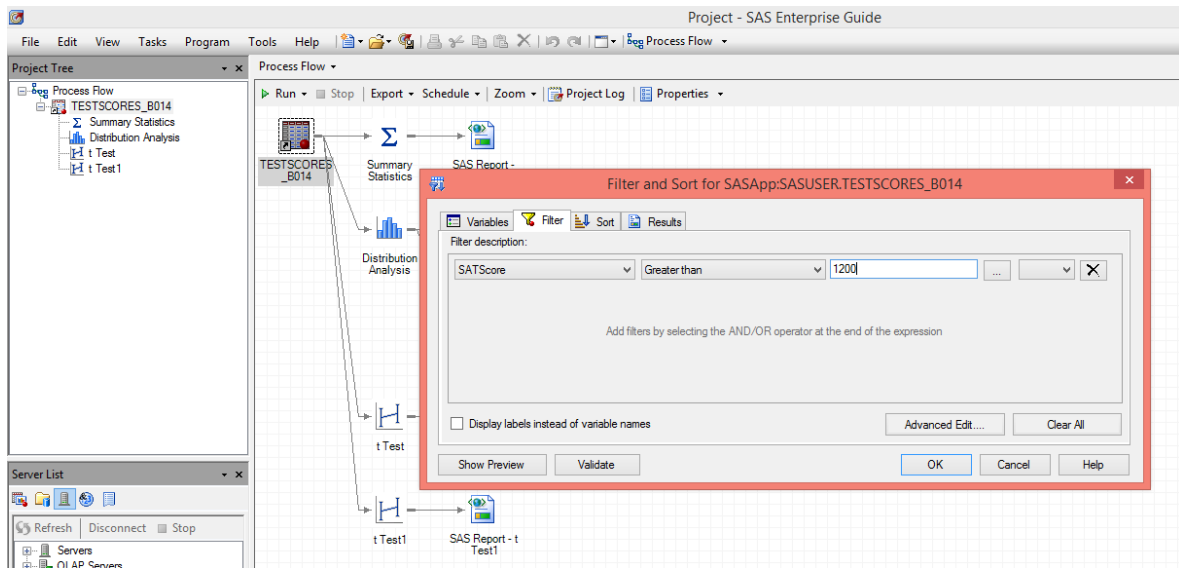| | Gender | SATScore | IDNumber |
|---|---|---|---|
| 1 | Male | 1170 | 61469897 |
| 2 | Female | 1090 | 33081197 |
| 3 | Male | 1240 | 68137597 |
| 4 | Female | 1000 | 37070397 |
| 5 | Male | 1210 | 64608797 |
| 6 | Female | 970 | 60714297 |
| 7 | Male | 1020 | 16907997 |
| 8 | Female | 1490 | 9589297 |
| 9 | Male | 1200 | 93891897 |
| 10 | Female | 1260 | 85859397 |
| 11 | Male | 1150 | 38152597 |
| 12 | Female | 1390 | 99108497 |
| 13 | Male | 1240 | 59666697 |
| 14 | Female | 1370 | 70847197 |
| 15 | Male | 1140 | 47613397 |
| 16 | Female | 1160 | 53750297 |
| 17 | Male | 1050 | 95948597 |
| 18 | Female | 1110 | 3873197 |
| 19 | Male | 1100 | 25756097 |
| 20 | Female | 1080 | 43493297 |
| 21 | Male | 1120 | 27543197 |
| 22 | Female | 1080 | 26212897 |
| 23 | Male | 1050 | 8945097 |
| 24 | Female | 1200 | 51799397 |
| 25 | Male | 1600 | 39196697 |
| 26 | Female | 1100 | 48154497 |
| 27 | Male | 1050 | 55189597 |
| 28 | Female | 1060 | 46028397 |
| 29 | Male | 1140 | 75332897 |
| 30 | Female | 1100 | 29520797 |
| 31 | Male | 1340 | 55983497 |
| 32 | Female | 1240 | 93236497 |
| 33 | Male | 1090 | 6975697 |
| 34 | Female | 1180 | 29686297 |
| 35 | Male | 1170 | 76815697 |
| 36 | Female | 1130 | 64045497 |
| 37 | Male | 1290 | 9880297 |
| 38 | Female | 1380 | 23048597 |
| 39 | Male | 1010 | 76058697 |
| 40 | Female | 1280 | 42586897 |
| 41 | Male | 1050 | 62688897 |
| 42 | Female | 1520 | 73461797 |
| 43 | Male | 1360 | 44327297 |
| 44 | Female | 1260 | 2854197 |

**Filter and Sort:**

1. Select data-set on the process window
2. Click on 'Tasks' in the menu bar
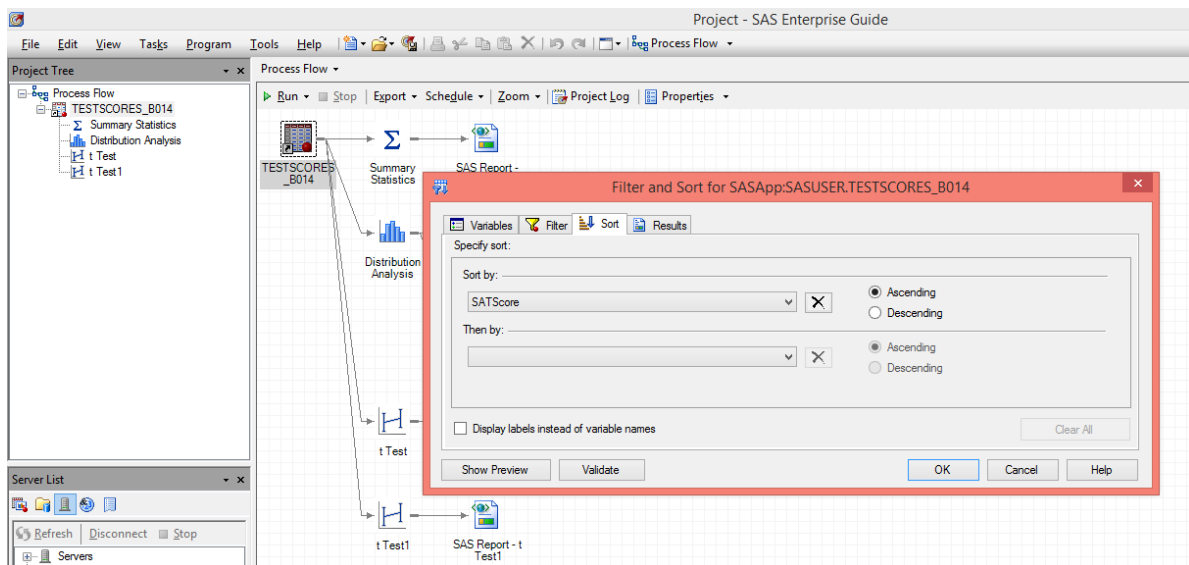   a. Select Data -> Filter and Sort (this opens a filter and sort wizard)



Select variables on the basis of which you want to filter the data-set

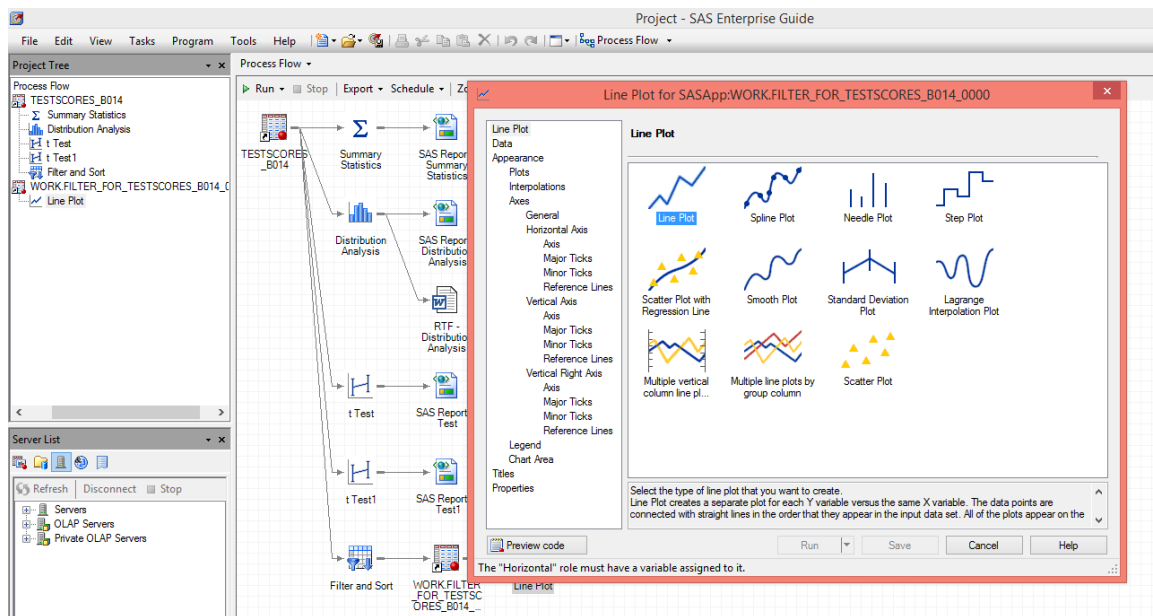Set filter conditions by selecting variable, filtering condition and value.



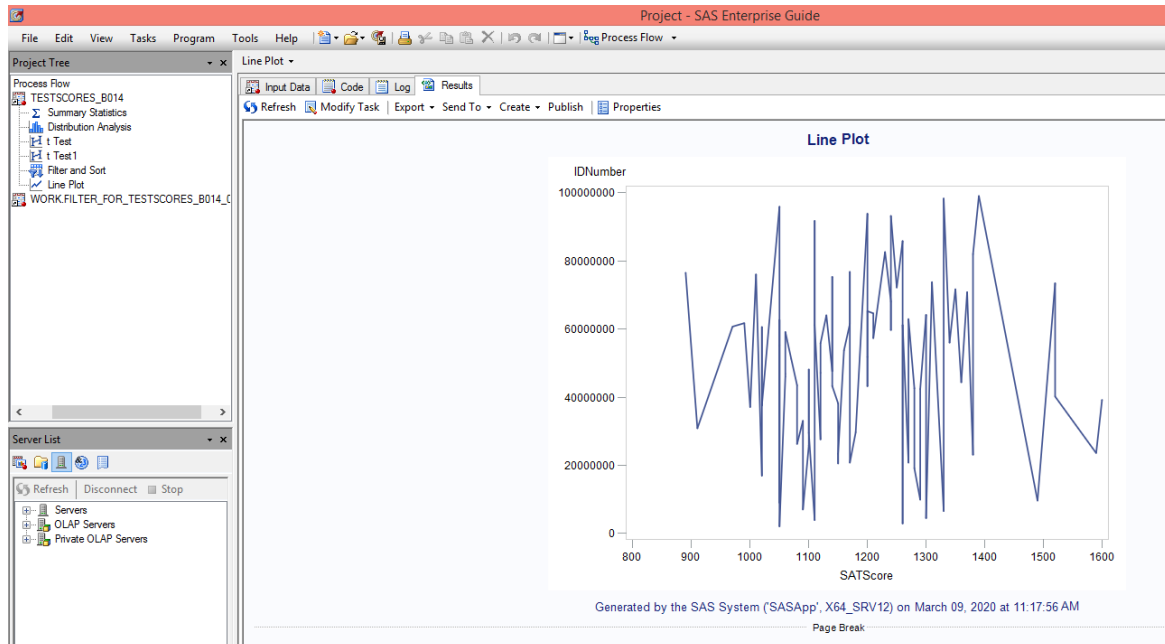Select sort and apply sorting conditions



When you are done, click 'OK' to apply the filters and sorting conditions to the data-set.

You will observe a new data-set being created in the process window that is derived from the main data-set after applying filter and sort function. Double click on this new data-set to view its contents.

**Line Plot:**



View line plot by double clicking the node on the process window.

## Summary Statistics:

Select data-set -> Tasks –> Describe -> Summary statistics
This will open a summary statistics wizard.
Select the variables for analysis

Summary statistics result shows the selected values, ie mean, standard deviation, variance, min,mac, range, quartiles and confidence levels for the selected data-set.

### Summary Statistics
### Results
### The MEANS Procedure

| | | | | | | | | | | Lower 95% | Upper 95% |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Mean | Std Dev | Variance | Minimum | Maximum | Range | N | Lower Quartile | Median | Upper Quartile | CL for Mean | CL for Mean |
| 1190.63 | 147.0584466 | 21626.19 | 890.0000000 | 1600.00 | 710.0000000 | 80 | 1085.00 | 1170.00 | 1280.00 | 1157.90 | 1223.35 |

Generated by the SAS System ('SASApp', X64_SRV12) on March 09, 2020 at 11:21:45 AM

The histogram plot shows the percentage distribution of SAT scores.



Box plot indicates the minimum, 1st quartile (25th percentile), median, 3rd quartile (75th percentile) and maximum values for the data-set.

## Distribution Analysis:

Select data-set -> Tasks –> Describe -> Distribution Analysis

# Distribution analysis of: SATScore BY B014

## The UNIVARIATE Procedure
### Variable: SATScore

| Basic Statistical Measures | | | |
|---|---|---|---|
| **Location** | | **Variability** | |
| Mean | 1190.625 | Std Deviation | 147.05845 |
| Median | 1170.000 | Variance | 21626 |
| Mode | 1050.000 | Range | 710.00000 |
| | | Interquartile Range | 195.00000 |

| Basic Confidence Limits Assuming Normality | | | |
|---|---|---|---|
| **Parameter** | **Estimate** | **95% Confidence Limits** | |
| Mean | 1191 | 1158 | 1223 |
| Std Deviation | 147.05845 | 127.27215 | 174.18670 |
| Variance | 21626 | 16198 | 30341 |

| Tests for Location: Mu0=0 | | | | |
|---|---|---|---|---|
| **Test** | | **Statistic** | **p Value** | |
| Student's t | t | 72.41525 | Pr > \|t\| | <.0001 |
| Sign | M | 40 | Pr >= \|M\| | <.0001 |
| Signed Rank | S | 1620 | Pr >= \|S\| | <.0001 |

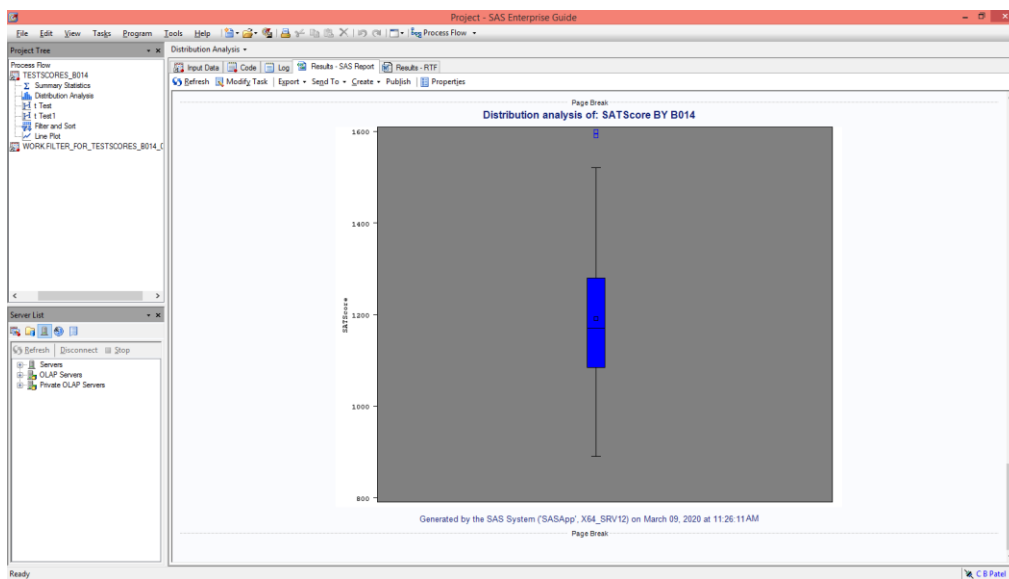Generated by the SAS System ('SASApp', X64_SRV12) on March 09, 2020 at 11:26:11 AM

Input Data | Code | Log | Results - SAS Report | Results - RTF
Refresh | Modify Task | Export | Send To | Create | Publish | Properties

Page Break

**Distribution analysis of: SATScore BY B014**

The UNIVARIATE Procedure

| N | 80 |
| Variance | 21626.19 |
| Skewness | 0.64202 |
| Kurtosis | 0.4241 |
| Number of Obs | 80 |

Generated by the SAS System ('SASApp', X64_SRV12) on March 09, 2020 at 11:26:11AM

Page Break

Page Break

# Distribution analysis of: SATScore BY B014

## The UNIVARIATE Procedure
### Fitted Normal Distribution for SATScore

| Parameters for Normal Distribution | | |
|---|---|---|
| Parameter | Symbol | Estimate |
| Mean | Mu | 1190.625 |
| Std Dev | Sigma | 147.0584 |

| Goodness-of-Fit Tests for Normal Distribution | | | | |
|---|---|---|---|---|
| Test | | Statistic | p Value | |
| Kolmogorov-Smirnov | D | 0.08382224 | Pr > D | >0.150 |
| Cramer-von Mises | W-Sq | 0.09964577 | Pr > W-Sq | 0.114 |
| Anderson-Darling | A-Sq | 0.70124822 | Pr > A-Sq | 0.068 |

| Quantiles for Normal Distribution | | |
|---|---|---|
| | Quantile | |
| Percent | Observed | Estimated |
| 1.0 | 890.000 | 848.516 |
| 5.0 | 995.000 | 948.735 |
| 10.0 | 1020.000 | 1002.162 |
| 25.0 | 1085.000 | 1091.436 |
| 50.0 | 1170.000 | 1190.625 |
| 75.0 | 1280.000 | 1289.814 |
| 90.0 | 1375.000 | 1379.088 |
| 95.0 | 1505.000 | 1432.515 |
| 99.0 | 1600.000 | 1532.734 |

Generated by the SAS System ('SASApp', X64_SRV12) on March 09, 2020 at 11:26:11AM

**T-Test:**

**One Sample**

Enter the null hypothesis condition (value of Ho) and specify confidence level.

# t Test

## The TTEST Procedure

### Variable: SATScore

| N | Mean | Std Dev | Std Err | Minimum | Maximum |
|---|---|---|---|---|---|
| 80 | 1190.6 | 147.1 | 16.4416 | 890.0 | 1600.0 |

| Mean | 95% CL Mean | | Std Dev | 95% CL Std Dev | |
|---|---|---|---|---|---|
| 1190.6 | 1157.9 | 1223.4 | 147.1 | 127.3 | 174.2 |

| DF | t Value | Pr > \|t\| |
|---|---|---|
| 79 | -0.57 | 0.5702 |



**Distribution of SATScore**
With 95% Confidence Interval for Mean

Generated by the SAS System ('SASApp', X64_SRV12) on March 09, 2020 at 11:31:25 AM

Page Break

T-test will either accept or fail to accept a claim (null hypothesis) based on the comparison of P value and alpha value (confidence interval)

## Two sample:

Ho is 0 to indicate that the null hypothesis is: values for the two variables are equal.



## t Test

### The TTEST Procedure

### Variable: SATScore

| Gender | N | Mean | Std Dev | Std Err | Minimum | Maximum |
|--------|-----|--------|---------|---------|---------|---------|
| Female | 40 | 1221.0 | 157.4 | 24.8864 | 910.0 | 1590.0 |
| Male | 40 | 1160.3 | 130.9 | 20.7008 | 890.0 | 1600.0 |
| Diff (1-2) | | 60.7500 | 144.8 | 32.3706 | | |

| Gender | Method | Mean | 95% CL Mean | | Std Dev | 95% CL Std Dev | |
|--------|--------|------|------|------|---------|------|------|
| Female | | 1221.0 | 1170.7 | 1271.3 | 157.4 | 128.9 | 202.1 |
| Male | | 1160.3 | 1118.4 | 1202.1 | 130.9 | 107.2 | 168.1 |
| Diff (1-2) | Pooled | 60.7500 | -3.6950 | 125.2 | 144.8 | 125.2 | 171.7 |
| Diff (1-2) | Satterthwaite | 60.7500 | -3.7286 | 125.2 | | | |

| Method | Variances | DF | t Value | Pr > \|t\| |
|--------|-----------|------|---------|---------|
| Pooled | Equal | 78 | 1.88 | 0.0643 |
| Satterthwaite | Unequal | 75.497 | 1.88 | 0.0644 |

| Equality of Variances | | | | |
|--------|---------|--------|---------|--------|
| Method | Num DF | Den DF | F Value | Pr > F |
| Folded F | 39 | 39 | 1.45 | 0.2545 |

Distribution of SATScore

Generated by the SAS System ('SASApp', X64_SRV12) on March 09, 2020 at 11:32:55 AM

Page Break

You either accept or fail to accept the null hypothesis based on comparison of the probability value obtained and the alpha value (confidence value).

**Final process flow:**