

Data Preparation

Look at the attribute type; e.g., nominal, ordinal or quantitative.

Categorical/Nominal Attributes:

1. State: Customer's state.
2. Area Code: Area code of customer.
3. Phone: Phone number of customer.
4. Intl_Plan: Binary indicator showing whether the customer has an international calling plan (Yes/No).
5. VMail_Plan: Indicator of voicemail plan (Yes/No).
6. Churn: Class attribute with binary values (True for churn and False for not churn).

Numerical/Quantitative Attributes:

1. Account Length: Integer number showing the duration of activity for the customer account.
2. CustServ Calls: The number of calls to customer support service.
3. Day Calls: Discrete attribute indicating the total number of calls during daytime.
4. Day Charge: Charges for using the service during daytime (continuous data type).
5. Day Mins: The number of minutes the customer used the service during daytime (continuous quantitative data type).
6. Eve Calls: The number of calls during evening time.
7. Eve Charge: Charges for using the service during evening time (continuous data type).
8. Eve Mins: The number of minutes the customer used the service during evening time (continuous quantitative data type).
9. Intl Calls: The number of international calls.
10. Intl Charge: Charges for international calls (continuous data type).
11. Intl Mins: The number of minutes the customer used the service to make international calls (continuous quantitative data type).
12. Night Calls: The number of calls during night time.
13. Night Charge: Charges for using the service during night time (continuous data type).
14. Night Mins: Number of minutes the customer used the service during night time (continuous quantitative data type).
15. VMail Message: The number of voicemail messages.

Determine any outlier values (records) for each of the attributes or attributes under consideration (min, max, std. dev, scatter plots, box plots or others can be used).

State – Categorical

Outliers in State Based on Frequency Distribution							
Obs	State	COUNT	PERCENT	_TYPE_	_FREQ_	mean_freq	std_freq
1	CA	34	1.02010	0	51	65.3529	11.8014
2	WV	106	3.18032	0	51	65.3529	11.8014

State	Count	Percent	Type	Freq	Mean	Std
AL	1	0.00309	0	51	65.3529	11.8014
AK	1	0.00309	0	51	65.3529	11.8014
AZ	1	0.00309	0	51	65.3529	11.8014
AR	1	0.00309	0	51	65.3529	11.8014
CA	34	1.02010	0	51	65.3529	11.8014
CO	1	0.00309	0	51	65.3529	11.8014
CT	1	0.00309	0	51	65.3529	11.8014
DE	1	0.00309	0	51	65.3529	11.8014
FL	1	0.00309	0	51	65.3529	11.8014
GA	1	0.00309	0	51	65.3529	11.8014
HI	1	0.00309	0	51	65.3529	11.8014
ID	1	0.00309	0	51	65.3529	11.8014
IL	1	0.00309	0	51	65.3529	11.8014
IN	1	0.00309	0	51	65.3529	11.8014
IA	1	0.00309	0	51	65.3529	11.8014
KS	1	0.00309	0	51	65.3529	11.8014
KY	1	0.00309	0	51	65.3529	11.8014
LA	1	0.00309	0	51	65.3529	11.8014
MA	1	0.00309	0	51	65.3529	11.8014
MD	1	0.00309	0	51	65.3529	11.8014
ME	1	0.00309	0	51	65.3529	11.8014
MI	1	0.00309	0	51	65.3529	11.8014
MN	1	0.00309	0	51	65.3529	11.8014
MO	1	0.00309	0	51	65.3529	11.8014
MS	1	0.00309	0	51	65.3529	11.8014
MT	1	0.00309	0	51	65.3529	11.8014
NE	1	0.00309	0	51	65.3529	11.8014
NH	1	0.00309	0	51	65.3529	11.8014
NJ	1	0.00309	0	51	65.3529	11.8014
NM	1	0.00309	0	51	65.3529	11.8014
NY	1	0.00309	0	51	65.3529	11.8014
NC	1	0.00309	0	51	65.3529	11.8014
ND	1	0.00309	0	51	65.3529	11.8014
OH	1	0.00309	0	51	65.3529	11.8014
OK	1	0.00309	0	51	65.3529	11.8014
OR	1	0.00309	0	51	65.3529	11.8014
PA	1	0.00309	0	51	65.3529	11.8014
RI	1	0.00309	0	51	65.3529	11.8014
SC	1	0.00309	0	51	65.3529	11.8014
SD	1	0.00309	0	51	65.3529	11.8014
TN	1	0.00309	0	51	65.3529	11.8014
TX	1	0.00309	0	51	65.3529	11.8014
UT	1	0.00309	0	51	65.3529	11.8014
VA	1	0.00309	0	51	65.3529	11.8014
VT	1	0.00309	0	51	65.3529	11.8014
WA	1	0.00309	0	51	65.3529	11.8014
WI	1	0.00309	0	51	65.3529	11.8014
WY	1	0.00309	0	51	65.3529	11.8014

- **Outlier Detection**
 - Outliers are identified using the mean and standard deviation. The typical range for normal data is within:
 - **Mean - 2 × Standard Deviation**
 - **Mean + 2 × Standard Deviation**
 - So, the outlier thresholds are:
 - **Lower Threshold:** Mean - 2 × Std Dev = 65.35 - 2 × 11.80 = 41.75
 - **Upper Threshold:** Mean + 2 × Std Dev = 65.35 + 2 × 11.80 = 89.95
 - Records that fall outside these thresholds are considered outliers.
- **Interpretation of Specific States**
 - **California (CA):**
 - **Count:** 34
 - **Percentage:** 1.02%
 - **Analysis:**
 - CA has a frequency count of 34, which is below the lower threshold of 41.75.
 - This indicates that CA's count is significantly lower than the average state frequency.
 - **Conclusion:** CA is an outlier with a notably lower frequency compared to other states.
 - **West Virginia (WV):**
 - **Count:** 106
 - **Percentage:** 3.18%
 - **Analysis:**
 - WV has a frequency count of 106, which is above the upper threshold of 89.95.
 - This indicates that WV's count is significantly higher than the average state frequency.
 - **Conclusion:** WV is an outlier with a notably higher frequency compared to other states

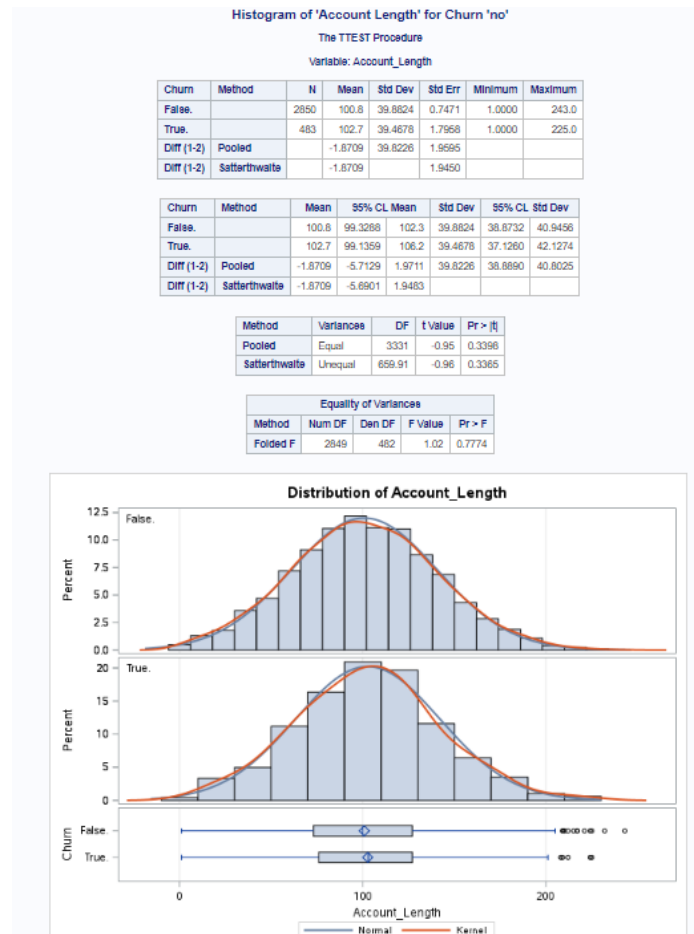
Summary

- **California (CA)** has fewer records than most other states, making it an outlier with a lower frequency.
- **West Virginia (WV)** has more records than most other states, making it an outlier with a higher frequency.

Analyze the distribution of numeric attributes (normal or other). Plot histograms for attributes of concern and analyze whether they have any influence on the class attribute. (NUMERICAL AND CLASS CORRELATIONS)

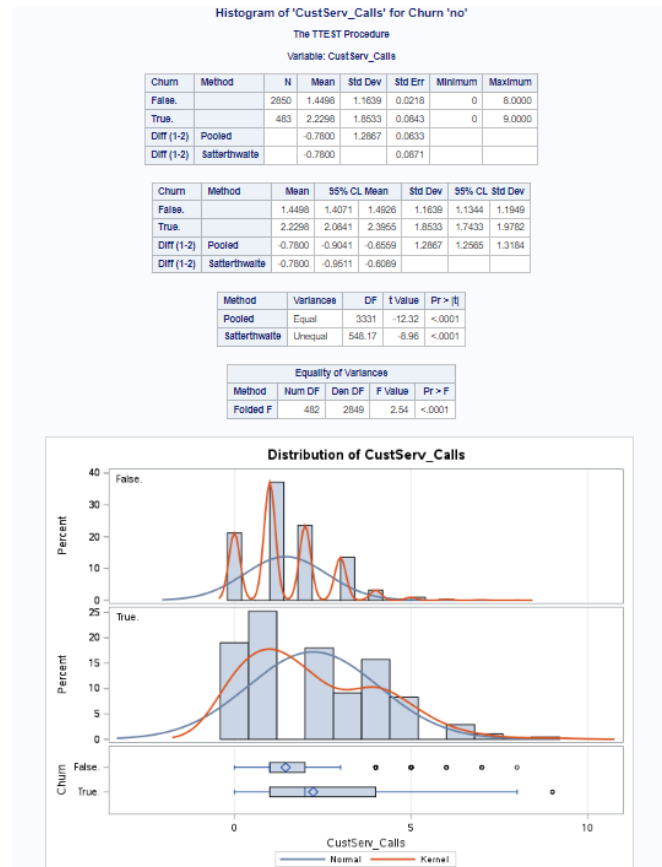
Account Length and Class Attribute

- Both the pooled and Satterthwaite t-tests show p-values of 0.3398 and 0.3365, respectively. These p-values are greater than the common significance level of 0.05, indicating that the difference in account length between the two groups is not statistically significant.
- The analysis suggests that there is no statistically significant difference in account length between clients who churned and those who did not churn. This means that account length alone is not a strong factor in predicting whether a client will churn.



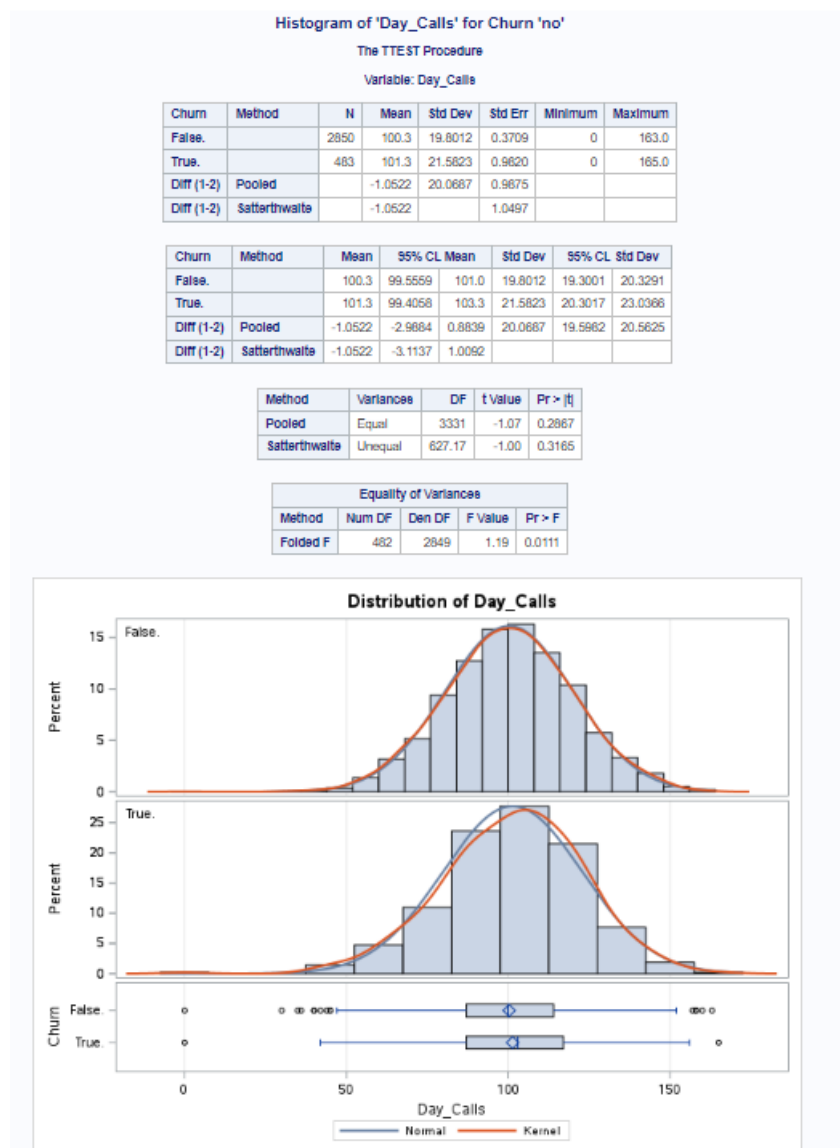
CustServ Calls & Class

- The mean number of customer service calls for clients who churned (yes) is higher (2.2298) compared to those who did not churn (no) (1.4498).
 - Both the pooled and Satterthwaite t-tests show highly significant p-values ($<.0001$), indicating that the difference in the number of customer service calls between the two groups is statistically significant.
 - The analysis reveals a statistically significant difference in the number of customer service calls between clients who have churned and those who have not. Clients who have churned tend to have a higher number of customer service calls compared to those who have not churned. This suggests that customer service calls are an influential factor in determining whether a client churns.



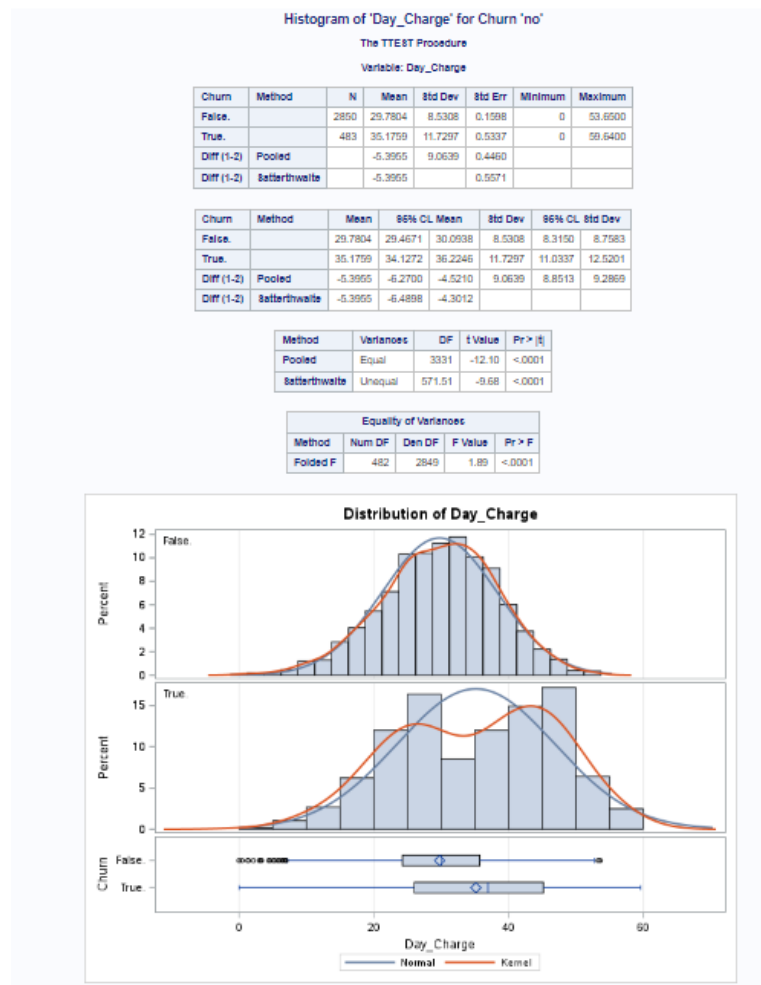
Day Calls & Class

- The mean number of Day_Calls for clients who churned (yes) is slightly higher (101.3) compared to those who did not churn (no) (100.3).
 - Both the pooled and Satterthwaite t-tests result in p-values of 0.2867 and 0.3165, respectively. These p-values are greater than the conventional significance level of 0.05, indicating that the difference in the number of Day_Calls between the churned and non-churned groups is not statistically significant.
 - The analysis suggests that there is no statistically significant difference in the number of Day_Calls between clients who churned and those who did not. This indicates that the number of Day_Calls may not be a strong factor in predicting churn status.



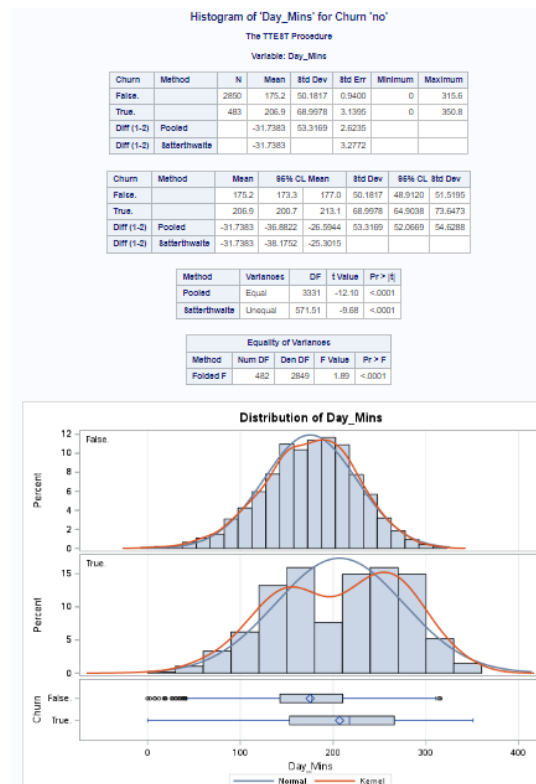
Day Charge & Class

- The mean Day_Charge for customers who did not churn (no) is lower (29.7804) compared to those who churned (yes) (35.1759).
 - Both the pooled and Satterthwaite t-tests show highly significant p-values ($<.0001$), indicating that the difference in the Day_Charge between the two groups is statistically significant.
 - The analysis suggests that there is a statistically significant difference in the Day_Charge between customers who churned and those who did not. Customers who churned (yes) tend to have higher Day_Charge compared to those who did not churn (no). This indicates that higher daytime charges may be associated with an increased likelihood of customer churn.



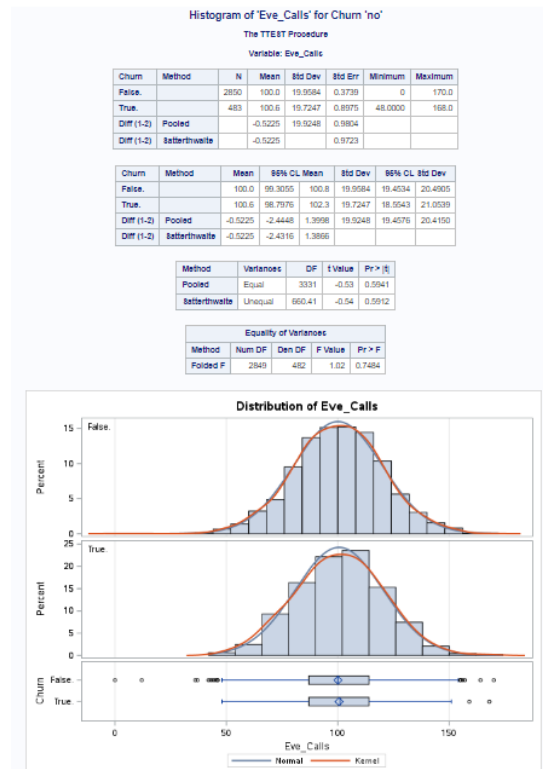
Day Mins & Class

- The mean Day_Mins for customers who churned (yes) is higher (206.9 minutes) compared to those who did not churn (no) (175.2 minutes).
 - Both the pooled and Satterthwaite t-tests show highly significant p-values ($<.0001$), indicating that the difference in Day_Mins between the two groups is statistically significant.
 - The analysis suggests that there is a statistically significant difference in the number of daytime minutes used between customers who churned and those who did not. Customers who churned tend to have a higher number of daytime minutes, implying that higher usage during the day could be a factor associated with customer churn.



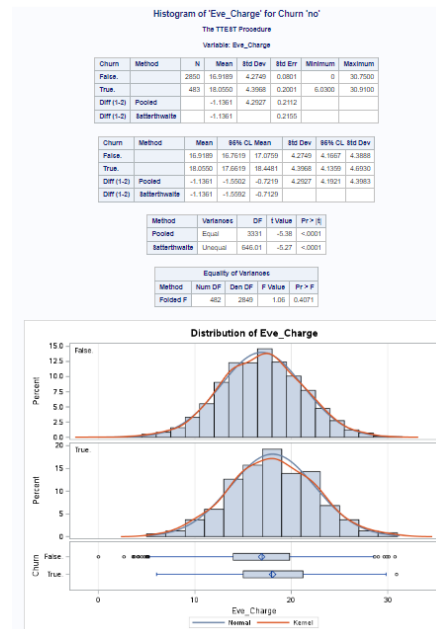
Eve Calls & Class

- The mean number of evening calls is slightly higher for clients who churned (yes) compared to those who did not churn (no), but the difference is very small (-0.5225) and not statistically significant.
 - The p-values from both the pooled and Satterthwaite t-tests are greater than 0.05, indicating that there is no strong evidence to suggest a significant difference in the number of evening calls between the churned and non-churned clients.
 - The analysis implies that the number of evening calls is not a strong predictor of churn status, as there is no significant difference between the two groups in terms of evening calls.



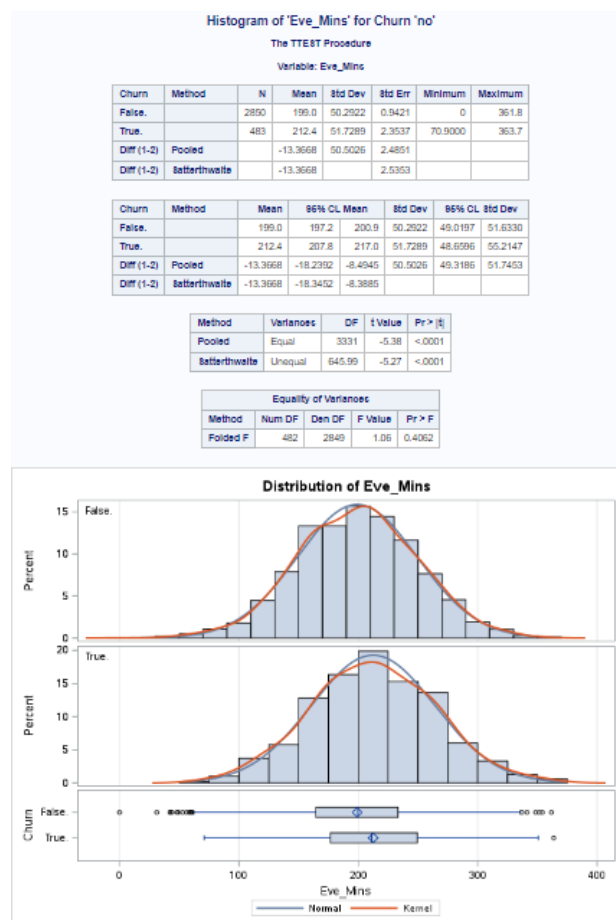
Eve Charge & Class

- The mean Eve_Charge is higher for clients who churned (yes) compared to those who did not churn (no). The mean Eve_Charge for churned clients is 18.0550, while for non-churned clients it is 16.9189, resulting in a difference of -1.1361. This difference is statistically significant.
 - The p-values from both the pooled and Satterthwaite t-tests are less than 0.0001, indicating that there is strong evidence to suggest a significant difference in Eve_Charge between churned and non-churned clients.
 - The analysis implies that Eve_Charge is a significant predictor of churn status, with higher Eve_Charge values associated with a higher likelihood of churn.



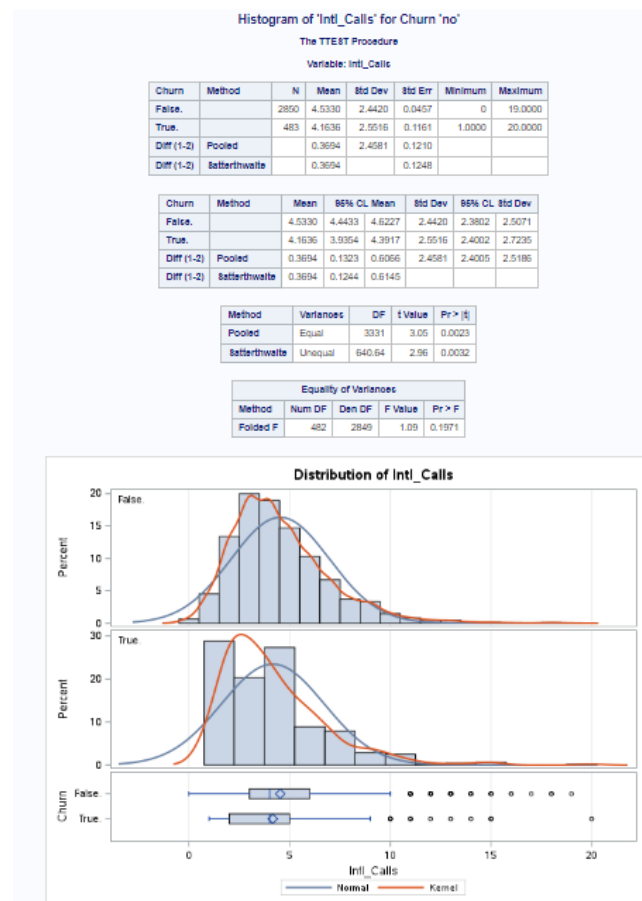
Eve Mins & Class

- The mean number of evening minutes is higher for clients who churned (yes) compared to those who did not churn (no). Specifically, clients who churned have an average of 212.4 minutes, while non-churned clients have an average of 199.0 minutes. The difference is approximately - 13.37 minutes.
 - The p-values from both the pooled and Satterthwaite t-tests are less than 0.0001, indicating that the difference in the number of evening minutes between churned and non-churned clients is statistically significant.
 - The analysis suggests that the number of evening minutes is a significant predictor of churn status. Clients who churn tend to use more evening minutes compared to those who do not churn, which implies that higher evening minutes are associated with a higher likelihood of churn.



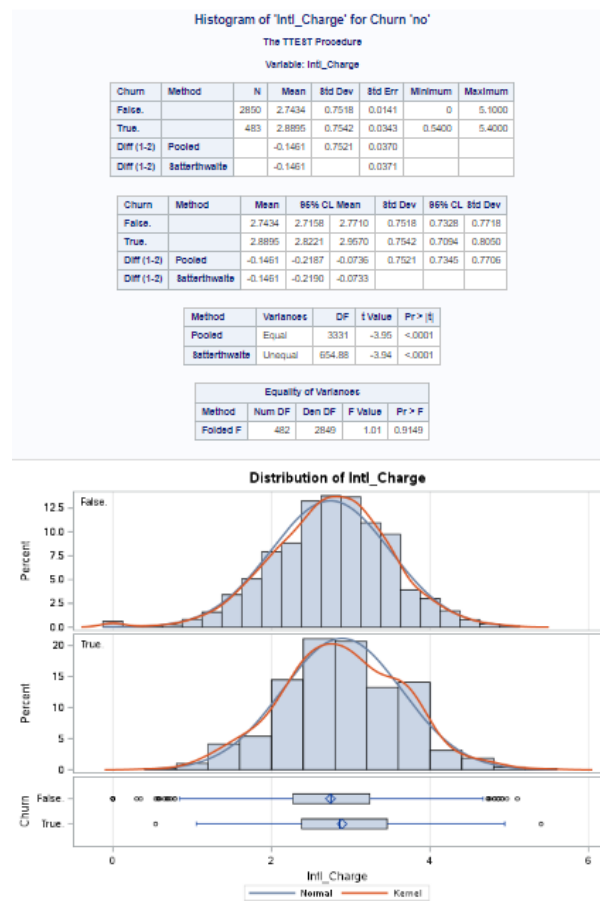
Intl Calls & Churn

- The mean number of international calls is slightly higher for clients who did not churn (no) compared to those who churned (yes), with a difference of 0.3694. This difference is statistically significant.
 - The p-values from both the pooled t-test (0.0023) and the Satterthwaite t-test (0.0032) are less than 0.05, indicating that there is a significant difference in the number of international calls between churned and non-churned clients.
 - The analysis suggests that international calls do have a significant relationship with churn status. Non-churned clients tend to make slightly more international calls compared to churned clients. This difference, though statistically significant, is relatively small, implying that while there is a difference, it may not be a strong predictor of churn on its own.



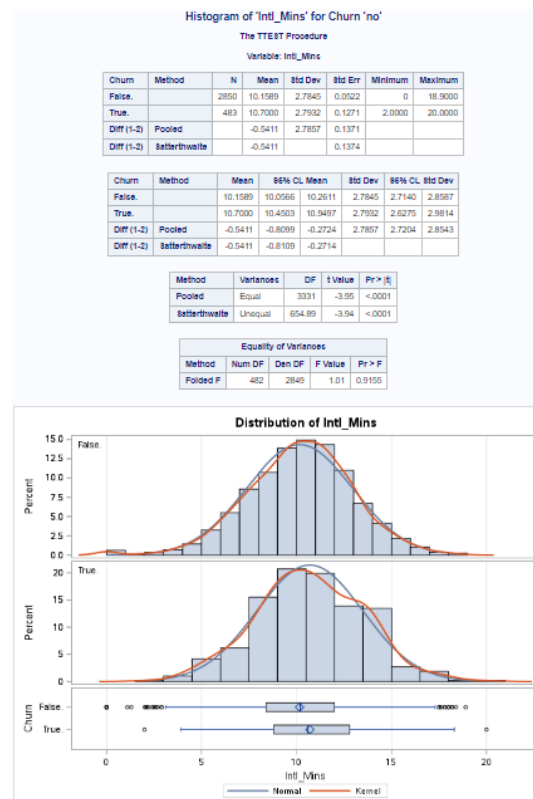
Intl_Charge & Class

- The mean Intl_Charge is slightly higher for clients who churned (yes) compared to those who did not churn (no). The mean Intl_Charge is 2.8895 for churned clients and 2.7434 for non-churned clients, with a difference of -0.1461.
 - The p-values from both the pooled and Satterthwaite t-tests are less than 0.05 (0.0001), indicating strong evidence to suggest a significant difference in the Intl_Charge between churned and non-churned clients.
 - The analysis implies that Intl_Charge is a significant predictor of churn status, as there is a statistically significant difference between the two groups in terms of Intl_Charge. Specifically, churned clients tend to have a slightly higher average Intl_Charge compared to non-churned clients.



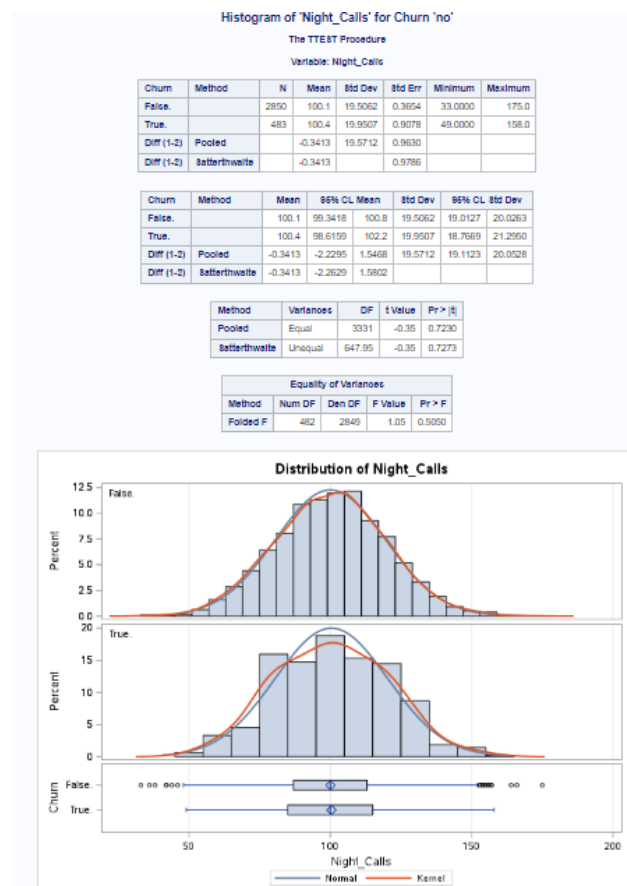
Intl Mins & Class

- The mean number of international minutes is slightly higher for clients who churned (yes) compared to those who did not churn (no), but the difference is statistically significant.
- The mean Intl_Mins for clients who churned is 10.70, compared to 10.16 for those who did not churn. The difference in means is -0.5411.
 - The p-values from both the pooled and Satterthwaite t-tests are less than 0.05 (0.0001), indicating a significant difference in the number of international minutes between churned and non-churned clients.
 - The analysis suggests that Intl_Mins does have a statistically significant relationship with churn status, implying that the number of international minutes could be a relevant factor in predicting customer churn.



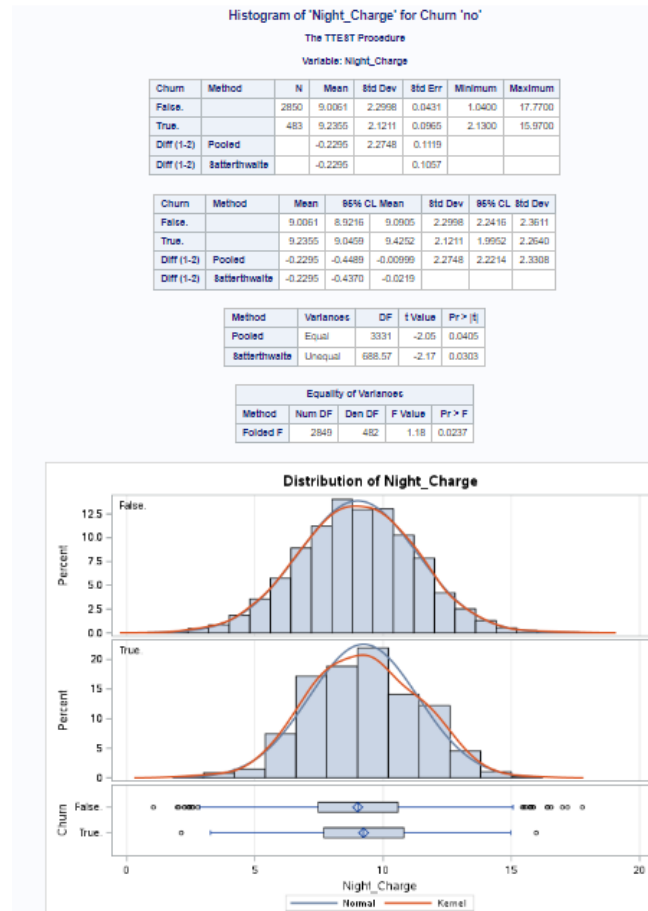
Night Calls & Class

- Mean Difference: The mean number of night calls is slightly higher for clients who churned (yes) compared to those who did not churn (no). However, the difference in means is very small, with a difference of -0.3413.
 - Statistical Significance:
 - The p-values from both the pooled t-test ($p = 0.7230$) and the Satterthwaite t-test ($p = 0.7273$) are greater than 0.05. This indicates that there is no strong evidence to suggest a significant difference in the number of night calls between clients who churned and those who did not.
 - The analysis implies that the number of night calls is not a strong predictor of churn status, as the observed differences in night calls between the two groups are not statistically significant.



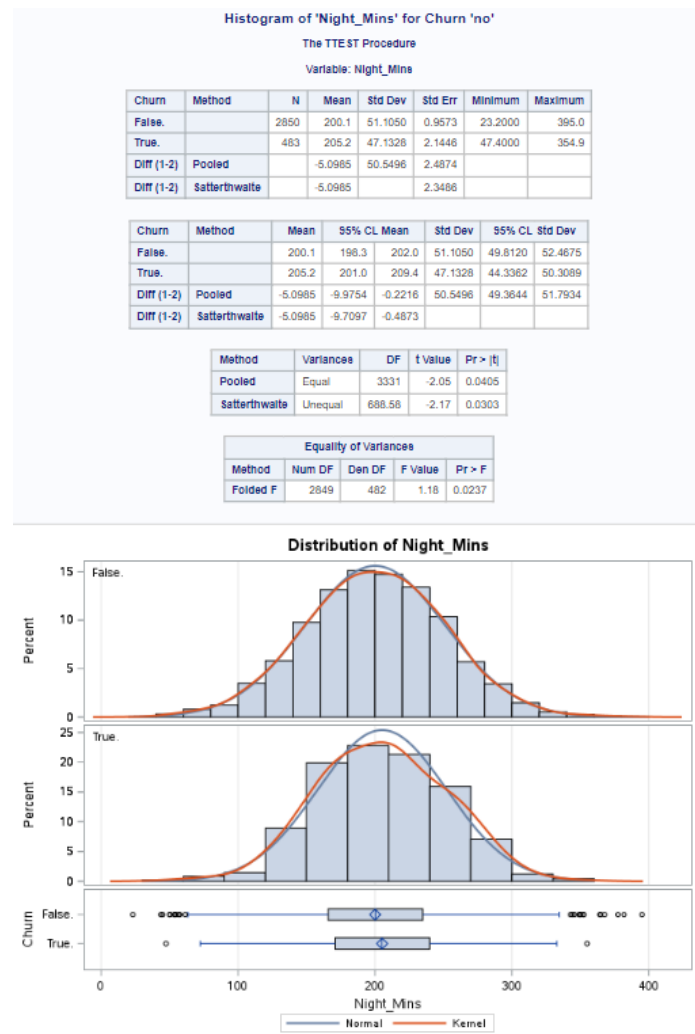
Night Charge & Class

- The mean night charge is slightly higher for clients who churned (9.24) compared to those who did not churn (9.01). The difference is small (-0.2295) but statistically significant at the 0.05 level.
 - The p-values from both the pooled (0.0405) and Satterthwaite (0.0303) t-tests are less than 0.05, indicating that there is evidence to suggest a significant difference in the night charge between the churned and non-churned clients.
 - Although the difference in night charge between the two groups is statistically significant, the magnitude of the difference is small, which may limit its practical significance as a predictor of churn status.



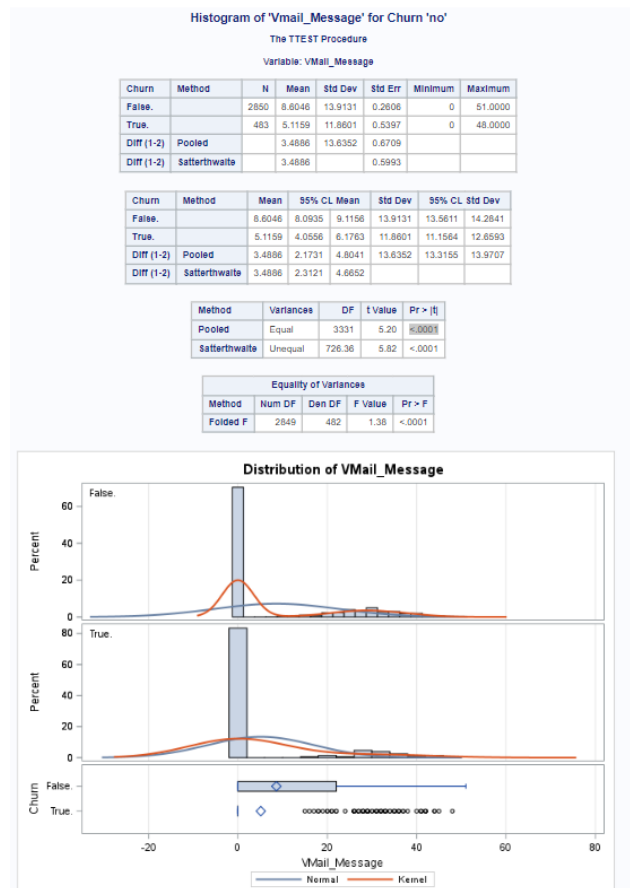
Night Mins & Class

- The mean number of night minutes is slightly higher for clients who churned (205.2) compared to those who did not churn (200.1), but the difference is small (-5.0985) and statistically significant.
 - The p-values from both the pooled ($p = 0.0405$) and Satterthwaite ($p = 0.0303$) t-tests are less than 0.05, indicating that there is strong evidence to suggest a significant difference in the number of night minutes between the churned and non-churned clients.
 - The analysis implies that the number of night minutes might be a predictor of churn status, as there is a statistically significant difference between the two groups in terms of night minutes. However, the practical significance of this difference should be considered, as the actual difference in means is relatively small.



Vmail Message & Class

- The mean number of voicemail messages is significantly higher for clients who did not churn (no) compared to those who did churn (yes), with a mean difference of 3.4886. This difference is statistically significant.
 - The p-values from both the pooled and Satterthwaite t-tests are less than 0.05, indicating strong evidence of a significant difference in the number of voicemail messages between churned and non-churned clients.
 - The analysis suggests that the number of voicemail messages is a significant predictor of churn status, as there is a substantial and statistically significant difference between the two groups in terms of voicemail messages.



(CATEGORICAL AND CLASS CORRELATIONS)

State & Class

- The p-value of 0.0023 indicates a significant association between the state and whether a client churned. This suggests that the likelihood of churn is not independent of the state a client is from.
- The Phi Coefficient and Cramer's V values are both 0.1578, which indicates a weak to moderate association between State and Churn. This implies that while there is a significant association, the strength of this relationship is not very strong.
- States with Higher Churn Rates:
 - West Virginia (WV): Has the highest proportion of churn with 14.49% of clients churning. This state shows a higher likelihood of churn compared to other states.
 - Mississippi (MS): This state has a relatively high churn rate at 21.54%, suggesting a greater likelihood of churn.
- States with Lower Churn Rates:
 - North Dakota (ND): Has the lowest proportion of churn with only 9.68% of clients churning. This state exhibits a lower likelihood of churn.
 - Nebraska (NE): Also has a lower churn rate at 8.20%, indicating a reduced likelihood of churn compared to other states.

Chi-Square Test for State and Churn Status
The PREG Procedure

Frequency
Expected
Percent
Row Pct
Col Pct

Table of State by Churn				
State	False	True	Total	
AK	49	3	52	
	44.464	7.5306		
	1.47	0.09	1.56	
	94.23	0.77		
	1.72	0.62		
AL	72	8	80	
	68.407	11.593		
	2.16	0.24	2.40	
	90.00	10.00		
	2.53	1.66		
AR	44	11	55	
	47.03	7.9703		
	1.32	0.33	1.65	
	80.00	20.00		
	1.54	2.28		
AZ	60	4	64	
	54.725	9.2745		
	1.80	0.12	1.92	
	93.75	6.25		
	2.11	0.63		
CA	25	9	34	
	29.073	4.9271		
	0.75	0.27	1.02	
	73.53	26.47		
	0.88	1.86		
CO	57	9	66	
	56.436	9.5644		
	1.71	0.27	1.98	
	86.36	13.64		
	2.00	1.86		
CT	62	12	74	
	63.276	10.724		
	1.86	0.36	2.22	
	83.78	16.22		
	2.18	2.40		
DC	49	5	54	
	46.175	7.8254		
	1.47	0.15	1.62	
	90.74	9.26		
	1.72	1.04		
DE	52	9	61	
	52.16	8.8388		
	1.56	0.27	1.83	
	85.25	14.75		
	1.82	1.86		
FL	55	8	63	
	53.87	9.1296		
	1.65	0.24	1.89	
	87.30	12.70		
	1.93	1.66		
GA	46	8	54	
	46.175	7.8254		
	1.38	0.24	1.62	
	85.19	14.81		
	1.61	1.66		
HI	50	3	53	
	45.32	7.6805		
	1.50	0.09	1.59	
	94.34	0.66		
	1.75	0.62		

IA	41	3	44	
	37.624	6.3762		
	1.23	0.09	1.32	
	93.18	6.82		
	1.44	0.62		
ID	64	9	73	
	62.421	10.579		
	1.92	0.27	2.19	
	87.67	12.33		
	2.05	1.86		
IL	53	5	58	
	49.055	8.945		
	1.59	0.15	1.74	
	91.38	8.62		
	1.86	1.04		
IN	62	9	71	
	60.111	10.289		
	1.76	0.27	2.13	
	87.52	12.68		
	2.18	1.86		
KS	57	13	70	
	59.856	10.144		
	1.71	0.39	2.10	
	81.42	18.57		
	2.00	2.69		
KY	51	8	59	
	50.465	8.535		
	1.53	0.24	1.77	
	86.64	13.36		
	1.79	1.66		
LA	47	4	51	
	43.809	7.1906		
	1.41	0.12	1.53	
	92.16	7.84		
	1.65	0.83		
MA	54	11	65	
	55.981	9.4194		
	1.62	0.33	1.95	
	83.08	16.92		
	1.89	2.28		
MD	53	17	70	
	59.856	10.144		
	1.59	0.51	2.10	
	75.71	24.29		
	1.86	3.52		
ME	49	13	62	
	53.015	8.9847		
	1.47	0.30	1.86	
	79.03	20.97		
	1.72	2.60		
MI	57	16	73	
	62.421	10.579		
	1.71	0.48	2.19	
	78.08	21.92		
	2.00	3.31		
MN	69	15	84	
	71.827	12.173		
	2.07	0.45	2.52	
	80.14	17.86		
	2.42	3.11		
MO	56	7	63	
	53.87	9.1296		
	1.68	0.21	1.89	
	89.89	11.11		
	1.96	1.45		
MS	51	14	65	
	55.981	9.4194		
	1.53	0.42	1.95	
	79.46	21.54		
	1.75	2.30		
MT	54	14	68	
	58.146	9.8542		
	1.62	0.42	2.04	
	75.41	20.59		
	1.89	2.30		
NC	57	11	68	
	58.146	9.8542		
	1.71	0.33	2.04	
	83.62	16.18		
	2.00	2.28		
ND	56	6	62	
	53.015	8.9847		
	1.68	0.18	1.86	
	90.32	9.68		
	1.96	1.24		
NE	56	5	61	
	52.16	8.8388		
	1.68	0.15	1.83	
	91.80	8.20		
	1.96	1.04		

NH	47	9	56	
	47.885	8.1152		
	1.41	0.27	1.68	
	83.93	16.07		
	1.65	1.86		
NJ	50	18	68	
	58.146	9.8542		
	1.50	0.34	2.04	
	73.53	26.47		
	1.75	3.73		
NM	56	6	62	
	53.015	8.9847		
	1.68	0.18	1.86	
	90.32	9.68		
	1.96	1.24		
NV	52	14	66	
	56.436	9.5644		
	1.56	0.42	1.98	
	78.79	21.21		
	1.82	2.90		
NY	68	15	83	
	70.972	12.028		
	2.04	0.45	2.49	
	81.93	18.07		
	2.39	3.11		
OH	68	10	78	
	66.897	11.303		
	2.04	0.30	2.34	
	87.18	12.82		
	2.39	2.07		
OK	52	9	61	
	52.16	8.8388		
	1.56	0.27	1.83	
	82.25	14.75		
	1.62	1.86		
OR	67	11	78	
	66.897	11.303		
	2.01	0.33	2.34	
	85.90	14.10		
	2.35	2.28		
PA	37	8	45	
	38.479	6.5212		
	1.11	0.24	1.35	
	82.22	17.78		
	1.30	1.66		
RI	59	6	65	
	55.981	9.4194		
	1.77	0.18	1.95	
	90.77	9.23		
	2.07	1.24		
SC	46	14	60	
	51.395	9.5940		
	1.38	0.42	1.80	
	76.67	23.33		
	1.61	2.90		
SD	52	8	60	
	51.395	9.5940		
	1.56	0.24	1.80	
	86.97	13.03		
	1.82	1.66		
TN	48	5	53	
	45.32	7.6805		
	1.44	0.15	1.59	
	90.11	9.43		
	1.68	1.04		
TX	54	18	72	
	61.396	10.434		
	1.62	0.54	2.16	
	75.00	25.00		
	1.89	3.73		
UT	62	10	72	
	61.396	10.434		
	1.86	0.30	2.16	
	86.11	13.89		
	2.18	2.07		
VA	72	5	77	
	65.362	11.198		
	1.56	0.15	2.31	
	93.51	6.49		
	2.33	1.04		
VT	65	8	73	
	62.421	10.579		
	1.95	0.24	2.19	
	89.04	10.96		
	2.28	1.66		
WA	52	14	66	
	56.436	9.5644		
	1.56	0.42	1.98	
	78.79	21.21		
	1.82	2.90		

WI	71	7	78	
	66.697	11.303		
	2.13	0.21	2.34	
	91.03	8.97		
	2.49	1.45		
WV	96	10	106	
	90.639	15.361		
	2.88	0.30	3.18	
	90.57	9.43		
	3.37	2.07		
WY	68	9	77	
	65.842	11.158		
	2.04	0.27	2.31	
	88.31	11.69		
	2.39	1.86		
Total	2850	483	3333	
	85.51	14.49	100.00	

Statistics for Table of State by Churn

Statistic	DF	Value	Prob
Chi-Square	50	83.0438	0.0023
Likelihood Ratio Chi-Square	50	83.1836	0.0022
Mantel-Haenszel Chi-Square	1	0.2017	0.6534
Phi Coefficient		0.1578	
Contingency Coefficient		0.1559	
Cramer's V		0.1578	

Sample Size = 3333

Area_Code & Class

- The p-value of 0.9151 indicates no significant association between Area_Code and whether a client churned. This suggests that the likelihood of churn is independent of the area code a client is from.
- The Phi Coefficient and Cramer's V values are both 0.0073, indicating a very weak association between Area_Code and Churn. This implies that while there may be some relationship, it is not strong.
- Area Codes with Higher Churn Rates:
 - Area Code 415: This area code has the highest number of churned clients with 236 out of 1655 clients, indicating a higher likelihood of churn compared to other area codes.
 - Area Code 510: This area code has 125 churned clients out of 840, suggesting a notable churn rate compared to the other area codes.
- Area Codes with Lower Churn Rates:
 - Area Code 408: This area code has 122 churned clients out of 1226, indicating a somewhat lower likelihood of churn compared to the higher churn rate area codes.
 - Area Code 838: This area code has 83 churned clients out of 331, showing a similar level of likelihood for churn compared to the other lower churn rate area codes.

Contingency Table for Area_Code and Churn Status

The FREQ Procedure

Frequency Percent Row Pct Col Pct	Table of Area_Code by Churn			
	Area_Code	Churn		Total
		False.	True.	
408		716	122	838
		21.48	3.66	25.14
		85.44	14.56	
		25.12	25.26	
415		1419	236	1655
		42.57	7.08	49.65
		85.74	14.26	
		49.79	48.66	
510		715	125	840
		21.45	3.75	25.20
		85.12	14.88	
		25.09	25.88	
Total		2850	483	3333
		85.51	14.49	100.00

Statistics for Table of Area_Code by Churn

Statistic	DF	Value	Prob
Chi-Square	2	0.1775	0.9151
Likelihood Ratio Chi-Square	2	0.1771	0.9153
Mantel-Haenszel Chi-Square	1	0.1270	0.7215
Phi Coefficient		0.0073	
Contingency Coefficient		0.0073	
Cramer's V		0.0073	

Sample Size = 3333

Phone & Class

- Chi-Square Test:
 - p-value: 0.4919
 - The p-value indicates that there is no significant association between Phone and Churn at the conventional significance level (e.g., 0.05). A p-value of 0.4919 suggests that the relationship between Phone and Churn is likely due to chance.
- Phi Coefficient/Cramer's V:
 - Phi Coefficient: 1.0000
 - Cramer's V: 1.0000
 - Both the Phi Coefficient and Cramer's V are at their maximum values. However, these values are not meaningful here due to the warning about expected counts. When 100% of cells have expected counts less than 5, these measures may not accurately represent the strength of association.

Statistics for Table of Phone by Churn			
Statistic	DF	Value	Prob
Chi-Square	3332	3333.0000	0.4919
Likelihood Ratio Chi-Square	3332	2758.2933	1.0000
Mantel-Haenszel Chi-Square	1	0.1055	0.7454
Phi Coefficient		1.0000	
Contingency Coefficient		0.7071	
Cramer's V		1.0000	
WARNING: 100% of the cells have expected counts less than 5. Chi-Square may not be a valid test.			

Intl_Plan & Class

- The p-value of <0.0001 indicates a highly significant association between the Intl_Plan status and whether a client churned. This suggests that the likelihood of churn is not independent of whether a client has an international plan.
- The Phi Coefficient and Cramer's V values are both 0.2599, which indicates a moderate association between Intl_Plan and Churn. This implies that while there is a significant relationship between having an international plan and the likelihood of churn, the strength of this relationship is moderate.
- Clients with Higher Churn Rates:
 - Clients with an international plan (Intl_Plan = yes):
 - 137 out of 323 clients with an international plan churned, which is a higher churn rate compared to those without an international plan.
- Clients with Lower Churn Rates:
 - Clients without an international plan (Intl_Plan = no):
 - 346 out of 3010 clients without an international plan churned, which is a lower churn rate compared to those with an international plan.

Contingency Table for Intl_Plan and Churn Status

The FREQ Procedure

Frequency Percent Row Pct Col Pct	Table of Intl_Plan by Churn		
	Intl_Plan	Churn	
		False.	True.
no		2664	346
		79.93	10.38
		88.50	11.50
		93.47	71.64
yes		186	137
		5.58	4.11
		57.59	42.41
		6.53	28.36
Total		2850	483
		85.51	14.49
			3333
			100.00

Statistics for Table of Intl_Plan by Churn

Statistic	DF	Value	Prob
Chi-Square	1	225.0541	<.0001
Likelihood Ratio Chi-Square	1	170.3998	<.0001
Continuity Adj. Chi-Square	1	222.5658	<.0001
Mantel-Haenszel Chi-Square	1	224.9888	<.0001
Phi Coefficient		0.2599	
Contingency Coefficient		0.2515	
Cramer's V		0.2599	

Fisher's Exact Test

Cell (1,1) Frequency (F)	2664
Left-sided Pr $\leq F$	1.0000
Right-sided Pr $\geq F$	<.0001
Table Probability (P)	<.0001
Two-sided Pr $\leq P$	<.0001

Sample Size = 3333

Vmail_Plan & Class

- The p-value of <0.0001 indicates a highly significant association between the Vmail_Plan status and whether a client churned. This suggests that the likelihood of churn is not independent of whether a client has a voicemail plan.
- The Phi Coefficient and Cramer's V values are both -0.1021 , which indicates a weak association between Vmail_Plan and Churn. This implies that while there is a significant relationship between having a voicemail plan and the likelihood of churn, the strength of this relationship is weak.
- Clients with Higher Churn Rates:
 - Clients without a voicemail plan (VMail_Plan = no):
 - Out of 2411 clients without a voicemail plan, 403 churned, resulting in a churn rate of approximately 16.73%.
- Clients with Lower Churn Rates:
 - Clients with a voicemail plan (VMail_Plan = yes):
 - Out of 922 clients with a voicemail plan, 80 churned, resulting in a churn rate of approximately 8.68%.
- This corrected analysis indicates that clients without a voicemail plan have a higher churn rate compared to those with a voicemail plan, suggesting that having a voicemail plan might be associated with lower churn

Contingency Table for Vmail_Plan and Churn Status

The FREQ Procedure

Frequency Percent Row Pct Col Pct	Table of VMail_Plan by Churn			
	VMail_Plan	Churn		Total
		False.	True.	
	no	2008 60.25 83.28 70.46	403 12.09 16.72 83.44	2411 72.34
	yes	842 25.26 91.32 29.54	80 2.40 8.68 16.56	922 27.66
	Total	2850 85.51	483 14.49	3333 100.00

Statistics for Table of VMail_Plan by Churn

Statistic	DF	Value	Prob
Chi-Square	1	34.7773	<.0001
Likelihood Ratio Chi-Square	1	37.9643	<.0001
Continuity Adj. Chi-Square	1	34.1317	<.0001
Mantel-Haenszel Chi-Square	1	34.7669	<.0001
Phi Coefficient		-0.1021	
Contingency Coefficient		0.1016	
Cramer's V		-0.1021	

Fisher's Exact Test

Cell (1,1) Frequency (F)	2008
Left-sided Pr $\leq F$	<.0001
Right-sided Pr $\geq F$	1.0000
Table Probability (P)	<.0001
Two-sided Pr $\leq P$	<.0001

Sample Size = 3333

(NUMERICAL CORRELATIONS)

Pearson Correlation Coefficients, N = 3333 Prob > r under H0: Rho=0															
	Account_Length	CustServ_Calls	Day_Calls	Day_Charge	Day_Mins	Eve_Calls	Eve_Charge	Eve_Mins	Intl_Calls	Intl_Charge	Intl_Mins	Night_Calls	Night_Charge	Night_Mins	VMail_Message
Account_Length	1.00000	-0.00380 0.8266	0.03847 0.0264	0.00621 0.7199	0.00622 0.7198	0.01926 0.2663	-0.00675 0.6971	-0.00676 0.6966	0.02066 0.2331	0.00955 0.5817	0.00951 0.5830	-0.01318 0.4470	-0.00896 0.6051	-0.00896 0.6053	-0.00463 0.7894
CustServ_Calls	-0.00380 0.8266	1.00000	-0.01894 0.2743	-0.01343 0.4384	-0.01342 0.4385	0.00242 0.8888	-0.01299 0.4535	-0.01298 0.4536	-0.01756 0.3108	-0.00967 0.5766	-0.00964 0.5780	-0.01280 0.4600	-0.00928 0.5924	-0.00929 0.5920	-0.01326 0.4440
Day_Calls	0.03847 0.0264	-0.01894 0.2743	1.00000	0.00675 0.6967	0.00675 0.6969	0.00646 0.7092	-0.02145 0.2157	-0.02145 0.2157	0.00457 0.7918	0.02167 0.2111	0.02156 0.2133	-0.01956 0.2590	0.02293 0.1857	0.02294 0.1855	-0.00955 0.5816
Day_Charge	0.00621 0.7199	-0.01343 0.4384	0.00675 0.6967	1.00000	1.00000 <.0001	0.01577 0.3628	0.00704 0.6847	0.00705 0.6841	0.00803 0.6430	-0.01009 0.5802	-0.01016 0.5578	0.02297 0.1849	0.00430 0.8040	0.00432 0.8029	0.00078 0.9643
Day_Mins	0.00622 0.7198	-0.01342 0.4385	0.00675 0.6969	1.00000 <.0001	1.00000	0.01577 0.3628	0.00703 0.6850	0.00704 0.6844	0.00803 0.6429	-0.01009 0.5803	-0.01015 0.5578	0.02297 0.1849	0.00430 0.8040	0.00432 0.8030	0.00078 0.9642
Eve_Calls	0.01926 0.2663	0.00242 0.8888	0.00646 0.7092	0.01577 0.3628	0.01577 0.3628	1.00000	-0.01142 0.5097	-0.01143 0.5095	0.01743 0.3143	0.00867 0.6167	0.00870 0.6155	0.00771 0.6564	-0.00206 0.9056	-0.00209 0.9039	-0.00586 0.7350
Eve_Charge	-0.00675 0.6971	-0.01299 0.4535	-0.02145 0.2157	0.00704 0.6847	0.00703 0.6850	-0.01142 0.5097	1.00000	1.00000 <.0001	0.00254 0.8834	-0.01107 0.5227	-0.01104 0.5239	0.00760 0.6611	-0.01260 0.4671	-0.01259 0.4674	0.01758 0.3103
Eve_Mins	-0.00676 0.6966	-0.01298 0.4536	-0.02145 0.2157	0.00705 0.6841	0.00704 0.6844	-0.01143 0.5095	1.00000	1.00000 <.0001	0.00254 0.8834	-0.01107 0.5230	-0.01103 0.5242	0.00759 0.6615	-0.01259 0.4674	-0.01258 0.4677	0.01756 0.3108
Intl_Calls	0.02066 0.2331	-0.01756 0.3108	0.00457 0.7918	0.00803 0.6430	0.00803 0.6429	0.01743 0.3143	0.00254 0.8834	0.00254 0.8834	1.00000	0.03237 0.0617	0.03230 0.0622	0.00030 0.9860	-0.01233 0.4767	-0.01235 0.4759	0.01396 0.4205
Intl_Charge	0.00955 0.5817	-0.00967 0.5766	0.02167 0.2111	-0.01009 0.5602	-0.01009 0.5603	0.00867 0.6167	-0.01107 0.5227	-0.01107 0.5230	0.03237 0.0617	1.00000	0.99999 <.0001	-0.01363 0.4315	-0.01519 0.3808	-0.01518 0.3810	0.00288 0.8678
Intl_Mins	0.00951 0.5830	-0.00964 0.5780	0.02156 0.2133	-0.01016 0.5578	-0.01015 0.5578	0.00870 0.6155	-0.01104 0.5239	-0.01103 0.5242	0.03230 0.0622	0.99999 <.0001	1.00000	-0.01360 0.4323	-0.01521 0.3799	-0.01521 0.3801	0.00286 0.8691
Night_Calls	-0.01318 0.4470	-0.01280 0.4600	-0.01956 0.2590	0.02297 0.1849	0.02297 0.1849	0.00771 0.6564	0.00760 0.6611	0.00759 0.6615	0.00030 0.9860	-0.01363 0.4315	-0.01360 0.4323	1.00000	0.01119 0.5185	0.01120 0.5179	0.00712 0.6810
Night_Charge	-0.00896 0.6051	-0.00928 0.5924	0.02293 0.1857	0.00430 0.8040	0.00430 0.8040	-0.00206 0.9056	-0.01260 0.4671	-0.01259 0.4674	-0.01233 0.4767	-0.01519 0.3808	-0.01521 0.3799	0.01119 0.5185	1.00000	1.00000 <.0001	0.00766 0.8583
Night_Mins	-0.00896 0.6053	-0.00929 0.5920	0.02294 0.1855	0.00432 0.8029	0.00432 0.8030	-0.00209 0.9039	-0.01259 0.4674	-0.01258 0.4677	-0.01235 0.4759	-0.01518 0.3810	-0.01521 0.3801	0.01120 0.5179	1.00000 <.0001	1.00000	0.00768 0.8576
VMail_Message	-0.00463 0.7894	-0.01326 0.4440	-0.00955 0.5816	0.00078 0.9643	0.00078 0.9642	-0.00586 0.7350	0.01758 0.3103	0.01756 0.3108	0.01396 0.4205	0.00288 0.8678	0.00286 0.8691	0.00712 0.6810	0.00766 0.8583	0.00768 0.8576	1.00000

(CATEGORICAL CORRELATIONS)

>Only showing Correlation = Yes

Area_Code & Intl_Plan

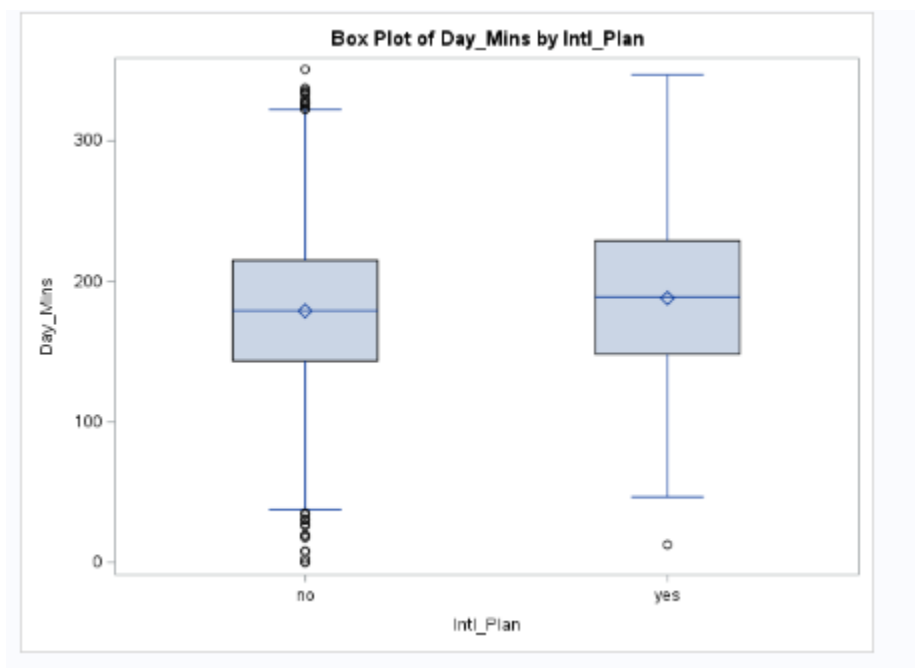
- p-value: The p-value is 0.0189, which is less than 0.05. This indicates a significant association between Area_Code and Intl_Plan. This suggests that the likelihood of having an international plan is not independent of the area code.
- Phi Coefficient and Cramer's V: Both values are 0.0488, which indicates a very weak association between Area_Code and Intl_Plan. This implies that while there is a statistically significant relationship, the strength of the association is quite weak.
- Frequency Distribution:
 - Clients with Intl_Plan = yes:
 - The percentage of clients with an international plan varies by area code:
 - Area Code 408: 8.47%
 - Area Code 415: 9.06%
 - Area Code 510: 12.14%
 - Clients with Intl_Plan = no:
 - The percentage of clients without an international plan:
 - Area Code 408: 91.53%
 - Area Code 415: 90.94%
 - Area Code 510: 87.86%

(CATEGORICAL and NUMERICAL CORRELATIONS)

>Only showing Correlation = Yes

Day_Mins & Intl_Plan

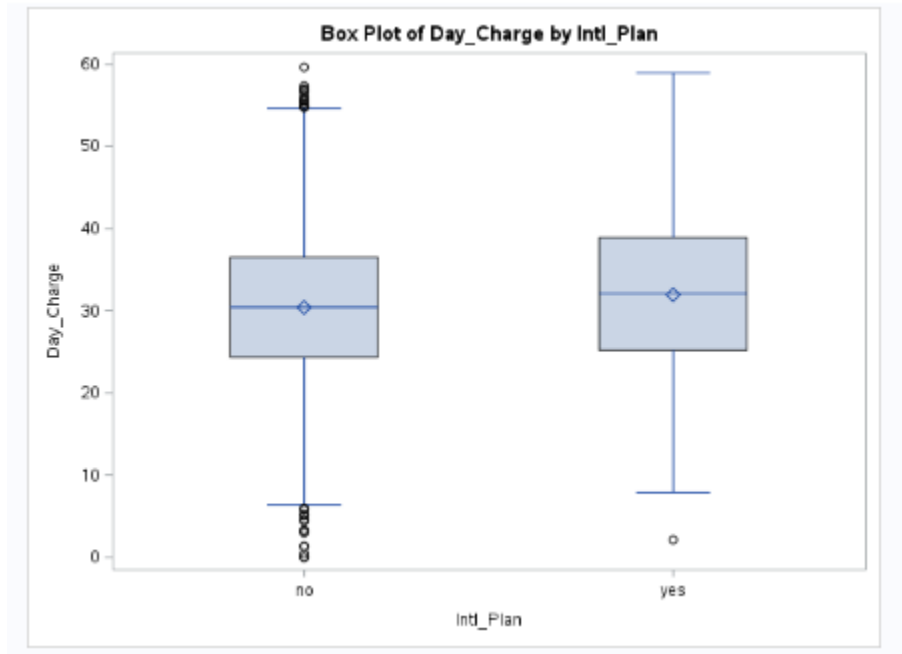
- The mean number of Day_Mins is slightly higher for clients with an international plan (Intl_Plan = yes) compared to those without (Intl_Plan = no), but the difference is relatively small and statistically significant.
 - The F-value is 8.15 with a p-value of 0.0043, indicating a statistically significant difference in Day_Mins between clients with and without an international plan.
 - The R-squared value is 0.002440, suggesting that the model explains a very small proportion of the variance in Day_Mins.
 - The p-value from the ANOVA is less than 0.05, which means there is sufficient evidence to suggest a significant difference in Day_Mins based on the Intl_Plan status.
 - The analysis indicates that while there is a significant difference in Day_Mins related to having an international plan, the practical significance might be limited due to the small R-squared value.



Day_Charge & Intl_Plan

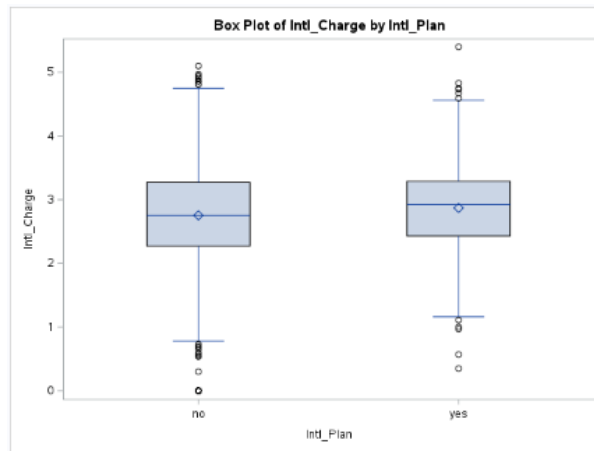
- The mean Day_Charge is slightly higher for clients with an international plan (Intl_Plan = yes) compared to those without (Intl_Plan = no), and this difference is statistically significant.
 - The F-value is 8.15 with a p-value of 0.0043, indicating a statistically significant difference in Day_Charge between clients with and without an international plan.
 - The R-squared value is 0.002440, which shows that the model explains a very small proportion of the variance in Day_Charge.

- The p-value from the ANOVA is less than 0.05, providing strong evidence of a significant difference in Day_Charge based on the Intl_Plan status.
- Despite the statistical significance, the small R-squared value suggests that the practical impact of Intl_Plan on Day_Charge might be limited, as the model does not explain much of the variability in the dependent variable.



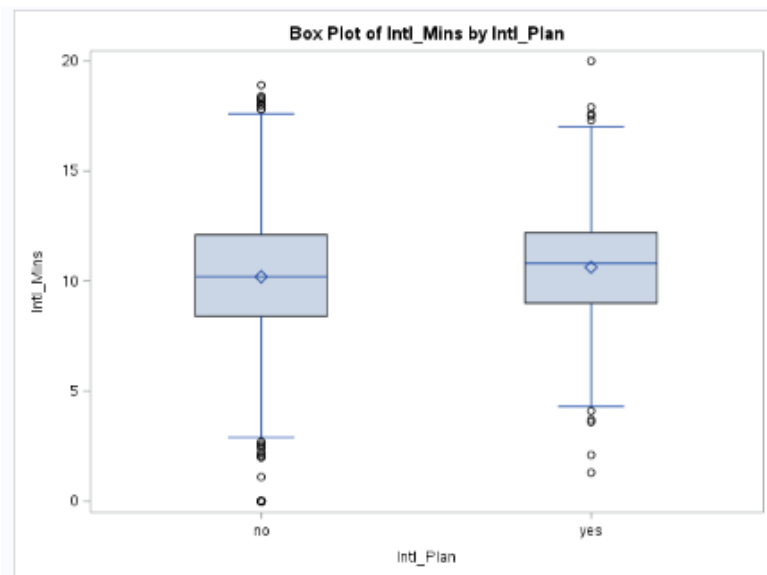
Intl_Charge & Intl_Plan

- The mean Intl_Charge is slightly higher for clients with an international plan (Intl_Plan = yes) compared to those without (Intl_Plan = no), and this difference is statistically significant.
 - The F-value is 7.00 with a p-value of 0.0082, indicating a statistically significant difference in Intl_Charge between clients with and without an international plan.
 - The R-squared value is 0.002096, which shows that the model explains a very small proportion of the variance in Intl_Charge.
 - The p-value from the ANOVA is less than 0.05, providing strong evidence of a significant difference in Intl_Charge based on the Intl_Plan status.
 - Despite the statistical significance, the small R-squared value suggests that the practical impact of Intl_Plan on Intl_Charge might be limited, as the model does not explain much of the variability in the dependent variable.



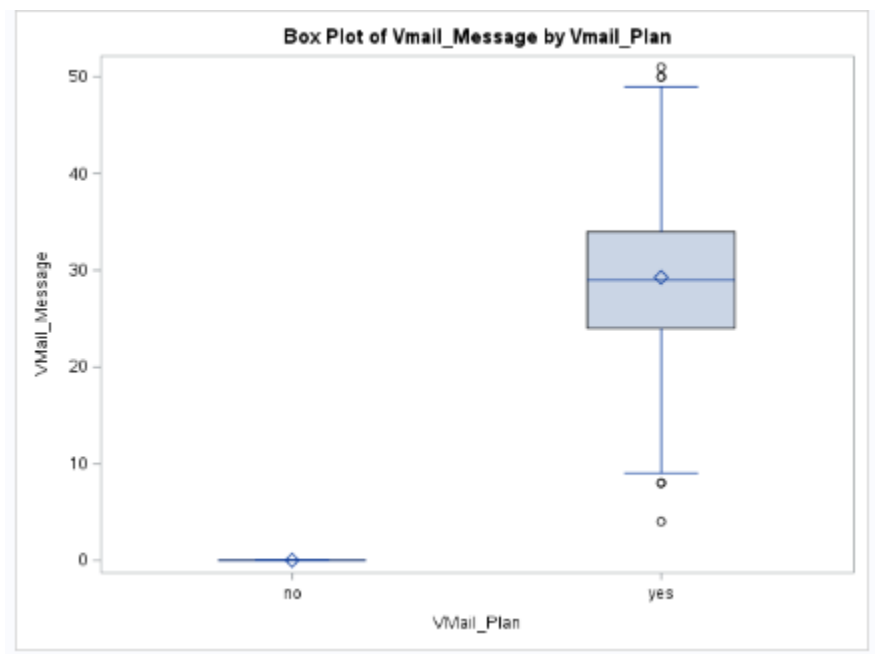
Intl_Mins & Intl_Plan

- The mean Intl_Mins is slightly higher for clients with an international plan (Intl_Plan = yes) compared to those without (Intl_Plan = no), and this difference is statistically significant.
 - The F-value is 7.02 with a p-value of 0.0081, indicating a statistically significant difference in Intl_Mins between clients with and without an international plan.
 - The R-squared value is 0.002104, which shows that the model explains a very small proportion of the variance in Intl_Mins.
 - The p-value from the ANOVA is less than 0.05, providing strong evidence of a significant difference in Intl_Mins based on the Intl_Plan status.
 - Despite the statistical significance, the small R-squared value suggests that the practical impact of Intl_Plan on Intl_Mins might be limited, as the model does not explain much of the variability in the dependent variable.



Vmail_Message & Vmail_Plan

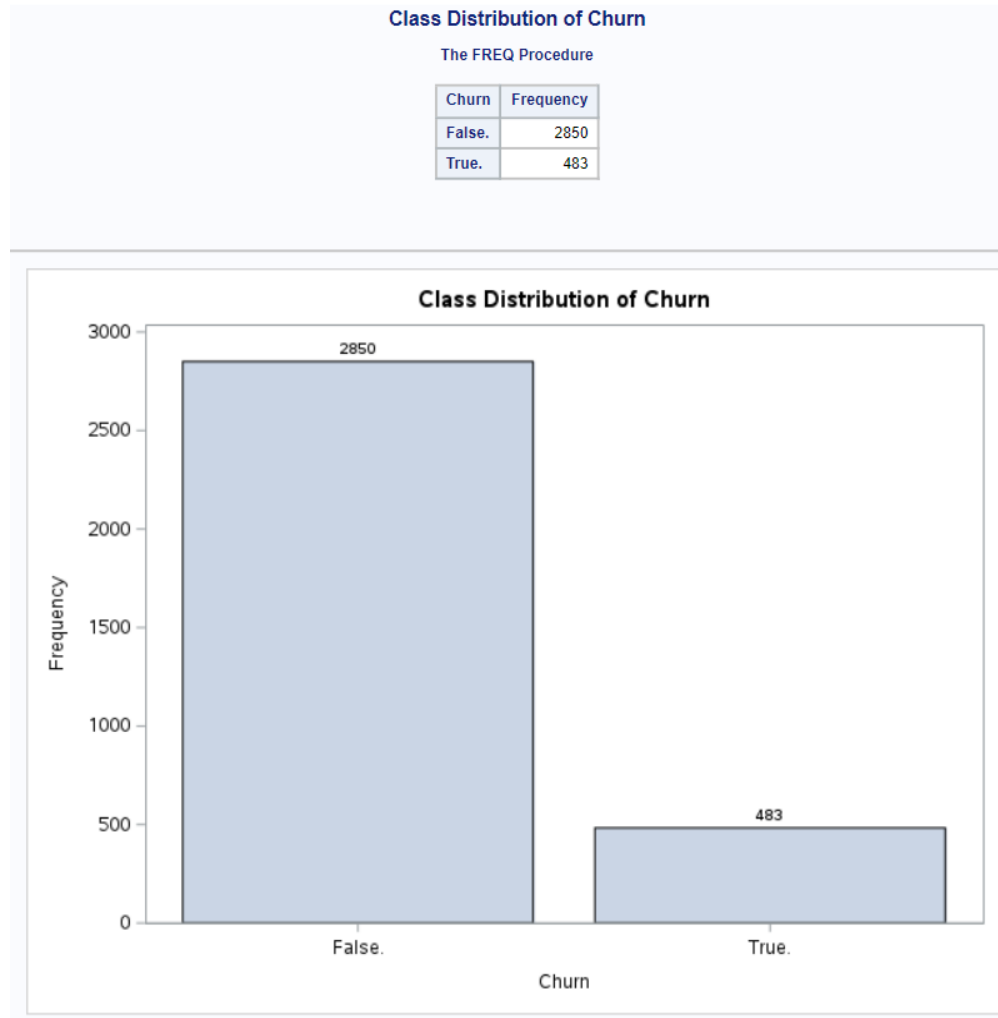
- The mean number of Vmail_Message is significantly different between clients with a voicemail plan (Vmail_Plan = yes) and those without (Vmail_Plan = no), and this difference is highly statistically significant.
 - The F-value is 36186.7 with a p-value of <0.0001 , indicating a highly significant difference in Vmail_Message between clients with and without a voicemail plan.
 - The R-squared value is 0.9157, which shows that the model explains a substantial proportion of the variance in Vmail_Message.
 - The p-value from the ANOVA is less than 0.05, providing very strong evidence of a significant difference in Vmail_Message based on the Vmail_Plan status.
 - The large R-squared value suggests that the presence of a voicemail plan has a considerable impact on the number of voicemail messages received, explaining a significant portion of the variability in Vmail_Message.



Which attributes do you think can be eliminated or included in the analysis? This can be a subjective decision or an objective decision based on statistical method

Selection Strategy included in Naïve Bayes and Decision Tree model

Determine whether the dataset has an imbalanced class distribution (same proportion of records of different types or not) and do you need to balance the dataset.



- The class distribution of the variable Churn (which represents whether clients have churned or not) is imbalanced:
 - The True class (churned) constitutes 483 records, which is approximately 14.6% of the total dataset.
 - The False class (not churned) constitutes 2850 records, which is approximately 85.4% of the total dataset.
- The True class is significantly smaller compared to the False class, indicating an imbalance. The percentage of the True class is under the 20-30% threshold, which highlights the imbalance in the dataset.

Balancing the Class Attribute:

In the predictive modeling efforts, I have addressed the class imbalance in the dataset, specifically focusing on the Churn attribute, which indicates whether clients have churned or not.

Reason for Balancing:

- The dataset exhibits a notable class imbalance, with the True class (churned) representing only 14.6% of the data, while the False class (not churned) accounts for 85.4%. This imbalance could result in a model that is biased towards predicting the majority class (False), which may impair its ability to accurately identify clients who have actually churned (True).

Impact of Balancing:

- Balancing the dataset is crucial as it allows the model to more effectively learn and predict the minority class (True). This approach enhances key performance metrics such as precision, recall, and F1-score for the churned class. By addressing the imbalance, the model becomes more robust and reliable in making accurate predictions, particularly improving its ability to identify clients who are at risk of churning.