BUAN 6337 Predictive Analytics using SAS Homework Assignment 1 Group 11

Prithvi Shivashankar
Soham Bhalerao
Harkirat kaur
Arjitkumar Sureshkumar
Anjana Sebastian
Sai Pratheek Banda

Q1. What is the distribution of gender, vehicle size, and vehicle class?

The SAS System

The FREQ Procedure

Gender	Frequency	Percent	Cumulative Frequency	Cumulative Percent
F	4658	51.00	4658	51.00
M	4476	49.00	9134	100.00

Vehicle_Size	Frequency	Percent	Cumulative Frequency	Cumulative Percent
Large	946	10.36	946	10.36
Medsize	6424	70.33	7370	80.69
Small	1764	19.31	9134	100.00

Vehicle_Class	Frequency	Percent	Cumulative Frequency	Cumulative Percent
Four-Door Car	4621	50.59	4621	50.59
Luxury Car	163	1.78	4784	52.38
Luxury SUV	184	2.01	4968	54.39
SUV	1796	19.66	6764	74.05
Sports Car	484	5.30	7248	79.35
Two-Door Car	1886	20.65	9134	100.00

Distribution of Gender: There is not much variation in terms of gender distribution for the given data.

Distribution of Vehicle Size: Medium-sized vehicles constitute 70% of the data.

Distribution of Vehicle Class: The vehicles are distributed widely. The top 3 dominant vehicle classes in the given data are Four-door car, two-door car and SUV respectively.

Q2. What is the average customer life time value of each level of gender, vehicle size, and vehicle class?

		The	e MEAN	S Proc	edure			
		Analysis Varia	ıble : Cu	stome	r_Lifetime	_Value		
Gender	Vehicle_Size	Vehicle_Class	N Obs	N	Mean	Std Dev	Minimum	Maximum
F	Large	Four-Door Car	249	249	6596.15	4753.13	2111.99	27564.74
		Luxury Car	7	7	13152.99	5183.70	7373.23	21435.88
		Luxury SUV	7	7	288 47.15	21238.57	7449.88	60556.19
		SUV	91	91	9441.19	7539.97	3853.47	51337.91
		Sports Car	30	30	11161.95	6318.59	4082.00	35537.85
		Two-Door Car	121	121	6637.54	51 18.68	2336.29	27528.31
	Medsize	Four-Door Car	1659	1659	6748.67	5503.89	1904.00	41787.90
		Luxury Car	55	55	14437.68	7992.32	6698.97	51428.25
		Luxury SUV	61	61	17888.00	13980.05	6991.25	73225.96
		SUV	660	660	10572.28	8322.12	3371.53	58753.88
		Sports Car	181	181	11542.64	90 10.80	3595.31	40132.01
		Two-Door Car	614	614	7028.99	5454.13	2147.68	38887.90
	Small	Four-Door Car	498	498	6820.34	5637.46	2004.35	38470.30
		Luxury Car	13	13	18922.65	7945.75	7255.14	25807.08
		Luxury SUV	15	15	16917.91	9972.78	6383.61	48770.98
		SUV	171	171	10438.55	7879.10	3451.10	510 16.07
		Sports Car	31	31	9801.49	6596.88	3884.86	26900.27
		Two-Door Car	195	195	6828.67	5781.18	1898.68	35186.26

M	Large	Four-Door Car	226	226	6075.99	4665.63	2052.95	35944.71
		Luxury Car	9	9	13478.59	6256.67	7126.60	22837.14
		Luxury SUV	11	11	16487.56	15022.67	6874.18	58207.13
		SUV	76	76	10147.42	9132.98	3123.08	48611.87
		Sports Car	19	19	9030.71	9463.37	3954.34	40636.67
		Two-Door Car	100	100	5853.44	3610.24	1940.98	22563.62
	Medsize	Four-Door Car	1578	1578	6804.89	4956.39	1994.77	32467.66
		Luxury Car	51	51	16551.64	12813.08	6191.40	74228.52
		Luxury SUV	64	64	15858.53	10308.01	6423.74	68025.75
		SUV	648	648	10387.80	7642.75	3099.54	49423.80
		Sports Car	185	185	10205.47	8339.33	3074.11	67907.27
		Two-Door Car	668	668	6535.13	5070.82	1898.01	35444.31
	Small	Four-Door Car	411	411	6361.32	4373.62	2030.78	29232.69
		Luxury Car	28	28	24361.32	19868.45	5886.22	83325.38
		Luxury SUV	26	26	16168.61	11739.64	6871.77	50568.26
		SUV	150	150	10883.60	7169.98	2864.82	44795.47
		Sports Car	38	38	10948.38	8764.80	3515.48	39561.08
		Two-Door Car	188	188	6277.78	4489.36	1918.12	29577.28

For both males and females,

Luxury SUV has higher Lifetime value under Large and Medium vehicle size and Luxury car has higher lifetime value under small vehicle size on an average.

Q3. Do Large cars have a higher lifetime value than medsize cars. Do a ttest and report on your findings.

				ine	HES	Proc	cedure				
			Vari	able: C	ustom	er_Li	fetime	_Va	lue		
	Vehicle	Size	N	Mean	Std [Dev	Std Er	r M	inimum	Maximu	ım
	Large		946	7545.0	662	625.4	215.4	1	1941.0	6055	5.2
	Medsize	•	6424	8050.7	683	33.1	85.254	0	1898.0	7422	3.5
	Diff (1-2)			-505.7	680	06.8	237.)			
Vehi	cle_Size	Meth	od	Me	an !	95% C	L Mea	n	Std Dev	95% CL	Std D
Larg	Large			7545	0.0 7	122.3	796	7.7	6625.4	6339.7	6938
Meds	size			8050	.7 78	883.5	821	7.8	6833.1 6806.8		6953
Diff (1-2)	Pool	ed	-505	.7 -9	970.3	-40.99	.9917			6918
Diff (1-2)	Satte	erthwai	te -505	.7 -9	960.2	-51.10	590			
		Meth	od	Var	iance	s	DF t	Valu	ue Pr>	t	
		Pool	ed	Equ	al	7	368	-2.	13 0.032	29	
		Satte	erthwa	ite Une	qual	125	59.7	-2.	18 0.029	12	
				Equa	ality o	f Vari	iances				
		M	ethod	Num I	OF D	en DF	F Va	lue	Pr > F		
		Fo	olded F	64	23	945	5	1.06	0.2183		

Average vehicle size for medium cars $-\mu 1$ Average vehicle size for large cars $-\mu 2$

Initial test for equality of variance:

- H0: μ1= μ2
- H1: μ1≠ μ2

From The equality of variances test we conclude that the variances are equal as the p-value is 0.2183

We can hence conclude at 5% level of significance that the null hypothesis that the variances are equal cannot be rejected.

H0: Large cars do not have higher lifetime value than medium size car (Medium size car>=Large size car)

H1: Large cars have higher lifetime value than medium size car (Medium size car<Large size car)

The T-value for the case of equality of variance (as tested above) is -2.13 and since this is a left tail test, the t-critical value is -1.64. As the calculated t-value is lesser than the t-critical value, we can reject the null hypothesis in favor of the alternative. Hence, we can conclude that large cars have higher lifetime value than medium sized cars.

Q4. Is there a significant difference between men and women in customer life time value?

				The	SAS	Sys	tem					
				The 1	TTEST	Proc	edure	е				
		Va	ariab	ole: Cu	ustome	r_Lif	etime	e_Val	ue			
G	ende	r N	Me	an S	td Dev	Std	Err	Minin	num	Max	imun	n
F		4658	8096	6.6	6956.1	10	1.9	18	98.7	73	3226.	0
M	l	4476	7909	9.6	6780.7	10	1.4	18	98.0	83	3325.	4
D	iff (1-	2)	187	7.1	6870.7	14	13.8					
Gender	Me	thod		Mean	95%	CL N	lean	Sto	l Dev	95%	CL :	Std De
F			8	096.6	789	6.8	8296.	4 6	956.1	681	7.6	7100.
M			7	909.6	771	0.9	8108.	3 6	780.7	664	3.1	6924
Diff (1-2)	Po	oled		187.1	-94.84	77	468.	9 6	870.7	677	2.5	6971.
Diff (1-2)	Sa	tterthwait	е	187.1	-94.70	43	468.	8				
		Method		Vari	iances		DF 1	t Valu	e Pr	> t		
		Pooled		Equ	al	9	132	1.3	0 0.	1934		
		Satterthw	aite	Une	qual	913	0.1	1.3	0 0.	1932		
				Faua	lity of	Varia	ances					
		Method	1 1	Num [alue	Pr >	F		
		Folded	F	46	57	4475		1.05	0.084	47		

Average Customer value for men – μ 1 Average Customer value for women – μ 2

Initial test for equality of variance:

• H0: μ1= μ2

H1: μ1≠ μ2

From The equality of variances test we conclude that the variances are equal as the pvalue is 0.0847. We can hence conclude at 5% level of significance that the null hypothesis that the variances are equal cannot be rejected.

H0: Customer life time value of men and women is not different.

H1: Customer life time value of men and women is different.

The p value for equal variances is 0.1934, so we do not reject the null hypothesis at 5% significance level and we don't have enough evidence to claim that customer lifetime value for men and women is different.

Q5. Use ANOVA to test whether there is difference in customer lifetime value across different sales channels. Which sales channel generates the highest lifetime value?



HO: Customer lifetime value across different sales channels is not different.

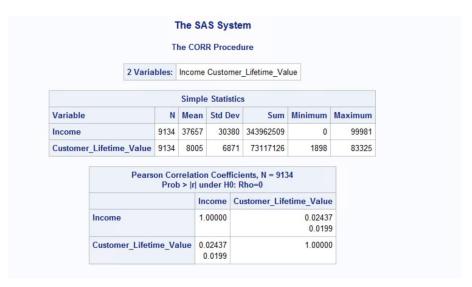
H1: Atleast customer lifetime value across one channel is different.

In The above ANOVA test, it is observed that the p value is 0.4503, Hence we do not reject the null hypothesis. Therefore, the data explains that customer lifetime value across different channels is not different.

As per the findings, it is observed that as a total (N * Mean), Agent sales channel generates the highest customer lifetime value.

Q6. What demographic factors (education, income, marital_status) affect customer lifetime value?

• INCOME and customer life time value:



We calculated the correlation to find how income affects customer life time value. Correlation of 0.0243 suggests that there's almost no correlation between income and CLV and since, the p-value is 0.0199, we conclude that the correlation is not zero. Hence the data explains Income isn't significantly affecting customer life time value.

• Education and customer life time value:

requency			Table	of clv	by	Educatio	n		
Percent Row Pct				1	Edu	cation			
Col Pct	clv	Bachelor	College	Docto	or	High Sch	ool or	Master	Total
·	1	692 7.58 30.30 25.18	680 7.44 29.77 25.36	1.1 4.4 29.5	11		629 6.89 27.54 23.99	182 1.99 7.97 24.56	2284 25.01
	2	1408 15.41 30.83 51.24	1339 14.66 29.32 49.94	1.7 3.5 46.7	75		1308 14.32 28.64 49.89	352 3.85 7.71 47.50	4567 50.00
	3	648 7.09 28.38 23.58	7.25 29.00 24.69	0.8 3.5 23.6	55		685 7.50 30.00 26.13	207 2.27 9.07 27.94	2283 24.99
	Total	2748 30.09	2681 29.35	34 3.7	-		2622 28.71	741 8.11	9134 100.00
	Stat	Statistics	for Table		by DF	Educatio	n Prot		
		Square			8	13.0316	2 080		
		elihood Rati	io Chi Sau	iaro	8	12.8480	100000000000000000000000000000000000000		
		itel-Haensz			1				
	718.0.014	Coefficient		aare		0.0378	0.025	,	
	1 111	Coemcient				100000000000000000000000000000000000000			
	Con	tingency Co	nefficient			0.0377			

Since Chi square is a procedure for testing if two categorical variables are related or not, we divided Customer Lifetime Value into three levels:

if Customer Lifetime Value < 3994 then clv=1

if 8963 < Customer Lifetime Value < 3994 then clv= 2

if Customer Lifetime Value >= 8963 then clv=3

Chi square test:

H0: Education and Customer life time value are independent. No relationship exists.

H1: Education and Customer life time value are dependent

As you can see, the p value of the chi square test is 0.1108, hence we fail to reject the null hypothesis. Therefore, there is no relationship between education and customer life time value.

• Marital status and customer life time value:

Percent Row Pct Col Pct Divorced Married Single Tota
Col Pct clv Divorced Married Single Total 1 353 1271 660 228 3.86 13.92 7.23 25.0 15.46 55.65 28.90 26.75 2 654 2693 1220 456
3.86 13.92 7.23 25.0 15.46 55.65 28.90 25.79 23.99 26.75 2 654 2693 1220 456
15.46 55.65 28.90 25.79 23.99 26.75 2 654 2693 1220 456
25.79 23.99 26.75 2 654 2693 1220 456
2 654 2693 1220 456
7.16 29.48 13.36 50.0
14.32 58.97 26.71
47.77 50.83 49.45
3 362 1334 587 228
3.96 14.60 6.43 24.9
15.86 58.43 25.71
26.44 25.18 23.79
Total 1369 5298 2467 913
14.99 58.00 27.01 100.0

Chi square test:

HO: Marital status and Customer life time value are independent. No relationship exists.

H1: Marital status and Customer life time value are dependent.

Contingency Coefficient

Phi Coefficient

Cramer's V

The p value here is 0.0340, hence we reject the null hypothesis at 5% significance level. Therefore, we conclude that marital status does affect customer lifetime value.

0.0338

0.0337

0.0239

Q7. Is there a relationship between renew_offer_type and response (use Chi-sq test)? Which offer type generates the highest response rate?

Frequency	Table of Renew_C	Offer	Туре	by R	esponse
Percent Row Pct			R	espo	nse
Col Pct	Renew_Offer_Type	е	No	Ye	s Total
	Offer1	1	3158 34.57 34.17 40.35		0 41.08 3
	Offer2		2242 24.55 76.62 28.65		9 32.03
	Offer3		1402 15.35 97.91 17.91	3 0.3 2.0 2.2	3 15.68 9
	Offer4	10	1024 11.21 00.00 13.08	0.0	0
	Total		7826 35.68	130	
Statistics for	or Table of Renew_C	Offer DF		by R	esponse
Chi-Squa	re	3	548.1	1645	<.0001
	re d Ratio Chi-Square	3			<.0001
Likelihoo				1675	
Likelihoo	d Ratio Chi-Square nenszel Chi-Square	3	751.4 242.3	1675	<.0001
Likelihoo Mantel-Ha Phi Coeffi	d Ratio Chi-Square nenszel Chi-Square	3	751.4 242.3 0.2	1675 3027	<.0001

H0: Renew Offer Type and Response rate are independent.

H1: Renew Offer Type and Response rate are dependent.

The p value here is less than 0.0001, hence we reject the null hypothesis at 5% significance level. Therefore, we conclude that renew offer type does affect response rate.

Offer type 1 generates the highest response rate.

Q8. Do different renew_offer_types have different lifetime values? Which offer type is the best?



The SAS System The MEANS Procedure Analysis Variable: Customer Lifetime Value Renew_Offer_Type N Obs N Mean Std Dev Minimum Maximum Offer1 3752 3752 8707.09 7336.98 83325.38 1898.01 Offer2 2926 2926 7396.75 6446.15 1994.77 61134.68 1432 7997.89 Offer3 1432 6669.59 1898.68 61850.19 Offer4 1024 1024 7179.95 6286.01 2121.31 56675.94

HO: Customer lifetime value across different renew offer types is not different.

H1: Customer lifetime value across different renew offer types is different.

The p value here is less than 0.0001, hence we reject the null hypothesis at 5% significance level. Therefore, we conclude that customer lifetime value across different renew offer types is different.

Offer 1 type is the best.

Q9. Is the effectiveness of renew_offer_type different across different states with respect to lifetime value?

				Cla	ss Level	Informati	n		
С	lass		Le	vels	Values	3			
S	tate			5	Arizona	California	Nevada (Oregon Was	hingto
R	enew_Offer	_Туре		4	Offer1	Offer2 Offe	3 Offer4		
		N	umb	er o	f Observ	ations Re	nd 9134		
		N	umb	er o	f Observ	ations Use	d 9134	ı	
								_	
				Т	he SAS	System			
				The	e ANOVA	Procedu	е		
		epend	dent	Var	iable: Cu	istomer_L	fetime_	Value	
So	ource	[OF .	Sum	of Squa	res Mea	Square	e F Value	Pr>
Mo	odel		19	4	07988168	33.7 214	30614.93	3 4.58	<.000
Er	ror	91	14	42	27090837	243 4686	0965.245	5	
Co	rrected Tot	al 91	33	43	31170718	927			
	R-Square	Coeff	Var	Ro	oot MSF	Custome	Lifetim	ne Value M	lean
	0.009462	85.5			845.507	Customo		_	.940
So	urce			DF	Anova	SS Mea	Square	e F Value	Pr >

ANOVA test:

Ho: Mean of Renew offer type with respect to lifetime value is equal across different states.

H1: Mean of at least one Renew offer type across different states with respect to lifetime value is different.

The p value here is <0.0001, hence we reject the null hypothesis. Therefore, we conclude that Renew offer type with respect to lifetime value is different across different states at 5% significance level.

Q10. What other interesting insights that are useful to the company in terms of action can be obtained from the data? Write any 3 and indicate which type of analysis is appropriate.

Insight 1: According to marketing theory, it is always wise to invest more in customers with high customer lifetime value. From the data, we can find out that premium renew type has more CLV than basic type. If we can prove that there is a significant difference between Premium Renew type and Basic

Renew Type with respect to CLV, which can be done using ANOVA, then we can recommend upgrading customers who are likely to churn from basic to premium for free because the cost is minuscule for the company and we can retain those customers.

Insight 2: People who work in urban areas tend to stay in suburban areas because of high cost of housing in the urban areas and factors such as lesser noise and stress. Such people tend to drive from the suburban areas to the urban areas and there is a higher probability that such people could encounter accidents on the road. We could first check the average claim amount for the three different location types and then run an ANOVA to check if there is a significant difference between the three groups and based on that, we would be able to see whether the people from suburban locations have a higher claim amount and then in such a case, we could charge a higher premium.

			The SAS			
	Ar	alysis	Variable : To	tal_Claim_An	nount	
Location_Code	N Obs	N	Mean	Std Dev	Minimum	Maximum
Rural	1773	1773	109.9050952	76.8463187	0.0990070	562.0885870
Suburban	5779	5779	562.1598701	275.1666070	292.8000000	2893.24
Urban	1582	1582	329.5723289	124.1755121	156.9212470	1065.05

Based on the above table, we can see that people from suburban areas have much higher claim amounts on average and then by running an ANOVA, we can check whether they are significantly different and then based on that, we can decide to charge a higher premium for the suburban people.

Insight 3: If we plot a trend line between customer lifetime value and Months since the last claim, we can find that there is inverse relationship i.e. as the number of months since the last claim increases CLV tends to decrease over time. If we can prove this by running a regression and get that the relationship is significant then we can get a segment of the customer, we can focus on.