

Αναφορά εργασίας εξαμήνου στο μάθημα  
Εφαρμοσμένα πληροφοριακά συστήματα II

Εργαλείο:SAS

<u>Ονόματα:</u>	<u>ΑΜ:</u>
Δημητριάδης Μιλτιάδης(*)	5288
Πατρώνη Σωτηρία	5399
Τσαρνάς Βασίλειος	5444

## Ερώτημα 1:

Για την ολοκλήρωση του καθαρισμού των δεδομένων του dataset χρειάστηκε αρχικά να ενωθούν τα δύο αρχεία δεδομένων Campaign1.csv και Campaign2.csv σε ένα dataset και έπειτα να ξεκινήσει ο καθαρισμός.

- Εισαγωγή στοιχείων του αρχείου Campaign1.csv στο dataset CAMP.campaign1

```
1 /*Import arxeiou Campaign1 */
2 /*Observations=1700, Variables=18*/
3
4 data CAMP.campaign1 ;
5 infile "&path/Campaign1.csv" dlm=';' dsd firstobs=2;
6 input id_ targetD:dollar10.2 GiftCnt36 GiftCntAll GiftCntCard36 GiftCntCardAll
7       GiftAvgLast:dollar10.2 GiftAvg36:dollar10.2 GiftAvgAll:dollar10.2
8       GiftAvgCard36:dollar10.2 GiftTimeLast GiftTimeFirst PromCnt12
9       PromCnt36 PromCntAll PromCntCard12 PromCntCard36 PromCntCardAll;
10 run;
11
12 proc print data=CAMP.campaign1;
13 run;
14
15 proc contents data=CAMP.campaign1;
16 run;
```

- Εισαγωγή στοιχείων του αρχείου Campaign2.csv στο dataset CAMP.campaign2

```
18 /*Import arxeiou Campaign2 */
19 /*Observations=1700, Variables=11*/
20 data CAMP.campaign2 ;
21 infile "&path/Campaign2.csv" dlm=';' dsd firstobs=2;
22 input id_ StatusCat96NK $ StatusCatStarAll DemCluster DemAge DeMGender $
23       DemHomeOwner $ DemIncomeGroup DemMedHomeValue:dollar10.2
24       DemPctVeterans TargetB;
25 run;
26
27 proc print data=CAMP.campaign2;
28 run;
29
30 proc contents data=CAMP.campaign2;
31 run;
```

- Ταξινόμηση των δύο dataset με βάση το id για την ένωση τους

```

28 /*Sorting */
29 proc sort data=CAMP.campaign1
30 out=CAMP.campaign1_sorted;
31 by id_;
32 run;
33
34 proc print data=CAMP.campaign1_sorted;
35 run;
36
37 proc sort data=CAMP.campaign2
38 out=CAMP.campaign2_sorted;
39 by id_;
40 run;
41
42 proc print data=CAMP.campaign2_sorted;
43 run;
44

```

- Ένωση των δύο dataset

```

51 /*Merging*/
52 data CAMP.campaign12 ;
53 merge CAMP.campaign1_sorted (in=a) CAMP.campaign2_sorted (in=b) ;
54 by id_;
55 if a=1 and b=1;
56 run;
57
58 proc print data=CAMP.campaign12;
59 run;
60
61 proc contents data=CAMP.campaign12;
62 run;
63
64 proc freq data=CAMP.campaign12;
65 tables _NUMERIC_ /missing;
66 tables _CHAR_ /missing;
67 run;
68

```

- Καθαρισμός του dataset μετά την ένωση

```

69 /*Cleaning the dataset*/
70 data CAMP.cleaned(drop= DemPctVeterans DemAge TargetB);
71 set CAMP.campaign12;
72
73 if DeMGender = 'Manan' then DeMGender = 'M';
74 if DeMGender = 'Man' then DeMGender = 'M';
75 if DemHomeOwner = 'Home' then DemHomeOwner = 'H';
76 if DemHomeOwner = 'NoHome' then DemHomeOwner = 'NH';
77 if DeMGender = 'Woman' then DeMGender = 'F';
78 if DeMGender = 'U' then delete;
79 if GiftAvgCard36 = '.' then delete;
80 if DemIncomeGroup = '.' then delete;
81 if TargetB = '0' then delete;
82 if DemMedHomeValue = '0' then delete;
83 if GiftAvgLast = '0' then delete;
84 if PromCntCard12 = '0' then delete;
85 if GiftCnt36 = '0' then delete;
86 if GiftCntCard36 = '0' then delete;
87 if GiftCntCardAll = '0' then delete;
88 if GiftAvg36 = '0' then delete;
89
90 run;

92 proc print data=CAMP.cleaned;
93 run;
94
95 proc contents data=CAMP.cleaned;
96 run;
97
98 proc freq data=CAMP.cleaned;
99 tables _NUMERIC_ /missing;
100 tables _CHAR_ /missing;
101 run;

```

Αρχικά, στην μεταβλητή DeMGender διορθώνονται όλες οι άτυπες τιμές και εκχωρούνται μόνο M (male) και F (female), έπειτα γίνεται το ίδιο και για τη μεταβλητή DemHomeOwner και τέλος διαγράφονται όλες οι λανθασμένες τιμές και τα missing values. Η τιμή U της μεταβλητής DeMGender και οι τιμές 0 των μεταβλητών DemMedHomeValue, GiftAvgLast, PromCntCard12, GiftCnt36, GiftCntCard36, GiftCntCardAll, GiftAvg36 είναι λανθασμένες τιμές, τα missing values βρίσκονται μόνο στις μεταβλητές GiftAvgCard36 και DemIncomeGroup και τέλος για την μεταβλητή TargetB η τιμή 0 είναι άχρηστη αφού δεν μας αφορούν οι πελάτες που δε θα πάρουν κάποιο δώρο. Επίσης, γίνεται διαγραφή των στηλών:

1. DemPctVeterans, αφού από την εκφώνηση είναι γνωστό ότι όλες οι μεταβλητές που υπολογίζουν κάποιο ποσοστό είναι άχρηστες και λανθασμένες,
2. DemAge, αφού το ποσοστό των λανθασμένων και των τιμών που απουσιάζουν είναι αρκετά μεγάλο για να αφαιρεθούν όλες αυτές οι εγγραφές οπότε γίνεται διαγραφή όλης της στήλης και
3. TargetB, αφού έχουν κρατηθεί όλες οι εγγραφές που είναι χρησιμες με βάση τη μεταβλητή αυτή (δηλαδή όλες οι εγγραφές των πελατών που θα λάβουν κάποιο δώρο).

Αποτελέσματα εκτέλεσης της διαδικασίας proc contents:

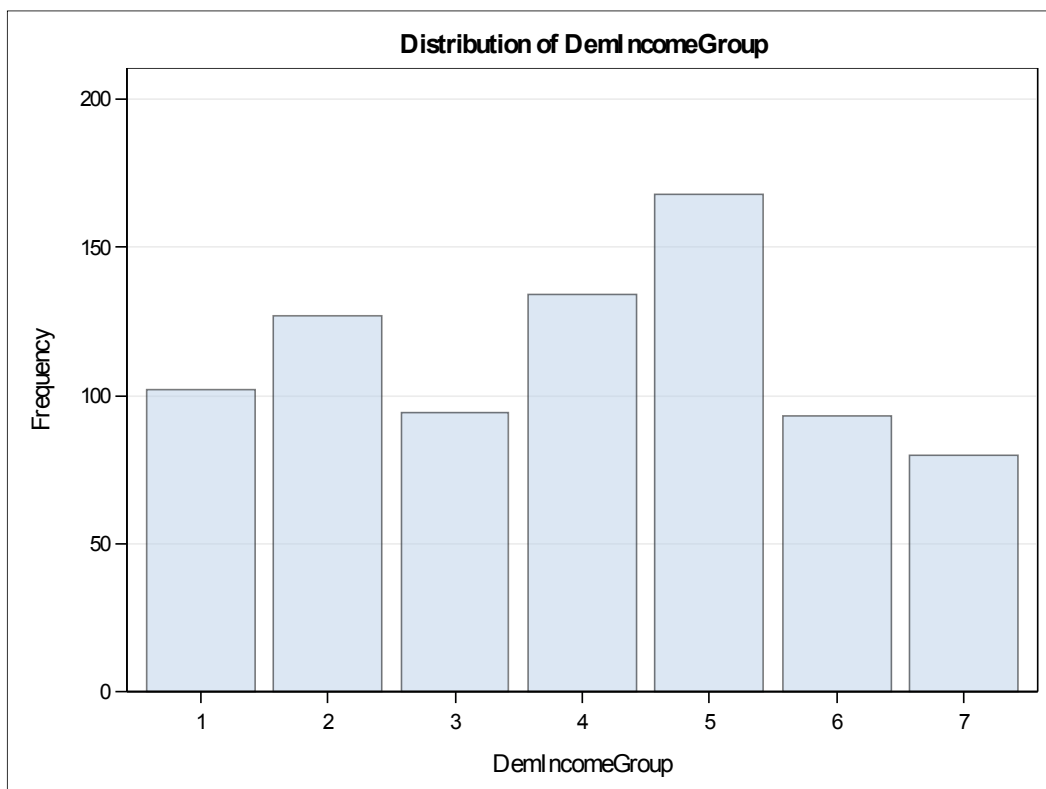
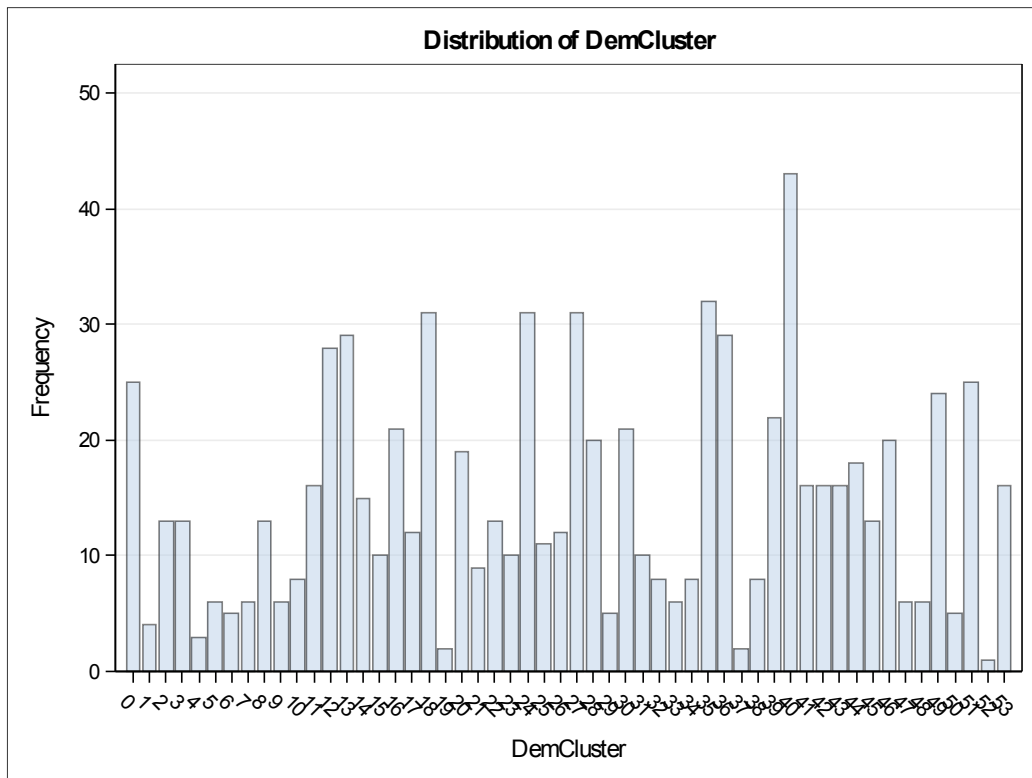
<b>Data Set Name</b>	CAMP.CLEANED	<b>Observations</b>	798
<b>Member Type</b>	DATA	<b>Variables</b>	25
<b>Engine</b>	V9	<b>Indexes</b>	0
<b>Created</b>	07/06/2015 09:59:51	<b>Observation Length</b>	200
<b>Last Modified</b>	07/06/2015 09:59:51	<b>Deleted Observations</b>	0
<b>Protection</b>		<b>Compressed</b>	NO
<b>Data Set Type</b>		<b>Sorted</b>	NO
<b>Label</b>			
<b>Data Representation</b>	SOLARIS_X86_64, LINUX_X86_64, ALPHA_TRU64, LINUX_IA64		
<b>Encoding</b>	utf-8 Unicode (UTF-8)		

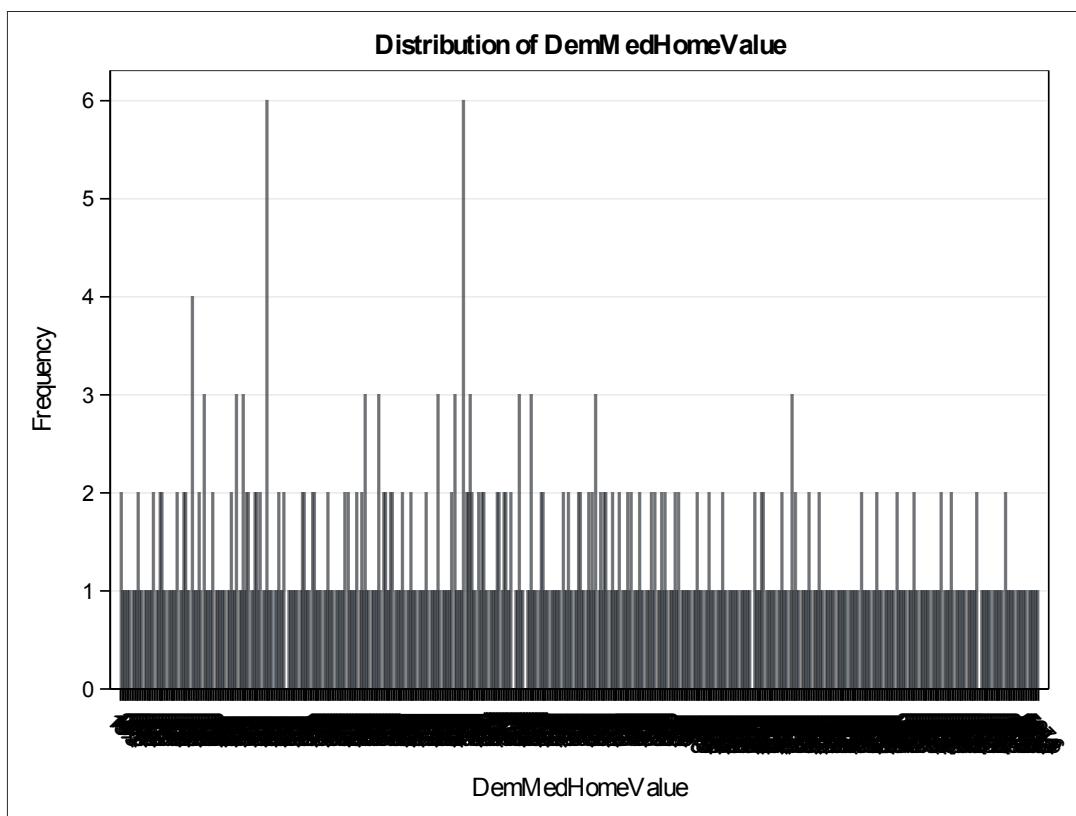
Engine/Host Dependent Information	
<b>Data Set Page Size</b>	65536
<b>Number of Data Set Pages</b>	3
<b>First Data Page</b>	1
<b>Max Obs per Page</b>	327
<b>Obs in First Data Page</b>	302
<b>Number of Data Set Repairs</b>	0
<b>Filename</b>	/folders/myfolders/myproject/cleaned.sas7bdat
<b>Release Created</b>	9.0401M2
<b>Host Created</b>	Linux
<b>Inode Number</b>	11554

Engine/Host Dependent Information	
Access Permission	rw-rw-rw-
Owner Name	root
File Size (bytes)	262144

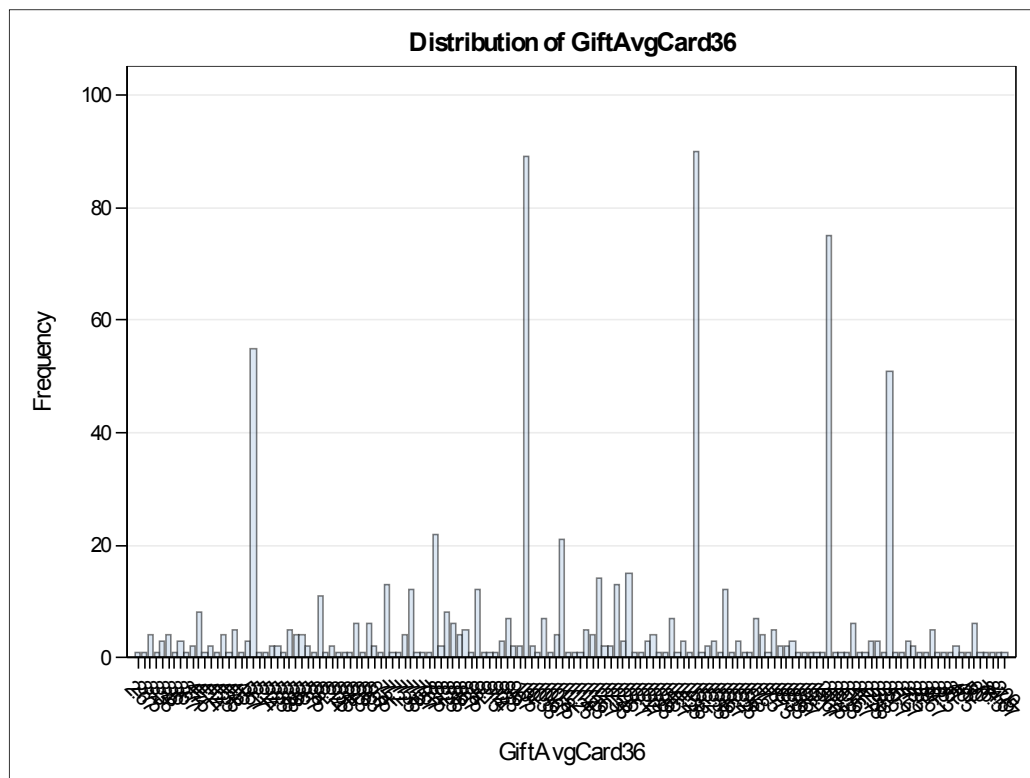
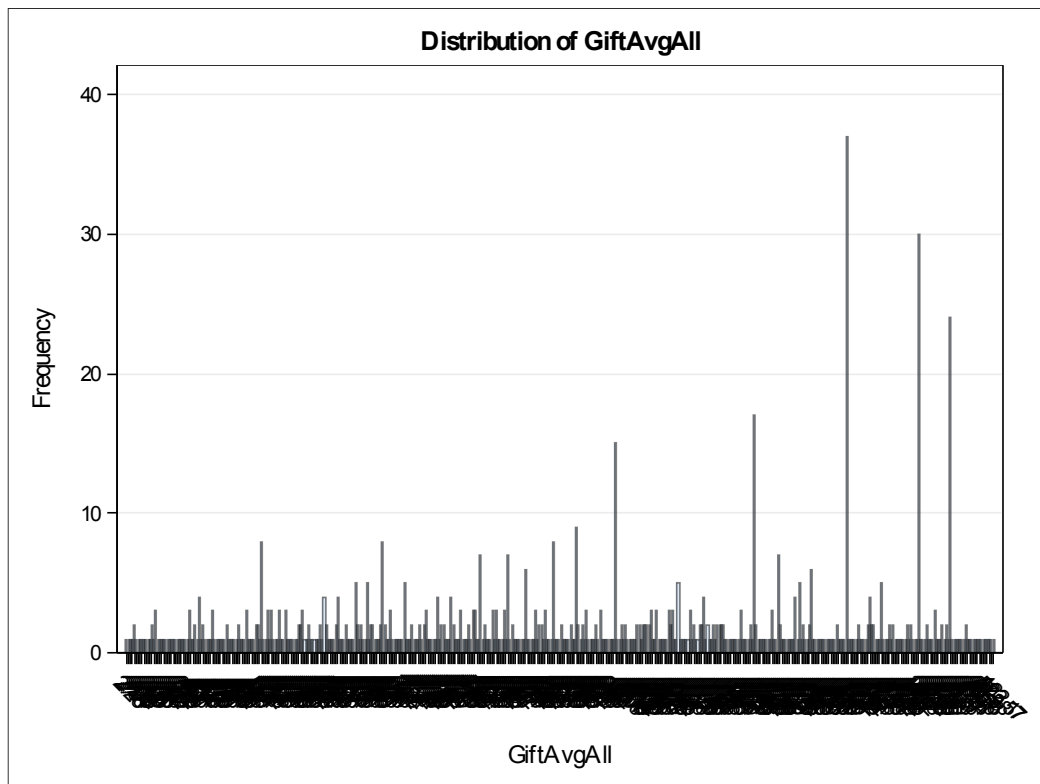
Alphabetic List of Variables and Attributes			
#	Variable	Type	Len
22	DeMGender	Char	8
21	DemCluster	Num	8
23	DemHomeOwner	Char	8
24	DemIncomeGroup	Num	8
25	DemMedHomeValue	Num	8
8	GiftAvg36	Num	8
9	GiftAvgAll	Num	8
10	GiftAvgCard36	Num	8
7	GiftAvgLast	Num	8
3	GiftCnt36	Num	8
4	GiftCntAll	Num	8
5	GiftCntCard36	Num	8
6	GiftCntCardAll	Num	8
12	GiftTimeFirst	Num	8
11	GiftTimeLast	Num	8
13	PromCnt12	Num	8
14	PromCnt36	Num	8
15	PromCntAll	Num	8
16	PromCntCard12	Num	8
17	PromCntCard36	Num	8
18	PromCntCardAll	Num	8
19	StatusCat96NK	Char	8
20	StatusCatStarAll	Num	8
1	id_	Num	8
2	targetD	Num	8

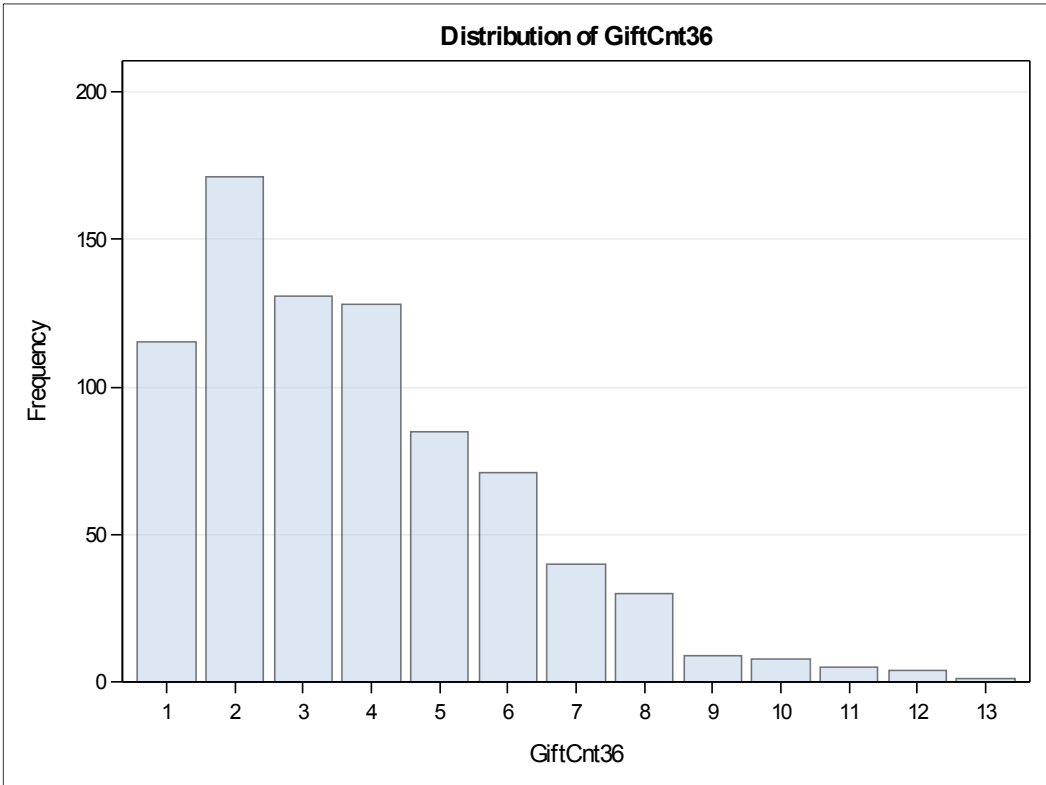
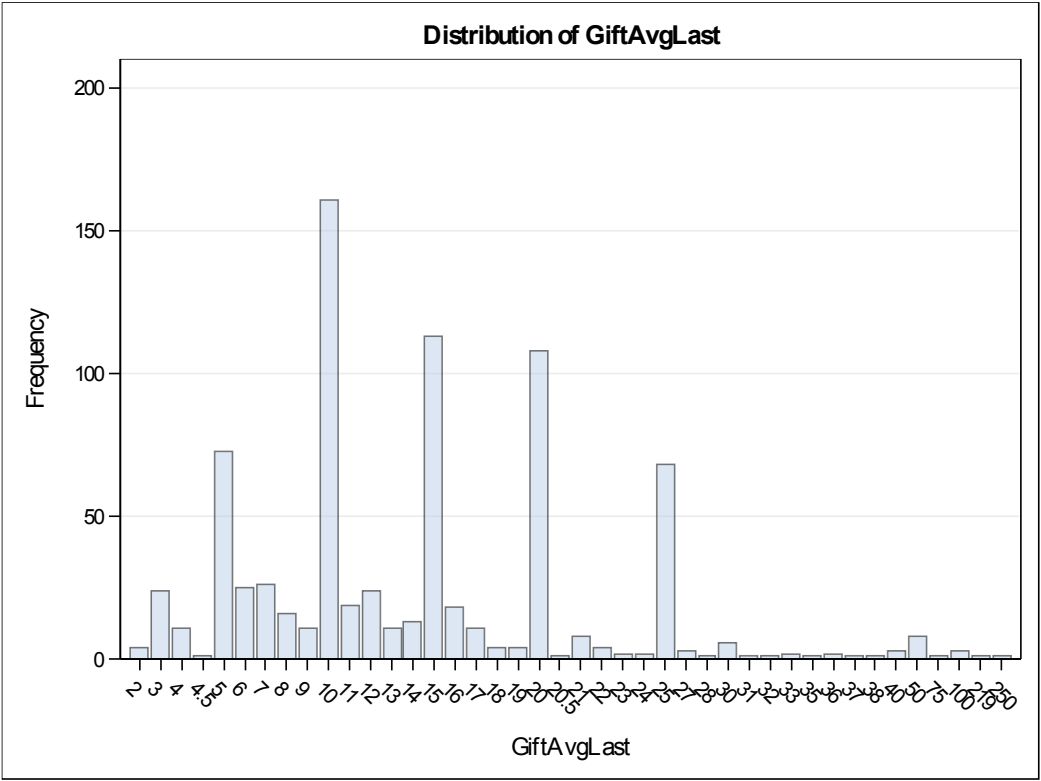
## Ιστογράμματα:

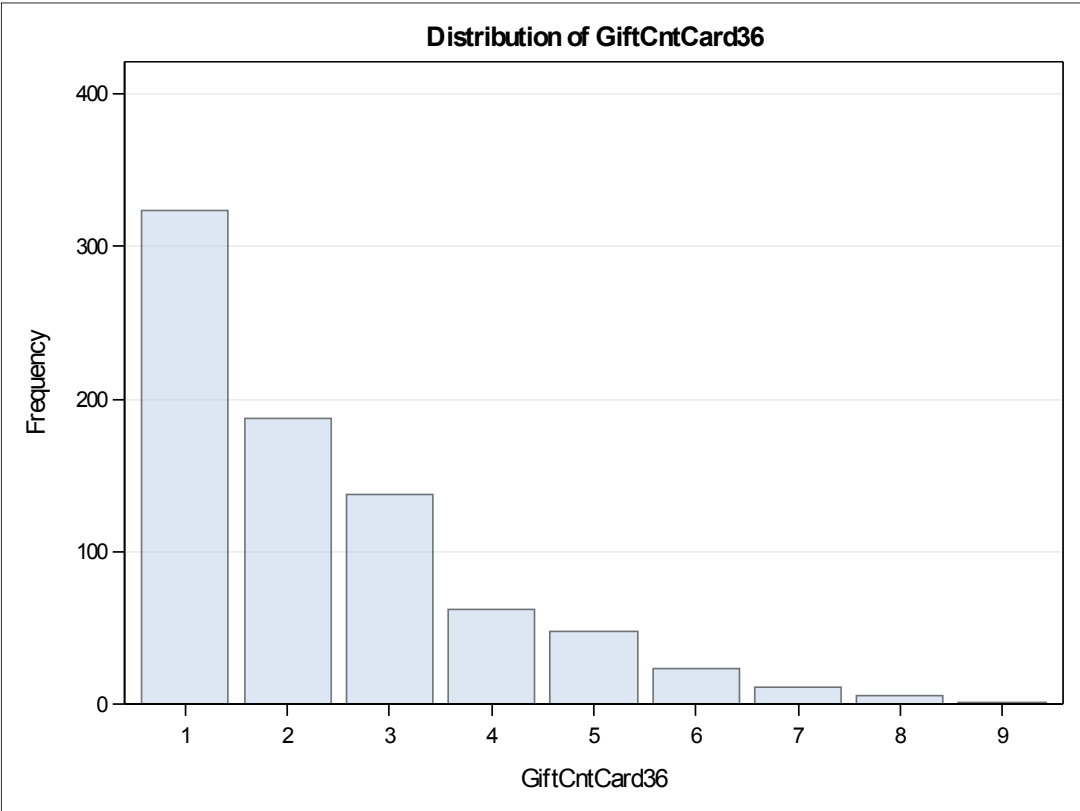
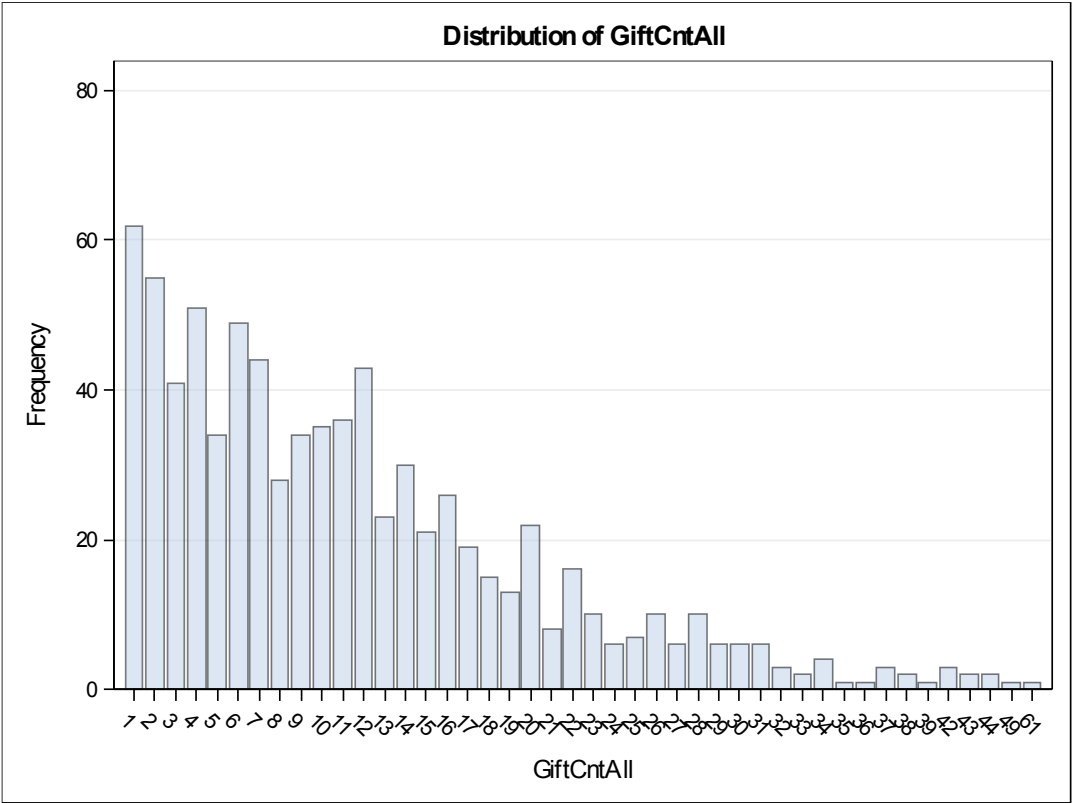


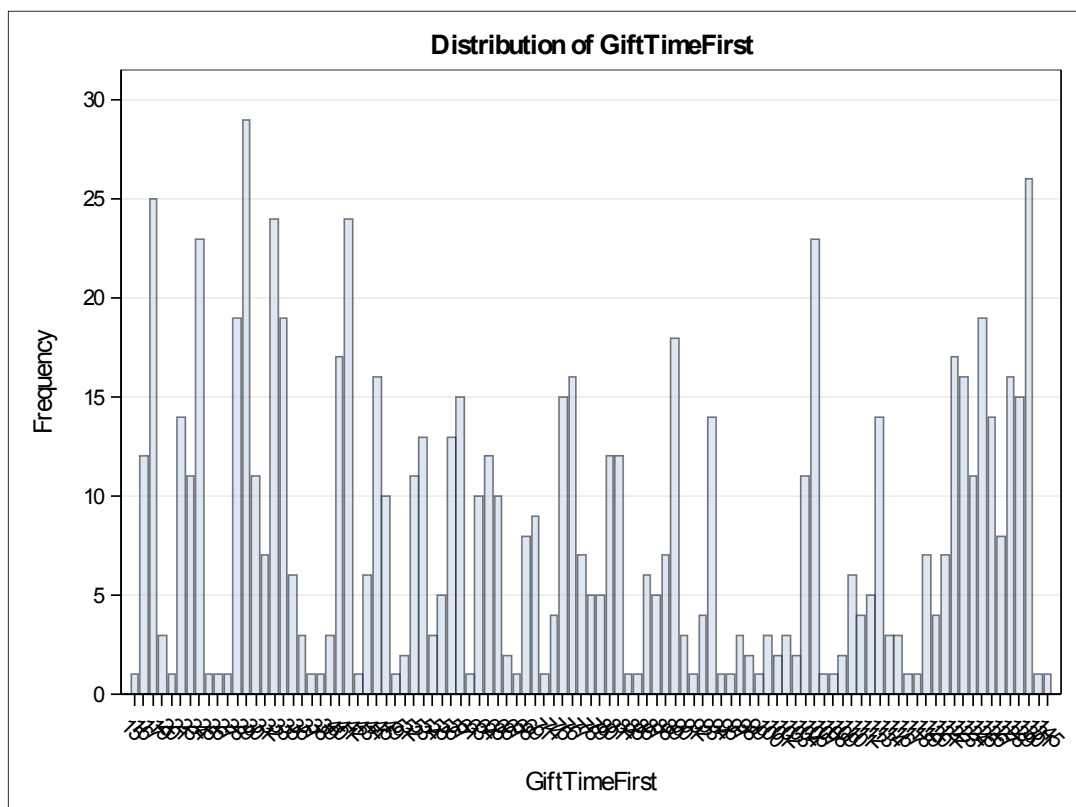
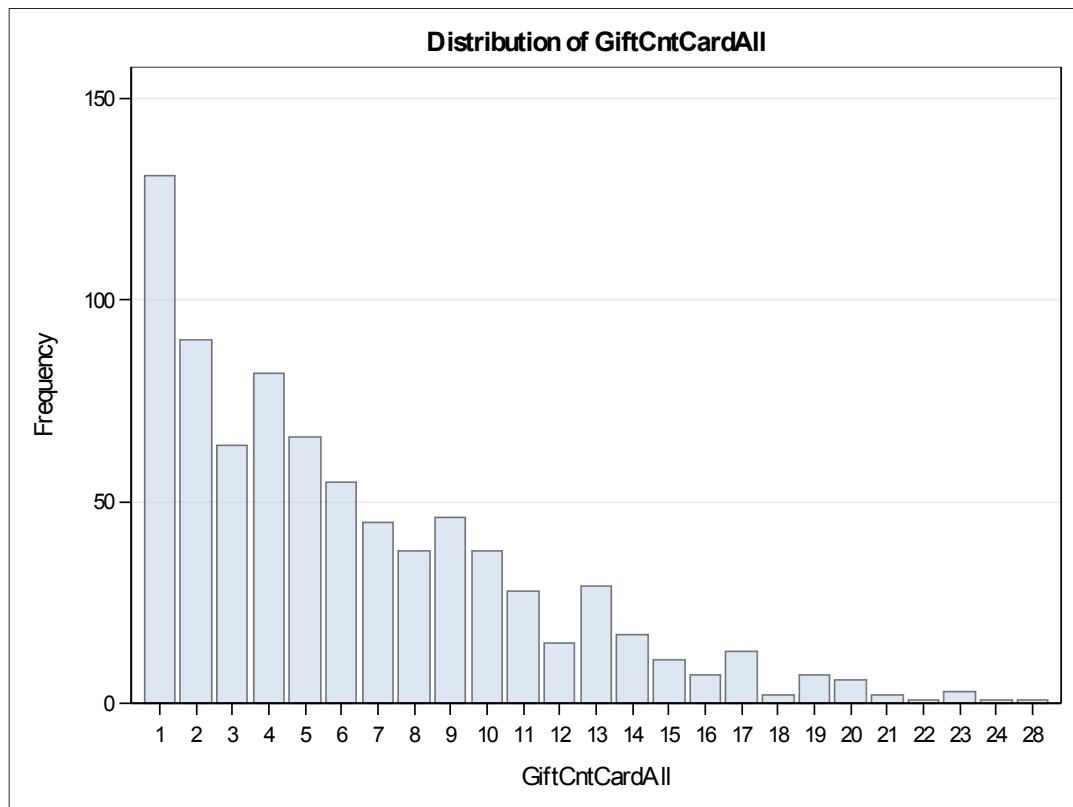


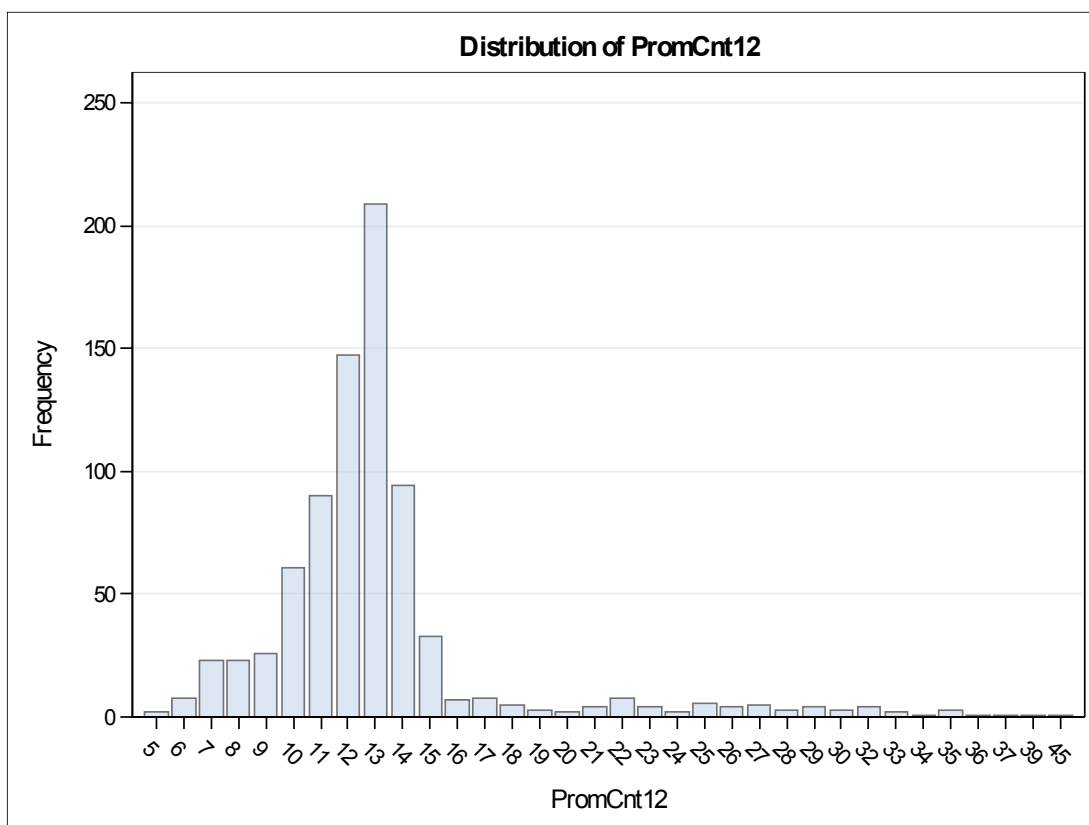
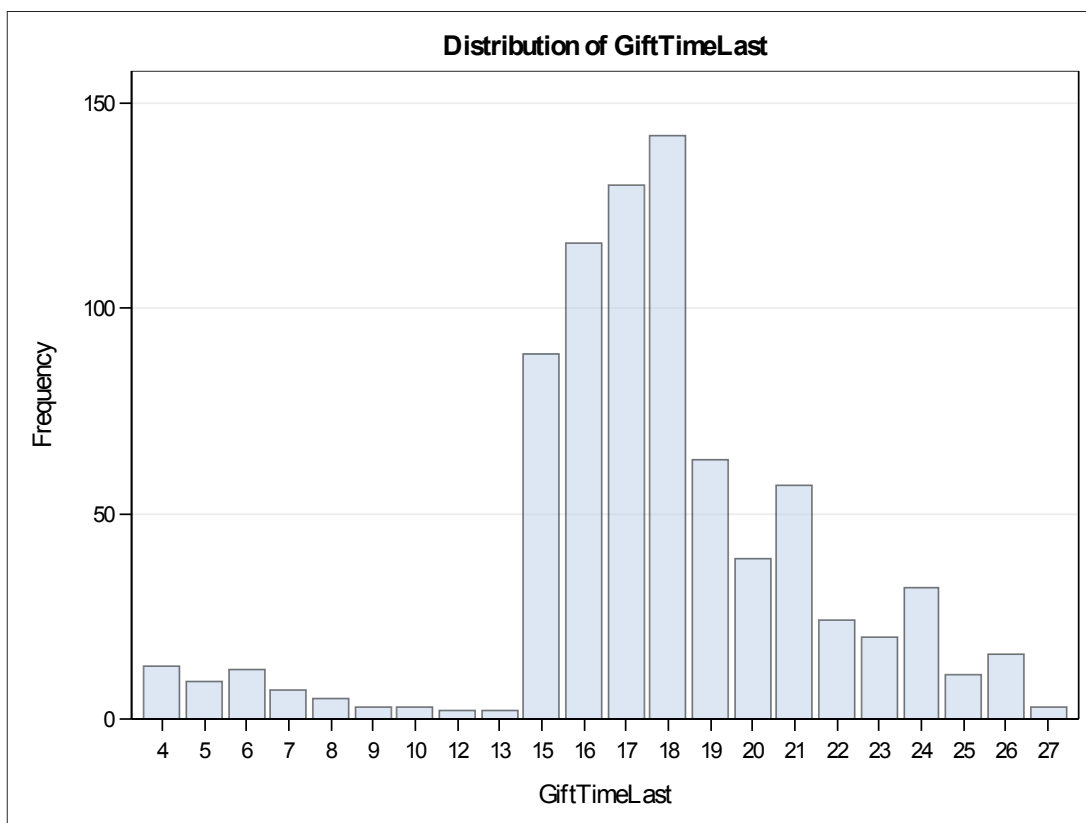


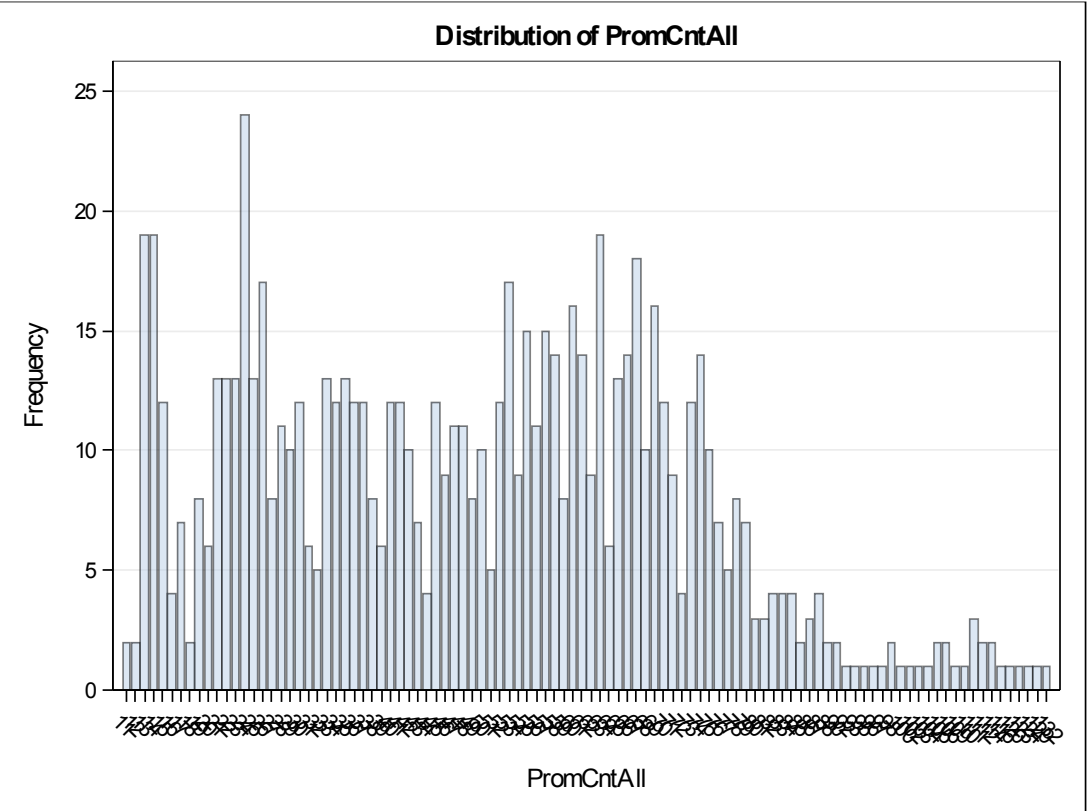
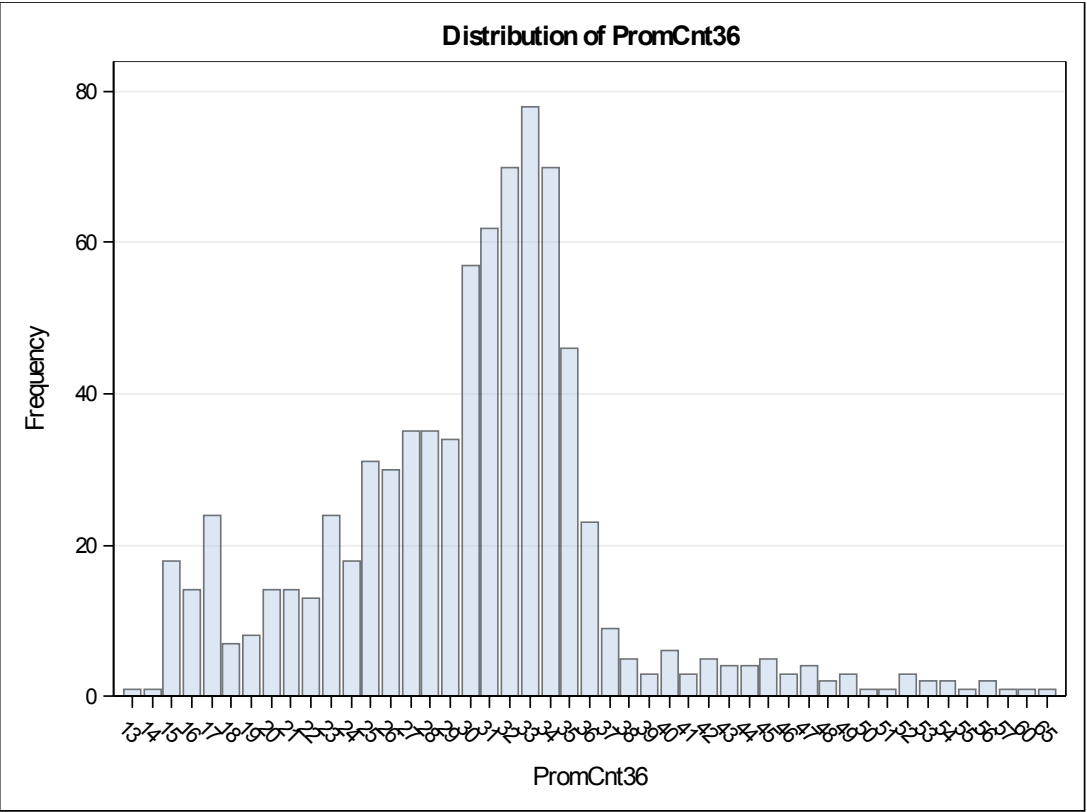


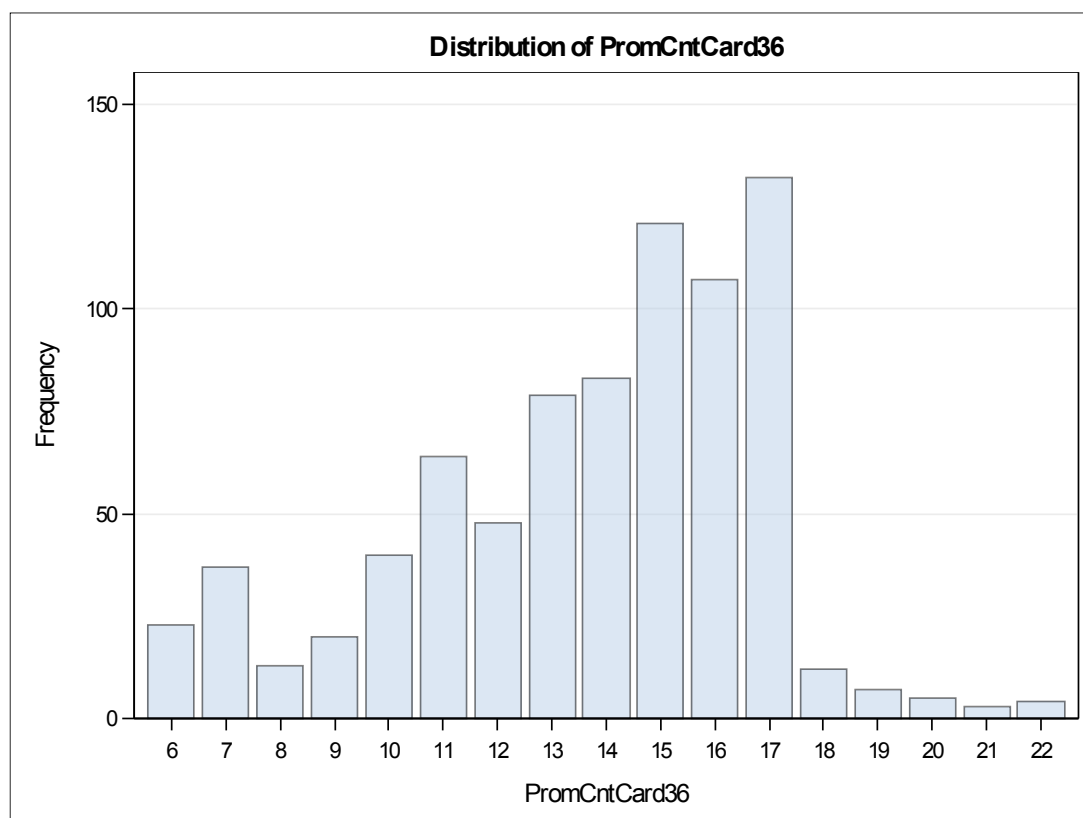
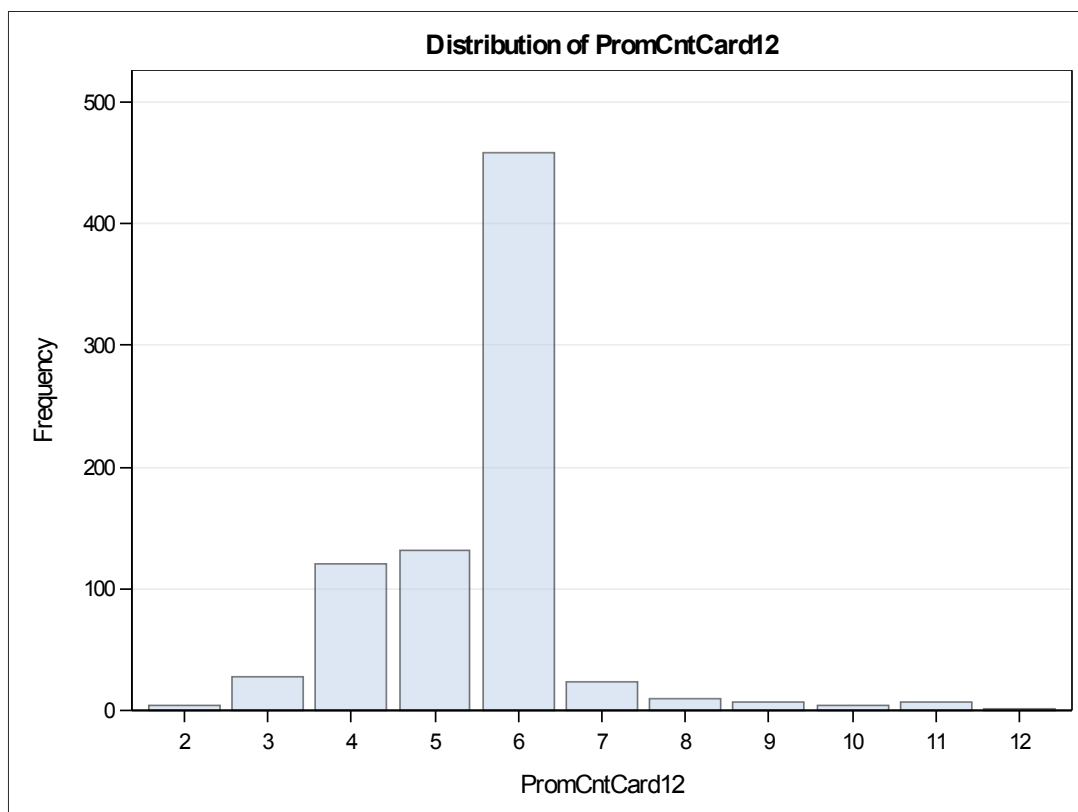


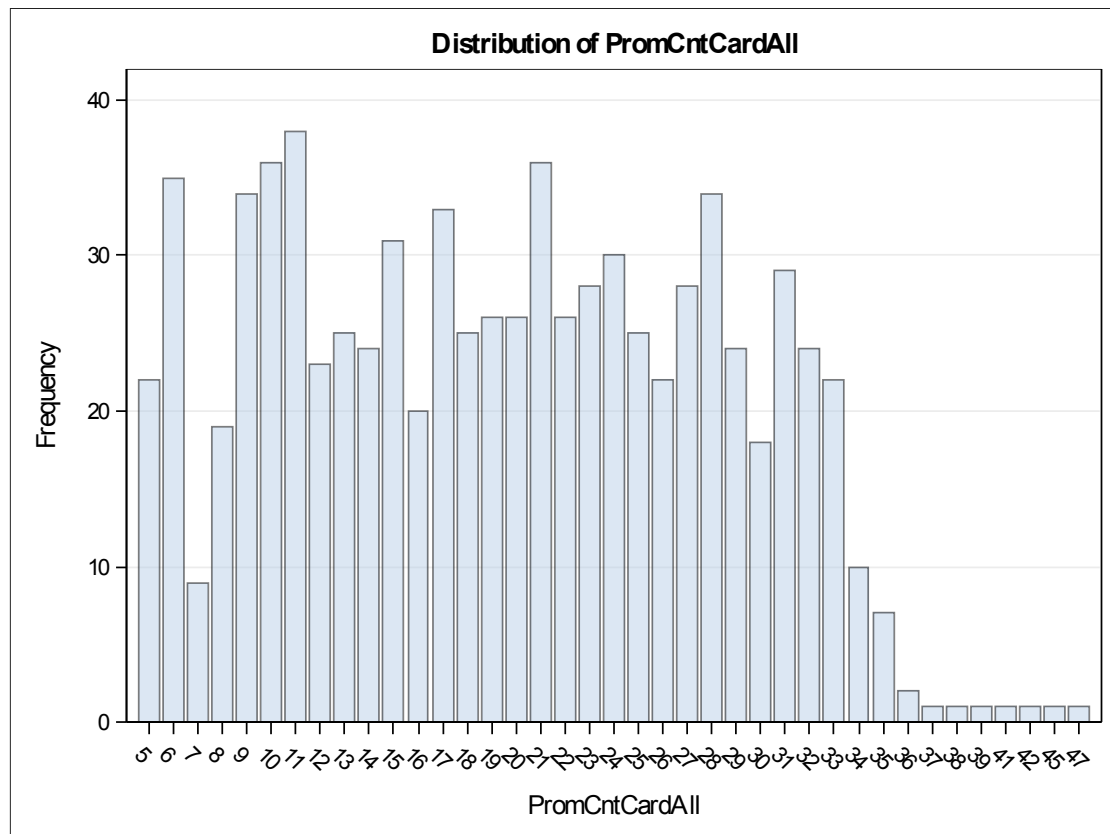












## Ερώτημα 2:

### Proc surveyselect και strata variables

Η surveyselect χρησιμοποιείται για την επιλογή δείγματος από ένα μεγάλο dataset, όπως αυτό της εργασίας. Το μέγεθος του δείγματος ισούται με το 70% του μεγέθους του καθαρισμένου dataset, άρα με 559 (από τα 798, στρογγυλοποίηση προς τα πάνω).

Για να κάνουμε stratified sampling στο εργαλείο SAS, πρέπει να χρησιμοποιηθεί το STRATA statement της surveyselect, το οποίο δέχεται μεταβλητές που εκφράζουν κατηγορία (όπως η DemHomeOwner στην δικιά μας περίπτωση). Με αυτό τον τρόπο, χωρίζουμε σε υποομάδες (**strata**) το δείγμα μας.

Το μέγεθος του κάθε stratum καθορίζεται από το ALLOC statement. Σε αυτήν την υλοποίηση επιλέχθηκε ALLOC=PROP, δηλαδή proportional allocation, που σημαίνει ότι αν (π.χ.) 60% του πληθυσμού έχει DemHomeOwner='H', τότε το 60% του δείγματος θα έχει DemHomeOwner='H'.



Ακολουθεί ο κώδικας για την δημιουργία του δείγματος.

```
113 /*sampling*/
114
115 title1 'Marketing campaign sampling (70% of population)';
116 title2 'Proportional Allocation';
117 proc surveyselect data=CAMP.cleaned_sorted
118     n=559 out=CAMP.sample;
119     strata DemHomeOwner / alloc=prop;
120 run;
```

Έπειτα, γίνεται εκτέλεση της διαδικασίας proc corr για να εμφανιστούν οι συσχετίσεις όλων των μεταβλητών με τη μεταβλητή targetD, όπως φαίνεται παρακάτω:

```
120 proc corr data=CAMP.sample;
121 var DemCluster DemIncomeGroup DemMedHomeValue GiftAvg36 GiftAvgAll GiftAvgCard36
122     GiftAvgLast GiftCnt36 GiftCntAll GiftCntCard36 GiftCntCardAll GiftTimeFirst
123     GiftTimeLast PromCnt12 PromCnt36 PromCntAll PromCntCard12 PromCntCard36
124     PromCntCardAll;
125 with targetD;
126 run;
```

Μετά την εκτέλεση του παραπάνω τμήματος προκύπτει ο παρακάτω πίνακας:

Pearson Correlation Coefficients, N = 559 Prob >  r  under H0: Rho=0							
	DemCluster	DemIncomeGroup	DemMedHomeValue	GiftAvg36	GiftAvgAll	GiftAvgCard36	
targetD	-0.11019	0.11320	0.09884	0.63304	0.53800	0.61434	
	0.0091	0.0074	0.0194	<.0001	<.0001	<.0001	
Pearson Correlation Coefficients, N = 559 Prob >  r  under H0: Rho=0							
	GiftAvgLast	GiftCnt36	GiftCntAll	GiftCntCard36	GiftCntCardAll	GiftTimeFirst	GiftTimeLast
targetD	0.76270	-0.26357	-0.18970	-0.19870	-0.16375	-0.11745	0.03967
	<.0001	<.0001	<.0001	<.0001	0.0001	0.0054	0.3492
Pearson Correlation Coefficients, N = 559 Prob >  r  under H0: Rho=0							
	PromCnt12	PromCnt36	PromCntAll	PromCntCard12	PromCntCard36	PromCntCardAll	
targetD	-0.02564	-0.09256	-0.10721	-0.12500	-0.14735	-0.13456	
	0.5452	0.0287	0.0112	0.0031	0.0005	0.0014	

Από αυτόν τον πίνακα συμπεραίνουμε ότι οι μεταβλητές με την μεγαλύτερη συσχέτιση με το targetD είναι όσες έχουν  $p\text{-value} < 0.0001$  και όσες είναι κοντά στο 1 κατ'απόλυτη τιμή. Αυτές οι μεταβλητές είναι οι GiftAvg36 GiftAvgAll GiftAvgCard36 GiftAvgLast GiftCnt36 GiftCntAll GiftCntCard36 και είναι αυτές που θα αποτελέσουν το μοντέλο πρόβλεψης που παρουσιάζεται παρακάτω:

```

127 proc reg data=CAMP.sample
128   outest=CAMP.estimates;
129
130 model targetD=GiftAvg36 GiftAvgAll GiftAvgCard36 GiftAvgLast
131   GiftCnt36 GiftCntAll GiftCntCard36/clm cli;
132
133 title "Regression of % targetD on sample";
134 run;
135
136 proc contents data=CAMP.estimates;
137 run;

```

Το αποτέλεσμα της παραπάνω εκτέλεσης φαίνεται στον παρακάτω πίνακα ANOVA:

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	7	55385	7912.07797	120.04	<.0001
Error	551	36316	65.91016		
Corrected Total	558	91701			

Root MSE	8.11851	R-Square	0.6040
Dependent Mean	14.47138	Adj R-Sq	0.5989
Coeff Var	56.10044		

Αποτελέσματα εκτέλεσης της διαδικασίας proc contents:

<b>Data Set Name</b>	CAMP.ESTIMATES	<b>Observations</b>	1
<b>Member Type</b>	DATA	<b>Variables</b>	13
<b>Engine</b>	V9	<b>Indexes</b>	0
<b>Created</b>	07/06/2015 13:33:48	<b>Observation Length</b>	13 6
<b>Last Modified</b>	07/06/2015 13:33:48	<b>Deleted Observations</b>	0
<b>Protection</b>		<b>Compressed</b>	N O
<b>Data Set Type</b>	EST	<b>Sorted</b>	N O
<b>Label</b>	Parameter Estimates and Statistics		
<b>Data Representation</b>	SOLARIS_X86_64, LINUX_X86_64, ALPHA_TRU64, LINUX_IA64		
<b>Encoding</b>	utf-8 Unicode (UTF-8)		

Engine/Host Dependent Information	
<b>Data Set Page Size</b>	65536
<b>Number of Data Set Pages</b>	1
<b>First Data Page</b>	1
<b>Max Obs per Page</b>	481
<b>Obs in First Data Page</b>	1
<b>Number of Data Set Repairs</b>	0
<b>Filename</b>	/folders/myfolders/myproject/estimates.sas7 bdat
<b>Release Created</b>	9.0401M2
<b>Host Created</b>	Linux
<b>Inode Number</b>	21785
<b>Access Permission</b>	rw-rw-rwx
<b>Owner Name</b>	root
<b>File Size (bytes)</b>	131072

Alphabetic List of Variables and Attributes				
#	Variable	Type	Len	Label
6	GiftAvg36	Num	8	
7	GiftAvgAll	Num	8	
8	GiftAvgCard36	Num	8	
9	GiftAvgLast	Num	8	
10	GiftCnt36	Num	8	
11	GiftCntAll	Num	8	
12	GiftCntCard36	Num	8	
5	Intercept	Num	8	Intercept
3	_DEPVAR_	Char	13	Dependent variable
1	_MODEL_	Char	32	Label of model
4	_RMSE_	Num	8	Root mean squared error
2	_TYPE_	Char	8	Type of statistics
13	targetD	Num	8	

### Ερώτημα 3:

Η ομαδοποίηση των δεδομένων με χρήση του k-means αλγόριθμου γίνεται με τη διαδικασία proc fastclus με βάση τις μεταβλητές που χρησιμοποιήθηκαν στην παραπάνω μοντελοποίηση, για k από 3 έως 7 (στον κώδικα το k ορίζεται ως maxclusters). Η εντολή maxiter=100 γίνεται για πιο ακριβές αποτέλεσμα εκτελώντας τον αλγόριθμο 100 φορές μέχρι να παραχθεί το αποτέλεσμα, ενώ γίνεται και μία ταξινόμηση του dataset με βάση το cluster για να μπορεί να τυπωθεί το αποτέλεσμα στον χρήστη. Ο κώδικας για την ομαδοποίηση είναι ο εξής:

```

139 /*clustering*/
140 proc fastclus data=CAMP.cleaned out=CAMP.Clust
141             maxclusters=7 maxiter=100;
142     var GiftAvg36 GiftAvgAll GiftAvgCard36 GiftAvgLast GiftCnt36 GiftCntAll
143         GiftCntCard36;
144 run;
145
146 proc sort data=CAMP.Clust;
147 by cluster;
148 run;
149 proc print;
150 by cluster;
151 run;

```

Μετά από πέντε εκτελέσεις του κώδικα αυτού για k από 3 έως 7 παρατηρήθηκε ότι καλύτερη ομαδοποίηση γίνεται για k=7, αφού γίνεται καλύτερη διαχείριση των outliers (απομόνωση σε ξεχωριστά clusters) με αποτέλεσμα να παράγεται ένα καλύτερο δείγμα δεδομένων.

Μετά από την επιλογή k=7 ομάδες επιλέγεται το cluster=7 λόγω περισσότερων παρατηρήσεων για το μοντέλο πρόβλεψης που θα δημιουργηθεί παρακάτω (Τα περισσότερα από τα υπόλοιπα clusters έχουν απομονώσει τα outliers και ως εκ τούτου δε μπορούν να χρησιμοποιηθούν για μοντελοποίηση). Αυτό παρατηρείται απο τον παρακάτω πίνακα:

Cluster Summary						
Cluster	Frequency	RMS Std Deviation	Maximum Distance from Seed to Observation	Radius Exceeded	Nearest Cluster	Distance Between Cluster Centroids
1	1	.	0		3	154.9
2	1	.	0		6	116.6
3	3	13.0043	33.3663		5	87.2145
4	302	3.4817	29.7666		7	18.6373
5	76	5.9202	51.5801		4	22.7041
6	1	.	0		2	116.6
7	414	4.3025	46.6018		4	18.6373

Ο κώδικας για τη συσχέτιση των μεταβλητών και τη μοντελοποίηση φαίνεται παρακάτω:

```

153 proc corr data=CAMP.Clust;
154 by cluster;
155 var DemCluster DemIncomeGroup DemMedHomeValue GiftAvg36 GiftAvgAll GiftAvgCard36
156     GiftAvgLast GiftCnt36 GiftCntAll GiftCntCard36 GiftCntCardAll GiftTimeFirst
157     GiftTimeLast PromCnt12 PromCnt36 PromCntAll PromCntCard12 PromCntCard36
158     PromCntCardAll;
159 with targetD;
160 run;
161

```

## Πίνακας συσχετίσεων:

Pearson Correlation Coefficients, N = 414 Prob >  r  under H0: Rho=0							
	DemCluster	DemIncomeGroup	DemMedHomeValue	GiftAvg36	GiftAvgAll	GiftAvgCard36	
targetD	-0.01153	0.02373	0.03524	0.40437	0.32760	0.39764	
	0.8151	0.6302	0.4746	<.0001	<.0001	<.0001	
Pearson Correlation Coefficients, N = 414 Prob >  r  under H0: Rho=0							
	GiftAvgLast	GiftCnt36	GiftCntAll	GiftCntCard36	GiftCntCardAll	GiftTimeFirst	GiftTimeLast
targetD	0.27508	-0.08317	-0.03349	-0.04636	-0.03560	0.00278	0.04905
	<.0001	0.0910	0.4967	0.3468	0.4701	0.9550	0.3194
Pearson Correlation Coefficients, N = 414 Prob >  r  under H0: Rho=0							
	PromCnt12	PromCnt36	PromCntAll	PromCntCard12	PromCntCard36	PromCntCardAll	
targetD	0.04155	0.05465	0.05347	0.04039	0.06319	0.03830	
	0.3991	0.2672	0.2777	0.4124	0.1995	0.4370	

Από τον πίνακα συμπεραίνουμε ότι οι μεταβλητές με την μεγαλύτερη συσχέτιση με το targetD είναι για το λόγο που αναφέρθηκε παραπάνω οι GiftAvg36 GiftAvgAll GiftAvgCard36 GiftAvgLast και είναι αυτές που θα αποτελέσουν το μοντέλο πρόβλεψης για το cluster=7 που παρουσιάζεται παρακάτω:

```

162 proc reg data=CAMP.Clust
163 outest=CAMP.estimates2;
164 by cluster;
165 model targetD=GiftAvg36 GiftAvgAll GiftAvgCard36 GiftAvgLast /clm cli;
166     title "Regression of % targetD in clusters";
167 run;
168
169 proc contents data=CAMP.estimates2;
170 run;
171

```

Το αποτέλεσμα της παραπάνω εκτέλεσης φαίνεται στον παρακάτω πίνακα ANOVA:

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	4	3089.73145	772.43286	20.59	<.0001
Error	409	15342	37.51215		
Corrected Total	413	18432			

Root MSE	6.12472	R-Square	0.1676
Dependent Mean	10.01087	Adj R-Sq	0.1595
Coeff Var	61.18066		

Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr >  t
Intercept	1	2.27451	0.95859	2.37	0.0181
GiftAvg36	1	0.63933	0.33156	1.93	0.0545
GiftAvgAll	1	-0.06934	0.21267	-0.33	0.7446
GiftAvgCard36	1	0.35997	0.26073	1.38	0.1681
GiftAvgLast	1	-0.00860	0.11283	-0.08	0.9393

Αποτελέσματα εκτέλεσης της διαδικασίας proc contents:

<b>Data Set Name</b>	CAMP.ESTIMATES2	<b>Observations</b>	7
<b>Member Type</b>	DATA	<b>Variables</b>	11
<b>Engine</b>	V9	<b>Indexes</b>	0
<b>Created</b>	07/06/2015 13:48:17	<b>Observation Length</b>	120
<b>Last Modified</b>	07/06/2015 13:48:17	<b>Deleted Observations</b>	0
<b>Protection</b>		<b>Compressed</b>	NO
<b>Data Set Type</b>	EST	<b>Sorted</b>	NO
<b>Label</b>	Parameter Estimates and Statistics		
<b>Data Representation</b>	SOLARIS_X86_64, LINUX_X86_64, ALPHA_TRU64, LINUX_IA64		
<b>Encoding</b>	utf-8 Unicode (UTF-8)		

Engine/Host Dependent Information	
Data Set Page Size	65536
Number of Data Set Pages	1
First Data Page	1
Max Obs per Page	545
Obs in First Data Page	7
Number of Data Set Repairs	0
Filename	/folders/myfolders/myproject/estimates2.sas7bdat
Release Created	9.0401M2
Host Created	Linux
Inode Number	24591
Access Permission	rw-rw-rw-
Owner Name	root
File Size (bytes)	131072

Alphabetic List of Variables and Attributes				
#	Variable	Type	Len	Label
1	CLUSTER	Num	8	Cluster
7	GiftAvg36	Num	8	
8	GiftAvgAll	Num	8	
9	GiftAvgCard36	Num	8	
10	GiftAvgLast	Num	8	
6	Intercept	Num	8	Intercept
4	_DEPVAR_	Char	13	Dependent variable
2	_MODEL_	Char	32	Label of model
5	_RMSE_	Num	8	Root mean squared error
3	_TYPE_	Char	8	Type of statistics
11	targetD	Num	8	

Ερώτημα 4:



Η αναγωγή του παραπάνω μοντέλου στο dataset μετά τη δειγματοληψία γίνεται ως εξής:

```
174 proc reg data=CAMP.sample  
175 outest=CAMP.estimates3;  
176 model targetD=GiftAvg36 GiftAvgAll GiftAvgCard36 GiftAvgLast /clm cli;  
177     title "Regression of % targetD on sample by model of cluster 6";  
178 run;  
179
```

Το αποτέλεσμα της παραπάνω εκτέλεσης φαίνεται στον παρακάτω πίνακα ANOVA:

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	4	54451	13613	202.45	<.0001
Error	554	37250	67.23889		
Corrected Total	558	91701			

Root MSE	8.19993	R-Square	0.5938
Dependent Mean	14.47138	Adj R-Sq	0.5909
Coeff Var	56.66311		

### Ερώτημα 5:

Παρατηρώντας του πίνακες ANOVA για τα δύο παραπάνω μοντέλα και ειδικότερα τιμές των πεδίων Sum of squares και R-Square τα οποία είναι μετρικές που πρέπει να λάβουμε υπόψη για να αποφασίσουμε ποιο από τα δύο μοντέλα είναι καλύτερο, συμπεραίνουμε ότι το αρχικό μοντέλο είναι καλύτερο από το δεύτερο. Το πεδίο R-square, το οποίο δείχνει πόσο καλά ταιριάζει ένα μοντέλο με την σχέση μεταβλητότητας ανεξάρτητης-εξαρτημένης μεταβλητής (το τετράγωνο της συσχέτισης), είναι ελαφρώς υψηλότερο στο πρώτο μοντέλο (r square υψηλότερο – μοντέλο καλύτερο). Το πεδίο Sum of squares όμως είναι μικρότερο στο δεύτερο μοντέλο (μικρότερο sum of squares σημαίνει καλύτερο μοντέλο). Κοιτώντας τα υπόλοιπα πεδία (γραμμή Error, το πρώτο μοντέλο έχει χαμηλότερες τιμές από το δεύτερο), αποφασίσαμε να χρησιμοποιήσουμε το πρώτο μοντέλο για την πρόβλεψη της δαπάνης.

### Ερώτημα 6:

Η εφαρμογή του καλύτερου μοντέλου είναι η εφαρμογή του αρχικού μοντέλου και η πρόβλεψη γίνεται με τον παρακάτω τρόπο:

- Εισαγωγή στοιχείων του αρχείου New\_Campaign.csv στο dataset CAMP.New\_Campaign

```
182 /*Import arxeiou New_Campaign */
183 /*Observations=50, Variables=18*/
184
185 data CAMP.New_Campaign ;
186 infile "&path/New_Campaign.csv" dlm=';' dsd firstobs=2;
187 input targetB id_ targetD:dollar10.2 GiftCnt36 GiftCntAll GiftCntCard36
188     GiftCntCardAll GiftAvgLast:dollar10.2 GiftAvg36:dollar10.2
189     GiftAvgAll:dollar10.2 GiftAvgCard36:dollar10.2 GiftTimeLast
190     GiftTimeFirst PromCnt12 PromCnt36 PromCntAll PromCntCard12
191     PromCntCard36 PromCntCardAll StatusCat96NK $ StatusCatStarAll
192     DemCluster DemAge DeMGender $ DemHomeOwner $ DemIncomeGroup
193     DemMedHomeValue:dollar10.2 DemPctVeterans;
194 run;
195
196 proc print data=CAMP.New_Campaign;
197 run;
198
```

- Πρόβλεψη του τελικού συνολικού ποσού δαπάνης

```
203 /*telikh problepsh posou*/
204
205 proc score data=CAMP.New_Campaign score=CAMP.estimates out=CAMP.scored
206     type=parms;
207     var GiftAvg36 GiftAvgAll GiftAvgCard36 GiftAvgLast
208         GiftCnt36 GiftCntAll GiftCntCard36;
209     title "Problepsh telikou posou";
210 run;
211
212 proc print data=CAMP.estimates;
213 run;
214
215 data CAMP.scored1;
216 set CAMP.scored;
217 if MODEL1 = '.' then delete;
218 run;
219
220 proc print data=CAMP.scored1;
221 var id_ MODEL1;
222 sum MODEL1;
223 run;
```

Πίνακας αποτελεσμάτων μετά την εκτέλεση της διαδικασίας proc  
score:

Obs	_MODEL_	_TYPE_	_DEPVAR_	_RMSE_	Intercept	GiftAvg36	GiftAvgAll
1	MODEL1	PARMS	targetD	8.11851	8.61576	-0.15486	-0.25077

Obs	GiftAvgCard36	GiftAvgLast	GiftCnt36	GiftCntAll	GiftCntCard36	targetD
1	0.31019	0.58109	-0.43353	-0.11673	0.38420	-1

Πίνακας αναλυτικού ποσού δαπάνης για κάθε πελάτη:

Obs	id_	MODEL 1
1	11961	19.127
2	12122	22.903
3	19020	11.889
4	20522	15.734
5	23989	14.413
6	27137	7.054
7	43492	14.085
8	53576	20.591
9	63135	7.424
10	65073	11.166
11	65189	12.859
12	76165	11.111
13	77523	13.269
14	87706	13.552
15	92972	25.907
16	96548	12.539
17	97596	15.158
18	102163	20.591
19	109858	15.352
20	112172	12.828
21	115166	9.527
22	120128	14.969

Obs	id_	MODEL 1
23	124935	20.751
24	125164	39.136
25	126787	21.028
26	135265	12.778
27	137256	19.937
28	139628	18.253
29	140294	14.793
30	142838	11.154
31	151907	11.364
32	155653	18.163
33	157427	6.560
34	162494	16.874
35	168275	19.989
36	171265	20.145
37	172167	15.145
38	172750	17.958
39	173772	10.423
40	177190	11.914
41	179863	12.020
42	181337	14.882
43	182519	15.734
44	187274	15.734
45	189725	11.400
		<b>698.185</b>

Άρα τελικά το συνολικό ποσό δαπάνης της εταιρίας για τα δώρα των πελατών αναμένεται στα 698.185\$.