

LoanLogisticRegModel

Laura Ahumada, Erin McClure-Price, Duy Nguyen

EDA

Looking into the data

Statistics

- ID can be dropped since it is not a useful predictor and just a unique identifier and will affect the results of the model
- Age mean is 45 years old, 23 has is the lowest and 67 is the highest. There seems to be a normal distribution.
- ***There are negative values in Experience which is odd because there can't be negative years of experience. This will be looked at next.***
- The mean income is 73 while the median is 64 showing skewness which seems normal representation of society. The minimum income is 8,000 and the max is 224,000. These are common Salaries
- ZIP.Code should not be numeric they should be changed to categorical. There are 467 distinct ZIP codes
- There seems to be an equal distribution for families size 1,2,3,4. Each are compose of about 25%
- The CCAVG, Average Spending per 1000 goes from 0 to 10,000, however the median is 1,5000. Here we can see that there are outliers.
- For Education, 41% have up to highschool, 28% up to under grad studies and interestingly 30% have up to grad school. That seems like a reasonable distribution.
- For mortgage we can see how skewed it is, the mean is 56.5 and the median is 0. Showing the extreme outliers. **This needs to be looked at**
- For the target variable, personal loan, we can see quite a difference, 90% without a loan and only 10% with loan. This makes sense because not many people take loan from banks
- Securities account also has a big difference where 90% does not have a security account while 10% does
- For CD.Account (Ceritificate deposit) once again we see the 94% does not have it while 6% has it.
- As for online (online banking capability), 40% doe not have it while 60% has it. That distribution is a little more balanced and makes sense.
- As for Credit card 70% does not have it while 30% has it.

ID	Age	Experience	Income	ZIP.Code
Min. : 1	Min. :23.00	Min. : -3.0	Min. : 8.00	Min. : 9307
1st Qu.:1251	1st Qu.:35.00	1st Qu.:10.0	1st Qu.: 39.00	1st Qu.:91911
Median :2500	Median :45.00	Median :20.0	Median : 64.00	Median :93437
Mean :2500	Mean :45.34	Mean :20.1	Mean : 73.77	Mean :93152
3rd Qu.:3750	3rd Qu.:55.00	3rd Qu.:30.0	3rd Qu.: 98.00	3rd Qu.:94608
Max. :5000	Max. :67.00	Max. :43.0	Max. :224.00	Max. :96651
Family	CCAVg	Education	Mortgage	
Min. :1.000	Min. : 0.000	Min. :1.000	Min. : 0.0	
1st Qu.:1.000	1st Qu.: 0.700	1st Qu.:1.000	1st Qu.: 0.0	
Median :2.000	Median : 1.500	Median :2.000	Median : 0.0	

```

Mean      :2.396   Mean      : 1.938   Mean      :1.881   Mean      : 56.5
3rd Qu.:3.000   3rd Qu.: 2.500   3rd Qu.:3.000   3rd Qu.:101.0
Max.      :4.000   Max.      :10.000   Max.      :3.000   Max.      :635.0
Personal.Loan  Securities.Account  CD.Account      Online
Min.      :0.000   Min.      :0.0000   Min.      :0.0000   Min.      :0.0000
1st Qu.:0.000   1st Qu.:0.0000   1st Qu.:0.0000   1st Qu.:0.0000
Median :0.000   Median :0.0000   Median :0.0000   Median :1.0000
Mean      :0.096   Mean      :0.1044   Mean      :0.0604   Mean      :0.5968
3rd Qu.:0.000   3rd Qu.:0.0000   3rd Qu.:0.0000   3rd Qu.:1.0000
Max.      :1.000   Max.      :1.0000   Max.      :1.0000   Max.      :1.0000
CreditCard
Min.      :0.000
1st Qu.:0.000
Median :0.000
Mean      :0.294
3rd Qu.:1.000
Max.      :1.000

```

Looking into questions obtained in the statistical analysis

- *Looking at the entries with negative experience*
- The 52 people with negative total experience are between 23 to 29 years old and salary median of 65,000. Sound like these could just young people that just started working. We can change their negative values to 1 since they have an income, thus are are working.

```

      Age      Income      Personal.Loan
Min.    :23.00   Min.    : 12.00   Min.    :0
1st Qu.:24.00   1st Qu.: 40.75   1st Qu.:0
Median :24.00   Median : 65.50   Median :0
Mean    :24.52   Mean    : 69.94   Mean    :0
3rd Qu.:25.00   3rd Qu.: 86.75   3rd Qu.:0
Max.    :29.00   Max.    :150.00   Max.    :0

```

- Setting categories as factors
- Creating a new variable called Exprience 2 due to Age and experience having a correlation of 97%. That way we can get rid of Age and Experience

PersonalLoan

```

11 Variables      5000 Observations
-----
Income
  n missing distinct      Info      Mean      Gmd      .05      .10
5000      0      162      1      73.77      50.91      18      22
.25      .50      .75      .90      .95
39      64      98      145      170

lowest :   8   9  10  11  12, highest: 203 204 205 218 224
-----
Family
  n missing distinct
5000      0      4

```

Value	1	2	3	4
Frequency	1472	1296	1010	1222
Proportion	0.294	0.259	0.202	0.244

CCAvg

n	missing	distinct	Info	Mean	Gmd	.05	.10
5000	0	108	0.999	1.938	1.794	0.1	0.3
.25	.50	.75	.90	.95			
0.7	1.5	2.5	4.3	6.0			

lowest : 0.0 0.1 0.2 0.3 0.4, highest: 8.8 8.9 9.0 9.3 10.0

Education

n	missing	distinct
5000	0	3

Value	1	2	3
Frequency	2096	1403	1501
Proportion	0.419	0.281	0.300

Mortgage

n	missing	distinct	Info	Mean	Gmd	.05	.10
5000	0	347	0.668	56.5	88.16	0	0
.25	.50	.75	.90	.95			
0	0	101	200	272			

lowest : 0 75 76 77 78, highest: 590 601 612 617 635

Personal.Loan

n	missing	distinct
5000	0	2

Value	No	Yes
Frequency	4520	480
Proportion	0.904	0.096

Securities.Account

n	missing	distinct
5000	0	2

Value	0	1
Frequency	4478	522
Proportion	0.896	0.104

CD.Account

n	missing	distinct
5000	0	2

Value	0	1
Frequency	4698	302
Proportion	0.94	0.06

Online

n	missing	distinct
---	---------	----------

```

5000      0      2

Value      0      1
Frequency  2016  2984
Proportion 0.403 0.597

```

```

CreditCard
  n missing distinct
5000      0      2

```

```

Value      0      1
Frequency  3530  1470
Proportion 0.706 0.294

```

```

Experience2
  n missing distinct  Info    Mean    Gmd    .05    .10
5000      0      178     1  0.404  0.1848  0.07678  0.14286
.25      .50      .75    .90    .95
0.28571  0.44444  0.53846  0.58730  0.60317

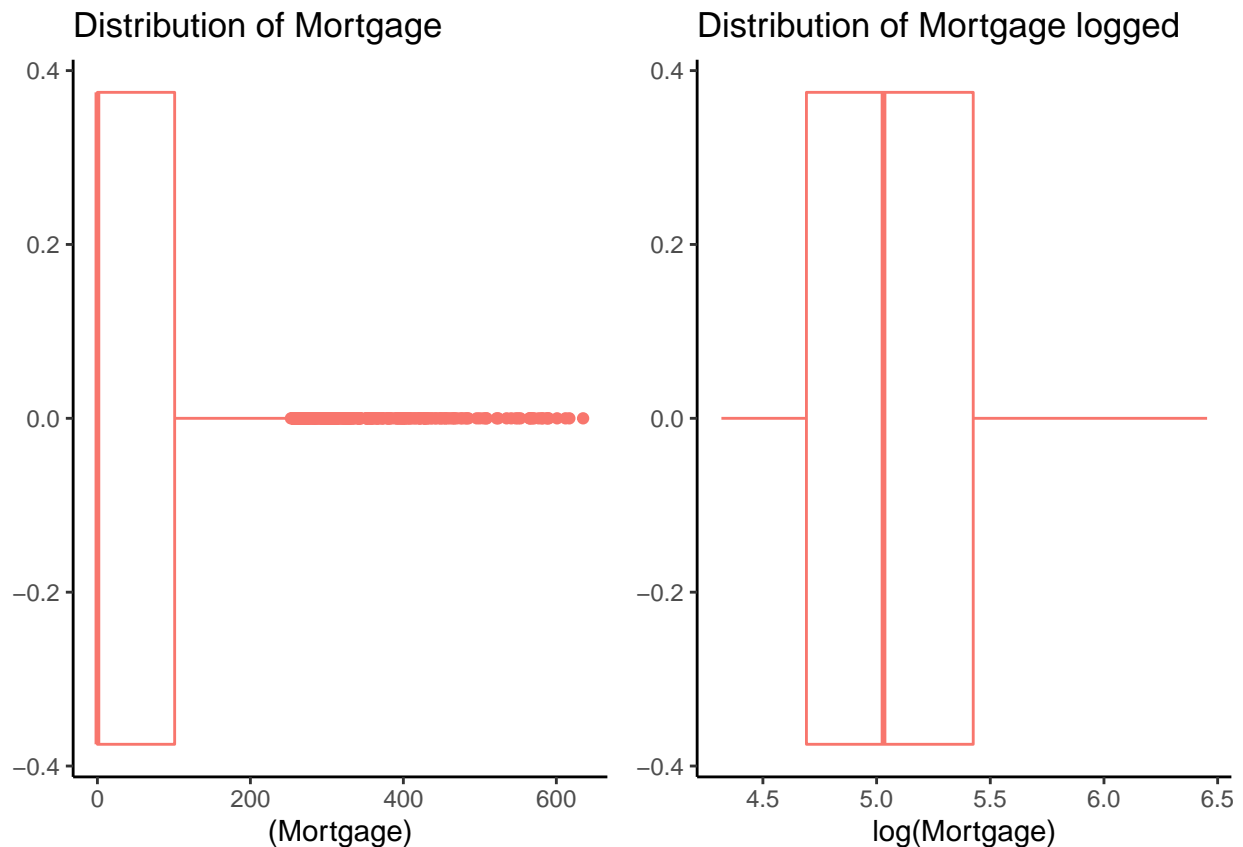
```

```

lowest : 0.00000000 0.03225806 0.03448276 0.03571429 0.03703704
highest: 0.62500000 0.62686567 0.63076923 0.63636364 0.64179104

```

- *Checking mortgage distribution*
- Mortgage may need to be logged as it is very skewed.

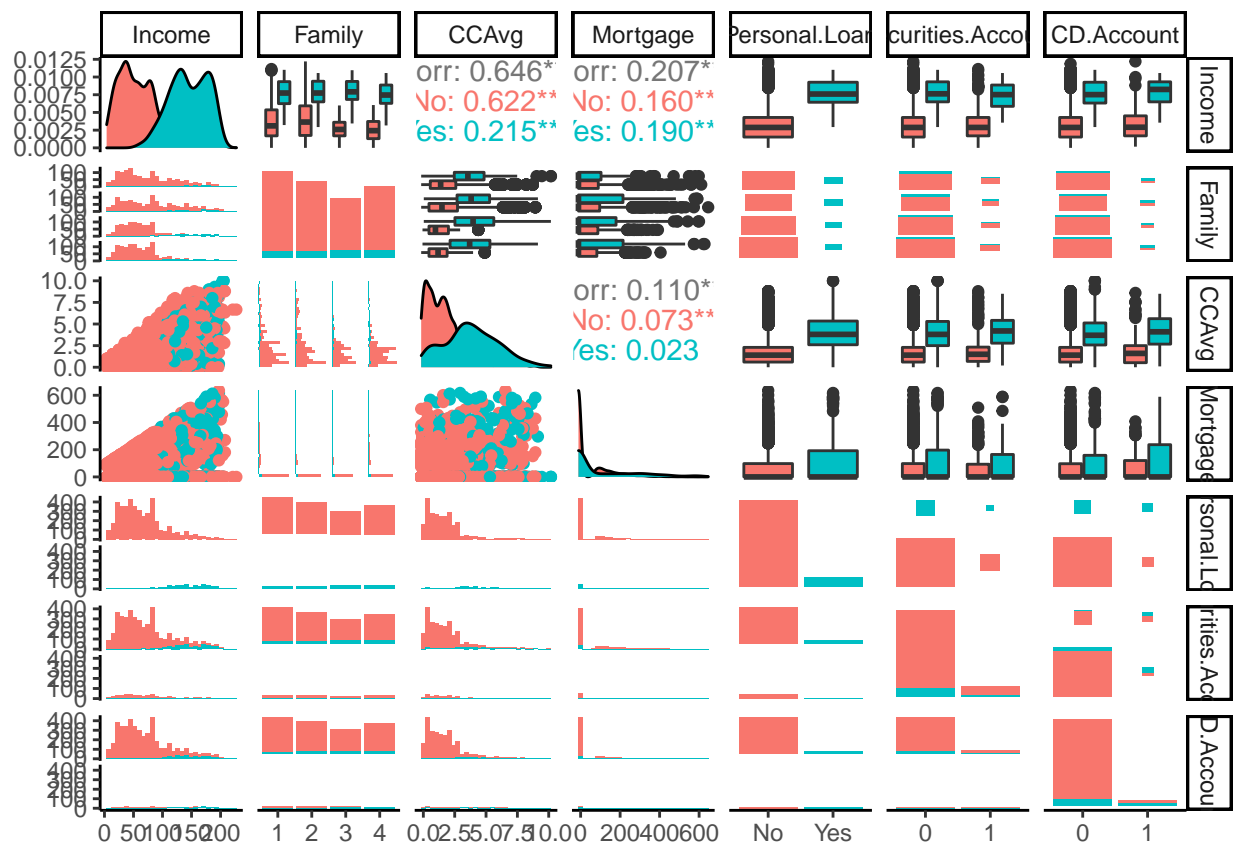


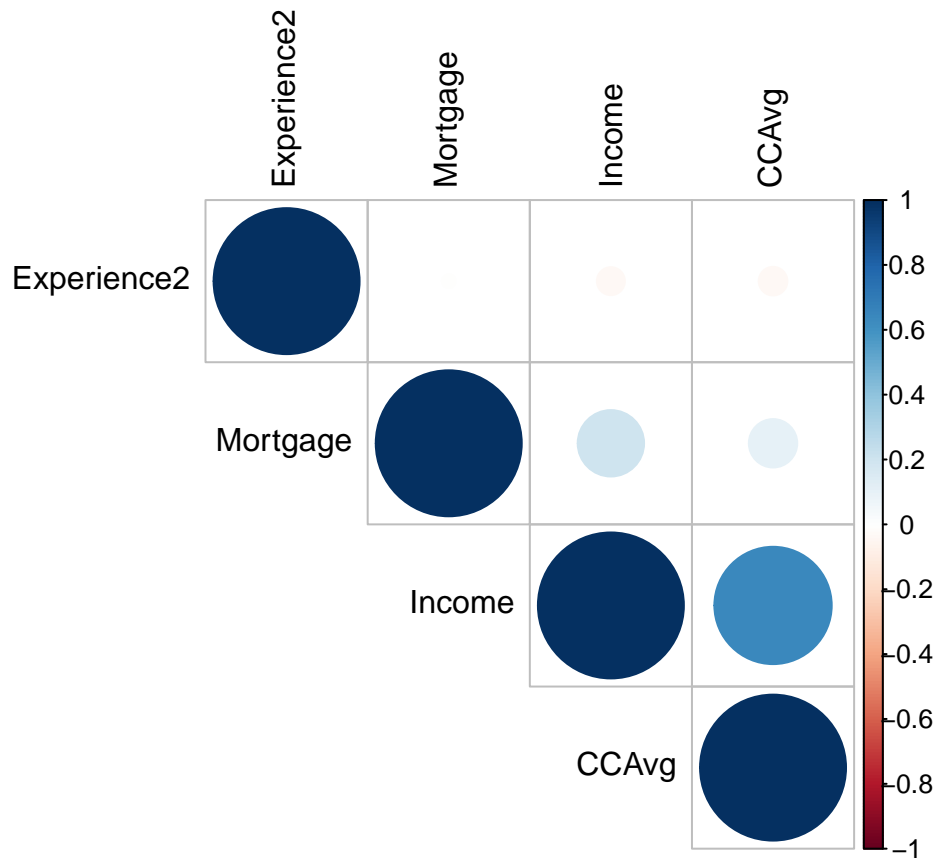
- Creating another data set with Mortgage logged since logging did improve the distribution. The zero's were replaced to with 1 in order to not get infitly when logging. Also, removing age since it has a 99% correlation with Experience, therefore only one is needed

```
#Creating a new data set with the modified attributes
PersonalLoan2=mutate(PersonalLoan)
#Updating values
PersonalLoan2$Mortgage[PersonalLoan2$Mortgage==0]=1
PersonalLoan2$MortgageLogged=log(PersonalLoan2$Mortgage)
#PersonalLoan2=PersonalLoan2 %>% select(-Mortgage)
```

Relationships and correlations

- Experience and Age has a Correlation of 99%.Too high. However, they seem to have no relationship with Loan as both; yes loan and no loan, have the same correlation with experience and age of 0.99
- CCAvg and Income has a 64% correlations and it does seems to have a relationship with Loan since there is 0.62 with no loan and .02 with yes loan.
- The rest of the explanatory variables do not seem to have relationship between each other





Checking relationship between Loan (response variable) and the rest of the predictors

Very high Relationship present

- *Income*, the more income the more changes to get a loan
- *Mortgage*, people high mortgages seem to ask for loans more so than those with lower mortgages

Relationship present

- *CD.Account* (certificate of deposit) seem to have a relationship with Personal Loan. Those with personal loan tend to have CD.Account more so than those with no CD.Account
- *Personal education*, people with up to highschool tend do not get loans as much as does with an education of higher than highschool
- `!!**there seems to be relationship with loan and education but that of education 2 and 3 seem to be`

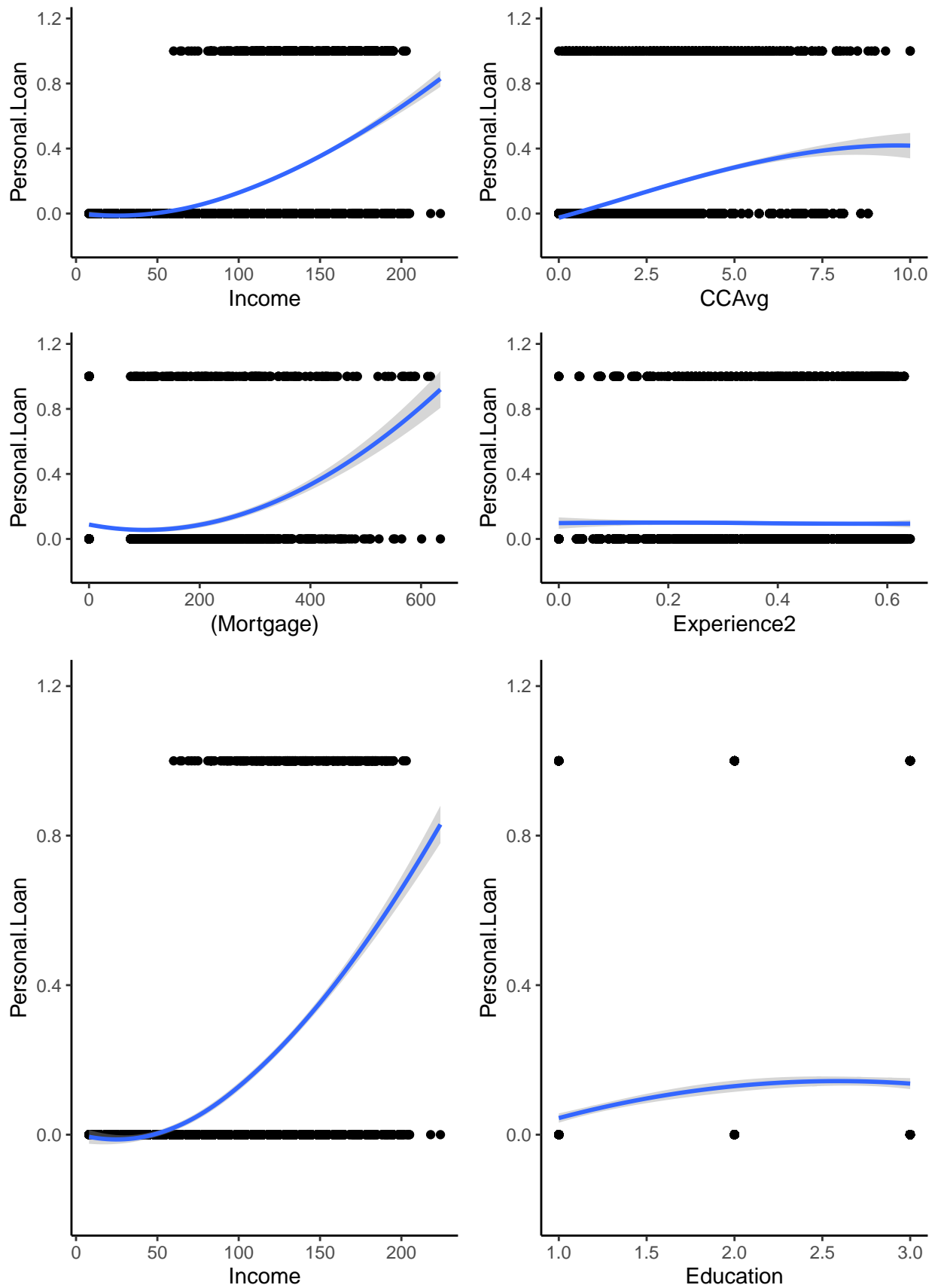
Very small relationship

- *Security account* and Personal Loan seem to have slight relationship

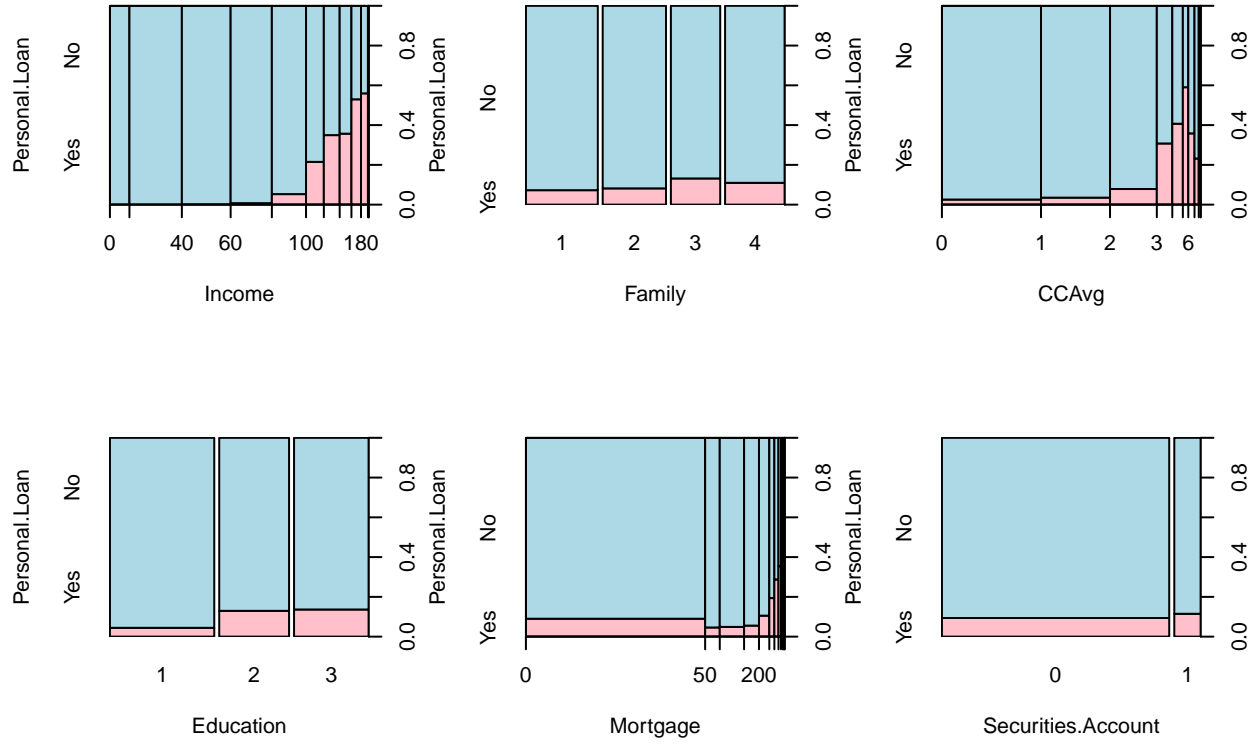
No relationship with response

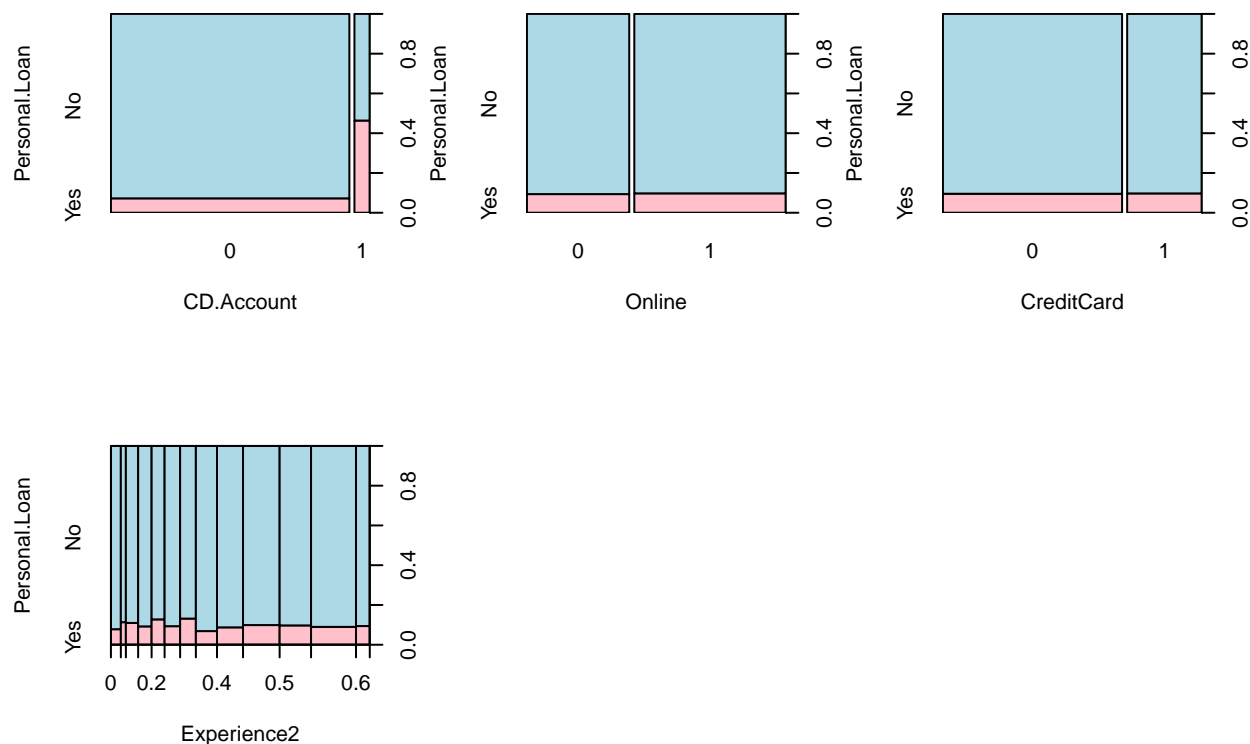
- *Online and Credit card* don't seem to have a relationship

Looking closer at each relationship to see if anything was missed



- More detailed graph on each variable with response





Model: Objective 1

```
# Train Test Split
set.seed(123)
index<-sample(1:dim(PersonalLoan)[1],round(.70 * dim(PersonalLoan)[1]))
train<-PersonalLoan[index,]
test<-PersonalLoan[-index,]

# Split Predict for lasso
dat.test.x = model.matrix(Personal.Loan ~ Income + Family + CCAvg-1 + Education + Securities.Account-1 +
dat.train.x = model.matrix(Personal.Loan ~ Income + Family + CCAvg-1 + Education + Securities.Account-1 +
dat.train.y = train$Personal.Loan
```

Performing model with all variables, some feature selection methods (forward, stepwise, LASSO) and another based on EDA + With the full model with all attributes it showed that the only important were Income, Family, CCAvg, Education, Securites.Account, CD.Account, Online, CreditCard

- Once Stepwise was added to the full model it selected all of those that appeared as significant in the full model: Income, Family, CCAvg, Education, Securities.Account, CD.Account, Online, CreditCard. It also selected Experience2 but that one was not significant as we had seen in the EDA.

- When the forward model was added to the full model it selected the same thing as stepwise but included CCAVG which was significant and Mortgage which was not significant .
- As for LASSO it selected all of the attributes by Stepwise and included CCAvg.

Call:

```
glm(formula = Personal.Loan ~ ., family = "binomial", data = train)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.8582	-0.1818	-0.0647	-0.0193	4.1124

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	-12.5251518	0.7051046	-17.764	< 0.0000000000000002	***
Income	0.0632191	0.0037351	16.926	< 0.0000000000000002	***
Family2	-0.2376211	0.2832713	-0.839	0.401556	
Family3	1.9082393	0.2950650	6.467	0.0000000000998	***
Family4	1.3733143	0.2907382	4.724	0.0000023177032	***
CCAvg	0.1330478	0.0558477	2.382	0.017203	*
Education2	4.0306143	0.3351091	12.028	< 0.0000000000000002	***
Education3	4.1978282	0.3350516	12.529	< 0.0000000000000002	***
Mortgage	0.0009203	0.0007186	1.281	0.200303	
Securities.Account1	-0.7246928	0.3493892	-2.074	0.038063	*
CD.Account1	3.5261823	0.4039920	8.728	< 0.0000000000000002	***
Online1	-0.8210820	0.2024400	-4.056	0.0000499357901	***
CreditCard1	-0.8793800	0.2592112	-3.393	0.000693	***
Experience2	-0.0876285	0.5689759	-0.154	0.877601	

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 2163.7 on 3499 degrees of freedom
Residual deviance: 792.4 on 3486 degrees of freedom
AIC: 820.4

Number of Fisher Scoring iterations: 8

```
[1] 0.0009615423
```

- Using different thresholds

Changing the threshold

```
[1] "All_Attributes"
```

```
[1] "Threshold | Accuracy| Sensitivity| Specificity"
```

```
$'0.3'
```

```
[1] 0.3000000 0.9566667 0.9747212 0.8000000
```

```

$'0.5'
[1] 0.5000000 0.9620000 0.9947955 0.6774194

$'0.7'
[1] 0.7000000 0.9553333 0.9992565 0.5741935

[1] "-----"

[1] "StepWiseAIC"

[1] "Threashhold | Accuracy| Sensitivity| Specificy"

$'0.3'
[1] 0.3000000 0.9573333 0.9747212 0.8064516

$'0.5'
[1] 0.5000000 0.9620000 0.9947955 0.6774194

$'0.7'
[1] 0.7000000 0.9553333 0.9992565 0.5741935

[1] "-----"

[1] "Threashhold | Accuracy| Sensitivity| Specificy"

[1] "ForwardModel"

$'0.3'
[1] 0.3000000 0.9566667 0.9747212 0.8000000

$'0.5'
[1] 0.5000000 0.9620000 0.9947955 0.6774194

$'0.7'
[1] 0.7000000 0.9553333 0.9992565 0.5741935

[1] "-----"

[1] "Threashhold | Accuracy| Sensitivity| Specificy"

[1] "LASSO"

$'0.3'
[1] 0.3000000 0.9560000 0.9747212 0.7935484

$'0.5'
[1] 0.5000000 0.9613333 0.9955390 0.6645161

$'0.7'
[1] 0.7000000 0.9553333 1.0000000 0.5677419

```

```

[1] "-----"

[1] "Intuition"

[1] "Threashhold | Accuracy| Sensitivity| Specificy"

$'0.3'
[1] 0.3000000 0.9553333 0.9762082 0.7741935

$'0.5'
[1] 0.5000000 0.9593333 0.9947955 0.6516129

$'0.7'
[1] 0.7000000 0.9540000 0.9992565 0.5612903

[1] "-----"

[1] "EDA"

[1] "Threashhold | Accuracy| Sensitivity| Specificy"

$'0.3'
[1] 0.3000000 0.9533333 0.9739777 0.7741935

$'0.5'
[1] 0.5000000 0.9600000 0.9955390 0.6516129

$'0.7'
[1] 0.7000000 0.9553333 0.9992565 0.5741935

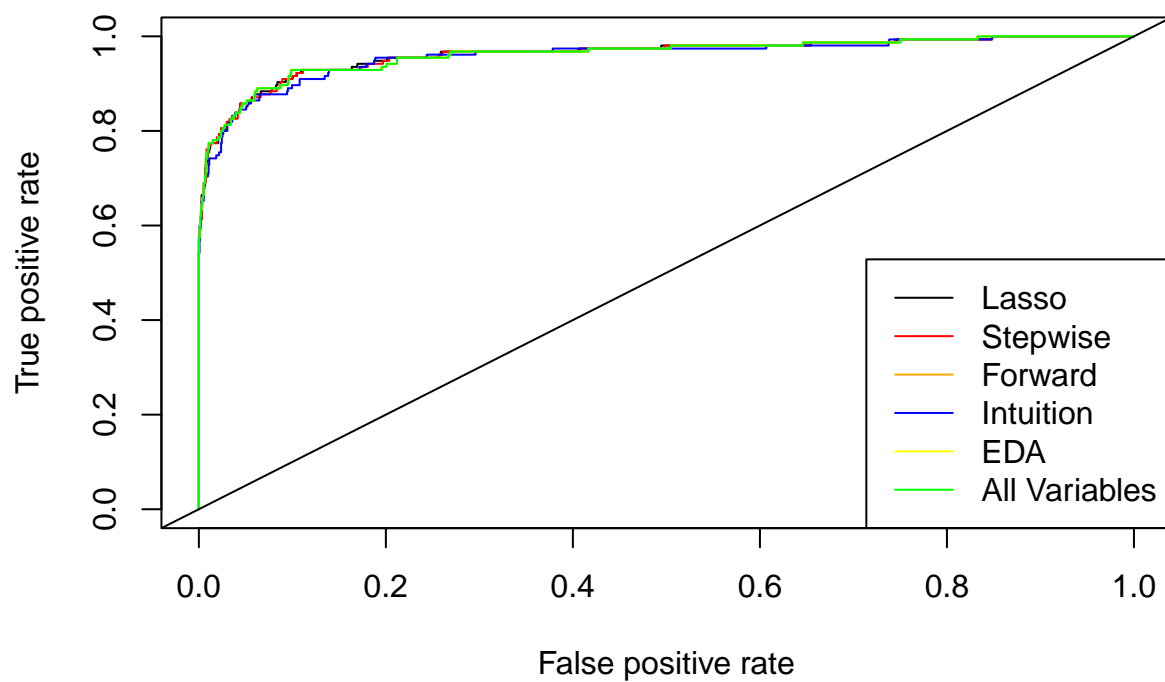
```

Choosing 0.3 threshold based of the threashold results. Criterion Comparison of all models

```
all_results
```

	Criterion	Full_Model	Step_Wise	Forward_Model	LASSO_model	Intuition	EDA
1	AIC	820.398	818.074	820.398	0.000	836.879	844.481
2	BIC	906.645	892.000	906.645	0.000	892.324	899.926
3	Accuracy	0.957	0.957	0.957	0.956	0.955	0.953
4	Sensitivity	0.975	0.975	0.975	0.975	0.976	0.974
5	Specificity	0.800	0.806	0.800	0.794	0.774	0.774

The ROC of models



Verify Proportions in test and train manually + Distribution in train and test do represent that of the whole data

```
[1] "All data"
```

	No	Yes
	0.904	0.096

```
[1] "Train"
```

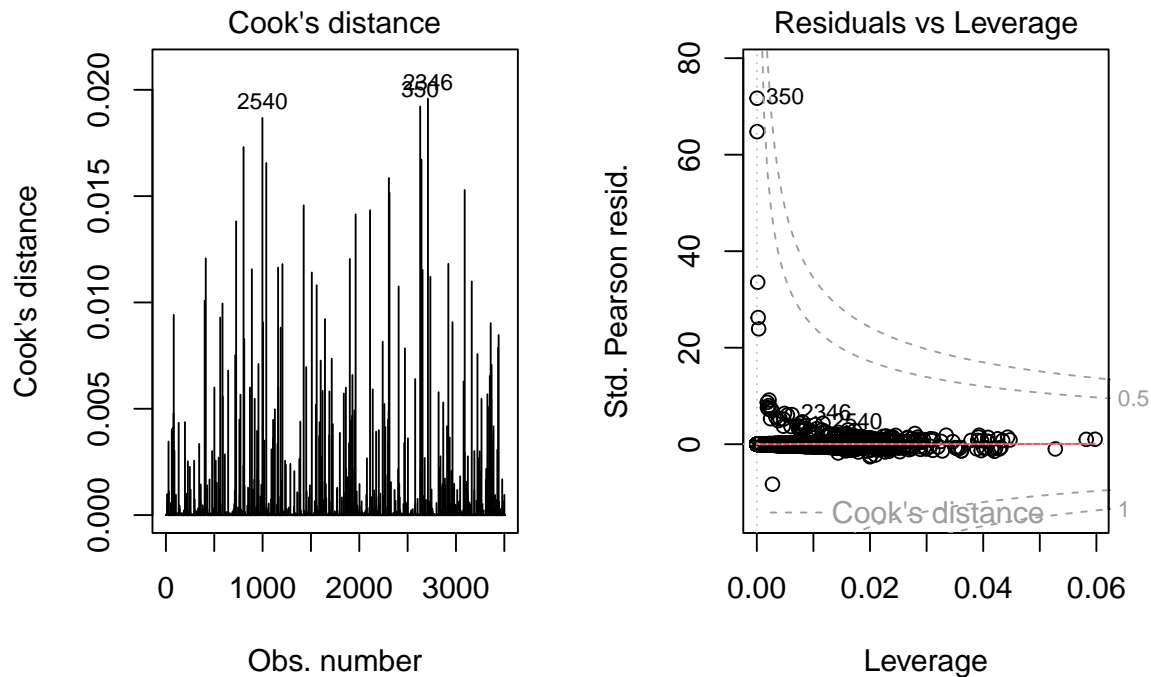
	No	Yes
	0.90714286	0.09285714

```
[1] "Test"
```

	No	Yes
	0.8966667	0.1033333

- Assumptions via PLOTS of selected model, Stepwise and checking VIF

- Plots look normal and there seems to be multicollinearity among variables based on VIF



	GVIF	Df	GVIF ^{1/(2*Df)}
Income	2.940809	1	1.714879
Family	1.529409	3	1.073381
CCAvg	1.516750	1	1.231564
Education	2.323075	2	1.234570
Securities.Account	1.291648	1	1.136507
CD.Account	1.936714	1	1.391659
Online	1.143566	1	1.069376
CreditCard	1.383602	1	1.176266

Conclusion from Part 1

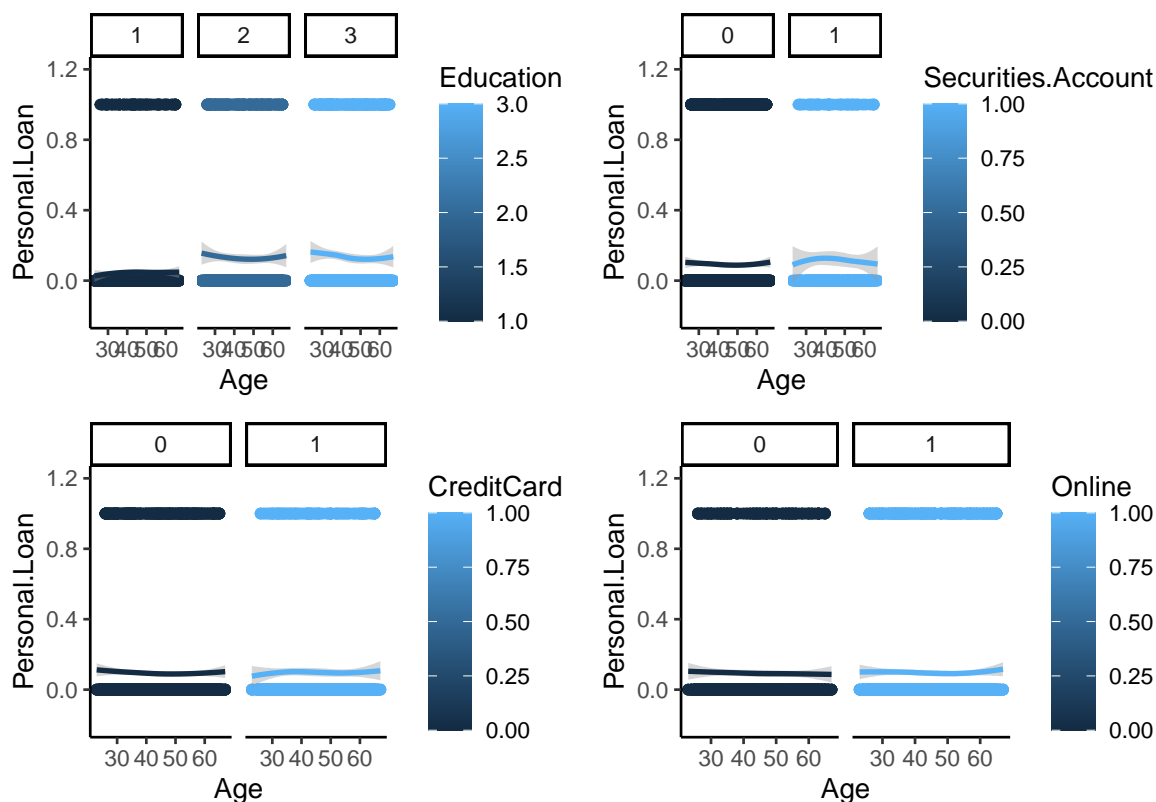
- The best model was step setting the threshold to 0.3 it gave an sensitivity of 94 and specificity of 72
- Due to the imbalance of amount of people with loan and without loan we do see we do see that the model favors no loan due to it but 72 compared to the 55 specificity was a great increase. This model is about trying to predict those who will say yes to Loan therefore Specificity is important.
- The attributes found useful were : Income, Family,CCAvg, Education,Securites.Account, CD.Account, Online, CreditCard
- The threshold was set to 0.3 and it lead to a Sensitivity of 0.96 and specificy of 0.71
- These were variables seen in the EDA as related to the loan.
- Coefficients results: For every unit increase in income the odd of getting a loan are $e^{1.06}$ times higher For every unit increase in Family the odd of getting a loan are $e^{0.698209}$ times higher For every unit increase in CCAvg the odd of getting a loan are $e^{0.120635}$ times higher For every unit increase in Education the odd of getting a loan are $e^{1.713690}$ times higher For every unit increase in

Securities.Account 1 the odd of getting a loan are $e^{-0.937183}$ times less likely For every unit increase in CD.Account1 the odd of getting a loan are $e^{3.840892}$ times higher For every unit increase in Online1 the odd of getting a loan are $e^{-0.673230}$ times less likely For every unit increase in CreditCard1 the odd of getting a loan are $e^{-1.122701}$ times higher

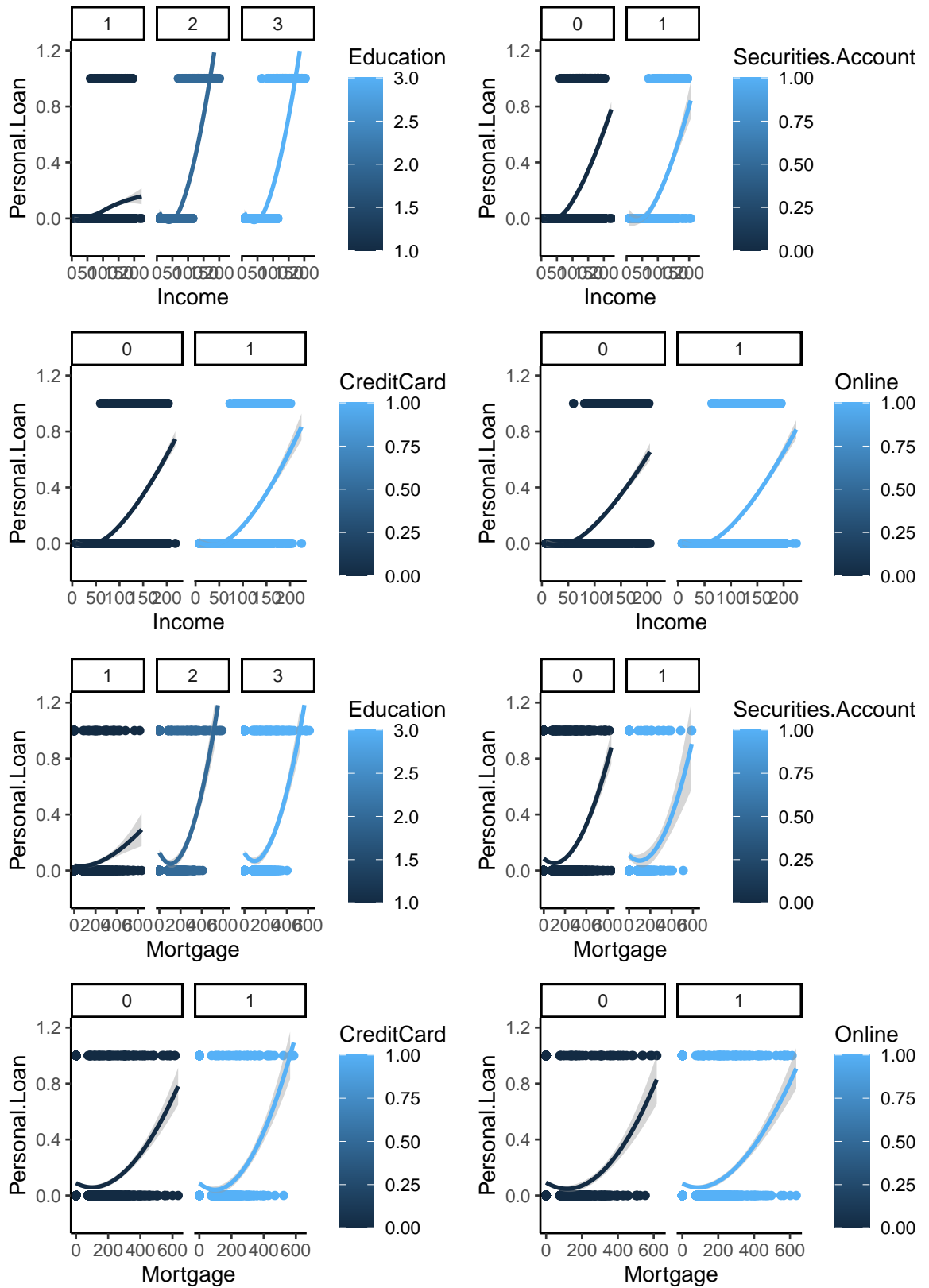
#####

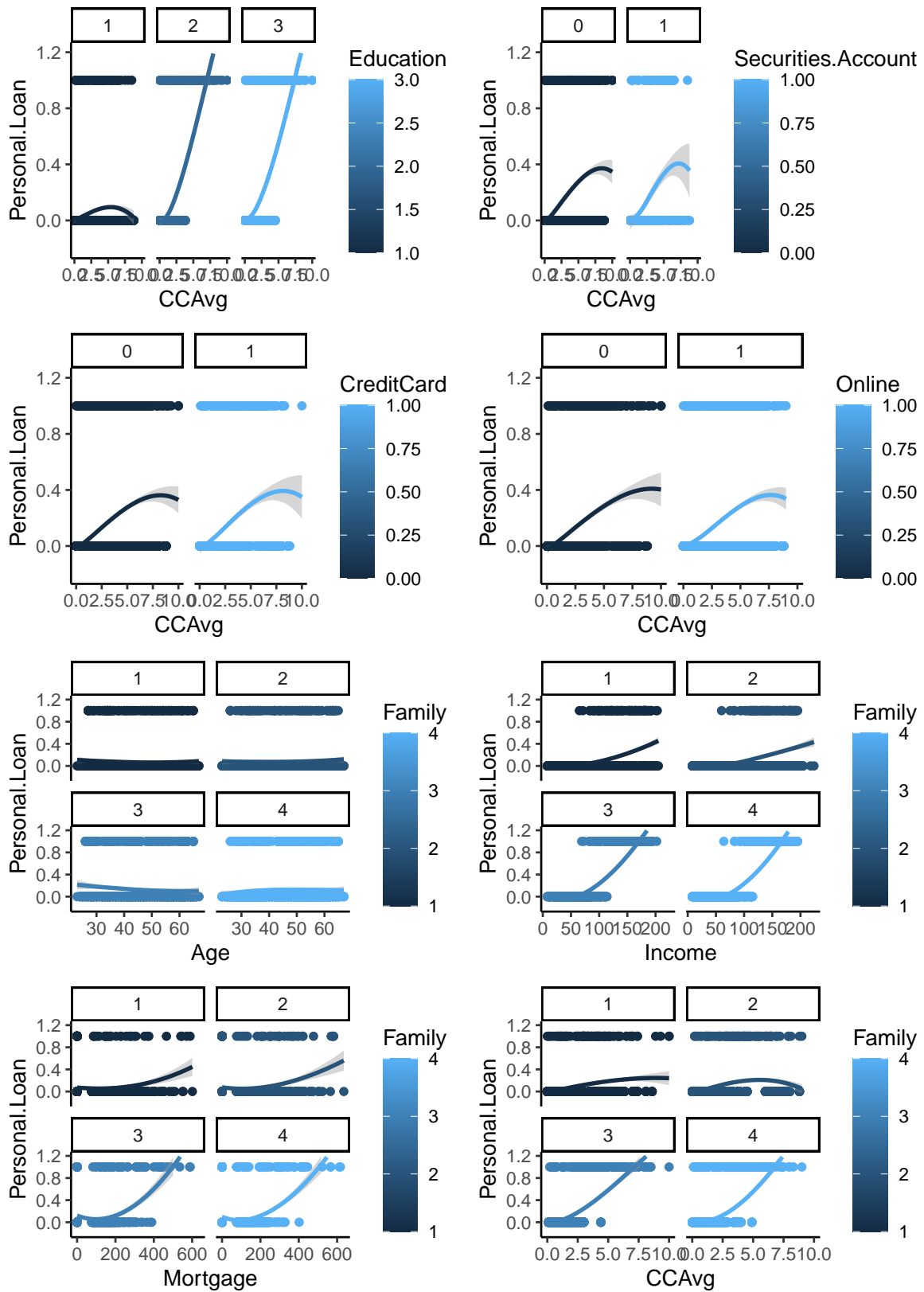
MODEL: Part 2

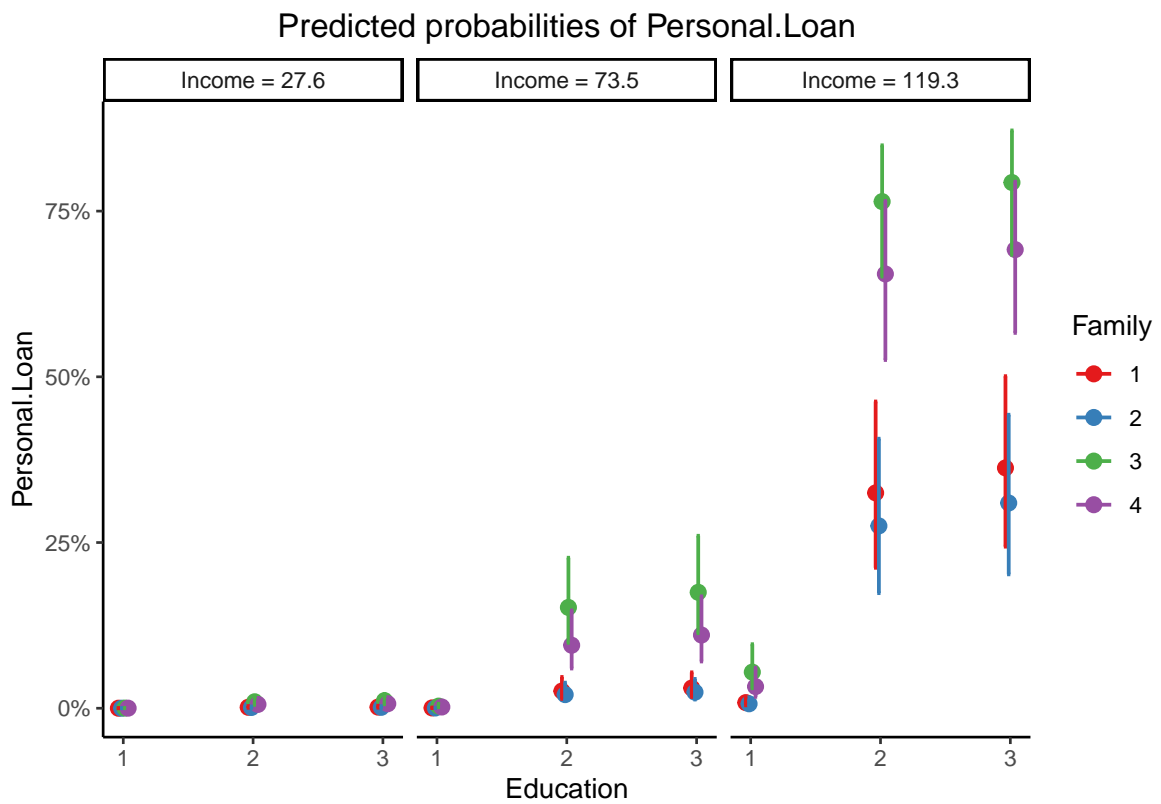
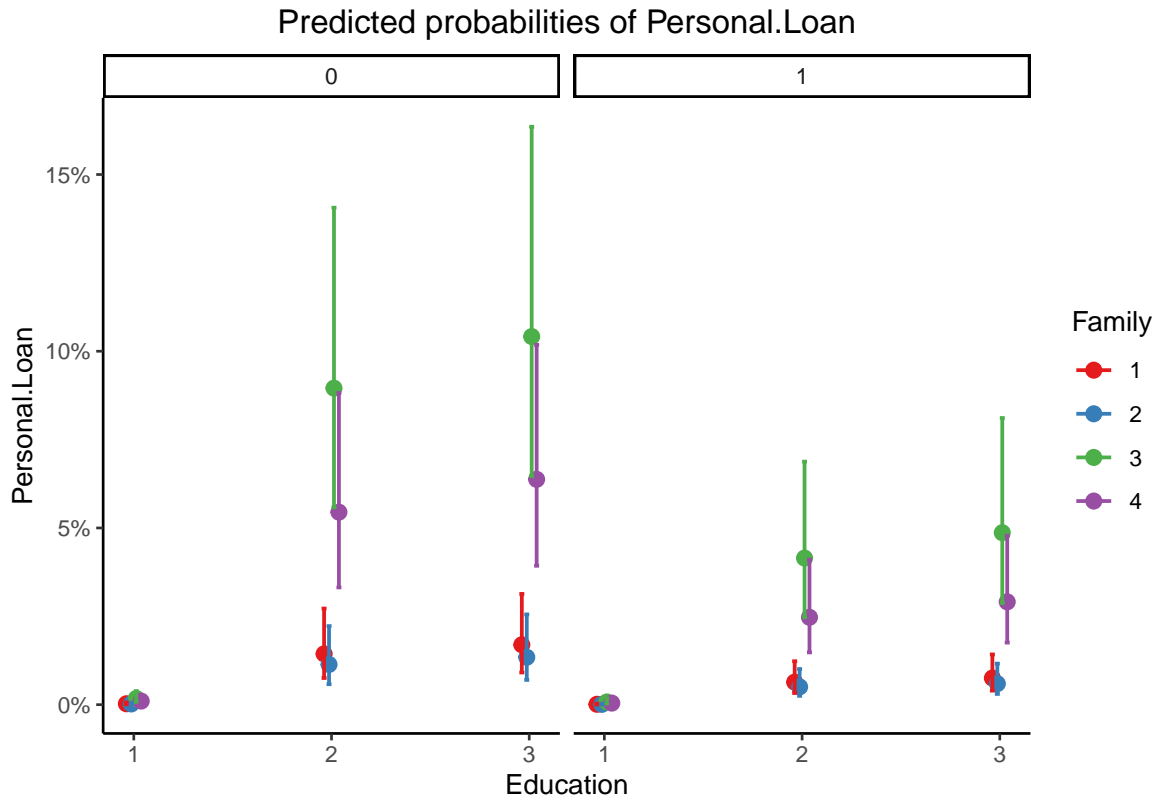
- Checking graphs for interactions
- Across all of the plots only one age with mortgage showed to maybe have an interaction

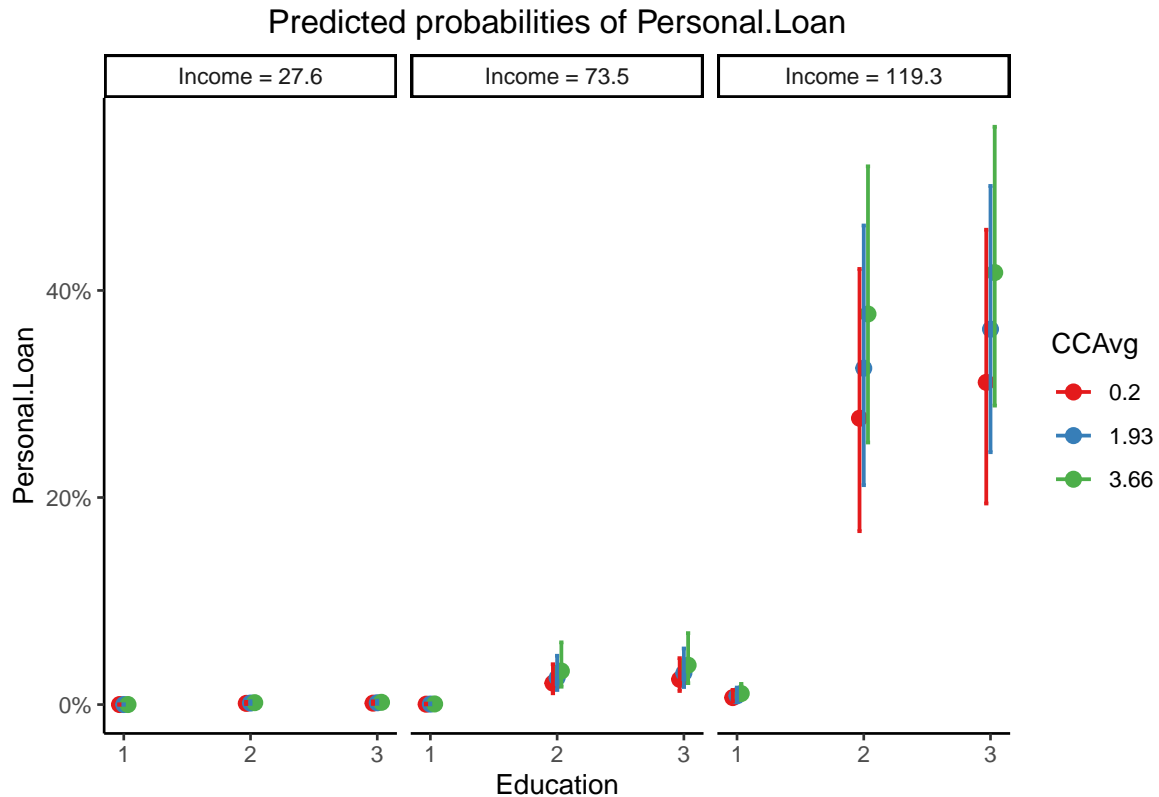


.









• Running the models

- In part 1 we saw that, Income, Family, CCAvg, CD.Account, Education and Credit Card were significant, therefore we leverage those variables and kept mortgage (just like in part 1) based on the EDA analysis, for Part 2

Interaction

- All variables: Income, Family, CCAvg, CD.Account, Education, and Credit Card and the interaction of family and mortgage were significant.
- When passing the model through stepwise and forward it kept the all variables including the interaction and thus kept the same significant variables showing signs that the interaction is indeed useful.
- When Anova was applied, it did show the interaction as significant
- The Hoslem test however showed that the model was a poor fit, yes again this is not reliable due to the size of the data (This I asked in class and Dr.Turner said we couldn't rely on this metric with big data)

Logged + We logged mortgage + All variables were significant: Income, Family, CCAvg, CD.Account, Education, and Credit Card including the log mortgage and log income + Once the models were mixed, having logged income and the interaction of logged mortgage with family lead to the interaction no longer being significant however all other variable remained significant.

Polynomials + When income was set to poly 2 all variables were significant + The income was set to 3 polys, the first poly income was significant but the rest were no longer significant. However, the rest of the variables are still significant + CCAvg set to poly2 showed all variables as significant just like poly 3 + when setting both CCAvg and income to poly 2 all variables showed importance however when setting both CCAvg and income to poly 3, income ² and income ³ did not show significance in the model while all

the rest of the attributes did show significance. + Interesting when MortgageLogged or mortgage was added it reduced the specificity and sensitivity + Therefore the attributes used where income, Family, CCAvg, CD.Account, Education, Credit Card . we brought up Securities.Account again as well as online since we had seen that they had importance in some models and EDA. Of course the difference within the models was having Incompo Poly or CCAvg or both to 2 and 3 polynomial.

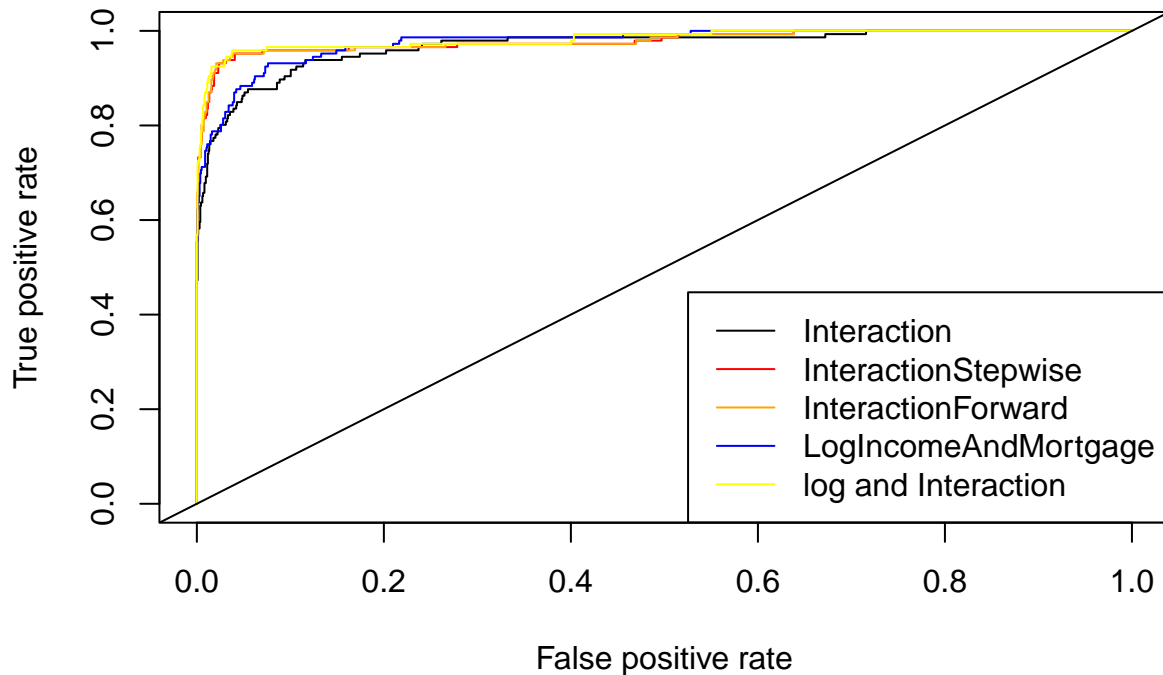
LDA and QDA + LDA mode and QDA model did poorly however between the two, LDA model outperformed QDA. This makes sense since the groups have a clear division of who gets loans and who doesn't which would be LDA. + Also this is a non parametric test therefore we can also explain why it did not do as good as the logisti regression models

```
Personal.Loan ~ Securities.Account + CD.Account + CreditCard +
  Education + Income + CCAvg + Family * Mortgage
```

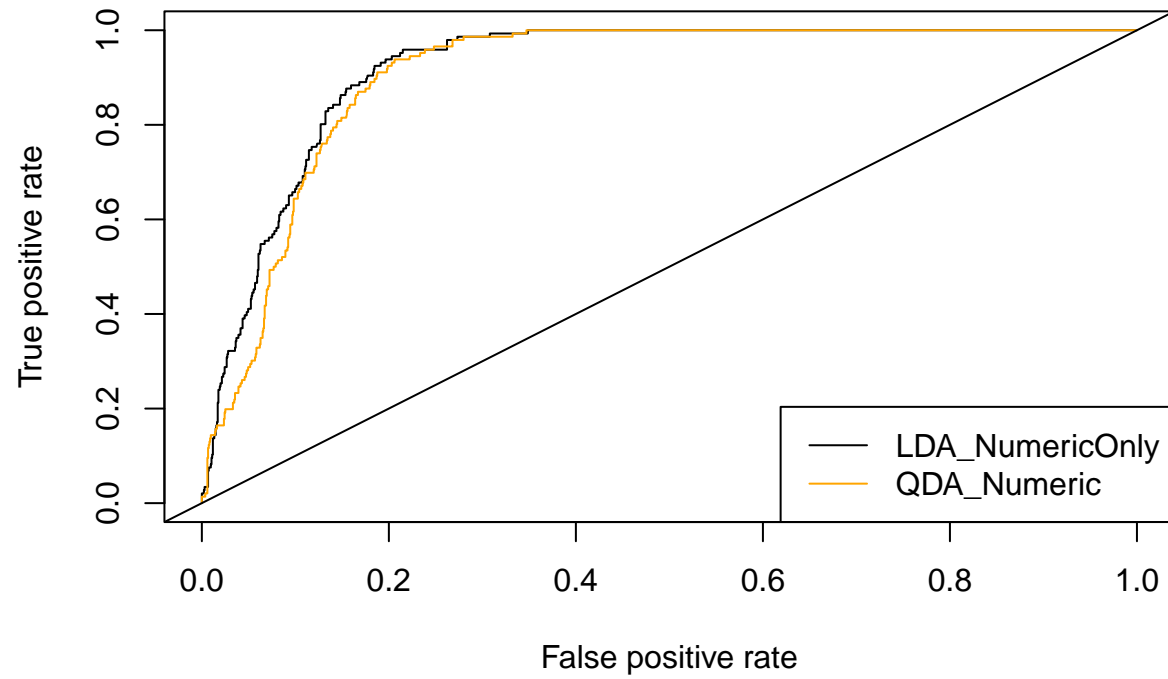
	Criterion	AllInteractions	LogIncomeMortgage	LogAndInteraction
1	AIC	628.246	771.857	620.167
2	BIC	751.456	845.784	786.501
3	Accuracy	0.974	0.955	0.977
4	Sensitivity	0.984	0.972	0.985
5	Specificity	0.877	0.801	0.904

The ROC of models

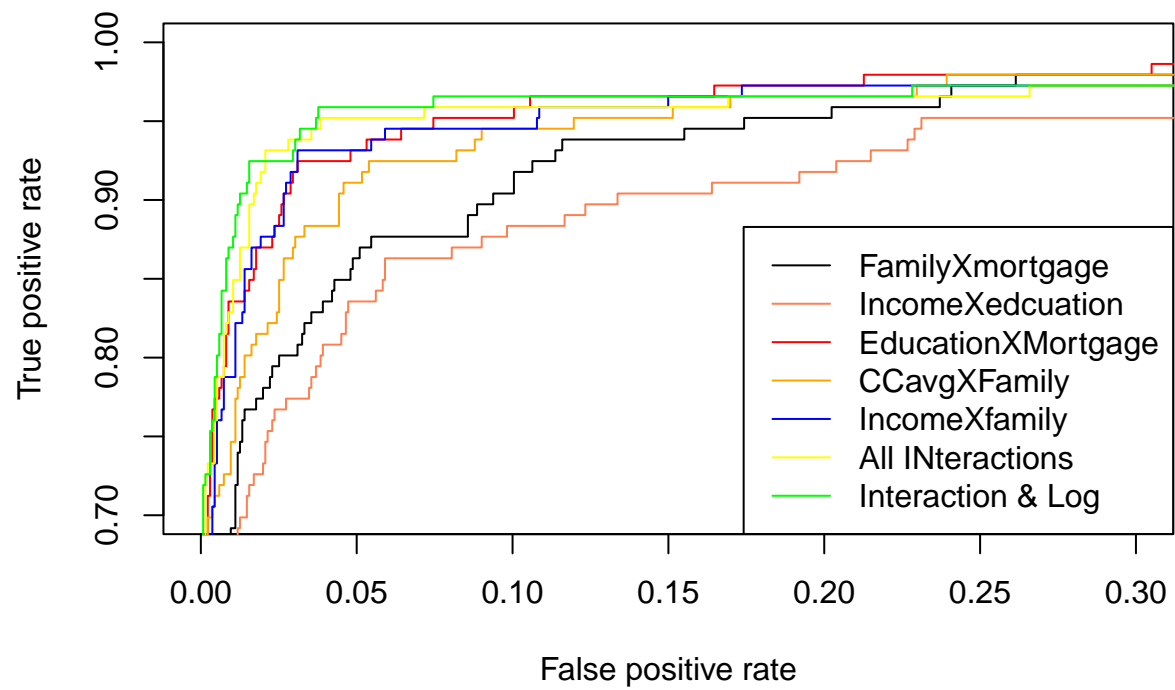
- *Interaction models*



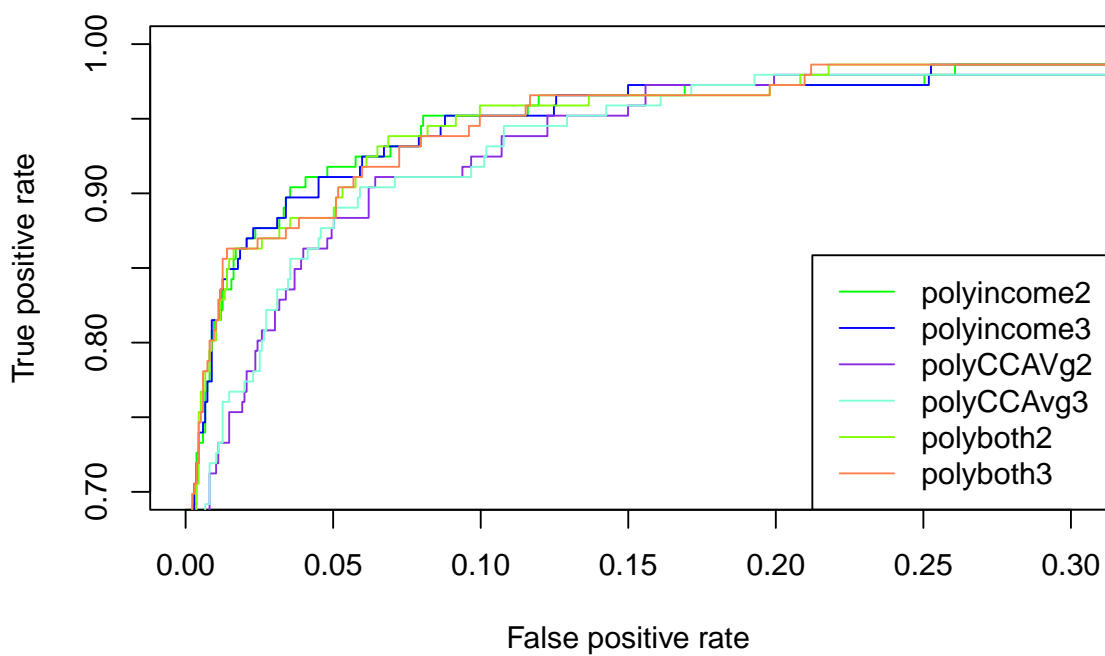
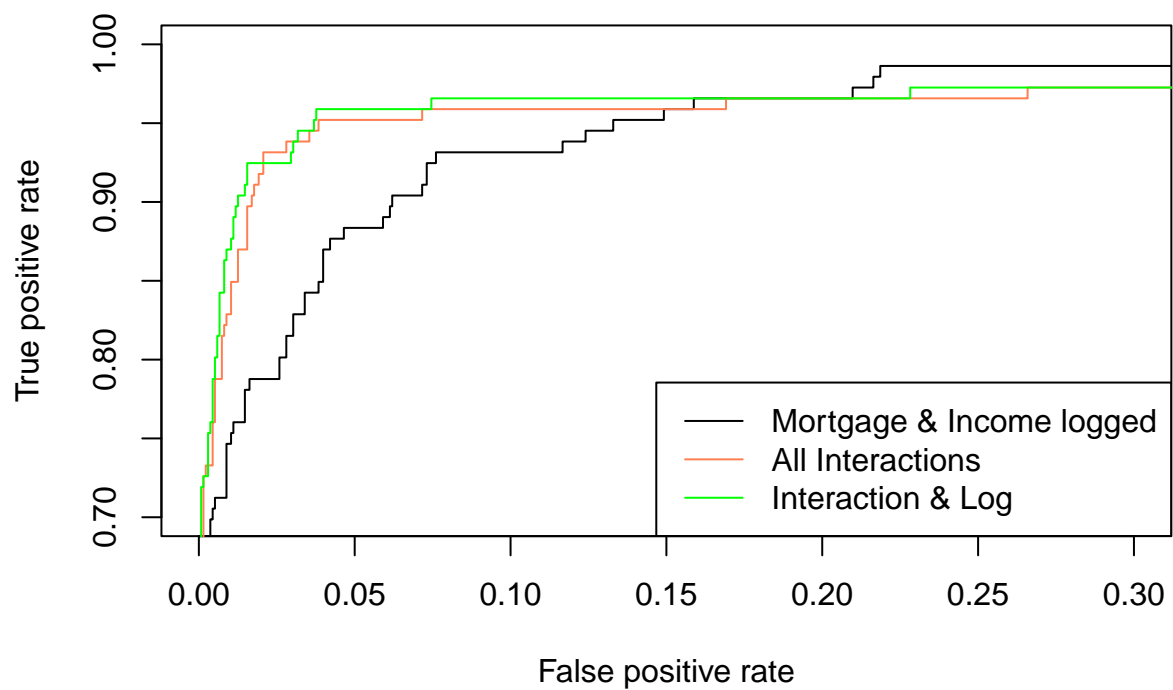
- *LDA & QDA*



- *Best of logs and interactions*

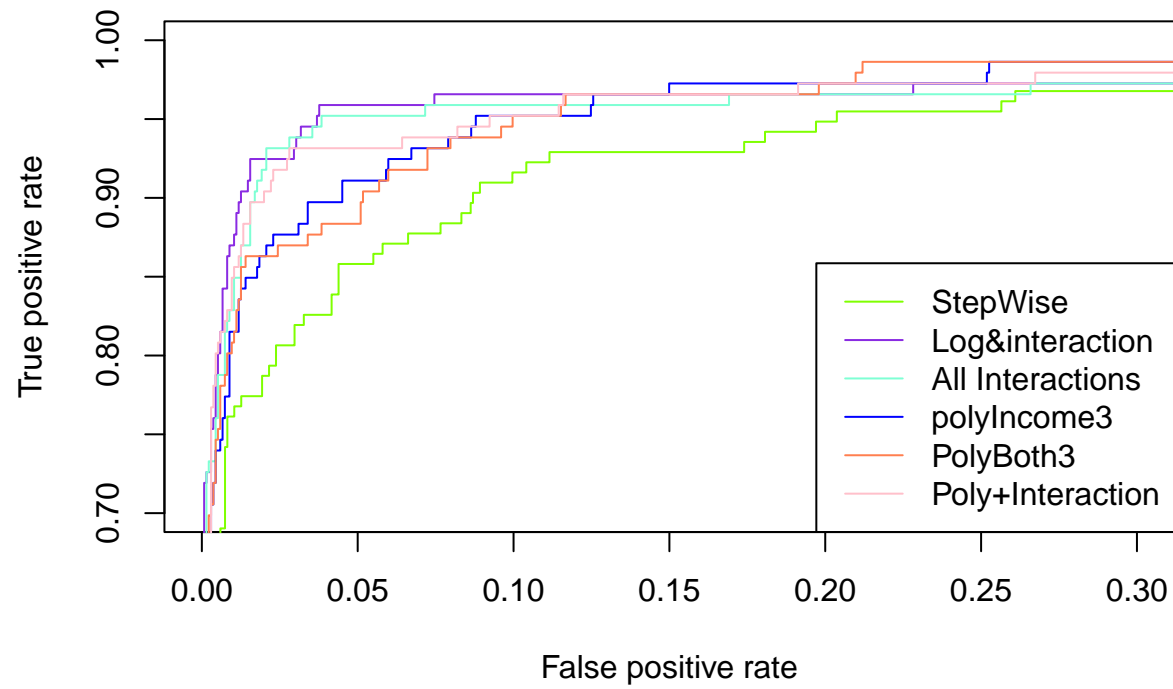


+ Lets get closer look at all the log and interaction , interaction with log and logging alone

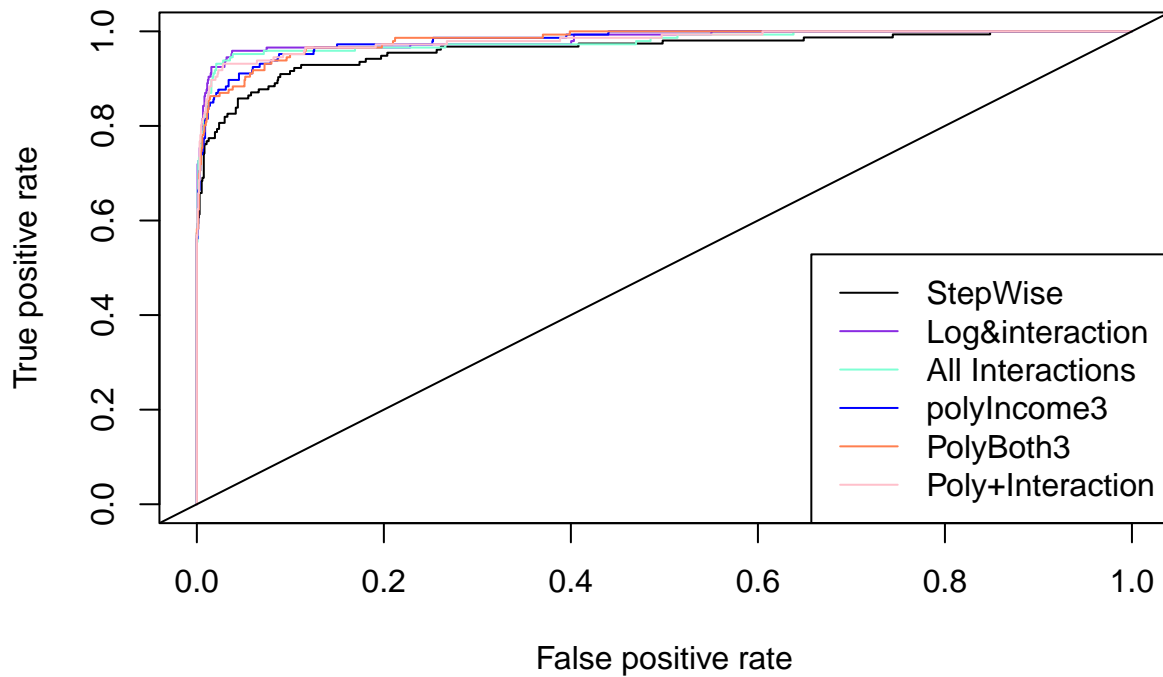


- *Poly models*

- *The best model from the first part vs the best models of the second set*



[1] "Closer view"



Criterion of best models

	Criterion	Step_Wise	LogAndInteraction	QDAnum	polyIncome3
1	AIC	818.074	620.167	0.000	695.393
2	BIC	892.000	786.501	0.000	781.640
3	Accuracy	0.957	0.977	0.875	0.969
4	Sensitivity	0.975	0.985	0.914	0.979
5	Specificity	0.806	0.904	0.514	0.870
	LogAndInteraction.1	AllInteractions	InteractinPolyLog		
1	620.167	628.246	603.767		
2	786.501	751.456	714.656		
3	0.977	0.974	0.974		
4	0.985	0.984	0.982		
5	0.904	0.877	0.897		

These could be for the extra info of criterion of all models

[1] "First basic models"

	Criterion	Full_Model	Step_Wise	Forward_Model	LASSO_model	Intuition	EDA
1	AIC	820.398	818.074	820.398	0.000	836.879	844.481
2	BIC	906.645	892.000	906.645	0.000	892.324	899.926
3	Accuracy	0.957	0.957	0.957	0.956	0.955	0.953
4	Sensitivity	0.975	0.975	0.975	0.975	0.976	0.974
5	Specificity	0.800	0.806	0.800	0.794	0.774	0.774

[1] "Interaction, Logginb and LDA/QDA results"

	Criterion	Famxmortgage	EduXMortgage	CCAvgxFam	IncomeXedu	FamxIncome
1	AIC	844.759	916.902	768.813	619.763	655.916
2	BIC	937.167	984.668	855.061	699.850	748.324
3	Accuracy	0.961	0.953	0.962	0.971	0.973
4	Sensitivity	0.980	0.973	0.977	0.984	0.984
5	Specificity	0.781	0.774	0.822	0.849	0.870
AllInteractions						
1		628.246				
2		751.456				
3		0.974				
4		0.984				
5		0.877				

[1] "LDA/QDA results"

	Criterion	AllInteractions	LogIncomeMortgage	LogAndInteraction	LDA	QDAnum
1	AIC	628.246	771.857	620.167	0.000	0.000
2	BIC	751.456	845.784	786.501	0.000	0.000
3	Accuracy	0.974	0.955	0.977	0.896	0.875
4	Sensitivity	0.984	0.972	0.985	0.945	0.914
5	Specificity	0.877	0.801	0.904	0.445	0.514

[1] "Polynomials results"

	Criterion	polyIncome2	polyIncome3	polyCCAvg2	polyCCAvg3	PolyBoth2	PolyBoth3
1	AIC	694.203	695.393	784.087	780.157	661.317	663.552
2	BIC	774.290	781.640	870.334	872.564	747.564	762.120
3	Accuracy	0.968	0.969	0.955	0.957	0.965	0.966
4	Sensitivity	0.979	0.979	0.971	0.975	0.976	0.977
5	Specificity	0.870	0.870	0.808	0.795	0.863	0.863
ExtremexModel InteractinPolyLog							
1		672.042	603.767				
2		795.252	714.656				
3		0.966	0.974				
4		0.978	0.982				
5		0.856	0.897				