

Submission 2

Nicholas Sherman, Kishor Kumar, Sridhar, Adjoa Adanledji, and Andrew Smith

11/18/2019

Data Cleaning

```
cData$CustomerId = as.character(cData$CustomerId)
cData$CreditScore = as.integer(cData$CreditScore)
cData$CreditScoreBins = cut(x = cData$CreditScore, breaks = c(349, seq(360,850,10)))
cData$Geography = as.factor(cData$Geography)
cData$Gender = as.factor(cData$Gender)
cData$Age = as.integer(cData$Age)
cData$AgeBins = cut(x = cData$Age, breaks = c(17, seq(25,100,5)))
cData$TenureFactor = factor(as.character(cData$Tenure), levels = c("0","1","2","3","4","5","6","7","8"))
cData$Tenure = as.factor(as.character(cData$Tenure))
cData$NumOfProducts = as.integer(cData$NumOfProducts)
cData$HasCrCard = as.factor(as.character(cData$HasCrCard))
levels(cData$HasCrCard) = c("inactive","active")
cData$IsActiveMember = as.factor(as.character(cData$IsActiveMember)) #assuming 1 = member, 0 = not
levels(cData$IsActiveMember) = c("inactive","active")
cData$Exited = as.factor(as.character(cData$Exited))
levels(cData$Exited) = c("retained","exited") #0 = retained, 1 = exited
cData = cData %>% select(-RowNumber)
```

Data Exploration

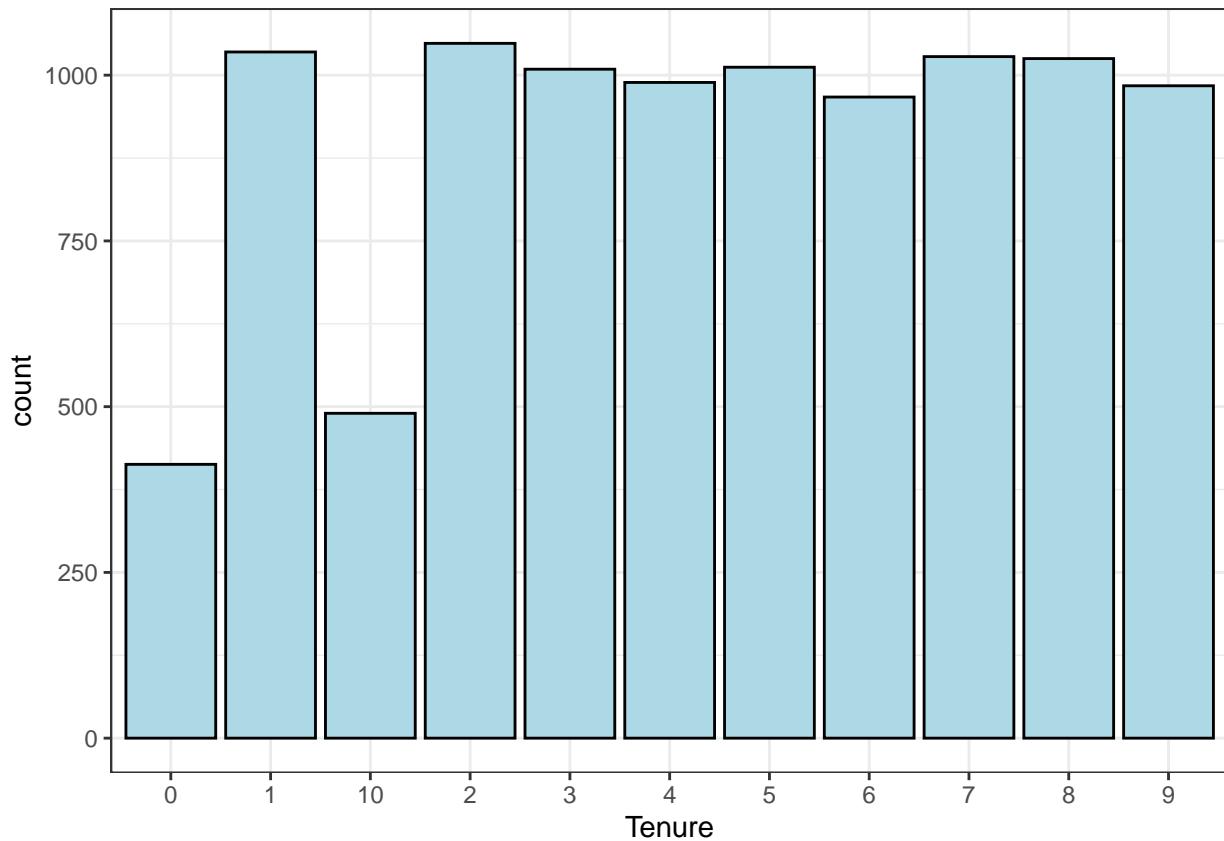
Difference between Active and Inactive

```
##           IsActiveMember
## Exited      inactive active
##   retained     3547    4416
##   exited       1302     735

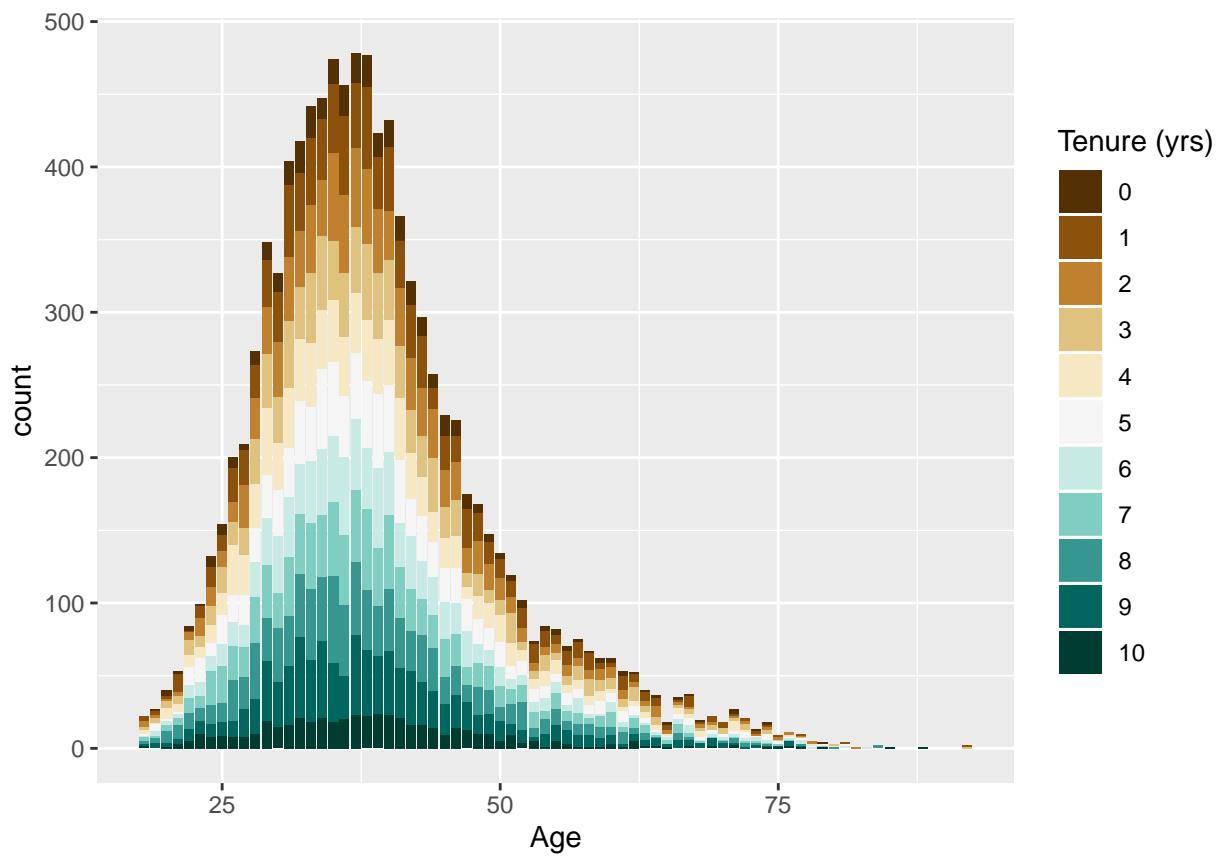
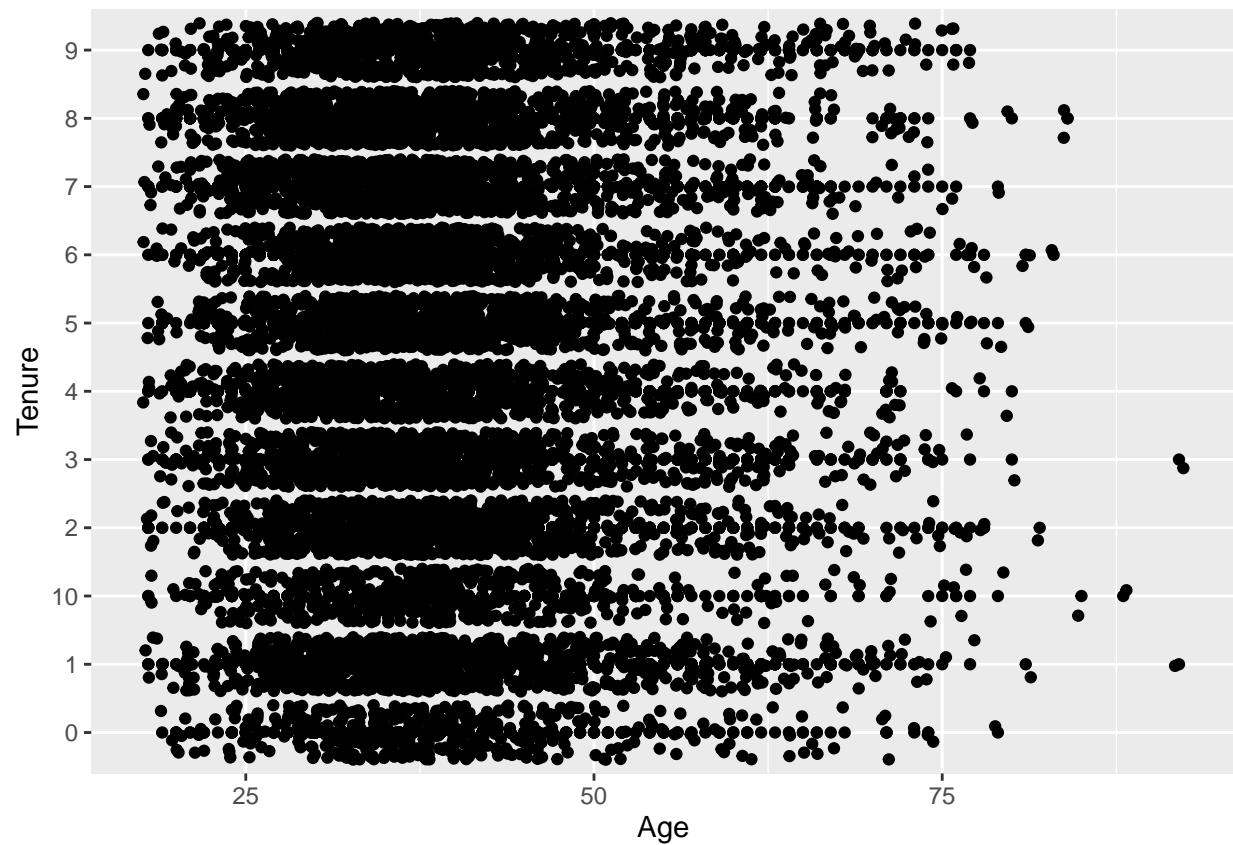
## [1] 10000
## [1] 2932
## [1] 7068
```

Possibility of duplicate values.

```
## Warning: Ignoring unknown parameters: binwidth, bins, pad
```

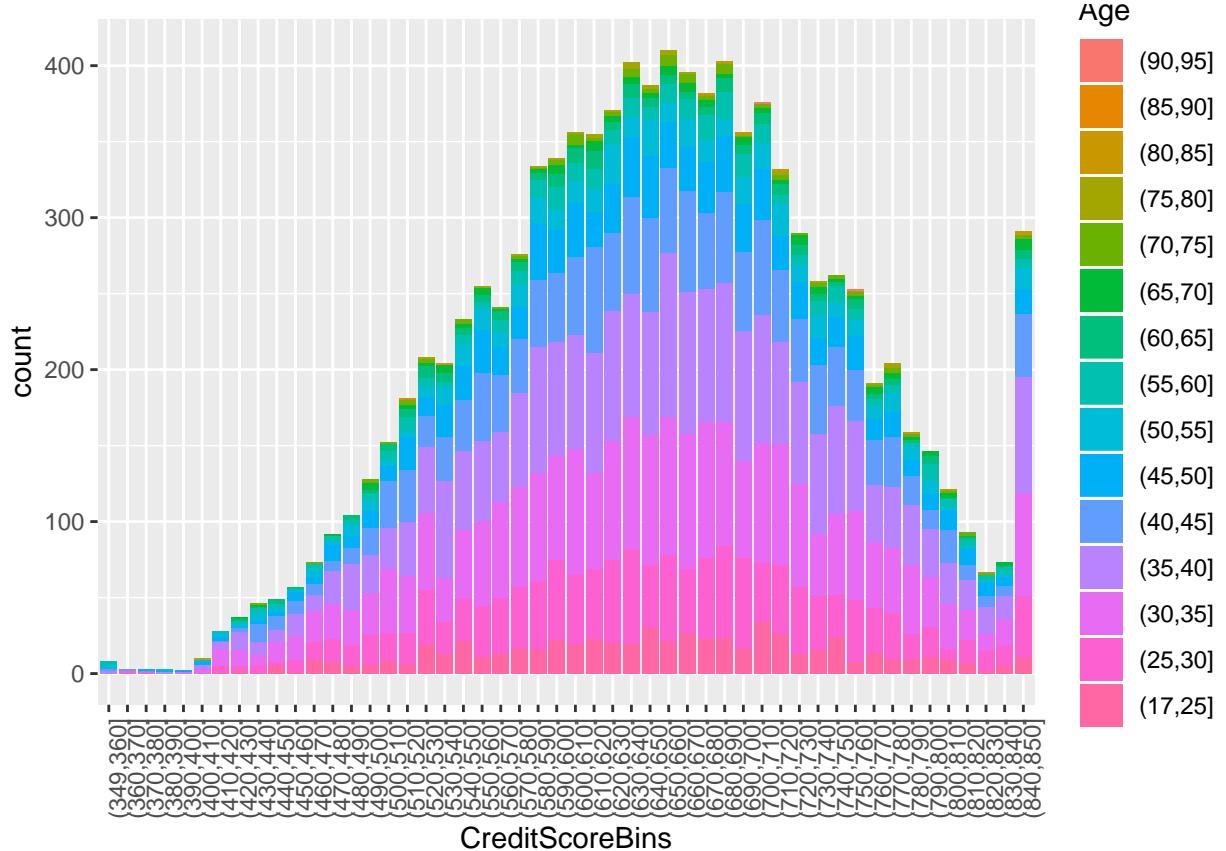


How do age and tenure compare?

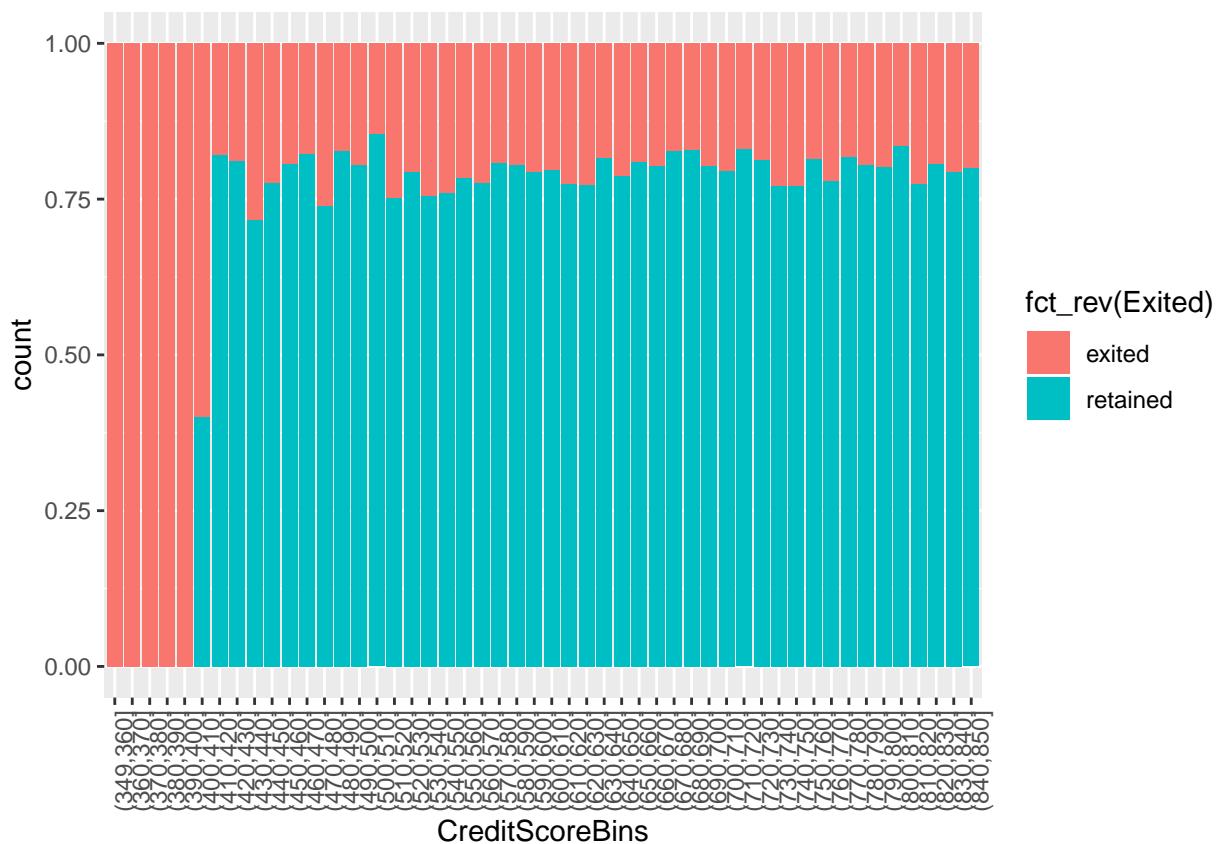


Looking into age groups

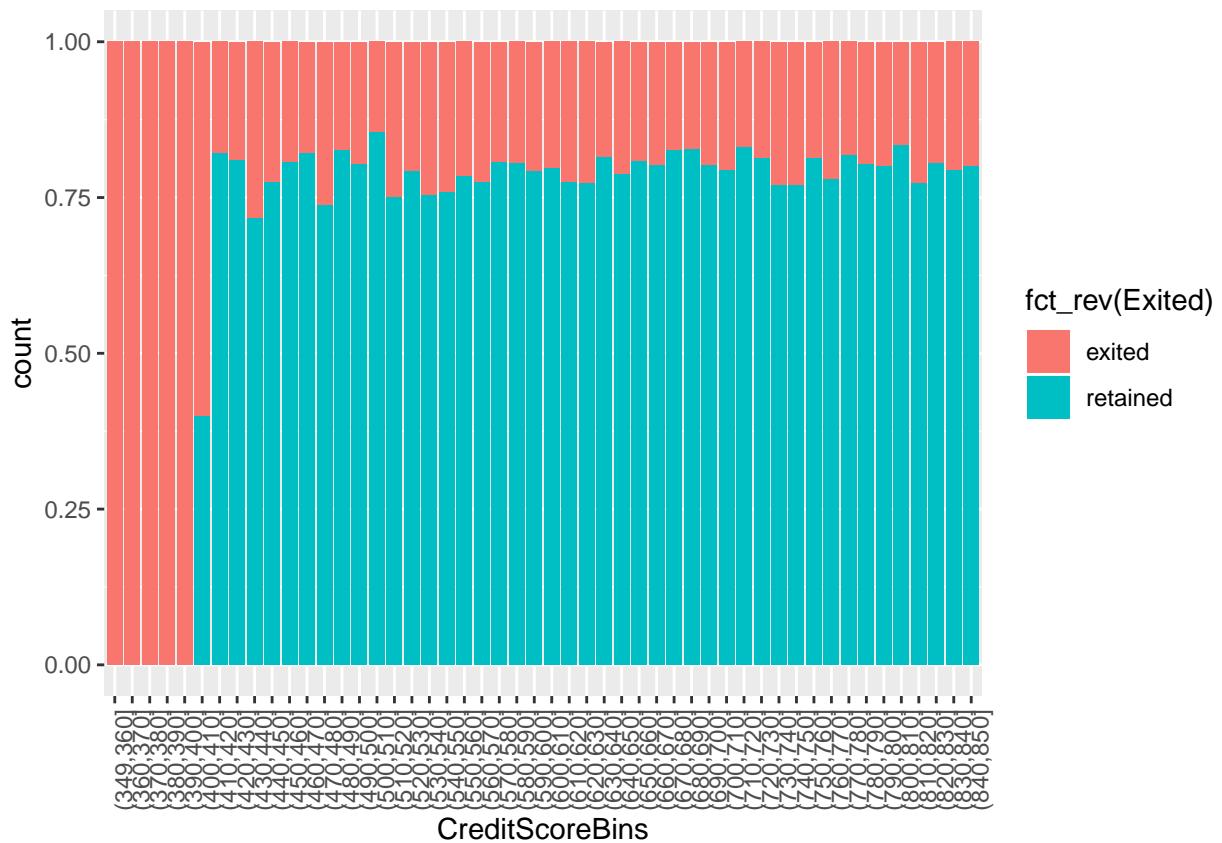
```
## 
##   (17,25]  (25,30]  (30,35]  (35,40]  (40,45]  (45,50]  (50,55]  (55,60]
##     611      1357     2185     2266     1470      850      461      336
##   (60,65]  (65,70]  (70,75]  (75,80]  (80,85]  (85,90]  (90,95]  (95,100]
##     200      131       88       33        9        1        2        0
```



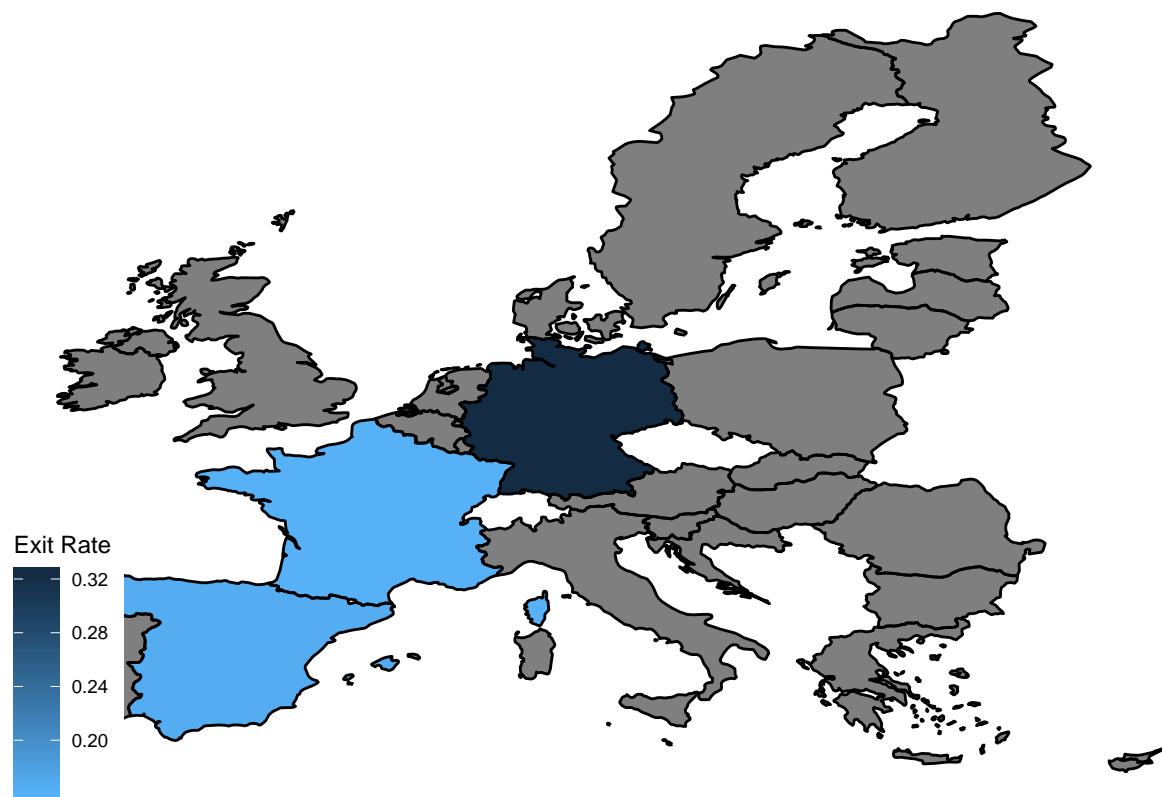
Data Exploration Response Variable: Exited



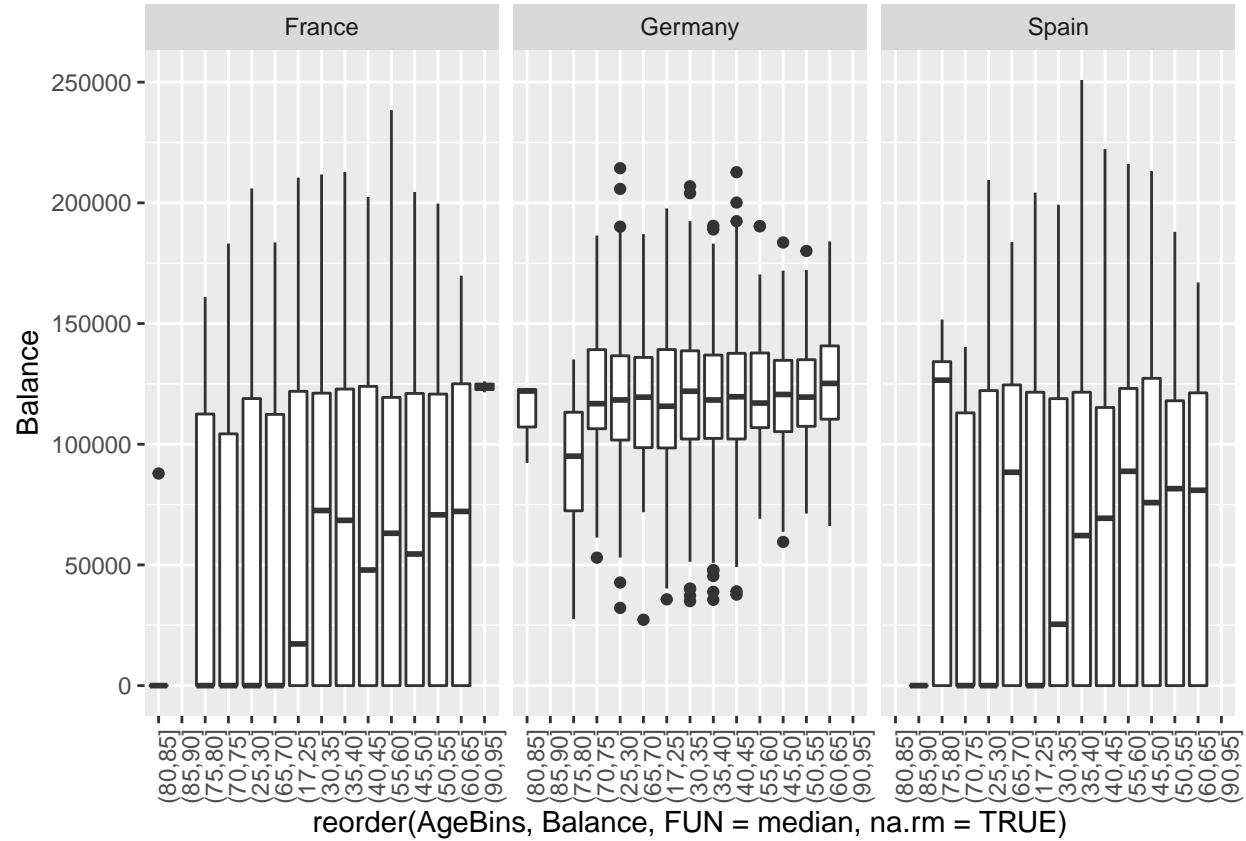
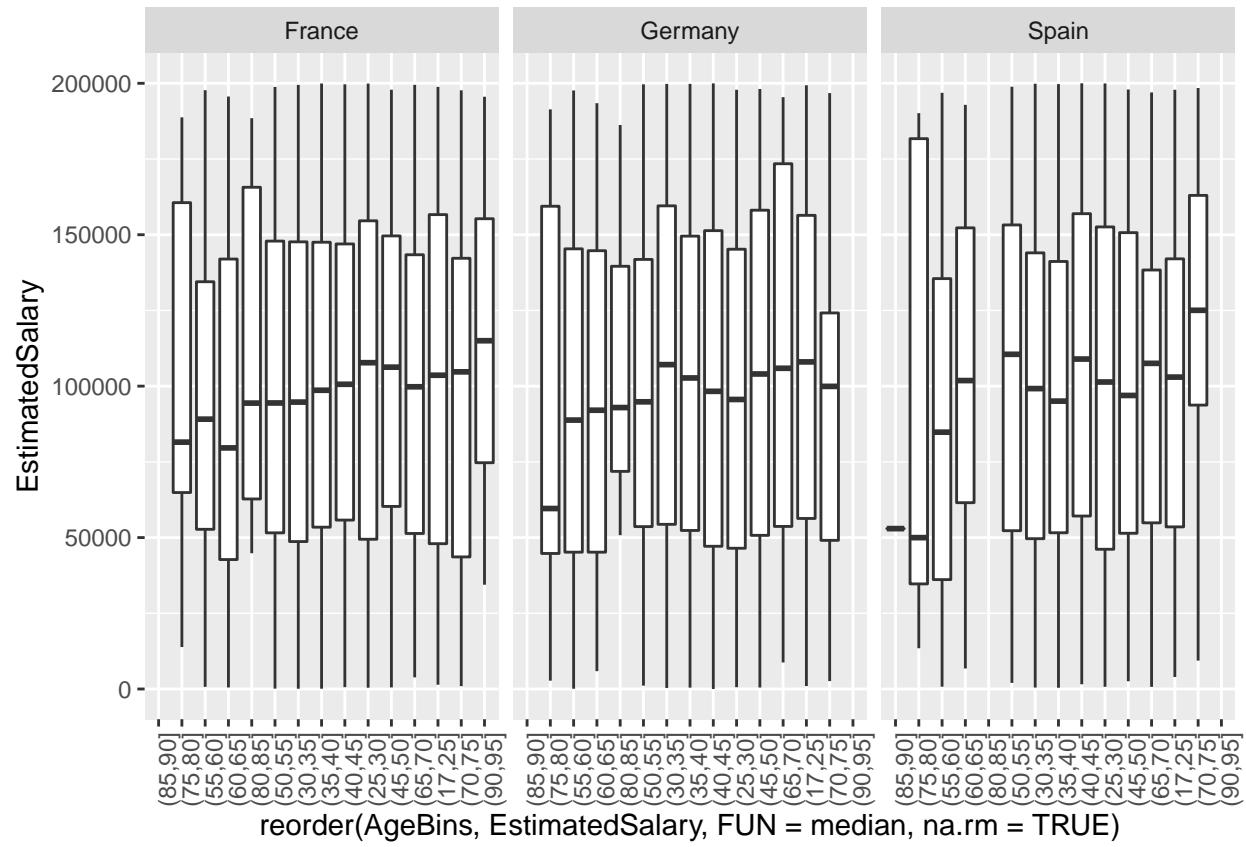
A low credit score maybe implies the bank didn't even accept them, or high fees and the customer didn't want to stay.



Geography

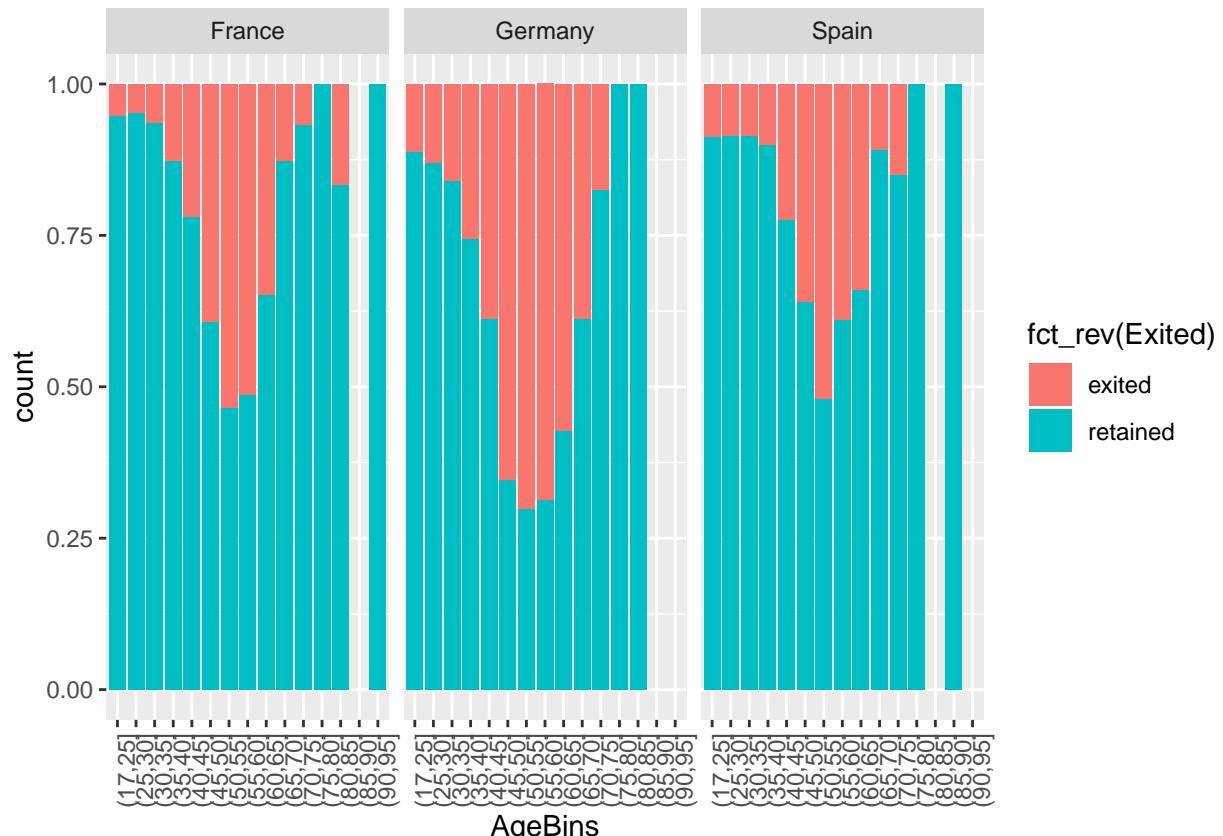
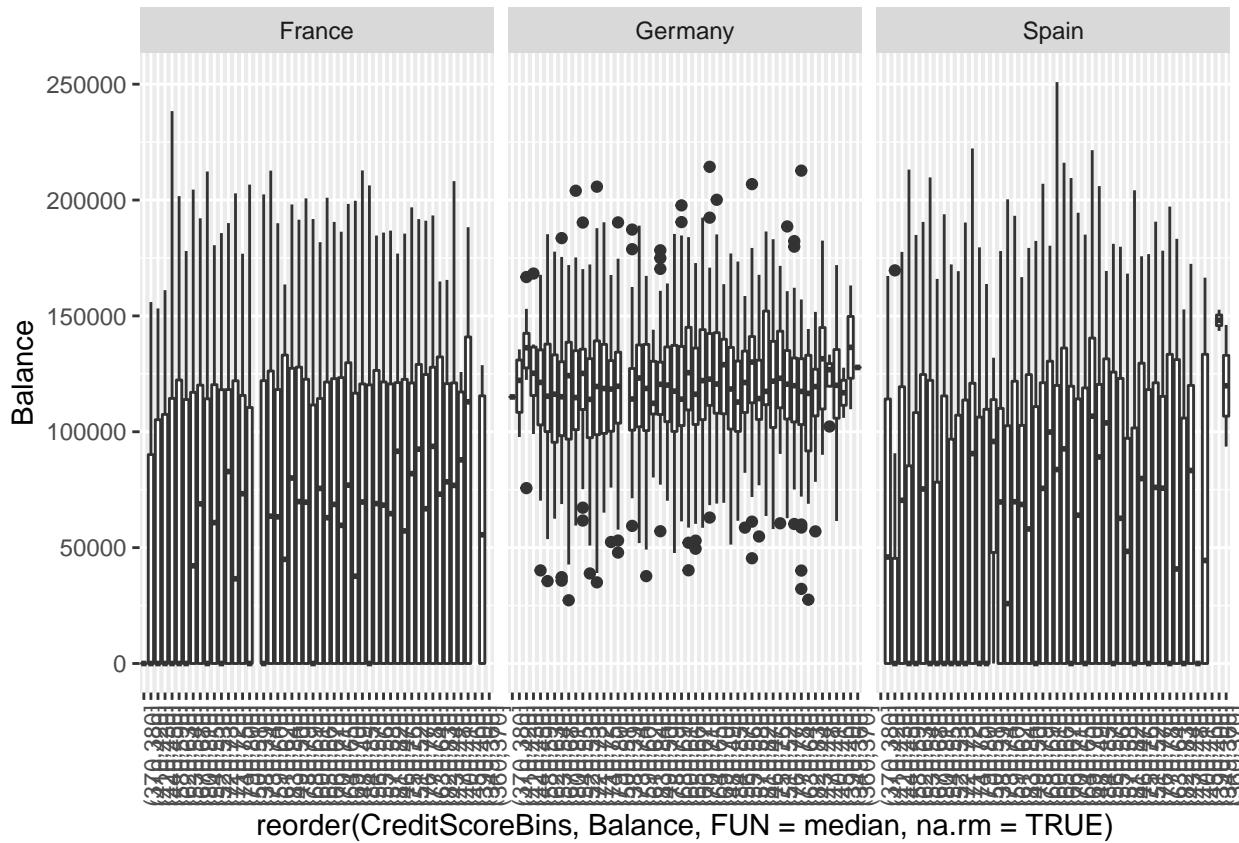


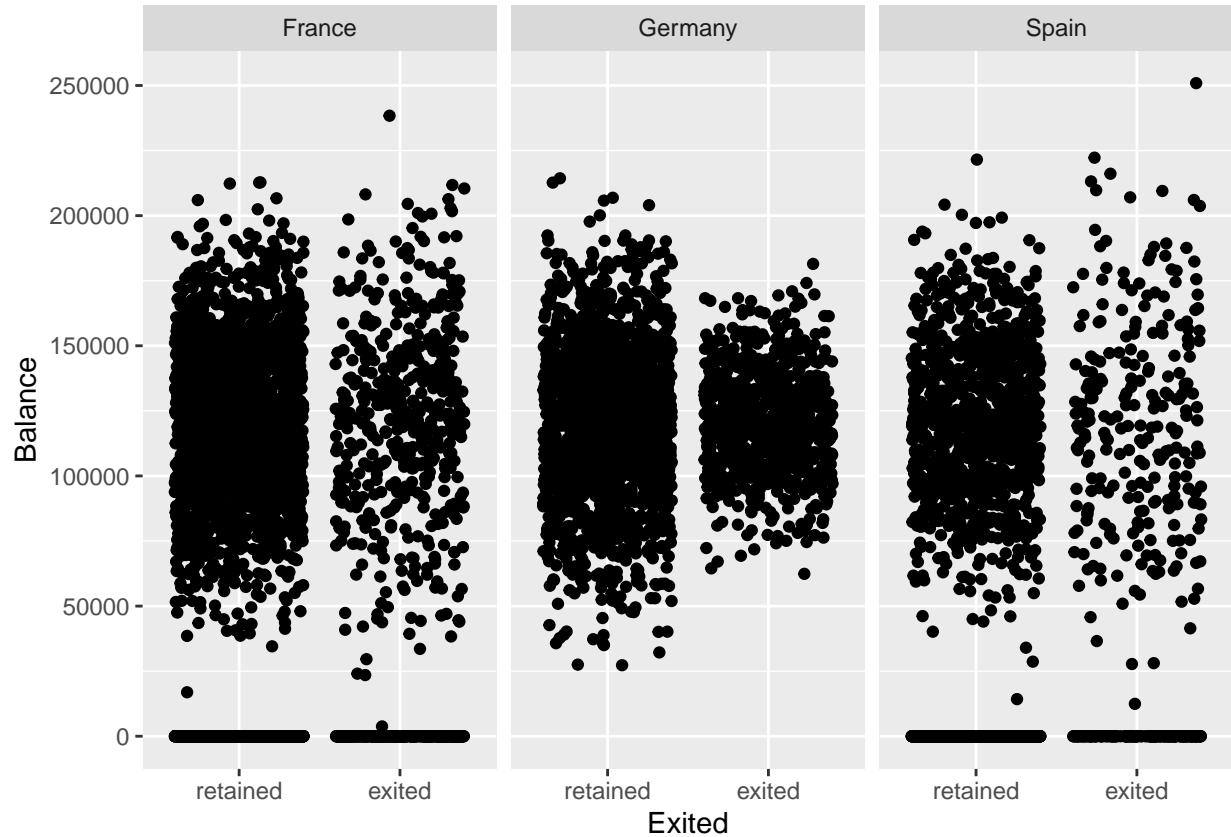
Germany

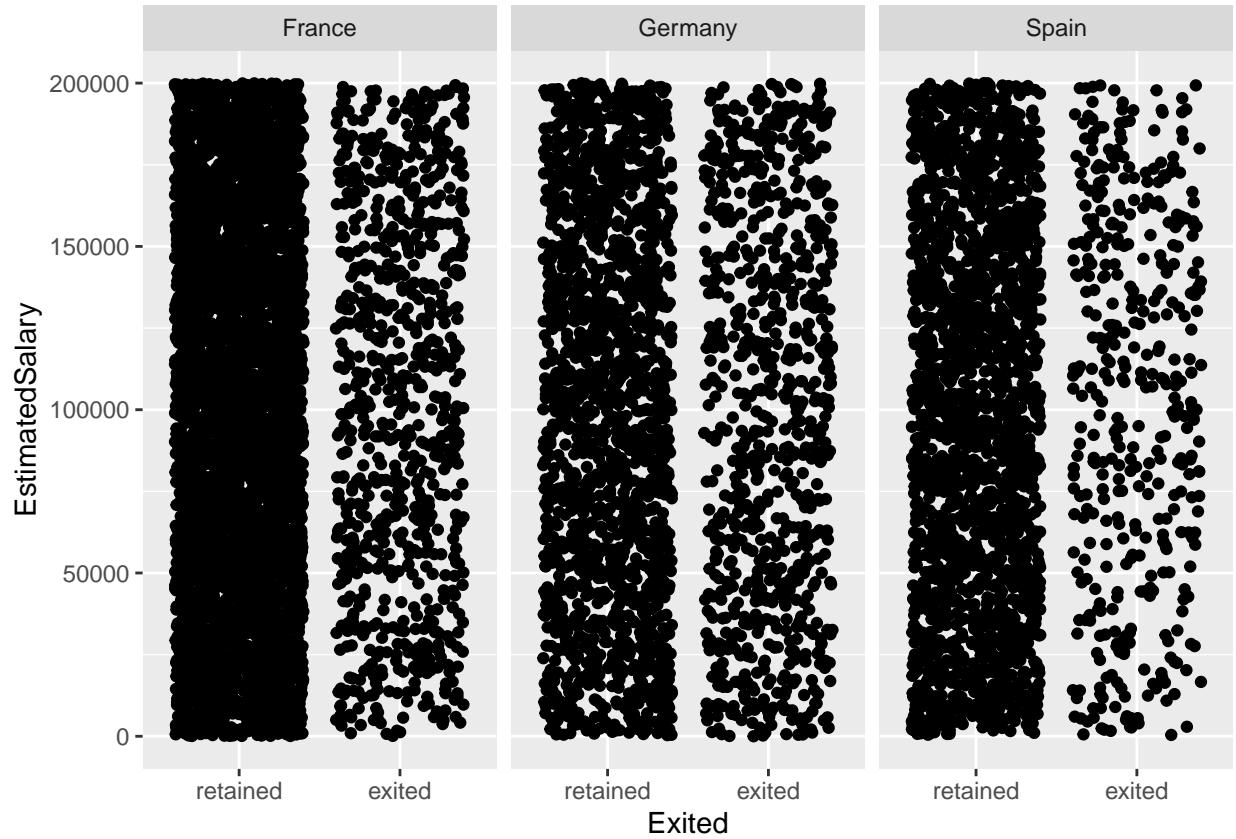


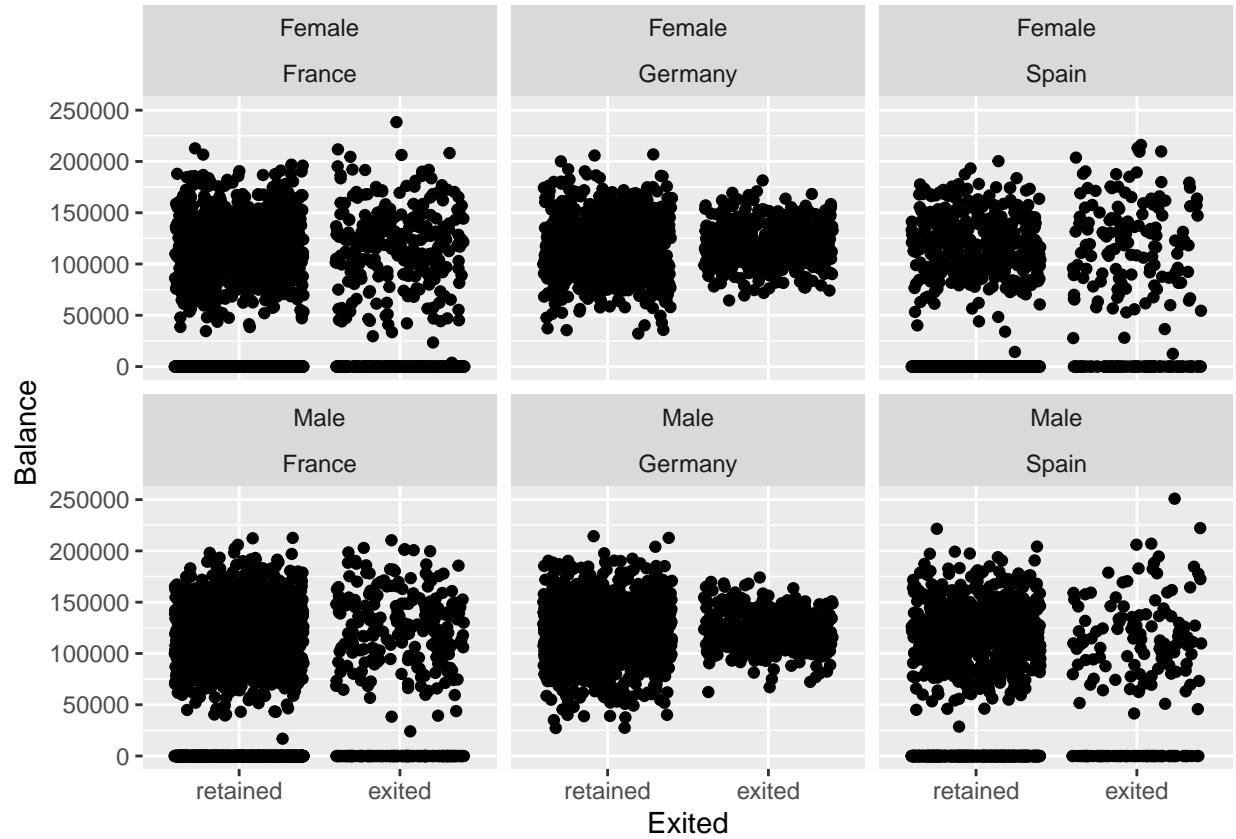
A lot of outliers are present in Germany compared to the other two countries. Same is the case if we create a boxplot with respect to balance and credit score bins.

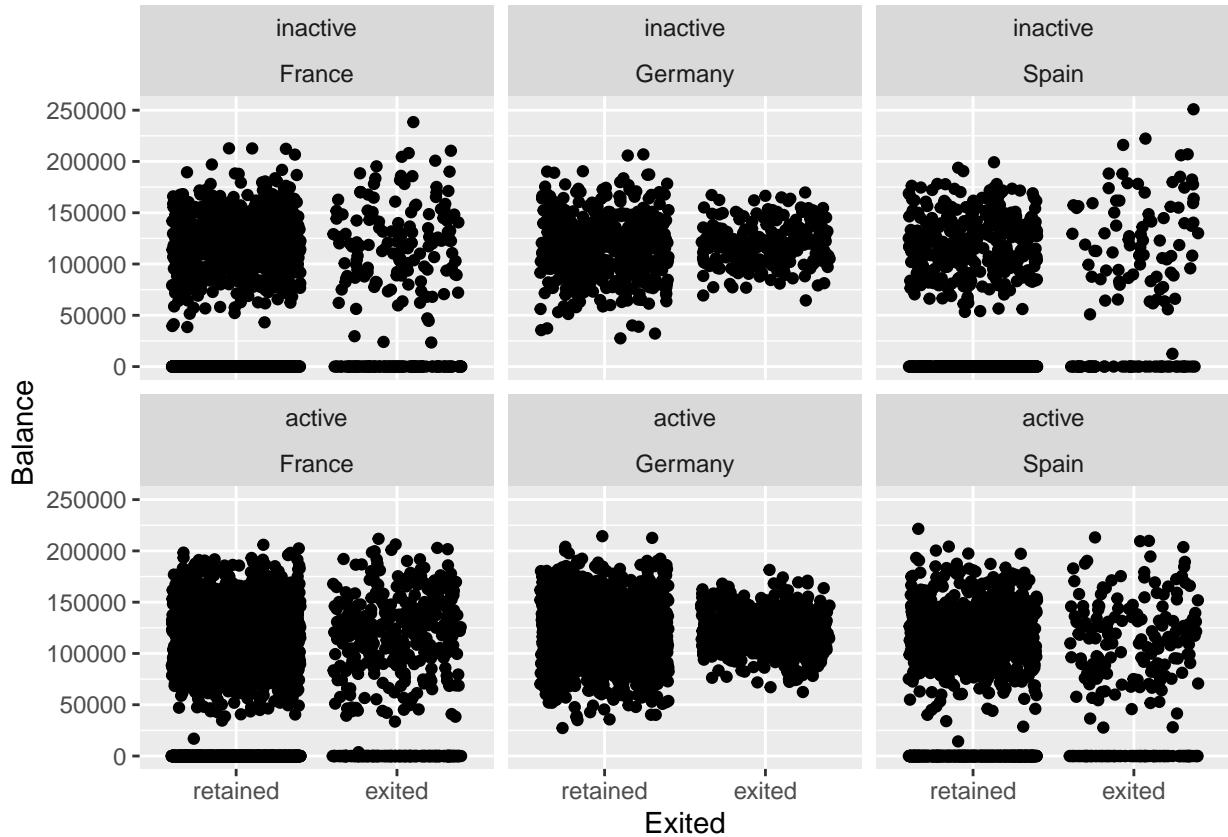
More Geography



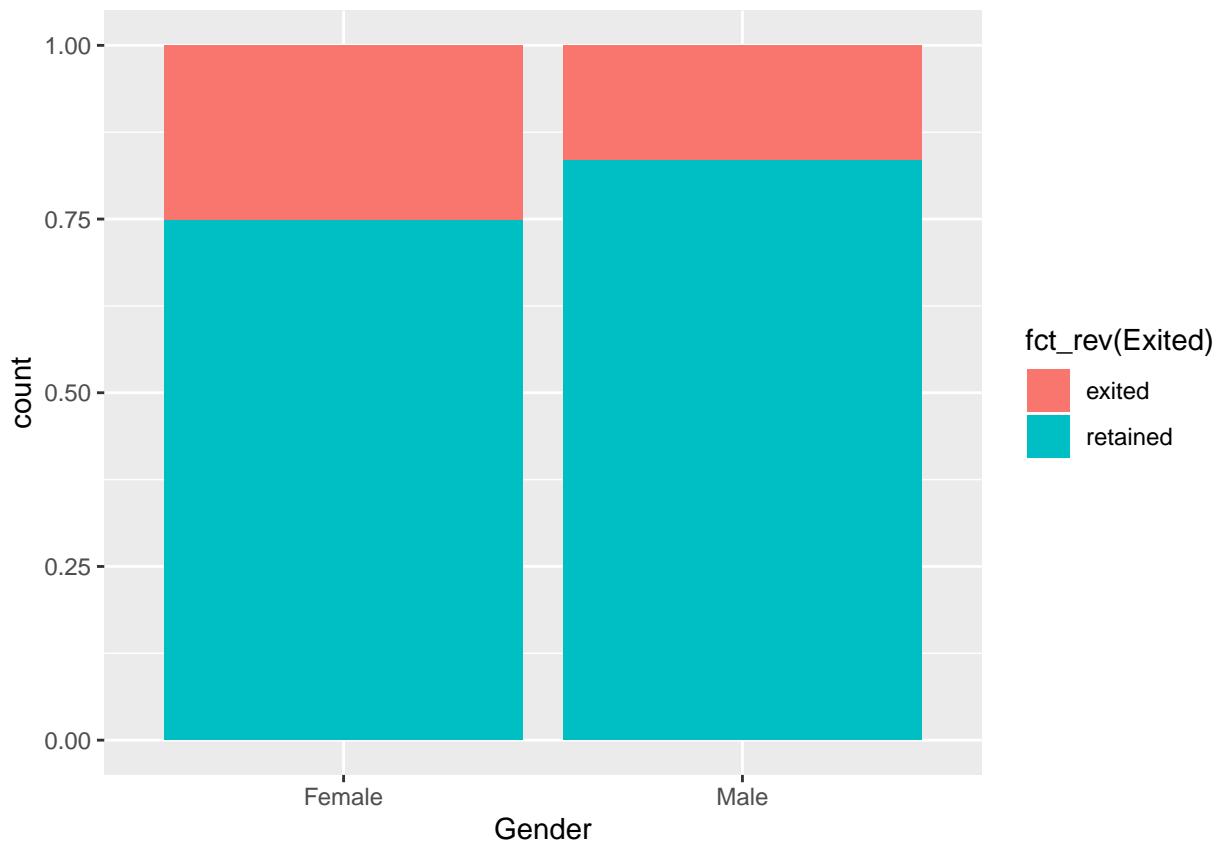




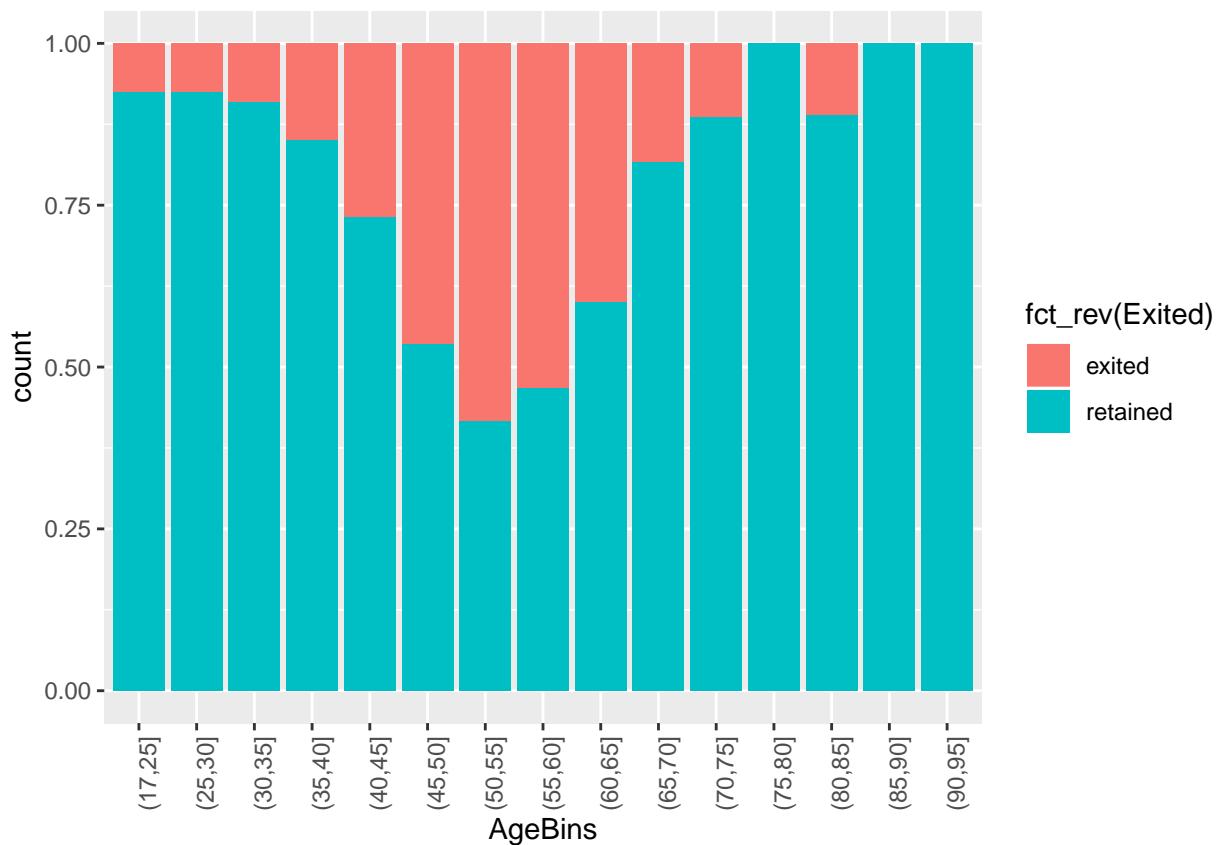




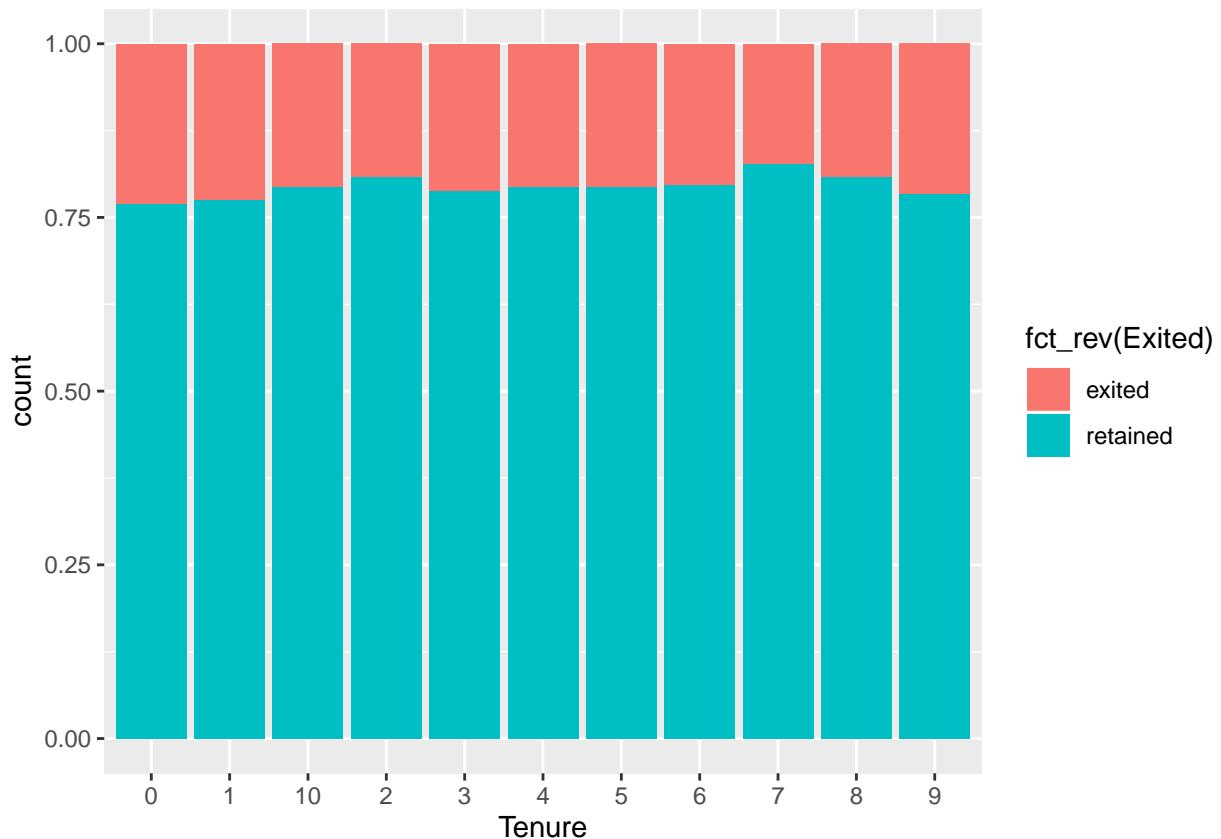
There doesn't seem to be much explanation with respect to the people who exited and their estimated salary. However, the spread of "balance" of the people who exited in Germany seems to be less and is hovering around 75000 to 175000, and this seems to be true for both the genders. But the number of exited people in Germany seems to be considerably higher for customers with Credit Card than without it comparing to the numbers in the other two countries #### GENDER



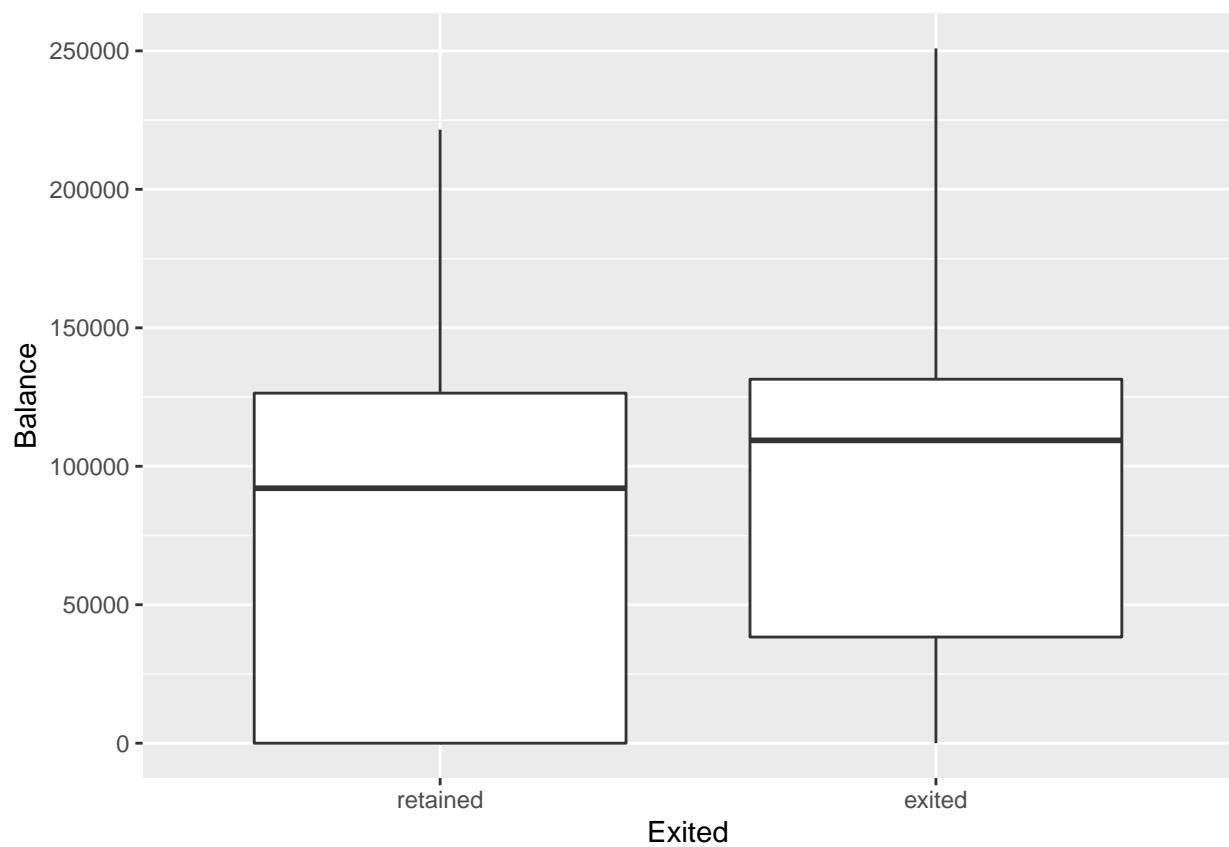
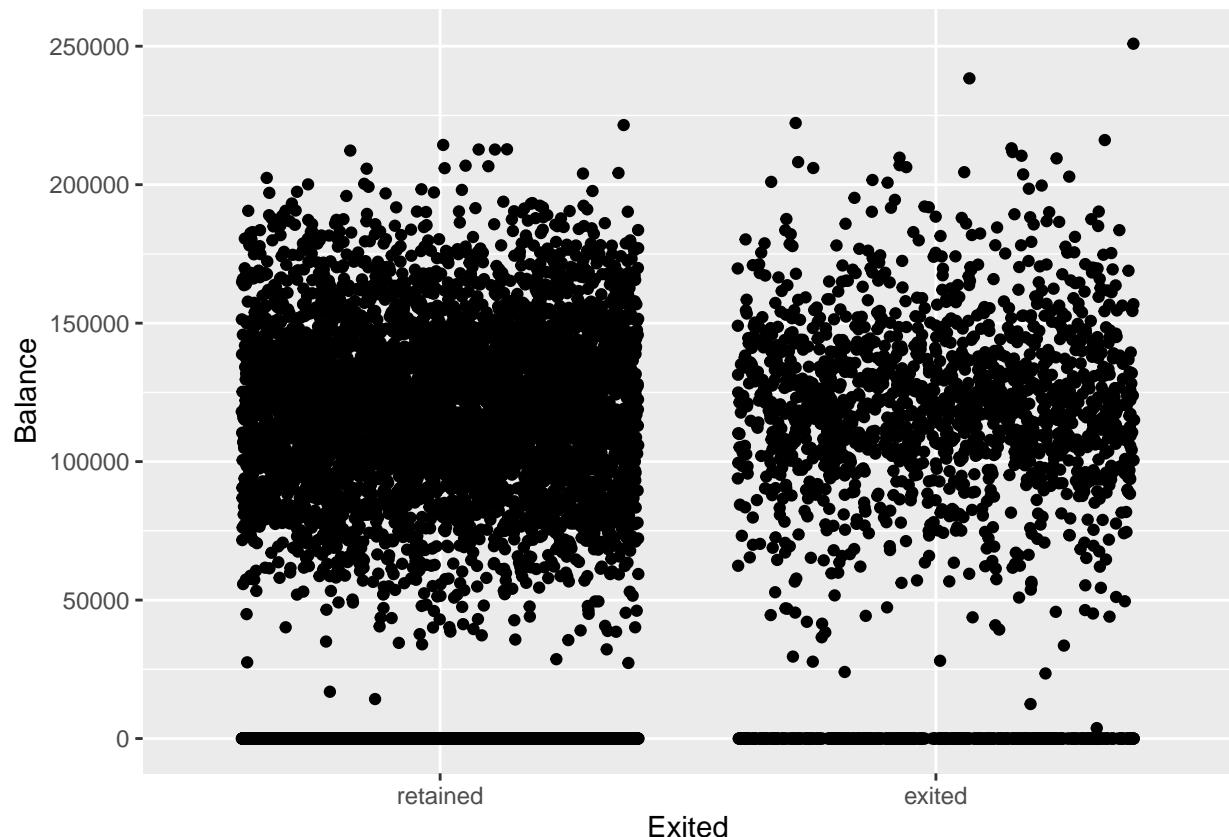
AGE



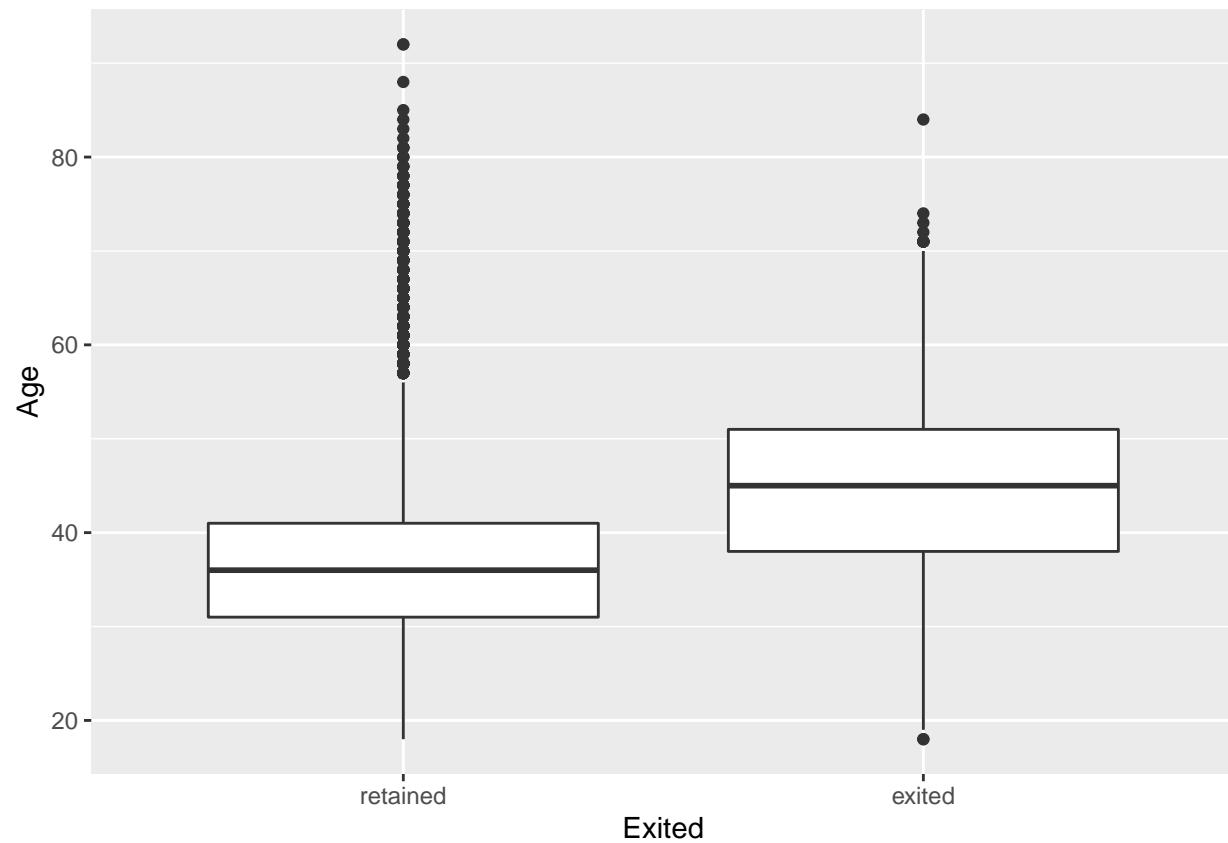
TENURE

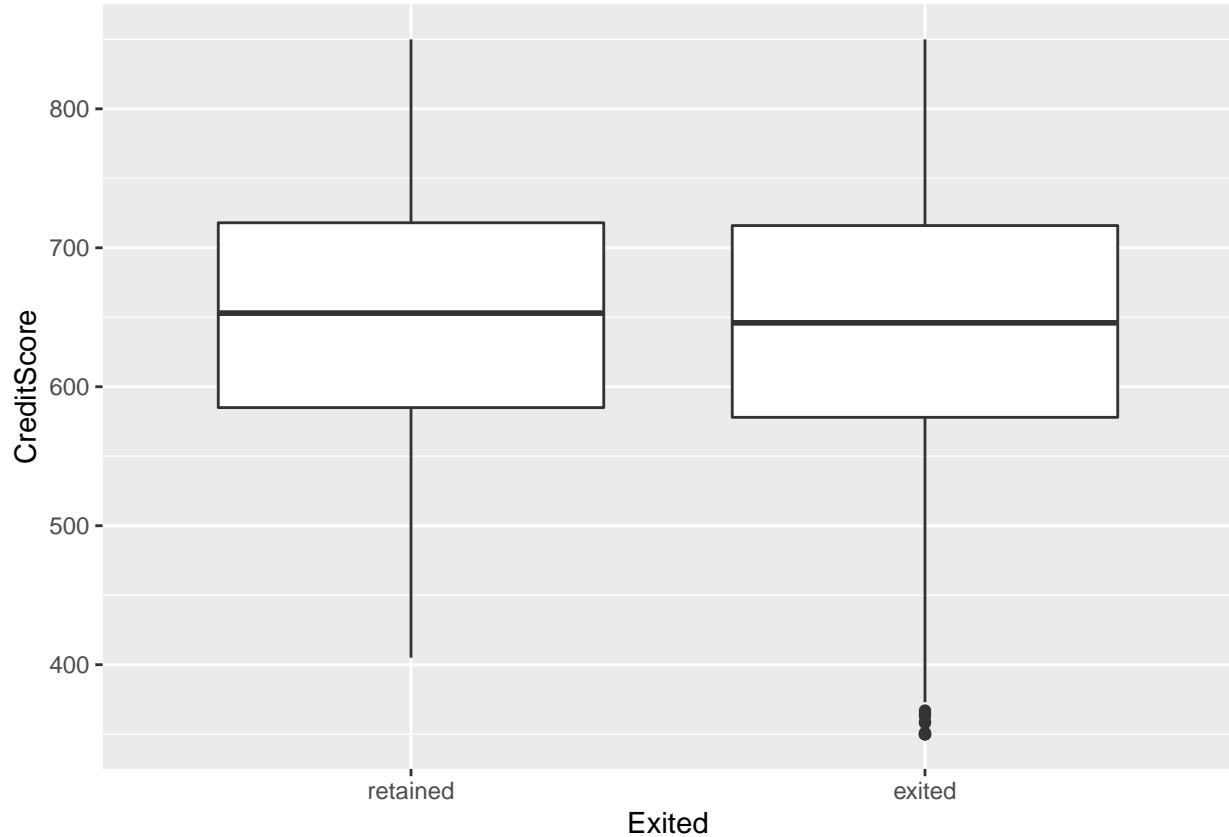


BALANCE

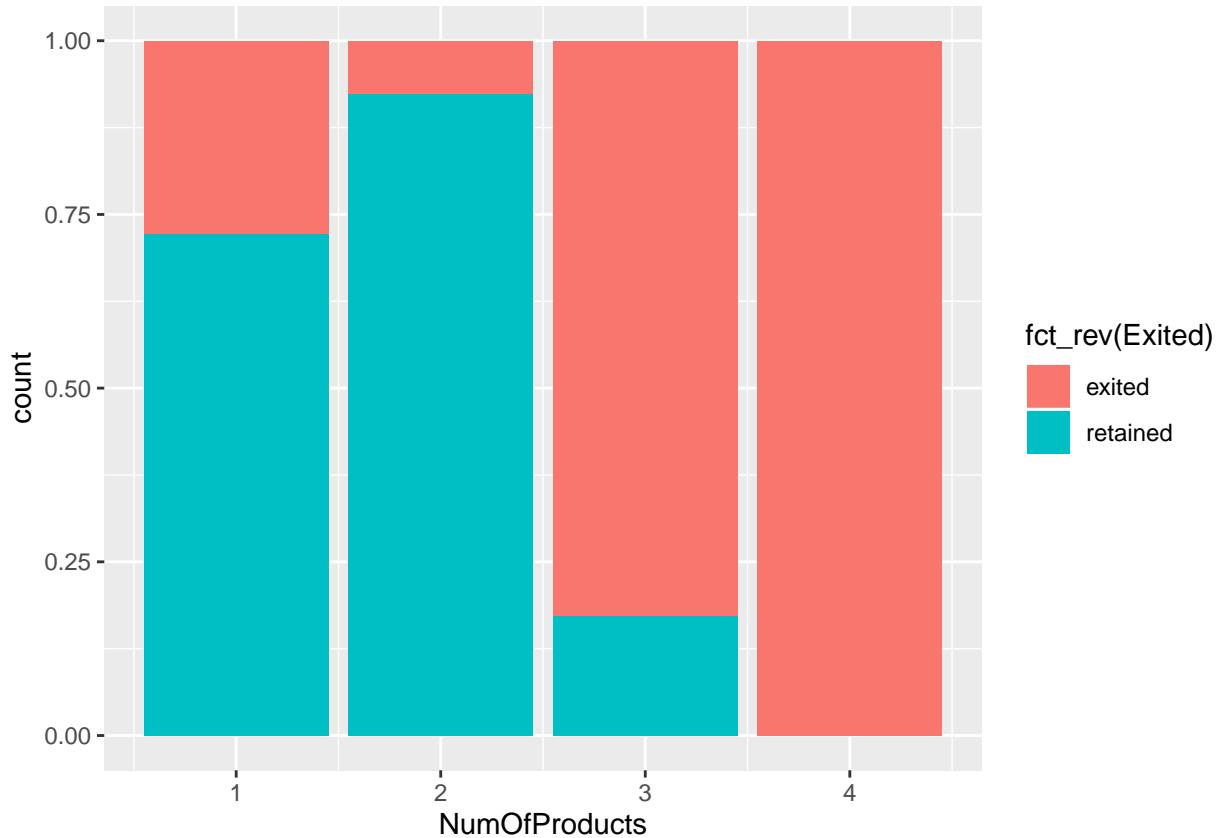


Low balance has a better probability of being retained? We should be looking at box plots for age, tenure, and credit score.





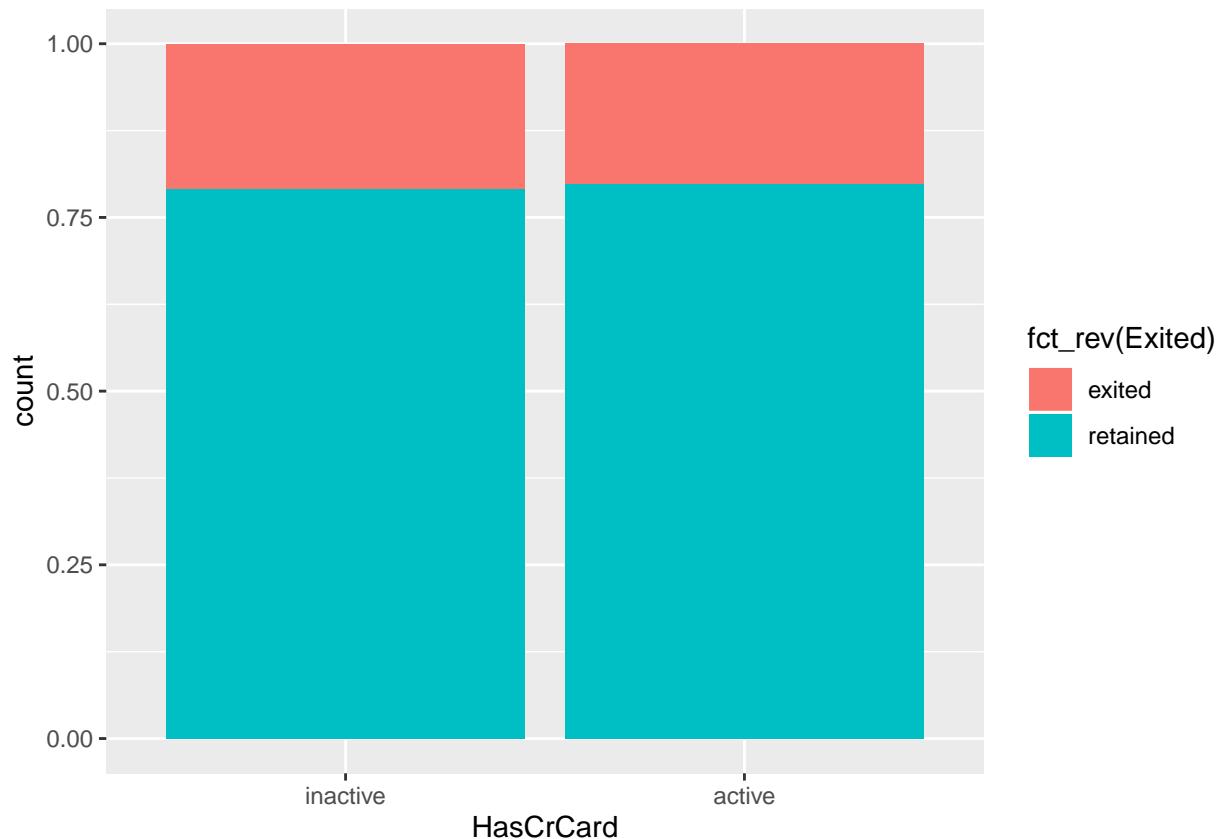
NUMBER OF PRODUCTS



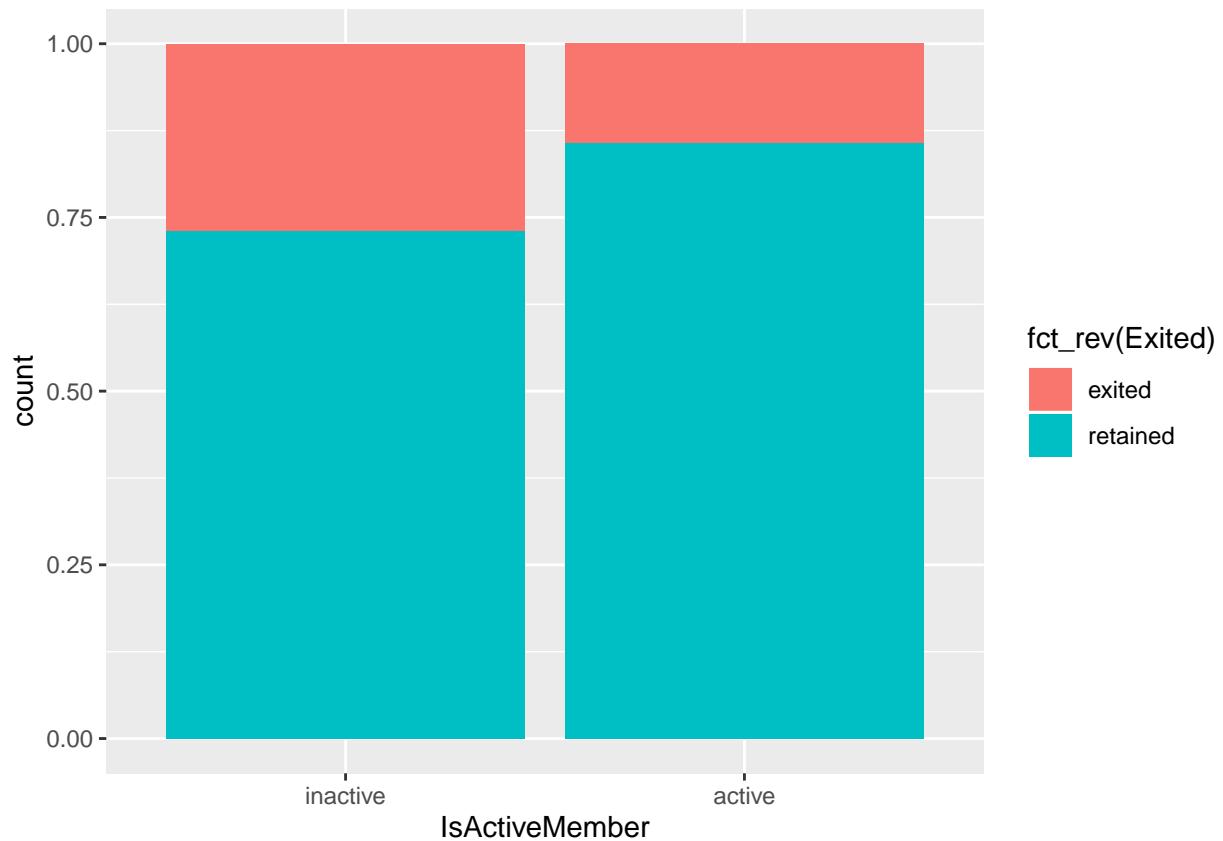
```
##  
##      1     2     3     4  
## 5084 4590  266   60
```

Groups are imbalanced, but still shows a high purity for some num of products based on this, it's more logical to treat NumOfProducts as a factor variable.

HAS CREDIT CARD

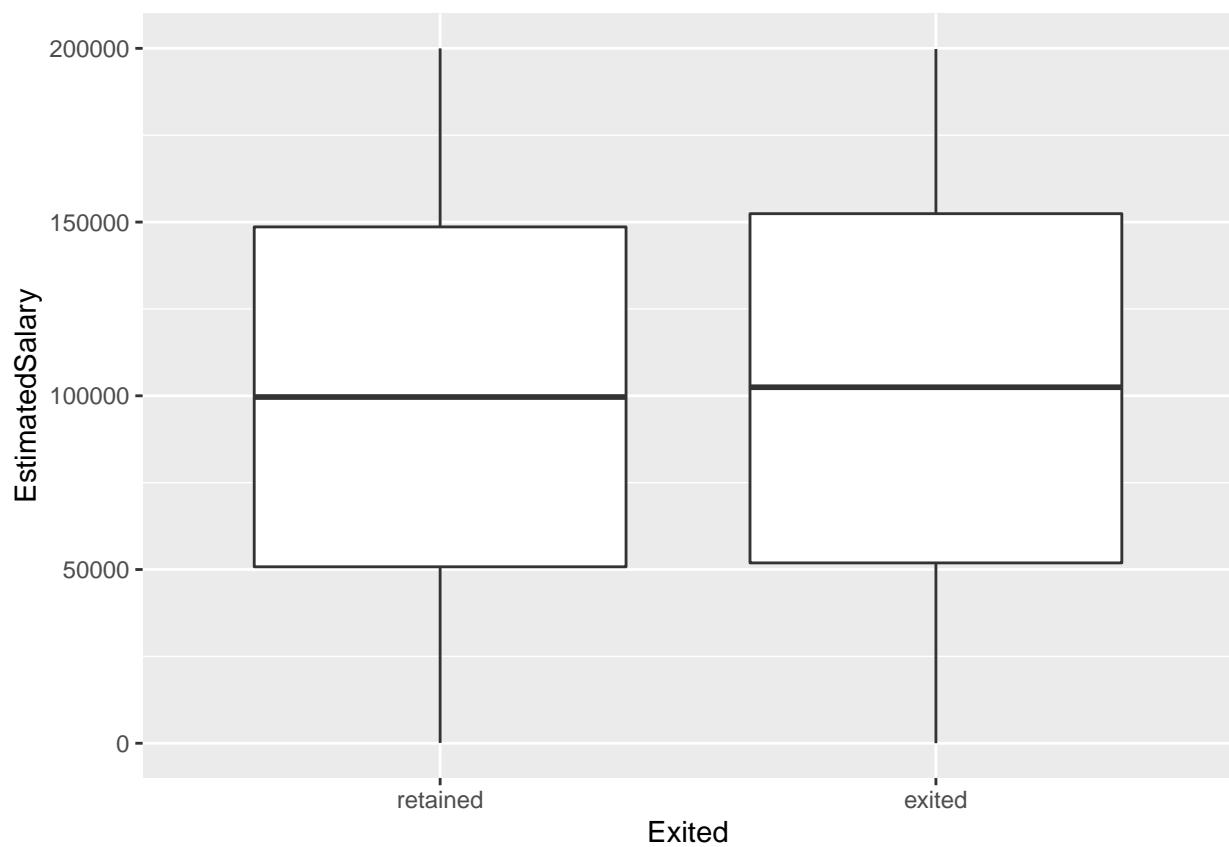
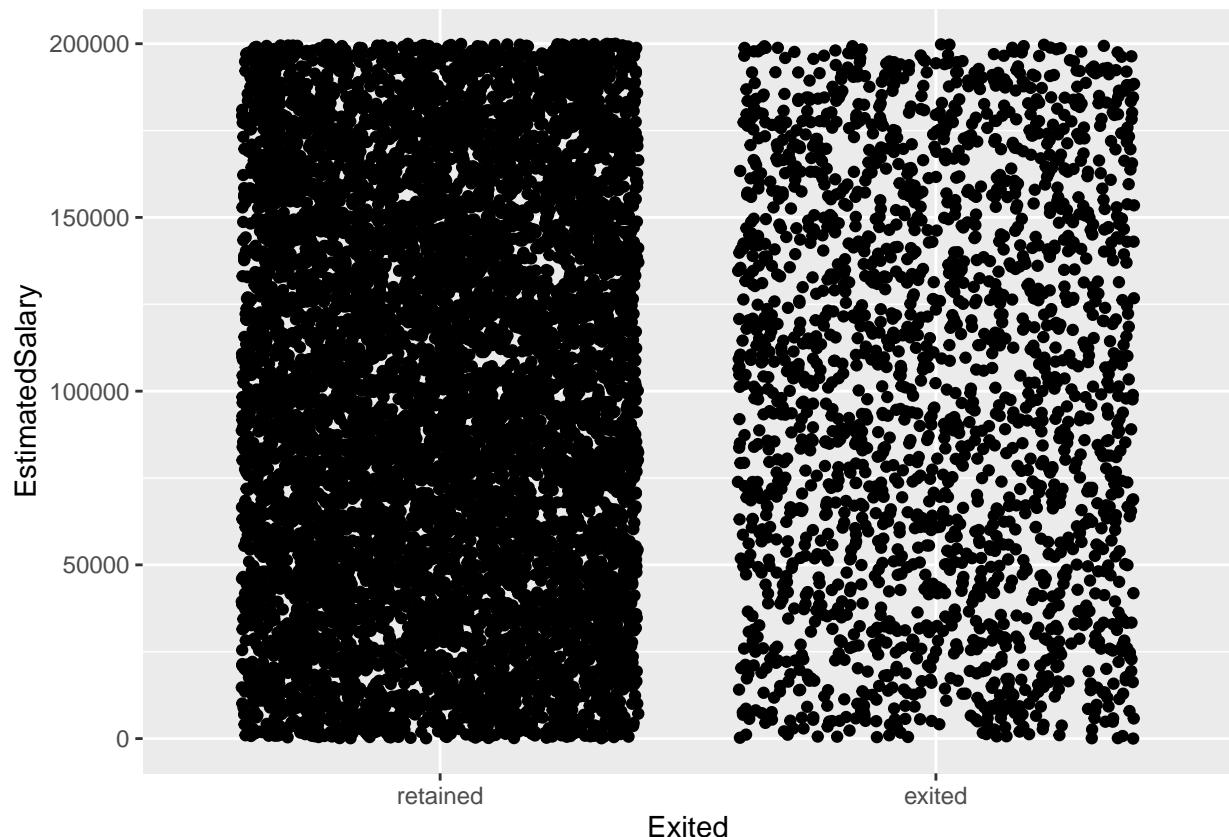


IS ACTIVE MEMBER



Does not provide a lot of insight or make much sense in the data set, will remove.

ESTIMATED SALARY



SURNAME

