

# DataVizA Tutorial: Plotting two variables: Solutions

*Department of Econometrics and Business Statistics, Monash University*

## *Tutorial 4*

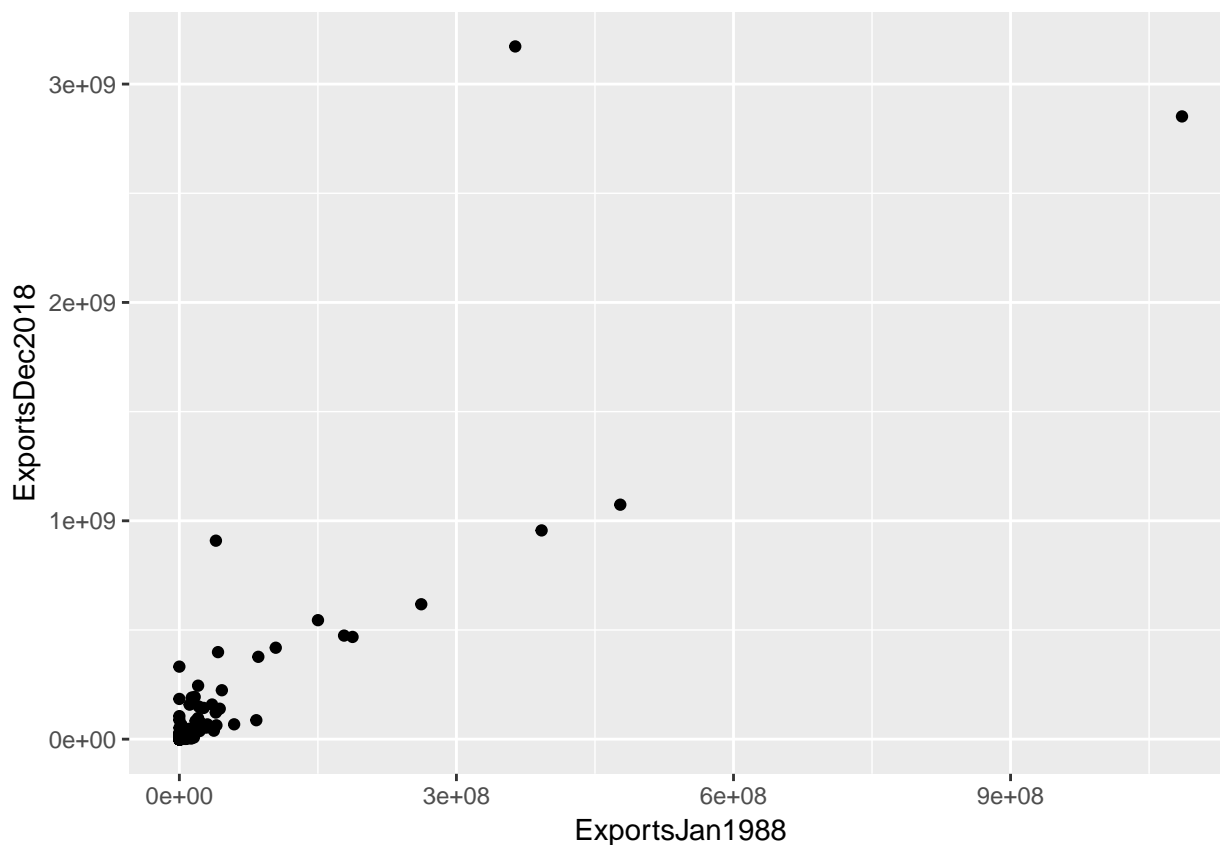
### Log Scales

1. Why might we display data using a log scale?

*Some data may be heavily skewed to the right. In such cases a log scale compresses together large values and stretches out small values.*

2. Using the Swiss Exports data plot a scatterplot of exports in January 1988 against exports in December 2018.

```
SwissExp<-readRDS('SwissExport.rds')
library(tidyverse)
ggplot(SwissExp,
  aes(x=ExportsJan1988,
      y=ExportsDec2018))+
  geom_point()
```



3. Do the same plot but using log scales.

```
SwissExp<-readRDS('SwissExport.rds')
library(tidyverse)
ggplot(SwissExp,
```

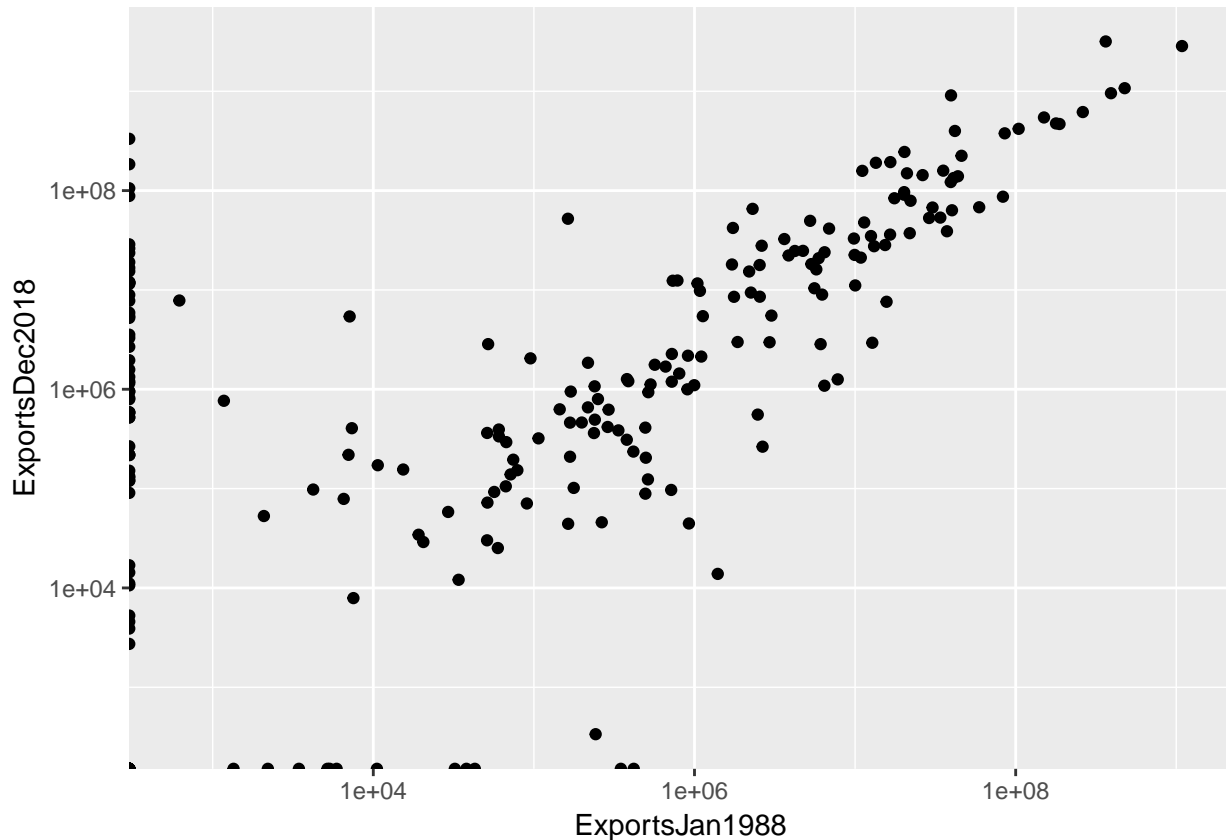
```

aes(x=ExportsJan1988,
    y=ExportsDec2018))+
  geom_point()+
  scale_x_log10()+scale_y_log10()

```

## Warning: Transformation introduced infinite values in continuous x-axis

## Warning: Transformation introduced infinite values in continuous y-axis



#### 4. Compare these two plots

*On the log log scale there is less overplotting. Rather than have all countries in the bottom left hand corner the points are spread out over the whole plot.*

#### 5. How do you understand the warning message that occurs when the log log scale is used?

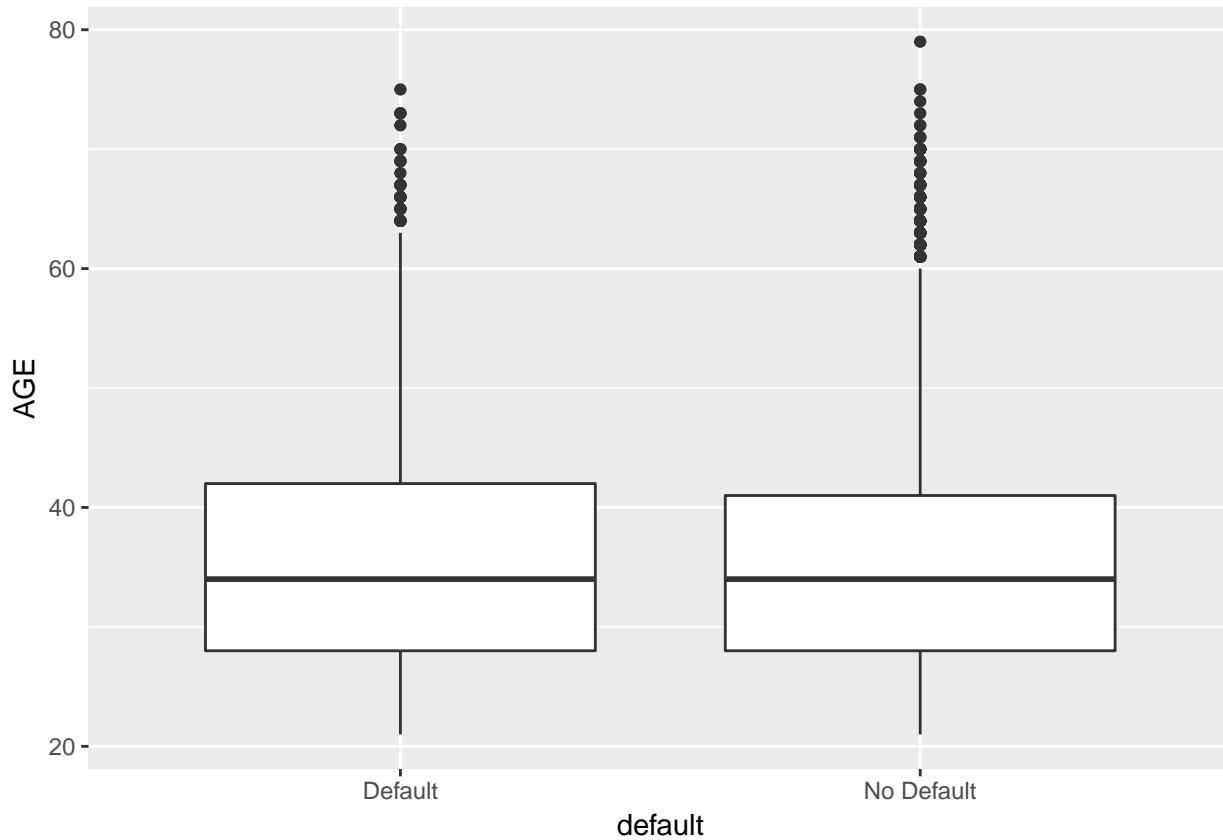
*There are a number of countries that receive no exports from Switzerland. These are values of 0. Mathematically, when taking the log of zero, the result is negative infinity. Note that ggplot has done something sensible and plotted these zero values along the margins of the plot. However, it may be better to either add a very small value to these zeroes or exclude them altogether.*

## Credit default data

The dataset *credit.rds* contains demographic information and repayment history for individuals who may have either defaulted on their credit card payments. The variable of interest is *default* which is equal to 1 for customers who fail to pay their credit card bill and 0 otherwise. More details on the dataset can be found [here](#)

6. Using box plots, explore whether the distribution of age is different for the default group and non-default group. Interpret your result.

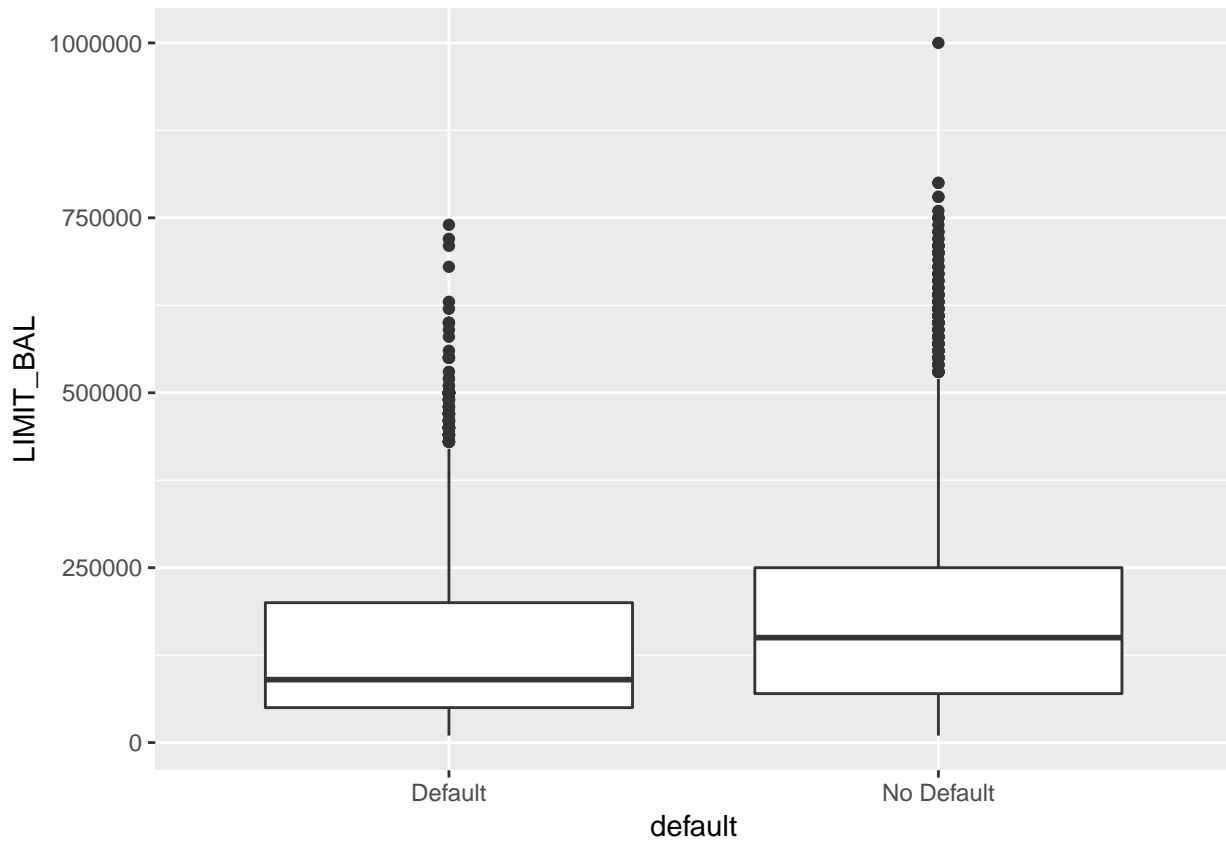
```
credit<-readRDS('Credit.rds')
ggplot(credit,
       aes(x=default,y=AGE))+
geom_boxplot()
```



*#The plot does not suggest differences between the age distribution  
#of the default and non-default groups. Age is not likely to be a  
#useful feature in discriminating between these groups.*

7. Using box plots, explore whether the distribution of the credit limit (LIMIT\_BAL) is different for the default group and non-default group. Interpret your result.

```
credit<-readRDS('Credit.rds')
ggplot(credit,
       aes(x=default,y=LIMIT_BAL))+
geom_boxplot()
```



*#The plot suggests that distribution of the credit limit is higher  
#for non-default groups.*

8. Suppose you work for the credit card company. In a business meeting a colleague suggests that the credit card limit could be useful in predicting when a customer defaults.

*The credit card limit could be useful in helping to explain default since a higher limit is associated with a lower probability of default.*

9. Suppose you work for the credit card company. In a business meeting a colleague looks at the plot and suggests that low credit card limits are causing defaults. They suggest to raise the credit card limit of all customers. Is this a good idea?

*This is probably a bad idea. Association does not imply causation. If anything the causation is likely to run the other way. Customers identified as being at risk of default for some other reason (unstable employment, low income) are only offered a low credit limit.*