# DataVizA Tutorial: Data Munging: Solutions

## Department of Econometrics and Business Statistics, Monash University

### Tutorial 5

### Swiss Exports: Full Data

The file *SwissExportsFull.csv* contains the full export data for Switzerland. Each row represents a different date. The first column is the date variable, the second column is the year only and each remaining column measures exports to a different country.

1. Read the data into R.

```
library(tidyverse)
SwissWide<-read_csv('SwissExportsFull.csv')
```

```
## Warning: Missing column names filled in: 'X154' [154]
```

```
#This works but with a quirky warning.  One country code is NA for
#Namibia, however R treats NA as a missing value.  It can be fixed
#with

SwissWide<-read_csv('SwissExportsFull.csv')%>%
  rename(`NA`=X154)
```

```
## Warning: Missing column names filled in: 'X154' [154]
```

```
SwissWide
```

2. Get the data into long form using the `pivot_longer` function.

```
library(tidyverse)
SwissLong<-pivot_longer(data = SwissWide,
                        cols=c(-Date,-Year), #Do not use these variables
                        names_to = 'Country', #Column names become variable
                        values_to =  'Exports') #All numbers are exports
SwissLong
```

3. Recall that in the previous tutorial, one issue was that monthly data were noisy. Using `group_by` and `summarise` create a new dataset of yearly aggregate exports to each country.

```
SwissLong%>% #This is much easier with long data
  group_by(Year,Country)%>%
  summarise(YearlyExports=sum(Exports))->SwissYearly
SwissYearly
```

4. Now produce a scatterplot on a log-log scale of 1988 exports against 2018 exports.

```
SwissYearly%>%
  filter(Year %in% c(1988,2018))%>% #Filter years
  pivot_wider(id_cols = Country,
              names_from = Year,
```
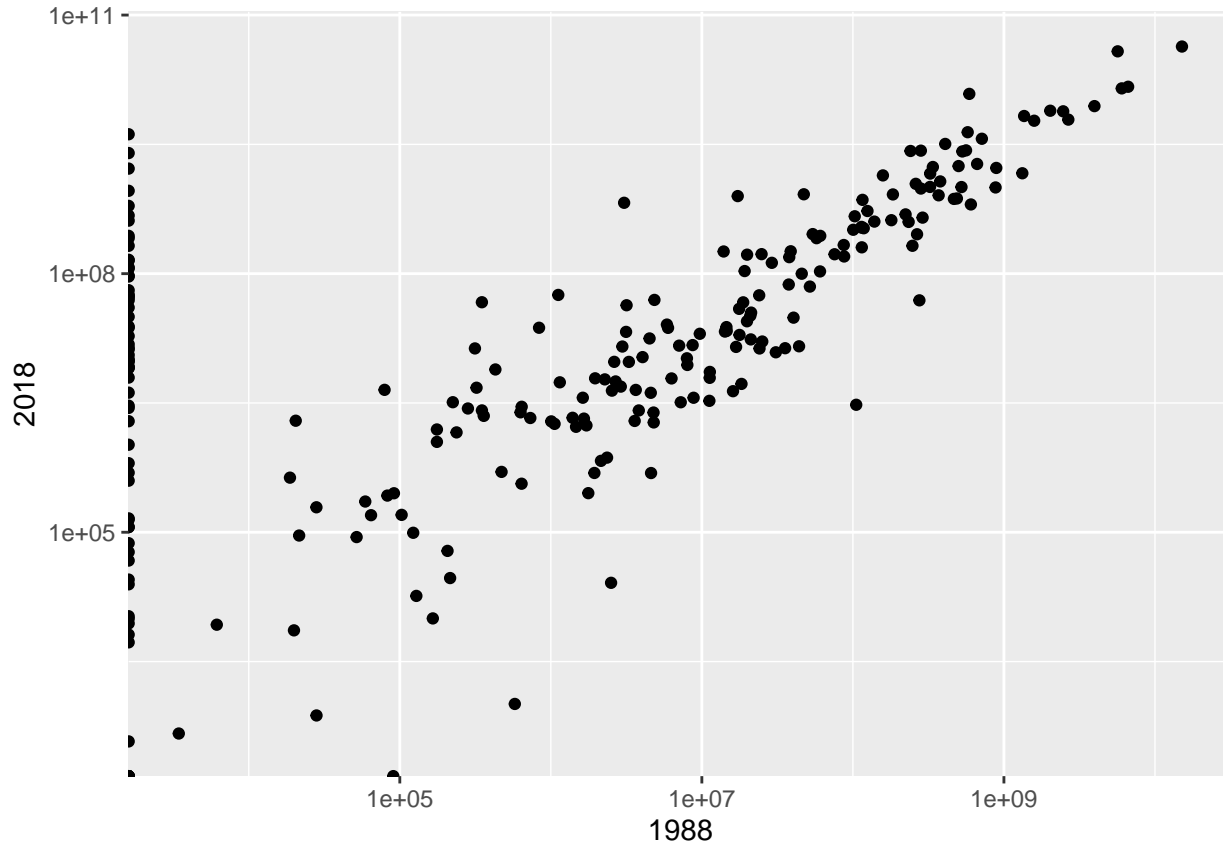
```
                values_from = YearlyExports)-> #Need to widen
  SwissYearlyWide

SwissYearlyWide%>%
  ggplot(aes(x=`1988`,y=`2018`))+
  geom_point()+
  scale_x_log10()+scale_y_log10()
```

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

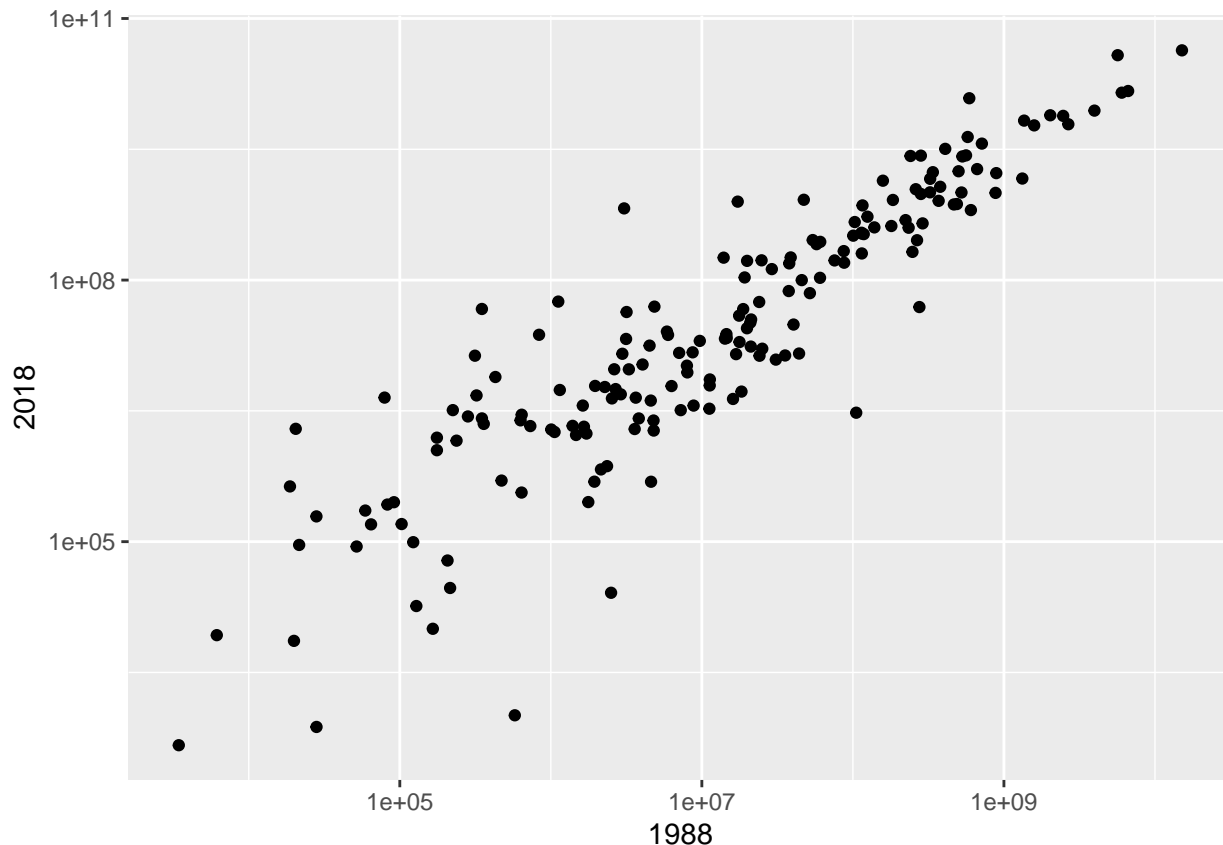## Warning: Transformation introduced infinite values in continuous y-axis



5. Produce the same plot but remove all countries for which exports are zero in either 1988 or 2018.

```
#Although it works, avoid the temptation to define the logical
#statement as it is written in words !((`1988`==0)|(`2018`==0))

SwissYearlyWide%>%
  filter((`1988`!=0)&(`2018`!=0))%>% #Filter years
  ggplot(aes(x=`1988`,y=`2018`))+
  geom_point()+
  scale_x_log10()+scale_y_log10()
```

2

## Options data

The following example uses Options data from Yahoo Finance. The owner of an put option has the right (but not the obligation) to sell stocks at a predetermined price (the Strike Price) on some fixed date (the Expiry date). A call is the same but gives the owner the right to buy stocks.

The objective of this exercise is to produce the well-known *volatility smile* result from finance. This result states that for a given Expiry date, a plot of Implied Volatilty against Strike Price is U-shaped. Implied Volatilty is the volatilty of a stock that is computed from stock option data assuming a specific pricing model. The exercise uses the data `apple_options.csv` which can be found on Moodle.

1. Read the data from this csv file into R.

```
apple<-read_csv('apple_options.csv')
```

2. The Implied Volatility has been imported as a character variable. To plot this is must be converted to a numeric variable. Create this using the `mutate` function.

Hint: The following code removes the percentage sign, converts to numeric and divides by 100.

```
gsub('%','','25%')%>%as.numeric()/100
```

```
## [1] 0.25
```

```
gsub('%','','1.32%')%>%as.numeric()/100
```

```
## [1] 0.0132
```

*The following code will create the new variable*

```
apple%>%
  mutate(ImpliedVol=gsub('%','',`Implied volatility`)%>%as.numeric()/100)
```

3. The volatility smile is best observed when options with a single expiry date are used. To use as much data as possible, find the expiry date that has the most put options.

```
apple%>%
  filter(Type=='Put')%>%
  group_by(`Expiry date`)%>%
  summarise(TotalOptions=n())
```

*The expiry date with the most options contracts is 2020-05-29 with 120 put options.*

4. Options that are very far *out of the money* (very low strike price for a put option) should be excluded from the analysis. Building on previous answers, construct a data frame that only keeps put options from the expiry date in your answer to Question 3, and that have a strike price above 250.

```
apple%>%
  mutate(ImpliedVol=gsub('%','',`Implied volatility`)%>%as.numeric()/100)%>%
  filter(Type=='Put',`Expiry date`=='2020-05-29',`Strike Price`>250)
```
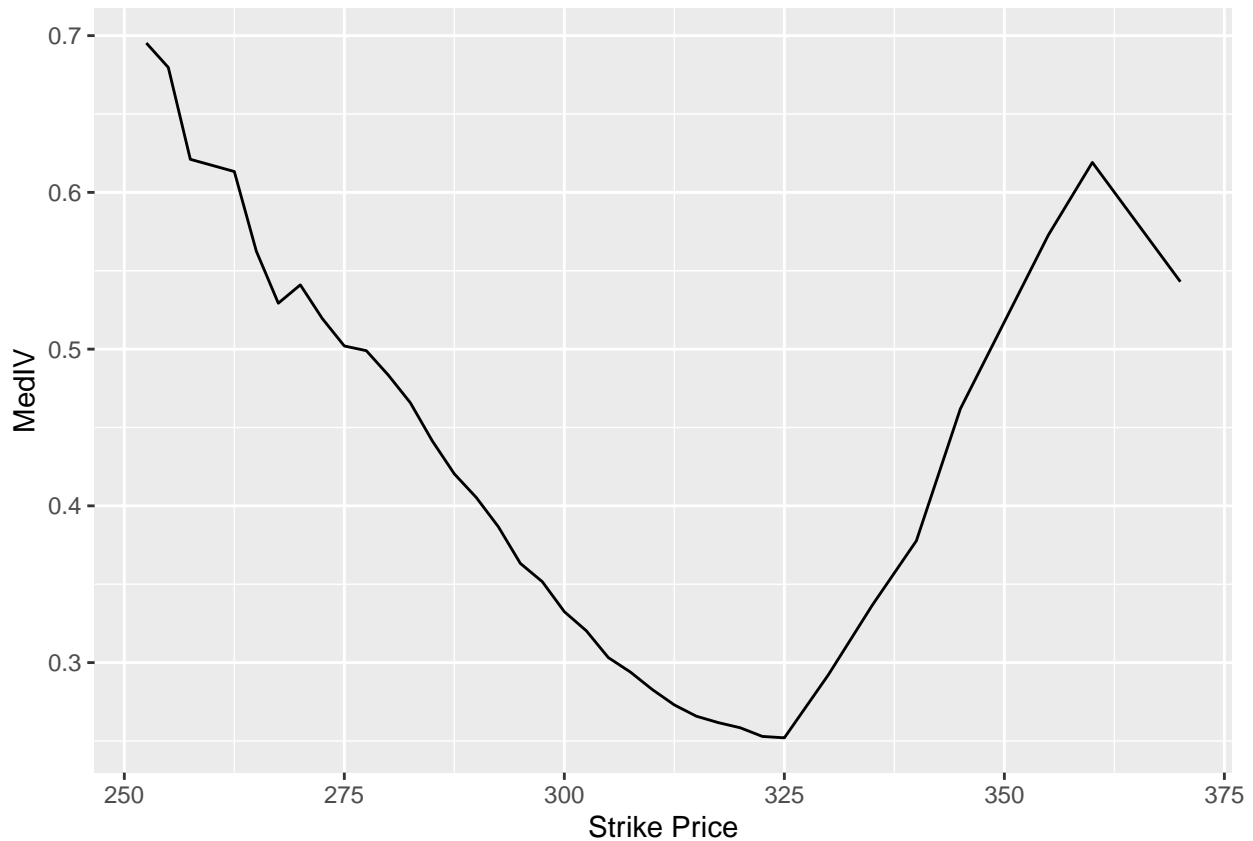
*Note that the filter function could use the AND operator as well.*

5. Using the data constructed in Q4, find the median value of Implied Volatility for each Strike Price.

```
apple%>%
  mutate(ImpliedVol=gsub('%','',`Implied volatility`)%>%as.numeric()/100)%>%
  filter(Type=='Put',`Expiry date`=='2020-05-29',`Strike Price`>250)%>%
  group_by(`Strike Price`)%>%
  summarise(MedIV=median(`ImpliedVol`,na.rm = T))
```

6. Plot Implied Volatility against Strike Price using a line plot.

```
apple%>%
  mutate(ImpliedVol=gsub('%','',`Implied volatility`)%>%as.numeric()/100)%>%
  filter(Type=='Put',`Expiry date`=='2020-05-29',`Strike Price`>250)%>%
  group_by(`Strike Price`)%>%
  summarise(MedIV=median(`ImpliedVol`,na.rm = T))%>%
  ggplot(aes(x=`Strike Price`,y=MedIV))+geom_line()
```

## Web scraping

The options data were obtained using web scraping, the following exercise helps you to scrape data for a single strike price. Go to this link and then click on the first strike price. Scrape the data that you find after clicking on a strike price.

*The exact link will depend on when the data is retrieved but will be similar to* https://au.finance.yahoo.com/quote/AAPL/options?strike=180&straddle=false

*The following code scrapes the data.*

```
#The following line should match the exact URL
url<-'https://au.finance.yahoo.com/quote/AAPL/options?strike=180&straddle=false'

url%>%
    read_html%>%
    html_table()->data
```

*This returns a list with first element:*

## First Normal Form

Discuss whether the following databases satisfy first normal form.

Database A:

| Name | Social Media Username |
| --- | --- |
| Jane Smith | Facebook: jsChampion |
| Kamal Usman | Twitter: kusman, LinkedIn: ku87 |
| Li Xiao | WeChat: lx99 |

Database B:

| Name | Social Media Username |
| --- | --- |
| Jane Smith | Facebook: jsChampion |
| Kamal Usman | Twitter: kusman |
| Kamal Usman | LinkedIn: ku87 |
| Li Xiao | WeChat: lx99 |

*Database A has the same issue as seen in lectures. Kamal Usman has two social media accounts so the entry is not atomic. The second Database resolves this issue. However arguably the variable Social Media Username is still not atomic. There are two separate pieces of information, the social media platform and the username. A better database would store these as two separate variables.*

*Another point worth mentioning is that even the variable name might not be considered to be atomic. For example in some contexts it may be important to serparate family name from given names.*