Mark Peters

# R Progammering Inlämningsuppgift

Once the game begins, who knows when it will stop?

# 1. Basic algorithmic thinking in R:

I started by designing each of the functions needed to run the program. It was stressed that it should be our focus, so I wanted to get it right.

Minimum function:

```
minimum <- function(x){
  min <- x[1]
  for (i in 2:length(x)) {
    if(x[i]< min)
      min <- x[i]
  }
  return(min)
}
```

The function accepts a vector and assigns the first value in the list to 'min'. It then cycles the rest of the values through a 'for' loop, from the second to the last in the list (length giving the function the point at which to stop), replacing the min value if a smaller one is found. This value is then returned to the parent program.

Maximum function:

```
9
10  maximum <- function(x){
11    max <- x[1]
12    for (i in 2:length(x)) {
13      if(x[i]> max)
14        max <- x[i]
15    }
16    return(max)
17  }
18
```

The maximum function applies almost the same process as the minimum except, naturally, replacing the 'less than' sign with the 'greater than'. It cycles through all supplied values and replaces 'max' with an element that is larger than it.

Mean function:

```
19  mean_me <- function(x){
20    sum <- x[1]
21    for (i in 2:length(x)) {
22      sum <- sum + x[i]
23    }
24    mean_me <- sum/length(x)
25    return(mean_me)
26  }
```

The mean function – which I christened 'mean-me' to avoid confusion – uses the same structure as the previous two functions (setting an initial value and then cycling through the rest with a for loop) but has an extra stage.

Sum becomes that number plus the next on the list. After it has collated all the numbers in the list it gets divided by the vector length. This gives the mean value for the entire vector.

Dataframe creation:

```
## created 10 vectors on 10 valu
temp1 <- rep(0, 10)      ## create
temp2 <- rep(0, 10)
temp3 <- rep(0, 10)
temp4 <- rep(0, 10)
temp5 <- rep(0, 10)
temp6 <- rep(0, 10)
temp7 <- rep(0, 10)
temp8 <- rep(0, 10)
temp9 <- rep(0, 10)
temp10 <- rep(0, 10)
## Create a dataframe of all ten
df <- data.frame(temp1, temp2, t
```

No doubt far from the most efficient way to create a dataframe, but I created each vector individually and filled them with zeros. I did this because my program asks for the city name first and I needed values in the vector for the name to 'stick'. All the vectors were placed into a dataframe named 'df'.

```
## A loop where user data is taken in and stored in a dataframe
for (i in 1:10) {
   name <- readline(prompt="Enter city name: ")
   names(df)[i] <- name
   print(paste("Enter temperatures for each day for ", names(df)[i]))
   df[ ,i] <- ReadVector() ## function to take in multiple values
}
```

Input User Data:

I used another 'for' loop to take input data from the user. I know this isn't a very 'user friendly' way to generate data, but for the purposes of this exercise it works.

The user enters the city name first which is added to the first vector with the 'names' function. The 'ReadVector' algorithm takes in each temperature value, one at a time, and fills the vector.

The for loop then moves on to the next city, and so on until all ten cities' values have been entered.

```
## A loop to generate max, min & mean for all entered temperatures
for (i in 1:10) {
   maxi <- maximum(df[ ,i])
   cat(paste("\nThe maximum value in column",i,"is",maxi,"\n"))
   mini <- minimum(df[ ,i])
   cat(paste("\nThe minimum value in column",i,"is",mini,"\n"))
   meani <- mean_me(df[ ,i])
   cat(paste("\nThe mean value for column",i,"is", meani,"\n"))
}
```

Print out max, min & mean:

Another 'for' loop (bit of a one-trick-pony I am) is used to generate min, max and mean values for each of the cities entered by the user. The loop takes the column number and then prints out the values for it (including the city's name) to the screen.

After printing the three values for the first city it moves onto the second until all ten have been printed to screen.

```
## Loop doing the same thing above but using inbuilt 'R' functions
for (i in 1:10) {
   maxi <- max(df[ ,i])
   cat(paste("\nThe R maximum value in column",i,"is",maxi,"\n"))
   mini <- min(df[ ,i])
   cat(paste("\nThe R minimum value in column",i,"is",mini,"\n"))
   meani <- mean(df[ ,i])
   cat(paste("\nThe R mean value for column",i,"is", meani,"\n"))
}
```
Same results but with inbuilt functions:

I repeat the same process as the previous but generate results with the inbuild R functions instead of my own.

Results of the program using a second test program:

This is a print out of all the values that came from my own functions:

```
The maximum value in column 1 is 30

The minimum value in column 1 is 20

The mean value for column 1 is 26.1

The maximum value in column 2 is 27

The minimum value in column 2 is 15

The mean value for column 2 is 20.8666666666667

The maximum value in column 3 is 36

The minimum value in column 3 is 26

The mean value for column 3 is 31.8666666666667
```

This is a printout of the values generated form the R functions:

```
The R maximum value in column 1 is 30

The R minimum value in column 1 is 20

The R mean value for column 1 is 26.1

The R maximum value in column 2 is 27

The R minimum value in column 2 is 15

The R mean value for column 2 is 20.8666666666667

The R maximum value in column 3 is 36

The R minimum value in column 3 is 26

The R mean value for column 3 is 31.8666666666667
>
```

Happily, the results were the same.

## 2.Data Pre-Processing:

For this question I worked with the suggested dataset, mainly because it needed a good amount of processing in order to be usable, and that wasn't the same for many of the ones that I came across.

Here is how the chronic kidney disease dataset was described on Kaggle:

*The data was taken over a 2-month period in India with 25 features ( eg, red blood cell count, white blood cell count, etc). The target is the 'classification', which is either 'ckd' or 'notckd' - ckd=chronic kidney disease.*

Here is a full list of each column in the dataset:

| | |
|---|---|
| age - age | sod - sodium |
| bp - blood pressure | pot - potassium |
| sg - specific gravity | hemo - hemoglobin |
| al - albumin | pcv - packed cell volume |
| su - sugar | wc - white blood cell count |
| rbc - red blood cells | rc - red blood cell count |
| pc - pus cell | htn - hypertension |
| pcc - pus cell clumps | dm - diabetes mellitus |
| ba - bacteria | cad - coronary artery disease |
| bgr - blood glucose random | appet - appetite |
| bu - blood urea | pe - pedal edema |
| sc - serum creatinine | ane - anemia |
| | class - classification |

I worked from left to right. I got the mean of all ages after removing the 'NA's' and them with the mean age of all participants.

```
> mean_age <- as.integer(mean(kidney_disease$age, na.rm = TRUE))
> kidney_disease$age[is.na(kidney_disease$age)] = mean_age
>
```

I repeated the same steps for each of the following columns until I got to 'rbc' (red blood cells).

| | id | age | bp | sg | al | su | rbc |
|---|---|---|---|---|---|---|---|
| 1 | 200 | 90 | 90 | 1.025 | 1 | 0 | NA |
| 2 | 171 | 83 | 70 | 1.020 | 3 | 0 | norma |
| 3 | 39 | 82 | 80 | 1.010 | 2 | 2 | norma |
| 4 | 160 | 81 | 60 | NA | NA | NA | NA |
| 5 | 194 | 80 | 70 | 1.010 | 2 | NA | NA |

In 'rbc' we meet 'characters' for the first time. My approach here was what was presented to us in class – that in order to make the code machine readable the bivalent values should be changed to either a numerical '1 and 0' or a logical 'True and False'.

I assigned 'normal' to '1' and 'abnormal' to a '0' and changed the values for the whole column, and subsequently assigned 'NA' to '0'.

```
> View(kidney_new)
> kidney_new$rbc = factor(kidney_new$rbc, levels = c('normal', 'abnorma
l'), labels = c(1, 0))
> kidney_new$rbc[is.na(kidney_new$rbc)] <- 0
>
```

This process was repeated for another four columns, the only difference being the nature of the character strings being replaced:

| | rbc | pc | pcc | ba | bgr |
|---|---|---|---|---|---|
| 0 | NA | normal | notpresent | notpresent | |
| 0 | normal | normal | notpresent | notpresent | |
| 2 | normal | NA | notpresent | notpresent | |
| NA | NA | NA | notpresent | notpresent | |
| NA | NA | abnormal | notpresent | notpresent | |

the 'pcc' and 'ba' columns (0,1) columns were 'notpresent' and 'present' respectively.

| | rbc | pc | pcc | ba | b |
|---|---|---|---|---|---|
| 0 | 0 | 1 | 0 | 0 | |
| 0 | 0 | 1 | 0 | 0 | |
| 3 | 1 | 1 | 0 | 0 | |

The next six columns were treated the same way as the first five, by calculating the mean value of each column and replacing the 'NA' values with the result.

| bgr | bu | sc | sod | pot | hemo | pc |
|---|---|---|---|---|---|---|
| 139 | 89.0 | 3.0 | 140 | 4.1 | 12.0 | 37 |
| 102 | 60.0 | 2.6 | 115 | 5.7 | 8.7 | 26 |
| 140 | 70.0 | 3.4 | 136 | 4.2 | 13.0 | 40 |
| 148 | 39.0 | 2.1 | 147 | 4.2 | 10.9 | 35 |

Interestingly, the 'pcv' column threw me a curveball, as though it was a list of integer values they were listed as being 'characters. I assume that was because they were entered into the dataset from user input. This meant I had to change all the values to 'numeric' before I could find the mean and so fill in all the empty boxes.

```
> kidney_new$pcv <- as.numeric(kidney_new$pcv)
Warning message:
NAs introduced by coercion
> summary(kidney_new$pcv)
   Min. 1st Qu.  Median   Mean 3rd Qu.   Max.   NA's
```

The next two columns were numerical and the remaining eight were all characters – all of them bivalent in nature. They had different permutations but were all either true or false in nature – 'yes/no', 'good/poor' and 'ckd/notckd'.

| yes | yes | no | good | no | no | ckd |
|---|---|---|---|---|---|---|
| yes | no | no | poor | no | yes | ckd |
| yes | yes | no | good | no | no | ckd |
| yes | yes | yes | poor | yes | no | ckd |

All columns were assigned ones and zeros and NA's were removed.

By the end of this process I had a dataset that was machine ready with all NA values removed and all columns being 'numeric' in form.
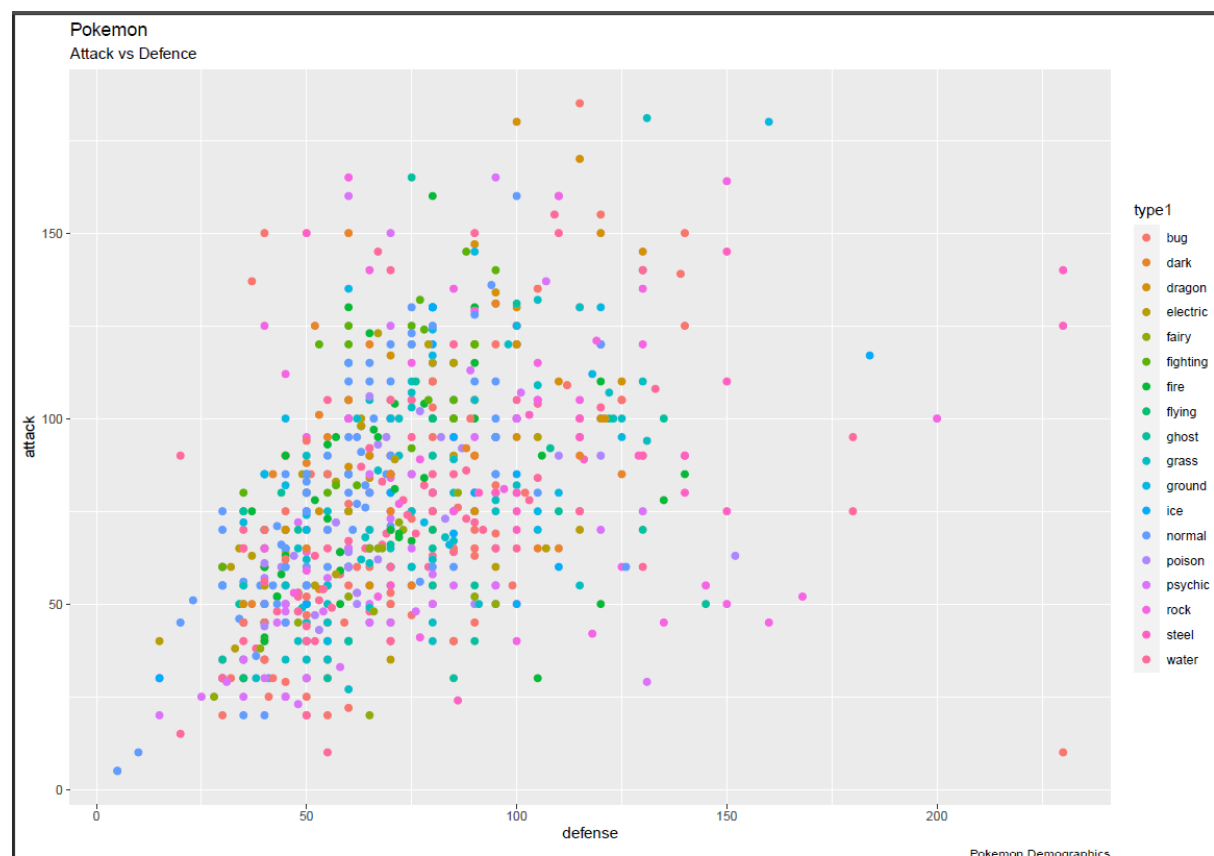
# 3. DataVisualization:

As I approach the end of a very long and intense couple of weeks, the goal very nearly in sight, I have already begun to critique my approach to the tasks. Maybe this is premature since I haven't finished yet, but I can already admit that I spent (and wasted) way too much time trying to find a dataset to work with, and even more time trying to decide what kind of graph would best portray the data. Hours were spent trying things out and then binning the results.

It's not clear but it seems you may have wanted to see all those abortive attempts to produce something of value. I regret my petulance in trashing the vast majority of what I did, if for no other reason than to illustrate the amount of time I've put into these tasks.

Not a time to stress the material I didn't keep and focus on what I did.

Practice:

I made a scatterplot of the attack vs defence capabilities of all (currently) discovered Pokemon for my kids. There were too many to name, so it's divided by colour into the various classes. They were mildly amused though not enlightened so I gave up on the dataset there and then.



(full-sized pdf versions can be found at the end of the document)

World Maps:

I eventually settled on a 'World' database which contained coordinates so I could plot on a world map. I wasted far too much time playing with all the various functions available, and I hope to develop this interest more when I have free time (somewhere around 2025 by my reckoning).

With the inbuilt view mode and 'polygon' function one can instantly create informative and interactive maps. I created a map showing life expectancy around the world which made for some sober viewing.
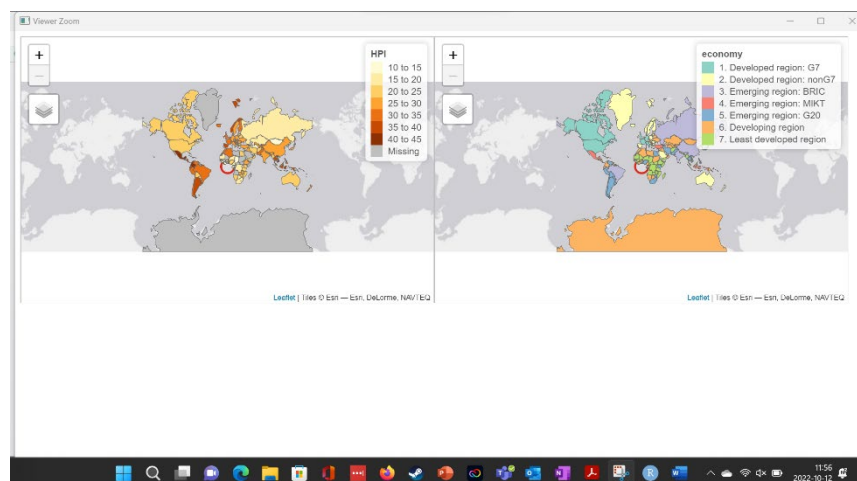
There were many fantastic functions, and the ones that stood out were the interactive ones, mainly as they were so easy to produce and are more pedagogic and memorable by their nature.



By using the 'tm_facets' function one can create maps that sit side by side and by 'syncing' them it's possible to access more information on each scan.

Scatterplot:

I wanted to create a scatterplot next in order to compare as many different factors as possible. In the limited time that I've been studying graphing with R, this has been the most powerful visual use of data that I've come across.

I chose to compare levels of inequality in countries to the HPI (happiness index) to see if racial or sexual discrimination was clearly reflected in the 'happiness' of each country's citizens. I added a legend for the income groups represented and the GDP of the countries as well.

```
gg2 <- ggplot(World, aes(HPI, inequality, col= income_grp, size= gdp_cap_est))+

  geom_point(alpha=0.5)+
  scale_size(range = c(.1, 24), name="GDP estimate")+

labs(title="Inequality Vs HPI", subtitle="indexed by GDP and Income Group",
     x="HPI (Happiness Index)",
     caption="World Demographics")
```



I really loved this graph. It clearly shows how without a stable economy and a relatively high standard of living it is hard for any country to achieve viable and visible equality (disregarding for now any cultural factors).

It is equally interested to note how spread out the HP index was among high earning countries. It raises many questions (for instance who is that country in the bottom left with lots of money, little inequality and with a miniscule amount of happiness to show for it. Tragic!

And the number of middle to high income countries that report very high indexes lends more credence to the cliché that 'money can't buy you everything'.

Boxplot:

Lastly, I created a boxplot to show the happiness of the different economic groups in the world. The world bank decides which group a country belongs to. Here is a map:



```
tmap_mode("plot")

tm_shape(World)+
  tm_polygons("economy")
```
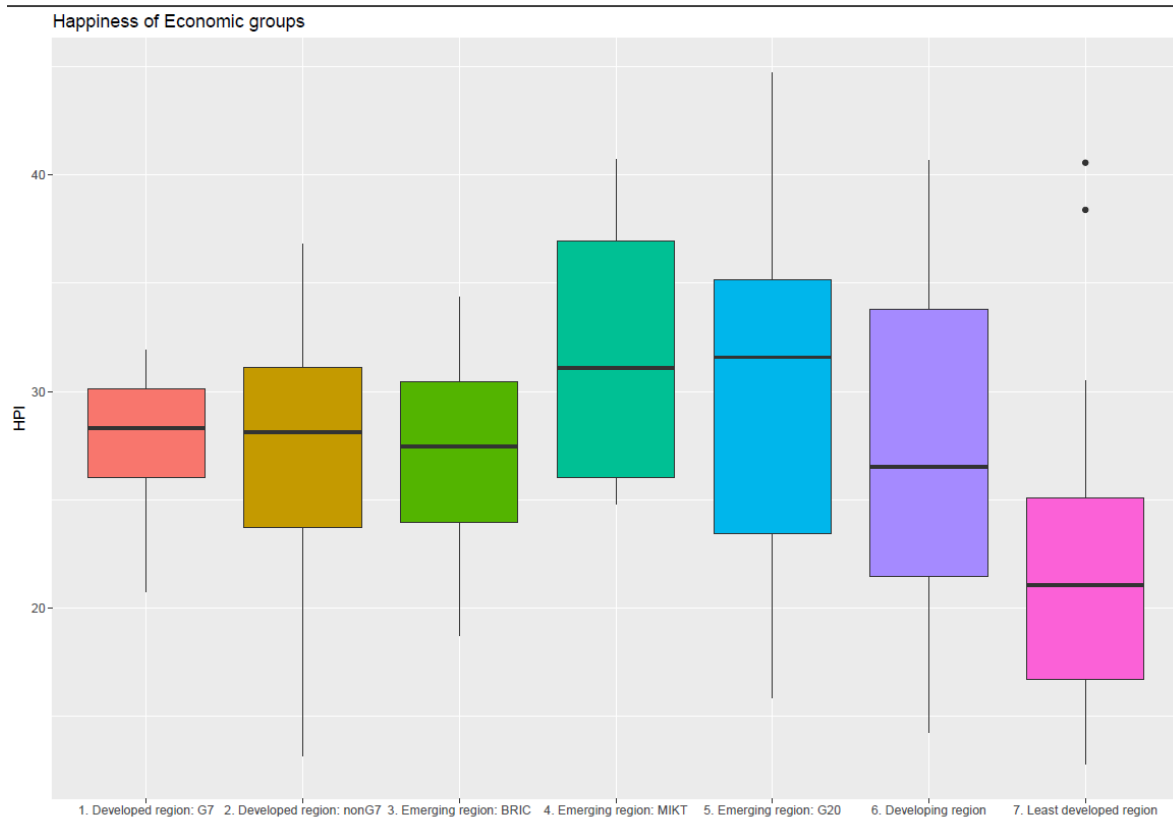
Again I used the tmap package to plot the map.

I created a boxplot with a focus on the economic groups, I wanted it to be as simple as possible.
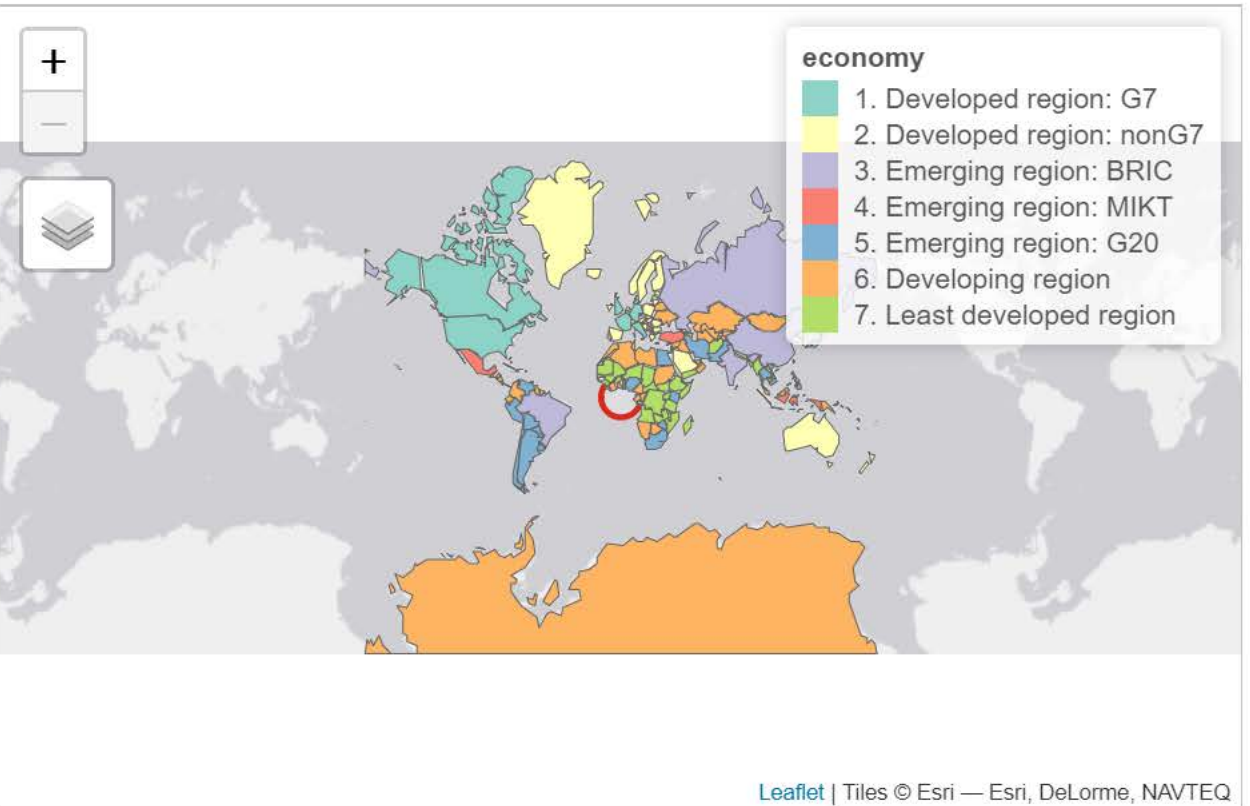
```
## create plot using economy as the x and happiness as the y, the fill
gg3 <- ggplot(World, aes(x=economy, y=HPI, fill=economy)) +

## call the boxplot through the geom function
  geom_boxplot()+
  labs(title="Happiness of Economic groups")+
xlab("Economic Group") +
  ## there is no legend so that the whole plot could fit on the page
  theme(legend.position="none") +
  xlab("")
```

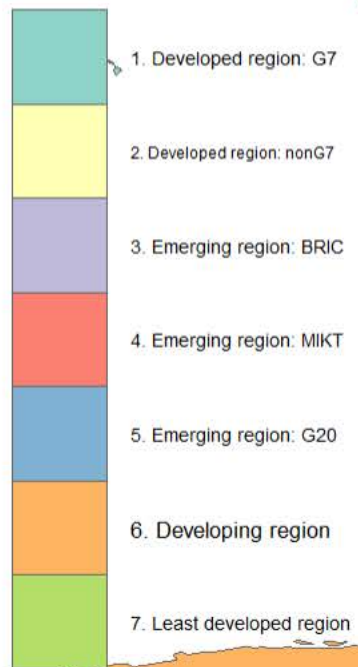Here is how the boxplot of happiness in these regions came out:



Not surprising that those with little money have little to be happy about, though it is surprising to see how middling and uninspiring the happiness measurement of the G7 countries was. There is something optimistic and heartening about seeing how high happiness is among emerging countries. If there is anything to conclude from this chart it is that how we measure happiness is fraught with difficulty, contradiction, and cultural relevance.
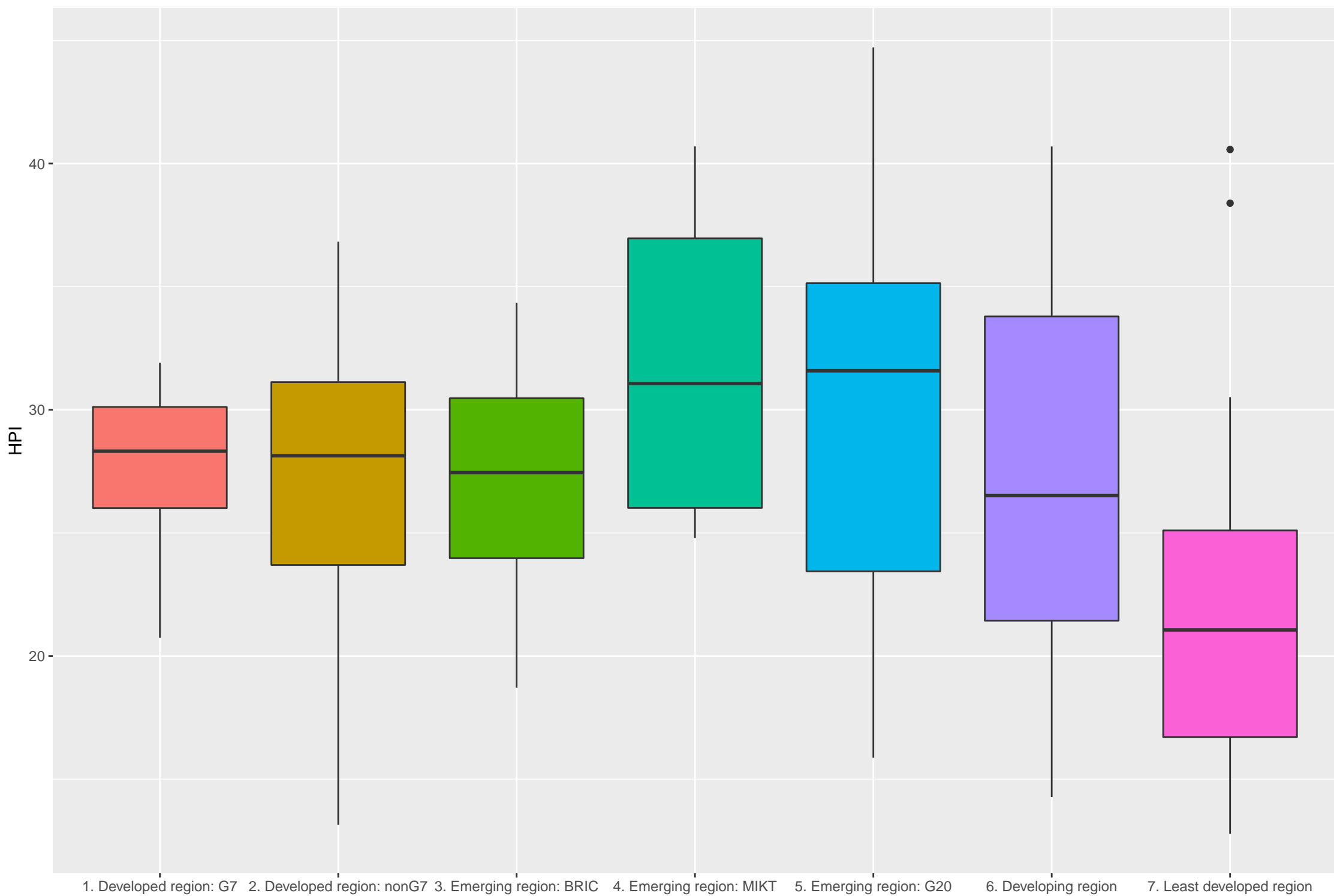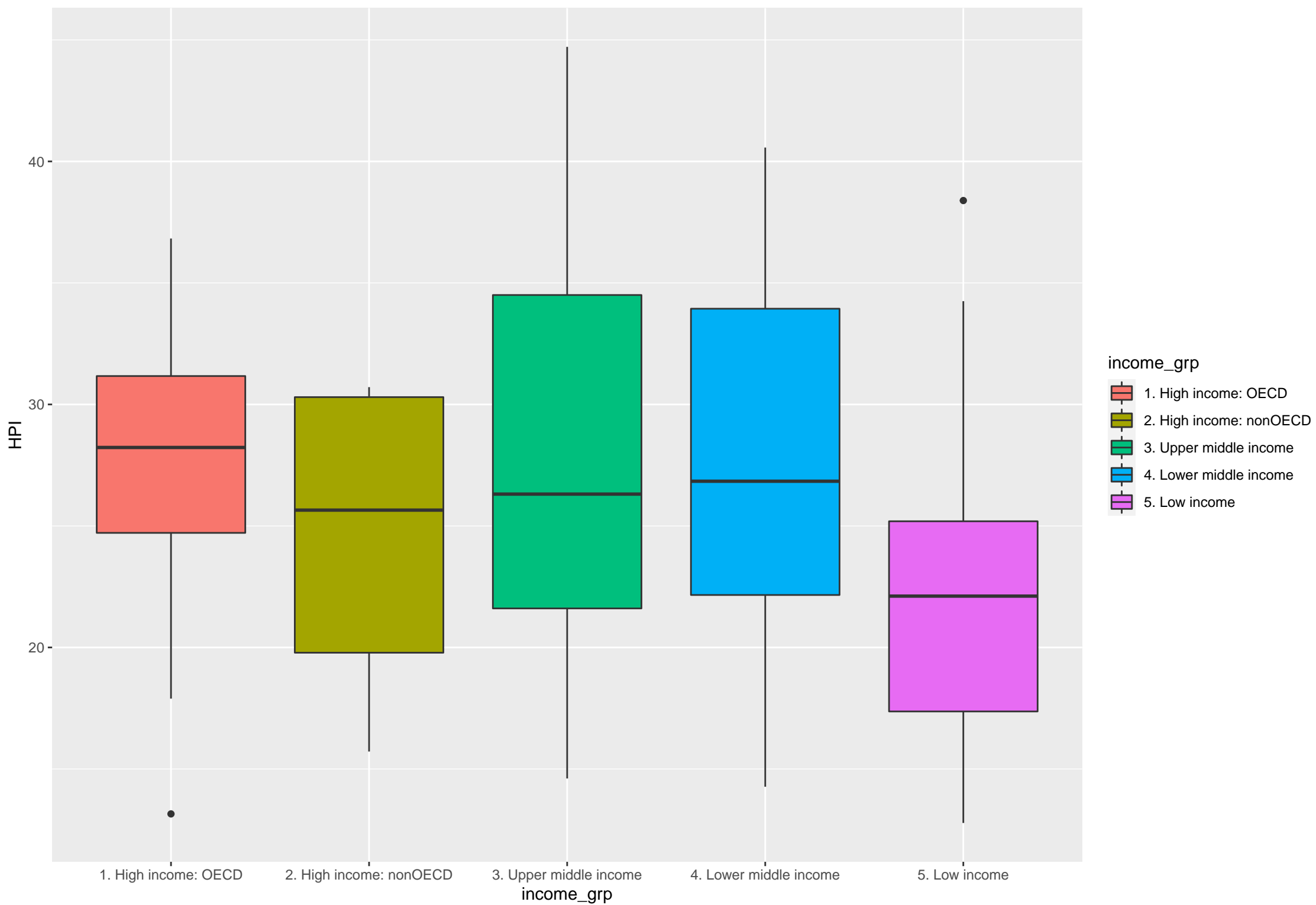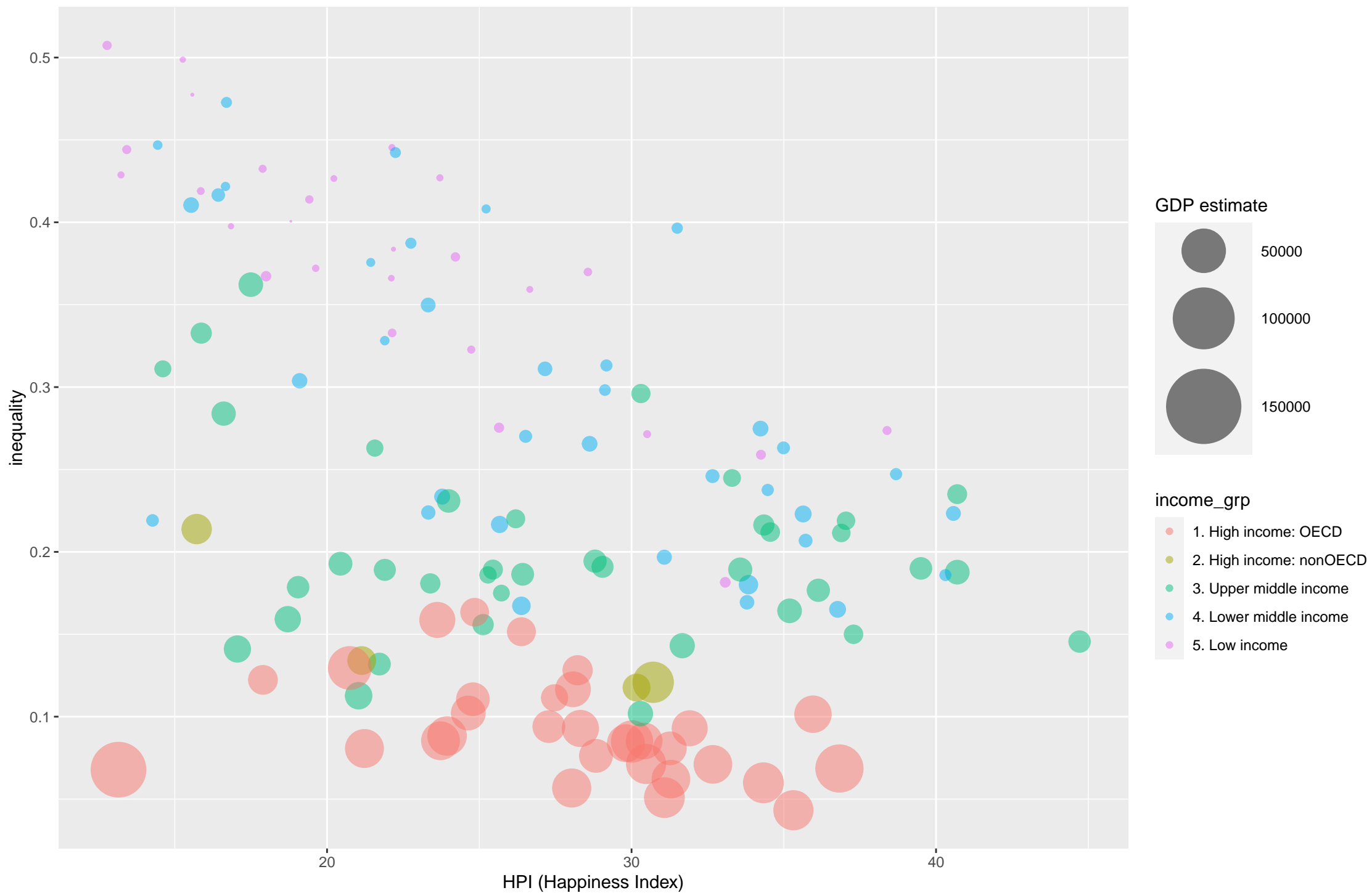
life_exp
- 45 to 50
- 50 to 55
- 55 to 60
- 60 to 65
- 65 to 70
- 70 to 75
- 75 to 80
- 80 to 85
- Missing

economy

1. Developed region: G7

2. Developed region: nonG7

3. Emerging region: BRIC

4. Emerging region: MIKT

5. Emerging region: G20

6. Developing region

7. Least developed region

Happiness of Income groups

## 4. Data Visualization – VG question:

I spent so long trying to decide upon a dataset and the graphs to go with it for question 3 that I had less time than I would have liked for question 4. After university I went to film school and have written screenplays and advertisements on and off since then. It seemed natural for me to search for a dataset within the field. I went with the top 1000 films on imdb, mainly because the length was sufficient for the brief, though I have not exactly been wowed by the data on offer. Given my self-enforced time restrictions I have had to make do.
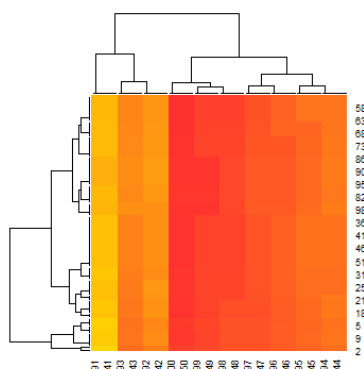
The dataset was not 'plot' ready, I found a number of rows with 'NA's'.

| | Meta_score | Director | Star1 | Star2 |
|---|---|---|---|---|
| ealthy erratic tippler, a dewy-eyed tramp ... | 99 | Charles Chaplin | Charles Chaplin | Virginia C |
| -year sentence for a violent crime, a 12-y... | 75 | Nadine Labaki | Zain Al Rafeea | Yordanos |
| e ravages of the Korean War, Sergeant Sü... | NA | Can Ulkay | Erdem Can | Çetin Tek |
| nse police officer, accompanied by Simo... | NA | Gayatri | Pushkar | Madhava |
| themselves linked in a bizarre way. When... | 79 | Makoto Shinkai | Ryûnosuke Kamiki | Mone Ka |
| havir Singh Phogat and his two wrestler ... | NA | Nitesh Tiwari | Aamir Khan | Sakshi Ta |

I used the same approach as in question 2, I got a mean for the row and replaced the missing values.
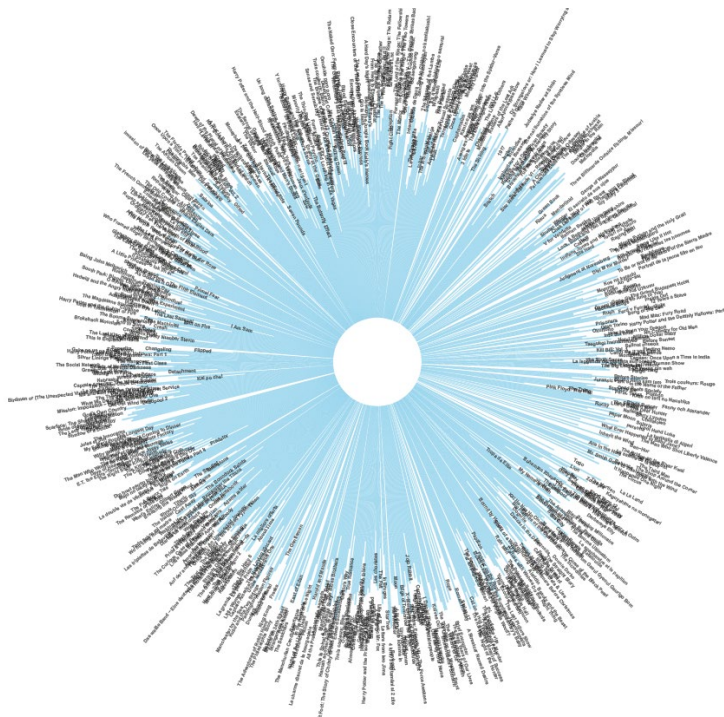
```
> summary(imdb_top_1000$Meta_score)
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
  28.00   70.00   79.00   77.97   87.00  100.00     157
> imdb_top_1000$Meta_score[is.na(imdb_top_1000$Meta_score)] <- 78
>
```

With a thousand rows of data I considered it wise to start with a landscape view, in order to get an idea of what I was dealing with. Here are a couple of valiant but failed attempts to discover insights from a mess of mass:
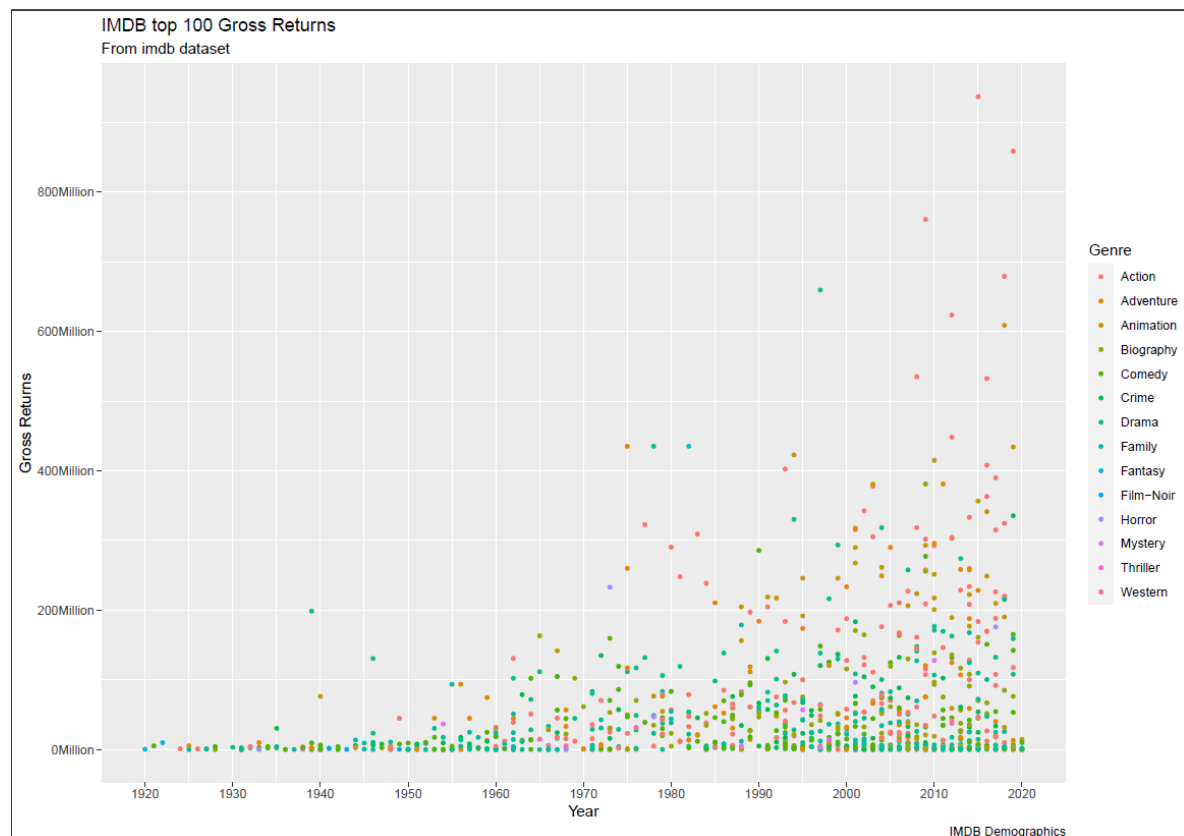


I started with a heatmap, but even with a thousand elements in the dataset it doesn't equate to much range when you consider the tens of thousands of rated films in the database.

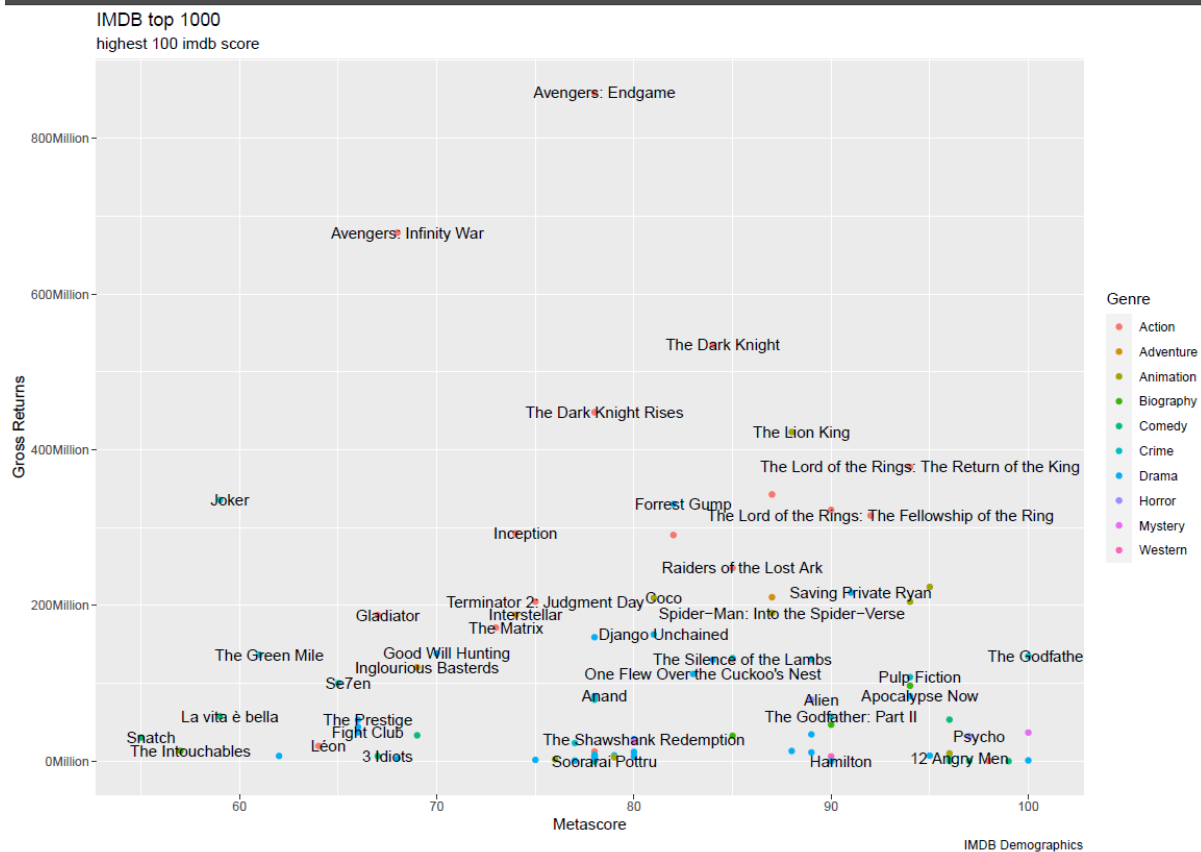Attempted a circular plot that came out eye strainingly mental:



I dare you to search for movie titles. Unfortunately, in the confusion of trying to figure out what I had created and the need to create something that made sense I deleted the code that formed this aberration.

So I turned to the tried and tested, and created a scatterplot when compared 'Gross profit' to 'Release Year' and using genre as the key. The idea of using actual film names for this graph was a non-starter. The problem with the result is that the profits of older films have not been adjusted according to inflation, so the results are skewed heavily in favour of films released in the last twenty years or so.

```
gg <- ggplot(newimdb, aes(x= Released_Year, y= Gross))+
geom_point(aes(col= Genre), size=1)+
labs(title="IMDB top 100 Gross Returns",
     subtitle="From imdb dataset", y="Gross Returns",
     x="Year", caption="IMDB Demographics")+
scale_x_continuous(breaks=seq(1920, 2020, 10))+
scale_y_continuous(breaks=seq(0, 1000000000,200000000), labels = function(x){paste0(x/1000000,'Million')})
```
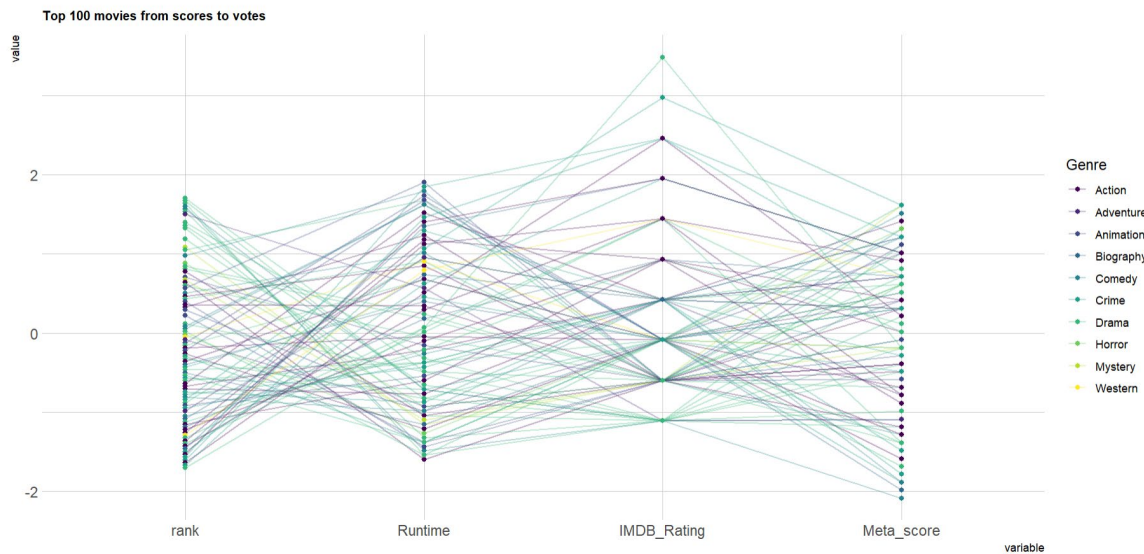
I did take the top hundred grossing films and plotted their gross profit to their 'Metascores' (a much more accurate measure of their critical worth than the 'IMDB rating', in my opinion). In this graph I could include the names of the films and again used 'genre' as a legend. I see it as the most natural and significant way to divide large groups of movies.



The graph clearly shows the dominance of action movies at the top of the box office, and within the genre how indomitable the comic book has been over the past decades.

```
gg2 <- ggplot(top100, aes(x=Meta_score, y=Gross)) +
  geom_point(aes(col=Genre))+
  labs(title="IMDB top 1000", subtitle="highest 100 imdb score", y="Gross Returns",
       x="Metascore", caption="IMDB Demographics")+
  geom_text(
    label=top100$Series_Title,
    nudge_x = 0.5, nudge_y = 0.5,
    check_overlap = T)+
  scale_y_continuous(breaks=seq(0, 1000000000,200000000),
                     labels = function(x){paste0(x/1000000,'Million')})
```

I wanted to find a way to show how people voting compared to critic's votes, which is the essential difference with the IMDB and metascore rating system. This was the graph that I tried out to show the difference. I also included a row for film length to see if that reflected any opinions on the films themselves. Genre again indexes, but the results are far from clear. What I think can be seen quite clearly is how critics tend to judge films more harshly than the average cinema goer.
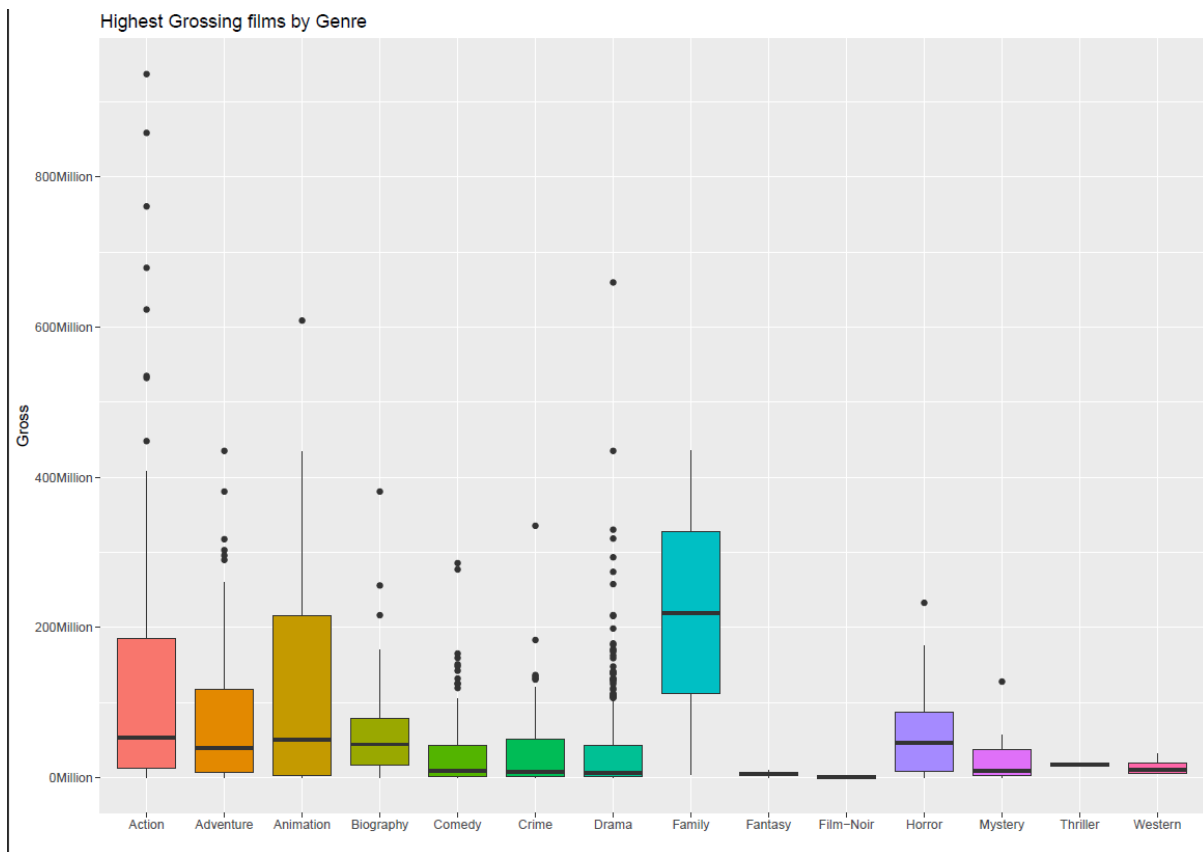
Top 100 movies from scores to votes

The third row is the public voting, the last row the critic's.

```
gg3 <- ggparcoord(rev_top100,
        columns = c(1,5,7,9), groupColumn = 6, order= "anyClass",
        showPoints = TRUE,
        title = "Top 100 movies from scores to votes",
        alphaLines = 0.3
) +
  scale_color_viridis(discrete=TRUE) +
  theme_ipsum()+
  theme(
    plot.title = element_text(size=10)
  )
```

'I will end my graphing of the thousand highest rated movies on IMDB with a boxplot, which shows much more clearly than the scatterplot how successful different genres have been over the decades. What stands out here is that although actions movies have outliers that are more successful on orders of magnitude from most other films, it is the family film which is the safest bet from an investment perspective, as they're most likely to break even or make a profit.
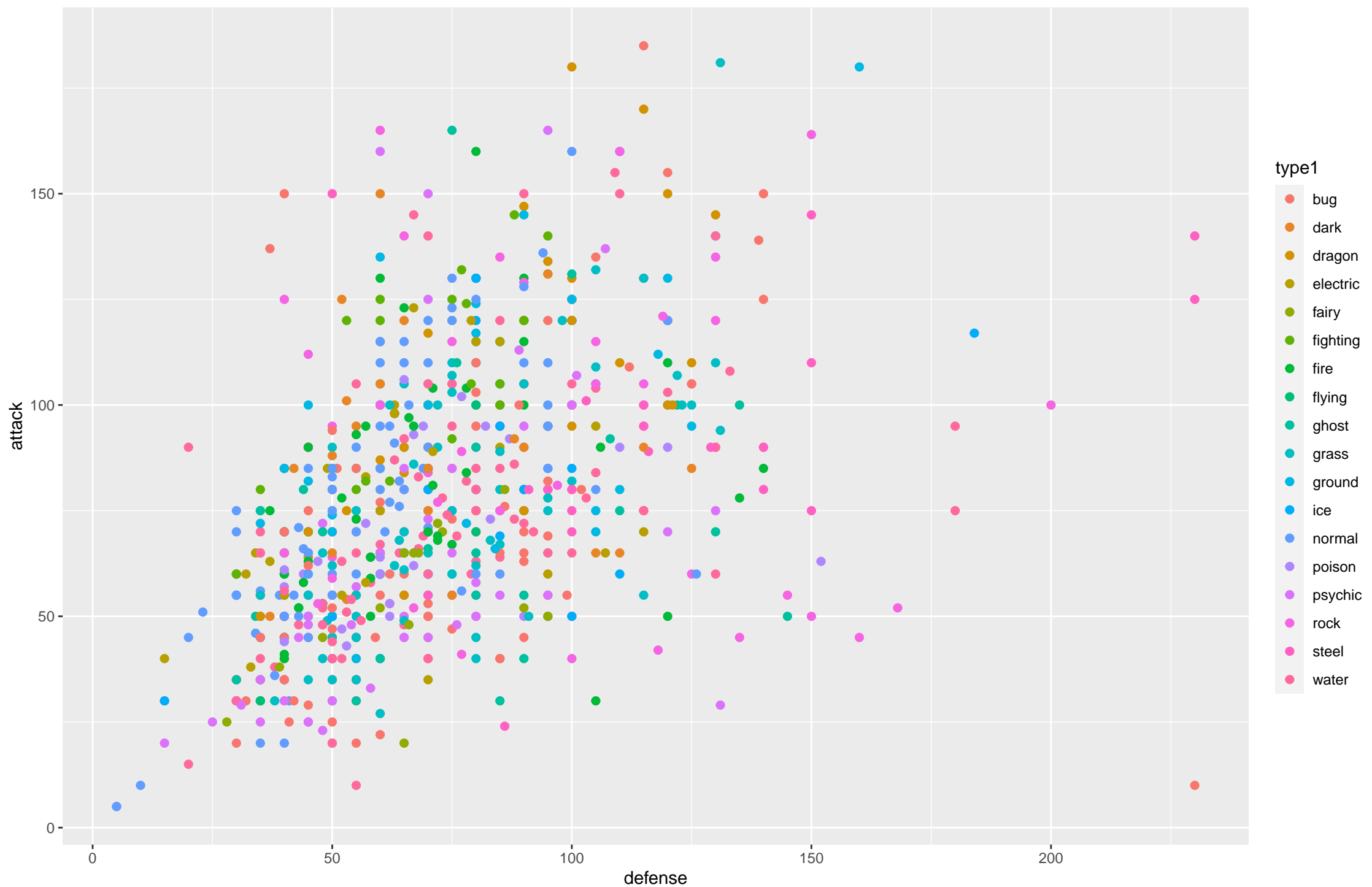
```
6  gg4 <- ggplot(newimdb, aes(x=Genre, y=Gross, fill=Genre)) +
7    geom_boxplot() +
8    xlab("Genre") +
9    theme(legend.position="none") +
0    xlab("") +
1    xlab("")+
2
3    labs(title="Highest Grossing films by Genre")+
4
5  scale_y_continuous(breaks=seq(0, 1000000000,200000000),
6                   labels = function(x){paste0(x/1000000,'Million')})
7
```

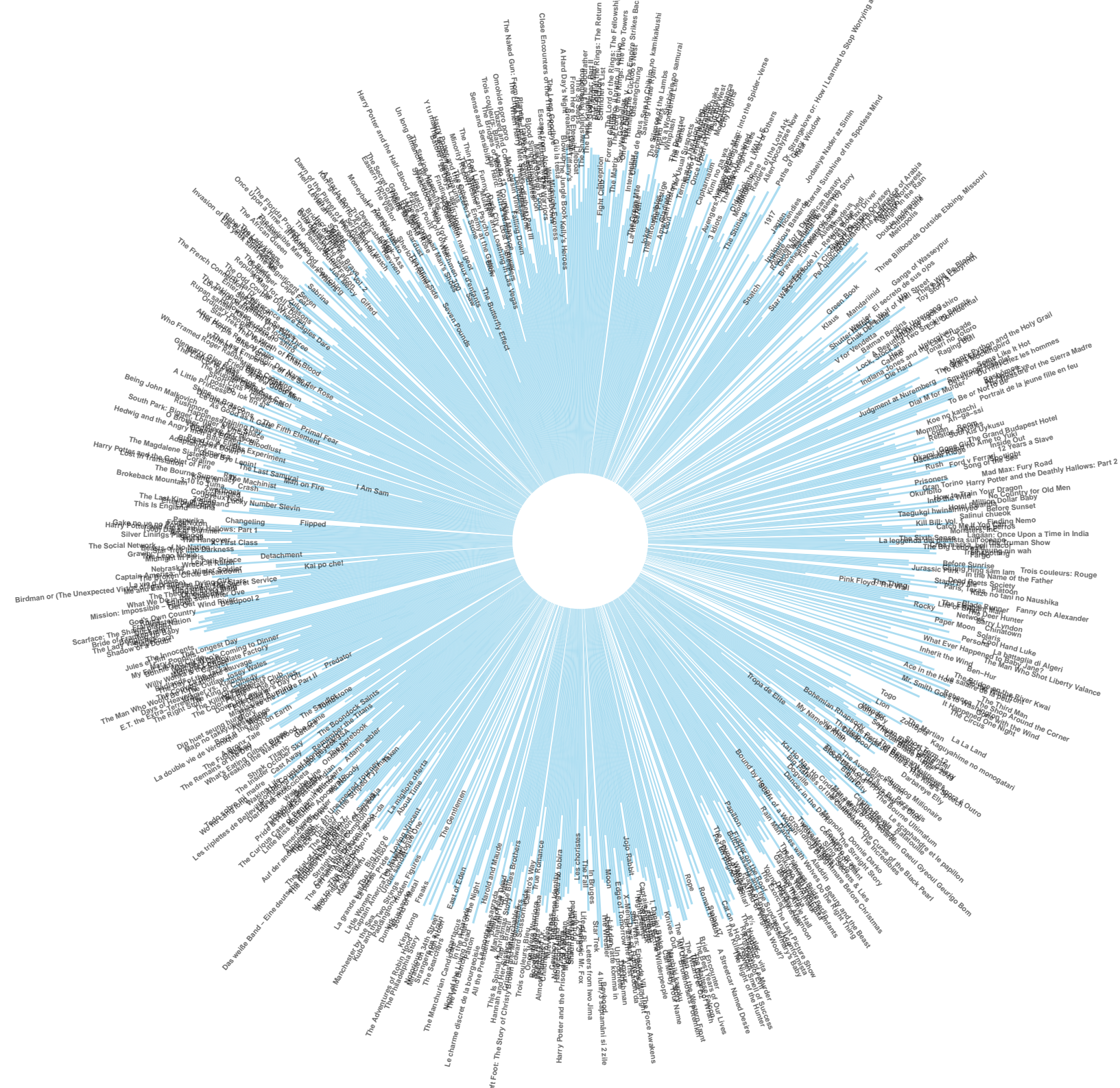Highest Grossing films by Genre

I hope this has been an enlightening look into how successful different genres of film have been throughout the world. It was incredibly interesting trying out different techniques and plotting styles. I just wish I'd had more time and was able to experiment more. This was a fun start though.

# Pokemon
## Attack vs Defence



Pokemon Demographics
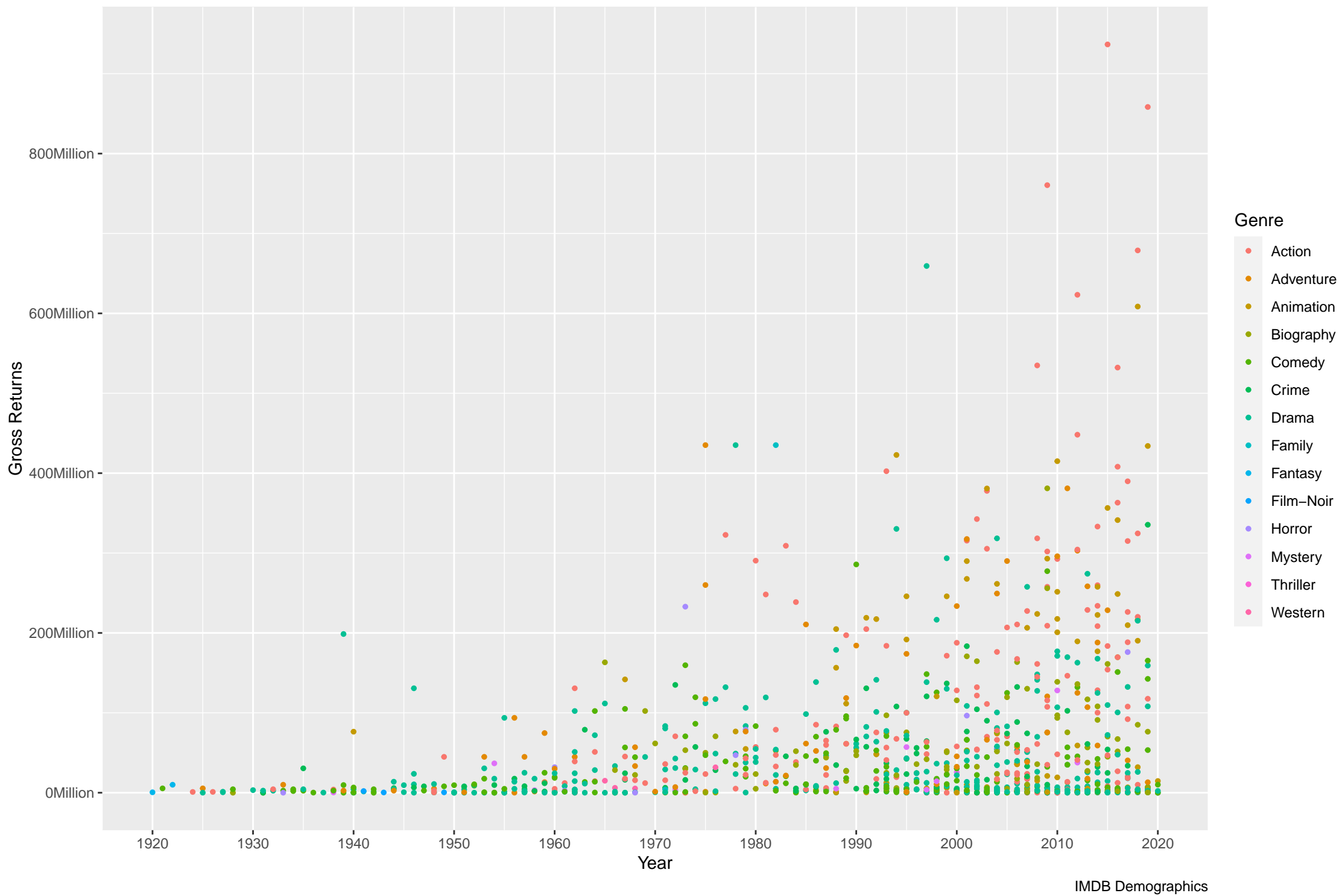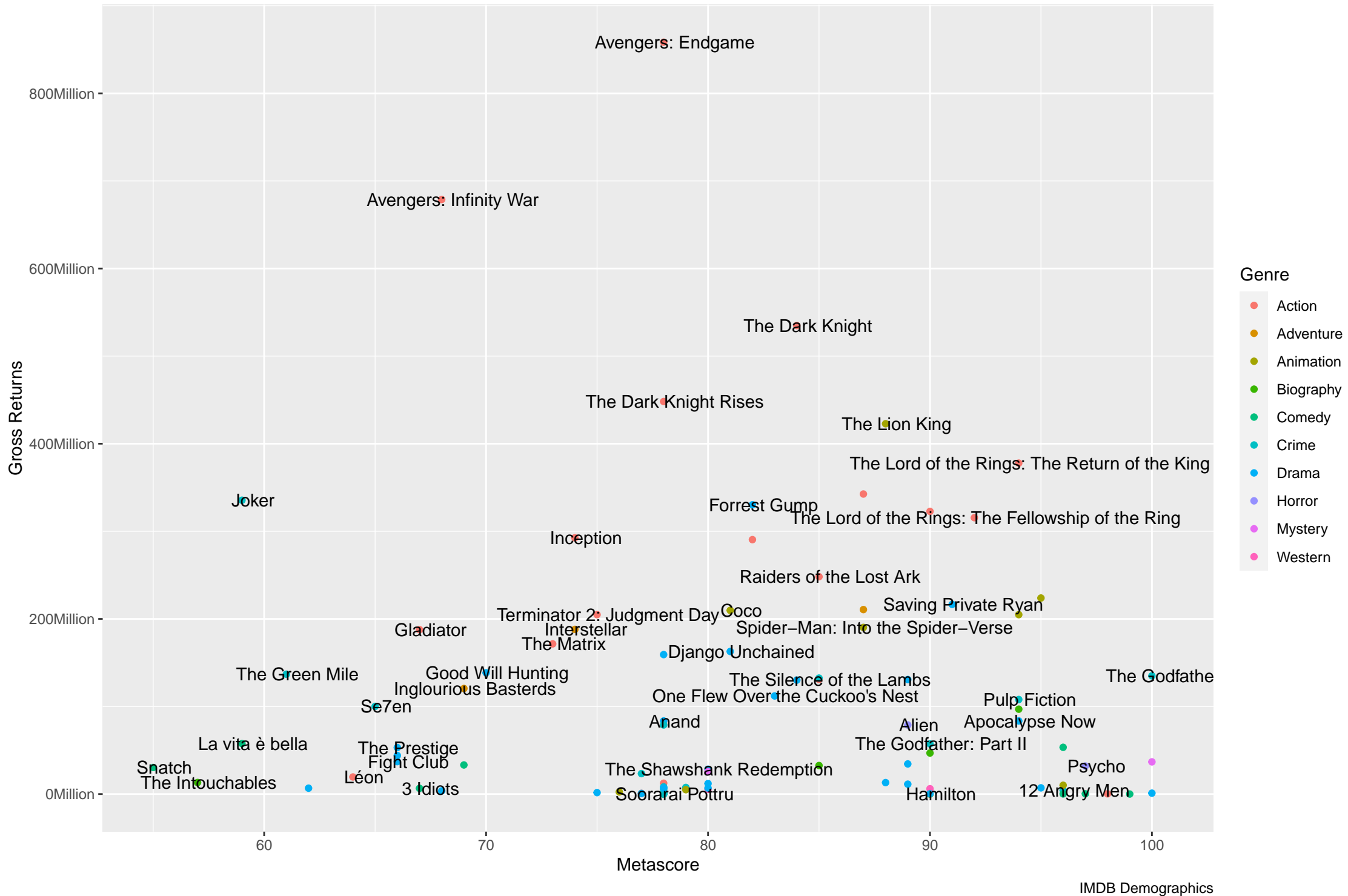
Better Circular Graph

# IMDB top 100 Gross Returns
From imdb dataset



**Genre**
- Action
- Adventure
- Animation
- Biography
- Comedy
- Crime
- Drama
- Family
- Fantasy
- Film−Noir
- Horror
- Mystery
- Thriller
- Western

IMDB Demographics

# IMDB top 1000
## highest 100 imdb score



**Gross Returns** (y-axis): 0Million, 200Million, 400Million, 600Million, 800Million

**Metascore** (x-axis): 60, 70, 80, 90, 100

**Genre**
- Action
- Adventure
- Animation
- Biography
- Comedy
- Crime
- Drama
- Horror
- Mystery
- Western

Avengers: Endgame
Avengers: Infinity War
The Dark Knight
The Dark Knight Rises
The Lion King
The Lord of the Rings: The Return of the King
Joker
Forrest Gump
The Lord of the Rings: The Fellowship of the Ring
Inception
Raiders of the Lost Ark
Terminator 2: Judgment Day
Coco
Saving Private Ryan
Gladiator
Interstellar
Spider-Man: Into the Spider-Verse
The Matrix
Django Unchained
The Green Mile
Good Will Hunting
The Silence of the Lambs
The Godfathe
Inglourious Basterds
One Flew Over the Cuckoo's Nest
Pulp Fiction
Se7en
Anand
Alien
Apocalypse Now
La vita è bella
The Godfather: Part II
The Prestige
Psycho
Snatch
Fight Club
The Shawshank Redemption
The Intouchables
Léon
3 Idiots
Hamilton
12 Angry Men
Soorarai Pottru

IMDB Demographics

**Top 100 movies from scores to votes**

# Highest Grossing films by Genre