FIFA Shiny App

Project Group 1

Maximilian Becker Stefan Ilic Romano Girardi Tobias Schär

Introduction

As we discussed what data we want to analyse and which questions we want to ask ourselves we came to the following research questions we had in our project proposal.

- Which section of performances do typically decrease/increase with age?
- How is heading precision related to weight/height and other factors?
- How does salary and the overall rating correlate?
- Which positions do have which average stats (preferred performance factors)

This and an additional question will be answered in the following document. Furthermore, we explain how we did this in our Shiny app with the corresponding code excerpts.

Overview

As you will see, we opted for a tab layout and a bootstrap css we took from bootswatch.org.

Every code excerpt from this document from the server part (where the calculations are done) only. We opted out of showing the UI here since there are only the Inputs (i.e.: Slider) and text paragraphs anyway.

In the Tab "Fifa Tabelle" is the raw dataset where you can get a general feeling what data we're dealing with.

Code:

```
fifaDataRaw <- read.csv("data/data.csv",stringsAsFactors=FALSE,encoding="UTF-8")
fifaDataRaw <- as_tibble(fifaDataRaw)</pre>
fifaDataRaw <- fifaDataRaw %>%
    {\tt separate(Height,\ c("HFeet",\ "HInches"),\ "'")}
fifaDataRaw <- fifaDataRaw %>%
    mutate(HeightMetric = ((as.numeric(HFeet)*30) + (as.numeric(HInches)*2.54)) )
fifaDataRaw <- fifaDataRaw %>%
    separate(Weight, c("Weight"), "lbs",convert=TRUE)
fifaDataRaw <- fifaDataRaw %>%
    mutate(WeightMetric = (as.numeric(Weight)/2.2) )
#Sperate the Value and Wage column with the euro sign and multiplicator (either M or K)
fifaDataRaw <- fifaDataRaw %>%
  separate(Wage,c(NA,"Wage"),"€",convert=TRUE)
fifaDataRaw <- fifaDataRaw %>%
  separate(Wage,c("Wage"),"K",convert=TRUE)
fifaDataRaw <- fifaDataRaw %>%
  separate(Value,c(NA,"Value"),"€",convert=TRUE)
fifaDataRaw <- fifaDataRaw %>%
  separate(Value,c("Value", "multiplicator"), "(?<=[0-9])(?=[A-Z])")
fifaDataRaw <- fifaDataRaw %>%
  mutate(
    Value = ifelse(multiplicator=="K", as.numeric(Value)*1000,
                   ifelse(multiplicator=="M",as.numeric(Value)*1000000,0))
newdata <- (fifaDataRaw %>% filter(Wage > quantile(Wage , 0.25 )))
newdata <- newdata[order(newdata$Wage),]</pre>
head(newdata)
print(newdata$Wage)
fifaDataRaw$WeightMetric <- as.numeric(as.character(fifaDataRaw$WeightMetric))</pre>
fifaDataRaw$HeightMetric <- as.numeric(as.character(fifaDataRaw$HeightMetric))</pre>
fifaDataRaw$HeadingAccuracy <- as.numeric(as.character(fifaDataRaw$HeadingAccuracy))</pre>
str(fifaDataRaw)
#Make Dataset Reactive
fifaData <- reactive({</pre>
  fifaData=fifaDataRaw
```

First the dataset is read and saved into a variable and converted to a Tibble.

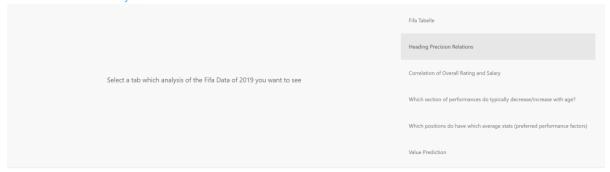
Then the columns for height is separated since the values were shown in an American format (i.e. 5'2") and converted into metric in a new column afterwards.

The same transformation was done for the weight.

Then, the Value had to be separated from the Euro-sign and the multiplicator had to be taken out for later calculations.

Q1: How is heading precision related to weight/height and other factors?





In order to visualize whether the heading performance is closely related to weight and height we need to see if these factors correlate. The easiest way to do that is by having a relationship between these factors. Luckily there already is an established factor for these values - it's called the BMI. So we added a column to our dataset and calculated each BMI for every player. The heading performance is on a scale from 0-100

Show 10 ~ entries			Search:	
	HeadingAccuracy 🔻	WeightMetric =	HeightMetric 🔷	вмі 💠
103	94	91.3636363636364	195.24	23.9682197947171
205	94	89.09090909091	192.7	23.9921938214764
500	93	86.36363636364	185.08	25.2122693675547
701	93	85	190.16	23.5061004801837
819	93	85	185.08	24.8141809038565
13	92	78.18181818182	185.08	22.8237385853653
317	92	78.18181818182	180	24.1301907968575
725	92	82.27272727273	190.16	22.7518940476645
9	91	82.2727272727273	180	25.3928170594837

Code (without the text paragraphs):

```
#Heading Performance
#Add BMI. BMI is a relationship between Weight and Height
headPerf <- fifaDataRaw[, c("HeadingAccuracy","WeightMetric","HeightMetric")]
headPerf <- headPerf %>% mutate(BMI = as.numeric(WeightMetric)/(as.numeric(HeightMetric)/100)^2)
#Output
output$headingPerfTable <- renderDataTable(</pre>
  headPerf
headPerf.agg <- aggregate(.~HeadingAccuracy,headPerf,FUN=mean,na.rm=TRUE,na.action=NULL)
#Output
output$headingPerfCorr <- renderPlot(</pre>
  ggcorr(headPerf.agg, label=TRUE, hjust=1, size=5)
#Output and gradient color
output$headingPerf <- renderPlot({</pre>
  palHP <- wes_palette("Zissou1", 100, type = "continuous")</pre>
  \texttt{ggplot}(\texttt{fifaData}(), \ \texttt{aes}(\texttt{x} = \texttt{WeightMetric}, \ \texttt{y} = \texttt{HeightMetric}, \ \texttt{fill} = \texttt{HeadingAccuracy})) \ + \\
    geom_tile() +
    scale_fill_gradientn(colours = palHP) +
    coord_equal()+
    xlab("Weight in kilogramms")+
    ylab("Height in meters")+
    geom_point(aes(x=input$weightM,y=input$heightM),size=3,colour="purple")
})
```

Here an additional column is added to the dataset. The column BMI is calculated and added for each row.

Then the table is aggregated to the HeadingAccuracy. We explicitly chose the mean function to have the outliers calculated too.

The correlation graph (ggcorr) is being put out.

The colour palette is being loaded (from a library) and the tile graph is created.

Q2: How does salary and the overall rating correlate?



Search:
Value \$
celona 110500000
rus 77000000
aint-Germain 118500000
nester United 72000000
nester City 102000000
93000000 93000000
ladrid 67000000
ch

Code:

```
#How does Value and the overall rating correlate?
output$valueTable <- renderDataTable(</pre>
  datatable(fifaData()[,c("Name","Age","Nationality","Overall","Club","Value")]
valueRat <- fifaDataRaw[,c("Name","Age","Nationality","Overall","Club","Value","Wage")]</pre>
valueRat <- valueRat %>%
  mutate( topTen=
     case_when
        Nationality == "Belgium" ~ TRUE,
Nationality == "France" ~ TRUE,
Nationality == "Brazil" ~TRUE,
        Nationality == "Brazil" ~TRUE,
Nationality == "England" ~TRUE,
Nationality == "Uruguay" ~TRUE,
Nationality == "Croatia" ~TRUE,
Nationality == "Portugal" ~TRUE,
Nationality == "Spain" ~TRUE,
Nationality == "Argentina" ~TRUE,
Nationality == "Colombia" ~TRUE,
Nationality == "Colombia" ~TRUE,
        TRUE ~ FALSE
valueRatR <- reactive({</pre>
  valueRatR=valueRat
Value.agg <- aggregate(.~Overall, valueRat[,c("Overall","Value","Wage")], FUN=mean, na.rm=TRUE, na.action=NULL)
output$corrGraphValue <- renderPlot(
  ggcorr(Value.agg,label=TRUE)
```

Here the raw data is first being put out with the relevant columns.

Then an additional column is added to either see from the nationality of a player whether it is a top ten country or not.

Afterwards, the data is being made reactive.

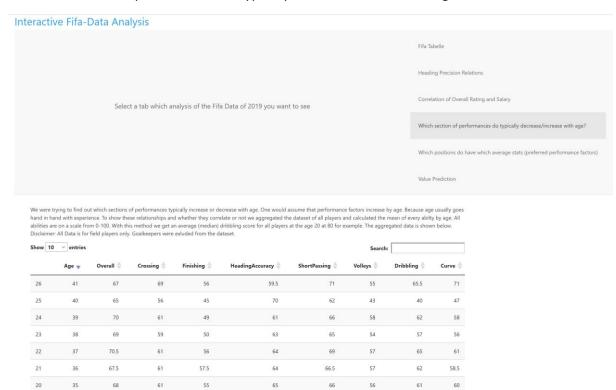
Then the data is being aggregated to the Overall-Rating.

```
output$ratingValue <- renderPlot({
       {\tt ggplot(valueRatR()~\%>\%~filter(Value>quantile(Value,input\$valueQuant/100,na.rm=~TRUE))}
                                   aes(x=Value,y=Overall,size=Wage,color=ifelse(topTen==TRUE,Nationality,"other"))) + (topTen==TRUE,Nationality,"other")) + (topTen==TRUE,Nationality,"other (topTen==TRUE,Nationali
                geom_point(alpha=0.5)
                scale_fill_viridis(discrete=TRUE, guide=FALSE, option="A")+
                theme_ipsum()+
                labs(colour="Fifa Top Ten (2020)",size="Wage in thousands")+scale_x_continuous(labels=comma,trans='log10')
output$densGraph <- renderPlot(</pre>
        ggplot(
                valueRatR()
               filter(
                       case_when(
                               input$selectValueWage == "Value" ~ Value>quantile(Value,input$valueQuant2/100,na.rm=TRUE),
                               input$selectValueWage == "Wage" ~ Wage>quantile(Wage,input$valueQuant2/100,na.rm=TRUE)
                 ,aes_string(input$selectValueWage))+
                geom_point(aes(y=Overall),alpha=0.5)+
theme_ipsum()+
               geom_density_2d(aes(y=0verall),colour="red")+
scale_x_continuous(labels=comma,trans='log10')
```

The plot is being generated with the chosen quantile input and the colours for the countries, the size for the wage and the points for the value.

The second graph is created with the chosen points and a density graph.

Q3: Which section of performances do typically decrease/increase with age?



Code:

```
#Question Which section of Performance...
#Omit Goalkeepers and select Age and Performance Factors
AgeStats <- fifaDataRaw %>% filter(Position != "GK")
#AgeStats <- AgeStats %>% filter(Age < quantile(Age , 0.90 ))</pre>
#Aggregate the mean of every ability by age
AgeStats.agg <- aggregate(.~Age, AgeStats, FUN=median, na.rm=TRUE, na.action=NULL)
#Output
output$AgeStatsTable <- renderDataTable(</pre>
 AgeStats.agg
#Correlation matrix
AgeStats.corr <-cor(AgeStats.agg)
#Output
output$AgeStatsCorr1 <- renderDataTable(</pre>
 AgeStats.corr
paletteHM1 = colorRampPalette(brewer.pal(3, "BuGn"))(20)
output$HeatmapCorr1 <- renderPlot(</pre>
 heatmap(x = AgeStats.corr, col = paletteHM1, symm = TRUE)
#print(AgeStats[order(-AgeStats$Age),])
#Correlation Matrix only by age, splitting the dataframe and putting
#it back together as tibble in order to create the geom_point graph
correlationAge <- AgeStats.corr[,'Age']</pre>
colAge <- colnames(AgeStats.corr)</pre>
Age.corr <- tibble(colAge, correlationAge)
output$CorrelationGraph1 <- renderPlot(</pre>
ggplot(data=Age.corr) +
  geom_point(aes(x=colAge, y=correlationAge)) + geom_hline(yintercept=0, color="blue", size=2)+
  theme(axis.text.x = element_text(size = 20, angle = 45, hjust = 1)) +
  xlab("Ability")+
 ylab("Corrleation Value")
```

Again, the dataset is being reduced for performance and only the relevant columns.

There is a lot going on here but basically the data is being aggregated again and put into correlation matrices and correlation graphs are being created.

```
#Output
output$CorrelationGraph2 <- renderPlot(</pre>
ggcorr(AgeStats.agg, label=TRUE, hjust=1, size=5)
#The mean of every ability by 'Overall' stat
OveStats.agg <- aggregate(.~Overall, AgeStats, FUN=median, na.rm=TRUE, na.action=NULL)
#Output
output$OveStatsTable <- renderDataTable(</pre>
    OveStats.agg
#Correlation matrix and heatmap
OveStats.corr <-cor(OveStats.agg)</pre>
output$OveStatsCorr1 <- renderDataTable(</pre>
    OveStats.corr
# #Output
# paletteHM2 = colorRampPalette(brewer.pal(3, "BuGn"))(20)
# output$HeatmapCorr2 <- renderPlot(</pre>
# heatmap(x = OveStats.corr, col = paletteHM2, symm = TRUE)
# )
#Output
output$CorrelationGraph3 <- renderPlot(</pre>
    ggcorr(OveStats.agg,label=TRUE,hjust=1,size=5)
print(OveStats.corr)
#Correlation data by Overall and create tibble again in order to create the geom_point
correlationOve <- OveStats.corr[,'Overall']</pre>
col <- colnames(OveStats.corr)</pre>
Ove.corr <- tibble(col, correlationOve)</pre>
#Output
output$CorrelationGraph4 <- renderPlot(</pre>
     ggplot(data=Ove.corr) +
         \label{eq:geom_point} geom\_point(aes(x=col, y=correlationOve)) + geom\_hline(yintercept=0, color="blue", size=2) + geom\_hline(yinterce
         theme(axis.text.x = element_text(size = 20, angle = 45, hjust = 1)+
         xlab("Ability")+
         ylab("Correlation Value")
)
#Output
output$AgeOverall <- renderPlot(</pre>
ggplot(AgeStats.agg,aes(x=Age,y=Overall))+
    geom_line()+
    geom_point(aes(x=AgeStats.agg[which.max(Overall),1],
                            y=AgeStats.agg[which.max(Overall),2]),
col="red",size=3,show.legend=TRUE)+
    )
```

Q4: Which positions do have which average stats (preferred performance factors)

		0 (1		'			
Interactive Fifa-Data	Analysis						
	Select a tab which analysis of the Fifa Data of 2019 you want to see			Fifa Tabelle			
			Heading Precision Relations				
		ant to see	ee	Correlation of Overall Rating and Salary			
				Which section of performances do typically decrease/increase with age?			
				Which positions do have which average stats (preferred performance factors)			
				Value Prediction			
To see who in which position has which average stats, we first had to have an underlying image of a football field. Then we had to extract the data and aggregate it in order to have the positions and enrich them with the coordinates on the football field. Then, each position is given a unique colour. Now, you can choose, if you want to see the mean of the overall rating, the players value or each unique ability of every player in that certain position. Or the number of players in that positions from our given data frame. Since the descriptions of each position would be too long, we added them below							
Select a value to be displayed a	at the corresponding positions						
Overall				<i>\</i> \$			
Value Crossing							
Crossing Finishing							
○ HeadingAccuracy							
○ ShortPassing							
○ Volleys							
Oribbling							
O Curve							
Count							
Code:							
###########	#######################	##					
##Stats on	preferred position						
posData <-	fifaDataRaw[,c("Positi	on","Overal	1","va	alue","Crossing",			
•	fifaDataRaw[,c("Positi "Finish	ina"."Headi	naÁcci	uracv"."ShortPassing".			
	"Vollev	s","Ďribbli	na"."(Curve"			
)]	,	,				
	7.3						
#Adding count column for later use in aggregat posData <- posData %>% mutate(count=1)							
	occurences of each pos						
	p1 <- aggregate(count		posDa	ata,sum,na.rm=TRUE)			
000	the mean of each other						
	p2 <- aggregate(. ~ Po	sition,posD	ata,me	ean)			
	mpty positions						
posData.tem	p1 <- posData.temp1 %>	% filter(Po	sition	n != "")			
#set count	column to null						
	p2\$count <- NULL						
	e aggregated dataframe	S					
	<- merge(posData.temp		emp2)				
1 9 9	9 - 4	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,					

- Only relevant Data
- The data is being counted and filtered

```
posData.agg <- posData.agg %>%
                                           posCoordY =
  mutate(
                                             case_when(
    posCoordX =
                                               Position == "LF" \sim 160,
      case_when(
                                               Position == "LB" ~ 80,
        Position == "LF" ~ 20,
                                               Position == "LCB" ~ 60,
        Position == "LB" ~ 20.
                                               Position == "LCM" \sim 100,
        Position == "LCB" \sim 40,
                                               Position == "LAM" ~ 120.
        Position == "LCM" ~ 40,
                                               Position == "LDM" \sim 80,
        Position == "LAM" ~ 40,
                                               Position == "LM" ~ 100,
        Position == "LDM" \sim 40,
                                               Position == "LS" ~ 140,
        Position == "LM" ~ 20,
                                               Position == "LW" \sim 150,
        Position == "LS" ~ 20,
                                               Position == "LWB" \sim 60,
        Position == "LW" ~ 20,
                                               Position == "RB" ~ 80,
        Position == "LWB" ~ 20,
                                               Position == "RWB" ~ 60,
        Position == "RB" ~ 80,
                                               Position == "RCB" ~ 80,
        Position == "RWB" ~ 80,
                                               Position == "RCM" \sim 100,
Position == "RAM" \sim 120,
        Position == "RCB" ~ 60,
        Position == "RCM" \sim 60,
                                               Position == "RDM" ~ 60,
        Position == "RAM" ~ 60,
                                               Position == "RM" \sim 100,
        Position == "RDM" ~ 60,
                                               Position == "RS" ~ 140,
        Position == "RM" \sim 80,
                                               Position == "RF" \sim 160,
        Position == "RS" \sim 80,
                                               Position == "RW" ~ 150,
        Position == "RF" ~ 80,
                                               Position == "CF" \sim 160,
        Position == "RW" ~ 80,
                                               Position == "CAM" ~ 120,
        Position == "CF" ~ 50,
                                               Position == "CB" ~ 80,
        Position == "CAM" \sim 50,
                                               Position == "CDM" \sim 80,
        Position == "CB" ~ 50,
        Position == "CDM" \sim 50,
                                               Position == "CM" \sim 100,
                                               Position == "ST" \sim 160.
        Position == "CM" \sim 50,
                                               Position == "GK" ~ 20,
        Position == "ST" ~ 50,
                                               TRUE \sim 0
        Position == "GK" \sim 50,
                                             )
        TRUE \sim 0
                                         )
      )
```

 Here, the positions and their respective coordinates are being added to the dataset for later usage on the football pitch image/graph

```
field <- readPNG("data/field2.png")</pre>
# grid <- rasterGrob(field, width=unit(1,"npc"), height=unit(1,"npc"))</pre>
grid <- rasterGrob(field, interpolate=TRUE, height = 1, width = 1)
# ggplot(posData.agg, aes(posCoordX,posCoordY,color=Position)) +
     annotation_custom(grid) +
    geom point(aes(size=Overall))
    scale_x_continuous(expand=c(0,0), lim=c(0,100)) + scale_y_continuous(expand=c(0,0), lim=c(0,200)) +
     theme_void() +
    theme(aspect.ratio = nrow(field)/ncol(field))+
    scale_fill_viridis(discrete=TRUE, guide=FALSE, option="A")+
guides(color=guide_legend(order=1),
            size=guide_legend(order=2))
#Make aggregated Data reactive
posDataR <- reactive({
 posDataR=posData.agg
#make input reactive
fieldChoice <- reactive(input$fieldChoice)</pre>
#make graph reactive for later if statement (scaling of values)
graph <- reactive({
  ggplot(posDataR(), aes(posCoordX,posCoordY,color=Position)) +
     annotation_custom(grid)
     geom_point(aes_string(size=fieldChoice()))
     scale_x\_continuous(expand=c(0,0), lim=c(0,100))
     scale_y = continuous(expand=c(0,0), lim=c(0,200)) +
     theme_void() +
    theme(aspect.ratio = nrow(field)/ncol(field),legend.key.size = unit(1, "cm"),legend.key.width = unit(1, "cm"))+ # scale_fill_viridis(discrete=TRUE, guide=FALSE, option="A")+
     guides(color=guide_legend(order=1),
            size=guide_legend(order=2))
})
```

- Here the picture is being loaded for later usage in the graph
- The graph is being created with the custom grid and the picture
- The positions are added and the size of a ball is the chosen value

```
output$avgStatsField <- renderPlot({</pre>
      if(fieldChoice()=="Value"){
           graph()+
                scale\_size(range = c(0.1, 20), labels=comma)
          graph()+
                 scale\_size(range = c(0.1, 15))
})
#Input Position
choicePosR <- reactive({
  choicePosR=input$choicePosition
#Data for Spiderchart and Rownames
spidDat <-
     posData.agg[,c("Overall","Crossing","Finishing","HeadingAccuracy","ShortPassing","Dribbling","Curve")]
 rownames(spidDat) <- posData.agg$Position
spidDat
#First two lines have to have minimun and maximum values for spiderchart
spidDat \leftarrow rbind(rep(100,7), rep(0,7), spidDat)
spidDatR <- reactive({</pre>
     spidDatR = spidDat
colors_borderSpid <- brewer.pal(3, "Set3"</pre>
colors_inSpid <- alpha(colors_borderSpid,0.5)</pre>
output$posSpider <- renderPlot({</pre>
      radarchart(spidDatR()%>%filter(rownames(spidDatR()) %in% c(1,2,choicePosR())), axistype=1 ,
                                  #custom polygon
                                  pcol=colors_borderSpid , pfcol=colors_inSpid , plwd=4 , plty=1,
                                  #custom the grid cglcol="grey", cglty=1, axislabcol="grey", caxislabels=seq(0,10,100), cglwd=0.8,
                                  #custom labels
                                  vlcex=0.8)
                                   # Add a legend
                                  legend(x=-\tilde{1}.9, y=0, legend = choicePosR(), bty = "n", pch=20 , col=colors\_borderSpid , text.col = "green text.col =
})
```

- The spider graph is being created with the colour palette and the respective positions.
- Each Positions can be chosen to be added

Value Prediction

This was a tricky one. We chose a linear regression model and added the most correlated variables one by one and things got tricky:

```
######--- Prediction Value
# Create training and test data
set.seed(69)
forecData <- fifaDataRaw[,c("Age","Value","Crossing","Finishing","HeadingAccuracy","ShortPassing","Vol
##Remove Value = NULL
forecData <- forecData %>% filter(Value!="")
##Remove outliers
outliers <- boxplot(forecData$Value, plot = TRUE)$out</pre>
x <- forecData
forecData <- x[-which(x$Value %in% outliers),]</pre>
#Train and Test Data
train.data<- sample_frac(forecData,0.7) #select 70% random samples</pre>
test.data <- setdiff(forecData,train.data)</pre>
#Create Correlation Table
forecData.corr<- cor(forecData,use = "complete.obs")</pre>
# view(forecData.corr[,c("Value")])
#View Correlations
ggcorr(forecData.corr,label=TRUE)
```

- Here again, only the relevant data. We wanted to find out which ability correlates most with the value and with those extrapolate a player's value.
- Here the massive outliers of the value have to be taken out in order to have usable data.

#wage would be the the most correlated, but since wage and value are naturally

• Train and test data is being created.

#strating with the most correlated variable as predictor

#dependend on each other (since the more value, the more wage)

Correlation table is calculated.

```
#and it would be meaningless to predict ones value from the wage we go from the 2nd highest
lm1 <- lm(Value~ShortPassing,data=train.data )</pre>
summary(1m1)
##All Coefficients are positive
C:/HSLU/DASB/FIFA_Project/ShinyFifaApp/
lm(formula = Value ~ ShortPassing, data = train.data)
Residuals:
     Min
                1Q
                     Median
                                    3Q
                                            Max
          -598332 -310361
-1288042
                               196813
                                        4921090
Coefficients:
               Estimate Std. Error t value Pr(>|t|)
                                               <2e-16 ***
(Intercept) -565401.8
                            37374.4
                                     -15.13
                                               <2e-16 ***
                26485.5
                              639.4
                                       41.42
ShortPassing
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 936700 on 10837 degrees of freedom
  (35 observations deleted due to missingness)
Multiple R-squared: 0.1367,
                                  Adjusted R-squared: 0.1366
F-statistic: 1716 on 1 and 10837 DF, p-value: < 2.2e-16
```

As soon as we added more variables, things got tricky:

In the next model, all variables were positive. But as soon as we added more, some became negative coefficients:

```
0.77
           lm2 <- lm(Value~ShortPassing+Curve,data=train.data)</pre>
  698
  699
           summary(1m2)
           ##All Coefficients are positive
  700
  701
  702
           lm3 <- lm(Value~ShortPassing+Curve+Dribbling,data=train.data)</pre>
  703
       # (Untitled) $
 689:5
Console
       Terminal × Jobs ×
C:/HSLU/DASB/FIFA_Project/ShinyFifaApp/ 🔊
      summary(1m2)
Call:
lm(formula = Value ~ ShortPassing + Curve, data = train.data)
Residuals:
     Min
               1Q
                    Median
                                  3Q
                                          Max
-1372195 -600885 -307900
                              198211 4854608
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
                          37802.7 -13.659 < 2e-16 ***
(Intercept)
            -516342.2
               20749.7
                                   21.278 < 2e-16 ***
ShortPassing
                             975.2
                                     7.774 8.27e-15 ***
                6132.6
                             788.8
Curve
Signif. codes: 0 '***' 0.001 '**' 0.05 '.' 0.1 ' ' 1
Residual standard error: 934100 on 10836 degrees of freedom
  (35 observations deleted due to missingness)
Multiple R-squared: 0.1415,
                                Adjusted R-squared: 0.1413
F-statistic: 892.7 on 2 and 10836 DF, p-value: < 2.2e-16
```

```
702
          lm3 <- lm(Value~ShortPassing+Curve+Dribbling,data=train.data)</pre>
 703
          summary(1m3)
  704
          ##The Dribbling Coefficient is negative. Which means greater value effects the value negativels
 705
706
689:5 # (Untitled) $
Console Terminal × Jobs ×
C:/HSLU/DASB/FIFA_Project/ShinyFifaApp/ > Summary(Imb)
Call:
lm(formula = Value ~ ShortPassing + Curve + Dribbling, data = train.data)
Residuals:
              1Q Median
067 -296076
                          196462 4850485
-1431066 -592067
Coefficients:
             Estimate Std. Error t value Pr(>|t|)
                        37943.9 -14.083 < 2e-16 ***
(Intercept)
            -534347.9
ShortPassing
              23823.4
                         1160.9 20.521 < 2e-16 ***
                          944.3
               8665.6
                                  9.176
                                        < 2e-16 ***
Dribbling
              -5074.9
                         1042.6 -4.867 1.15e-06 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 933100 on 10835 degrees of freedom
  (35 observations deleted due to missingness)
Multiple R-squared: 0.1433,
                              Adjusted R-squared: 0.1431
F-statistic: 604.3 on 3 and 10835 DF, p-value: < 2.2e-16
  706
            lm4 <- lm(Value~ShortPassing+Curve+Dribbling+Volleys,data=train.data)</pre>
  707
            summary(1m4)
  708
            ##Dribbling Coeficcient is still negative
  709
  710
            lm5 <- lm(Value~ShortPassing+Curve+Dribbling+Volleys+Crossing,data=train.data)</pre>
  711
            summary(1m5)
  712
 708:5
       # (Untitled) $
Console Terminal × Jobs ×
 C:/HSLU/DASB/FIFA_Project/ShinyFifaApp/ 🗇
       summary(1m4)
call:
lm(formula = Value ~ ShortPassing + Curve + Dribbling + Volleys,
    data = train.data)
Residuals:
     Min
                 1Q
                      Median
                                     3Q
                                               Max
-1518744
          -590341
                     -292105
                                 190262 4832977
Coefficients:
                Estimate Std. Error t value Pr(>|t|)
              -546563.1
                             37915.5 -14.415 < 2e-16 ***
(Intercept)
ShortPassing
                 24077.8
                              1159.3 20.769
                                                 < 2e-16 ***
                                        5.969 2.47e-09 ***
                  6090.7
                              1020.4
Curve
                 -7796.8
                                       -6.963 3.53e-12 ***
Dribbling
                               1119.8
                                        6.583 4.84e-11 ***
Volleys
                  6331.9
                                961.9
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 931300 on 10834 degrees of freedom
  (35 observations deleted due to missingness)
                                    Adjusted R-squared: 0.1464
Multiple R-squared: 0.1467,
F-statistic: 465.8 on 4 and 10834 DF, p-value: < 2.2e-16
```

- And so on.
- So we played around with the variables until we got a model with the most variables where none of them is a negative coefficient:

```
playing around with the least correlated variables and the variable which were a negative coefficient earli
  719
  720
               lm.red <- lm(Value~.-HeightMetric-WeightMetric-Age-Dribbling-Crossing-Finishing,data=train.data)
  721
               summary(1m.red)
  722
               ##Here HeadingAccuracy, ShortPassing, Volleys and Curve are coefficient and are positive
  723
  724
               #let's cpompare the models obtained, in a structured way
  725
              anova(lm1, lm2, lm3, lm4,lm5, lm.red , lm.tot)
  726 4
                                                                                                                                                       R Script
 712:13 ## (Untitled) $
 Console Terminal × Jobs ×
 C:/HSLU/DASB/FIFA_Project/ShinyFifaApp/
lm(formula = Value ~ . - HeightMetric - WeightMetric - Age -
Dribbling - Crossing - Finishing, data = train.data)
Residuals:
Min 1Q Median 3Q Max
-1365589 -600351 -316821 209310 4846585
Coefficients:
                      Estimate Std. Error t value Pr(>|t|)
                    -529069.9
1812.0
                                    38016.9 -13.917 < 2e-16 ***
709.8 2.553 0.010704 *
1145.5 16.103 < 2e-16 ***
(Intercept)
HeadingAccuracy
ShortPassing
                       18446.5
Volleys
                                                   3.619 0.000297 ***
                         3327.6
                                        919.5
                                                 4.314 1.62e-05 ***
                         4251.6
                                        985.6
Curve
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 933100 on 10834 degrees of freedom
(35 observations deleted due to missingness)
Multiple R-squared: 0.1434, Adjusted R-squared: 0.1431
F-statistic: 453.6 on 4 and 10834 DF, p-value: < 2.2e-16
```

According to the anova function, this model is among the best(nr. 6)

```
Model 1: Value ~ ShortPassing
Model 2: Value ~ ShortPassing + Curve
Model 3: Value ~ ShortPassing + Curve + Dribbling
Model 4: Value ~ ShortPassing + Curve + Dribbling + Volleys
Model 5: Value ~ ShortPassing + Curve + Dribbling + Volleys + Crossing
Model 6: Value ~ (Age + Crossing + Finishing + HeadingAccuracy + ShortPassing +
    Volleys + Dribbling + Curve + HeightMetric + WeightMetric) -
    HeightMetric - WeightMetric - Age - Dribbling - Crossing -
    Finishing
Model 7: Value ~ Age + Crossing + Finishing + HeadingAccuracy + ShortPassing +
    Volleys + Dribbling + Curve + HeightMetric + WeightMetric
                         Sum of Sq
  Res.Df
                RSS Df
                                              Pr(>F)
  10837 9.5076e+15
   10836 9.4549e+15 1 5.2737e+13 63.6286 1.654e-15 ***
   10835 9.4342e+15 1 2.0628e+13 24.8884 6.169e-07 ***
   10834 9.3967e+15 1 3.7582e+13 45.3429 1.738e-11 ***
   10833 9.3929e+15 1 3.7710e+12 4.5498
                                           0.03295 *
6
   10834 9.4330e+15 -1 -4.0148e+13 48.4390 3.604e-12 ***
   10828 8.9746e+15 6 4.5845e+14 92.1879 < 2.2e-16 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
>
```

So we chose this one and added it to our calculations.