Regression Project Amirali Khatib 7/3/2022 Introduction In the this classification and regression project, we'll work with players\_csv file which consist of detailed attribute for each football player, we're following 2 purpose in this project. Regression goal: We fit regression models to our data in order to predict players' value in the transfer market. So, football clubs' owners will know how much they should pay to sign a contract with a player that has a special combination of skills. Classification goal: Our purpose in fitting classification models to our data is to build a high performance classifier that be able to recommend the head coach the squad position for each player with specific sort of abilities. Dataset features: To know some more detailes and explanation about the features that exist in the data set, I recommend you to check these two link bellow: Attributes expalation • Example : Luka Modric attributes **Intialization** Import required libraries library(tidyverse) ## Warning: package 'tidyverse' was built under R version 4.1.3 ## -- Attaching packages ----- tidyverse 1.3.1 --## v ggplot2 3.3.6 v purrr 0.3.4 ## v tibble 3.1.7 v dplyr 1.0.9 ## v tidyr 1.2.0 v stringr 1.4.0 ## v readr 2.1.2 v forcats 0.5.1 ## Warning: package 'ggplot2' was built under R version 4.1.3 ## Warning: package 'tibble' was built under R version 4.1.3 ## Warning: package 'tidyr' was built under R version 4.1.3 ## Warning: package 'dplyr' was built under R version 4.1.3 ## Warning: package 'forcats' was built under R version 4.1.3 ## -- Conflicts ----- tidyverse\_conflicts() --## x dplyr::filter() masks stats::filter() ## x dplyr::lag() masks stats::lag() library(MASS) ## Attaching package: 'MASS' ## The following object is masked from 'package:dplyr': ## select library(caret) ## Warning: package 'caret' was built under R version 4.1.3 ## Loading required package: lattice ## Attaching package: 'caret' ## The following object is masked from 'package:purrr': lift ## library(rpart) library(rpart.plot) ## Warning: package 'rpart.plot' was built under R version 4.1.3 library(knitr) ## Warning: package 'knitr' was built under R version 4.1.3 library(kableExtra) ## Warning: package 'kableExtra' was built under R version 4.1.3 ## Attaching package: 'kableExtra' ## The following object is masked from 'package:dplyr': ## group\_rows library(e1071) ## Warning: package 'e1071' was built under R version 4.1.3 library(nnet) library(pROC) ## Warning: package 'pROC' was built under R version 4.1.3 ## Type 'citation("pROC")' for a citation. ## Attaching package: 'pROC' ## The following objects are masked from 'package:stats': ## cov, smooth, var library(ellipse) ## Warning: package 'ellipse' was built under R version 4.1.3 ## Attaching package: 'ellipse' ## The following object is masked from 'package:graphics': ## pairs library(lares) ## Warning: package 'lares' was built under R version 4.1.3 ## Attaching package: 'lares' ## The following object is masked from 'package:e1071': ## ## impute library(reshape2) ## Warning: package 'reshape2' was built under R version 4.1.3 ## Attaching package: 'reshape2' ## The following object is masked from 'package:tidyr': ## smiths library(ROSE) ## Warning: package 'ROSE' was built under R version 4.1.3 ## Loaded ROSE 0.0-4 Read dataset fifa = read.csv(file = 'D:\\Amirali\\University\\Applied Statistical Analysis\\R-tutorial\\Project\\players\_22.cs v', header = TRUE) Data pre-processing fifa = na.omit(fifa)fifa = fifa %>% dplyr::select(-position, everything()) %>% dplyr::select(-value\_eur, everything())  $fifa1.0_reg = fifa$ fifa1.0\_reg[,c("sofifa\_id", "short\_name", "potential", "club\_name", "league\_name", "club\_jersey\_number", "nationality\_n ame", "weight\_kg", "height\_cm")] <- list(NULL)</pre> colnames(fifa1.0\_reg) ## [1] "overall" "potential\_growth" ## [3] "age" "international\_reputation" ## [5] "crossing" "finishing" ## [7] "heading\_accuracy" "short\_passing" ## [9] "volleys" "dribbling" "long\_passing" ## [11] "curve" ## [13] "ball\_control" "sprint\_speed" ## [15] "agility" "reactions" "jumping" ## [17] "shot\_power" ## [19] "stamina" "strength" ## [21] "long\_shots" "aggression" ## [23] "interceptions" "positioning" ## [25] "vision" "composure" ## [27] "marking\_awareness" "standing\_tackle" "gk\_diving" ## [29] "sliding\_tackle" ## [31] "gk\_handling" "gk\_kicking" "gk\_reflexes" ## [33] "gk\_positioning" "value\_eur" ## [35] "position" Data spliting set.seed(23) # test-train split idx = sample(nrow(fifa1.0\_reg), size = 0.8 \* nrow(fifa1.0\_reg)) fifa\_trn = fifa1.0\_reg[idx,] fifa\_tst = fifa1.0\_reg[-idx,] # estimation-validation split idx1 = sample(nrow(fifa\_trn), size = 0.8 \* nrow(fifa\_trn)) fifa\_est = fifa\_trn[idx1,] fifa\_val = fifa\_trn[-idx1,] # check data length(fifa1.0\_reg) ## [1] 36 unique(fifa\$position) "midfielder" "goalkeeper" "defender" ## [1] "forward" fifa\_est\_c = fifa\_est fifa\_est\_c[,c('position')] <- list(NULL)</pre> cormat = cor(fifa\_est\_c) cormat = round(cormat, 2) **EDA** melted\_cormat = melt(cormat) ggplot(data = melted\_cormat, aes(x=Var1, y=Var2, fill=value)) + geom\_tile()+ theme(axis.text.x = element\_text(angle = 90))+ scale\_fill\_gradient2(low = "blue", high = "red", mid = "white", midpoint = 0, limit = c(-1,1), space = "Lab", name="Pearson\nCorrelation") composure -Pearson Correlation 0.5 0.0 -0.5 -1.0 crossing international\_reputation age potential\_growth overall Var1 corr\_cross(fifa\_est, max\_pvalue = 0.05, top = 30) ## Returning only the top 30. You may override with the 'top' argument ## Warning in .font\_global(font, quiet = FALSE): Font 'Arial Narrow' is not ## installed, has other name, or can't be found **Ranked Cross-Correlations** 30 most relevant .25 .5 .75 0 short\_passing + long\_passing long\_shots + positioning volleys + long\_shots overall + reactions ball control + positioning potential\_growth + age crossing + dribbling volleys + positioning Correlations with p-value < 0.05 corr\_var(fifa\_est, value\_eur, top = 30) **Correlations of value\_eur** 30 largest correlation variables (original & dummy) .4 international\_reputation overall .557 reactions composure vision short\_passing .304 .299 shot\_power long passing .289 ball\_control .251 curve dribbling .242 .235 crossing volleys .231 long\_shots positioning finishing .228 .223 .212 .205 stamina aggression agility .18 .171 heading\_accuracy .166 sprint\_speed marking\_awareness interceptions .156 .155 standing\_tackle .135 strength .127 jumping .12 sliding\_tackle .118 potential\_growth position\_midfielder -.113 .0457  $ggplot(data = fifa_est, aes(x = position, y = value_eur, group = position, color = position))+$  $geom_point(size = 4)$ 2.0e+08 -1.5e+08 position 11.0e+08 defender forward goalkeeper midfielder 5.0e+07 -0.0e+00 defender forward goalkeeper midfielder table(fifa\_est\$position) ## defender forward goalkeeper midfielder ## 2897 454 1787 7127 **RMSE Function** rmse =function(predicted, actual){sqrt(mean((actual - predicted)^2))} Regression Models Linear models Training on estimation set lm\_model\_list = list( lm\_1 = lm(formula = value\_eur ~ ., data = fifa\_est),  $lm_2 = lm(formula = value_eur \sim .^2, data = fifa_est),$ lm\_3 = step(lm(formula = value\_eur ~ ., data = fifa\_est), trace = FALSE, direction = 'both'),  $lm_4 = lm(formula = value_eur \sim poly(overall, 2) + poly(age, 2) + poly(international_reputation, 2) + poly(international_reputa$ poly(finishing,2)+poly(long\_passing,2)+poly(ball\_control,2)+ poly(sprint\_speed, 2)+poly(strength, 2)+poly(interceptions, 2)+ poly(vision, 2)+poly(sliding\_tackle, 2)+poly(gk\_diving, 2), data = fifa\_est), lm\_5 = lm(formula = value\_eur ~ poly(overall, 3)+poly(age, 3)+poly(international\_reputation, 3)+ poly(finishing, 3)+poly(long\_passing, 3)+poly(ball\_control, 3)+ poly(sprint\_speed,3)+poly(strength,3)+poly(interceptions,3)+ poly(vision,3)+poly(sliding\_tackle,3)+poly(gk\_diving,3),data = fifa\_est)) Predict value and RMSE value of linear models fifa\_est\_lm\_predicted\_list = lapply(lm\_model\_list, predict, fifa\_est) fifa\_val\_lm\_predicted\_list = lapply(lm\_model\_list, predict, fifa\_val) fifa\_est\_lm\_rmse\_vector = sapply(fifa\_est\_lm\_predicted\_list, rmse, fifa\_est\$value\_eur) fifa\_val\_lm\_rmse\_vector = sapply(fifa\_val\_lm\_predicted\_list, rmse, fifa\_val\$value\_eur) rmse\_lm = data.frame(fifa\_est\_lm\_rmse\_vector, fifa\_val\_lm\_rmse\_vector) colnames(rmse\_lm) = c('Est','Val') rownames(rmse\_lm) =  $c('lm_1', 'lm_2', 'lm_3', 'lm_4', 'lm_5')$ rmse\_lm = data.matrix(rmse\_lm) rmse\_lm = melt(rmse\_lm) colnames(rmse\_lm) = c('model', 'data\_type', 'RMSE')  $ggplot(data = rmse_lm, aes(x = model, y = RMSE, group = data_type, color = data_type)) + geom_point(size = 2.5) + geom_$ line(size = 0.75)5000000 4000000 data\_type - Est **→** Val 3500000 -3000000 -2500000 lm\_3 lm\_1 lm\_2 lm\_4 lm\_5 model KNN models Feature scaling Estimation set fifa\_est\_scaled = data.frame(rep(0, nrow(fifa\_est))) center\_vec = vector() scale\_vec = vector() for (i in 1:(length(fifa\_est)-2)){ scaled\_feature = scale(fifa\_est[,i]) fifa\_est\_scaled[,i] = scaled\_feature center\_vec = append(center\_vec, attr(scaled\_feature, 'scaled:center')) scale\_vec = append(scale\_vec, attr(scaled\_feature, 'scaled:scale')) fifa\_est\_scaled[,length(fifa\_est\_scaled)+1] = fifa\_est\$position fifa\_est\_scaled[,length(fifa\_est\_scaled)+1] = fifa\_est\$value\_eur colnames(fifa\_est\_scaled) = colnames(fifa\_est) Validation fifa\_val\_scaled = data.frame(rep(0, nrow(fifa\_val))) for (i in 1:(length(fifa\_val)-2)){ scaled\_feature = scale(fifa\_val[,i], scale = scale\_vec[i], center = center\_vec[i]) scaled\_feature = as.vector(scaled\_feature) fifa\_val\_scaled[,i] = scaled\_feature} class(scaled\_feature) ## [1] "numeric" fifa\_val\_scaled[,length(fifa\_val\_scaled)+1] = fifa\_val\$position fifa\_val\_scaled[,length(fifa\_val\_scaled)+1] = fifa\_val\$value\_eur colnames(fifa\_val\_scaled) = colnames(fifa\_val) fifa\_knn\_model\_list = list()  $k_{list} = seq(10, 100, 10)$ for(i in 1:length(k\_list)){fifa\_knn\_model\_list[[i]] = knnreg(formula = value\_eur ~ . , data = fifa\_est,  $k = k_list[i])$ fifa\_knn\_predicted\_est\_list = lapply(fifa\_knn\_model\_list, predict, fifa\_est) fifa\_rmse\_est\_list = sapply(fifa\_knn\_predicted\_est\_list, rmse, fifa\_est\$value\_eur) fifa\_knn\_predicted\_val\_list = lapply(fifa\_knn\_model\_list, predict, fifa\_val) fifa\_rmse\_val\_list = sapply(fifa\_knn\_predicted\_val\_list, rmse, fifa\_val\$value\_eur) fifa\_scaled\_knn\_model\_list = list() for(i in 1:length(k\_list)){fifa\_scaled\_knn\_model\_list[[i]] = knnreg(formula = value\_eur ~ . , data = fifa\_est\_scaled,  $k = k_list[i]$ fifa\_scaled\_knn\_predicted\_est\_list = lapply(fifa\_scaled\_knn\_model\_list, predict, fifa\_est\_scaled) fifa\_scaled\_knn\_rmse\_est\_list = sapply(fifa\_scaled\_knn\_predicted\_est\_list, rmse, fifa\_est\_scaled\$value\_eur) fifa\_scaled\_knn\_predicted\_val\_list = lapply(fifa\_scaled\_knn\_model\_list, predict, fifa\_val\_scaled) fifa\_scaled\_knn\_rmse\_val\_list = sapply(fifa\_scaled\_knn\_predicted\_val\_list, rmse, fifa\_val\_scaled\$value\_eur) knn\_report = data.frame(fifa\_rmse\_est\_list, fifa\_rmse\_val\_list, fifa\_scaled\_knn\_rmse\_est\_list, fifa\_scaled\_knn\_rmse\_val\_list) colnames(knn\_report) = c('est', 'val', 'scaled\_est', 'scaled\_val') rownames(knn\_report) = c('k=10', 'k=20', 'k=30', 'k=40', 'k=50', 'k=60', 'k=70', 'k=80', 'k=90', 'k=100')knn\_report val scaled\_est scaled\_val est ## k=10 3546599 3847711 3236364 3415462 3669803 3599962 ## k=20 3957183 3928143 3878460 ## k=30 4136475 4024404 3740975 3904888 ## k=40 4312672 4052335 4024642 ## k=50 4449114 4143601 4134414 4019128 ## k=60 4569169 4239056 4260066 4056857 4350526 ## k=70 4676300 4293656 4094884 ## k=80 4756708 4359516 4441625 4162836 ## k=90 4838993 4420409 4507762 4224761 ## k=100 4913750 4472461 4560098 4267872 knn\_report = melt(data.matrix(knn\_report)) colnames(knn\_report) = c('k', 'data\_type', 'RMSE')  $ggplot(data = knn\_report, aes(x = k, y = RMSE, group = data\_type, color = data\_type))+geom\_point(size = 2)+geom\_l$ ine(size = 0.75)4500000 data\_type - est W 40000000 scaled est scaled\_val 3500000 k=10 k=20 k=30 k=40 k=50 k=60 k=70 k=80 k=90 k=100 Regression tree models  $cp_list = c(0.1, 0.01, 0.001)$  $minsplit_list = seq(10, 100, 30)$ fifa\_regtr\_list = list( tr1 = rpart(value\_eur ~ ., data = fifa\_est, cp = 0.1, minsplit = 10), tr2 = rpart(value\_eur ~ ., data = fifa\_est, cp = 0.1, minsplit = 40), tr3 = rpart(value\_eur ~ ., data = fifa\_est, cp = 0.1, minsplit = 70), tr4 = rpart(value\_eur ~ ., data = fifa\_est, cp = 0.1, minsplit = 100), tr5 = rpart(value\_eur ~ ., data = fifa\_est, cp = 0.01, minsplit = 10), tr6 = rpart(value\_eur ~ ., data = fifa\_est, cp = 0.01, minsplit = 40), tr7 = rpart(value\_eur ~ ., data = fifa\_est, cp = 0.01, minsplit = 70), tr8 = rpart(value\_eur ~ ., data = fifa\_est, cp = 0.01, minsplit = 100), tr9 = rpart(value\_eur ~ ., data = fifa\_est, cp = 0.001, minsplit = 10), tr10 = rpart(value\_eur ~ ., data = fifa\_est, cp = 0.001, minsplit = 40), tr11 = rpart(value\_eur ~ ., data = fifa\_est, cp = 0.001, minsplit = 70), tr12 = rpart(value\_eur ~ ., data = fifa\_est, cp = 0.001, minsplit = 100)) fifa\_regtr\_predicted\_est\_list = lapply(fifa\_regtr\_list, predict, fifa\_est) fifa\_regtr\_rmse\_est\_list = sapply(fifa\_regtr\_predicted\_est\_list, rmse, fifa\_est\$value\_eur) fifa\_regtr\_predicted\_val\_list = lapply(fifa\_regtr\_list, predict, fifa\_val) fifa\_regtr\_rmse\_val\_list = sapply(fifa\_regtr\_predicted\_val\_list, rmse, fifa\_val\$value\_eur) regtr\_report = data.frame(fifa\_regtr\_rmse\_est\_list,fifa\_regtr\_rmse\_val\_list) colnames(regtr\_report) = c('Est', 'Val') rownames(regtr\_report) = c('tr1', 'tr2', 'tr3', 'tr4', 'tr5', 'tr6', 'tr7', 'tr8', 'tr9', 'tr10', 'tr11', 'tr12') regtr\_report Est Val ## tr1 4212303 4438735 ## tr2 4212303 4438735 ## tr3 4212303 4438735 ## tr4 4212303 4438735 ## tr5 2495107 2694443 ## tr6 2731512 2813335 ## tr7 2961682 3478820 ## tr8 2961682 3478820 ## tr9 1425771 1754542 ## tr10 2054955 2140377 ## tr11 2488345 3022428 ## tr12 2565220 3127878 regtr\_report = melt(data.matrix(regtr\_report)) colnames(regtr\_report) = c('model', 'data\_type', 'RMSE')  $ggplot(data = regtr\_report, aes(x = model, y = RMSE, group = data\_type, color = data\_type))+$  $geom_point(size = 2)+$  $geom\_line(size = 0.75)$ 4e+06 data\_type 3e+06-- Est **→** Val 2e+06 tr10 tr11 tr12 model rpart.plot(fifa\_regtr\_list\$tr9) 2.9e+6 100% 1.8e+6 96% 9.2e+6 7% 28e+6 3% 7.2e+6 5% 21e+6 2% 1.3e+6 29% 3.5e+6 3% 3.4e+6 1% 5e+6 0% 11e+6 1% 25e+6 0% 18e+6 1% 32e+6 0% 20e+6 0% 31e+6 0% 80e+6 0% 52e+6 0% 516e+3 43% (2.1e+6) 7% (3.2e+8) 4% (6.4e+6) 1% 6.6e+6 2% 12e+8 1% 17e+6 1% 8.3e+6 0% 24e+8 1% Fit final model on train data and calculate RMSE in predicting test set. tr = rpart(value\_eur ~ ., data = fifa\_trn, cp = 0.001, minsplit = 10) rmse(predict(tr, fifa\_tst), fifa\_tst\$value\_eur) ## [1] 2143079 rpart.plot(tr) 6.9e+6 4% 8.8e+6 1% 16e+6 1% 7.9e+8 0% 22e+6 0% \$159a+2 2.1a+0 3.2a+0 (6.4a+0 6.4a+0 156a+0 25e+6 0%