

Assignment4

Prahlad

6/13/2020

Load Packages

```
suppressMessages(library(ggplot2))
suppressMessages(library(caret))
suppressMessages(library(e1071))
library(randomForest)

## randomForest 4.6-14

## Type rfNews() to see new features/changes/bug fixes.

##
## Attaching package: 'randomForest'

## The following object is masked from 'package:ggplot2':
##
##     margin
```

Read Data

```
setwd("C:/Users/SPRAHLA2/Desktop/ml_rcourse")
#Load DataSet
df_training <- read.csv("C:/Users/SPRAHLA2/Desktop/ml_rcourse/pm1-
training.csv",stringsAsFactors = FALSE)
df_testing <- read.csv("C:/Users/SPRAHLA2/Desktop/ml_rcourse/pm1-
testing.csv",stringsAsFactors = FALSE)
```

Understanding Data

dimension of data

```
#colnames of data
dim(df_training)

## [1] 19622 160
```

160 variables

no of unique users

#no of unique users

```
unique(df_training$user_name)
```

```
## [1] "carlitos" "pedro" "adelmo" "charles" "eurico" "jeremy"
```

6 users in dataset

no of classes in classe

#classe dependent variable

```
unique(df_training$classe)
```

```
## [1] "A" "B" "C" "D" "E"
```

6 classe classes and number of users != number of classe

range of data

#data recording start and end time

```
min(df_training$cvtd_timestamp)
```

```
## [1] "02/12/2011 13:32"
```

```
max(df_training$cvtd_timestamp)
```

```
## [1] "30/11/2011 17:12"
```

Data is from november 30th to december 2nd

Are There Any Missing Value In Data?

```
which(sapply(df_training,function(x) sum(is.na(x))/nrow(df_training)*100)>95)
```

```
##          max_roll_belt          max_pitch_belt          min_roll_belt
##              18              19              21
##      min_pitch_belt      amplitude_roll_belt      amplitude_pitch_belt
##              22              24              25
##      var_total_accel_belt          avg_roll_belt          stddev_roll_belt
##              27              28              29
##          var_roll_belt          avg_pitch_belt          stddev_pitch_belt
##              30              31              32
##          var_pitch_belt          avg_yaw_belt          stddev_yaw_belt
##              33              34              35
##          var_yaw_belt          var_accel_arm          avg_roll_arm
##              36              50              51
##      stddev_roll_arm          var_roll_arm          avg_pitch_arm
##              52              53              54
##      stddev_pitch_arm          var_pitch_arm          avg_yaw_arm
##              55              56              57
```

```
##          stddev_yaw_arm          var_yaw_arm          max_roll_arm
##              58              59              75
##          max_picth_arm          max_yaw_arm          min_roll_arm
##              76              77              78
##          min_pitch_arm          min_yaw_arm          amplitude_roll_arm
##              79              80              81
##          amplitude_pitch_arm          amplitude_yaw_arm          max_roll_dumbbell
##              82              83              93
##          max_picth_dumbbell          min_roll_dumbbell          min_pitch_dumbbell
##              94              96              97
##          amplitude_roll_dumbbell          amplitude_pitch_dumbbell          var_accel_dumbbell
##              99              100              103
##          avg_roll_dumbbell          stddev_roll_dumbbell          var_roll_dumbbell
##              104              105              106
##          avg_pitch_dumbbell          stddev_pitch_dumbbell          var_pitch_dumbbell
##              107              108              109
##          avg_yaw_dumbbell          stddev_yaw_dumbbell          var_yaw_dumbbell
##              110              111              112
##          max_roll_forearm          max_picth_forearm          min_roll_forearm
##              131              132              134
##          min_pitch_forearm          amplitude_roll_forearm          amplitude_pitch_forearm
##              135              137              138
##          var_accel_forearm          avg_roll_forearm          stddev_roll_forearm
##              141              142              143
##          var_roll_forearm          avg_pitch_forearm          stddev_pitch_forearm
##              144              145              146
##          var_pitch_forearm          avg_yaw_forearm          stddev_yaw_forearm
##              147              148              149
##          var_yaw_forearm
##              150
```

**Yes there are missing values.67 variables have missing value % greater than 95%.
We have to be careful of variable selection for model building**

Taking only set of variables which are not having high percentage of missing value and leaving out statistical derived variables like (max,min,avg,skewness,kurtosis,variance,standard deviation) for analysis

```
colnames(df_training)
```

```
## [1] "X"          "user_name"
## [3] "raw_timestamp_part_1" "raw_timestamp_part_2"
## [5] "cvtd_timestamp"      "new_window"
## [7] "num_window"          "roll_belt"
## [9] "pitch_belt"          "yaw_belt"
## [11] "total_accel_belt"     "kurtosis_roll_belt"
## [13] "kurtosis_picth_belt"  "kurtosis_yaw_belt"
```

## [15]	"skewness_roll_belt"	"skewness_roll_belt.1"
## [17]	"skewness_yaw_belt"	"max_roll_belt"
## [19]	"max_picth_belt"	"max_yaw_belt"
## [21]	"min_roll_belt"	"min_pitch_belt"
## [23]	"min_yaw_belt"	"amplitude_roll_belt"
## [25]	"amplitude_pitch_belt"	"amplitude_yaw_belt"
## [27]	"var_total_accel_belt"	"avg_roll_belt"
## [29]	"stddev_roll_belt"	"var_roll_belt"
## [31]	"avg_pitch_belt"	"stddev_pitch_belt"
## [33]	"var_pitch_belt"	"avg_yaw_belt"
## [35]	"stddev_yaw_belt"	"var_yaw_belt"
## [37]	"gyros_belt_x"	"gyros_belt_y"
## [39]	"gyros_belt_z"	"accel_belt_x"
## [41]	"accel_belt_y"	"accel_belt_z"
## [43]	"magnet_belt_x"	"magnet_belt_y"
## [45]	"magnet_belt_z"	"roll_arm"
## [47]	"pitch_arm"	"yaw_arm"
## [49]	"total_accel_arm"	"var_accel_arm"
## [51]	"avg_roll_arm"	"stddev_roll_arm"
## [53]	"var_roll_arm"	"avg_pitch_arm"
## [55]	"stddev_pitch_arm"	"var_pitch_arm"
## [57]	"avg_yaw_arm"	"stddev_yaw_arm"
## [59]	"var_yaw_arm"	"gyros_arm_x"
## [61]	"gyros_arm_y"	"gyros_arm_z"
## [63]	"accel_arm_x"	"accel_arm_y"
## [65]	"accel_arm_z"	"magnet_arm_x"
## [67]	"magnet_arm_y"	"magnet_arm_z"
## [69]	"kurtosis_roll_arm"	"kurtosis_picth_arm"
## [71]	"kurtosis_yaw_arm"	"skewness_roll_arm"
## [73]	"skewness_pitch_arm"	"skewness_yaw_arm"
## [75]	"max_roll_arm"	"max_picth_arm"
## [77]	"max_yaw_arm"	"min_roll_arm"
## [79]	"min_pitch_arm"	"min_yaw_arm"
## [81]	"amplitude_roll_arm"	"amplitude_pitch_arm"
## [83]	"amplitude_yaw_arm"	"roll_dumbbell"
## [85]	"pitch_dumbbell"	"yaw_dumbbell"
## [87]	"kurtosis_roll_dumbbell"	"kurtosis_picth_dumbbell"
## [89]	"kurtosis_yaw_dumbbell"	"skewness_roll_dumbbell"
## [91]	"skewness_pitch_dumbbell"	"skewness_yaw_dumbbell"
## [93]	"max_roll_dumbbell"	"max_picth_dumbbell"
## [95]	"max_yaw_dumbbell"	"min_roll_dumbbell"
## [97]	"min_pitch_dumbbell"	"min_yaw_dumbbell"
## [99]	"amplitude_roll_dumbbell"	"amplitude_pitch_dumbbell"
## [101]	"amplitude_yaw_dumbbell"	"total_accel_dumbbell"
## [103]	"var_accel_dumbbell"	"avg_roll_dumbbell"
## [105]	"stddev_roll_dumbbell"	"var_roll_dumbbell"
## [107]	"avg_pitch_dumbbell"	"stddev_pitch_dumbbell"
## [109]	"var_pitch_dumbbell"	"avg_yaw_dumbbell"
## [111]	"stddev_yaw_dumbbell"	"var_yaw_dumbbell"
## [113]	"gyros_dumbbell_x"	"gyros_dumbbell_y"

```
## [115] "gyros_dumbbell_z"      "accel_dumbbell_x"
## [117] "accel_dumbbell_y"      "accel_dumbbell_z"
## [119] "magnet_dumbbell_x"     "magnet_dumbbell_y"
## [121] "magnet_dumbbell_z"     "roll_forearm"
## [123] "pitch_forearm"         "yaw_forearm"
## [125] "kurtosis_roll_forearm" "kurtosis_pitch_forearm"
## [127] "kurtosis_yaw_forearm"  "skewness_roll_forearm"
## [129] "skewness_pitch_forearm" "skewness_yaw_forearm"
## [131] "max_roll_forearm"      "max_pitch_forearm"
## [133] "max_yaw_forearm"       "min_roll_forearm"
## [135] "min_pitch_forearm"     "min_yaw_forearm"
## [137] "amplitude_roll_forearm" "amplitude_pitch_forearm"
## [139] "amplitude_yaw_forearm" "total_accel_forearm"
## [141] "var_accel_forearm"     "avg_roll_forearm"
## [143] "stddev_roll_forearm"   "var_roll_forearm"
## [145] "avg_pitch_forearm"     "stddev_pitch_forearm"
## [147] "var_pitch_forearm"     "avg_yaw_forearm"
## [149] "stddev_yaw_forearm"    "var_yaw_forearm"
## [151] "gyros_forearm_x"       "gyros_forearm_y"
## [153] "gyros_forearm_z"       "accel_forearm_x"
## [155] "accel_forearm_y"       "accel_forearm_z"
## [157] "magnet_forearm_x"      "magnet_forearm_y"
## [159] "magnet_forearm_z"      "classe"
```

data variable selection

```
df_training_selected <-
df_training[,c("user_name", "cvtd_timestamp", "new_window", "num_window", "roll_b
elt", "pitch_belt", "yaw_belt", "total_accel_belt", "roll_arm", "pitch_arm", "yaw_a
rm", "total_accel_arm", "roll_dumbbell", "pitch_dumbbell", "yaw_dumbbell", "total_
accel_dumbbell", "roll_forearm", "pitch_forearm", "yaw_forearm", "total_accel_for
earm", "classe")]
```

check missing value %

```
sapply(df_training_selected, function(x) sum(is.na(x))/nrow(df_training)*100)
```

```
##          user_name      cvtd_timestamp      new_window
##              0              0              0
##      num_window      roll_belt      pitch_belt
##              0              0              0
##      yaw_belt      total_accel_belt      roll_arm
##              0              0              0
##      pitch_arm      yaw_arm      total_accel_arm
##              0              0              0
##      roll_dumbbell      pitch_dumbbell      yaw_dumbbell
##              0              0              0
## total_accel_dumbbell      roll_forearm      pitch_forearm
##              0              0              0
##      yaw_forearm      total_accel_forearm      classe
##              0              0              0
```

Exploratory Analysis(Only 5 plots are recommended so i am picking just one variable in each of arm,dumbbell,forearm,belt from selected variable done in the above step)

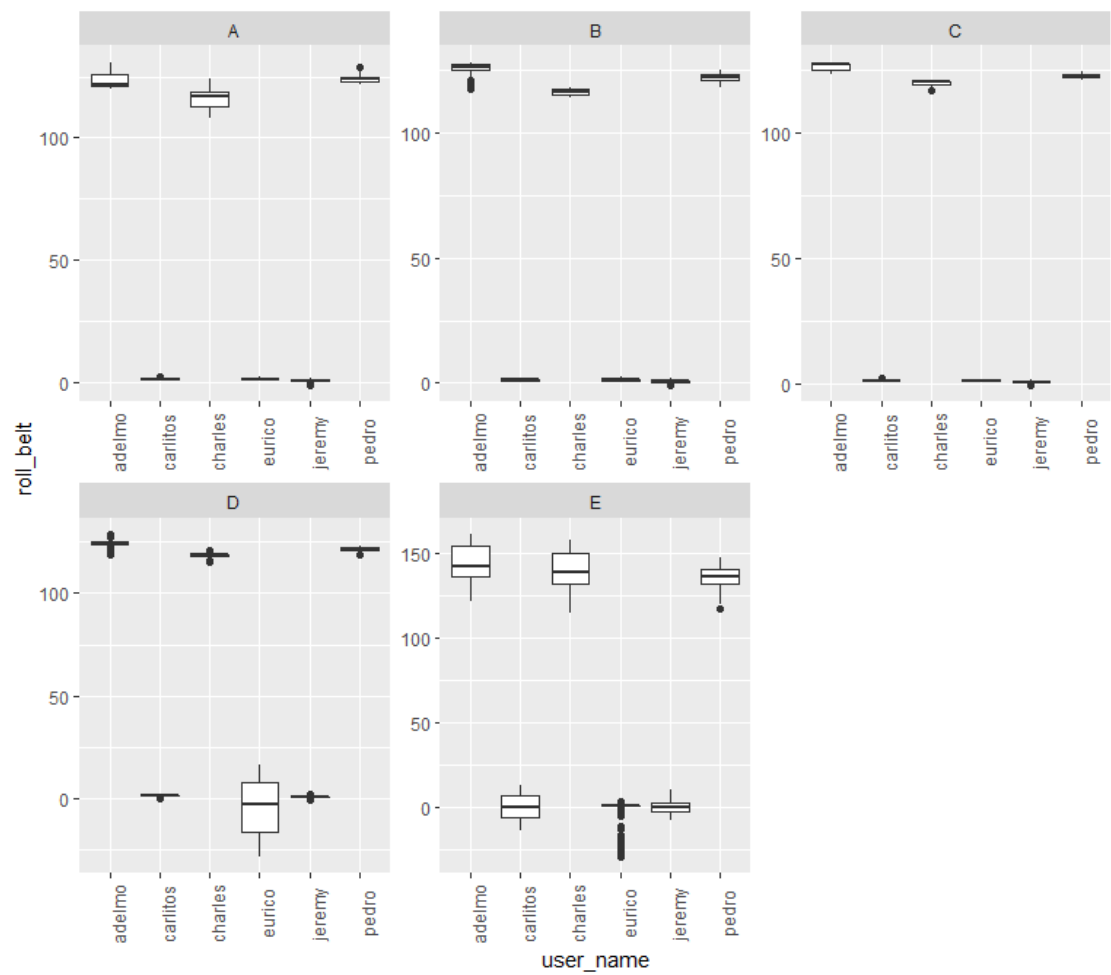
Usage of belt,arm,dumbbell,forearm among users each day(4 days!)

```
df_training_selected$day <-  
as.character(as.Date(df_training_selected$cvtd_timestamp, "%d/%m/%Y"))  
table(df_training_selected$user_name,df_training_selected$day)
```

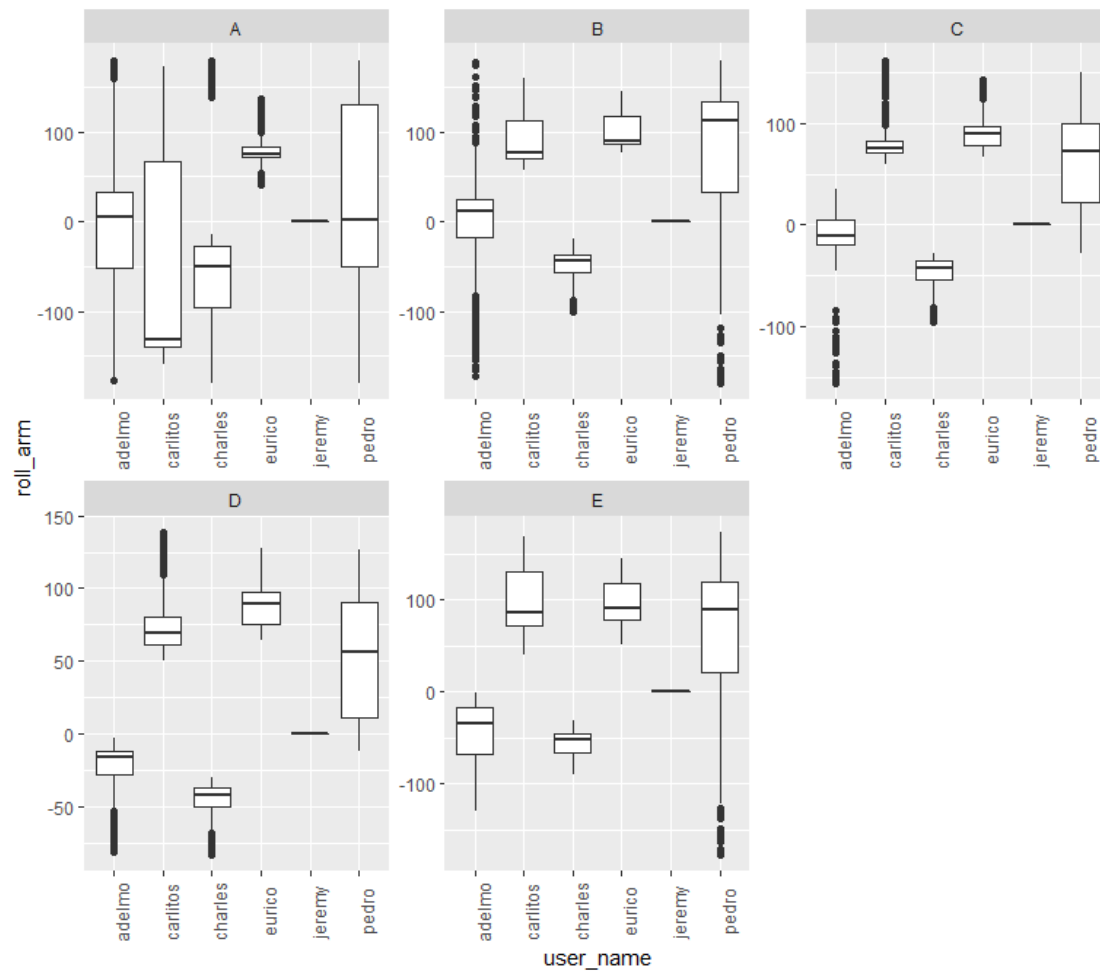
```
##  
##      2011-11-28 2011-11-30 2011-12-02 2011-12-05  
##  adelmo          0          0        3892          0  
##  carlitos         0          0          0        3112  
##  charles          0          0        3536          0  
##  eurico        3070          0          0          0  
##  jeremy          0        3402          0          0  
##  pedro           0          0          0        2610
```

In distribution of 4 days usage each user has tried on different days

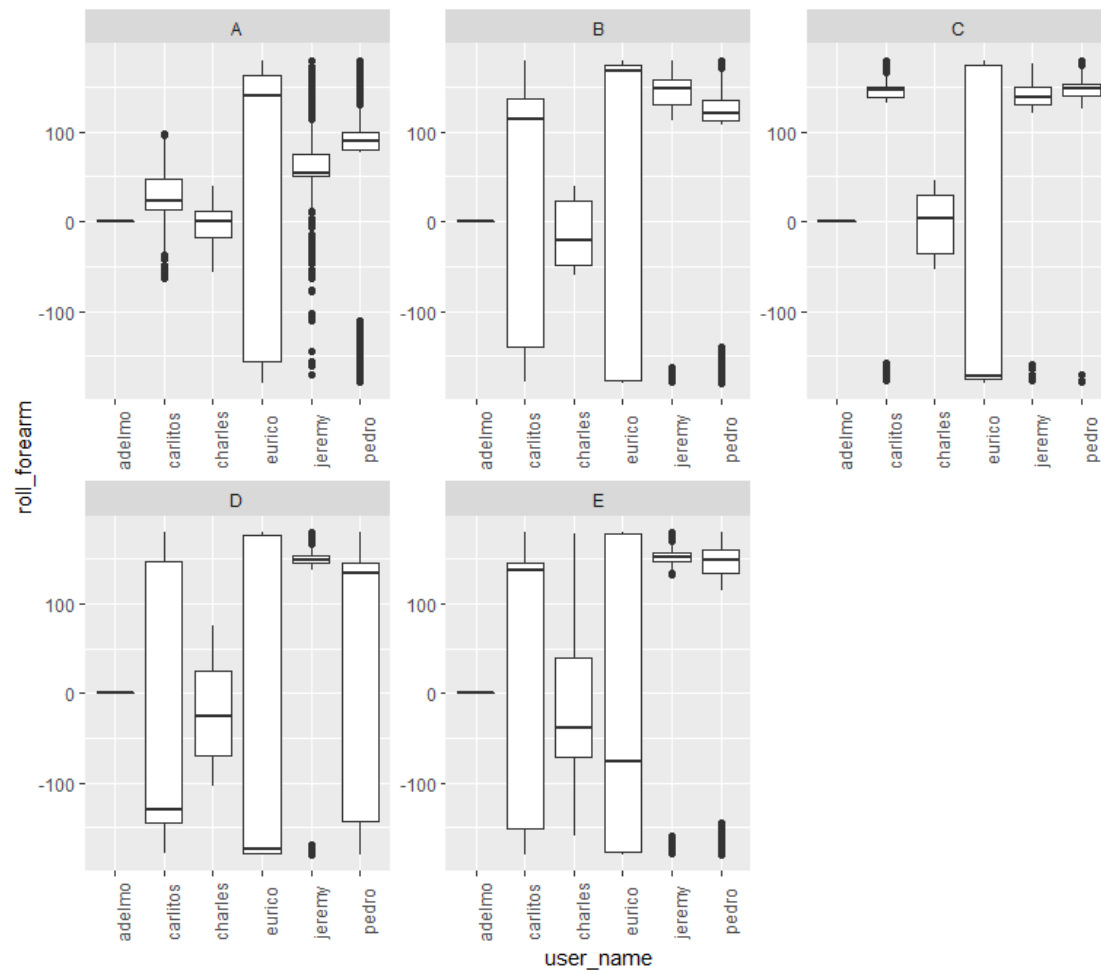
Performance of classe exercise around belt



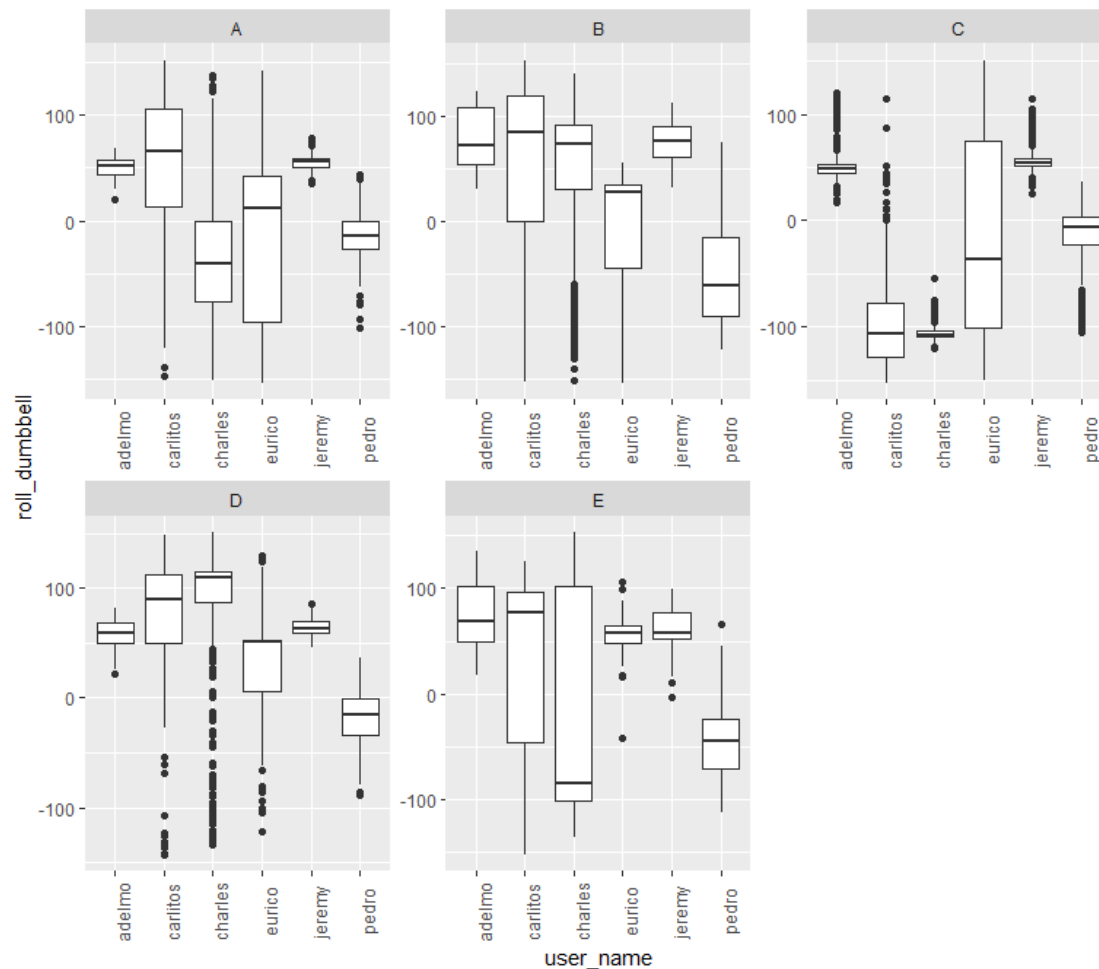
Performance of classe exercise around arm



Performance of classe exercise around fore arm



Performance of classe exercise around dumbbell



There is variability in each user for each classe

Check For Multicollinearity Between Numerical Variables

```
suppressMessages(library(pander))
p <- sapply(df_training_selected, function(x) is.numeric(x))
t <- df_training_selected[,p]
#Remove Na's
q1 <- na.omit(t)
c <- as.data.frame(cor(q1))
panderOptions('table.split.table', Inf)
pander(c)
```

	r	y														
nu	o	pi	a	tot	r	pi	y	tot	rol	pit	ya	total	ro	pit	ya	tota
m_	ll	tc	w	al_	ol	tc	a	al_	l_d	ch_	w_	_acc	ll_	ch_	w_	l_ac
wi	_	h_	_	acc	l_	h_	w	acc	u	du	du	el_d	fo	_fo	for	cel_f
nd	b	b	b	el_	a	ar	_a	el_	m	mb	mb	umb	re	re	ea	ore
ow	el	el	el	bel	r	m	r	ar	bb	bel	bel	bell	ar	ar	rm	arm

			7	5		3	3										
						6	5										
yaw	-	-	0.	-	-	0.	-	1	-	-	0.0	0.1	-	0.	0.0	0.	-
_ar	0.	0.	1	0.	0.2	4	0.		0.0	0.1	86	44	0.06	05	60	18	0.11
m	04	2	4	2	10	0	0		14	25	96	3	311	71	93	45	29
	79	2	3	2	4	7	8		91					6			
	8	5	4	8			2										
		9		6			3										
							5										
total	-	-	0.	-	-	0.	0.	-	1	-	-	-	0.13	0.	-	0.	-
_acc	0.	0.	0	0.	0.2	0	0	0.		0.0	0.0	0.0	7	02	0.1	09	0.07
el_a	06	2	9	2	66	5	3	0		05	72	28		37	86	28	364
rm	45	7	3	4		7	0	1		29	26	83		5	7	8	
	5	8	0	0		7	8	4		3							
		3	2	7		1	9	9									
								1									
roll_	0.	-	-	0.	-	-	0.	-	-	1	0.2	-	0.35	-	-	-	0.18
dum	03	0.	0.	0	0.1	0.	0	0.	0.0		12	0.2	99	0.	0.1	0.	71
bbel	40	1	3	9	31	1	2	1	05		9	77		02	08	00	
l	7	2	5	7	4	4	1	2	29			1		39	3	02	
		6	0	1		9	5	5	3					2		38	
		5	2	3		2	4									1	
pitc	-	0.	0.	-	0.0	0.	-	0.	-	0.2	1	0.5	-	-	0.2	-	0.01
h_du	0.	0	2	0.	58	2	0.	0	0.0	12		17	0.40	0.	90	0.	484
mbb	15	6	3	0	95	2	1	8	72	9		4	52	07	8	09	
ell	27	4	1	9		8	5	6	26					23		26	
		0	6	3		6	2	9						6		4	
		3		1			6	6									
				2													
yaw	-	0.	0.	-	0.0	0.	-	0.	-	-	0.5	1	-	0.	0.2	-	-
_du	0.	0	6	0.	60	2	0.	1	0.0	0.2	17		0.52	01	80	0.	0.18
mbb	10	2	6	3	43	7	1	4	28	77	4		67	27	4	05	94
ell	84	5		4		9	7	4	83	1				6		57	
		6		7		6	8	3								9	
		4		7			4										
total	0.	-	-	0.	-	-	0.	-	0.1	0.3	-	-	1	0.	-	0.	0.12
_acc	08	0.	0.	0	0.1	0.	1	0.	37	59	0.4	0.5		21	0.3	22	57
el_d	87	1	3	4	71	1	9	0		9	05	26		9	44	89	
umb	3	9	1	5	4	5	7	6			2	7			9		
bell		2	7	7		5	2	3									
		1	5	6		4		1									
								1									
roll_	-	-	0.	-	-	0.	0.	0.	0.0	-	-	0.0	0.21	1	-	0.	-
fore	0.	0.	1	0.	0.1	0	0	0	23	0.0	0.0	12	9		0.0	34	0.08

arm	01	1	4	1	12	7	2	5	75	23	72	76			54	67	146
	11	5	5	8	7	5	2	7		92	36				41		
	3	0	2	0		8	8	1									
		2		7		8	7	6									
pitc	-	0.	0.	-	0.1	0.	-	0.	-	-	0.2	0.2	-	-	1	-	-
h_fo	0.	1	2	0.	84	1	0.	0	0.1	0.1	90	80	0.34	0.		0.	0.17
rear	04	7	5	0	5	6	1	6	86	08	8	4	49	05		22	93
m	69	4	3	3		5	5	0	7	3				44		08	
	7	6	6	0		7	4	9						1			
				7			5	3									
				1													
yaw	-	-	0.	-	-	0.	0.	0.	0.0	-	-	-	0.22	0.	-	1	0.19
_for	0.	0.	0	0.	0.2	2	0	1	92	0.0	0.0	0.0	89	34	0.2		98
ear	08	2	4	1	38	3	7	8	88	00	92	55		67	20		
m	41	6	8	9	7	6	9	4		23	64	79			8		
	3	6	5	5		5	7	5		81							
		4	1	7			7										
total	-	0.	-	0.	0.0	-	0.	-	-	0.1	0.0	-	0.12	-	-	0.	1
_acc	0.	0	0.	2	32	0.	1	0.	0.0	87	14	0.1	57	0.	0.1	19	
el_fo	02	6	3	2	22	0	0	1	73	1	84	89		08	79	98	
rear	65	5	2	7		1	8	1	64			4		14	3		
m	1	5	9	8		1	9	2						6			
		6	3			9		9									
						1											

Multi collinearity exists but since we are building tree based models including rf where each tree is different.we can try to remove them in further iterations

Model Building

distribution of classe variable

```
table(df_training_selected$classe)
```

```
##
##      A      B      C      D      E
## 5580 3797 3422 3216 3607
```

Equal distribution of classes. We can go for accuracy as metric

Building a default RF with no hyper parameters and OOB is a CV metric

1) Since the user names are same in both train and test we can use user name

2) dropping new window as test set only has value "no"

```
#convert dependent to factor
df_training_selected$classe <- as.factor(df_training_selected$classe)

# dropping unwanted column

df_training_selected$day <- NULL

#character to factor
df_training_selected$new_window <- as.factor(df_training_selected$new_window)
df_training_selected$user_name <- as.factor(df_training_selected$user_name)

set.seed(3457)

#training
model_rf <- randomForest(classe~.-cvtd_timestamp-
new_window,data=df_training_selected)

model_rf

##
## Call:
## randomForest(formula = classe ~ . - cvtd_timestamp - new_window,
data = df_training_selected)
##
##           Type of random forest: classification
##           Number of trees: 500
## No. of variables tried at each split: 4
##
##           OOB estimate of  error rate: 0.09%
## Confusion matrix:
##      A    B    C    D    E  class.error
## A 5580     0     0     0     0 0.0000000000
## B   1 3795     0     1     0 0.0005267316
## C   0   4 3416     2     0 0.0017533606
## D   0   1   1 3213     1 0.0009328358
## E   0   0   3   3 3601 0.0016634322
```

The accuracy is 0.99 and the oob error is 1% which is equivalent to validation test error

Varimp

```
varImp(model_rf)
```

```
##              Overall
## user_name      457.9678
## num_window    2942.2121
## roll_belt     1936.1633
## pitch_belt    1253.7556
## yaw_belt      1541.9041
## total_accel_belt  538.2202
## roll_arm      524.8097
## pitch_arm     272.3909
## yaw_arm       437.8473
## total_accel_arm  233.3761
## roll_dumbbell  843.3334
## pitch_dumbbell 451.4661
## yaw_dumbbell   626.2025
## total_accel_dumbbell 689.6410
## roll_forearm   930.6746
## pitch_forearm  1225.7259
## yaw_forearm    384.6047
## total_accel_forearm 220.3053
```

**we could see from model_rf that variables
user_name,num_window,roll_belt,pitch_belt,yaw_belt have high importance
compared to other predictors**

Prediction on test data

```
df_testing_selected <-  
df_testing[,c("user_name","cvtd_timestamp","new_window","num_window","roll_belt",  
"pitch_belt","yaw_belt","total_accel_belt","roll_arm","pitch_arm","yaw_arm",  
"total_accel_arm","roll_dumbbell","pitch_dumbbell","yaw_dumbbell","total_accel_dumbbell",  
"roll_forearm","pitch_forearm","yaw_forearm","total_accel_forearm")]
```

#character to factor

```
df_testing_selected$new_window <- as.factor(df_testing_selected$new_window)  
df_testing_selected$user_name <- as.factor(df_testing_selected$user_name)
```

```
predictions <- predict(model_rf,df_testing_selected)
```

```
predictions
```

```
##  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20  
##  B  A  B  A  A  E  D  B  A  A  B  C  B  A  E  E  A  B  B  B  
## Levels: A B C D E
```