

Kobe Shot Selection



**KOBE
BRYANT**

#10 USA BASKETBALL

Content

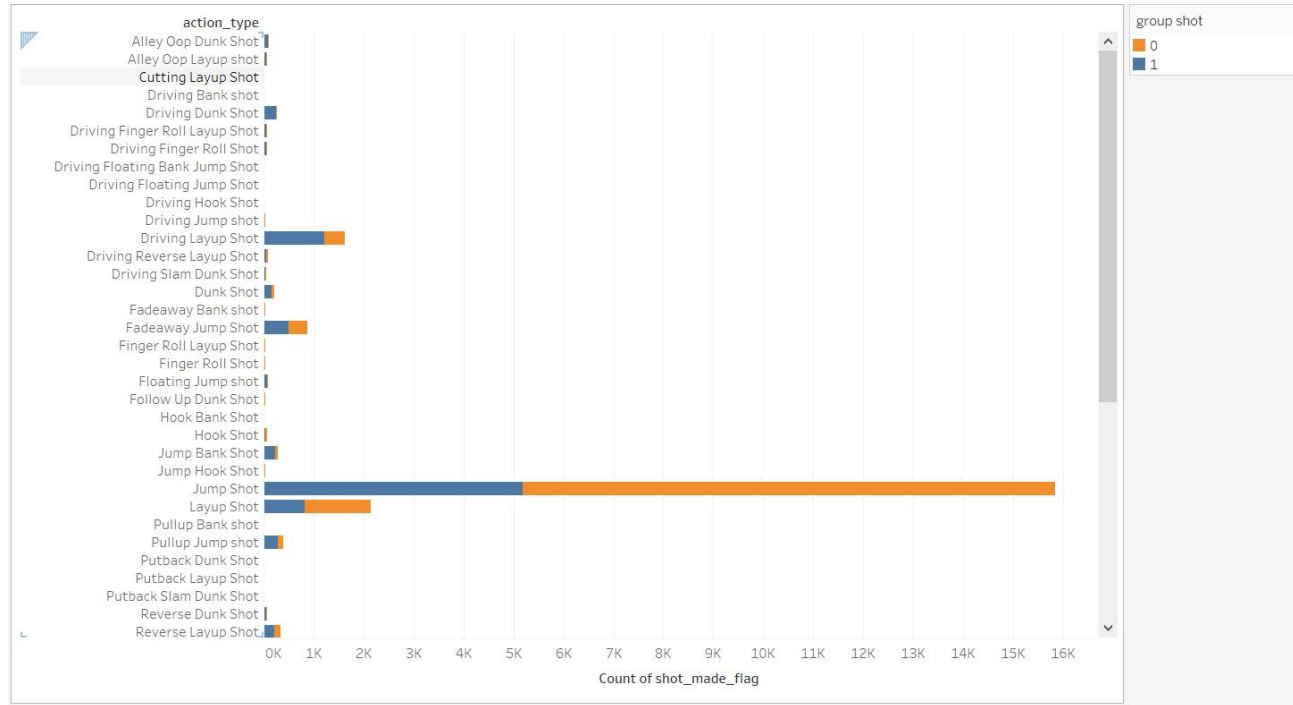
- Data exploration
- Data preparation
- Data modeling
- Model evaluation
- Result

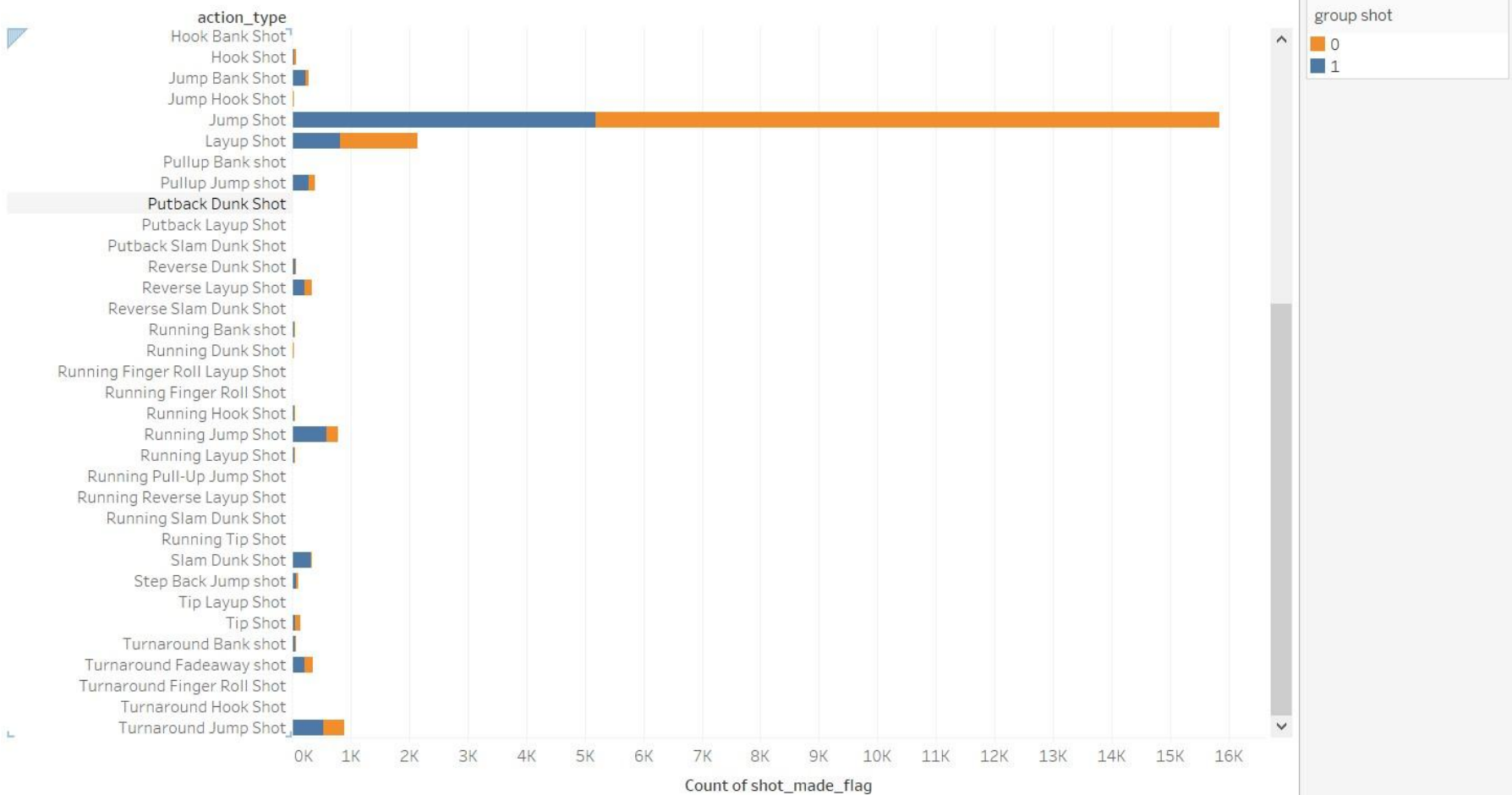
Data Exploration

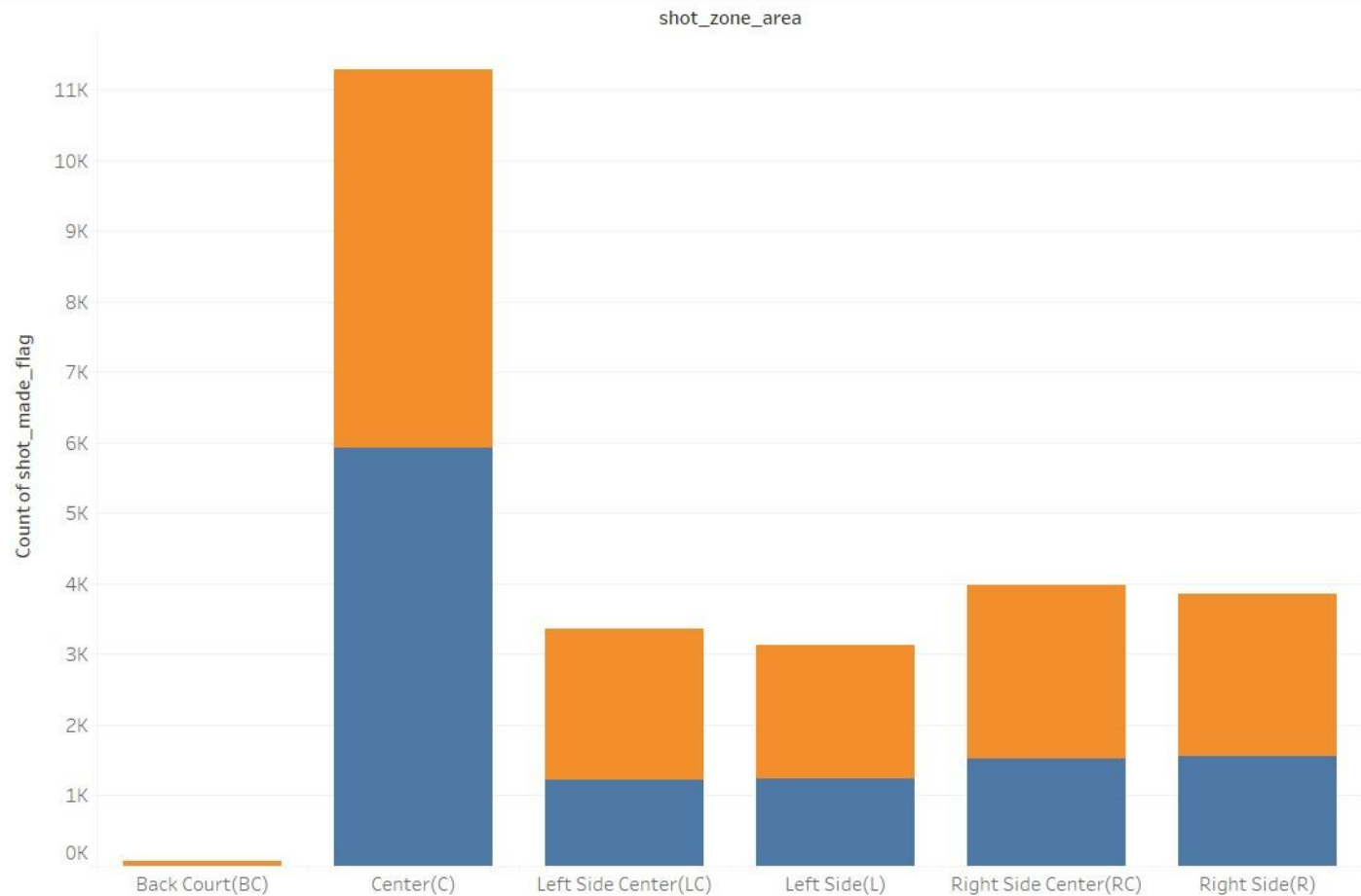
- We are going to explore the data of the selection data set.
- We list out the feature that we will use for the prediction model.

```
## [1] "action_type"      "combined_shot_type" "game_event_id"
## [4] "game_id"          "lat"                "loc_x"
## [7] "loc_y"            "lon"                "minutes_remaining"
## [10] "period"           "playoffs"           "season"
## [13] "seconds_remaining" "shot_distance"      "shot_made_flag"
## [16] "shot_type"        "shot_zone_area"     "shot_zone_basic"
## [19] "shot_zone_range"  "team_id"            "team_name"
## [22] "game_date"        "matchup"            "opponent"
## [25] "shot_id"
```

Let see kobe accuracy with each of this features by graph.



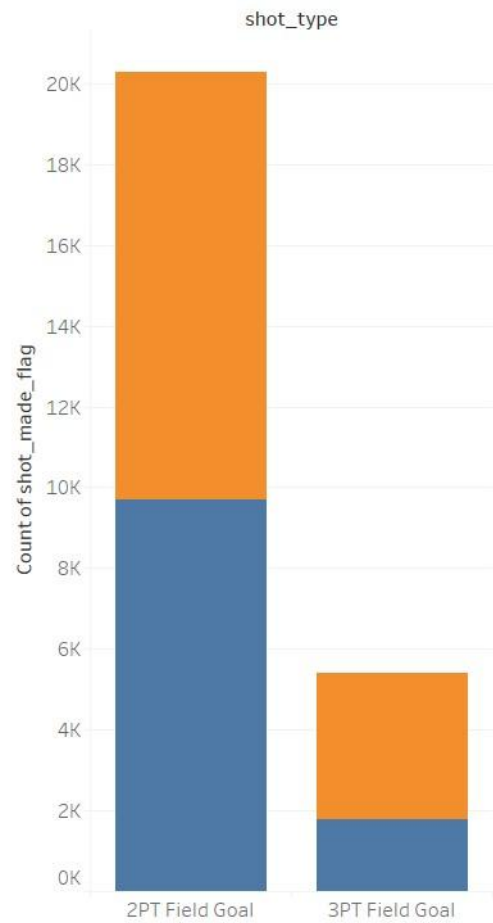




group shot

0

1

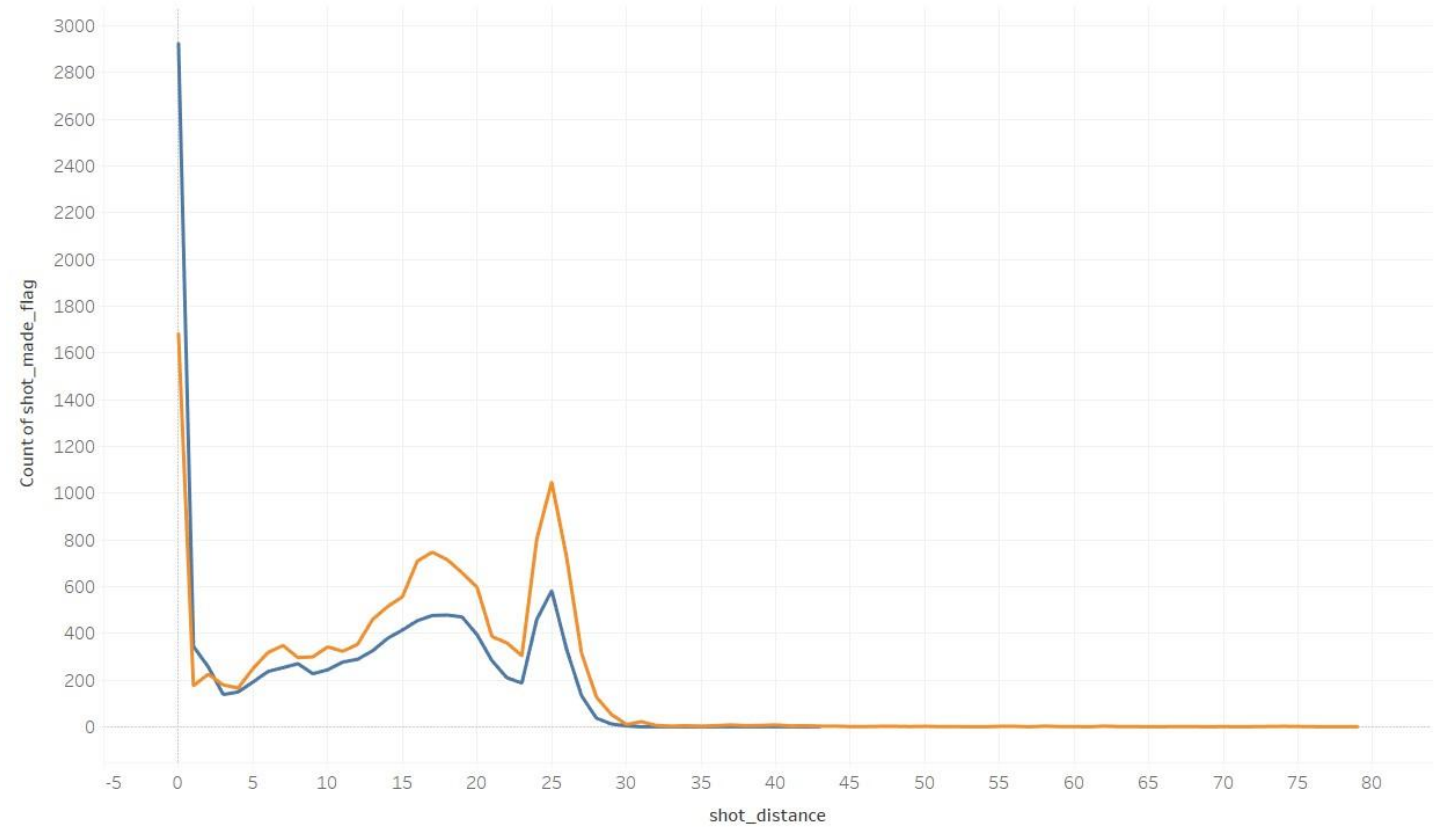


group shot

0

1

Sheet 6

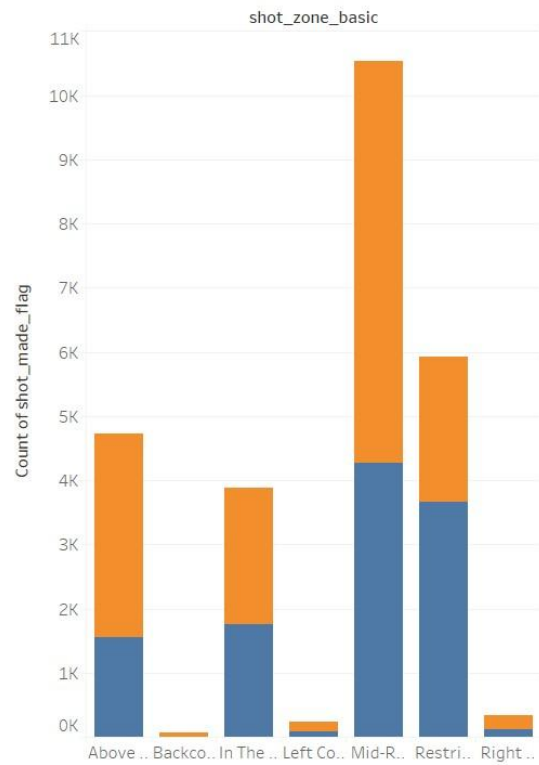


group shot

0

1

Sheet 7

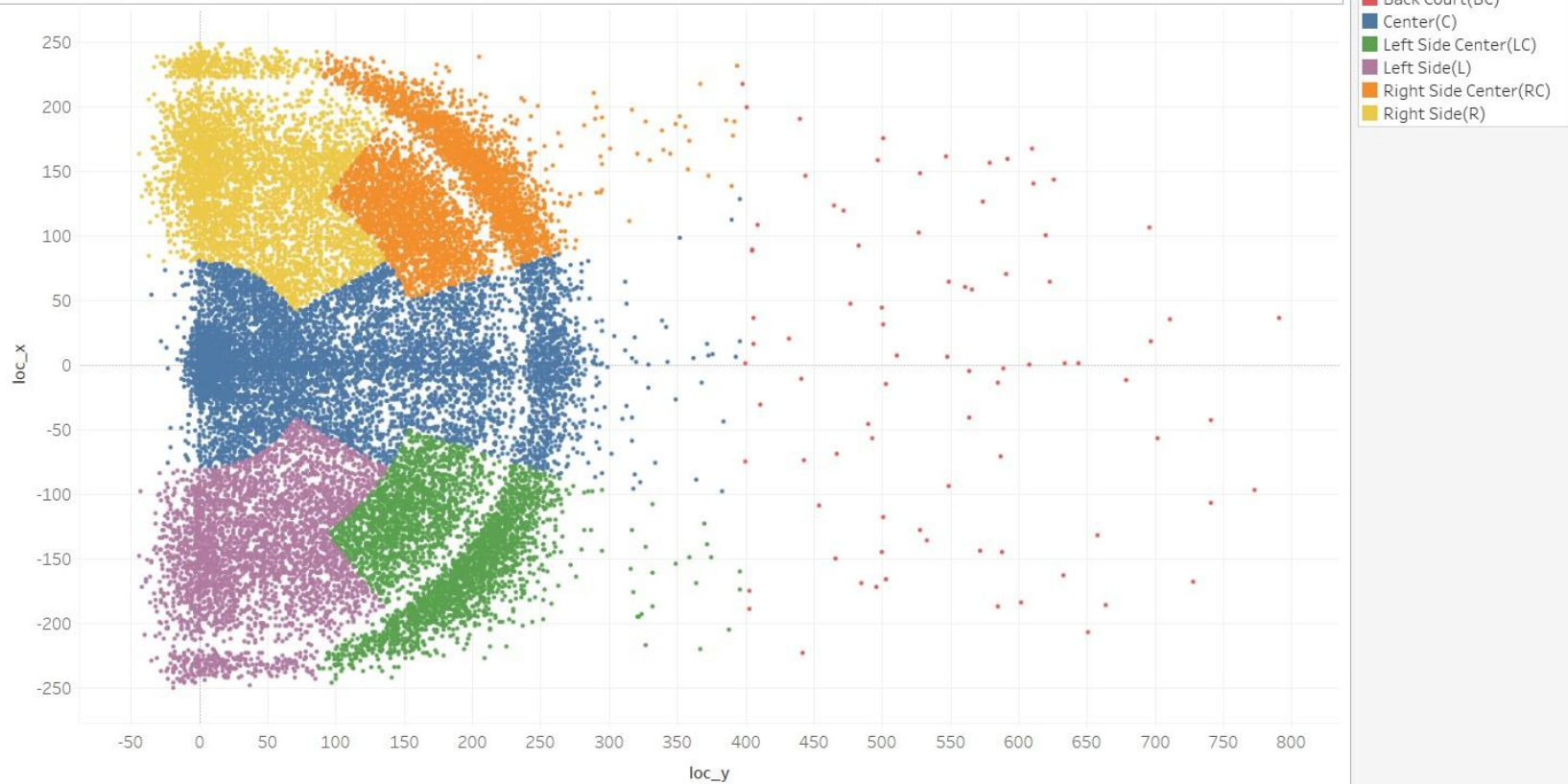


group shot

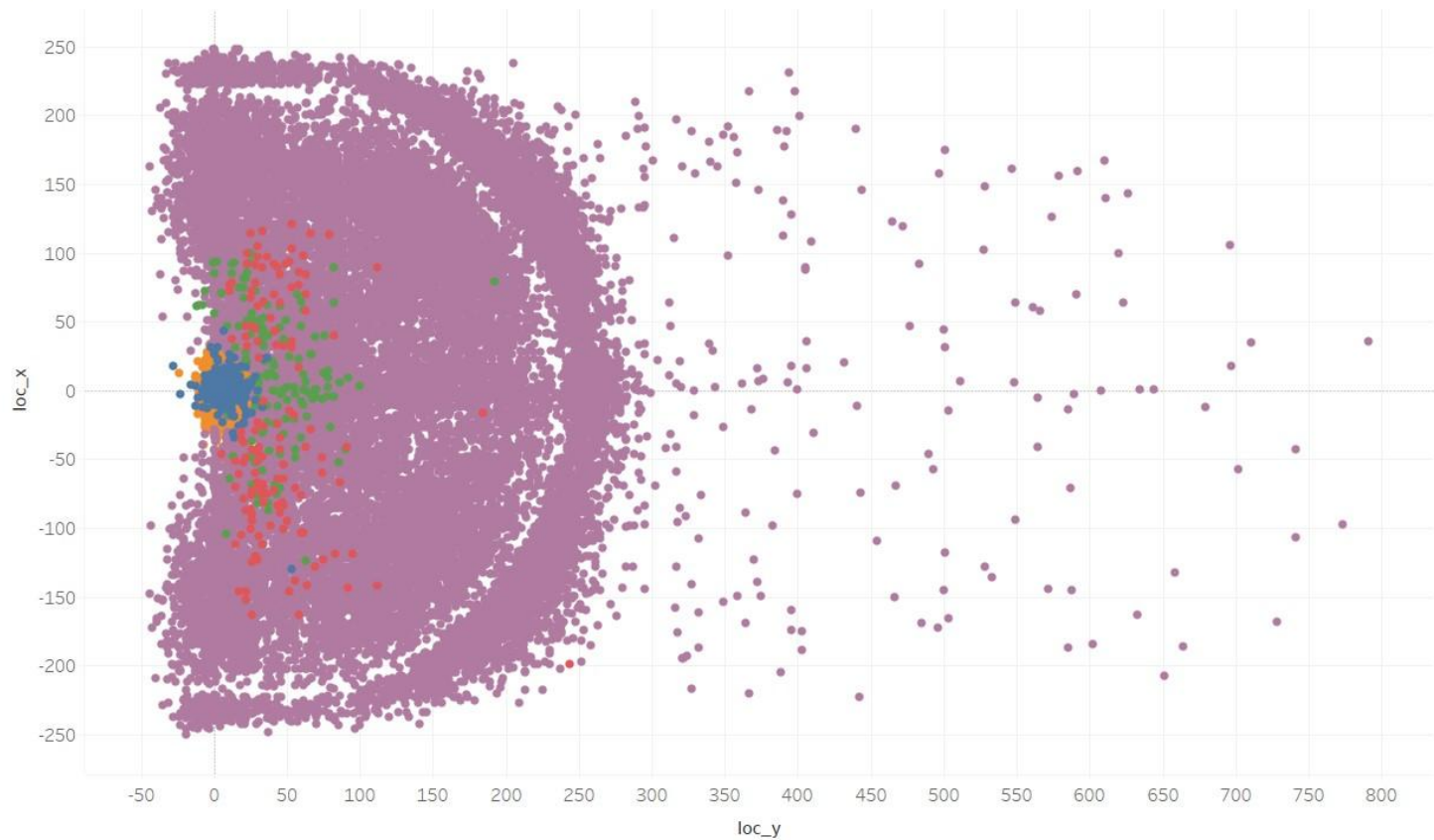
0

1

Sheet 2



Sheet 1



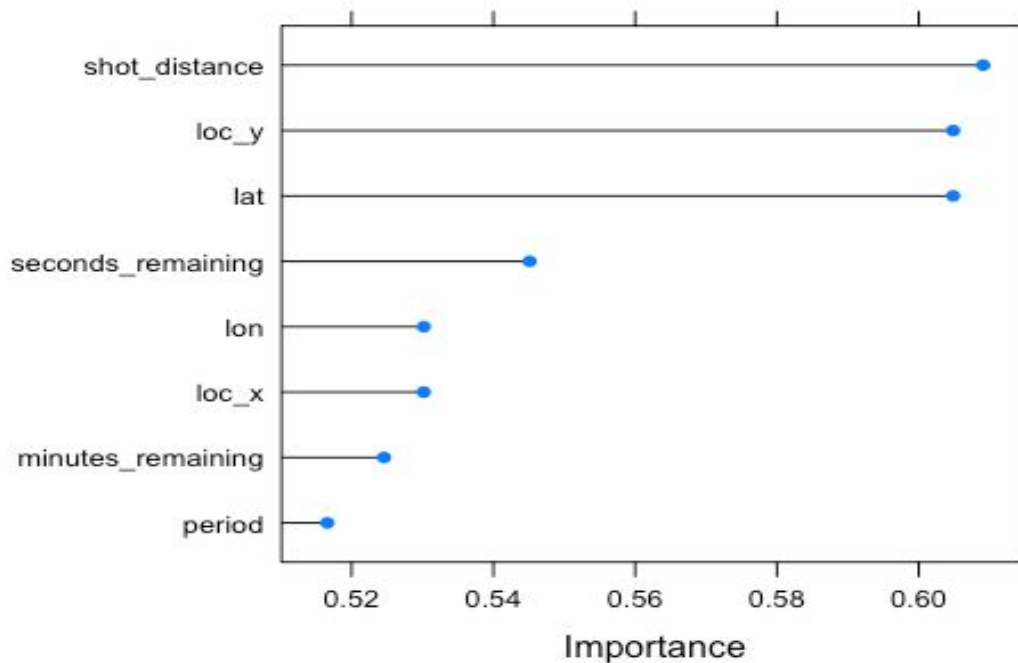
combined_shot_type

- Bank Shot
- Dunk
- Hook Shot
- Jump Shot
- Layup
- Tip Shot

Data preparation

- We are going to prepare the data of the selection data set for modeling.
- There are some method for the preparation of the selected prediction model.
- We also use Learning Vector Quantization model to help us in finding suitable features for the prediction model.

Learning Vector Quantization



Method to preparation.

- The first step we convert the date(yyyy-mm-dd) to weekday, month, and year.

```
11 # Convert date to weekday and other
12 | install.packages("lubridate") #remove comment here if you didn't install this package yet
13 library(lubridate)
14 train$week_day <- wday(train$game_date)
15 train$month <- month(train$game_date)
16 train$year <- year(train$game_date)
17
```

Method to preparation.

- The next step we will prepare the data by normalize the data into range 0:1.

```
6  # Function : normalize data into range 0:1
7  normalize <- function(data) {
8    data <- (data - min(data)) / (max(data) - min(data))
9  }
10
```

```
18 # Normalize data
19 train$loc_x <- normalize(train$loc_x)
20 train$loc_y <- normalize(train$loc_y)
21 train$shot_distance <- normalize(train$shot_distance)
22 train$time_remaining <- (train$minutes_remaining*60) + train$seconds_remaining
23 train$time_remaining <- normalize(train$time_remaining)
24 train$period <- normalize(train$period)
25 train$matchup <- ifelse( grepl("@",train$matchup), train$matchup <- "AWAY", train$matchup <- "HOME" )
26 train$week_day <- normalize(train$week_day)
27 train$month <- normalize(train$month)
28 train$year <- normalize(train$year)
```

Method to preparation.

- The next step we will remove unused features.

```
9  # Remove unused attributes (Useless for data modeling such as _id)
1  train$seconds_remaining <- NULL # use time_remaining instead
2  train$minutes_remaining <- NULL # use time_remaining instead
3  train$game_id <- NULL
4  train$game_event_id <- NULL
5  train$game_date <- NULL          # use weekday,month,year
6  train$team_id <- NULL
7  train$team_name <- NULL          # Kobe NEVER play for other teams
8  train$season <- NULL
9  train$shot_id <- NULL
10 train$lat <- NULL
11 train$lon <- NULL
12
```


Method to preparation.

- The next step we used the library(dummies)
- We used this library to convert categorical data to binary matrix so it is useable in the model.

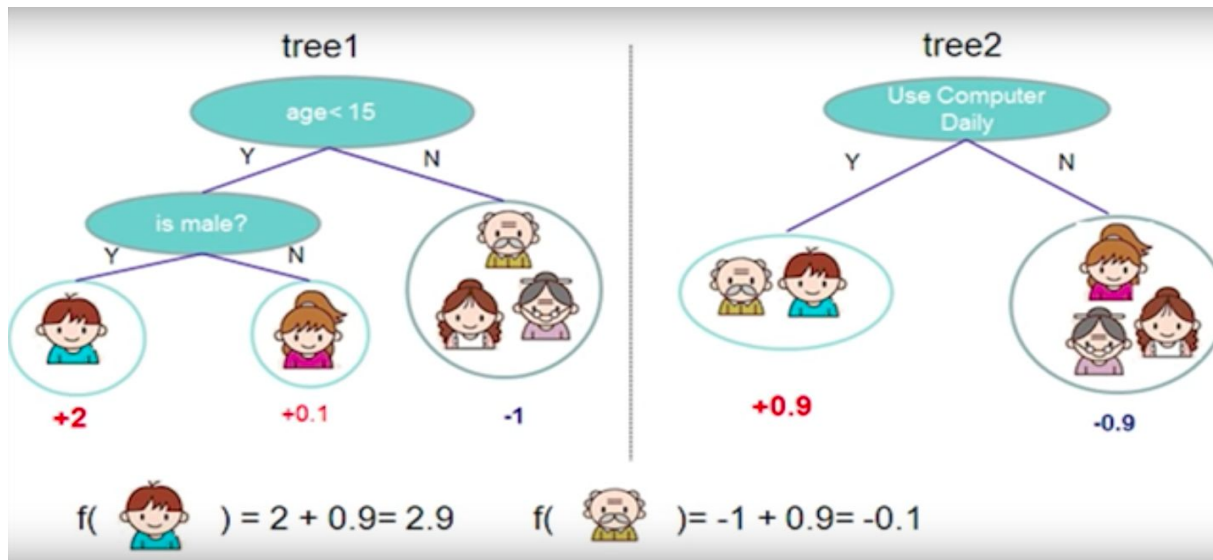
```
43 # Prepare for model
44 install.packages("dummies") #remove comment here if you didn't install this package yet
45 library(dummies)
46 df<- dummy.data.frame(train , names = c('action_type','combined_shot_type' , 'shot_type' ,
47 | 'shot_zone_area' , 'playoffs' , 'shot_zone_basic' , 'shot_zone_range' , 'matchup' , 'opponent') , sep='_')
48 X_train <- df[(!is.na(df$shot_made_flag)),]
49 X_test <- df[(is.na(df$shot_made_flag)),]
50 y <- X_train$shot_made_flag
51
```

Data Modeling

25,697 records as train dataset and 5,000 as test dataset

Xgboost - execute quickly and get good model. Based on **Decision tree**

We tune the objective parameter to **binary:logistic (logistic regression)** because our problem is classification. Whether the shot success or not (1 or 0)



Model Evaluation

We test our model accuracy by sampling data from train dataset 75% remains as model and 25% for testing the model

<u>Confusion matrix</u>	Actual Positive (shot_made_flag = 1)	Actual Negative (shot_made_flag = 0)
Model say YES (shot_made_flag >= 0.5)	TP = 1371	FP = 507
Model say No (shot_made_flag < 0.5)	FN = 1543	TN = 3004

Accuracy ~ 70% Precision ~73% Recall ~ 50%

What is the problem ? - Prediction in range ~ 0.35-0.65














Another evaluation

However, we got **34th** rank in the Kaggle leaderboard with score 0.5998

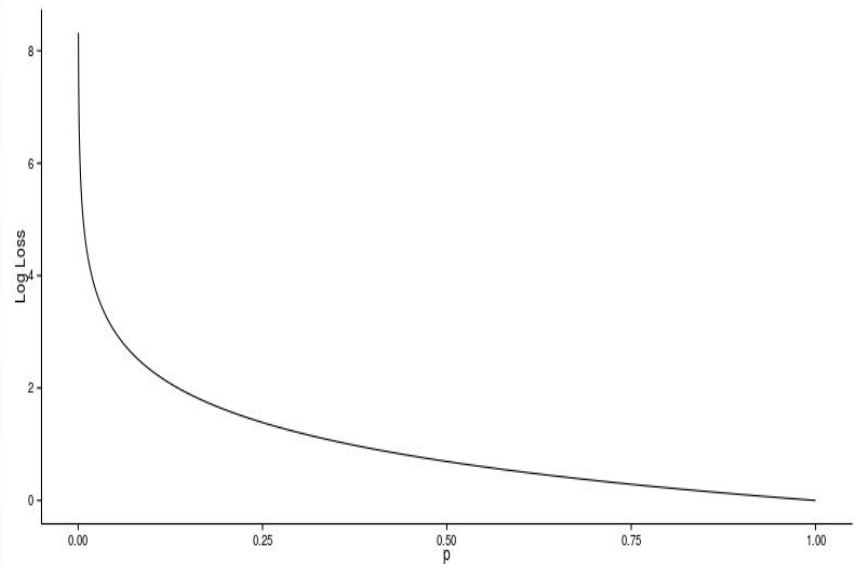
Name	Submitted	Wait time	Execution time	Score
kobe_prediction_result.csv	9 hours ago	0 seconds	0 seconds	0.59986

Complete

[Jump to your position on the leaderboard](#) ▾

1	—	Humberto Brandão		0.56529	50	1y
2	—	y130038		0.56977	15	1y
3	—	The Black Mamba Team	 	0.57222		
4	—	kshain		0.57414		
5	—	Pavel Shashkin		0.57509		
6	—	LeBronLearnsML	  	0.58713		
7	—	Jean Souquet		0.58943		
8	—	derekatwood		0.58949		
9	—	Juan Aguilera		0.59134		
10	—	Gaspard Jallat		0.59741		

Log Loss metric



Result

