# Baseball Case Study



## Problem Statement:

This dataset utilizes data from 2014 Major League Baseball seasons in order to develop an algoirthm that predicts the number of wins for a given team in the 2015 season based on several different indicators of success. This model is used to select best team based on best input feature. .There are 16 different features that will be used as the inputs to the machine learning and the output will be a value that represents the number of wins.

**Input features**:

- R: Runs-times reached home plate legally and safely
- AB: At Bats-plate appearances, not including bases on balls, being hit by pitch, sacrifices, interference, or obstruction
- H: Hits- reaching base because of a batted, fair ball without error by the defence
- 2B: Doubles-hits on which the batter reaches second base safely without the contribution of a fielding error
- 3B: Triples- hits on which the batter reaches third base safely without the contribution of a fielding error
- HR: Homeruns
- BB: Walks-times pitching four balls, allowing the batter to take first base
- SO: Strikeouts
- SB: Stolen Bases- number of bases advanced by the runner while the ball is in the possession of the defence
- RA: Runs Allowed
- ER: Earned Runs- number of runs that did not occur as a result of errors or passed balls
- ERA: Earned Run Average (ERA)- total number of earned runs (see "ER" above), multiplied by 9, divided by innings pitched
- SO: Shutouts- number of complete games pitched with no runs allowed
- SV: Saves- number of games where the pitcher enters a game led by the pitcher's team, finishes the game without surrendering the lead, is not the winning pitcher, and either (a) the lead was three runs or fewer when the pitcher entered the game; (b) the

potential tying run was on base, at bat, or on deck; or (c) the pitcher pitched three or more innings

- ➢ CG: Complete Games-number of games where player was the only pitcher for their team
- ➢ E: Errors- the judgment of the official scorer, of a fielder misplaying a ball in a manner that allows a batter or base runner to advance one or more bases

**Output**: Number of predicted wins (W)

**Importing the Libraries**

import pandas as pd
import numpy as np
import seaborn as sns
from scipy import stats
import warnings
warnings.filterwarnings('ignore')
from matplotlib import pyplot as plt
from scipy.stats import zscore

**Loading Dataset:**

```
In [3]:    1  df=pd.read_csv("E:\\baseball.csv")
           2  df.head(5)
```

Out[3]:

| | W | R | AB | H | 2B | 3B | HR | BB | SO | SB | RA | ER | ERA | CG | SHO | SV | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 95 | 724 | 5575 | 1497 | 300 | 42 | 139 | 383 | 973 | 104 | 641 | 601 | 3.73 | 2 | 8 | 56 | 88 |
| 1 | 83 | 696 | 5467 | 1349 | 277 | 44 | 156 | 439 | 1264 | 70 | 700 | 653 | 4.07 | 2 | 12 | 45 | 86 |
| 2 | 81 | 669 | 5439 | 1395 | 303 | 29 | 141 | 533 | 1157 | 86 | 640 | 584 | 3.67 | 11 | 10 | 38 | 79 |
| 3 | 76 | 622 | 5533 | 1381 | 260 | 27 | 136 | 404 | 1231 | 68 | 701 | 643 | 3.98 | 7 | 9 | 37 | 101 |
| 4 | 74 | 689 | 5605 | 1515 | 289 | 49 | 151 | 455 | 1259 | 83 | 803 | 746 | 4.64 | 7 | 12 | 35 | 86 |

loaded the abalone dataset.

```
n [201]:   1  df.shape
ut[201]: (30, 17)
```

It has 17 coloumns and 30 rows

```
n [202]:   1
           2  df.isnull().sum()
```

ut[202]:
```
W       0
R       0
AB      0
H       0
2B      0
3B      0
HR      0
BB      0
SO      0
SB      0
RA      0
ER      0
ERA     0
CG      0
SHO     0
SV      0
E       0
dtype: int64
```

```
1  From the above observation this dataset has no null values
```
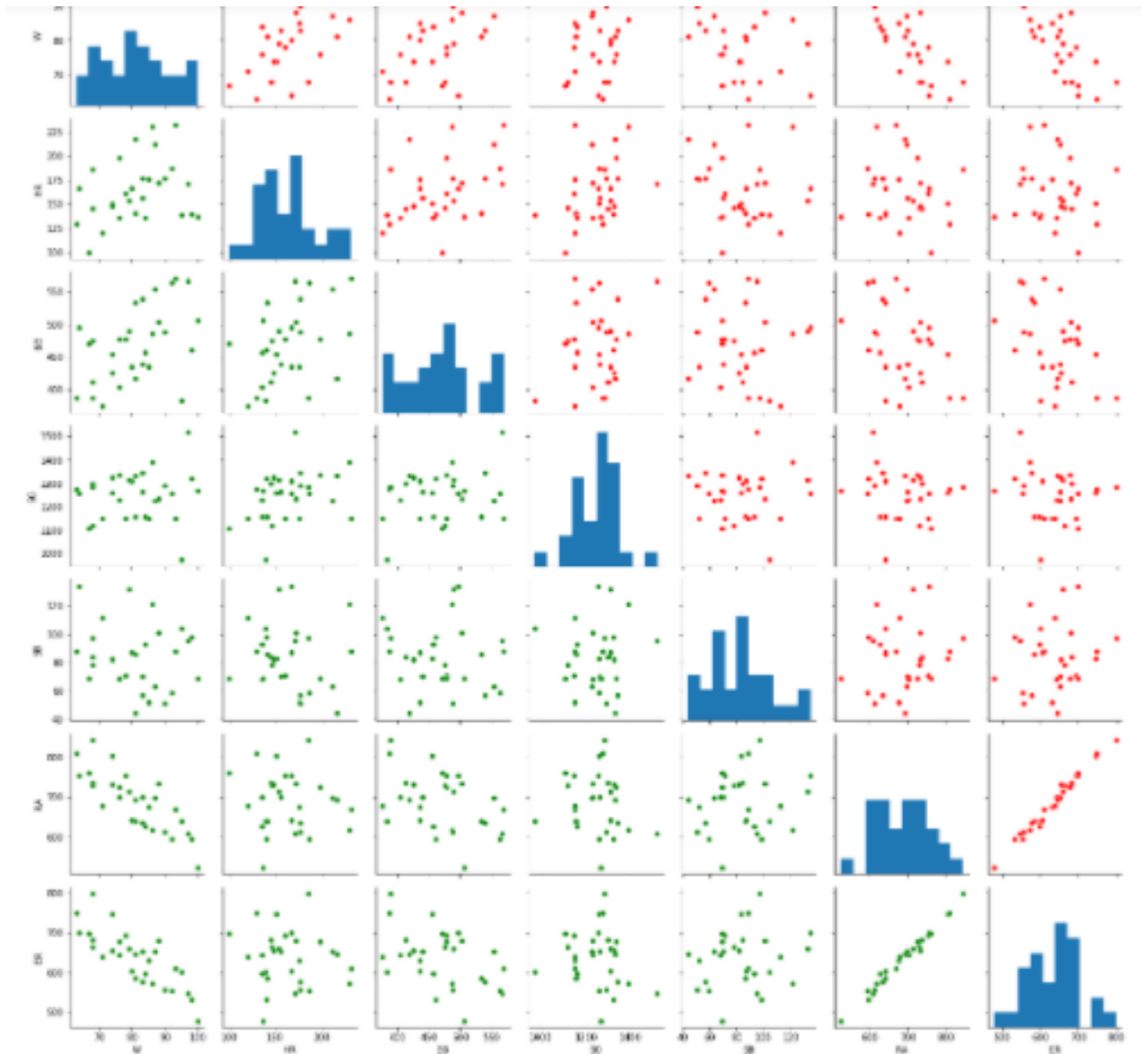
This dataset has no null values and it has 30 rows and 17 columns

## CHECKING FOR EMPTY SPACE IN DATASET

```
3]:   1  col=['W', 'R', 'AB', 'H', '2B', '3B', 'HR', 'BB', 'SO', 'SB', 'RA', 'ER','ERA', 'CG', 'SHO', 'SV', 'E']
      2  for i in col:
      3      print(df.loc[df[i]==""])
```
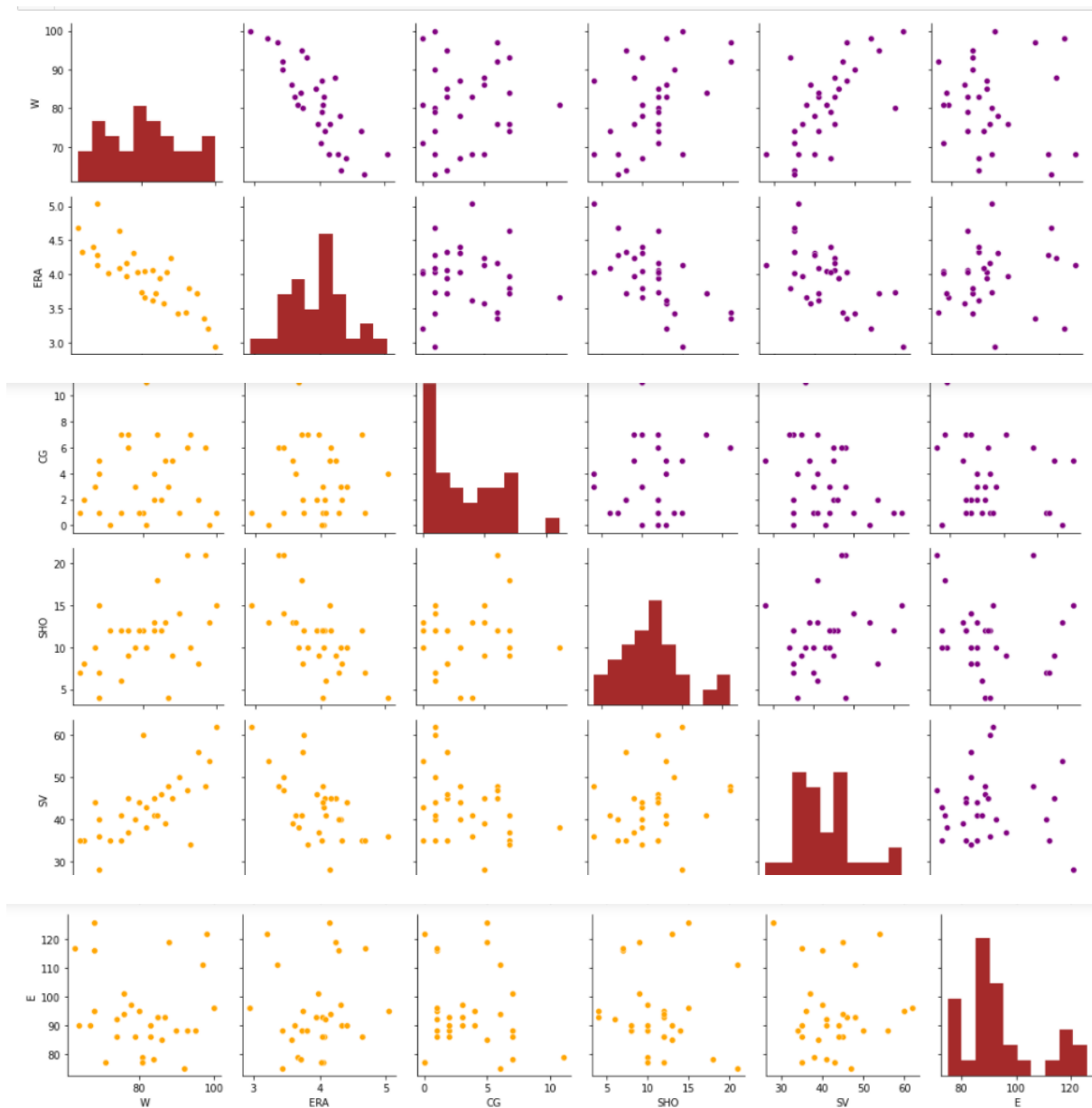
```
Empty DataFrame
Columns: [W, R, AB, H, 2B, 3B, HR, BB, SO, SB, RA, ER, ERA, CG, SHO, SV, E]
Index: []
Empty DataFrame
Columns: [W, R, AB, H, 2B, 3B, HR, BB, SO, SB, RA, ER, ERA, CG, SHO, SV, E]
Index: []
Empty DataFrame
Columns: [W, R, AB, H, 2B, 3B, HR, BB, SO, SB, RA, ER, ERA, CG, SHO, SV, E]
Index: []
Empty DataFrame
Columns: [W, R, AB, H, 2B, 3B, HR, BB, SO, SB, RA, ER, ERA, CG, SHO, SV, E]
Index: []
Empty DataFrame
Columns: [W, R, AB, H, 2B, 3B, HR, BB, SO, SB, RA, ER, ERA, CG, SHO, SV, E]
Index: []
Empty DataFrame
Columns: [W, R, AB, H, 2B, 3B, HR, BB, SO, SB, RA, ER, ERA, CG, SHO, SV, E]
Index: []
Empty DataFrame
Columns: [W, R, AB, H, 2B, 3B, HR, BB, SO, SB, RA, ER, ERA, CG, SHO, SV, E]
Index: []
Empty DataFrame
Columns: [W, R, AB, H, 2B, 3B, HR, BB, SO, SB, RA, ER, ERA, CG, SHO, SV, E]
Index: []
Empty DataFrame
Columns: [W, R, AB, H, 2B, 3B, HR, BB, SO, SB, RA, ER, ERA, CG, SHO, SV, E]
Index: []
Empty DataFrame
Columns: [W, R, AB, H, 2B, 3B, HR, BB, SO, SB, RA, ER, ERA, CG, SHO, SV, E]
Index: []
Empty DataFrame
Columns: [W, R, AB, H, 2B, 3B, HR, BB, SO, SB, RA, ER, ERA, CG, SHO, SV, E]
Index: []
Empty DataFrame
Columns: [W, R, AB, H, 2B, 3B, HR, BB, SO, SB, RA, ER, ERA, CG, SHO, SV, E]
Index: []
```

This dataset has no empty space as value

The data points of the features R and AB showing an uphill pattern as you move from left to right so it has positive relationship. The data points of the features R and w showing an uphill pattern as you move from left to right so it has positive relationship.

The data points of ER and RA showing an uphill pattern as you move from left to right so it has positive relationship.
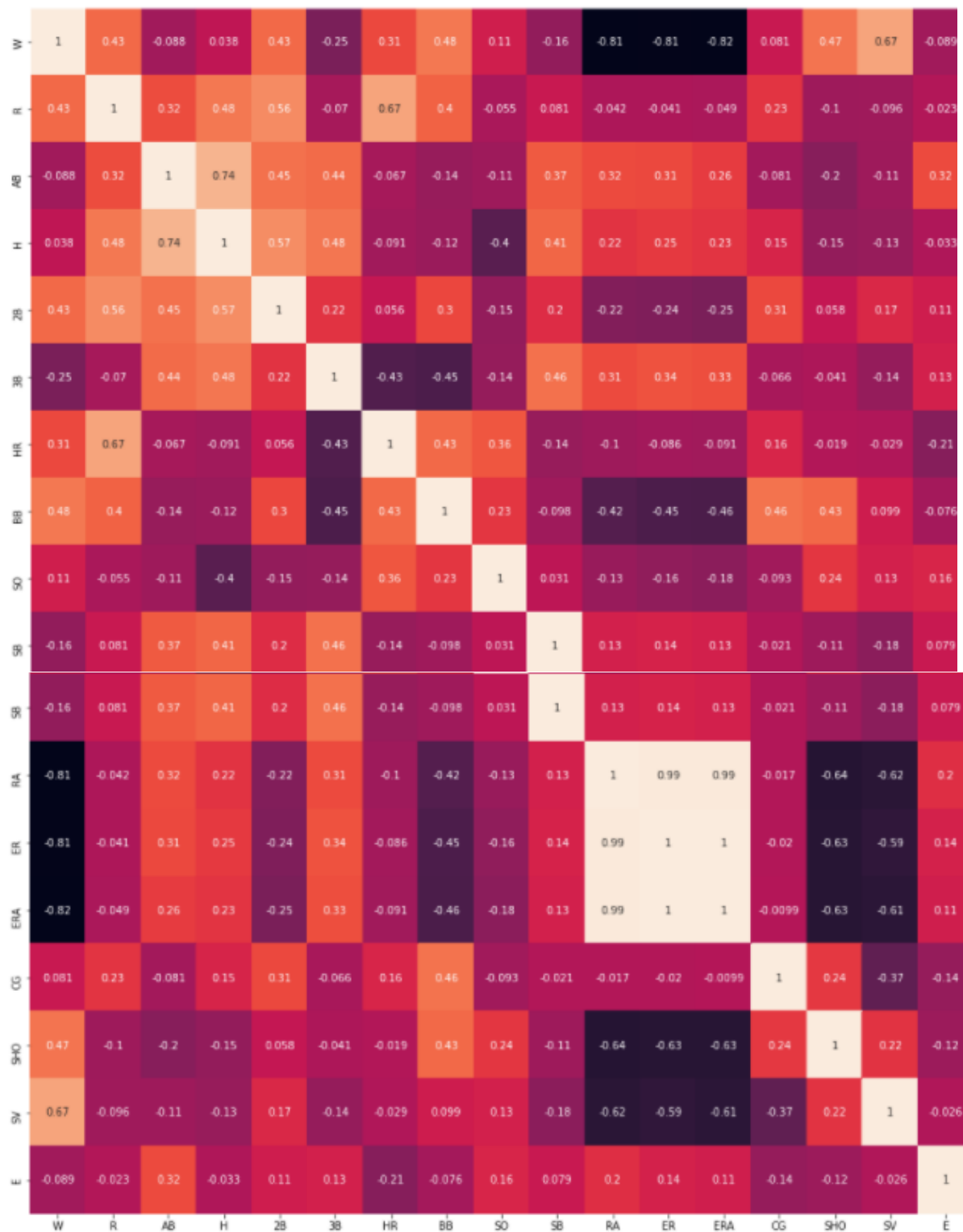
The data points of ER and w showing an downhill pattern as you move from left to right so it has negative relationship.

The data points of RA and w showing an downhill pattern as you move from left to right so it has negative relationship.

The data points of the features ERA and w showing an downhill pattern as you move from left to right so it has negative relationship.

The data points of the features SV and W showing an uphill pattern as you move from left to right so it has positive relationship.

## CORRELATION MATRIX



> ➢ RA,ER,ERA features are highly negatively correlated with wins.

> ➢ RA,ER,ERA features are highly negatively correlated with BB.

➢ RA,ER,ERA features are highly negatively correlated with sv and sho.

➢ Sho,bb,2b,SV,HR are positively correlated with win

➢ ER andERA are posively correlated with RA

➢ wins feature is highly and positively correlated with saves

➢ RA,ER,ERA features are highly negatively correlated with saves.
➢ h,cg having very less correlation with winsRA,ER,ERA features are highly negatively correlated with wins.

➢ RA,ER,ERA features are highly negatively correlated with BB.

➢ RA,ER,ERA features are highly negatively correlated with sv and sho.

➢ Sho,bb,2b,SV,HR are positively correlated with win

➢ ER andERA are posively correlated with RA

➢ wins feature is highly and positively correlated with saves

➢ RA,ER,ERA features are highly negatively correlated with saves.
➢ h,cg having very less correlation with wins

**Finding whether all features which are positively correlated with Wins feature have linear relationship:**

```
1  plt.figure(figsize=(7,5))
2  plt.title('Wins vs Saves')
3  sns.set_style="white_grid"
4  sns.regplot(x="W",y="SV",data=df,scatter_kws = {'color': 'purple'}, line_kws = {'color': 'r'})
5  plt.show()
6
7
```



Wins vs Saves

In this reg plot it is apparent that SV is positively correlated with W data.most of the

datapoints are almost close to best fit line it means that the feature SV and W having linear relationship

```
[282]:  1  plt.figure(figsize=(7,5))
        2  plt.title('Wins vs Walks')
        3  sns.set_style="white_grid"
        4  sns.regplot(x="W",y="BB",data=df,scatter_kws = {'color': 'purple'}, line_kws = {'color': 'r'})
        5
        6
        7
```
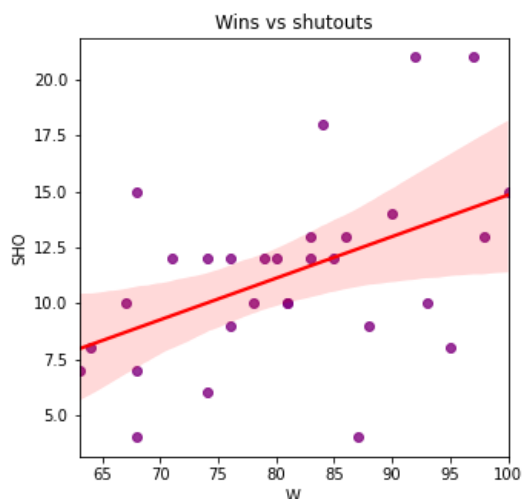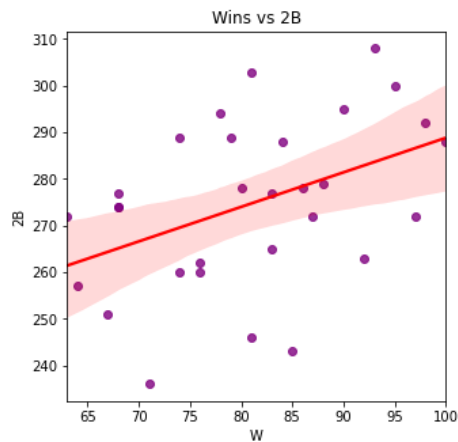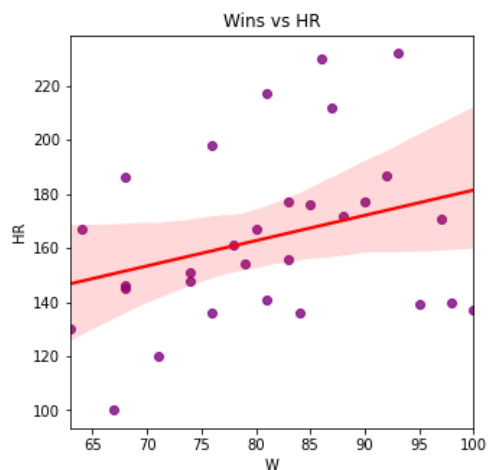
[282]: <AxesSubplot:title={'center':'Wins vs Walks'}, xlabel='W', ylabel='BB'>



In this reg plot it is apparent that BB is positively correlated with W data. Only few datapoints are close to best fit line it means that the feature BB and W are not having linear relationship. If the value of bb is between 400 to 450 there is more chance to win

```
:   1  plt.figure(figsize=(5,5))
    2  plt.title('Wins vs shutouts')
    3  sns.set_style="white_grid"
    4  sns.regplot(x="W",y="SHO",data=df,scatter_kws = {'color': 'purple'}, line_kws = {'color': 'r'})
```

: <AxesSubplot:title={'center':'Wins vs shutouts'}, xlabel='W', ylabel='SHO'>
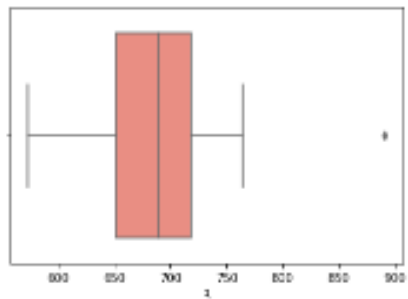


In this reg plot it is apparent that SHO is positively correlated with W data.only few datapoints are close to best fit line it means that the feature SHO and W are not having linear relationship

```
In [5]:   1  plt.figure(figsize=(5,5))
          2  plt.title('Wins vs 2B')
          3  sns.set_style="white_grid"
          4  sns.regplot(x="W",y="2B",data=df,scatter_kws = {'color': 'purple'}, line_kws = {'color': 'r'})
```

Out[5]: <AxesSubplot:title={'center':'Wins vs 2B'}, xlabel='W', ylabel='2B'>



In this reg plot it is apparent that 2B is positively correlated with W data. Only few data points are close to best fit line it means that the feature 2B and W are not having linear relationship

```
[6]:   1  plt.figure(figsize=(5,5))
       2  plt.title('Wins vs HR')
       3  sns.set_style="white_grid"
       4  sns.regplot(x="W",y="HR",data=df,scatter_kws = {'color': 'purple'}, line_kws = {'color': 'r'})
```

[6]: <AxesSubplot:title={'center':'Wins vs HR'}, xlabel='W', ylabel='HR'>



In this reg plot it is apparent that HR is positively correlated with W data. Only few datapoints are close to best fit line it means that the feature HR and W are not having linear relationship

# IDENTIFYING OUTLIERS:

```
In [209]:   1  sns.boxplot(df['R'],color="salmon")
```

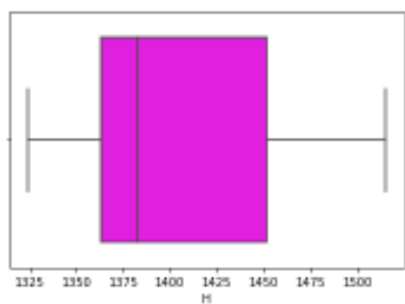Out[209]: <AxesSubplot:xlabel='R'>



feature R has outlier

```
In [210]:   1  sns.boxplot(df['AB'],color="brown")
```
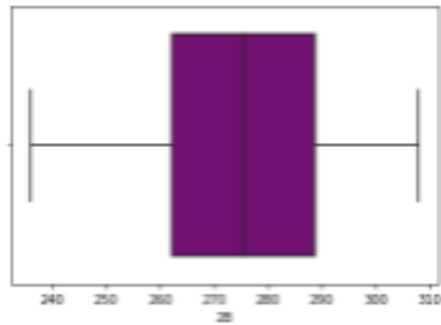
Out[210]: <AxesSubplot:xlabel='AB'>



```
In [211]:   1  sns.boxplot(df['H'],color=(1,0,1))
```
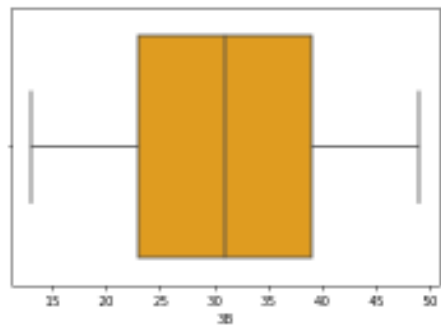
Out[211]: <AxesSubplot:xlabel='H'>

```
In [212]:   1  sns.boxplot(df['2B'],color="Purple")
```
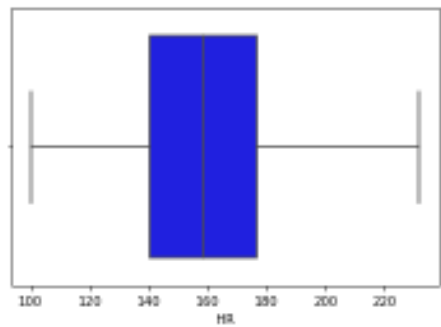
Out[212]:  <AxesSubplot:xlabel='2B'>



```
In [213]:   1  sns.boxplot(df['3B'],color="orange")
```

Out[213]:  <AxesSubplot:xlabel='3B'>



```
In [214]:   1  sns.boxplot(df['HR'],color=(0,0,1))
```
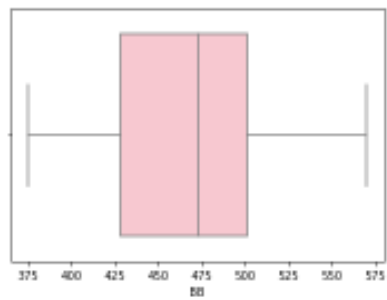
Out[214]:  <AxesSubplot:xlabel='HR'>



```
In [215]:   1  sns.boxplot(df['BB'],color="pink")
```
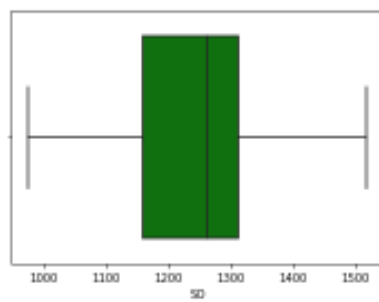
Out[215]:  <AxesSubplot:xlabel='BB'>

```
In [215]:   1  sns.boxplot(df['BB'],color="pink")
```
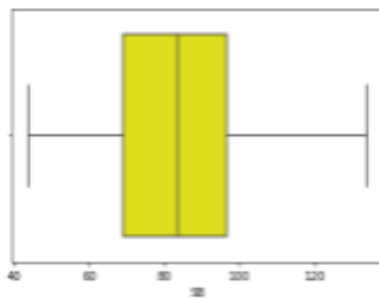
Out[215]: <AxesSubplot:xlabel='BB'>



```
In [221]:   1  sns.boxplot(df['SO'],color="green")
```

Out[221]: <AxesSubplot:xlabel='SO'>
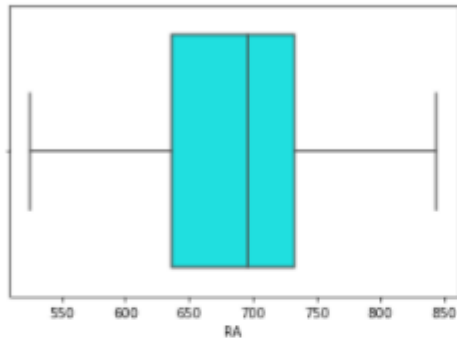


```
In [217]:   1  sns.boxplot(df['SB'],color="yellow")
```

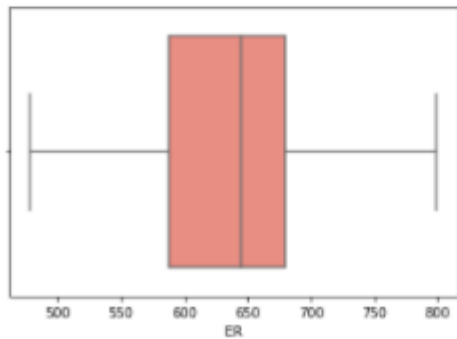Out[217]: <AxesSubplot:xlabel='SB'>

```
In [218]:    1  sns.boxplot(df['RA'],color=(0,1,1))
```

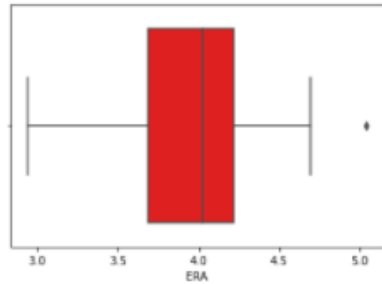Out[218]: <AxesSubplot:xlabel='RA'>



```
In [219]:    1  sns.boxplot(df['ER'],color="salmon")
```

Out[219]: <AxesSubplot:xlabel='ER'>



```
In [220]:    1  sns.boxplot(df['ERA'],color="red")
```
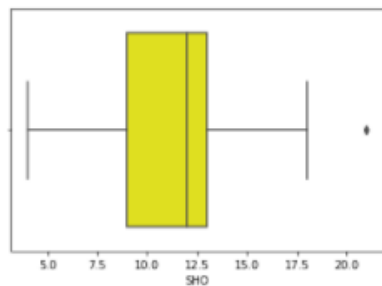
Out[220]: <AxesSubplot:xlabel='ERA'>



feature ERA has outlier

```
In [143]:    1  sns.boxplot(df['SHO'],color=(1,1,0))
```

Out[143]: <AxesSubplot:xlabel='SHO'>
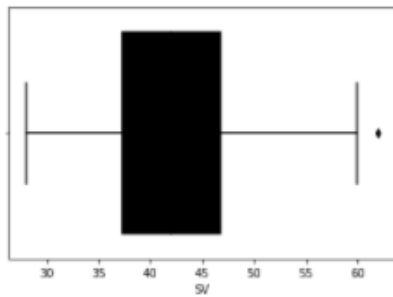


feature SHO has outlier

feature SHO has outlier

In [144]: 1 sns.boxplot(df['SV'],color=(0,0,0))

Out[144]: <AxesSubplot:xlabel='SV'>



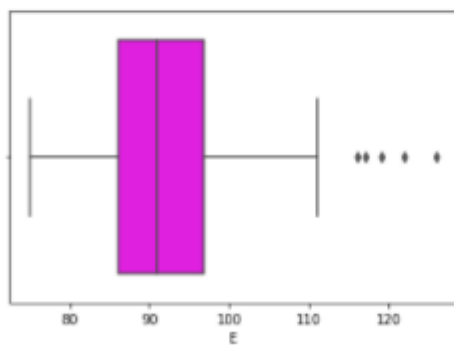feature SV has outlier

In [145]: 1 sns.boxplot(df['CG'],color=(0,1,1))

Out[145]: <AxesSubplot:xlabel='CG'>



In [146]: 1 sns.boxplot(df['E'],color=(1,0,1))

Out[146]: <AxesSubplot:xlabel='E'>



The features R,ERA,SHO,SV,E are having outliers

**REMOVING OUTLIERS**

```
In [7]:    1  z = np.abs(zscore(df))
           2  print(np.where(z > 3))

(array([5], dtype=int64), array([1], dtype=int64))

In [8]:    1  data_mod = df[(z<3).all(axis=1)]
           2  data_mod.shape

Out[8]:  (29, 17)

In [9]:    1  ((30-29)/30)*100

Out[9]:  3.3333333333333335
```

After using z score it has only 3.3333% data loss

IQR TO REMOVE OUTLIERS:

```
In [230]:   1  data=df
            2  Q1=data.quantile(0.25)
            3  Q3=data.quantile(0.75)
            4  IQR=Q3-Q1
            5  df_new=data[~((data<(Q1-1.5*IQR))| (data>(Q1+1.5*IQR))).any(axis=1)]

In [26]:    1  df_new.shape

Out[26]:  (6, 17)

In [27]:    1  ((30-6)/30)*100

Out[27]:  80.0
```
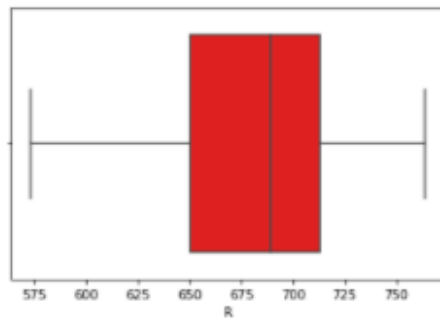
80% data loss so I'm choosing zscore method

AFTER USING ZSCORE:

```
In [229]:   1  sns.boxplot(data_mod['E'],color="purple")
Out[229]:   <AxesSubplot:xlabel='E'>
```



After using z_score the outliers in the features R and ERA are removed but the features SV,SHO,E are still having outliers

**HANDLING SKEWNESS**:

```
In [10]:   1  for col in enumerate(list(data_mod.columns.values)):
           2      print(col[1],"=",data_mod[col[1]].skew())

W = 0.11901344569985461
R = -0.21536363420992782
AB = 0.16957316834729352
H = 0.7837722117274881
2B = -0.335303936110201
3B = 0.09012434653848651
HR = 0.45086158125803544
BB = 0.15119282971519954
SO = -0.2338149185462262
SB = 0.4949657663368456
RA = 0.018155177145956613
ER = 0.018460990156758887
ERA = 0.016693217783651695
CG = 0.8549795901105167
SHO = 0.5269430585305683
SV = 0.6274804879503074
E = 0.8402711976867623
```

the features H,CG,SHO,SV,E having skewness

 To stabilize variance, make the data more normal distribution like, improve the validity of measures of association and to remove skewness I have used power transformation

```
1  col_s=['H','CG','SHO','SV','E']
2  data_clean=data_mod
3  from sklearn.preprocessing import power_transform
4  data_clean[col_s]=power_transform (data_mod[col_s])
```

```
1  for col in enumerate(list(data_clean.columns.values)):
2      print(col[1],"=",data_clean[col[1]].skew())
```

```
W = 0.11901344569985461
R = -0.21536363420992782
AB = 0.16957316834729352
H = 0
2B = -0.335303936110201
3B = 0.09012434653848651
HR = 0.45086158125803544
BB = 0.15119282971519954
SO = -0.2338149185462262
SB = 0.4949657663368456
RA = 0.0181551771145956613
ER = 0.018460990156758887
ERA = 0.016693217783651695
CG = -0.045947323970913174
SHO = 0.0005293650356868707
SV = -0.0009249344497408174
E = 0.06558547868786976
```

Skewness is completely removed

I have used Power transform on skewed data to make it symmetric, and then fit it to a symmetric distribution

```
1  sns.boxplot(data_clean['E'],color="purple")
```

: <AxesSubplot:xlabel='E'>



Using zscore the outliers in the feature E is completely removed and the data loss after removing outlier is 3.333333%

```
In [38]:    1  sns.boxplot(data_clean['SHO'],color="blue")
```

Out[38]:  &lt;AxesSubplot:xlabel='SHO'&gt;



```
In [39]:    1  z3 = np.abs(stats.zscore(data_clean['SV']))
            2  data_clean['SV'] = data_clean[(z3<3)]
            3  sns.boxplot(data_clean['SV'],color="brown")
```

Out[39]:  &lt;AxesSubplot:xlabel='SV'&gt;



I have used zscore to remove outliers in the features SV, SHO and E.

**SCALING:**

___feature scaling transforming un scaled data into scaled data using min max scalining technique

```
:   1  from sklearn.preprocessing import MinMaxScaler
    2  scaler=MinMaxScaler()
    3  scaled = scaler.fit_transform(x1)
```

**MODELLING**:

___MODELING

```
In [108]:   1  from sklearn.neighbors import KNeighborsRegressor
            2  from sklearn.svm import SVR
            3  from sklearn.tree import DecisionTreeRegressor
            4  from sklearn.linear_model import LinearRegression
            5  from sklearn.linear_model import Ridge
            6  from sklearn.linear_model import Lasso
            7  from sklearn.ensemble import RandomForestRegressor
            8  from sklearn.ensemble import GradientBoostingRegressor
            9  import xgboost as xgb
           10  from sklearn.model_selection import train_test_split
           11
```

```
In [109]:   1  x=scaled
            2  y=y1
```

```
In [110]:   1  xtrain,xtest,ytrain,ytest=train_test_split(x,y,test_size=0.30,random_state=7)
            2  from sklearn.metrics import mean_absolute_error,mean_squared_error,r2_score
            3  models=[KNeighborsRegressor(),SVR(),DecisionTreeRegressor(),LinearRegression(),Lasso(),Ridge(),
            4         RandomForestRegressor(),GradientBoostingRegressor(),xgb.XGBRegressor(objective="reg:squarederror")]
            5  maelist=[]
            6  mselist=[]
            7  rmselist=[]
            8  r2list=[]
```

```
In [111]:   1  def create_model(model):
            2      m=model
            3      m.fit(xtrain,ytrain)
            4      p=m.predict(xtest)
            5
            6      mae=mean_absolute_error(p,ytest)
            7      mse=mean_squared_error(p,ytest)
            8      rmse=np.sqrt(mean_squared_error(p,ytest))
            9      r2=r2_score(ytest,p)
           10
           11      maelist.append(mae)
           12      mselist.append(mse)
           13      rmselist.append(rmse)
           14      r2list.append(r2)
```

```
KNeighborsRegressor()
Mean absolute error 5.266666666666666
Mean squared error 35.76444444444443
Root Mean squared error 5.980338154690287
R2 Score 0.6688477366255146
------------------------------------------------------------------------------
SVR()
Mean absolute error 7.641601405900634
Mean squared error 76.28996804828544
Root Mean squared error 8.734412862252702
R2 Score 0.293611406960031995
------------------------------------------------------------------------------
DecisionTreeRegressor()
Mean absolute error 1.7777777777777777
Mean squared error 5.111111111111111
Root Mean squared error 2.260776661041756
R2 Score 0.9526748971193416
------------------------------------------------------------------------------
LinearRegression()
Mean absolute error 7.894919286223336e-15
Mean squared error 1.12193550964766613e-28
Root Mean squared error 1.0592145720521698e-14
R2 Score 1.0
------------------------------------------------------------------------------
Lasso()
Mean absolute error 3.171414003610454
Mean squared error 15.237323453470554
Root Mean squared error 3.9035014350542454
R2 Score 0.8589136717271245
------------------------------------------------------------------------------
Ridge()
Mean absolute error 2.8099060243455765
Mean squared error 8.601635421789048
Root Mean squared error 2.9328544835687036
R2 Score 0.9203552275760273
------------------------------------------------------------------------------
RandomForestRegressor()
Mean absolute error 3.034444444444443
Mean squared error 12.24765555555555
Root Mean squared error 3.4996650633389974
R2 Score 0.8865957818930041

------------------------------------------------------------------------------
GradientBoostingRegressor()
Mean absolute error 2.3271601042110572
Mean squared error 7.157694017274431
Root Mean squared error 2.6753867042493935
R2 Score 0.9337250553956071
------------------------------------------------------------------------------
XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
             colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,
             importance_type='gain', interaction_constraints='',
             learning_rate=0.300000012, max_delta_step=0, max_depth=6,
             min_child_weight=1, missing=nan, monotone_constraints='()',
             n_estimators=100, n_jobs=4, num_parallel_tree=1, random_state=0,
             reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
             tree_method='exact', validate_parameters=1, verbosity=None)
Mean absolute error 2.9089228312174478
Mean squared error 10.716167470886527
Root Mean squared error 3.2735557839888
R2 Score 0.900776227121421
------------------------------------------------------------------------------
Minimum Mean Absolute error is shown by  LinearRegression() 7.894919286223336e-15
Minimum Mean squared error is shown by  LinearRegression() 1.12193550964766613e-28
Minimum Root Mean squared error is shown by  LinearRegression() 1.0592145720521698e-14
Maximun R2 Score is shown by  LinearRegression() 1.0
```

**Cross Validation:**

In order to avoid over fitting, Cross-validation is used to estimate the skill of a machine learning model on unseen data.

```
n [139]:   1  from sklearn.model_selection import cross_val_score
           2  k=KNeighborsRegressor()
           3  scores=cross_val_score(k,x,y,scoring='r2',cv=5)
           4  scorel.append(scores)
           5  scores
```

```
ut[139]: array([ 0.5065    , -0.22736842,  0.56332915,  0.59501639,  0.90259434])
```

```
n [140]:   1  from sklearn.model_selection import cross_val_score
           2  svr=SVR()
           3  scores=cross_val_score(svr,x,y,scoring='r2',cv=5)
           4  scorel.append(scores)
           5  scores
```

```
ut[140]: array([ 0.19800765, -0.43552919,  0.14123261, -0.19247998,  0.41517121])
```

```
n [141]:   1  from sklearn.model_selection import cross_val_score
           2  dt=DecisionTreeRegressor()
           3  scores=cross_val_score(dt,x,y,scoring='r2',cv=5)
           4  scorel.append(scores)
           5  scores
```

```
ut[141]: array([0.88068182, 0.18421053, 0.87787537, 0.89016393, 0.89976415])
```

```
n [142]:   1  from sklearn.model_selection import cross_val_score
           2  lr=LinearRegression()
           3  scores=cross_val_score(lr,x,y,scoring='r2',cv=5)
           4  scorel.append(scores)
           5  scores
```

```
ut[142]: array([1., 1., 1., 1., 1.])
```

```
n [143]:   1  from sklearn.model_selection import cross_val_score
           2  l=Lasso()
           3  scores=cross_val_score(l,x,y,scoring='r2',cv=5)
           4  scorel.append(scores)
           5  scores
```

```
ut[143]: array([0.88553127, 0.84887851, 0.72013556, 0.64974046, 0.88639659])
```

```
In [144]:  1  from sklearn.model_selection import cross_val_score
           2  rid=Ridge()
           3  scores=cross_val_score(rid,x,y,scoring='r2',cv=5)
           4  scorel.append(scores)
           5  scores

Out[144]:  array([0.8984672 , 0.76298028, 0.90904629, 0.92935865, 0.99225316])

In [145]:  1  from sklearn.model_selection import cross_val_score
           2  rf=RandomForestRegressor()
           3  scores=cross_val_score(rf,x,y,scoring='r2',cv=5)
           4  scorel.append(scores)
           5  scores

Out[145]:  array([0.9720642 , 0.95058684, 0.859889  , 0.86196803, 0.9760625 ])

In [146]:  1  from sklearn.model_selection import cross_val_score
           2  gb=GradientBoostingRegressor()
           3  scores=cross_val_score(gb,x,y,scoring='r2',cv=5)
           4  scorel.append(scores)
           5  scores

Out[146]:  array([0.98639233, 0.8546094 , 0.92123603, 0.90768891, 0.98418366])

In [147]:  1  from sklearn.model_selection import cross_val_score
           2  xb=xgb.XGBRegressor()
           3  scores=cross_val_score(xb,x,y,scoring='r2',cv=5)
           4  scorel.append(scores)
           5  scores

Out[147]:  array([0.82751376, 0.94214528, 0.88634985, 0.89240543, 0.94730646])
```

**Difference of predicted model and cross validation score:**

- KNeighborsRegressor - 0.2337466
- SVR()-0.12155981
- DecisionTreeRegressor()-0.05291075
- LinearRegression() -0
- Lasso - 0.02748291
- Ridge()-0.07189794
- RandomForestRegressor()-0.08946672
- GradientBoostingRegressor()-0.0504586
- XGBRegressor-0.046530

from the observation Linear regression model model has least difference so I'm selecting
Linear regression as best model

**Hyper Tuning:**

___Hyper tuning

```
In [150]:  1  from sklearn.model_selection import GridSearchCV
           2  parameters = { "normalize":[True, False], "copy_X":[True, False],
           3                 "fit_intercept": [True, False]
           4               }
           5  grid = GridSearchCV(LinearRegression(), param_grid = parameters, cv = 5, scoring = "r2")
```

```
In [151]:  1  grid.fit(xtrain,ytrain)
           2
           3  print("Best_parameters",grid.best_params_)
           4
```

Best_parameters {'copy_X': True, 'fit_intercept': True, 'normalize': True}

**Best parameters**: {'copy_X': True, 'fit_intercept': True, 'normalize': True}
**Modelling using best parameter and best model:**

```
In [154]:  1  xtrain,xtest,ytrain,ytest=train_test_split(x,y,test_size=0.25,random_state=1)
           2  from sklearn.metrics import mean_absolute_error,mean_squared_error,r2_score,accuracy_score
           3  model=LinearRegression(copy_X=True,fit_intercept=True,normalize=True)
           4  model.fit(xtrain,ytrain)
           5  p=model.predict(xtest)
           6  acc=model.score(xtest,ytest)
           7  mae=mean_absolute_error(p,ytest)
           8  mse=mean_squared_error(p,ytest)
           9  rmse=np.sqrt(mean_squared_error(p,ytest))
          10  r2=r2_score(ytest,p)
          11  print('Accuracy',acc)
          12  print('Mean absolute error',mae)
          13  print('Mean squared error',mse)
          14  print('Root Mean squared error',rmse)
          15  print('r2 score',r2)
          16
```

```
Accuracy 1.0
Mean absolute error 1.0658141036401503e-14
Mean squared error 2.0194839173657902e-28
Root Mean squared error 1.4210854715202004e-14
r2 score 1.0
```

Final model after hyper tuning its retaining 100% accuracy and error values got reduced

**Conclusion:**

I have developed a model to predict number of wins with 100% accuracy

**Saving the model**

```
In [79]:  1  from joblib import dump
          2  dump(model, 'model_baseball.joblib')
```

Out[79]: ['model_baseball.joblib']

```
In [80]:  1  from joblib import load
          2  loaded = load('model_baseball.joblib')
```