

Author: Ganesh Kale

Week#3

Date: Dec 14, 2021

Improve Marketing Promotion of MLB Game:

Major League Baseball (MLB) is a professional baseball organization and the oldest major professional sports league in the world. As of 2021, a total of 30 teams play in Major League Baseball—15 teams in the National League (NL) and 15 in the American League (AL)—with 29 in the United States and 1 in Canada.

Problem Statement:

- What is the best time to run Marketing promotion to increase attendance for Dodgers Game?
- Propose best date or day of the week or Day & month to run marketing promotion to increase audience.

Data Information:

The data collected here is from Dodgers MLB game from 2012. This data have different features such as -

- month: Month of the Game - String format
- day: The day/date of the Game - Number
- attend: Number of audience attended to see the game - numerical data type
- day_of_week: The day of the week - String format- String format
- temp: Temperature of the game day - numerical data type
- skies: Sky condition of the game day - String format- String format
- day_night: Game was at day or night - String format- String format
- cap: Cap distributed - Boolean
- shirt: shirt distributed - Boolean
- fireworks: fireworks distributed/happened - Boolean
- bobblehead: bobblehead distributed - Boolean

import required packages

```
In [1]: import numpy as np # for numeric operations
import pandas as pd # for data manipulation
import matplotlib.pyplot as plt # for data visualization
import seaborn as sns # for data visualization
from sklearn.preprocessing import LabelEncoder

from IPython.core.interactiveshell import InteractiveShell
InteractiveShell.ast_node_interactivity = 'all'
```

```
In [2]: # import the dataset

dodgers = pd.read_csv("Data/dodgers.csv")
```

```
In [3]: # display shape and head of df

dodgers.shape
dodgers.head()
```

```
Out[3]: (81, 12)
```

	month	day	attend	day_of_week	opponent	temp	skies	day_night	cap	shirt	fireworks	bobblehead
0	APR	10	56000	Tuesday	Pirates	67	Clear	Day	NO	NO	NO	NO
1	APR	11	29729	Wednesday	Pirates	58	Cloudy	Night	NO	NO	NO	NO
2	APR	12	28328	Thursday	Pirates	57	Cloudy	Night	NO	NO	NO	NO
3	APR	13	31601	Friday	Padres	54	Cloudy	Night	NO	NO	YES	NO
4	APR	14	46549	Saturday	Padres	57	Cloudy	Night	NO	NO	NO	NO

Data Exploratory Analysis:

```
In [4]: # statistical information about data

dodgers.describe(include='all')
```

```
Out[4]:
```

	month	day	attend	day_of_week	opponent	temp	skies	day_night	cap	shirt	fireworks	bobblehead
count	81	81	81000000	81	81	81	81	81000000	81	81	81	81
unique	7	NaN	NaN	7	17	NaN	2	2	2	2	2	2
top	MAY	NaN	NaN	Tuesday	Giants	NaN	Clear	NO	NO	NO	NO	NO
freq	18	NaN	NaN	13	9	NaN	62	66	79	78	67	70
mean	NaN	NaN	6135802	41040.074074	NaN	NaN	73.148148	NaN	NaN	NaN	NaN	NaN
std	NaN	9.605666	8297.539460	NaN	NaN	8.317318	NaN	NaN	NaN	NaN	NaN	NaN
min	NaN	1.000000	24312.000000	NaN	NaN	54.000000	NaN	NaN	NaN	NaN	NaN	NaN
25%	NaN	NaN	8000000	34493.000000	NaN	NaN	67.000000	NaN	NaN	NaN	NaN	NaN
50%	NaN	NaN	15000000	40284.000000	NaN	NaN	73.000000	NaN	NaN	NaN	NaN	NaN
75%	NaN	NaN	25000000	46568.000000	NaN	NaN	79.000000	NaN	NaN	NaN	NaN	NaN
max	NaN	31.000000	56000.000000	NaN	NaN	85.000000	NaN	NaN	NaN	NaN	NaN	NaN

```
In [5]: # data type information

dodgers.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 81 entries, 0 to 80
Data columns (total 12 columns):
 # Column            Non-Null Count  Dtype
---  --
 0 month              81 non-null    object
 1 day                81 non-null    int64
 2 attend             81 non-null    int64
 3 day_of_week        81 non-null    object
 4 opponent           81 non-null    object
 5 temp               81 non-null    int64
 6 skies              81 non-null    object
 7 day_night          81 non-null    object
 8 cap                81 non-null    object
 9 shirt              81 non-null    object
10 fireworks          81 non-null    object
11 bobblehead         81 non-null    object
dtypes: int64(3), object(9)
memory usage: 7.7+ KB
```

```
In [6]: # see if any missing values in the data

import missingno as msn
msno.bar(dodgers, color='b')
```

```
Out[6]:
```

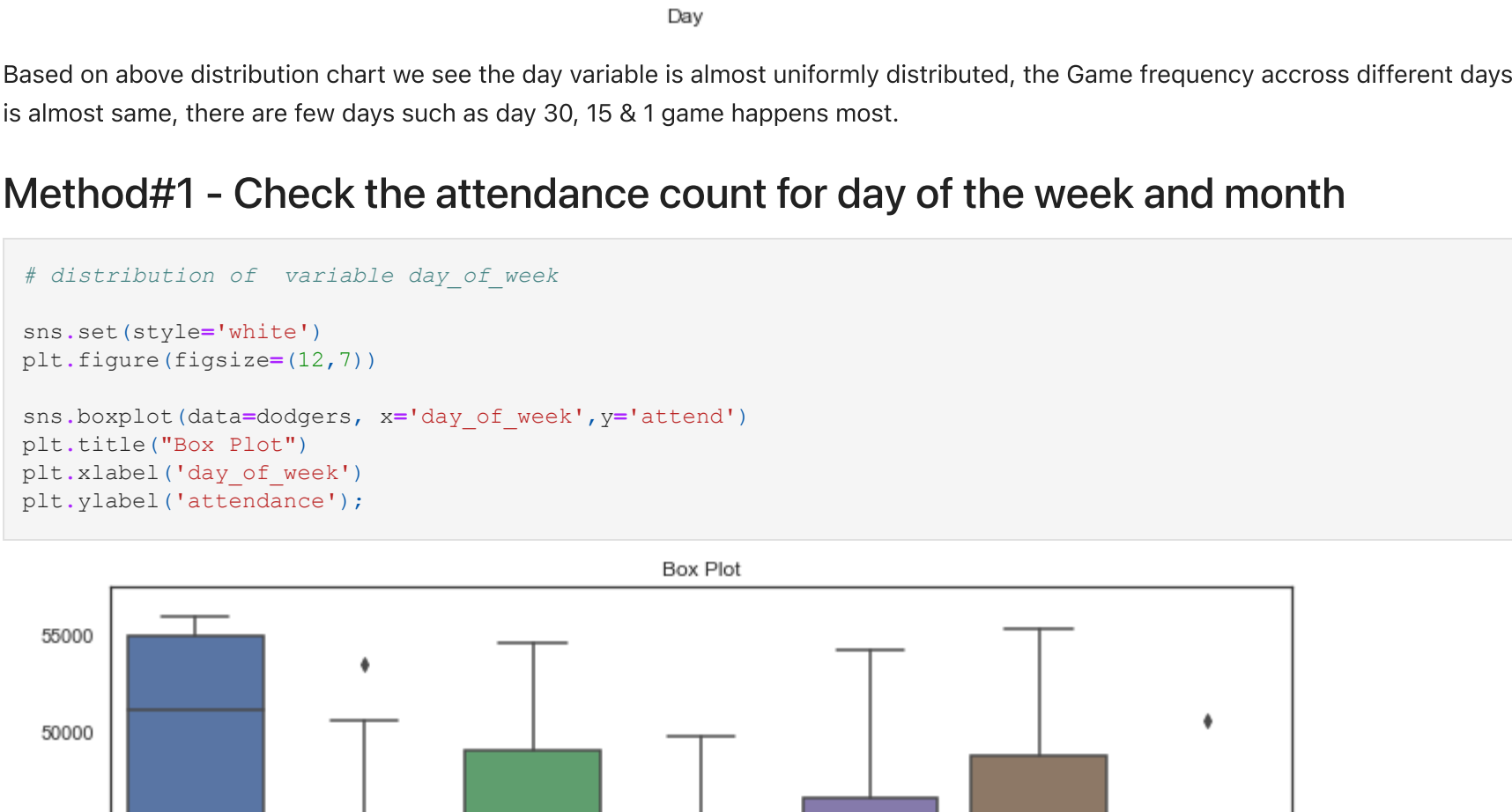
Based on above bar charts we can see that there are no missing values in the data set

```
In [7]: # distribution of independent variable attend

sns.set(style='white')
plt.figure(figsize=(12,7))

sns.distplot(a=dodgers['attend'],color='g', bins=10)

plt.title("Distribution of Audience")
plt.xlabel("Audience attendance")
plt.ylabel("Count");
```



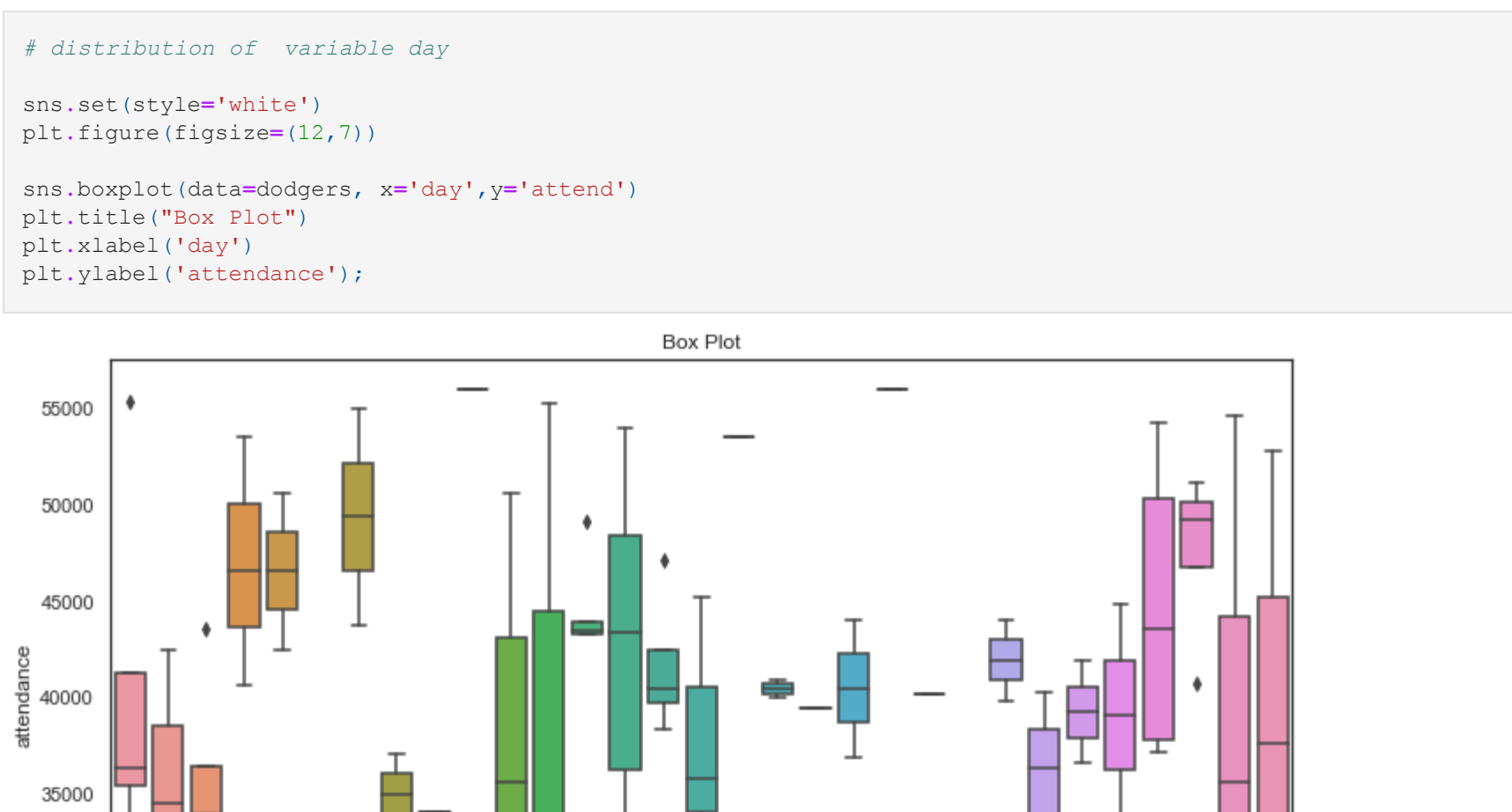
Based on above distribution plot, we see that the attendance is normally distributed and average attendance is around 40K in each game. Also, we do not see any outliers in the attendance.

```
In [8]: # distribution of variable day

sns.set(style='white')
plt.figure(figsize=(12,7))

sns.distplot(a=dodgers['day'],color='g', bins=10)

plt.title("Distribution of Day")
plt.xlabel("Day")
plt.ylabel("Count");
```



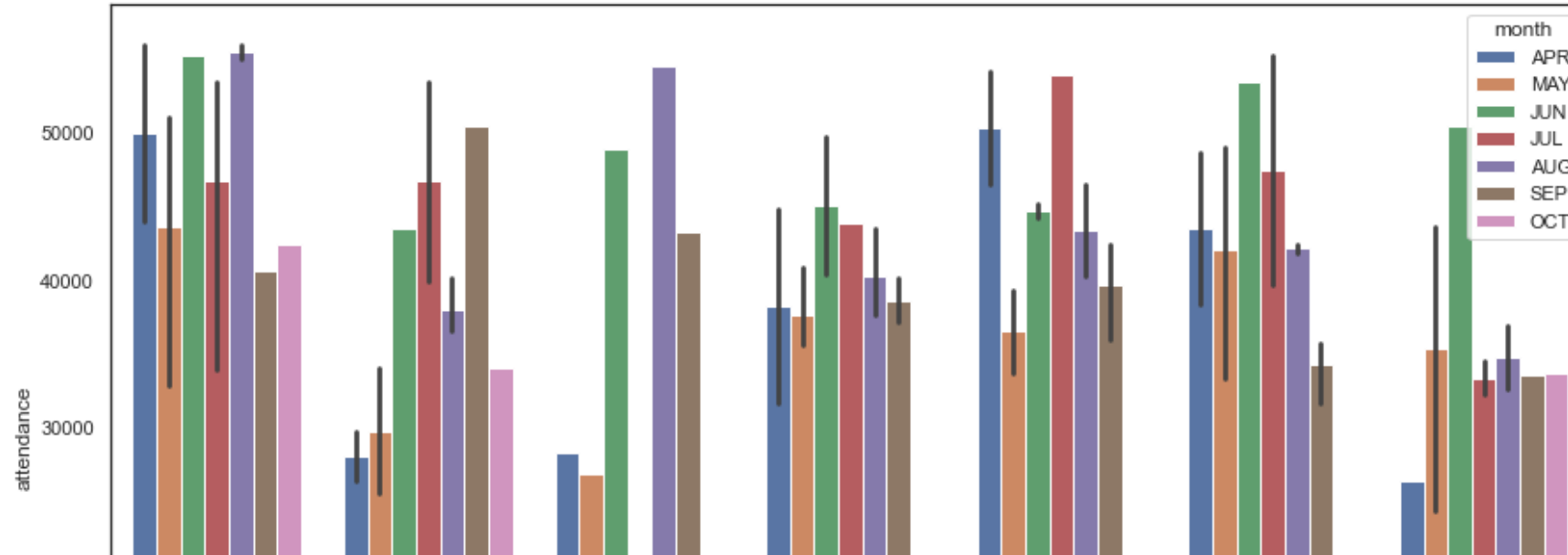
Based on above distribution chart we see the day variable is almost uniformly distributed, the Game frequency across different days is almost same, there are few days such as day 30, 15 & 1 game happens most.

Method#1 - Check the attendance count for day of the week and month

```
In [9]: # distribution of variable day_of_week

sns.set(style='white')
plt.figure(figsize=(12,7))

sns.boxplot(data=dodgers, x='day_of_week', y='attend')
```

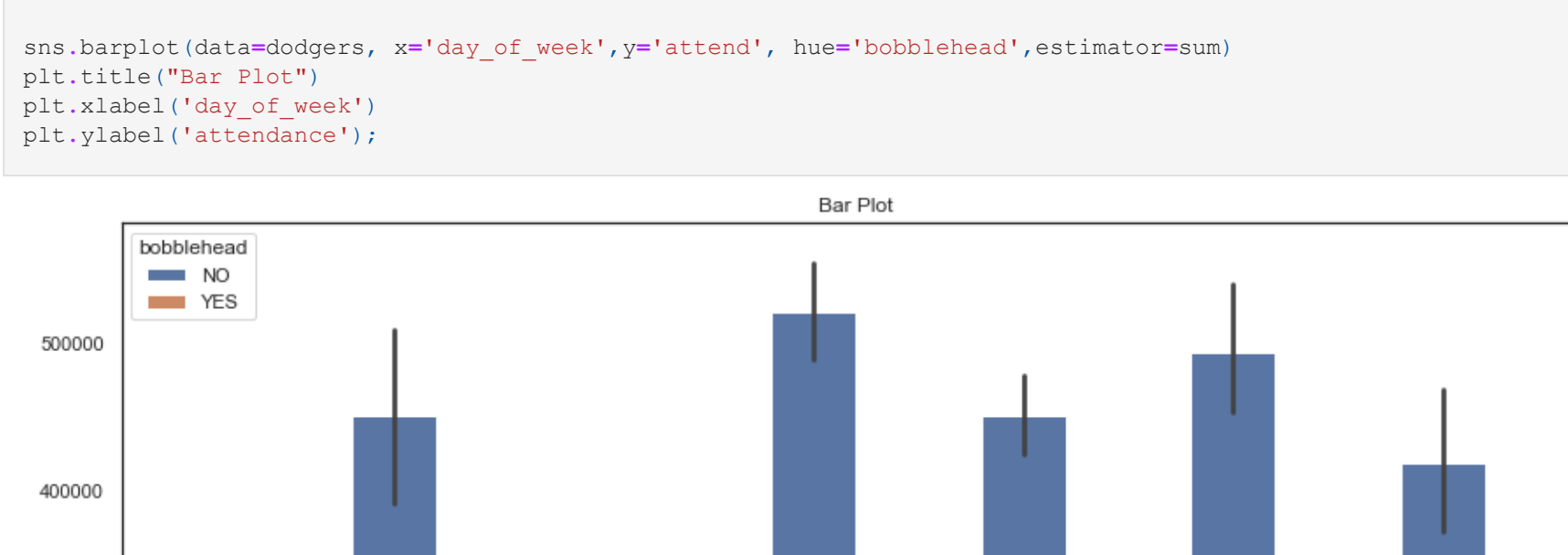


Based on above box plots, we see that Tuesday is the day wherein median attendance is highest compared to other days and overall attendance is higher on Tuesdays. The 25 percentile attendance on Tuesdays is higher than all other days median or 50 percentile attendance.

```
In [10]: # distribution of variable day

sns.set(style='white')
plt.figure(figsize=(12,7))

sns.boxplot(data=dodgers, x='day', y='attend')
```



Based on above box chart the day 7th which has highest median value of attendance compared to other days and minimum number of attendance on this day is more than most of the days median or 75% percentile of other days. This is the day where attendance is more.

```
In [11]: # distribution of variable day_of_week & month

sns.set(style='white')
plt.figure(figsize=(12,9))

sns.barplot(data=dodgers, x='day_of_week', y='attend', hue='month', estimator=np.mean)
```

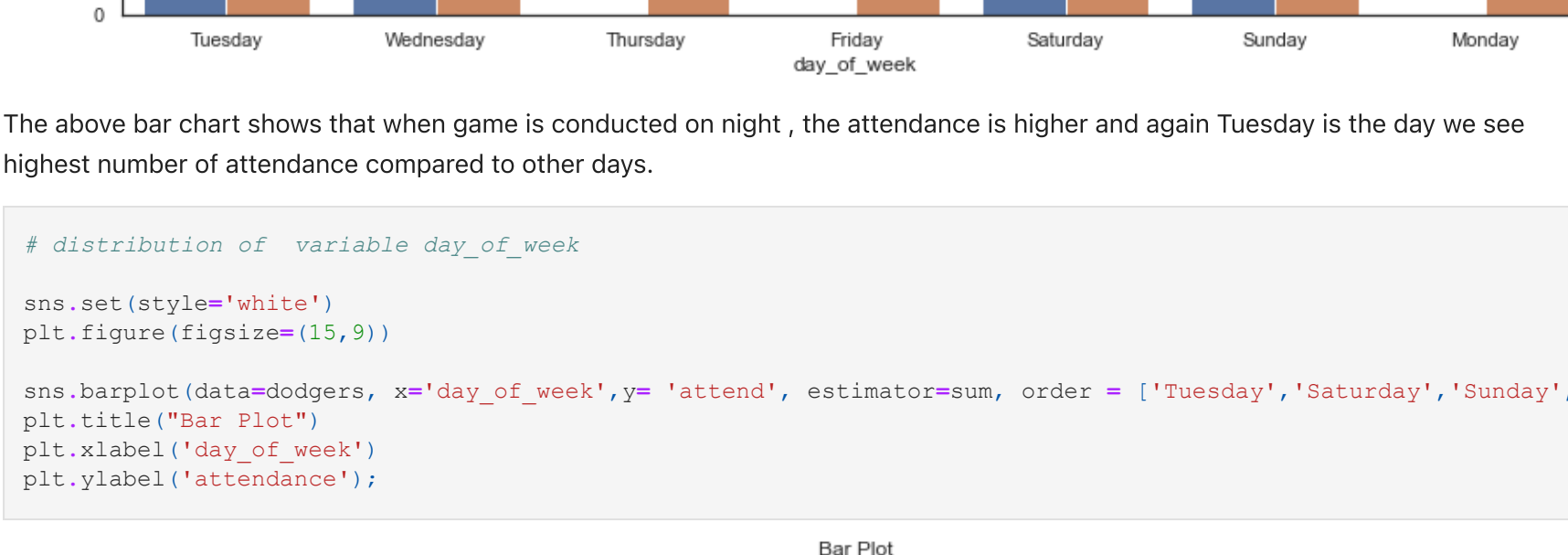


Based on above bar charts we see that Tuesday is the day wherein average attendance is highest across all the month when compared to different days of the week and months. Tuesdays got more than average 40K attendance every month which is above average of other days.

```
In [12]: # day of week distribution to see attendance

sns.set(style='white')
plt.figure(figsize=(12,9))

sns.barplot(data=dodgers, x='day_of_week', y='attend', hue='bobblehead', estimator=sum)
```

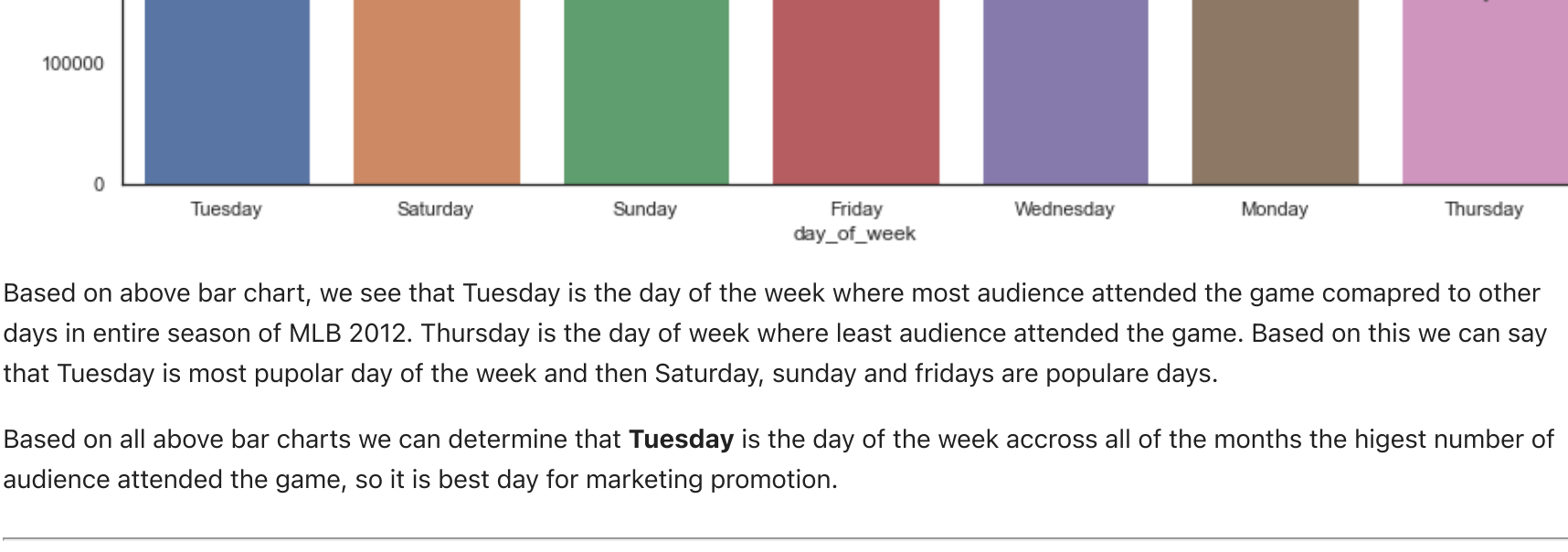


Based on above bar chart, we see that when bobble head is distributed or offered the attendance is higher so it is playing important role in marketing promotion. On Tuesday we see the attendance is highest when bobble head is offered.

```
In [13]: # day of week distribution to see attendance

sns.set(style='white')
plt.figure(figsize=(12,9))

sns.barplot(data=dodgers, x='day_of_week', y='attend', hue='day_night', estimator=sum)
```



The above bar chart shows that when game is conducted on night, the attendance is higher and again Tuesday is the day where we see highest number of attendance compared to other days.

```
In [14]: # distribution of variable day_of_week

sns.set(style='white')
plt.figure(figsize=(12,9))

sns.barplot(data=dodgers, x='day_of_week', y='attend', estimator=sum, order = ['Tuesday','Saturday','Sunday'],'attend')
```



Based on above bar chart, we see that Tuesday is the day of the week where most audience attended the game compared to other days in entire season of MLB 2012. Thursday is the day of week where least audience attended the game. Based on this we can say that Tuesday is most popular day of the week and then Saturday, Sunday and Fridays are popular days.

Based on all above bar charts we can determine that Tuesday is the day of the week across all of the months the highest number of audience attended the game, so it is best day for marketing promotion.

Method#2 - Finding factors that impact attendance increase, correlating factors with attendance and see what are the different factors correlate and their strength and direction of correlation.

```
In [15]: # scatter plot to see correlations

sns.set(style='white')
plt.figure(figsize=(12,7))

sns.scatterplot(data=dodgers, y='attend', x='day')
```



Based on above chart we do not see strong correlation but as day increases the attendance as well.

```
In [16]: # scatter plot to see correlations

sns.set(style='white')
plt.figure(figsize=(12,7))

sns.scatterplot(data=dodgers, y='attend', x='temp')
```



Based on above scatter plot we can say that as temperatures become warmer audience increases but as temperature gets hotter above 80F the audience decreases.

Feature Engineering:

```
In [17]: # create new variables by combining day of the week and month, day, day of week and month

dodgers['day_of_week_month'] = dodgers['month']+'-'+dodgers['day_of_week']
```

```
In [18]: # display head of the df

dodgers.head()
```

```
Out[18]:
```

	month	day	attend	day_of_week	opponent	temp	skies	day_night	cap	shirt	fireworks	bobblehead	day_of_week_month
0	APR	10	56000	Tuesday	Pirates	67	Clear	Day	NO	NO	NO	NO	NO
1	APR	11	29729	Wednesday	Pirates	58	Cloudy	Night	NO	NO	NO	NO	APR-Wednesday
2	APR	12	28328	Thursday	Pirates	57	Cloudy	Night	NO	NO	NO	NO	APR-Thursday
3	APR	13	31601	Friday	Padres	54	Cloudy	Night	NO	NO	YES	NO	APR-Friday
4	APR	14	46549	Saturday	Padres	57	Cloudy	Night	NO	NO	NO	NO	APR-Saturday

```
In [19]: # encode the categorical features

from sklearn.preprocessing import StandardScaler, OrdinalEncoder

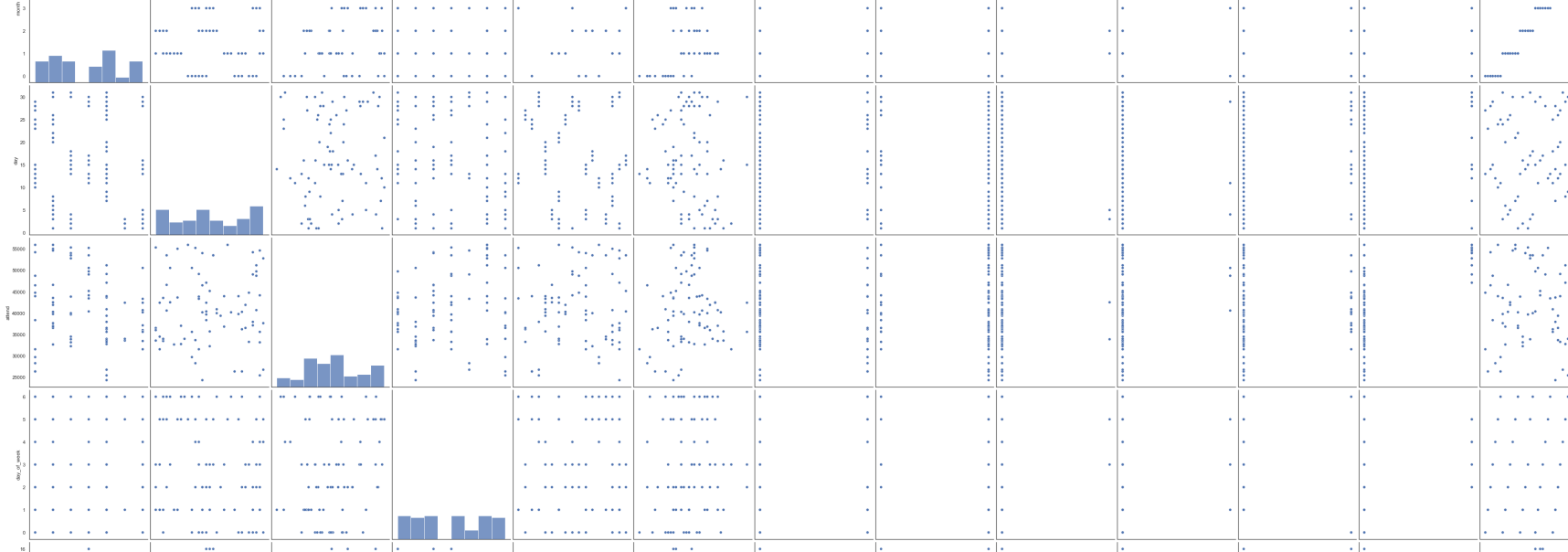
enc = OrdinalEncoder()

sc = StandardScaler()
```

```
In [20]: # scatter plot of day of week enc and attend

sns.set(style='white')
plt.figure(figsize=(12,7))

sns.scatterplot(data=dodgers, x='day_of_week', y='attend')
```

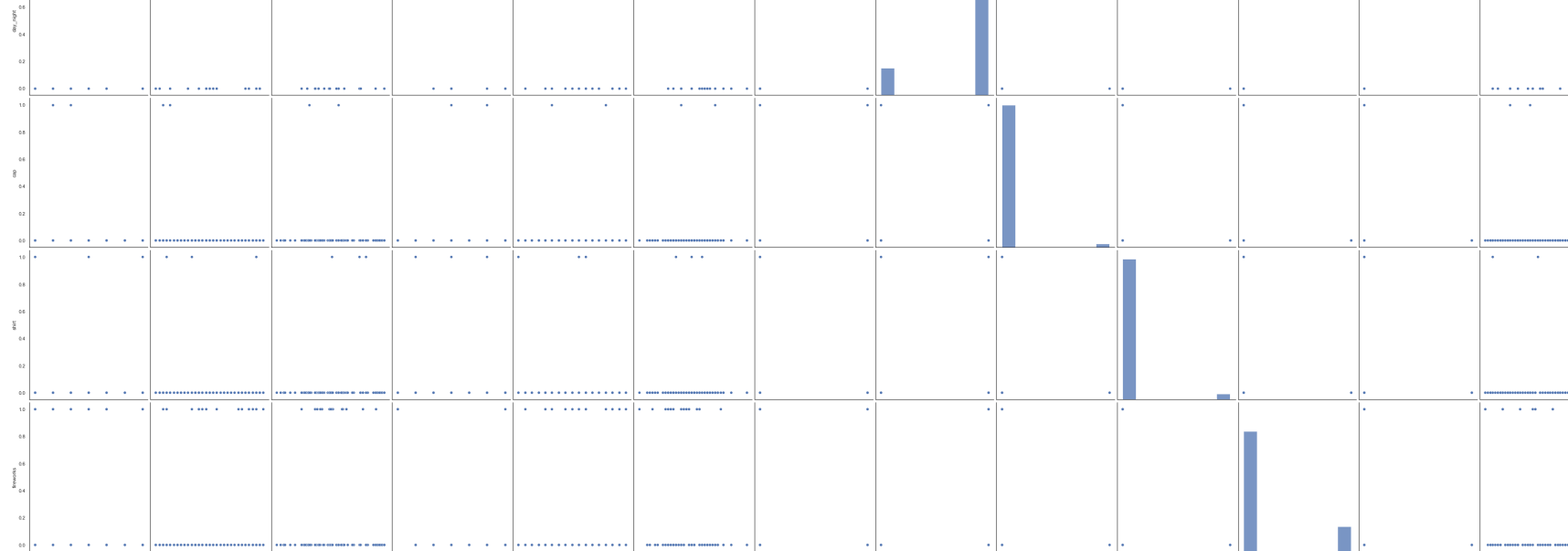


Based on above scatter plot we do not see any strong correlation of day of the week with attendance.

```
In [21]: # scatter plot of day of week enc and attend

sns.set(style='white')
plt.figure(figsize=(12,7))

sns.scatterplot(data=dodgers, x='day_of_week_month', y='attend')
```



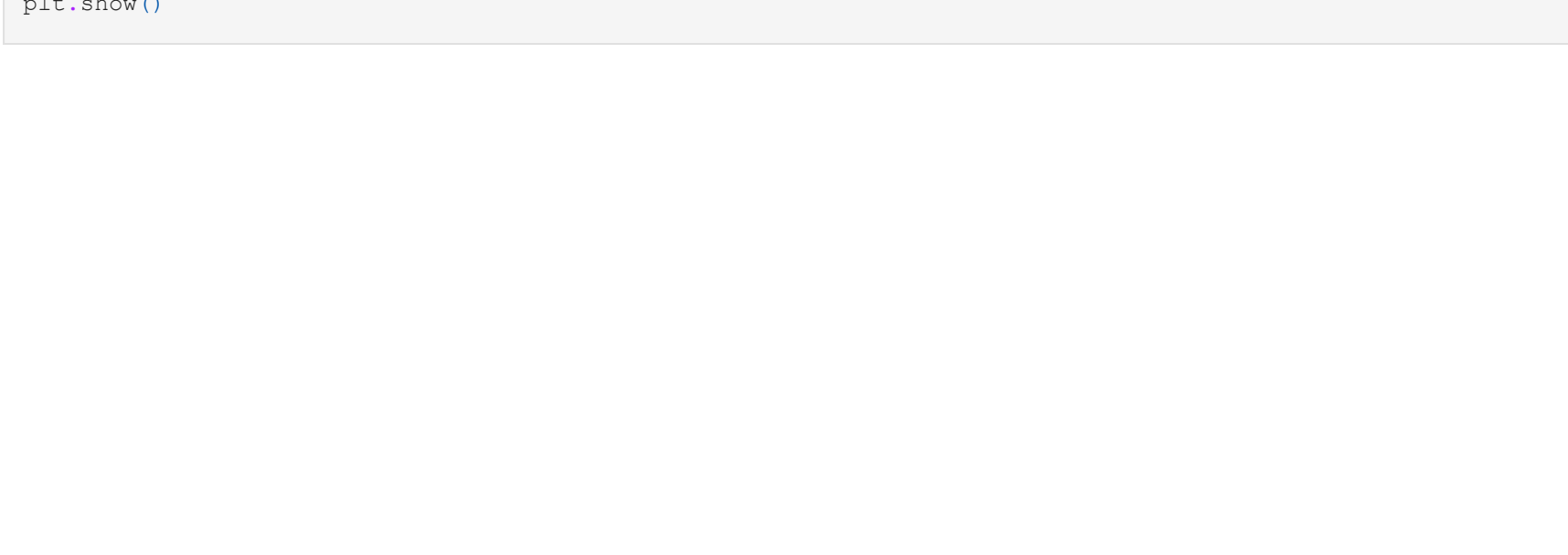
Based on above scatter plot we don't see any correlation between attendance and day and month.

```
In [22]: # pair plot to see pair wise correlations

sns.set(style='white')

sns.pairplot(data=dodgers,height=5, aspect=0.8)

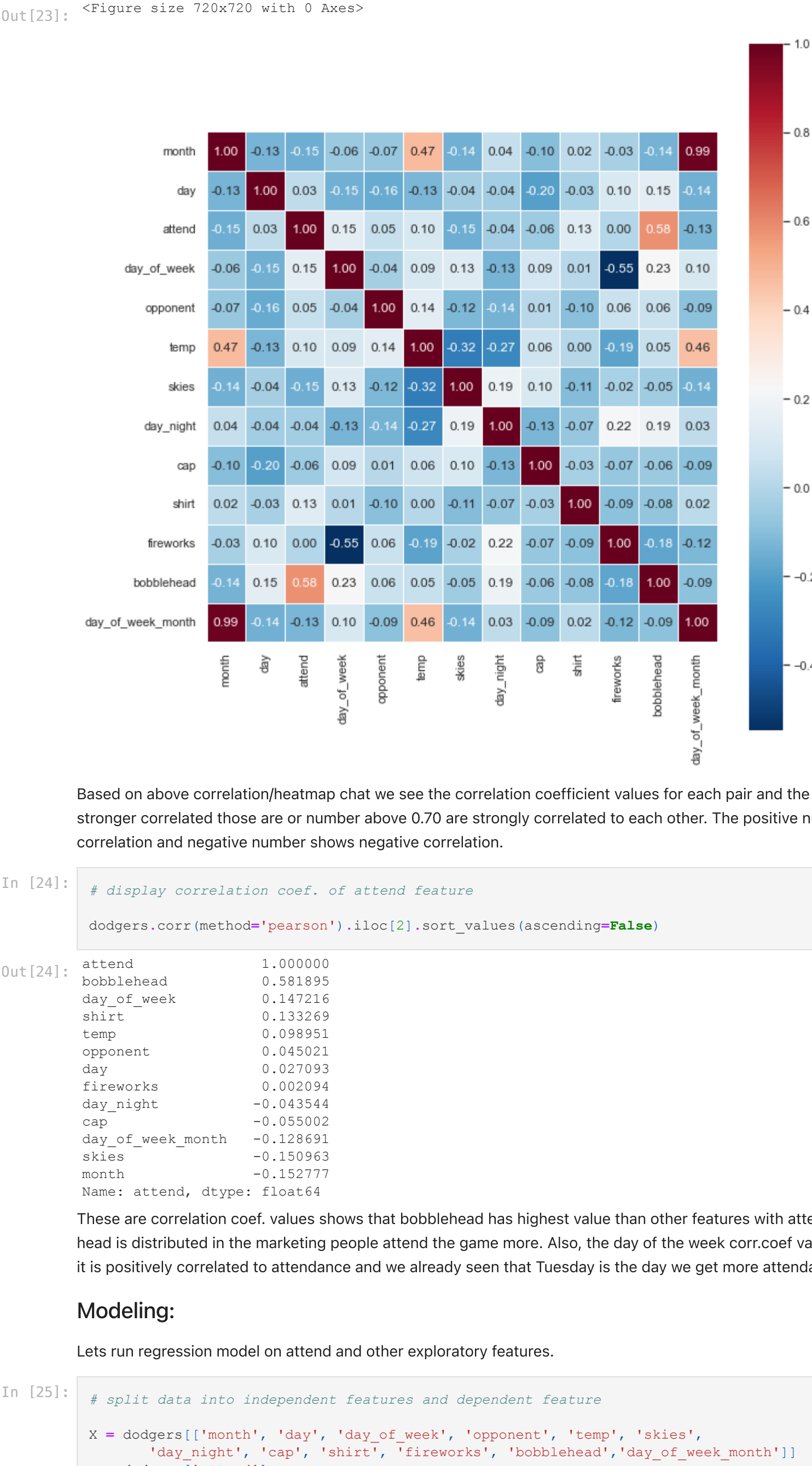
plt.title("Pair Plot");
```



```
In [23]:

plt.figure(figsize=(10,10))

sns.heatmap(dodgers.get_numeric_data().astype(float).corr(),
            square=True, cmap='RdBu_r', linewidth=5,
            annot=True, fmt='.2f').figure.tight_layout()
```

Based on above correlation/heatmap chat we see the correlation coefficient values for each pair and the dark red is color more stronger correlated those are or number above 0.70 are strongly correlated to each other. The positive number shows positive correlation and negative number shows negative correlation.

```
In [24]: # display correlation coef. of attend feature
dodgers.corr(method="pearson").iloc[2].sort_values(ascending=False)
```

```
Out[24]:
attend          1.000000
bobblehead      0.581895
day_of_week     0.147216
shirt           0.133269
temp            0.098951
opponent        0.045021
day             0.027093
fireworks       0.002094
day_night       -0.043544
cap             -0.055002
day_of_week_month 0.128691
skies           -0.150963
month           -0.152777
Name: attend, dtype: float64
```

These are correlation coef. values shows that bobblehead has highest value than other features with attend. That means if bobble head is distributed in the marketing people attend the game more. Also, the day of the week corr.coef value is 0.15 which positive so it is positively correlated to attendance and we already seen that Tuesday is the day we get more attendance.

Modeling:

Lets run regression model on attend and other exploratory features.

```
In [25]: # split data into independent features and dependent feature
X = dodgers[['month', 'day', 'day_of_week', 'opponent', 'temp', 'skies',
'day_night', 'cap', 'shirt', 'fireworks', 'bobblehead', 'day_of_week_month']]
y = dodgers['attend']
```

```
In [26]: # save feature names
feature_names = ['month', 'day', 'day_of_week', 'opponent', 'temp', 'skies',
'day_night', 'cap', 'shirt', 'fireworks', 'bobblehead', 'day_of_week_month']
```

```
In [27]: # transform the data to std scaler
X = sc.fit_transform(X)
y = sc.fit_transform(y.values.reshape(-1,1))
```

```
In [28]: X_features = pd.DataFrame(X,columns=feature_names) # to see features names in model summary
```

```
In [29]: # use statsmodel to run linear regression model
import statsmodels.api as sm
model = sm.OLS(y, sm.add_constant(X_features)).fit()
```

```
In [30]: # print the summary of model result
model.summary()
```

influencing the attendance. Alos, attendance is influenced by promotion

Afterall, to improve attendace it would be great idea to run marketing pro when temperature is best suitable for the game and people who prefer t

END

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

Regression Result:

Based on above regression result the value of the R-squared is 0.46, this means that the features we considered here are accounting for 45% variation in the attendance of the game and there are still 65% factors that may influence attendance of the game.

The features such as - day_of_week which has positive coefficient value means that this feature is positively correlated with attendance. Other features such as shirt, fireworks, bobblehead, temperature, month also shows positive relationship with attendance, that means these are the factors influence attendance of the game. Rest of the features which have negative coefficient values means those are negatively correlated with attendance.

When we consider the p-value, only the features- bobblehead & fireworks having p-value less than 0.05 which indicates that these features making significant contribution to the model.

Summary:

Overall, the analysis performed so far on the Dodgers games data we see that the **Tuesday** is the day we saw highest attendance in the game. June and August months Tuesdays we saw highest average attendance that means these months people prefer to attend the game.

Regression model supports the analysis we have done and suggest that day of week, months and temperatures are significantly influencing the attendance. Also, attendance is influenced by promotions such as bobble head, fireworks and shirts.

Afterall, to improve attendance it would be great idea to run marketing promotion on Tuesday of the summer months June, july, august when temperature is best suitable for the game and people who prefer to attend game and offer bobblehead or fireworks or shirts.

END