

```
In [ ]: import pandas as pd
import statsmodels.tsa.arima.model import ARIMA
from statsmodels.graphics import tsaplots
from statsmodels.tsa.arima import auto_arima
import matplotlib.pyplot as plt
import matplotlib.pyplot as pylab

# mape
fig, s1=tsaplots.metrics import mean_absolute_percentage_error

# dickey fuller
from statsmodels.tsa.stattools import adfuller

In [ ]: # change font properties for plots
# font = ('size', 12)
# plt.rc('font', "font")
params = {'legend.fontsize': 'small',
          'figure.figsize': (15, 5),
          'axes.labelsize': 'medium',
          'axes.titlesize': 'x-large',
          'xtick.labelsize': 'medium',
          'ytick.labelsize': 'medium'}

pylab.rcParams.update(params)

In [ ]: df = pd.read_csv('seasonSummaries.csv')

In [ ]: df.index = pd.period_range('1903', '2022', freq='Y')
df = df[['Rate']]

In [ ]: df.info()
<class 'pandas.core.frame.DataFrame'>
PeriodIndex: 120 entries, 1903 to 2022
Freq: A-BE
Data columns (total 1 columns):
# column Non-Null Count Dtype
...
0 Rate 120 non-null float64
dtypes: float64(1)
memory usage: 1.9 KB

In [ ]: # plot strikeout rate
plt.clf()
fig = plt.figure(figsize=(8,5))
ax = fig.add_subplot(111)
plt.plot(np.arange(1903,2023,1),df['Rate'],c='#D58032',linewidth=2)
plt.xlabel('Season')
plt.ylabel('Ks per 100 AB')
plt.title('MLB Strikeout Rate')
ax.text(y=97,x=5,s=1903 - 2022,transform=ax.transAxes, ha='center')
ax.spines['top'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.spines['left'].set_visible(False)
ax.spines['bottom'].set_visible(False)
ax.tick_params(bottom=True, top=False, left=True, right=False)
plt.xticks(ticks=np.arange(5,26,5))
plt.xticks(np.arange(1900,2021,20))
plt.xticks(rotation=45)
plt.show()

<Figure size 1080x360 with 0 Axes>
```

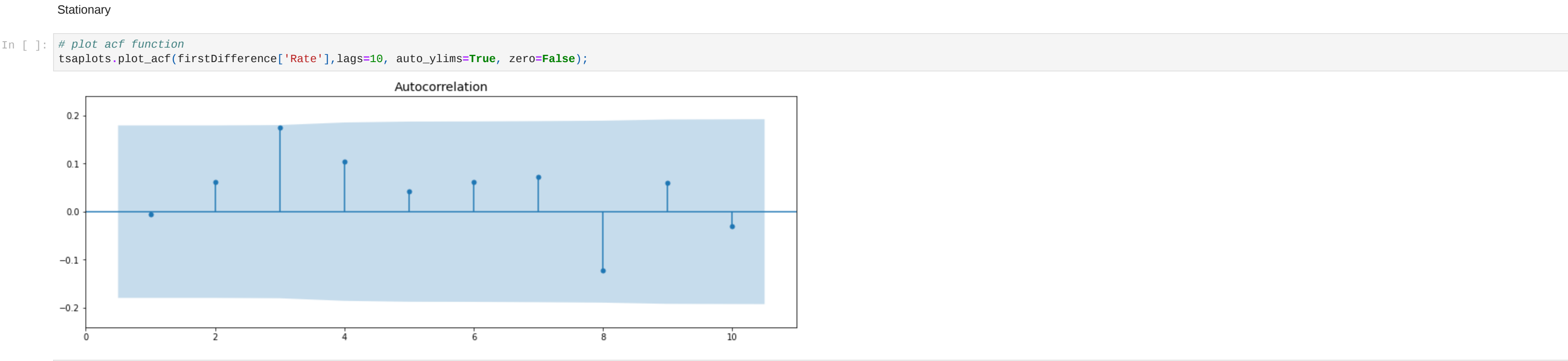


```
In [ ]: # check if process is stationary
dickey = adfuller(df['Rate'])
print(f"Test Statistic: {dickey[0]}VNP-Value: {dickey[1]}")
Test Statistic: 0.803940105582916
P-Value: 0.991695682284964
not stationary

In [ ]: firstDifference = pd.DataFrame({'Rate':df['Rate'].diff()})
firstDifference = firstDifference.iloc[1:]

In [ ]: # plot strikeout rate
plt.clf()
fig = plt.figure(figsize=(8,5))
ax = fig.add_subplot(111)
plt.plot(np.arange(1904,2023,1),firstDifference['Rate'],c='#D58032',linewidth=2)
plt.xlabel('Season')
plt.ylabel('Ks per 100 AB')
plt.title('MLB Strikeout Rate (First Difference)')
ax.text(y=97,x=5,s=1904 to 2022,transform=ax.transAxes, ha='center')
ax.spines['top'].set_visible(False)
ax.spines['right'].set_visible(False)
ax.spines['left'].set_visible(False)
ax.spines['bottom'].set_visible(False)
ax.tick_params(bottom=True, top=False, left=True, right=False)
plt.xticks(ticks=np.arange(2,25,1))
plt.xticks(np.arange(1900,2021,20))
plt.xticks(rotation=45)
plt.show()

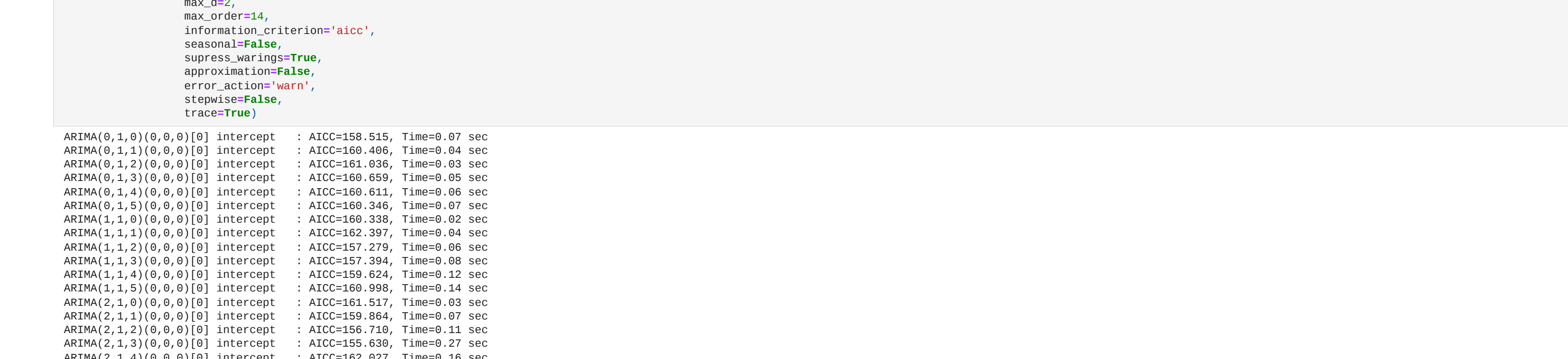
<Figure size 1080x360 with 0 Axes>
```



```
In [ ]: # check if process is stationary
dickey = adfuller(firstDifference['Rate'])
print(f"Test Statistic: {dickey[0]}VNP-Value: {dickey[1]}")
Test Statistic: -4.72795668289665
P-Value: 7.46611910757497e-05
Stationary

In [ ]: # plot acf function
tsaplots.plot_acf(firstDifference['Rate'],lags=10, auto_yaxis=True, zero=False);

Autocorrelation
```



```
In [ ]: # plot pacf function
tsaplots.plot_pacf(firstDifference['Rate'], method='ymw',lags=10, auto_yaxis=True, zero=False);

Partial Autocorrelation
```



```
In [ ]: # split into a test and train dataset
train = df.iloc[:int(len(df)*.7),:]
test = df.iloc[int(len(df)*.7):,:]

In [ ]: auto = auto_arima(train,
                      max_p=5,
                      max_d=5,
                      max_order=14,
                      information_criterion='aicc',
                      seasonal=False,
                      suppress_warnings=True,
                      approximation=False,
                      error_action='warn',
                      stepsize=False,
                      trace=True)

ARIMA(0,1,0)(0,0,0)[0] intercept : AIC=158.515, Time=0.07 sec
ARIMA(0,1,1)(0,0,0)[0] intercept : AIC=159.486, Time=0.04 sec
ARIMA(0,1,2)(0,0,0)[0] intercept : AIC=161.836, Time=0.03 sec
ARIMA(0,1,3)(0,0,0)[0] intercept : AIC=166.659, Time=0.05 sec
ARIMA(1,1,4)(0,0,0)[0] intercept : AIC=166.611, Time=0.06 sec
ARIMA(0,1,5)(0,0,0)[0] intercept : AIC=168.346, Time=0.07 sec
ARIMA(1,1,6)(0,0,0)[0] intercept : AIC=169.336, Time=0.02 sec
ARIMA(1,1,1)(0,0,0)[0] intercept : AIC=162.397, Time=0.04 sec
ARIMA(1,1,2)(0,0,0)[0] intercept : AIC=157.279, Time=0.06 sec
ARIMA(1,1,3)(0,0,0)[0] intercept : AIC=157.284, Time=0.08 sec
ARIMA(1,1,4)(0,0,0)[0] intercept : AIC=159.624, Time=0.12 sec
ARIMA(1,1,5)(0,0,0)[0] intercept : AIC=160.986, Time=0.14 sec
ARIMA(2,1,6)(0,0,0)[0] intercept : AIC=161.517, Time=0.03 sec
ARIMA(1,1,7)(0,0,0)[0] intercept : AIC=159.864, Time=0.07 sec
ARIMA(1,2)(0,0,0)[0] intercept : AIC=156.719, Time=0.11 sec
ARIMA(2,1,3)(0,0,0)[0] intercept : AIC=155.630, Time=0.27 sec
ARIMA(2,1,4)(0,0,0)[0] intercept : AIC=162.627, Time=0.18 sec
ARIMA(2,1,5)(0,0,0)[0] intercept : AIC=163.510, Time=0.28 sec
ARIMA(3,1,6)(0,0,0)[0] intercept : AIC=159.428, Time=0.04 sec
ARIMA(3,1,1)(0,0,0)[0] intercept : AIC=159.443, Time=0.05 sec
ARIMA(3,1,2)(0,0,0)[0] intercept : AIC=inf, Time=0.22 sec
ARIMA(3,1,3)(0,0,0)[0] intercept : AIC=inf, Time=0.28 sec
ARIMA(3,1,4)(0,0,0)[0] intercept : AIC=inf, Time=0.33 sec
ARIMA(4,1,5)(0,0,0)[0] intercept : AIC=inf, Time=0.43 sec
ARIMA(4,1,6)(0,0,0)[0] intercept : AIC=160.443, Time=0.05 sec
ARIMA(4,1,1)(0,0,0)[0] intercept : AIC=162.171, Time=0.19 sec
ARIMA(4,1,2)(0,0,0)[0] intercept : AIC=157.812, Time=0.28 sec
ARIMA(4,1,3)(0,0,0)[0] intercept : AIC=inf, Time=0.32 sec
ARIMA(4,1,4)(0,0,0)[0] intercept : AIC=inf, Time=0.38 sec
ARIMA(4,1,5)(0,0,0)[0] intercept : AIC=inf, Time=0.42 sec
ARIMA(5,1,6)(0,0,0)[0] intercept : AIC=162.218, Time=0.07 sec
ARIMA(5,1,1)(0,0,0)[0] intercept : AIC=164.383, Time=0.11 sec
ARIMA(5,1,2)(0,0,0)[0] intercept : AIC=166.622, Time=0.26 sec
ARIMA(5,1,3)(0,0,0)[0] intercept : AIC=inf, Time=0.35 sec
ARIMA(5,1,4)(0,0,0)[0] intercept : AIC=inf, Time=0.39 sec
ARIMA(5,1,5)(0,0,0)[0] intercept : AIC=inf, Time=0.53 sec

Best model: ARIMA(2,1,3)(0,0,0)[0] intercept
Total fit time: 6.363 seconds
```

```
In [ ]: model = ARIMA(train,order=auto.get_params()['order'])
results = model.fit(method_kwargs={"warn_convergence": False})

In [ ]: results.summary()
```

SARIMAX Results				
Dep. Variable:	Rate	No. Observations:	84	
Model:	ARIMA(2,1,3)	Log Likelihood:	-70.196	
Date:	Thu, 20 Oct 2022	AIC	152.391	
Time:	13:59:10	BIC	166.904	
Sample:	12-31-1903	HQIC	158.222	
	12-31-1986			
Covariance Type:	opg			
	coef	std err	z	P> z [0.025 0.975]
ar.L1	1.8384	0.083	19.836	0.000 1.657 2.020
ar.L2	-0.9202	0.093	-9.908	0.000 -1.102 0.738
ma.L1	-2.0862	0.489	-4.267	0.000 -3.045 -1.128
ma.L2	1.4404	0.641	2.246	0.025 0.184 2.697
ma.L3	-0.2406	0.183	-1.313	0.189 -0.600 0.118
sigma2	0.2999	0.172	1.748	0.080 -0.036 0.636
Ljung-Box (L1) [Q]:	0.01	Jarque-Bera (JB):	7.68	
	Prob(Q):	0.94	Prob(B):	0.02
Heteroskedasticity (H):	0.93	Skew:	0.18	
Prob(H) (two-sided):	0.84	Kurtosis:	4.48	

```
Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).
```

```
In [ ]: nAhead = len(test)
forecast = results.get_forecast(steps=nAhead)
alphaVal = .1
forecastFrame = forecast.summary_frame(alpha=alphaVal)
forecastFrame.columns = ['mean','se','lower','upper']
forecastFrame.head()

Out [ ]:
mean se lower upper
1987 17.370278 0.552176 16.462030 18.278527
1988 17.755143 0.695907 16.620477 18.899808
1989 18.101209 0.836009 16.726157 19.476381
1990 18.383414 0.986717 16.760409 20.006419
1991 18.935987 1.154534 16.684548 20.462625
```

```
In [ ]: forecastFrame.tail()

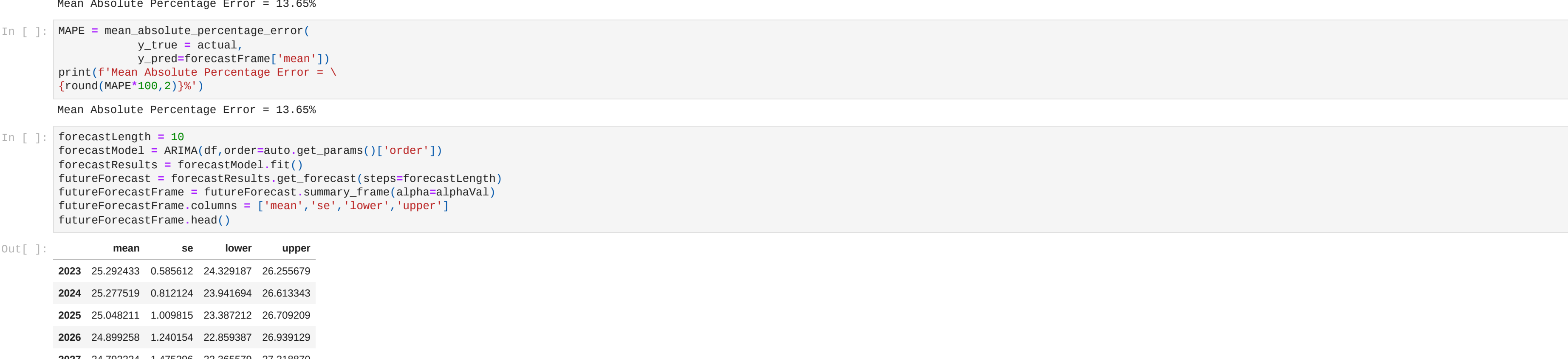
Out [ ]:
mean se lower upper
2019 17.789534 4.364382 10.601765 24.959304
2020 17.678707 4.446265 10.365252 24.992162
2021 17.583216 4.525319 10.139728 25.026704
2022 17.503373 4.600584 9.934085 25.068861
2023 17.438789 4.671483 9.754883 25.122696
```

```
In [ ]: trainYears = np.arange(1903,1987,1)
testYears = np.arange(1986,2023,1)
totalYears = np.arange(1903,2023,1)

In [ ]: # plot forecast over the test data set with all other seasons plotted too
plt.clf()
fig = plt.figure(figsize=(10,5))
ax = fig.add_subplot(111)
plt.plot(trainYears,train['Rate'],label='Observed', color='#002072')
plt.plot(testYears,test['Rate'][0:nAhead], color='#002072',label='Observed')
plt.plot(testYears,forecastFrame['mean'],color='#D58032',label='Forecast')
plt.fill_between(testYears,forecastFrame['lower'],forecastFrame['upper'],color='#FFB364',label=f'((1.0-alphaVal)*100,.07% Confidence Interval')
plt.xlabel('Season')
plt.ylabel('Ks Per 100 AB')

plt.xticks(ticks=np.arange(5,26,5))
plt.xticks(rotation = 45)
plt.legend(loc='upper left')
plt.title('MLB Strikeout Rate',loc='center')
ax.text(y=97,x=5,s='Forecast vs Observed',transform=ax.transAxes, ha='center')
ax.text(y=97,x=5,s='Forecast vs Observed',transform=ax.transAxes, ha='center')
ax.spines['top'].set_visible(False)
ax.spines['right'].set_visible(False)
plt.show()

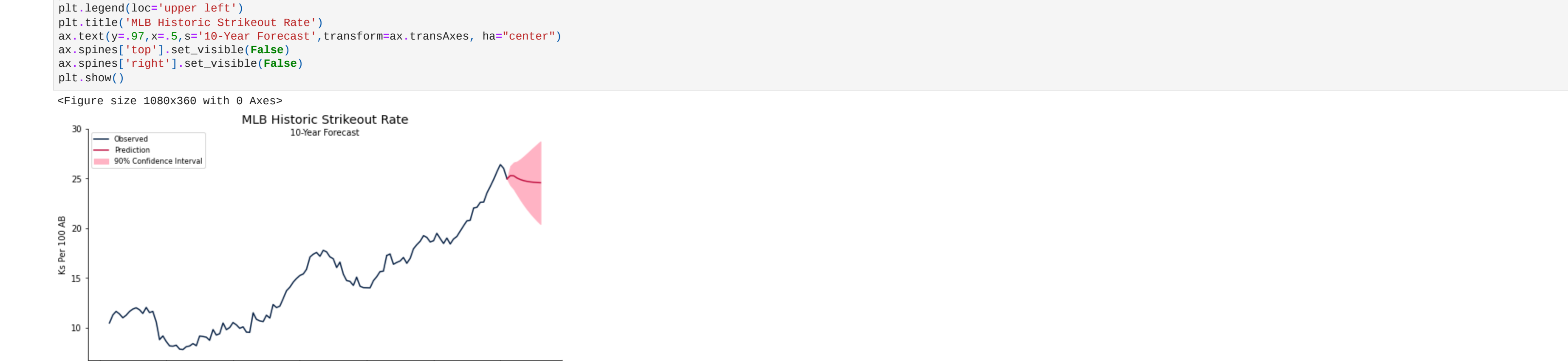
<Figure size 1080x360 with 0 Axes>
```



```
In [ ]: # plot forecast over the test data set with all other seasons plotted too
plt.clf()
fig = plt.figure(figsize=(10,5))
ax = fig.add_subplot(111)
plt.plot(trainYears,train['Rate'][0:nAhead], color='#002072',label='Observed')
plt.plot(testYears,test['Rate'][0:nAhead], color='#002072',label='Observed')
plt.plot(futureYears,futureForecastFrame['mean'],color='#D58032',label='Forecast')
plt.fill_between(testYears,futureForecastFrame['lower'],futureForecastFrame['upper'],color='#FFB364',label=f'((1.0-alphaVal)*100,.07% Confidence Interval')
plt.xlabel('Season')
plt.ylabel('Ks Per 100 AB')

plt.xticks(ticks=np.arange(10,31,5))
plt.xticks(rotation = 45)
plt.legend(loc='upper left')
plt.title('MLB Strikeout Rate')
ax.text(y=95,x=5,s='10-Year Forecast',transform=ax.transAxes, ha='center')
ax.spines['top'].set_visible(False)
ax.spines['right'].set_visible(False)
plt.show()

<Figure size 1080x360 with 0 Axes>
```



```
In [ ]: # plot forecast over the test data set with all other seasons plotted too
plt.clf()
fig = plt.figure(figsize=(10,5))
ax = fig.add_subplot(111)
plt.plot(trainYears,train['Rate'],label='Observed', color='#041e42')
plt.plot(testYears,test['Rate'][0:nAhead],label='Observed', color='#041e42')
plt.plot(futureYears,futureForecastFrame['mean'],color='#D58032',label='Forecast')
plt.fill_between(testYears,futureForecastFrame['lower'],futureForecastFrame['upper'],color='#FFB364',label=f'((1.0-alphaVal)*100,.07% Confidence Interval')
plt.xlabel('Season')
plt.ylabel('Ks Per 100 AB')

plt.xticks(ticks=np.arange(20,31,2))
plt.xticks(ticks=np.arange(2022,2023,2))
plt.xticks(rotation = 45)
plt.legend(loc='upper left')
plt.title('MLB Strikeout Rate')
ax.text(y=95,x=5,s='10-Year Forecast',transform=ax.transAxes, ha='center')
ax.spines['top'].set_visible(False)
ax.spines['right'].set_visible(False)
plt.show()

<Figure size 1080x360 with 0 Axes>
```

