# How to Create an ARIMA Model for Time Series Forecasting in Python

by **Jason Brownlee** on January 9, 2017 in **Time Series**

Last Updated on December 10, 2020

A popular and widely used statistical method for time series forecasting is the ARIMA model.

ARIMA is an acronym that stands for AutoRegressive Integrated Moving Average. It is a class of model that captures a suite of different standard temporal structures in time series data.

In this tutorial, you will discover how to **develop an ARIMA model for time series** forecasting in Python.

After completing this tutorial, you will know:

- About the ARIMA model the parameters used and assumptions made by the model.
- How to fit an ARIMA model to data and use it to make forecasts.
- How to configure the ARIMA model on your time series problem.

**Kick-start your project** with my new book Time Series Forecasting With Python, including *step-by-step tutorials* and the *Python source code* files for all **examples**.

Let's get started.

- **Updated Apr/2019**: Updated the link to dataset.
- **Updated Sep/2019**: Updated examples to use latest API.
- **Updated Dec/2020**: Updated examples to use latest API.

## Autoregressive Integrated Moving Average Model

An ARIMA model is a class of statistical models for analyzing and forecasting time series data.

It explicitly caters to a suite of standard structures in time series data, and as such provides a simple yet powerful method for making skillful time series forecasts.

ARIMA is an acronym that stands for AutoRegressive Integrated Moving Average. It is a generalization of the simpler AutoRegressive Moving Average and adds the notion of integration.

This acronym is descriptive, capturing the key aspects of the model itself. Briefly, they are:

- **AR**: *Autoregression*. A model that uses the dependent relationship between an observation and some number of lagged observations.
- **I**: *Integrated*. The use of differencing of raw observations (e.g. subtracting an observation from an observation at the previous time step) in order to make the time series stationary.
- **MA**: *Moving Average*. A model that uses the dependency between an observation and a residual error from a moving average model applied to lagged observations.

Each of these components are explicitly specified in the model as a parameter. A standard notation is used of ARIMA(p,d,q) where the parameters are substituted with integer values to quickly indicate the specific ARIMA model being used.

The parameters of the ARIMA model are defined as follows:

- **p**: The number of lag observations included in the model, also called the lag order.
- **d**: The number of times that the raw observations are differenced, also called the degree of differencing.
- **q**: The size of the moving average window, also called the order of moving average.

A linear regression model is constructed including the specified number and type of terms, and the data is prepared by a degree of differencing in order to make it stationary, i.e. to remove trend and seasonal structures that negatively affect the regression model.

A value of 0 can be used for a parameter, which indicates to not use that element of the model. This way, the ARIMA model can be configured to perform the function of an ARMA model, and even a simple AR, I, or MA model.

Adopting an ARIMA model for a time series assumes that the underlying process that generated the observations is an ARIMA process. This may seem obvious, but helps to motivate the need to confirm the assumptions of the model in the raw observations and in the residual errors of forecasts from the model.

Next, let's take a look at how we can use the ARIMA model in Python. We will start with loading a simple univariate time series.

---

---

## Shampoo Sales Dataset

This dataset describes the monthly number of sales of shampoo over a 3 year period.

The units are a sales count and there are 36 observations. The original dataset is credited to Makridakis, Wheelwright, and Hyndman (1998).

- Download the dataset.

Download the dataset and place it in your current working directory with the filename "*shampoo-sales.csv*".

Below is an example of loading the Shampoo Sales dataset with Pandas with a custom function to parse the date-time field. The dataset is baselined in an arbitrary year, in this case 1900.
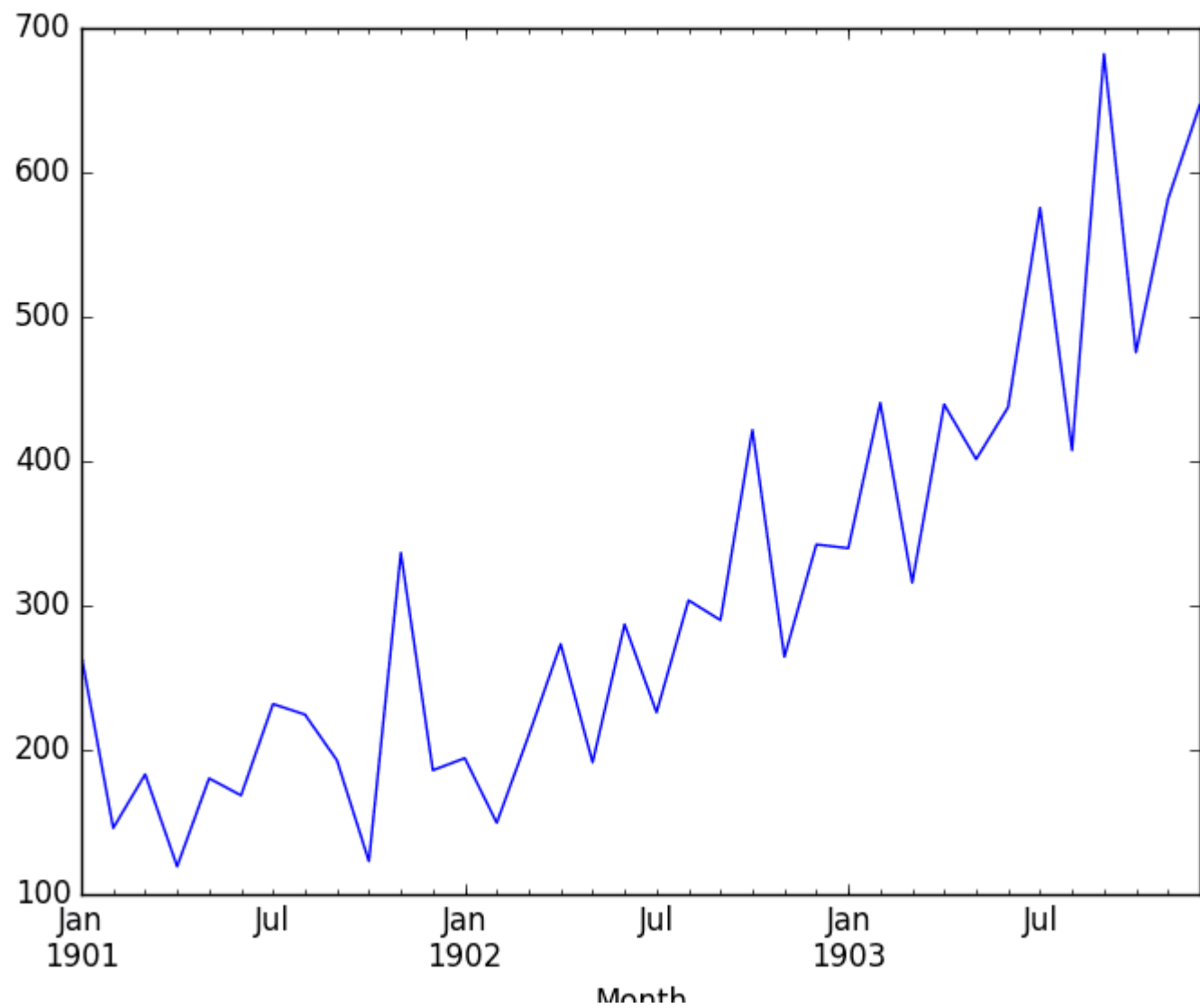
```
1  from pandas import read_csv
2  from pandas import datetime
3  from matplotlib import pyplot
4
5  def parser(x):
6      return datetime.strptime('190'+x, '%Y-%m')
7
8  series = read_csv('shampoo-sales.csv', header=0, parse_dates=[0], index_col=0, squeeze=True, date_parser=parser)
```

```
 9  print(series.head())
10  series.plot()
11  pyplot.show()
```

Running the example prints the first 5 rows of the dataset.

```
1  Month
2  1901-01-01 266.0
3  1901-02-01 145.9
4  1901-03-01 183.1
5  1901-04-01 119.3
6  1901-05-01 180.3
7  Name: Sales, dtype: float64
```

The data is also plotted as a time series with the month along the x-axis and sales figures on the y-axis.

Shampoo Sales Dataset Plot

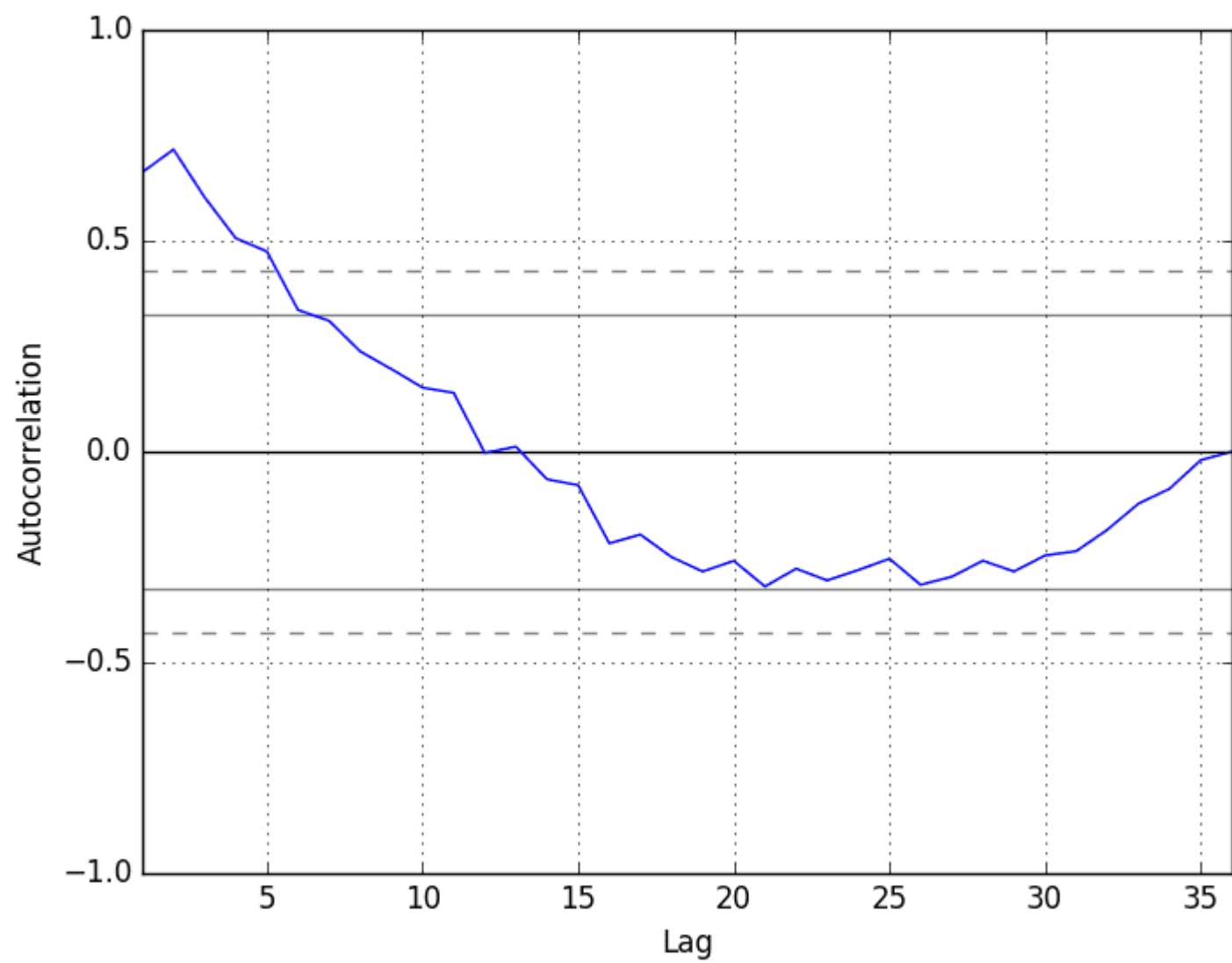We can see that the Shampoo Sales dataset has a clear trend.

This suggests that the time series is not stationary and will require differencing to make it stationary, at least a difference order of 1.

Let's also take a quick look at an autocorrelation plot of the time series. This is also built-in to Pandas. The example below plots the autocorrelation for a large number of lags in the time series.

```
1   from pandas import read_csv
2   from pandas import datetime
3   from matplotlib import pyplot
4   from pandas.plotting import autocorrelation_plot
5
6   def parser(x):
7       return datetime.strptime('190'+x, '%Y-%m')
8
9   series = read_csv('shampoo-sales.csv', header=0, parse_dates=[0], index_col=0, squeeze=True, date_parser=parser)
10  autocorrelation_plot(series)
11  pyplot.show()
```

Running the example, we can see that there is a positive correlation with the first 10-to-12 lags that is perhaps significant for the first 5 lags.

A good starting point for the AR parameter of the model may be 5.

Autocorrelation Plot of Shampoo Sales Data

# ARIMA with Python

The statsmodels library provides the capability to fit an ARIMA model.

An ARIMA model can be created using the statsmodels library as follows:

1. Define the model by calling ARIMA() and passing in the *p*, *d*, and *q* parameters.
2. The model is prepared on the training data by calling the fit() function.
3. Predictions can be made by calling the predict() function and specifying the index of the time or times to be predicted.

Let's start off with something simple. We will fit an ARIMA model to the entire Shampoo Sales dataset and review the residual errors.

First, we fit an ARIMA(5,1,0) model. This sets the lag value to 5 for autoregression, uses a difference order of 1 to make the time series stationary, and uses a moving average model of 0.

```
1  # fit an ARIMA model and plot residual errors
2  from pandas import datetime
3  from pandas import read_csv
4  from pandas import DataFrame
5  from statsmodels.tsa.arima.model import ARIMA
6  from matplotlib import pyplot
7  # load dataset
8  def parser(x):
9      return datetime.strptime('190'+x, '%Y-%m')
10 series = read_csv('shampoo-sales.csv', header=0, index_col=0, parse_dates=True, squeeze=True, date_parser=parser)
11 series.index = series.index.to_period('M')
12 # fit model
13 model = ARIMA(series, order=(5,1,0))
14 model_fit = model.fit()
15 # summary of fit model
16 print(model_fit.summary())
17 # line plot of residuals
18 residuals = DataFrame(model_fit.resid)
```
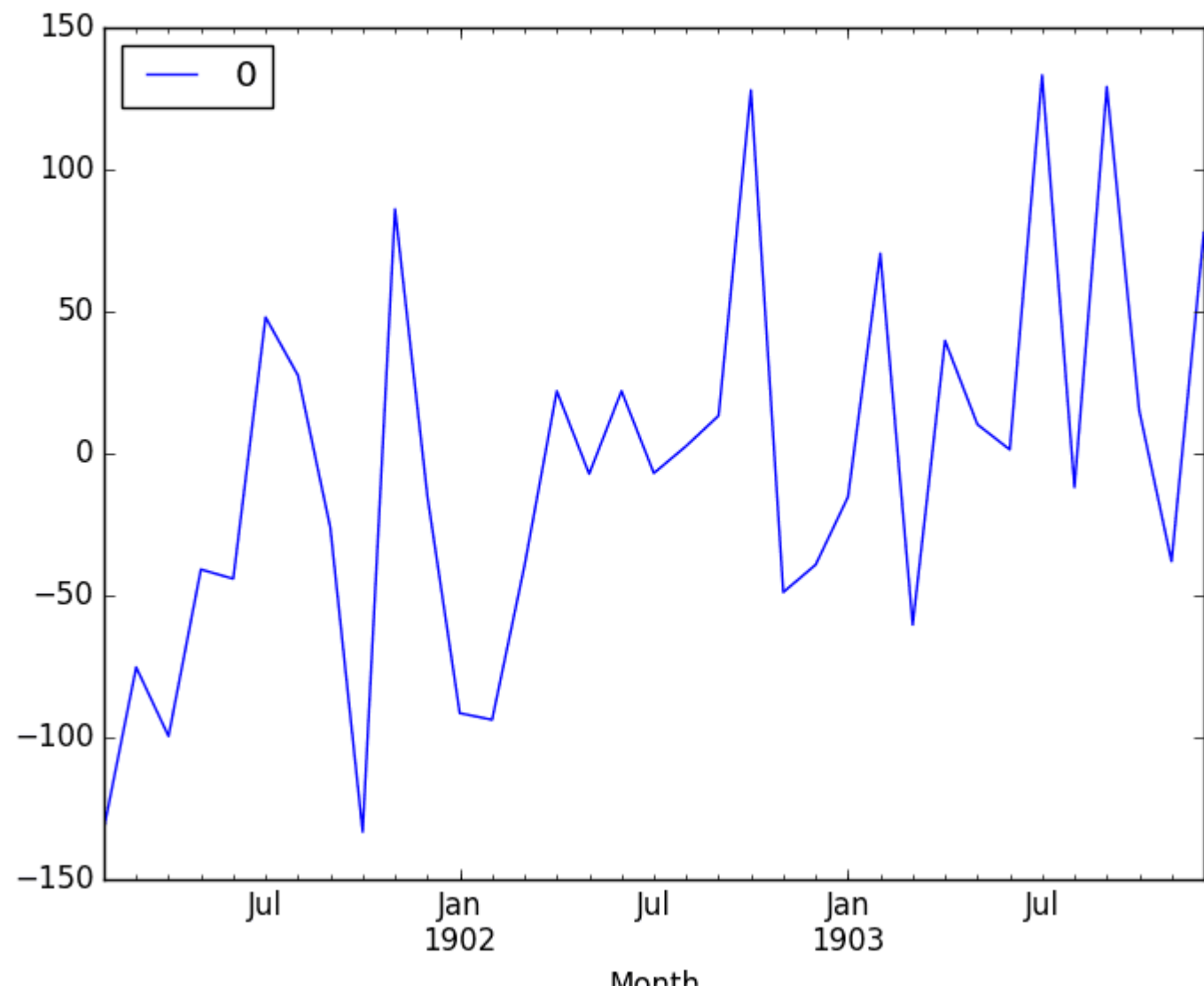
```
19  residuals.plot()
20  pyplot.show()
21  # density plot of residuals
22  residuals.plot(kind='kde')
23  pyplot.show()
24  # summary stats of residuals
25  print(residuals.describe())
```

Running the example prints a summary of the fit model. This summarizes the coefficient values used as well as the skill of the fit on the on the in-sample observations.

```
1                              SARIMAX Results
2   ==============================================================================
3   Dep. Variable:                    Sales   No. Observations:                 36
4   Model:                   ARIMA(5, 1, 0)   Log Likelihood              -198.485
5   Date:                  Thu, 10 Dec 2020   AIC                          408.969
6   Time:                          09:15:01   BIC                          418.301
7   Sample:                      01-31-1901   HQIC                         412.191
8                              - 12-31-1903
9   Covariance Type:                    opg
10  ==============================================================================
11                   coef    std err          z      P>|z|      [0.025      0.975]
12  ------------------------------------------------------------------------------
13  ar.L1         -0.9014      0.247     -3.647      0.000      -1.386      -0.417
14  ar.L2         -0.2284      0.268     -0.851      0.395      -0.754       0.298
15  ar.L3          0.0747      0.291      0.256      0.798      -0.497       0.646
16  ar.L4          0.2519      0.340      0.742      0.458      -0.414       0.918
17  ar.L5          0.3344      0.210      1.593      0.111      -0.077       0.746
18  sigma2      4728.9608   1316.021      3.593      0.000    2149.607    7308.314
19  ===================================================================================
20  Ljung-Box (L1) (Q):                0.61   Jarque-Bera (JB):                 0.96
21  Prob(Q):                           0.44   Prob(JB):                         0.62
22  Heteroskedasticity (H):            1.07   Skew:                             0.28
23  Prob(H) (two-sided):               0.90   Kurtosis:                         2.41
24  ===================================================================================
```
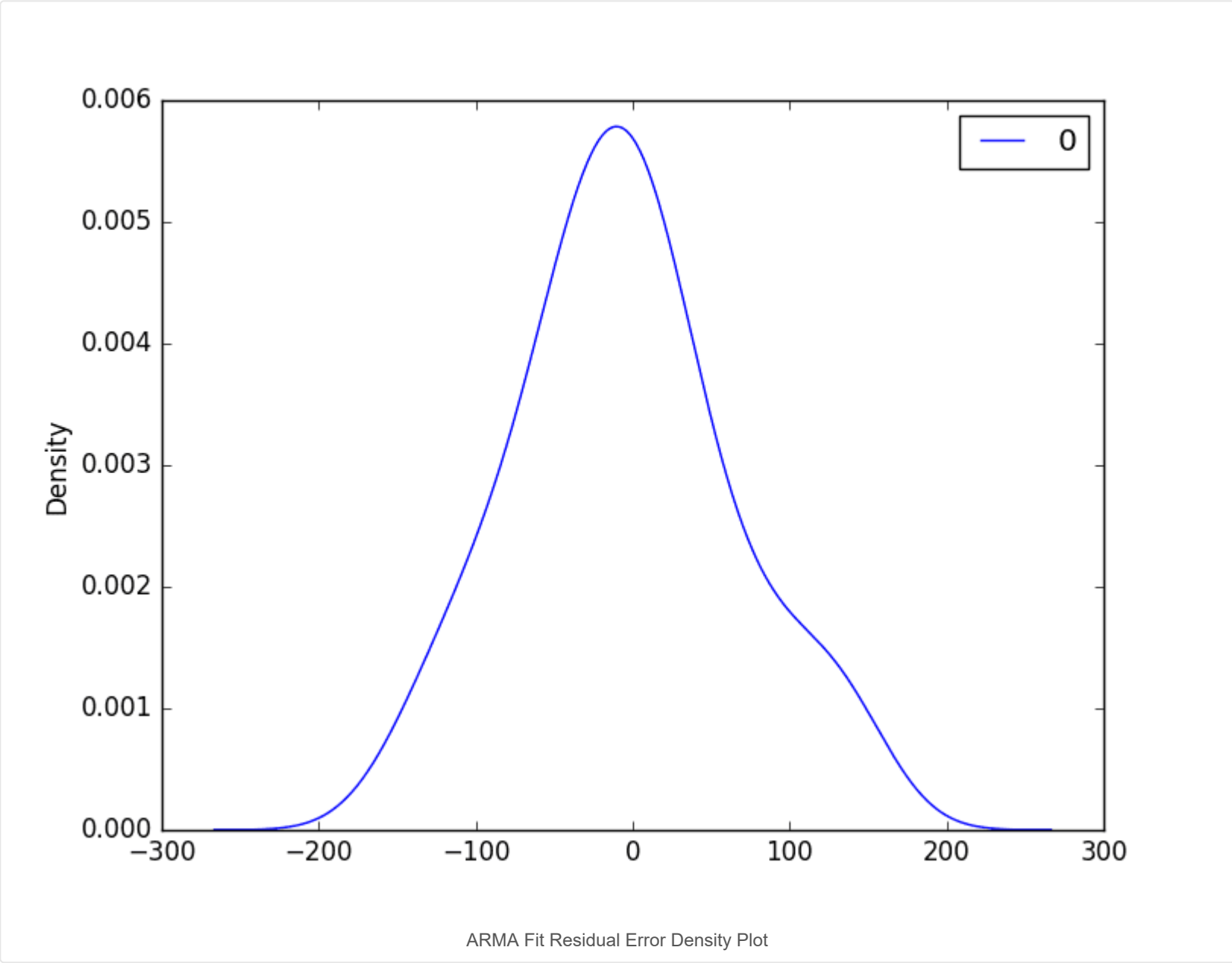
First, we get a line plot of the residual errors, suggesting that there may still be some trend information not captured by the model.

ARMA Fit Residual Error Line Plot

Next, we get a density plot of the residual error values, suggesting the errors are Gaussian, but may not be centered on zero.

ARMA Fit Residual Error Density Plot

The distribution of the residual errors is displayed. The results show that indeed there is a bias in the prediction (a non-zero mean in the residuals).

```
1  count    36.000000
2  mean     21.936144
3  std      80.774430
```

```
4  min    -122.292030
5  25%     -35.040859
6  50%      13.147219
7  75%      68.848286
8  max     266.000000
```

Note, that although above we used the entire dataset for time series analysis, ideally we would perform this analysis on just the training dataset when developing a predictive model.

Next, let's look at how we can use the ARIMA model to make forecasts.

## Rolling Forecast ARIMA Model

The ARIMA model can be used to forecast future time steps.

We can use the predict() function on the ARIMAResults object to make predictions. It accepts the index of the time steps to make predictions as arguments. These indexes are relative to the start of the training dataset used to make predictions.

If we used 100 observations in the training dataset to fit the model, then the index of the next time step for making a prediction would be specified to the prediction function as *start=101, end=101*. This would return an array with one element containing the prediction.

We also would prefer the forecasted values to be in the original scale, in case we performed any differencing (*d>0* when configuring the model). This can be specified by setting the *typ* argument to the value *'levels'*: *typ='levels'*.

Alternately, we can avoid all of these specifications by using the forecast() function, which performs a one-step forecast using the model.

We can split the training dataset into train and test sets, use the train set to fit the model, and generate a prediction for each element on the test set.

A rolling forecast is required given the dependence on observations in prior time steps for differencing and the AR model. A crude way to perform this rolling forecast is to re-create the ARIMA model after each new observation is received.

We manually keep track of all observations in a list called history that is seeded with the training data and to which new observations are appended each iteration.

Putting this all together, below is an example of a rolling forecast with the ARIMA model in Python.
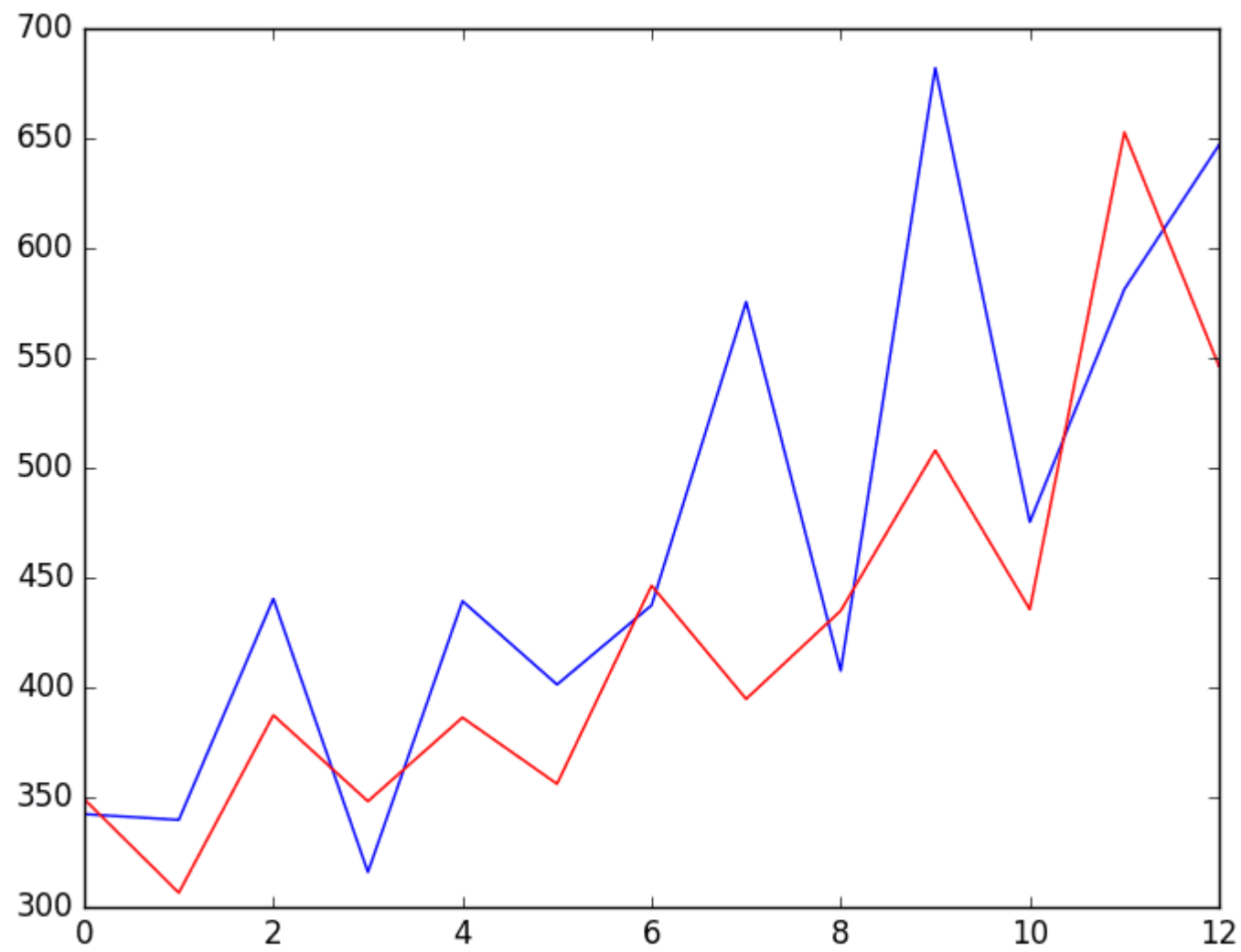
```python
# evaluate an ARIMA model using a walk-forward validation
from pandas import read_csv
from pandas import datetime
from matplotlib import pyplot
from statsmodels.tsa.arima.model import ARIMA
from sklearn.metrics import mean_squared_error
from math import sqrt
# load dataset
def parser(x):
    return datetime.strptime('190'+x, '%Y-%m')
series = read_csv('shampoo-sales.csv', header=0, index_col=0, parse_dates=True, squeeze=True, date_parser=parser)
series.index = series.index.to_period('M')
# split into train and test sets
X = series.values
size = int(len(X) * 0.66)
train, test = X[0:size], X[size:len(X)]
history = [x for x in train]
predictions = list()
# walk-forward validation
for t in range(len(test)):
    model = ARIMA(history, order=(5,1,0))
    model_fit = model.fit()
    output = model_fit.forecast()
    yhat = output[0]
    predictions.append(yhat)
    obs = test[t]
    history.append(obs)
    print('predicted=%f, expected=%f' % (yhat, obs))
# evaluate forecasts
rmse = sqrt(mean_squared_error(test, predictions))
print('Test RMSE: %.3f' % rmse)
# plot forecasts against actual outcomes
pyplot.plot(test)
pyplot.plot(predictions, color='red')
pyplot.show()
```

Running the example prints the prediction and expected value each iteration.

We can also calculate a final root mean squared error score (RMSE) for the predictions, providing a point of comparison for other ARIMA configurations.

```
1   predicted=343.272180, expected=342.300000
2   predicted=293.329674, expected=339.700000
3   predicted=368.668956, expected=440.400000
4   predicted=335.044741, expected=315.900000
5   predicted=363.220221, expected=439.300000
6   predicted=357.645324, expected=401.300000
7   predicted=443.047835, expected=437.400000
8   predicted=378.365674, expected=575.500000
9   predicted=459.415021, expected=407.600000
10  predicted=526.890876, expected=682.000000
11  predicted=457.231275, expected=475.300000
12  predicted=672.914944, expected=581.300000
13  predicted=531.541449, expected=646.900000
14  Test RMSE: 89.021
```

A line plot is created showing the expected values (blue) compared to the rolling forecast predictions (red). We can see the values show some trend and are in the correct scale.

ARIMA Rolling Forecast Line Plot

The model could use further tuning of the p, d, and maybe even the q parameters.

# Configuring an ARIMA Model

The classical approach for fitting an ARIMA model is to follow the Box-Jenkins Methodology.

This is a process that uses time series analysis and diagnostics to discover good parameters for the ARIMA model.

In summary, the steps of this process are as follows:

1. **Model Identification**. Use plots and summary statistics to identify trends, seasonality, and autoregression elements to get an idea of the amount of differencing and the size of the lag that will be required.
2. **Parameter Estimation**. Use a fitting procedure to find the coefficients of the regression model.
3. **Model Checking**. Use plots and statistical tests of the residual errors to determine the amount and type of temporal structure not captured by the model.

The process is repeated until either a desirable level of fit is achieved on the in-sample or out-of-sample observations (e.g. training or test datasets).

The process was described in the classic 1970 textbook on the topic titled Time Series Analysis: Forecasting and Control by George Box and Gwilym Jenkins. An updated 5th edition is now available if you are interested in going deeper into this type of model and methodology.

Given that the model can be fit efficiently on modest-sized time series datasets, grid searching parameters of the model can be a valuable approach.

For an example of how to grid search the hyperparameters of the ARIMA model, see the tutorial:

- How to Grid Search ARIMA Model Hyperparameters with Python

# Summary

In this tutorial, you discovered how to develop an ARIMA model for time series forecasting in Python.

Specifically, you learned:

- About the ARIMA model, how it can be configured, and assumptions made by the model.
- How to perform a quick time series analysis using the ARIMA model.
- How to use an ARIMA model to forecast out of sample predictions.

**Do you have any questions about ARIMA, or about this tutorial?**
Ask your questions in the comments below and I will do my best to answer.