

ARIMA Model – Complete Guide to Time Series Forecasting in Python

by [Selva Prabhakaran](https://www.machinelearningplus.com/author/selva86/) (https://www.machinelearningplus.com/author/selva86/)

[f](#) (/#facebook) [t](#) (/#twitter) [w](#) (/#whatsapp) [in](#) (/#linkedin) [r](#) (/#reddit)

[G](#) (/#google_bookmarks)

[+](https://www.addtoany.com/share?url=https%3A%2F%2Fwww.machinelearningplus.com%2Ftime-series%2Farima-model-time-series-forecasting-python%2F&title=ARIMA%20Model%20%E2%80%93%20Complete%20Guide%20to%20Time%20Series%20Forecasting%20in%20Python) (https://www.addtoany.com/share?url=https%3A%2F%2Fwww.machinelearningplus.com%2Ftime-series%2Farima-model-time-series-forecasting-python%2F&title=ARIMA%20Model%20%E2%80%93%20Complete%20Guide%20to%20Time%20Series%20Forecasting%20in%20Python)

Using ARIMA model, you can forecast a time series using the series past values. In this post, we build an optimal ARIMA model from scratch and extend it to Seasonal ARIMA (SARIMA) and SARIMAX models. You will also see how to build autoarima models in python



ARIMA Model – Time Series Forecasting. Photo by Cerquiera

Recommended

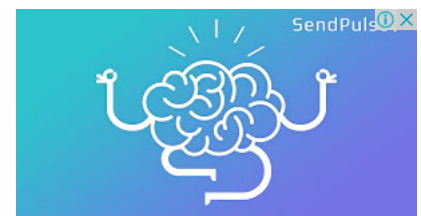
Time Series Analysis in Python – A Comprehensive Guide with Examples (https://www.machinelearningplus.com/time-series/time-series-analysis-python/)

Vector Autoregression (VAR) – Comprehensive Guide with Examples in Python (https://www.machinelearningplus.com/time-series/vector-autoregression-examples-python/)

ARIMA Model – Complete Guide to Time Series Forecasting in Python (https://www.machinelearningplus.com/time-series/arima-model-time-series-forecasting-python/)

Augmented Dickey Fuller Test (ADF Test) – Must Read Guide (https://www.machinelearningplus.com/time-series/augmented-dickey-fuller-test/)

KPSS Test for Stationarity (https://www.machinelearningplus.com/time-series/kpss-test-for-stationarity/)



Email Sending Service #1

OPEN

[ezoic](https://www.ezoic.com/what-is-ezoic/) (https://www.ezoic.com/what-is-ezoic/)

More Recommended

report this ad

Gradient Boosting – A Concise Introduction from Scratch (https://www.machinelearningplus.com/machine-learning/gradient-boosting/)

Bias Variance Tradeoff – Clearly Explained (https://www.machinelearningplus.com/machine-learning/bias-variance-tradeoff/)

Caret Package – A Practical Guide to Machine Learning in R (https://www.machinelearningplus.com/machine-learning/caret-package/)

Linear Regression in Julia



Contents

1. Introduction to Time Series Forecasting
2. Introduction to ARIMA Models
3. What does the p, d and q in ARIMA model mean?
4. What are AR and MA models
5. How to find the order of differencing (d) in ARIMA model
6. How to find the order of the AR term (p)
7. How to find the order of the MA term (q)
8. How to handle if a time series is slightly under or over differenced
9. How to build the ARIMA Model
10. How to do find the optimal ARIMA model manually using Out-of-Time Cross validation
11. Accuracy Metrics for Time Series Forecast
12. How to do Auto Arima Forecast in Python
13. How to interpret the residual plots in ARIMA model
14. How to automatically build SARIMA model in python
15. How to build SARIMAX Model with exogenous variable
16. Practice Exercises
17. Conclusion

1. Introduction to Time Series Forecasting

A time series is a sequence where a metric is recorded over regular time intervals.

Depending on the frequency, a time series can be of yearly (ex: annual budget), quarterly (ex: expenses), monthly (ex: air traffic), weekly (ex: sales qty), daily (ex: weather), hourly (ex: stocks price), minutes (ex: inbound calls in a call center) and even seconds wise (ex: web traffic).

We have already seen the steps involved in a previous post on [Time Series Analysis](https://www.machinelearningplus.com/time-series-analysis-python/) (<https://www.machinelearningplus.com/time-series-analysis-python/>). If you haven't read it, I highly encourage you to do so.

Forecasting is the next step where you want to predict the future values the series is going to take.

But why forecast?

Because, forecasting a time series (like demand and sales) is often of tremendous commercial value.

In most manufacturing companies, it drives the fundamental business planning, procurement and production activities. Any errors in the forecasts will ripple down throughout the supply chain or any business context for that matter. So it's important to get the forecasts accurate in order to save on costs and is critical to success.

Not just in manufacturing, the techniques and concepts behind time series forecasting are applicable in any business.

Now forecasting a time series can be broadly divided into two types.

If you use only the previous values of the time series to predict its future values, it is called **Univariate Time Series Forecasting**

And if you use predictors other than the series (a.k.a exogenous variables) to forecast it is called **Multi Variate Time Series Forecasting**

This post focuses on a particular type of forecasting method called **ARIMA** modeling.

(<https://www.machinelearningplus.com/linear-regression-in-julia/>)

ARIMA Model – Complete Guide to Time Series Forecasting in Python
(<https://www.machinelearningplus.com/time-series/arima-model-time-series-forecasting-python/>)

How Naive Bayes Algorithm Works? (with example and full code)
(<https://www.machinelearningplus.com/predictive-modeling/how-naive-bayes-algorithm-works-with-example-and-full-code/>)

Principal Component Analysis (PCA) – Better Explained
(<https://www.machinelearningplus.com/machine-learning/principal-components-analysis-pca-better-explained/>)

Mahalanobis Distance – Understanding the math with examples (python)
(<https://www.machinelearningplus.com/statistics/mahalanobis-distance/>)

Investor's Portfolio Optimization with Python using Practical Examples
(<https://www.machinelearningplus.com/machine-learning/portfolio-optimization-python-example/>)

Augmented Dickey Fuller Test (ADF Test) – Must Read Guide
(<https://www.machinelearningplus.com/time-series/augmented-dickey-fuller-test/>)

Complete Introduction to Linear Regression in R
(<https://www.machinelearningplus.com/machine-learning/complete-introduction-linear-regression-r/>)

Cosine Similarity – Understanding the math and how it works (with python codes)
(<https://www.machinelearningplus.com/nlp/cosine-similarity/>)

Feature Selection – Ten Effective Techniques with Examples
(<https://www.machinelearningplus.com/machine-learning/feature-selection/>)

Gensim Tutorial – A Complete Beginners Guide
(<https://www.machinelearningplus.com/nlp/gensim-tutorial/>)

K-Means Clustering Algorithm from Scratch
(<https://www.machinelearningplus.com/predictive-modeling/k-means-clustering/>)

KPSS Test for Stationarity
(<https://www.machinelearningplus.com/time-series/kpss-test-for-stationarity/>)

Lemmatization Approaches with Examples in Python
(<https://www.machinelearningplus.com/nlp/lemmatization-examples-python/>)

Python Numpy – Introduction to ndarray [Part 1]
(<https://www.machinelearningplus.com/python/numpy-tutorial-part1-array-python-examples/>)

Numpy Tutorial Part 2 – Vital Functions for Data Analysis
(<https://www.machinelearningplus.com/python/numpy-tutorial-python-part2/>)

ARIMA, short for 'AutoRegressive Integrated Moving Average', is a forecasting algorithm based on the idea that the information in the past values of the time series can alone be used to predict the future values.

2. Introduction to ARIMA Models

So what exactly is an ARIMA model?

ARIMA, short for 'Auto Regressive Integrated Moving Average' is actually a class of models that 'explains' a given time series based on its own past values, that is, its own lags and the lagged forecast errors, so that equation can be used to forecast future values.

Any 'non-seasonal' time series that exhibits patterns and is not a random white noise can be modeled with ARIMA models.

An ARIMA model is characterized by 3 terms: p , d , q

where,

p is the order of the AR term

q is the order of the MA term

d is the number of differencing required to make the time series stationary

If a time series, has seasonal patterns, then you need to add seasonal terms and it becomes SARIMA, short for 'Seasonal ARIMA'. More on that once we finish ARIMA.

So, what does the 'order of AR term' even mean? Before we go there, let's first look at the ' d ' term.

3. What does the p , d and q in ARIMA model mean

The first step to build an ARIMA model is to ~~make the time series stationary~~ (<https://www.machinelearningplus.com/stationary-time-series/>).

Why?

Because, term 'Auto Regressive' in ARIMA means it is a linear regression model (<https://www.machinelearningplus.com/machine-learning/complete-introduction-linear-regression-r/>) that uses its own lags as predictors. Linear regression models, as you know, work best when the predictors are not correlated and are independent of each other.

So how to make a series stationary?

The most common approach is to difference it. That is, subtract the previous value from the current value. Sometimes, depending on the complexity of the series, more than one differencing may be needed.

The value of d , therefore, is the minimum number of differencing needed to make the series stationary. And if the time series is already stationary, then $d = 0$.

P-Value – Understanding from Scratch
(<https://www.machinelearningplus.com/statistics/what-is-p-value/>)

Vector Autoregression (VAR) – Comprehensive Guide with Examples in Python
(<https://www.machinelearningplus.com/time-series/vector-autoregression-examples-python/>)

Time Series Analysis in Python – A Comprehensive Guide with Examples
(<https://www.machinelearningplus.com/time-series/time-series-analysis-python/>)

Top 15 Evaluation Metrics for Classification Models
(<https://www.machinelearningplus.com/machine-learning/evaluation-metrics-classification-models-r/>)

Topic modeling visualization – How to present the results of LDA models?
(<https://www.machinelearningplus.com/nlp/topic-modeling-visualization-how-to-present-results-lda-models/>)

Topic Modeling with Gensim (Python)
(<https://www.machinelearningplus.com/nlp/topic-modeling-gensim-python/>)

Building chatbot with Rasa and spaCy
(<https://www.machinelearningplus.com/nlp/chatbot-with-rasa-and-spacy/>)

Top Posts & Pages

ARIMA Model - Complete Guide to Time Series Forecasting in Python
(<https://www.machinelearningplus.com/time-series/arma-model-time-series-forecasting-python/>)

Parallel Processing in Python - A Practical Guide with Examples
(<https://www.machinelearningplus.com/python/parallel-processing-python/>)

Time Series Analysis in Python - A Comprehensive Guide with Examples
(<https://www.machinelearningplus.com/time-series/time-series-analysis-python/>)

Top 50 matplotlib Visualizations - The Master Plots (with full python code)
(<https://www.machinelearningplus.com/plots/top-50-matplotlib-visualizations-the-master-plots-python/>)

Topic Modeling with Gensim (Python)
(<https://www.machinelearningplus.com/nlp/topic-modeling-gensim-python/>)

Machine Learning Better Explained!
(<https://www.machinelearningplus.com/>)

Cosine Similarity - Understanding the math and how it works (with python codes)
(<https://www.machinelearningplus.com/nlp/cosine-similarity/>)

Matplotlib Histogram - How to Visualize Distributions in Python
(<https://www.machinelearningplus.com/plots/matplotlib-histogram-python-examples/>)

Next, what are the 'p' and 'q' terms?

'p' is the order of the 'Auto Regressive' (AR) term. It refers to the number of lags of Y to be used as predictors. And 'q' is the order of the 'Moving Average' (MA) term. It refers to the number of lagged forecast errors that should go into the ARIMA Model.

4. What are AR and MA models

So what are AR and MA models? what is the actual mathematical formula for the AR and MA models?

A pure **Auto Regressive (AR only) model** is one where Y_t depends only on its own lags. That is, Y_t is a function of the 'lags of Y_t '.

(<https://www.machinelearningplus.com/wp-content/uploads/2019/02/Equation-1-min.png>)

where, Y_{t-1} is the lag1 of the series, β_1 is the coefficient of lag1 that the model estimates and α is the intercept term, also estimated by the model.

Likewise a pure **Moving Average (MA only) model** is one where Y_t depends only on the lagged forecast errors.

(<https://www.machinelearningplus.com/wp-content/uploads/2019/02/Equation-2-min.png>)

where the error terms are the errors of the autoregressive models of the respective lags. The errors E_t and E_{t-1} are the errors from the following equations :

(<https://www.machinelearningplus.com/wp-content/uploads/2019/02/Equation-3-min.png>)

[Error: The beta coefficients in the second equation above is incorrect.]

That was AR and MA models respectively.

So what does the equation of an ARIMA model look like?

Python Logging - Simplest Guide with Full Code and Examples
(<https://www.machinelearningplus.com/python/python-logging-guide/>)

101 Pandas Exercises for Data Analysis
(<https://www.machinelearningplus.com/python/101-pandas-exercises-python/>)

Recent Posts

Matplotlib Plotting Tutorial – Complete overview of Matplotlib library
(<https://www.machinelearningplus.com/plots/matplotlib-plotting-tutorial/>)

How to implement Linear Regression in TensorFlow
(<https://www.machinelearningplus.com/deep-learning/linear-regression-tensorflow/>)

Brier Score – How to measure accuracy of probabilistic predictions
(<https://www.machinelearningplus.com/statistics/brier-score/>)

Modin – How to speedup pandas by changing one line of code
(<https://www.machinelearningplus.com/python/modin-speedup-pandas/>)

Dask – How to handle large dataframes in python using parallel computing
(<https://www.machinelearningplus.com/python/dask-tutorial/>)

Text Summarization Approaches for NLP – Practical Guide with Generative Examples
(<https://www.machinelearningplus.com/nlp/text-summarization-approaches-nlp-example/>)

Bias Variance Tradeoff – Clearly Explained
(<https://www.machinelearningplus.com/machine-learning/bias-variance-tradeoff/>)

Gradient Boosting – A Concise Introduction from Scratch
(<https://www.machinelearningplus.com/machine-learning/gradient-boosting/>)

Complete Guide to Natural Language Processing (NLP) – with Practical Examples
(<https://www.machinelearningplus.com/nlp/natural-language-processing-guide/>)

Portfolio Optimization with Python using Efficient Frontier with Practical Examples
(<https://www.machinelearningplus.com/machine-learning/portfolio-optimization-python-example/>)

What does Python Global Interpreter Lock – (GIL) do?
(<https://www.machinelearningplus.com/python/python-global-interpreter-lock-gil/>)

Logistic Regression in Julia – Practical Guide with Examples
(<https://www.machinelearningplus.com/julia/logistic-regression-in-julia-practical-guide-with-examples/>)

An ARIMA model is one where the time series was differenced at least once to make it stationary and you combine the AR and the MA terms. So the equation becomes:

(<https://www.machinelearningplus.com/wp-content/uploads/2019/02/Equation-4-min.png>)

ARIMA model in words:

Predicted Y_t = Constant + Linear combination Lags of Y (upto p lags) + Linear Combination of Lagged forecast errors (upto q lags)

The objective, therefore, is to identify the values of p , d and q . But how?

Let's start with finding the 'd'.

5. How to find the order of differencing (d) in ARIMA model

The purpose of differencing it to make the time series stationary.

But you need to be careful to not over-difference the series. Because, an over differenced series may still be stationary, which in turn will affect the model parameters.

So how to determine the right order of differencing?

The right order of differencing is the minimum differencing required to get a near-stationary series which roams around a defined mean and the ACF plot reaches to zero fairly quick.

If the autocorrelations are positive for many number of lags (10 or more), then the series needs further differencing. On the other hand, if the lag 1 autocorrelation itself is too negative, then the series is probably over-differenced.

In the event, you can't really decide between two orders of differencing, then go with the order that gives the least standard deviation in the differenced series.

Let's see how to do it with an example.

First, I am going to check if the series is stationary using the Augmented Dickey Fuller test (<https://www.machinelearningplus.com/time-series/augmented-dickey-fuller-test/>)(`adfuller()`), from the statsmodels package.

Why?

Because, you need differencing only if the series is non-stationary. Else, no differencing is needed, that is, $d=0$.

The null hypothesis of the ADF test is that the time series is non-stationary. So, if the p-value of the test is less than the significance level (0.05) then you reject the null hypothesis and infer that the time series is indeed stationary.

Tags

Chatbot

(<https://www.machinelearningplus.com/tag/chatbot/>)

Classification

(<https://www.machinelearningplus.com/tag/classification/>)

Confidence Interval

(<https://www.machinelearningplus.com/tag/confidence-interval/>)

dask

(<https://www.machinelearningplus.com/tag/dask/>)

data.table

(<https://www.machinelearningplus.com/tag/data-table/>)

Data Manipulation

(<https://www.machinelearningplus.com/tag/data-manipulation/>)

Debugging

(<https://www.machinelearningplus.com/tag/debugging/>)

Evaluation Metrics

(<https://www.machinelearningplus.com/tag/evaluation-metrics/>)

Exercises

(<https://www.machinelearningplus.com/tag/exercises/>)

FastText

(<https://www.machinelearningplus.com/tag/fasttext/>)

Gensim

(<https://www.machinelearningplus.com/tag/gensim/>)

HuggingFace

(<https://www.machinelearningplus.com/tag/huggingface/>)

Julia

(<https://www.machinelearningplus.com/tag/julia/>)

Julia Packages

(<https://www.machinelearningplus.com/tag/julia-packages/>)

LDA

(<https://www.machinelearningplus.com/tag/lda/>)

Lemmatization

(<https://www.machinelearningplus.com/tag/lemmatization/>)

Linear Regression

(<https://www.machinelearningplus.com/tag/linear-regression/>)

Logistic

(<https://www.machinelearningplus.com/tag/logistic/>)

Loop

(<https://www.machinelearningplus.com/tag/loop/>)

Machine Learning

(<https://www.machinelearningplus.com/tag/machine-learning/>)

Matplotlib

(<https://www.machinelearningplus.com/tag/matplotlib/>)

NLP

(<https://www.machinelearningplus.com/tag/nlp/>)

NLTK

(<https://www.machinelearningplus.com/tag/nltk/>)

Numpy

(<https://www.machinelearningplus.com/tag/numpy/>)

P-Value

(<https://www.machinelearningplus.com/tag/p-value/>)

plots

(<https://www.machinelearningplus.com/tag/plots/>)

Practice Exercise

(<https://www.machinelearningplus.com/tag/practice-exercise/>)

Python

(<https://www.machinelearningplus.com/tag/python/>)

So, in our case, if P Value > 0.05 we go ahead with finding the order of differencing.

```
from statsmodels.tsa.stattools import adfuller
from numpy import log
result = adfuller(df.value.dropna())
print('ADF Statistic: %f' % result[0])
print('p-value: %f' % result[1])
```

```
ADF Statistic: -2.464240
p-value: 0.124419
```

Since P-value is greater than the significance level, let's difference the series and see how the autocorrelation plot looks like.

```
import numpy as np, pandas as pd
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
import matplotlib.pyplot as plt
plt.rcParams.update({'figure.figsize':(9,7), 'figure.dpi':120})

# Import data
df = pd.read_csv('https://raw.githubusercontent.com/selva86/datasets/master/wwusage.csv (https://raw.

# Original Series
fig, axes = plt.subplots(3, 2, sharex=True)
axes[0, 0].plot(df.value); axes[0, 0].set_title('Original Series')
plot_acf(df.value, ax=axes[0, 1])

# 1st Differencing
axes[1, 0].plot(df.value.diff()); axes[1, 0].set_title('1st Order Differencing')
plot_acf(df.value.diff().dropna(), ax=axes[1, 1])

# 2nd Differencing
axes[2, 0].plot(df.value.diff().diff()); axes[2, 0].set_title('2nd Order Differencing')
plot_acf(df.value.diff().diff().dropna(), ax=axes[2, 1])

plt.show()
```

R

(<https://www.machinelearningplus.com/tag/>

Regex (<https://www.machinelearningplus.com/tag/regex/>)

Regression

(<https://www.machinelearningplus.com/tag/regression/>)

Residual Analysis

(<https://www.machinelearningplus.com/tag/residual-analysis/>)

Significance Tests

(<https://www.machinelearningplus.com/tag/significance-tests/>)

Soft Cosine Similarity

(<https://www.machinelearningplus.com/tag/soft-cosine-similarity/>)

spaCy

(<https://www.machinelearningplus.com/tag/spacy/>)

Stationarity

(<https://www.machinelearningplus.com/tag/stationarity/>)

Statistics

(<https://www.machinelearningplus.com/tag/statistics/>)

Tensorflow

(<https://www.machinelearningplus.com/tag/tensorflow/>)

TextBlob

(<https://www.machinelearningplus.com/tag/textblob/>)

TextSummarization

(<https://www.machinelearningplus.com/tag/textsummarization/>)

Text Summarization

(<https://www.machinelearningplus.com/tag/text-summarization/>)

Time Series

(<https://www.machinelearningplus.com/tag/timeseries/>)

Topic Modeling

(<https://www.machinelearningplus.com/tag/topic-modeling/>)

T Test

(<https://www.machinelearningplus.com/tag/t-test/>)

Visualization

(<https://www.machinelearningplus.com/tag/visualization>.

(https://www.machinelearningplus.com/wp-content/uploads/2019/02/output_5_0-min.png).

Order of Differencing

For the above series, the time series reaches stationarity with two orders of differencing. But on looking at the autocorrelation plot for the 2nd differencing the lag goes into the far negative zone fairly quick, which indicates, the series might have been over differenced.

So, I am going to tentatively fix the order of differencing as 1 even though the series is not perfectly stationary (weak stationarity).

```
from pmdarima.arima.utils import ndiffs
df = pd.read_csv('https://raw.githubusercontent.com/selva86/datasets/master/wwwusage.csv') (https://raw.
y = df.value

## Adf Test
ndiffs(y, test='adf') # 2

# KPSS test
ndiffs(y, test='kpss') # 0

# PP test:
ndiffs(y, test='pp') # 2
```

202

6. How to find the order of the AR term (p)

The next step is to identify if the model needs any AR terms. You can find out the required number of AR terms by inspecting the Partial Autocorrelation (PACF) plot.

But what is PACF?

Partial autocorrelation can be imagined as the correlation between the series and its lag, after excluding the contributions from the intermediate lags. So, PACF sort of conveys the pure correlation between a lag and the series. That way, you will know if that lag is needed in the AR term or not.

So what is the formula for PACF mathematically?

Partial autocorrelation of lag (k) of a series is the coefficient of that lag in the autoregression equation of Y.

$$Y_t = \alpha_0 + \alpha_1 Y_{t-1} + \alpha_2 Y_{t-2} + \alpha_3 Y_{t-3}$$

That is, suppose, if Y_t is the current series and Y_{t-1} is the lag 1 of Y , then the partial autocorrelation of lag 3 (Y_{t-3}) is the coefficient α_3 of Y_{t-3} in the above equation.

Good. Now, how to find the number of AR terms?

Any autocorrelation in a stationarized series can be rectified by adding enough AR terms. So, we initially take the order of AR term to be equal to as many lags that crosses the significance limit in the PACF plot.

```
# PACF plot of 1st differenced series
plt.rcParams.update({'figure.figsize':(9,3), 'figure.dpi':120})

fig, axes = plt.subplots(1, 2, sharex=True)
axes[0].plot(df.value.diff()); axes[0].set_title('1st Differencing')
axes[1].set_ylim=(0,5)
plot_pacf(df.value.diff().dropna(), ax=axes[1])

plt.show()
```

```
/Library/Frameworks/Python.framework/Versions/3.7/lib/python3.7/site-packages/statsmodels/regression/linear_model.p
return rho, np.sqrt(sigmassq)
```

(https://www.machinelearningplus.com/wp-content/uploads/2019/02/output_9_1-min.png).

You can observe that the PACF lag 1 is quite significant since is well above the significance line. Lag 2 turns out to be significant as well, slightly managing to cross the significance limit (blue region). But I am going to be conservative and tentatively fix the p as 1.

7. How to find the order of the MA term (q)

Just like how we looked at the PACF plot for the number of AR terms, you can look at the ACF plot for the number of MA terms. An MA term is technically, the error of the lagged forecast.

The ACF tells how many MA terms are required to remove any autocorrelation in the stationarized series.

Let's see the autocorrelation plot of the differenced series.

```
import pandas as pd
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
import matplotlib.pyplot as plt
plt.rcParams.update({'figure.figsize':(9,3), 'figure.dpi':120})

# Import data
df = pd.read_csv('https://raw.githubusercontent.com/selva86/datasets/master/austa.csv')

fig, axes = plt.subplots(1, 2, sharex=True)
axes[0].plot(df.value.diff()); axes[0].set_title('1st Differencing')
axes[1].set_ylim=(0,1.2)
plot_acf(df.value.diff().dropna(), ax=axes[1])

plt.show()
```

(https://www.machinelearningplus.com/wp-content/uploads/2019/02/output_12_0-min.png).

Order of MA Term

Couple of lags are well above the significance line. So, let's tentatively fix q as 2. When in doubt, go with the simpler model that sufficiently explains the Y.

8. How to handle if a time series is slightly under or over differenced

It may so happen that your series is slightly under differenced, that differencing it one more time makes it slightly over-differenced.

How to handle this case?

If your series is slightly under differenced, adding one or more additional AR terms usually makes it up. Likewise, if it is slightly over-differenced, try adding an additional MA term.

9. How to build the ARIMA Model

Now that you’ve determined the values of p, d and q, you have everything needed to fit the ARIMA model. Let’s use the `ARIMA()` implementation in `statsmodels` package.

```
from statsmodels.tsa.arima_model import ARIMA

# 1,1,2 ARIMA Model
model = ARIMA(df.value, order=(1,1,2))
model_fit = model.fit(dis=0)
print(model_fit.summary())
```

ARIMA Model Results						
=====						
Dep. Variable:	D.value	No. Observations:	99			
Model:	ARIMA(1, 1, 2)	Log Likelihood	-253.790			
Method:	css-mle	S.D. of innovations	3.119			
Date:	Wed, 06 Feb 2019	AIC	517.579			
Time:	23:32:56	BIC	530.555			
Sample:	1	HQIC	522.829			
=====						
	coef	std err	z	P> z	[0.025	0.975]

const	1.1202	1.290	0.868	0.387	-1.409	3.649
ar.L1.D.value	0.6351	0.257	2.469	0.015	0.131	1.139
ma.L1.D.value	0.5287	0.355	1.489	0.140	-0.167	1.224
ma.L2.D.value	-0.0010	0.321	-0.003	0.998	-0.631	0.629
Roots						
=====						
	Real	Imaginary	Modulus	Frequency		

AR.1	1.5746	+0.0000j	1.5746	0.0000		
MA.1	-1.8850	+0.0000j	1.8850	0.5000		
MA.2	545.3515	+0.0000j	545.3515	0.0000		

The model summary reveals a lot of information. The table in the middle is the coefficients table where the values under ‘coef’ are the weights of the respective terms.

Notice here the coefficient of the MA2 term is close to zero and the P-Value in ‘P>|z|’ column is highly insignificant. It should ideally be less than 0.05 for the respective X to be significant.

So, let’s rebuild the model without the MA2 term.

1,1,1 ARIMA Model

```
model = ARIMA(df.value, order=(1,1,1))
model_fit = model.fit(dispatch=0)
print(model_fit.summary())
```

ARIMA Model Results

```
=====
Dep. Variable:          D.value  No. Observations:          99
Model:                 ARIMA(1, 1, 1)  Log Likelihood      -253.790
Method:                css-mle  S.D. of innovations        3.119
Date:                 Sat, 09 Feb 2019  AIC                  515.579
Time:                 12:16:06  BIC                   525.960
Sample:                1  HQIC                   519.779
=====
```

```
=====
coef  std err      z    P>|z|    [0.025    0.975]
-----
const      1.1205    1.286    0.871    0.386   -1.400    3.641
ar.L1.D.value  0.6344    0.087    7.317    0.000    0.464    0.804
ma.L1.D.value  0.5297    0.089    5.932    0.000    0.355    0.705
=====
```

Roots

```
=====
Real      Imaginary      Modulus      Frequency
-----
AR.1      1.5764      +0.0000j      1.5764      0.0000
MA.1      -1.8879      +0.0000j      1.8879      0.5000
=====
```

The model AIC has reduced, which is good. The P Values of the AR1 and MA1 terms have improved and are highly significant (< 0.05).

Let's plot the residuals to ensure there are no patterns (that is, look for constant mean and variance).

Plot residual errors

```
residuals = pd.DataFrame(model_fit.resid)
fig, ax = plt.subplots(1,2)
residuals.plot(title="Residuals", ax=ax[0])
residuals.plot(kind='kde', title='Density', ax=ax[1])
plt.show()
```

(https://www.machinelearningplus.com/wp-content/uploads/2019/02/output_20_0-min.png)

Residuals Density

The residual errors seem fine with near zero mean and uniform variance. Let's plot the actuals against the fitted values using `plot_predict()`.

```
# Actual vs Fitted
model_fit.plot_predict(dynamic=False)
plt.show()
```

(https://www.machinelearningplus.com/wp-content/uploads/2019/02/output_22_0-min.png)

Actual vs Fitted

When you set `dynamic=False` the in-sample lagged values are used for prediction.

That is, the model gets trained up until the previous value to make the next prediction. This can make the fitted forecast and actuals look artificially good.

So, we seem to have a decent ARIMA model. But is that the best?

Can't say that at this point because we haven't actually forecasted into the future and compared the forecast with the actual performance.

So, the real validation you need now is the Out-of-Time cross-validation.

10. How to do find the optimal ARIMA model manually using Out-of-Time Cross validation

In Out-of-Time cross-validation, you take few steps back in time and forecast into the future to as many steps you took back. Then you compare the forecast against the actuals.

To do out-of-time cross-validation, you need to create the training and testing dataset by splitting the time series into 2 contiguous parts in approximately 75:25 ratio or a reasonable proportion based on time frequency of series.

Why am I not sampling the training data randomly you ask?

That's because the order sequence of the time series should be intact in order to use it for forecasting.

```
from statsmodels.tsa.stattools import acf
```

```
# Create Training and Test
```

```
train = df.value[:85]
```

```
test = df.value[85:]
```

You can now build the ARIMA model on training dataset, forecast and plot it.

```
# Build Model
```

```
# model = ARIMA(train, order=(3,2,1))
```

```
model = ARIMA(train, order=(1, 1, 1))
```

```
fitted = model.fit(dispatch=-1)
```

```
# Forecast
```

```
fc, se, conf = fitted.forecast(15, alpha=0.05) # 95% conf
```

```
# Make as pandas series
```

```
fc_series = pd.Series(fc, index=test.index)
```

```
lower_series = pd.Series(conf[:, 0], index=test.index)
```

```
upper_series = pd.Series(conf[:, 1], index=test.index)
```

```
# Plot
```

```
plt.figure(figsize=(12,5), dpi=100)
```

```
plt.plot(train, label='training')
```

```
plt.plot(test, label='actual')
```

```
plt.plot(fc_series, label='forecast')
```

```
plt.fill_between(lower_series.index, lower_series, upper_series,  
                 color='k', alpha=.15)
```

```
plt.title('Forecast vs Actuals')
```

```
plt.legend(loc='upper left', fontsize=8)
```

```
plt.show()
```

(https://www.machinelearningplus.com/wp-content/uploads/2019/02/output_28_0-min.png).

Forecast vs Actuals

From the chart, the ARIMA(1,1,1) model seems to give a directionally correct forecast. And the actual observed values lie within the 95% confidence band. That seems fine.

But each of the predicted forecasts is consistently below the actuals. That means, by adding a small constant to our forecast, the accuracy will certainly improve. So, there is definitely scope for improvement.

While doing this, I keep an eye on the P values of the AR and MA terms in the model summary. They should be as close to zero, ideally, less than 0.05.

```
# Build Model
model = ARIMA(train, order=(3, 2, 1))
fitted = model.fit(dispatch=-1)
print(fitted.summary())

# Forecast
fc, se, conf = fitted.forecast(15, alpha=0.05) # 95% conf

# Make as pandas series
fc_series = pd.Series(fc, index=test.index)
lower_series = pd.Series(conf[:, 0], index=test.index)
upper_series = pd.Series(conf[:, 1], index=test.index)

# Plot
plt.figure(figsize=(12,5), dpi=100)
plt.plot(train, label='training')
plt.plot(test, label='actual')
plt.plot(fc_series, label='forecast')
plt.fill_between(lower_series.index, lower_series, upper_series,
                 color='k', alpha=.15)
plt.title('Forecast vs Actuals')
plt.legend(loc='upper left', fontsize=8)
plt.show()
```

ARIMA Model Results						
=====						
Dep. Variable:	D2.value	No. Observations:	83			
Model:	ARIMA(3, 2, 1)	Log Likelihood	-214.248			
Method:	css-mle	S.D. of innovations	3.153			
Date:	Sat, 09 Feb 2019	AIC	440.497			
Time:	12:49:01	BIC	455.010			
Sample:	2	HQIC	446.327			
=====						
	coef	std err	z	P> z	[0.025	0.975]

const	0.0483	0.084	0.577	0.565	-0.116	0.212
ar.L1.D2.value	1.1386	0.109	10.399	0.000	0.924	1.353
ar.L2.D2.value	-0.5923	0.155	-3.827	0.000	-0.896	-0.289
ar.L3.D2.value	0.3079	0.111	2.778	0.007	0.091	0.525
ma.L1.D2.value	-1.0000	0.035	-28.799	0.000	-1.068	-0.932
Roots						
=====						
	Real	Imaginary	Modulus	Frequency		

AR.1	1.1557	-0.0000j	1.1557	-0.0000		
AR.2	0.3839	-1.6318j	1.6763	-0.2132		
AR.3	0.3839	+1.6318j	1.6763	0.2132		
MA.1	1.0000	+0.0000j	1.0000	0.0000		

(https://www.machinelearningplus.com/wp-content/uploads/2019/02/output_30_1-min.png).

Revised Forecast vs Actuals

The AIC has reduced to 440 from 515. Good. The P-values of the X terms are less than 0.05, which is great.

So overall it's much better.

Ideally, you should go back multiple points in time, like, go back 1, 2, 3 and 4 quarters and see how your forecasts are performing at various points in the year.

Here's a great practice exercise: Try to go back 27, 30, 33, 36 data points and see how the forecasts perform. The forecast performance can be judged using various accuracy metrics discussed next.

11. Accuracy Metrics for Time Series Forecast

The commonly used accuracy metrics to judge forecasts are:

1. Mean Absolute Percentage Error (MAPE)
2. Mean Error (ME)
3. Mean Absolute Error (MAE)
4. Mean Percentage Error (MPE)
5. Root Mean Squared Error (RMSE)
6. Lag 1 Autocorrelation of Error (ACF1)
7. Correlation between the Actual and the Forecast (corr)
8. Min-Max Error (minmax)

Typically, if you are comparing forecasts of two different series, the MAPE, Correlation and Min-Max Error can be used.

Why not use the other metrics?

Because only the above three are percentage errors that vary between 0 and 1. That way, you can judge how good is the forecast irrespective of the scale of the series.

The other error metrics are quantities. That implies, an RMSE of 100 for a series whose mean is in 1000's is better than an RMSE of 5 for series in 10's. So, you can't really use them to compare the forecasts of two different scaled time series.

Accuracy metrics

```
def forecast_accuracy(forecast, actual):  
    mape = np.mean(np.abs(forecast - actual)/np.abs(actual)) # MAPE  
    me = np.mean(forecast - actual) # ME  
    mae = np.mean(np.abs(forecast - actual)) # MAE  
    mpe = np.mean((forecast - actual)/actual) # MPE  
    rmse = np.mean((forecast - actual)**2)**.5 # RMSE  
    corr = np.corrcoef(forecast, actual)[0,1] # corr  
    mins = np.amin(np.hstack([forecast[:,None],  
                             actual[:,None]]), axis=1)  
    maxs = np.amax(np.hstack([forecast[:,None],  
                             actual[:,None]]), axis=1)  
    minmax = 1 - np.mean(mins/maxs) # minmax  
    acf1 = acf(fc-test)[1] # ACF1  
    return({'mape':mape, 'me':me, 'mae': mae,  
           'mpe': mpe, 'rmse':rmse, 'acf1':acf1,  
           'corr':corr, 'minmax':minmax})
```

```
forecast_accuracy(fc, test.values)
```

```
#> {'mape': 0.02250131357314834,  
#> 'me': 3.230783108990054,  
#> 'mae': 4.548322194530069,  
#> 'mpe': 0.016421001932706705,  
#> 'rmse': 6.373238534601827,  
#> 'acf1': 0.5105506325288692,  
#> 'corr': 0.9674576513924394,  
#> 'minmax': 0.02163154777672227}
```

Around 2.2% MAPE implies the model is about 97.8% accurate in predicting the next 15 observations.

Now you know how to build an ARIMA model manually.

But in industrial situations, you will be given a lot of time series to be forecasted and the forecasting exercise be repeated regularly.

So we need a way to automate the best model selection process.

12. How to do Auto Arima Forecast in Python

Like R's popular `auto.arima()` function, the `pmdarima` package provides `auto_arima()` with similar functionality.

`auto_arima()` uses a stepwise approach to search multiple combinations of p,d,q parameters and chooses the best model that has the least AIC.


```

from statsmodels.tsa.arima_model import ARIMA
import pmdarima as pm

df = pd.read_csv('https://raw.githubusercontent.com/selva86/datasets/master/wwwusage.csv (https://raw.

model = pm.auto_arima(df.value, start_p=1, start_q=1,
    test='adf',      # use adftest to find optimal 'd'
    max_p=3, max_q=3, # maximum p and q
    m=1,            # frequency of series
    d=None,         # let model determine 'd'
    seasonal=False, # No Seasonality
    start_P=0,
    D=0,
    trace=True,
    error_action='ignore',
    suppress_warnings=True,
    stepwise=True)

print(model.summary())

#> Fit ARIMA: order=(1, 2, 1); AIC=525.586, BIC=535.926, Fit time=0.060 seconds
#> Fit ARIMA: order=(0, 2, 0); AIC=533.474, BIC=538.644, Fit time=0.005 seconds
#> Fit ARIMA: order=(1, 2, 0); AIC=532.437, BIC=540.192, Fit time=0.035 seconds
#> Fit ARIMA: order=(0, 2, 1); AIC=525.893, BIC=533.648, Fit time=0.040 seconds
#> Fit ARIMA: order=(2, 2, 1); AIC=515.248, BIC=528.173, Fit time=0.105 seconds
#> Fit ARIMA: order=(2, 2, 0); AIC=513.459, BIC=523.798, Fit time=0.063 seconds
#> Fit ARIMA: order=(3, 2, 1); AIC=512.552, BIC=528.062, Fit time=0.272 seconds
#> Fit ARIMA: order=(3, 2, 0); AIC=515.284, BIC=528.209, Fit time=0.042 seconds
#> Fit ARIMA: order=(3, 2, 2); AIC=514.514, BIC=532.609, Fit time=0.234 seconds
#> Total fit time: 0.865 seconds
#>
#> ARIMA Model Results
#> =====
#> Dep. Variable:      D2.y  No. Observations:      98
#> Model:      ARIMA(3, 2, 1)  Log Likelihood      -250.276
#> Method:      css-mle  S.D. of innovations      3.069
#> Date:      Sat, 09 Feb 2019  AIC      512.552
#> Time:      12:57:22  BIC      528.062
#> Sample:      2  HQIC      518.825
#>
#> =====
#>      coef  std err      z  P>|z|  [0.025  0.975]
#> -----
#> const      0.0234   0.058   0.404   0.687   -0.090   0.137
#> ar.L1.D2.y   1.1586   0.097  11.965   0.000   0.969   1.348
#> ar.L2.D2.y  -0.6640   0.136  -4.890   0.000  -0.930  -0.398
#> ar.L3.D2.y   0.3453   0.096   3.588   0.001   0.157   0.534
#> ma.L1.D2.y  -1.0000   0.028  -36.302   0.000  -1.054  -0.946
#>
#> Roots
#> =====
#>      Real    Imaginary    Modulus    Frequency
#> -----
#> AR.1      1.1703    -0.0000j     1.1703    -0.0000
#> AR.2      0.3763    -1.5274j     1.5731    -0.2116
#> AR.3      0.3763    +1.5274j     1.5731     0.2116
#> MA.1      1.0000    +0.0000j     1.0000     0.0000
#> -----

```

13. How to interpret the residual plots in ARIMA model

Let's review the residual plots using `stepwise_fit`.

```
model.plot_diagnostics(figsize=(7,5))  
plt.show()
```

(https://www.machinelearningplus.com/wp-content/uploads/2019/02/output_39_0-min.png).

Residuals Chart

So how to interpret the plot diagnostics?

Top left: The residual errors seem to fluctuate around a mean of zero and have a uniform variance.

Top Right: The density plot suggest normal distribution with mean zero.

Bottom left: All the dots should fall perfectly in line with the red line. Any significant deviations would imply the distribution is skewed.

Bottom Right: The Correlogram, aka, ACF plot shows the residual errors are not autocorrelated. Any autocorrelation would imply that there is some pattern in the residual errors which are not explained in the model. So you will need to look for more X's (predictors) to the model.

Overall, it seems to be a good fit. Let's forecast.

```

# Forecast
n_periods = 24
fc, confint = model.predict(n_periods=n_periods, return_conf_int=True)
index_of_fc = np.arange(len(df.value), len(df.value)+n_periods)

# make series for plotting purpose
fc_series = pd.Series(fc, index=index_of_fc)
lower_series = pd.Series(confint[:, 0], index=index_of_fc)
upper_series = pd.Series(confint[:, 1], index=index_of_fc)

# Plot
plt.plot(df.value)
plt.plot(fc_series, color='darkgreen')
plt.fill_between(lower_series.index,
                 lower_series,
                 upper_series,
                 color='k', alpha=.15)

plt.title("Final Forecast of WWW Usage")
plt.show()

```

(https://www.machinelearningplus.com/wp-content/uploads/2019/02/output_41_0-min.png).

Final Forecast of WWW Usage

14. How to automatically build SARIMA model in python

The problem with plain ARIMA model is it does not support seasonality.

If your time series has defined seasonality, then, go for SARIMA which uses seasonal differencing.

Seasonal differencing is similar to regular differencing, but, instead of subtracting consecutive terms, you subtract the value from previous season.

So, the model will be represented as SARIMA(p,d,q)x(P,D,Q), where, P, D and Q are SAR, order of seasonal differencing and SMA terms respectively and 'x' is the frequency of the time series.

If your model has well defined seasonal patterns, then enforce D=1 for a given frequency 'x'.

Here's some practical advice on building SARIMA model:

As a general rule, set the model parameters such that D never exceeds one. And the total differencing 'd + D' never exceeds 2. Try to keep only either SAR or SMA terms if your model has seasonal components.

Let's build an SARIMA model on 'a10' – the drug sales dataset.

```
# Import
data = pd.read_csv('https://raw.githubusercontent.com/selva86/datasets/master/a10.csv')

# Plot
fig, axes = plt.subplots(2, 1, figsize=(10,5), dpi=100, sharex=True)

# Usual Differencing
axes[0].plot(data[:], label='Original Series')
axes[0].plot(data[:].diff(1), label='Usual Differencing')
axes[0].set_title('Usual Differencing')
axes[0].legend(loc='upper left', fontsize=10)

# Seasonal Differencing
axes[1].plot(data[:], label='Original Series')
axes[1].plot(data[:].diff(12), label='Seasonal Differencing', color='green')
axes[1].set_title('Seasonal Differencing')
plt.legend(loc='upper left', fontsize=10)
plt.suptitle('a10 - Drug Sales', fontsize=16)
plt.show()
```

(https://www.machinelearningplus.com/wp-content/uploads/2019/02/output_43_0-min.png).

Seasonal Differencing

As you can clearly see, the seasonal spikes is intact after applying usual differencing (lag 1). Whereas, it is rectified after seasonal differencing.

Let's build the SARIMA model using `pmdarima`'s `auto_arima()`. To do that, you need to set `seasonal=True`, set the frequency `m=12` for month wise series and enforce `D=1`.

```
# !pip3 install pyramid-arma
```

```
import pmdarima as pm
```

```
# Seasonal - fit stepwise auto-ARIMA
```

```
smodel = pm.auto_arma(data, start_p=1, start_q=1,  
                      test='adf',  
                      max_p=3, max_q=3, m=12,  
                      start_P=0, seasonal=True,  
                      d=None, D=1, trace=True,  
                      error_action='ignore',  
                      suppress_warnings=True,  
                      stepwise=True)
```

```
smodel.summary()
```

```
Fit ARIMA: order=(1, 0, 1) seasonal_order=(0, 1, 1, 12); AIC=534.818, BIC=551.105, Fit time=1.742 seconds  
Fit ARIMA: order=(0, 0, 0) seasonal_order=(0, 1, 0, 12); AIC=624.061, BIC=630.576, Fit time=0.028 seconds  
Fit ARIMA: order=(1, 0, 0) seasonal_order=(1, 1, 0, 12); AIC=596.004, BIC=609.034, Fit time=0.683 seconds  
Fit ARIMA: order=(0, 0, 1) seasonal_order=(0, 1, 1, 12); AIC=611.475, BIC=624.505, Fit time=0.709 seconds  
Fit ARIMA: order=(1, 0, 1) seasonal_order=(1, 1, 1, 12); AIC=557.501, BIC=577.046, Fit time=3.687 seconds  
(...TRUNCATED...)  
Fit ARIMA: order=(3, 0, 0) seasonal_order=(1, 1, 1, 12); AIC=554.570, BIC=577.372, Fit time=2.431 seconds  
Fit ARIMA: order=(3, 0, 0) seasonal_order=(0, 1, 0, 12); AIC=554.094, BIC=570.381, Fit time=0.220 seconds  
Fit ARIMA: order=(3, 0, 0) seasonal_order=(0, 1, 2, 12); AIC=529.502, BIC=552.305, Fit time=2.120 seconds  
Fit ARIMA: order=(3, 0, 0) seasonal_order=(1, 1, 2, 12); AIC=nan, BIC=nan, Fit time=nan seconds  
Total fit time: 31.613 seconds
```

The model has estimated the AIC and the P values of the coefficients look significant. Let's look at the residual diagnostics plot.

The best model SARIMAX(3, 0, 0)x(0, 1, 1, 12) has an AIC of 528.6 and the P Values are significant.

Let's forecast for the next 24 months.

```
# Forecast
n_periods = 24
fitted, confint = smodel.predict(n_periods=n_periods, return_conf_int=True)
index_of_fc = pd.date_range(data.index[-1], periods = n_periods, freq='MS')

# make series for plotting purpose
fitted_series = pd.Series(fitted, index=index_of_fc)
lower_series = pd.Series(confint[:, 0], index=index_of_fc)
upper_series = pd.Series(confint[:, 1], index=index_of_fc)

# Plot
plt.plot(data)
plt.plot(fitted_series, color='darkgreen')
plt.fill_between(lower_series.index,
                 lower_series,
                 upper_series,
                 color='k', alpha=.15)

plt.title("SARIMA - Final Forecast of a10 - Drug Sales")
plt.show()
```

There you have a nice forecast that captures the expected seasonal demand pattern.

15. How to build SARIMAX Model with exogenous variable

The SARIMA model we built is good. I would stop here typically.

But for the sake of completeness, let's try and force an external predictor, also called, 'exogenous variable' into the model. This model is called the SARIMAX model.

The only requirement to use an exogenous variable is you need to know the value of the variable during the forecast period as well.

For the sake of demonstration, I am going to use the seasonal index from the [classical seasonal decomposition \(https://www.machinelearningplus.com/time-series/time-series-analysis-python/\)](https://www.machinelearningplus.com/time-series/time-series-analysis-python/) on the latest 36 months of data.

Why the seasonal index? Isn't SARIMA already modeling the seasonality, you ask?

You are correct.

But also, I want to see how the model looks if we force the recent seasonality pattern into the training and forecast.

Secondly, this is a good variable for demo purpose. So you can use this as a template and plug in any of your variables into the code. The seasonal index is a good exogenous variable because it repeats every frequency cycle, 12 months in this case.

So, you will always know what values the seasonal index will hold for the future forecasts.

```
# Import Data
data = pd.read_csv('https://raw.githubusercontent.com/selva86/datasets/master/a10.csv', https://raw.githubusercontent.com/selva86/datasets/master/a10.csv')
```

Let's compute the seasonal index so that it can be forced as a (exogenous) predictor to the SARIMAX model.

```

# Compute Seasonal Index
from statsmodels.tsa.seasonal import seasonal_decompose
from dateutil.parser import parse

# multiplicative seasonal component
result_mul = seasonal_decompose(data['value'][-36:], # 3 years
                                model='multiplicative',
                                extrapolate_trend='freq')

seasonal_index = result_mul.seasonal[-12:].to_frame()
seasonal_index['month'] = pd.to_datetime(seasonal_index.index).month

# merge with the base data
data['month'] = data.index.month
df = pd.merge(data, seasonal_index, how='left', on='month')
df.columns = ['value', 'month', 'seasonal_index']
df.index = data.index # reassign the index.

```

The exogenous variable (seasonal index) is ready. Let's build the SARIMAX model.

```

import pmdarima as pm

# SARIMAX Model
sxmodel = pm.auto_arima(df[['value']], exogenous=df[['seasonal_index']],
                        start_p=1, start_q=1,
                        test='adf',
                        max_p=3, max_q=3, m=12,
                        start_P=0, seasonal=True,
                        d=None, D=1, trace=True,
                        error_action='ignore',
                        suppress_warnings=True,
                        stepwise=True)

sxmodel.summary()

```

```

Fit ARIMA: order=(1, 0, 1) seasonal_order=(0, 1, 1, 12); AIC=536.818, BIC=556.362, Fit time=2.083 seconds
Fit ARIMA: order=(0, 0, 0) seasonal_order=(0, 1, 0, 12); AIC=626.061, BIC=635.834, Fit time=0.033 seconds
Fit ARIMA: order=(1, 0, 0) seasonal_order=(1, 1, 0, 12); AIC=598.004, BIC=614.292, Fit time=0.682 seconds
Fit ARIMA: order=(0, 0, 1) seasonal_order=(0, 1, 1, 12); AIC=613.475, BIC=629.762, Fit time=0.510 seconds
Fit ARIMA: order=(1, 0, 1) seasonal_order=(1, 1, 1, 12); AIC=559.530, BIC=582.332, Fit time=3.129 seconds
(...Truncated...)
Fit ARIMA: order=(3, 0, 0) seasonal_order=(0, 1, 0, 12); AIC=556.094, BIC=575.639, Fit time=0.260 seconds
Fit ARIMA: order=(3, 0, 0) seasonal_order=(0, 1, 2, 12); AIC=531.502, BIC=557.562, Fit time=2.375 seconds
Fit ARIMA: order=(3, 0, 0) seasonal_order=(1, 1, 2, 12); AIC=nan, BIC=nan, Fit time=nan seconds
Total fit time: 30.781 seconds

```


(<https://www.machinelearningplus.com/wp-content/uploads/2019/02/Model-Summary-2-min.png>).

So, we have the model with the exogenous term. But the coefficient is very small for x_1 , so the contribution from that variable will be negligible. Let's forecast it anyway.

We have effectively forced the latest seasonal effect of the latest 3 years into the model instead of the entire history.

Alright let's forecast into the next 24 months. For this, you need the value of the seasonal index for the next 24 months.

```

# Forecast
n_periods = 24
fitted, confint = sxmodel.predict(n_periods=n_periods,
                                  exogenous=np.tile(seasonal_index.value, 2).reshape(-1,1),
                                  return_conf_int=True)

index_of_fc = pd.date_range(data.index[-1], periods = n_periods, freq='MS')

# make series for plotting purpose
fitted_series = pd.Series(fitted, index=index_of_fc)
lower_series = pd.Series(confint[:, 0], index=index_of_fc)
upper_series = pd.Series(confint[:, 1], index=index_of_fc)

# Plot
plt.plot(data['value'])
plt.plot(fitted_series, color='darkgreen')
plt.fill_between(lower_series.index,
                 lower_series,
                 upper_series,
                 color='k', alpha=.15)

plt.title("SARIMAX Forecast of a10 - Drug Sales")
plt.show()

```

(https://www.machinelearningplus.com/wp-content/uploads/2019/02/output_59_0-min.png)

SARIMAX Forecast

16. Practice Exercises

In the AirPassengers dataset, go back 12 months in time and build the SARIMA forecast for the next 12 months.

- 1. Is the series stationary? If not what sort of differencing is required?
- 2. What is the order of your best model?
- 3. What is the AIC of your model?
- 4. What is the MAPE achieved in OOT cross-validation?
- 5. What is the order of the best model predicted by `auto_arima()` method?

17. Conclusion

Congrats if you reached this point. Give yourself a BIG hug if you were able to solve the practice exercises.

I really hope you found this useful?

We have covered a lot of concepts starting from the very basics of forecasting, AR, MA, ARIMA, SARIMA and finally the SARIMAX model. If you have any questions please write in the comments section. Meanwhile, I will work on the next article.

Happy Learning!

ALSO ON MACHINELEARNINGPLUS.COM

Understanding the math and how it works? ...

2 years ago • 15 comments

Cosine similarity is a metric used to measure how similar the documents are ...

How to Visualize Distributions in Python

2 years ago • 1 comment

Matplotlib histogram is used to visualize the frequency distribution of numeric ...

A Comprehensive Guide with Examples

2 years ago • 25 comments

Time series is a sequence of observations recorded at regular time intervals. This ...

Confidence

3 months ago • 1

Confidence interval is a measure to quantify the uncertainty in ...

What do you think?
110 Responses

👍 Upvote

😍 Love

41 Comments machinelearningplus.com Disqus' Privacy Policy Login

Recommend 11 Tweet Share Sort by Newest

Join the discussion...

LOG IN WITH OR SIGN UP WITH DISQUS ?

Name

Giordano Pazzaglia • 10 days ago

In coding part :

```
# Make as pandas series
fc_series = pd.Series(fc, index=test.index)
lower_series = pd.Series(conf[:, 0], index=test.index)
upper_series = pd.Series(conf[:, 1], index=test.index)
```

The ouput has : `TypeError: 'builtin_function_or_method' object is not iterable`

Can someone help me ? Thanks you !

```
TypeError                                Traceback (most recent call last)
<ipython-input-63-95101ec585a9> in <module>
    14
    15 # Make as pandas series
```

```

---> 16 fc_series = pd.Series(fc, index=test.index)
      17 lower_series = pd.Series(conf[:, 0], index= test.index)
      18 upper_series = pd.Series(conf[:, 1], index= test.index)

~\anaconda3\lib\site-packages\pandas\core\series.py in __init__(self, data, index, dtype, name, copy, fastpath)
      213
      214         if index is not None:
--> 215             index = ensure_index(index)
      216
      217         if data is None:

```

[see more](#)

^ | v • Reply • Share ›



Smart Talk • 4 months ago

thanks for detailed explanation. But i want to ask you one thing, The Pacf is quite high where the y axes is more than one, how could the partial autocorrelation may reach number more than one, thank you

^ | v • Reply • Share ›



HBeing3 • 7 months ago

Thank you for the detailed step-by-step instruction.
I've followed this to the last bit.
Just wondering the last section of code

should `exogenous=np.tile(seasonal_index.value, 2).reshape(-1,1)`,
be like this
`exogenous=np.tile(seasonal_index.seasonal, 2).reshape(-1,1)`,

^ | v • Reply • Share ›



Monia • 8 months ago

hello ,
I applied the auto-arima on my data which does not carry seasonality (univariate-time series)however the best model provided was sarimax (3,1,2)

^ | v • Reply • Share ›



Mike Heimlich • 9 months ago

Did anyone else had the problem that in each of the (partial) autocorr and normal autocorrelation graphs only the first 20 values were used?

^ | v • Reply • Share ›



Dmitry Serebryanskiy → Mike Heimlich • 8 months ago

Just specify lags parameter:
`plot_acf(df.value, ax=axes[0, 1], lags=80)`

Or if you want it to be exactly as in the article:
`plot_acf(df.value, ax=axes[0, 1], lags=df.value.shape[0] - 1)`
`plot_acf(df.value.diff().dropna(), ax=axes[1, 1], lags=df.value.shape[0] - 2)`
`plot_acf(df.value.diff().diff().dropna(), ax=axes[2, 1], lags=df.value.shape[0] - 3)`

The idea is that lags parameter should equals the length of the plotted series.

^ | v • Reply • Share ›



Davide Manniello → Mike Heimlich • 9 months ago

yes, but I don't know how to solve it

^ | v • Reply • Share ›



Thushara De Silva • 10 months ago

I follow the codes to fit a SARIMAX. However, I didn't get the 'x1' term in my model. I got the same SARIMA model. Any comments on this?

^ | v • Reply • Share ›



Asad • 10 months ago

`plot_diagnostics` function giving error: `AttributeError: 'ARMAResults' object has no attribute 'plot_diagnostics'`. I checked, I installed and imported all the necessary libraries and modules.

^ | v • Reply • Share ›



Saurabh • a year ago • edited

Hi! this is an amazing article! However, I noticed that

- 1) you confined your PACF plot to (0,5) by setting the y limit - why did you do that?
- 2) Secondly, I thought PACF values should only be between 1 and - 1 but for your plot, some of them are >1, could you explain that too?
- 3) My PACF lag 1 is -0.5 and significant: does this suggest that I use an AR(1) model?

Thanks again!

^ | v • Reply • Share ›



Mustafa Kathawala → Saurabh • 6 months ago

I have the same questions! Did you get the answer?

^ | v • Reply • Share ›



Tukiya • a year ago

Great and thank you

^ | v • Reply • Share ›



Antoine Blais • a year ago

Hi, your article is amazing!! I have a question regarding the SARIMAX model. You say that the only requirement is to use an exogenous variable. Could we use montly average tempreature or even snowfall in the region where the data are from? I am trying to forecast sales of products that are strongly tied to seasons and snowfall.

And since the snowfall / temperature are "easy to forecast" since it is somehow constant, I was wondering if we could call this exogenous.

Thanks for your help!

^ | v • Reply • Share ›



Tanjilul Amin • a year ago

Hello,

when executing the line "from pmdarima.arima.utils import ndiffs"
I am getting the error: "module 'scipy' has no attribute 'stats'"

I cannot solve this issue. Even using "from scipy import stats" does not have any effect.

Do you have nay advise on how to troubleshoot this?

Thanks

^ | v • Reply • Share ›



TXN • a year ago

Many Thanks

^ | v • Reply • Share ›



pankaj • a year ago

Thanks for this beatifull artical on Time Series! :)

^ | v • Reply • Share ›



Amr a • a year ago

Hello,

Thank you very much for this article, it's great, in your graphs the forecast line does not connect to the actual data. Do you know how I can fix this?

^ | v • Reply • Share ›



Abdellah sebti • a year ago

I am truly thankful for this amazing article and your effort,
it was one of the best informative article I found

^ | v • Reply • Share ›



Selva Prabhakaran Mod → Abdellah sebti • a year ago

Thanks very much :)

^ | v • Reply • Share ›



Asto Maa Sadgamy • a year ago

Hi,

Just wondering, Can these model are suitable for long data training and forecasting. let say to train 40K data and forecast next 30k data ? if not could you suggest any time series method for that ?

^ | v • Reply • Share ›



Selva Prabhakaran Mod → Asto Maa Sadgamy • a year ago

Depends on the type of data and the problem you are dealing with. Can't suggest a method without context.

^ | v • Reply • Share ›



Asto Maa Sadgamy → Selva Prabhakaran • a year ago

I am working on wind speed time series data that covers 10 years of data of 10 mins average.
So using these 10 years of wind speed data, can I forecast next 10 years wind speed data ?

^ | v • Reply • Share ›



mjos • 2 years ago

Hi, I'm not sure if the numbering of beta coefficients in the equations for errors (see "The errors E_t and E_{t-1} are the errors from the following equation:") are correct. Could you please have a look at it?

^ | v • Reply • Share ›



Selva Prabhakaran Mod → mjos • a year ago

Oh, yeah. Its wrong. Thanks for finding this.. and apologies for such a late reply

^ | v • Reply • Share ›



Vũ Trung Nghĩa • 2 years ago

The original serie may be changed many time because differencing, detrend,... So when my model predict 100 next value, do I need to transform it to the original serie by the inverse process that I do when differencing, detrend the series or is there another the way to do that?

^ | v • Reply • Share ›



Selva Prabhakaran Mod → Vũ Trung Nghĩa • a year ago

The algo will give do it internally.. just specify the p,d and q. Not required if you use autoarima.

^ | v • Reply • Share ›



Dídac Cortiada • 2 years ago • edited

Great post! Really educational and helpful. Thanks :)

^ | v • Reply • Share ›



Selva Prabhakaran Mod → Dídac Cortiada • a year ago

Thanks :)

^ | v • Reply • Share ›



Wael Abu Rezeq • 2 years ago

I would like to thank you for this AMAZING and comprehensive article. I hereby have two questions:

1- If we are not considering seasonality, then $M=1$ (since the data is on monthly basis) in `auto_arima` exercise. However, If we have seasonality then $M=12$ to look at seasonality inside every 12 months. My questions is what if the dataset frequency in days not month, what should these variables hold?

2- We are forecasting the dataset on time series aggregations, what If we want to forecast including other attributes like Region or Sales Channel incorporated with time series (so suppose we want to do the forecast for a dataset that has the attributes of TimeSeries, Region, Value) per region?

Looking forward to hearing from you

Thanks for your great support

^ | v • Reply • Share ›



abu bakar siddique • 2 years ago

The computed initial AR coefficients are not stationary

You should induce stationarity, choose a different model order, or you can pass your own `start_params`...

What should i do now...?

1 ^ | v • Reply • Share ›



Julien → abu bakar siddique • 2 years ago • edited

I have the same issue...

Running this code seems to provide a `ValueError`:

```
from statsmodels.tsa.arima_model import ARIMA
# 1,1,2 ARIMA Model
model = ARIMA(df.value, order=(1,1,2))
model_fit = model.fit(dispatch=0)
print(model_fit.summary())
```

^ | v • Reply • Share ›



esko1779 → Julien • a year ago

You should try `order=(0,1,0)` or `order=(1,2,1)` for example (but 0,1,0 is the best, I think).

Date seems to be truncated somehow

^ | v • Reply • Share ›



Selva Prabhakaran Mod → Julien • a year ago

Probably some problem with the data you're using

^ | v 1 • Reply • Share ›



Phuc Coi • 2 years ago



This is very cool :D

^ | v • Reply • Share ›



Selva Prabhakaran Mod ➔ Phuc Coi • a year ago

Thanks :)

^ | v • Reply • Share ›



Akshita Gupta • 2 years ago

This article is quite informative. But I would like to ask what needs to be done if the data presents the double seasonality.

1 ^ | v • Reply • Share ›



Selva Prabhakaran Mod ➔ Akshita Gupta • a year ago

You might want to look at multiseasonal time series forecasting algo. The 'forecast' package in R has an implementation.

^ | v • Reply • Share ›



Betül Okan • 2 years ago

That's a great article. Everything is explained very well. Thank you very much.

^ | v • Reply • Share ›



Selva Prabhakaran Mod ➔ Betül Okan • 2 years ago

Welcome :)

^ | v • Reply • Share ›



Eric Kumar • 2 years ago

Great article, very concise and current. One request. it would be great if the data-sets are linked. If the data sets are linked somewhere, they are hard to find.

Data set for the #10

https://raw.githubusercontent.com/ezoi/ezoi/master/data/ezoi_data.csv

^ | v • Reply • Share ›



Selva Prabhakaran Mod ➔ Eric Kumar • 2 years ago • edited

Thanks for pointing it out Eric. I will make them available via links.

EDIT: Updated!

^ | v • Reply • Share ›

[Subscribe](#) [Add Disqus to your site](#)[Add Disqus](#) [Do Not Sell My Data](#)

(<https://www.ezoic.com/what-is-ezoic/>)

[report this ad](#)

Search ...

Search

Subscribe to Blog

Enter your email address to receive notifications of new posts by email.

Email Address

Subscribe

Copyright ML+ (<https://www.machinelearningplus.com/>). All rights reserved.

[Home \(https://www.machinelearningplus.com/\)](https://www.machinelearningplus.com/) [Contact Us \(https://www.machinelearningplus.com/contact-us/\)](https://www.machinelearningplus.com/contact-us/)

[Privacy Policy \(https://www.machinelearningplus.com/privacy-policy/\)](https://www.machinelearningplus.com/privacy-policy/) [About Selva \(https://www.machinelearningplus.com/about/\)](https://www.machinelearningplus.com/about/)

[Terms and Conditions \(https://www.machinelearningplus.com/terms-of-use/\)](https://www.machinelearningplus.com/terms-of-use/)