☰ | Navigation

Start Here     Blog     Books     About     Contact

Search...                                                    🔍

# Sensitivity Analysis of History Size to Forecast Skill with ARIMA in Python

by **Jason Brownlee** on March 27, 2017 in **Time Series**

How much history is required for a time series forecast model?

This is a problem-specific question that we can investigate by designing an experiment.

In this tutorial, you will discover the effect that history size has on the skill of an ARIMA forecast model in Python.
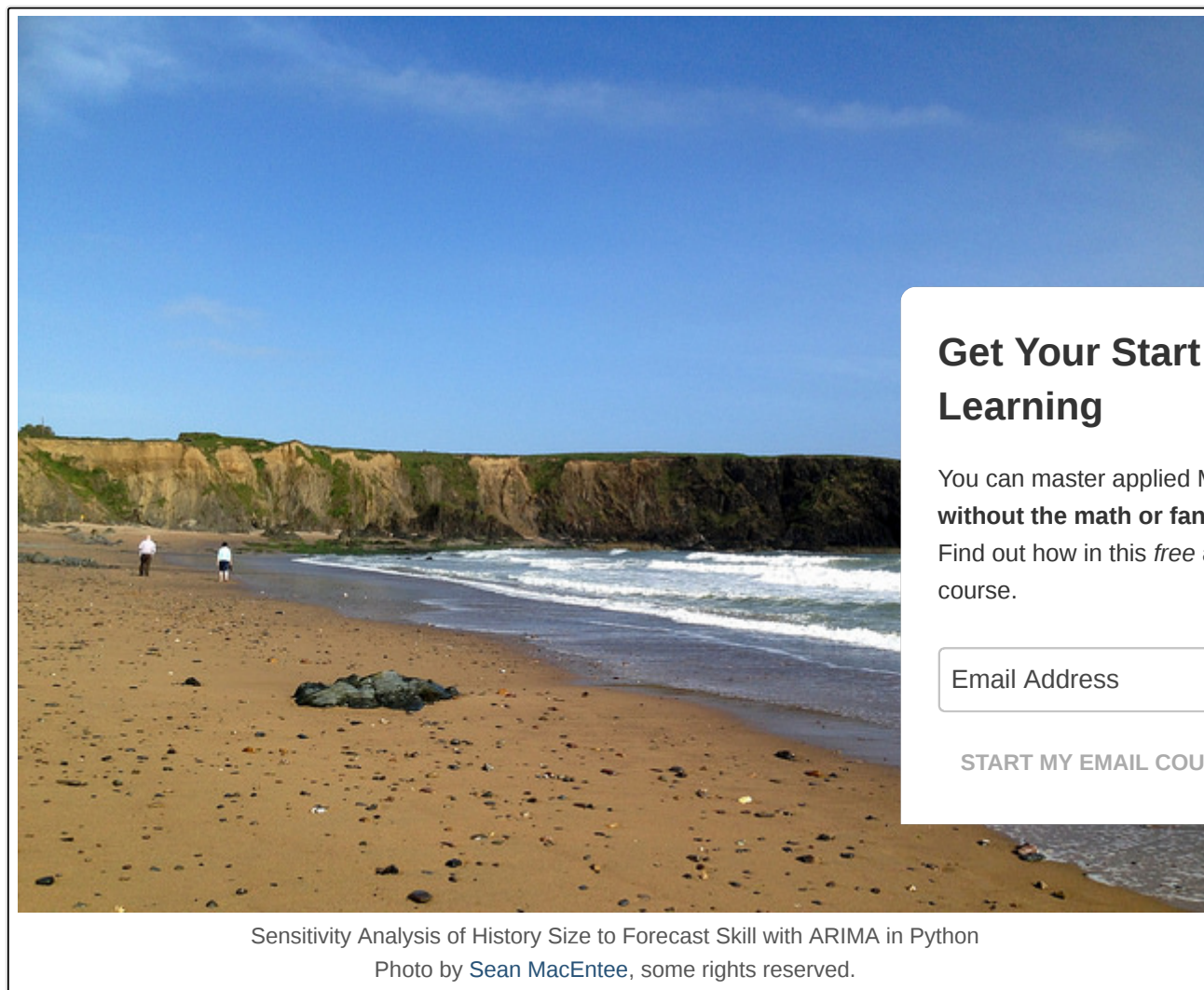
Specifically, in this tutorial, you will:

- Load a standard dataset and fit an ARIMA model.
- Design and execute a sensitivity analysis of the number of years of historic data to model skill.
- Analyze the results of the sensitivity analysis.

This will provide a template for performing a similar sensitivity analysis of historical data set size on y

**Get Your Start in Machine Learning**

Let's get started.

- **Update Aug/2017**: Fixed a bug where the models were constructed on the raw data instead of the seasonally differenced version of the data. Thanks David Ravnsborg!



**Get Your Start in Machine Learning**

You can master applied Machine Learning **without the math or fancy degree**.
Find out how in this *free* and *practical* email course.

Email Address

START MY EMAIL COURSE

Sensitivity Analysis of History Size to Forecast Skill with ARIMA in Python
Photo by Sean MacEntee, some rights reserved.

## Minimum Daily Temperatures Dataset

This dataset describes the minimum daily temperatures over 10 years (1981-1990) in the city of Mel

Get Your Start in Machine Learning

The units are in degrees Celsius and there are 3,650 observations. The source of the data is credited as the Australian Bureau of Meteorology.

Download the dataset and save it in your current working directory with the filename "*daily-minimum-temperatures.csv*".

Note: The downloaded file contains some question mark ("?") characters that must be removed before you can use the dataset. Open the file in a text editor and remove the "?" characters. Also, remove any footer information in the file.

The example below loads the dataset as a Pandas Series.

```
1  # line plot of time series
2  from pandas import Series
3  from matplotlib import pyplot
4  # load dataset
5  series = Series.from_csv('daily-minimum-temperatures.csv', header=0)
6  # display first few rows
7  print(series.head(20))
8  # line plot of dataset
9  series.plot()
10 pyplot.show()
```

Running the example prints the first 20 rows of the loaded file.

```
1  Date
2  1981-01-01 20.7
3  1981-01-02 17.9
4  1981-01-03 18.8
5  1981-01-04 14.6
6  1981-01-05 15.8
7  1981-01-06 15.8
8  1981-01-07 15.8
9  1981-01-08 17.4
10 1981-01-09 21.8
11 1981-01-10 20.0
12 1981-01-11 16.2
13 1981-01-12 13.3
14 1981-01-13 16.7
15 1981-01-14 21.5
16 1981-01-15 25.0
17 1981-01-16 20.7
18 1981-01-17 20.6
19 1981-01-18 24.8
20 1981-01-19 17.7
21 1981-01-20 15.5
```

**Get Your Start in Machine Learning**  ✕

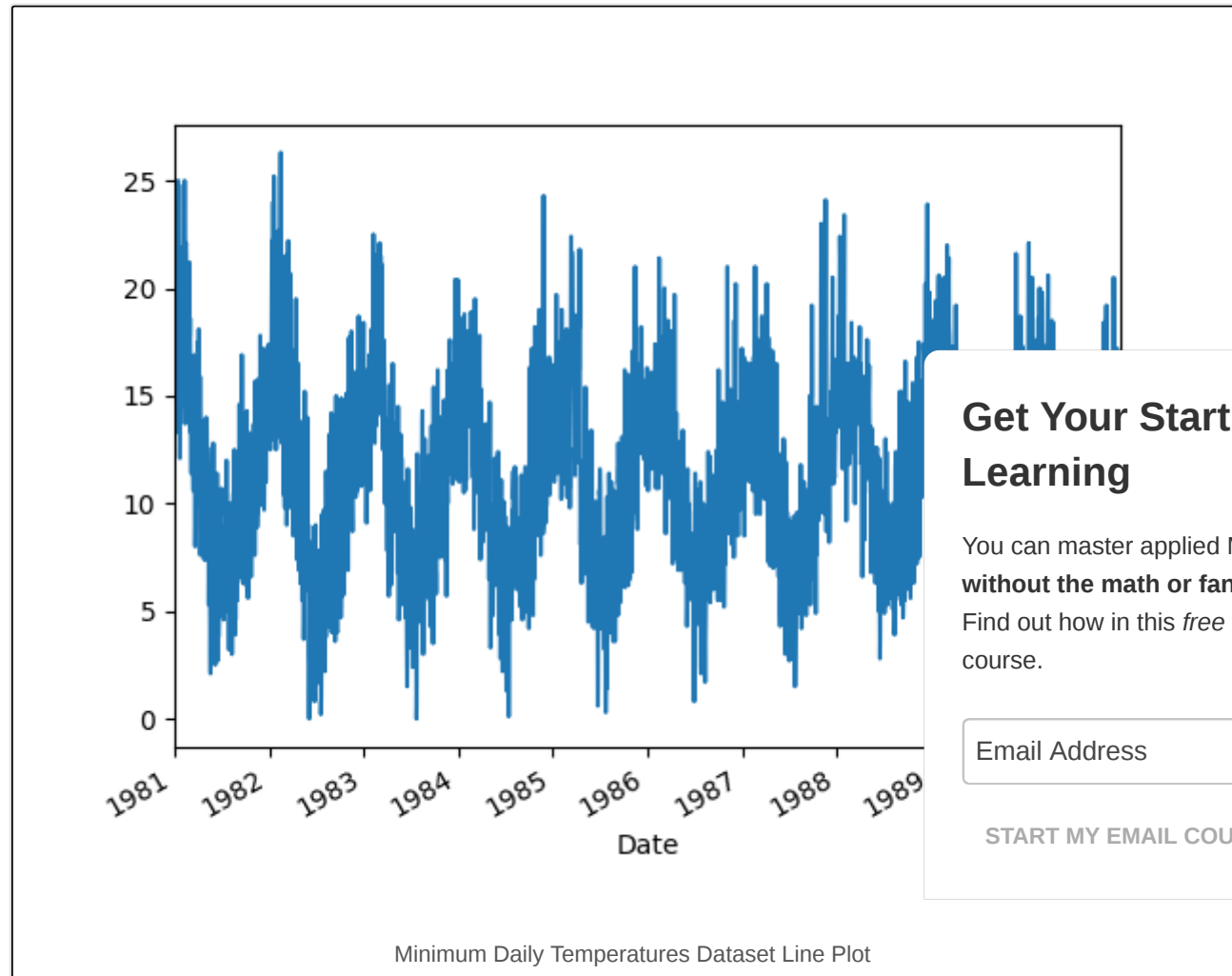You can master applied Machine Learning **without the math or fancy degree**.
Find out how in this *free* and *practical* email course.

Email Address

START MY EMAIL COURSE

Get Your Start in Machine Learning

```
22  Name: Temp, dtype: float64
```

Then the data is graphed as a line plot showing the seasonal pattern.



Minimum Daily Temperatures Dataset Line Plot

## ARIMA Forecast Model

In this section, we will fit an ARIMA forecast model to the data.

The parameters of the model will not be tuned, but will be skillful.

**Get Your Start in Machine Learning**

You can master applied Machine Learning **without the math or fancy degree**.
Find out how in this *free* and *practical* email course.

Email Address

START MY EMAIL COURSE

The data contains a one-year seasonal component that must be removed to make the data stationary and suitable for use with an ARIMA model.

We can take the seasonal difference by subtracting the observation from one year ago (365 days). This is rough in that it does not account for leap years. It also means that the first year of data will be unavailable for modeling as there is no data one year before to difference the data.

```
1  # seasonal difference
2  differenced = series.diff(365)
3  # trim off the first year of empty data
4  differenced = differenced[365:]
```

We will fit an ARIMA(7,0,0) model to the data and print the summary information. This demonstrates that the model is stable.

```
1  # fit model
2  model = ARIMA(differenced, order=(7,0,0))
3  model_fit = model.fit(trend='nc', disp=0)
4  print(model_fit.summary())
```

Putting this all together, the complete example is listed below.

```
1   # fit an ARIMA model
2   from pandas import Series
3   from matplotlib import pyplot
4   from statsmodels.tsa.arima_model import ARIMA
5   # load dataset
6   series = Series.from_csv('daily-minimum-temperatures.csv', header=0)
7   # seasonal difference
8   differenced = series.diff(365)
9   # trim off the first year of empty data
10  differenced = series[365:]
11  # fit model
12  model = ARIMA(differenced, order=(7,0,0))
13  model_fit = model.fit(trend='nc', disp=0)
14  print(model_fit.summary())
```

Running the example provides a summary of the fit ARIMA model.

```
1                     ARMA Model Results
2  ==============================================================================
3  Dep. Variable:              Temp   No. Observations:              3285
4  Model:                 ARMA(7, 0)   Log Likelihood             -8690.089
5  Method:                   css-mle   S.D. of innovations            3.409
6  Date:           Fri, 25 Aug 2017   AIC                        17396.178
7  Time:                   15:02:59   BIC                        17444.955
8  Sample:                01-01-1982   HQIC                       17413.643
```

```
 9                      - 12-31-1990
10  ==============================================================
11               coef    std err        z      P>|z|     [0.025     0.975]
12  --------------------------------------------------------------
13  ar.L1.Temp    0.5278    0.017    30.264    0.000     0.494     0.562
14  ar.L2.Temp   -0.1099    0.020    -5.576    0.000    -0.149    -0.071
15  ar.L3.Temp    0.0286    0.020     1.441    0.150    -0.010     0.067
16  ar.L4.Temp    0.0307    0.020     1.549    0.122    -0.008     0.070
17  ar.L5.Temp    0.0090    0.020     0.456    0.648    -0.030     0.048
18  ar.L6.Temp    0.0164    0.020     0.830    0.407    -0.022     0.055
19  ar.L7.Temp    0.0272    0.017     1.557    0.120    -0.007     0.061
20                              Roots
21  ==============================================================
22               Real          Imaginary          Modulus         Frequency
23  --------------------------------------------------------------
24  AR.1        1.3305         -0.0000j           1.3305          -0.0000
25  AR.2        0.9936         -1.1966j           1.5553          -0.1397
26  AR.3        0.9936         +1.1966j           1.5553           0.1397
27  AR.4       -0.2067         -1.7061j           1.7186          -0.2692
28  AR.5       -0.2067         +1.7061j           1.7186           0.2692
29  AR.6       -1.7536         -0.8938j           1.9683          -0.4250
30  AR.7       -1.7536         +0.8938j           1.9683           0.4250
31  --------------------------------------------------------------
```

## Model History Sensitivity Analysis

In this section, we will explore the effect that history size has on the skill of the fit model.

The original data has 10 years of data. Seasonal differencing leaves us with 9 years of data. We will perform walk-forward validation across this final year. The day-by-day forecasts will be collected and

The day-by-day forecasts will be collected and a root mean squared error (RMSE) score will be calc

The snippet below separates the seasonally adjusted data into training and test datasets.

```
1  train, test = differenced[differenced.index < '1990'], differenced['1990']
```

It is important to choose an interval that makes sense for your own forecast problem.

We will evaluate the skill of the model with the previous 1 year of data, then 2 years, all the way back through the 8 available years of historical data.

A year is a good interval to test for this dataset given the seasonal nature of the data, but other intervals could be tested, such as month-wise or multi-year intervals.

The snippet below shows how we can step backwards by year and cumulatively select all available observations.

For example

- Test 1: All data in 1989
- Test 2: All data in 1988 to 1989

And so on.

```
1  # split
2  train, test = differenced[differenced.index < '1990'], differenced['1990']
3  years = ['1989', '1988', '1987', '1986', '1985', '1984', '1983', '1982']
4  for year in years:
5      # select data from 'year' cumulative to 1989
6      dataset = train[train.index >= year]
```

The next step is to evaluate an ARIMA model.

We will use walk-forward validation. This means that a model will be constructed on the selected his[...]
1990). The real observation for that time step will be added to the history, a new model constructed, [...]

The forecasts will be collected together and compared to the final year of observations to give an err[...]
scores and will be in the same scale as the observations themselves.

**Get Your Start in Machine Learning** ✕

You can master applied Machine Learning **without the math or fancy degree**.
Find out how in this *free* and *practical* email course.

Email Address

START MY EMAIL COURSE

```
1   # walk forward over time steps in test
2   values = dataset.values
3   history = [values[i] for i in range(len(values))]
4   predictions = list()
5   test_values = test.values
6   for t in range(len(test_values)):
7       # fit model
8       model = ARIMA(history, order=(7,0,0))
9       model_fit = model.fit(trend='nc', disp=0)
10      # make prediction
11      yhat = model_fit.forecast()[0]
12      predictions.append(yhat)
13      history.append(test_values[t])
```

Get Your Start in Machine Learning

```
14  rmse = sqrt(mean_squared_error(test_values, predictions))
15  print('%s-%s (%d values) RMSE: %.3f' % (years[0], year, len(values), rmse))
```

Putting this all together, the complete example is listed below.

```
1   # fit an ARIMA model
2   from pandas import Series
3   from matplotlib import pyplot
4   from statsmodels.tsa.arima_model import ARIMA
5   from sklearn.metrics import mean_squared_error
6   from math import sqrt
7   # load dataset
8   series = Series.from_csv('daily-minimum-temperatures.csv', header=0)
9   # seasonal difference
10  differenced = series.diff(365)
11  # trim off the first year of empty data
12  differenced = differenced[365:]
13  # split
14  train, test = differenced[differenced.index < '1990'], differenced['1990']
15  years = ['1989', '1988', '1987', '1986', '1985', '1984', '1983', '1982']
16  for year in years:
17      # select data from 'year' cumulative to 1989
18      dataset = train[train.index >= year]
19      # walk forward over time steps in test
20      values = dataset.values
21      history = [values[i] for i in range(len(values))]
22      predictions = list()
23      test_values = test.values
24      for t in range(len(test_values)):
25          # fit model
26          model = ARIMA(history, order=(7,0,0))
27          model_fit = model.fit(trend='nc', disp=0)
28          # make prediction
29          yhat = model_fit.forecast()[0]
30          predictions.append(yhat)
31          history.append(test_values[t])
32      rmse = sqrt(mean_squared_error(test_values, predictions))
33      print('%s-%s (%d values) RMSE: %.3f' % (years[0], year, len(values), rmse))
```

Running the example prints the interval of history, number of observations in the history, and the RMSE skill of the model trained with that history.

The example does take awhile to run as 365 ARIMA models are created for each cumulative interval of historic training data.

```
1   1989-1989 (365 values) RMSE: 3.120
2   1989-1988 (730 values) RMSE: 3.109
```

```
3  1989-1987 (1095 values) RMSE: 3.104
4  1989-1986 (1460 values) RMSE: 3.108
5  1989-1985 (1825 values) RMSE: 3.107
6  1989-1984 (2190 values) RMSE: 3.103
7  1989-1983 (2555 values) RMSE: 3.099
8  1989-1982 (2920 values) RMSE: 3.096
```
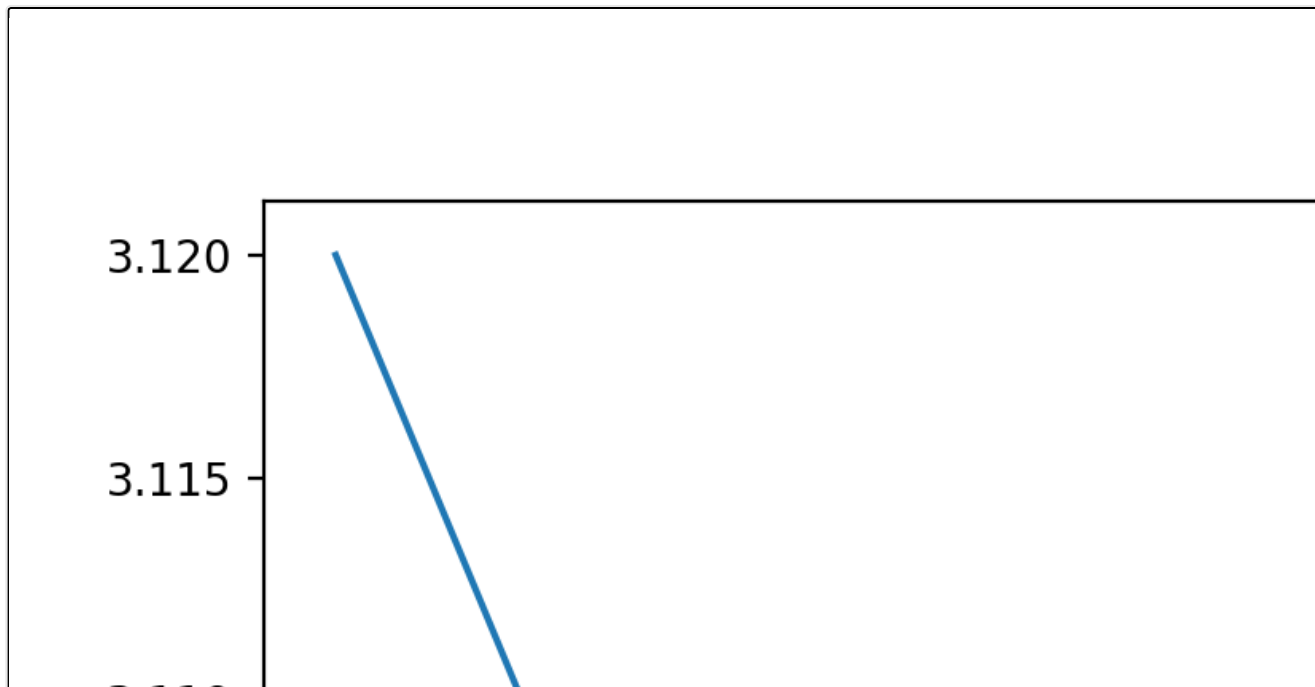
The results show that as the size of the available history is increased, there is a decrease in model error, but the trend is not purely linear.

We do see that there may be a point of diminishing returns at 2-3 years. Knowing that you can use fewer years of data is useful in domains where data availability or long model training time is an issue.

We can plot the relationship between ARIMA model error and the number of training observations.

```
1  from matplotlib import pyplot
2  x = [365, 730, 1095, 1460, 1825, 2190, 2555, 2920]
3  y = [3.120, 3.109, 3.104, 3.108, 3.107, 3.103, 3.099, 3.096]
4  pyplot.plot(x, y)
5  pyplot.show()
```

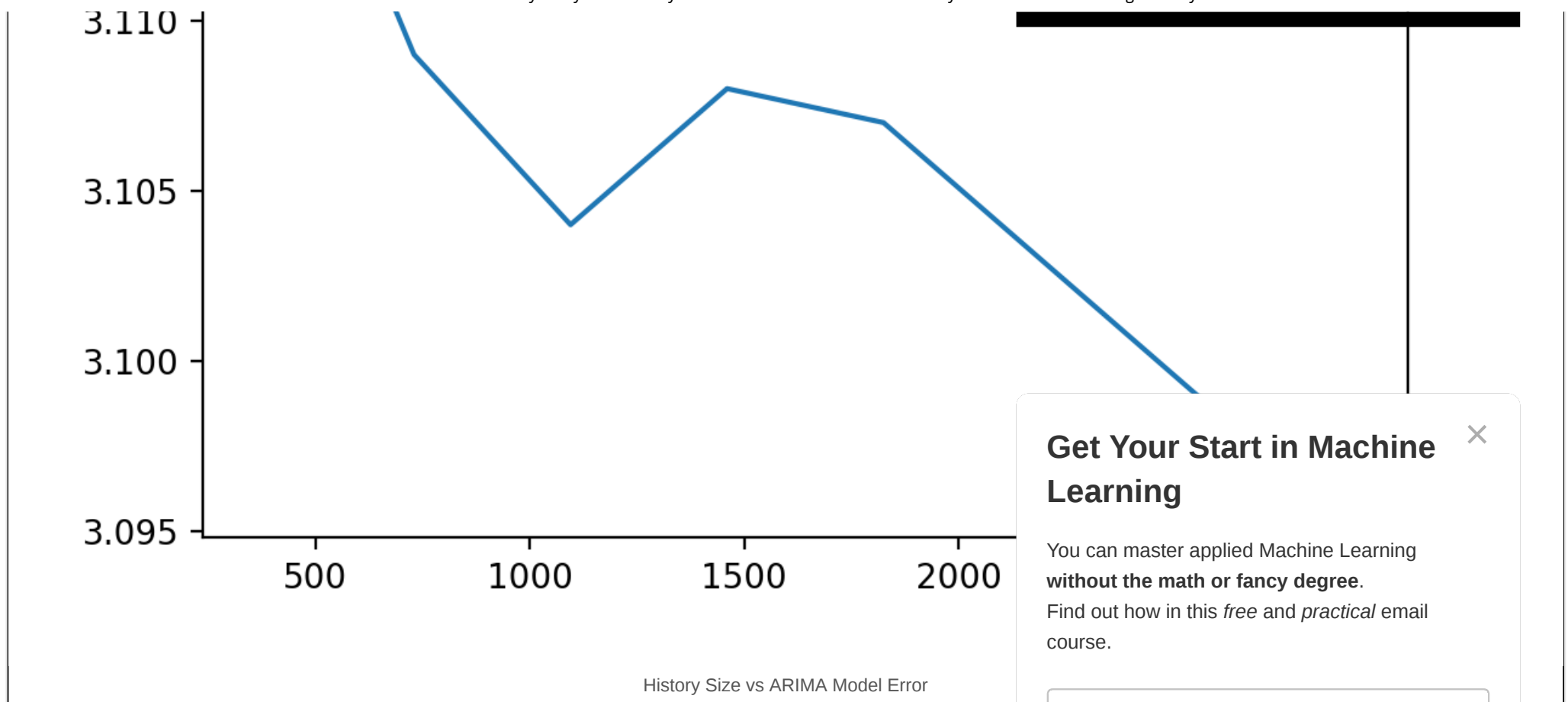Running the example creates a plot that almost shows a linear trend down in error as training sampl

**Get Your Start in Machine Learning**

You can master applied Machine Learning **without the math or fancy degree**.
Find out how in this *free* and *practical* email course.

Email Address

START MY EMAIL COURSE



Get Your Start in Machine Learning

History Size vs ARIMA Model Error

This is generally expected, as more historical data means that the coefficients may be better optimiz[...] om more years of data, for the most part.

There is also a counter-intuition. One may expect the performance of the model to increase with more history, as the data from the most recent years may be more like the data next year. This intuition is perhaps more valid in domains subjected to greater concept drift.

## Extensions

This section discusses limitations and extensions to the sensitivity analysis.

**Get Your Start in Machine Learning**

You can master applied Machine Learning **without the math or fancy degree**.
Find out how in this *free* and *practical* email course.

Email Address

START MY EMAIL COURSE

- **Untuned Model**. The ARIMA model used in the example is by no means tuned to the problem. Ideally, a sensitivity analysis of the size of training history would be performed with an already tuned ARIMA model or a model tuned to each case.
- **Statistical Significance**. It is not clear whether the difference in model skill is statistically significant. Pairwise statistical significance tests can be used to tease out whether differences in RMSE are meaningful.
- **Alternate Models**. The ARIMA uses historical data to fit coefficients. Other models may use the increasing historical data in other ways. Alternate nonlinear machine learning models may be investigated.
- **Alternate Intervals**. A year was chosen to joint the historical data, but other intervals may be used. A good interval might be weeks or months within one or two years of historical data for this dataset, as the extreme recency may bias the coefficients in useful ways.

## Summary

In this tutorial, you discovered how you can design, execute, and analyze a sensitivity analysis of the amount of history used to fit a time series forecast model.

Do you have any questions?
Ask your questions in the comments and I'll do my best to answer.

---

## Want to Develop Time Series Forecasts

### Develop Your Own Forecasts in Minutes

...with just a few lines of python code

Discover how in my new Ebook:
[Introduction to Time Series Forecasting With Python](#)

It covers **self-study tutorials** and **end-to-end projects** on topics like:
*Loading data, visualization, modeling, algorithm tuning,* and much more...

### Finally Bring Time Series Forecasting to Your Own Projects
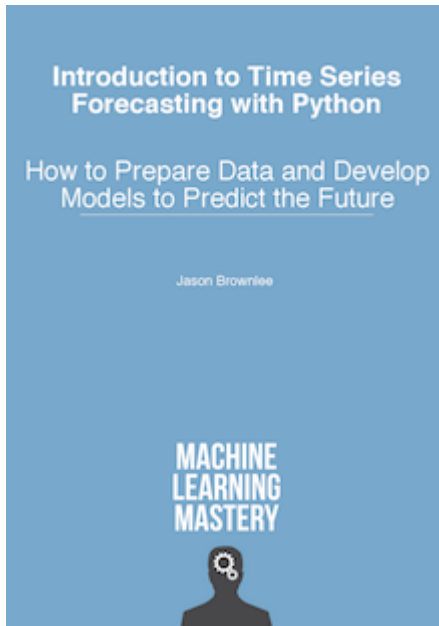
Skip the Academics. Just Results.

**Get Your Start in Machine Learning** ✕

You can master applied Machine Learning **without the math or fancy degree**.
Find out how in this *free* and *practical* email course.

Email Address

**START MY EMAIL COURSE**

**Get Your Start in Machine Learning**

Click to learn more.

### About Jason Brownlee

Dr. Jason Brownlee is a husband, proud father, academic researcher, author, professional devel[...] to helping developers get started and get good at applied machine learning. Learn more.

View all posts by Jason Brownlee →

## Get Your Start in Machine Learning

You can master applied Machine Learning **without the math or fancy degree**.
Find out how in this *free* and *practical* email course.

Email Address

**START MY EMAIL COURSE**

‹ How to Make Out-of-Sample Forecasts with ARIMA in Python

Feature Selection for Time Series Forecasting with Python ›

Get Your Start in Machine Learning

# 8 Responses to *Sensitivity Analysis of History Size to Forecast Skill with ARIMA in Python*

**sura** April 7, 2017 at 3:39 am #                                                     REPLY ↩

thank you ! but, can you give me a download address of dataset? i want to try again!

thank!

**Jason Brownlee** April 9, 2017 at 2:46 pm #                                          REPLY ↩

Sorry, here it is:

https://datamarket.com/data/set/2324/daily-minimum-temperatures-in-melbourne-australia-1981-199

**Eugeniy** May 5, 2017 at 6:44 pm #

Good afternoon!
Thank you for your article.
Tell me please, if there are more formal and mathematical definitions of sensitivity analysis of history size
Or is this usually only an experimental way of determining?
Thank you!

**Jason Brownlee** May 6, 2017 at 7:40 am #                                            REPLY ↩

I'm sure you can analyze the effect of history size on the model analytically.

A sensitivity analysis seeks to answer the question empirically.

---

**Get Your Start in Machine Learning**                                        ✕

You can master applied Machine Learning **without the math or fancy degree**.
Find out how in this *free* and *practical* email course.

| Email Address |

**START MY EMAIL COURSE**

---

Get Your Start in Machine Learning

**David Ravnsborg** August 18, 2017 at 9:36 am #

REPLY ↩

Why do you declare "differenced" and then immediately write over it without using it?

**Jason Brownlee** August 18, 2017 at 4:37 pm #

REPLY ↩

Good question, that looks like a bug to me. I'll add a note to trello to fix it up.

**David Ravnsborg** August 18, 2017 at 5:48 pm #

I looked into it a little further on my end. I think it was just a typo where:

```
# seasonal difference
differenced = series.diff(365)
# trim off the first year of empty data
differenced = series[365:]
```

should have been…

```
# seasonal difference
differenced = series.diff(365)
# trim off the first year of empty data
differenced = differenced[365:]
```

But it completely changes the results for the worse 🙁 Any chance you could cover this? It would make a great follow-up. Here are the results I get:
model.py:496: ConvergenceWarning: Maximum Likelihood optimization failed to converge. Check mle_retvals
"Check mle_retvals", ConvergenceWarning)
1989-1989 (365 values) RMSE: 3.120
1989-1988 (730 values) RMSE: 3.109
model.py:496: ConvergenceWarning: Maximum Likelihood optimization failed to converge. Check mle_retvals
"Check mle_retvals", ConvergenceWarning)
1989-1987 (1095 values) RMSE: 3.104

## Get Your Start in Machine Learning

You can master applied Machine Learning **without the math or fancy degree**.
Find out how in this *free* and *practical* email course.

Email Address

**START MY EMAIL COURSE**

Get Your Start in Machine Learning

1989-1986 (1460 values) RMSE: 3.108

1989-1985 (1825 values) RMSE: 3.107

model.py:496: ConvergenceWarning: Maximum Likelihood optimization failed to converge. Check mle_retvals

"Check mle_retvals", ConvergenceWarning)

1989-1984 (2190 values) RMSE: 3.103

1989-1983 (2555 values) RMSE: 3.099

1989-1982 (2920 values) RMSE: 3.096

**Jason Brownlee** August 25, 2017 at 3:15 pm #

REPLY ↰

I have updated the post, thanks again David.

We still see the same linearly downward trend in error.

Remember that the RMSE scores are in fact in the units of seasonally differenced temperatu

If you're interested in better results, you can try using a grid search on the ARIMA paramete
whether performing a seasonal difference results in better final RMSE on this problem.

## Leave a Reply

Name (required)

Email (will not be published) (required)

Website

SUBMIT COMMENT

**Welcome to Machine Learning Mastery**

Hi, I'm Dr. Jason Brownlee.
My goal is to make practitioners like YOU awesome at applied machine learning.

Read More

**Get Good at Time Series Forecasting**

Need visualizations and forecast models?
Looking for step-by-step tutorials?
Want end-to-end projects?

Get Started with Time Series Forecasting in Python!
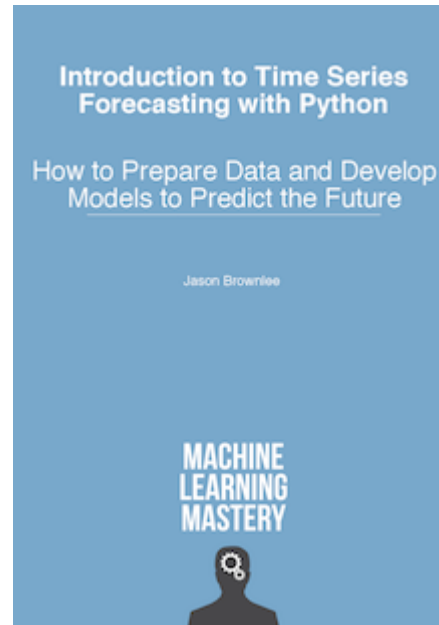
**Get Your Start in Machine Learning**

You can master applied Machine Learning
**without the math or fancy degree**.
Find out how in this *free* and *practical* email course.

Email Address

START MY EMAIL COURSE

Introduction to Time Series
Forecasting with Python

How to Prepare Data and Develop
Models to Predict the Future

Jason Brownlee

MACHINE
LEARNING
MASTERY

POPULAR

**Time Series Prediction with LSTM Recurrent Neural Networks in Python with Keras**
JULY 21, 2016

**Your First Machine Learning Project in Python Step-By-Step**
JUNE 10, 2016

**Develop Your First Neural Network in Python With Keras Step-By-Step**
MAY 24, 2016

**Sequence Classification with LSTM Recurrent Neural Networks in Python with Keras**
JULY 26, 2016

**How to Setup a Python Environment for Machine Learning and Deep Learning with Anaconda**
MARCH 13, 2017

# Get Your Start in Machine Learning

You can master applied Machine Learning
**without the math or fancy degree**.
Find out how in this *free* and *practical* email
course.

Email Address

START MY EMAIL COURSE

**Get Your Start in Machine Learning**

**Time Series Forecasting with the Long Short-Term Memory Network in Python**

APRIL 7, 2017

**Multi-Class Classification Tutorial with the Keras Deep Learning Library**

JUNE 2, 2016

**Regression Tutorial with the Keras Deep Learning Library in Python**

JUNE 9, 2016

**Multivariate Time Series Forecasting with LSTMs in Keras**

AUGUST 14, 2017

**How to Implement the Backpropagation Algorithm From Scratch In Python**

NOVEMBER 7, 2016

Privacy | Contact | About

## Get Your Start in Machine Learning

You can master applied Machine Learning **without the math or fancy degree**.
Find out how in this *free* and *practical* email course.

Email Address

**START MY EMAIL COURSE**

**Get Your Start in Machine Learning**