

Estimating Monthly Power Prices Using Time Series Models

Vladas Aleusius

IBM Machine Learning Professional Certificate

Specialized Models: Time Series And Survival Analysis

Coursera

2022-01-03

Main Objective

The objective of this analysis is to find out which time series prediction model does the best job predicting monthly energy prices for 2020 using the prices from 2000 to 2019 as the training dataset. Our focus is to select the model with the best *accuracy*, which will be assessed with Root Mean Squared Error (RMSE) as a metric. The models will be Exponential Smoothing, ARIMA and Simple Recurrent Neural Network (Simple RNN).

For Exponential Smoothing and ARIMA, 2000-2019 prices will be used as the training dataset. However, the prices from 2000 to 2018 will need to be used to train the Simple RNN model.

In this study, Exponential Smoothing and ARIMA models will include seasonality assumptions, in particular - that monthly prices have a 12-month seasonality. However, as we will see below, there is no clear trend pattern, so the seasonality pattern might be more complex than we assume.

Data Description

According to Wikipedia, "Nord Pool AS is a European power exchange owned by Euronext and the continental Nordic and Baltic countries' Transmission system operators (TSOs). Nord Pool delivers power trading across Europe. Nord Pool offers day-ahead and intraday trading, clearing and settlement, data and compliance, as well as consultancy services. More than 360 countries trade on Nord Pool today." Nord Pool website will be used as a source for historic System (SYS) monthly energy prices (euros per megawatt hour - EUR/MWh) that will be used to forecast monthly prices for 2021.

We will start with importing necessary modules and functions:

```
In [2]: import pandas as pd
import seaborn as sns
from datetime import datetime
from scipy.stats import normtest, norm, uniform, boxcox
from pandas.api.types import is_datetime64_any_dtype as is_datetime
import numpy as np
import matplotlib.pyplot as plt
import statsmodels.api as sm
from statsmodels.tsa.statespace.sarimax import SARIMAX
import matplotlib.pyplot as plt
import numpy as np
from statsmodels.tsa.stattools import adfuller
from statsmodels.tsa.api import ExponentialSmoothing
from sklearn.metrics import mean_squared_error
from math import sqrt
import keras
from keras.models import Sequential
from keras.layers import Dense, SimpleRNN, LSTM, Activation, Dropout
```

Historic monthly prices from 2000 to 2020 will be extracted from Nord Pool website using the links below:

```
In [3]: price_links = [
    "https://www.nordpoolgroup.com/49bfe7/globalassets/marketedata-excel-files/elspot-prices_2000_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48e484/globalassets/marketedata-excel-files/elspot-prices_2001_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2002_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2003_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2004_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2005_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2006_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2007_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2008_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2009_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2010_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2011_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2012_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2013_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2014_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2015_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2016_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2017_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2018_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2019_monthly_eur.xls",
    "https://www.nordpoolgroup.com/48b947/globalassets/marketedata-excel-files/elspot-prices_2020_monthly_eur.xls"
]
```

Below, the data is split into training and testing subsets. For Exponential Smoothing and ARIMA models, the training data will contain values from 2000 to 2019, and the testing data - prices for 2020.

```
In [5]: price_train = pd.DataFrame()
price_test = pd.DataFrame()
for i in range(len(price_links)):
    prices = pd.read_html(price_links[i], header=2, decimal=',', thousands=',')[0]
    prices['Month'] = i * 12 + prices['Unnamed: 0'].str[1:2].str[1:2] + prices['Unnamed: 0'].str[3:4].str[1:2]
    prices.set_index('Month', inplace=True)
    prices.index = pd.to_datetime(prices.index)
    prices = prices['SYS']
    price_trainpd.concat([price_train, prices], axis=0)
    price_testpd.concat([price_test, prices], axis=0)
else:
    price_testpd.concat([price_test, prices], axis=0)
```

The final training dataset consists of a month index and SYS prices from 2000 to 2019:

In [4]: price_train

Out [4]:

SYS

Month

2000-01-01 16.22

2000-02-01 12.89

2000-03-01 11.78

2000-04-01 9.80

2000-05-01 12.51

...

2019-08-01 36.11

2019-09-01 32.92

2019-10-01 37.10

2019-11-01 42.15

2019-12-01 36.79

240 rows x 1 columns

Data Exploration

As we can see from the plot and summary tables below, monthly energy prices fluctuate a lot, from less than 10 EUR/MWh in 2000 to more than 80 EUR/MWh in 2010. The average monthly price is about 34 EUR/MWh, and the standard deviation is about 13 EUR/MWh:

In [5]: price_train.describe()

Out [5]:

SYS

count 240.000000

mean 33.679708

std 12.877149

min 6.350000

25% 25.317500

50% 31.690000

75% 40.550000

max 81.650000

In [6]: print("Minimum Price")

price_train.loc[price_train.idxmin(), :]

Minimum Price

Out [6]:

SYS

Month

2000-07-01 6.35

In [7]: print("Maximum Price")

price_train.loc[price_train.idxmax(), :]

Maximum Price

Out [7]:

SYS

Month

2010-12-01 81.65

The plot below shows sharp price swings throughout the years:

In [22]: sns.lineplot(data=price_train, y='SYS', x=price_train.index).set_title("Monthly Prices")

Out [22]:

Text(0.5, 1.0, 'Monthly Prices')

Monthly Prices

Y-axis: SYS

X-axis: Month

2000 2004 2008 2012 2016 2020

Also, average price levels change throughout the years, which might indicate a need to difference data to make it stationary. The plot below shows that average price levels are more constant for the differenced data, although extremely high and low values for some months still remain:

In [23]: sns.lineplot(data=price_train.diff(), y='SYS', x=price_train.index).set_title("Monthly Price Non-Seasonal Differ")

Out [23]:

Text(0.5, 1.0, 'Monthly Price Non-Seasonal Differences')

Monthly Price Non-Seasonal Differences

Y-axis: SYS

X-axis: Month

2000 2004 2008 2012 2016 2020

Monthly price differences fluctuate around 0, but the highest absolute differences are for increases:

In [10]: price_train.diff().mean()

Out [10]:

SYS 0.080667

dtype: float64

The violin plot below demonstrates that although the highest price values are for winter months, typical price values are not much different throughout the year. So our initial assumption of 12-month seasonality might be put into question.

In [12]: # violinplot of months to determine variance and range

violinplot(price_train.index.month, yprice_train.SYS)

plt.grid(b=True)

violinplot

Y-axis: SYS

X-axis: Month

1 2 3 4 5 6 7 8 9 10 11 12

Below, autocorrelation and partial autocorrelation functions are plotted for the training dataset. Autocorrelation increases for lags 20-30 again confirm that the initial data is not stationary.

In [34]: print("Monthly Data Autocorrelation Plots")

Autocorrelation and Partial Autocorrelation Functions for Monthly Data

acf_plot_2 = plot_acf(price_train['SYS'], lags=50, title='Autocorrelation in Power Monthly Price Data')

pacf_plot_2 = plot_pacf(price_train['SYS'], lags=50, title='Partial Autocorrelation in Power Monthly Price Data')

Monthly Data Autocorrelation Plots

Autocorrelation in Power Monthly Price Data

Y-axis: SYS

X-axis: Lag

0 10 20 30 40 50

Partial Autocorrelation in Power Monthly Price Data

Y-axis: SYS

X-axis: Lag

0 10 20 30 40 50

The plots below show that for non-seasonally differenced data, both ACF and PACF values are small and patternless. This is a sign of not over-differencing. However, the plots below do not help to select the best number of either AR or MA components.

In [25]: print("Monthly Data diff Autocorrelation Plots")

Autocorrelation and Partial Autocorrelation Functions for Monthly Data

acf_plot_2 = plot_acf(price_train['SYS'].diff(), lags=50, title='Autocorrelation in Power Monthly Price Data')

pacf_plot_2 = plot_pacf(price_train['SYS'].diff(), lags=50, title='Partial Autocorrelation in Power Monthly Price Data')

Monthly Data diff Autocorrelation Plots

Autocorrelation in Power Monthly Price Diffs

Y-axis: SYS

X-axis: Lag

0 10 20 30 40 50

Partial Autocorrelation in Power Monthly Price Diffs

Y-axis: SYS

X-axis: Lag

0 10 20 30 40 50

Augmented Dickey-Fuller test confirms that monthly prices are stationary:

In [13]: adf, pvalue, usedlag, nobs, critical_values, icbest = adfuller(price_train['SYS'].diff().dropna())

Out [13]:

print(pvalue)

2.86955266287494e-12

Seasonal differencing seems to be unnecessary, since it produces strongly negative autocorrelation values (which indicates over-differencing):

In [26]: print("Monthly Data Seasonal diff Autocorrelation Plots")

Autocorrelation and Partial Autocorrelation Functions for Monthly Data

acf_plot_2 = plot_acf(price_train['SYS'].diff().dropna(), lags=50, title='Autocorrelation in Power Monthly Price Data')

pacf_plot_2 = plot_pacf(price_train['SYS'].diff().dropna(), lags=50, title='Partial Autocorrelation in Power Monthly Price Data')

Monthly Data Seasonal diff Autocorrelation Plots

Autocorrelation in Power Monthly Price Seasonal Diffs

Y-axis: SYS

X-axis: Lag

0 10 20 30 40 50

Partial Autocorrelation in Power Monthly Price Seasonal Diffs

Y-axis: SYS

X-axis: Lag

0 10 20 30 40 50

In [27]: print("Monthly Data Autocorrelation Plots")

Autocorrelation and Partial Autocorrelation Functions for Monthly Data

acf_plot_2 = plot_acf(price_train['SYS'].diff().dropna(), lags=50, title='Autocorrelation in Power Monthly Price Data')

pacf_plot_2 = plot_pacf(price_train['SYS'].diff().dropna(), lags=50, title='Partial Autocorrelation in Power Monthly Price Data')

Monthly Data Autocorrelation Plots

Autocorrelation in Power Monthly Price Seasonal and Non-Seasonal Diffs

Y-axis: SYS

X-axis: Lag

0 10 20 30 40 50

Partial Autocorrelation in Power Monthly Price Seasonal and Non-Seasonal Diffs

Y-axis: SYS

X-axis: Lag

0 10 20 30 40 50

The distribution of prices is approximately bell shaped, although somewhat skewed:

In [34]: price_train['SYS'].hist()

Out [34]:

Text(0.5, 1.0, 'Distribution Of Monthly Prices')

Distribution Of Monthly Prices

Y-axis: SYS

X-axis: Price

10 20 30 40 50 60 70 80

Differences are distributed more symmetrically, however:

In [35]: price_train['SYS'].diff().hist()

Out [35]:

Text(0.5, 1.0, 'Distribution Of Monthly Price Differences')

Distribution Of Monthly Price Differences

Y-axis: SYS

X-axis: Price

-20 -10 0 10 20 30

The average residual value is still almost 0:

In [39]: pd.Series(estimated_residual).mean()

Out [39]:

0.00705043859493289

The model is fit below - we assume an additive model for both trend and seasonal components, and 12-month seasonality:

In [10]: triple = ExponentialSmoothing(price_train['SYS'],

seasonal='additive',

seasonal_periods=12, freq='MS').fit(optimized=True)

triple_preds = triple.forecast(triple.index)

triple_preds_df = pd.DataFrame(triple_preds, columns=['SYS'], index=triple.index)

C:\Users\tenov\anaconda3\lib\site-packages\statsmodels\tsa\holtwinters\model.py:524: FutureWarning: After 0.13 initialization method will be handled at model creation

warnings.warn("No frequency information was")

Model summary is printed below. Coefficient values show that the forecasts mainly depend on seasonal components, and the trend has only a minor impact. Seasonal component values show that middle year prices are expected to be lower:

In [11]: triple.summary()

Out [11]:

ExponentialSmoothing Model Results

Dep. Variable: SYS No. Observations: 240

Model: ExponentialSmoothing SSE 8937.826

Optimized: True AIC 950.178

Trend: Additive BIC 955.868

Seasonal: Additive AICC 903.273

Seasonal Periods: 12 Date: Sat, 01 Jan 2022

Box-Cox: False Time: 1714.12

Box-Cox Coeffs: None

coeff code optimized

smoothing_level 0.999389 alpha True

smoothing_trend 2.6875e-07 beta True

smoothing_seasonal 8.9648e-16 gamma True

initial_level 36.573169 l0 True

initial_trend 0.087389 b.0 True

initial_seasons 0 -20.52435 s.0 True

initial_seasons.1 -21.791752 s.1 True

initial_seasons.2 -24.371459 s.2 True

initial_seasons.3 -25.702214 s.3 True

initial_seasons.4 -29.048464 s.4 True

initial_seasons.5 -28.900650 s.5 True

initial_seasons.6 -29.590462 s.6 True

initial_seasons.7 -26.487169 s.7 True

initial_seasons.8 -25.499835 s.8 True

initial_seasons.9 -24.634499 s.9 True

initial_seasons.10 -22.008141 s.10 True

initial_seasons.11 -21.009110 s.11 True

Price predictions for 2020 are printed below:

In [12]: triple_preds_df

Out [12]:

SYS

Month

2020-01-01 37.362548

2020-02-01 36.182531

2020-03-01 33.690214

2020-04-01 32.428849

2020-05-01 29.187989

2020-06-01 29.423193

2020-07-01 28.820771

2020-08-01 32.011453

2020-09-01 33.086177

2020-10-01 34.038903

2020-11-01 36.752551

2020-12-01 37.839072

The plot below compares testing subset prices and price predictions - the latter are much higher, which shows that the model has done not a very good job in predicting:

In [33]: plt.plot(price_train.index, price_train['SYS'], 'b--', label='train')

plt.plot(price_train.index, price_test, 'r--', label='test')

plt.plot(price_train.index, triple_preds_df['SYS'], 'b--', label='predictions')

plt.legend(loc='upper left')

plt.title("Power Monthly Prices")

plt.grid(alpha=0.3)

plt.show()

Power Monthly Prices

Y-axis: SYS

X-axis: Month

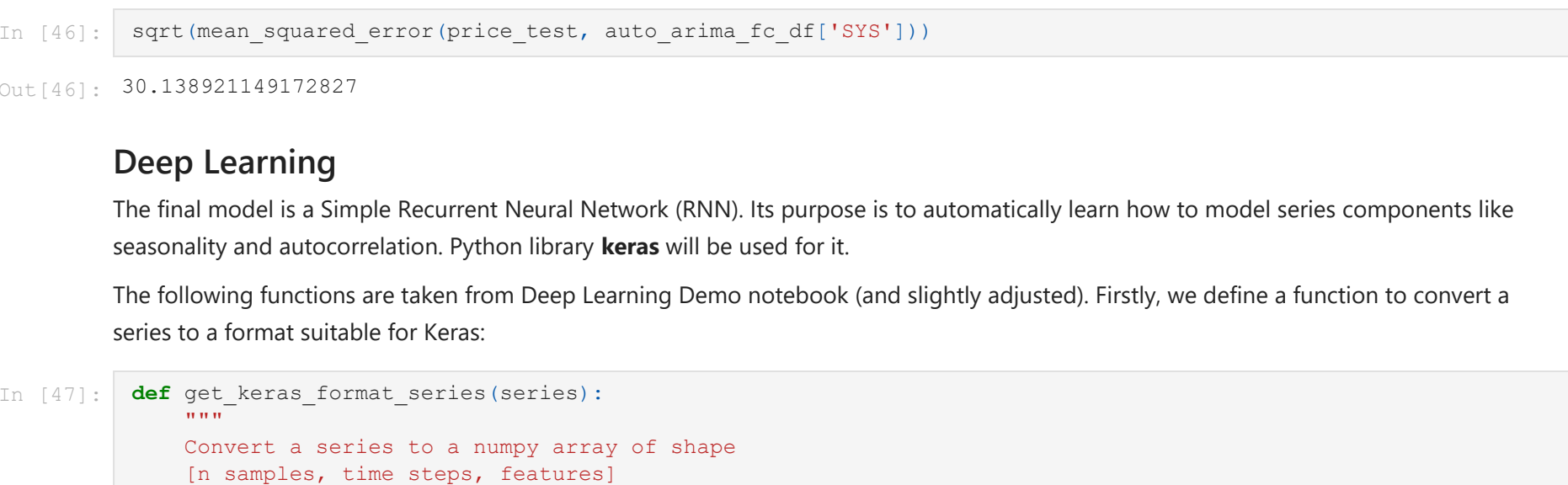
2000 2004 2008 2012 2016 2020

Root Mean Squared Error (RMSE


```
In [42]: auto_arima_fc = setuse_model.predict(len(price_test))
Out[43]: auto_arima_fc_df = DataFrame(auto_arima_fc, index=price_test.index, columns=["SYS"])
```

SYS	
Month	
2020-01-01	35.933251
2020-02-01	37.964997
2020-03-01	39.250499
2020-04-01	38.907914
2020-05-01	38.572125
2020-06-01	41.158908
2020-07-01	41.870328
2020-08-01	41.977700
2020-09-01	41.919301
2020-10-01	40.717496
2020-11-01	41.137459
2020-12-01	42.286685

Test subset price values and ARIMA predictions are compared below. As for Exponential Smoothing, predictions are much higher than actual 2020 prices. However, Exponential Smoothing predicted lower mid-year values (the same price fall was actually observed), while ARIMA predictions tend to moderately increase throughout the year:



RMSE is even higher for ARIMA than Exponential Smoothing and exceeds 30. This shows that selected ARIMA model predictive power is even lower than for Exponential Smoothing:

```
In [46]: sqrt(mean_squared_error(price_test, auto_arima_fc_df['SYS']))
Out[46]: 30.139921149172827
```

Deep Learning

The final model is a Simple Recurrent Neural Network (RNN). Its purpose is to automatically learn how to model series components like seasonality and autocorrelation. Python library **keras** will be used for it.

The following functions are taken from Deep Learning Demo notebook (and slightly adjusted). Firstly, we define a function to convert a series to a format suitable for Keras:

```
In [47]: def get_keras_format_series(series):
    """
    Convert a series to a numpy array of shape
    [n_samples, time_steps, features]
    """
    series = np.array(series)
    return series.reshape(series.shape[0], series.shape[1], 1)
```

The below function splits the total dataset into training and testing subsets. For the training subset, many training samples will be generated and added to train_X list, each sample of a selected length (in our case, they will be of length 12), and the following price value will be added to train_y list. A sample gap is used to make the samples not follow each other. For the testing subset, the same number of first values (12 in our case) will be used as X-s and used for the model based prediction, and the rest - as y-s.

```
In [48]: def get_train_test_data(df, series_name, input_months,
    test_months, sample_gap=3):
    """
    Utility processing function that splits a monthly time series into
    train and test with keras-friendly format, according to user-specified
    choice of shape.

    arguments
    -----
    df (dataframe): dataframe with time series columns
    series_name (string): column name in df
    input_months (int): length of sequence input to network
    test_months (int): length of held-out terminal sequence
    sample_gap (int): step size between start of train sequences; default 3

    returns
    -----
    tuple: train_X, test_X_init, train_y, test_y
    """
    train = df[series_name][:test_months] # training data is remaining months until amount of test_months
    test = df[series_name][test_months:] # test data is the remaining test_months

    train_X, train_y = [], []

    # range 0 through # of train samples - input_months by sample_gap.
    for i in range(0, train.shape[0]-input_months, sample_gap):
        train_X.append(train[i:i+input_months]) # each training sample is of length input_months
        train_y.append(train[i+input_months]) # each y is just the next step after training sample

    train_X = get_keras_format_series(train_X) # format our new training set to keras format
    train_y = np.array(train_y) # make sure y is an array to work properly with keras

    # The set that we held out for testing (must be same length as original train input)
    test_X_init = test[:input_months]
    test_y = test[input_months:] # test_y is remaining values from test set

    return train_X, test_X_init, train_y, test_y
```

As already mentioned, training sequences will be of 12 month data each, and the testing subset will consist of the last two years. The testing subset will again be split into two parts, 12 months each, the first of which will be used as input:

```
In [49]: input_months = 12
test_months = 24
train_X, test_X_init, train_y, test_y = \
    get_train_test_data(price_total, 'SYS', input_months, test_months)
```

A Simple RNN model with one layer is defined and fit below:

```
In [50]: def fit_simpleRNN(train_X, train_y, cell_units, epochs):
    """
    Fit Simple RNN to data train_X, train_y

    arguments
    -----
    train_X (array): input sequence samples for training
    train_y (list): next step in sequence targets
    cell_units (int): number of hidden units for RNN cells
    epochs (int): number of training epochs

    """
    # initialize model
    model = Sequential()

    # construct an RNN layer with specified number of hidden units
    # per cell and desired sequence input format
    model.add(SimpleRNN(cell_units, input_shape=train_X.shape[1],1))

    # add an output layer to make final predictions
    model.add(Dense(1))

    # define the loss function / optimization strategy, and fit
    # the model with the desired number of passes over the data (epochs)
    model.compile(loss='mean_squared_error', optimizer='adam')
    model.fit(train_X, train_y, epochs=epochs, verbose=0)

    return model
```

```
In [51]: model = fit_simpleRNN(train_X, train_y, cell_units=50, epochs=50000)
```

The model above was trained to predict only one future time step. For multi-step forecasting, we will iteratively generate one prediction, append it to the end of the input sequence (and shift that sequence forward by one step), then feed the new sequence back to the model.

We stop once we have generated all the time step predictions we need. In our case, the first 12 values of testing subset will be used as the initial set, and predictions will be generated for the rest 12 months one by one, shifting the input sequence.

```
In [52]: def predict(X_init, n_steps, model):
    """
    Given an input series matching the model's expected format,
    generates model's predictions for next n_steps in the series
    """
    X_init = np.array(X_init).copy().reshape(1,-1,1)
    preds = []

    # iteratively take current input sequence, generate next step pred,
    # and shift input sequence forward by a step (to end with latest pred).
    for i in range(n_steps):
        pred = model.predict(X_init)
        preds.append(pred)
        X_init[-1,-1,:] = X_init[-1,1,1] # replace first 11 values with 2nd through 12th
        X_init[-1,1,1] = pred # replace 12th value with prediction

    preds = np.array(preds).reshape(-1,1)

    return preds
```

```
In [53]: y_preds_rnn = predict(test_X_init, len(test_y), model)
```

Finally, a comparison of RNN predictions and actual data is plotted below. True values are much lower than the predictions. Furthermore, RNN predicted a price increase for the mid-year when the price actually fell:

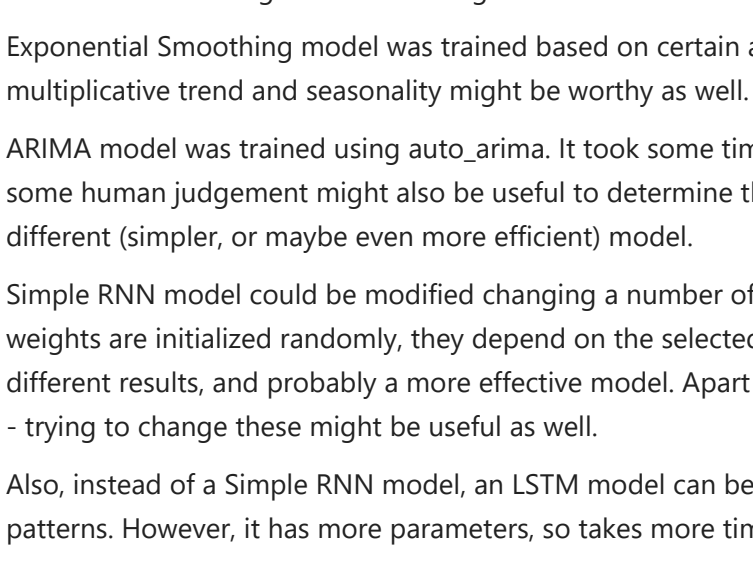
```
In [59]: def predict_and_plot(title):
    """
    Plots model's predictions against the ground truth

    arguments
    -----
    title (string): plot title

    """
    #using our ranges we plot X_init
    plt.plot(test_X_init.index, test_X_init, color='blue', linestyle="--")
    #and test and actual prices
    plt.plot(test_y.index, test_y, color='orange', linestyle="--")
    plt.plot(test_y.index, y_preds_rnn, color='red', linestyle="--")

    plt.title(title)
    plt.legend(['train', 'test', 'predictions'])
    plt.xticks(rotation=45)
```

```
In [70]: predict_and_plot('PM Series: Test Data and Simple RNN Predictions')
```



RMSE for RNN predictions exceeds 33 and is the worst among examined models:

```
In [71]: sqrt(mean_squared_error(test_y, y_preds_rnn))
Out[71]: 33.59419182852007
```

The model summary below shows that it had more than 2000 parameters, so still not very complex:

```
In [57]: model.summary()
Model: "sequential"
Layer (type) Output Shape Param #
-----
simple_rnn (SimpleRNN) (None, 50) 2600
Dense (Dense) (None, 1) 51
-----
Total params: 2,651
Trainable params: 2,651
Non-trainable params: 0
```

Final Recommendation

In this analysis, we tried running three time series prediction models for monthly energy price forecast: Exponential Smoothing, ARIMA and Recurrent Neural Network (RNN). Exponential Smoothing forecasts are the most in line with the actual values, and the model took the least time to train. RNN took the longest to model and provided the worst results in terms of RMSE, so it is the worst model of the three. ARIMA is in between in both regards.

So the best model in our case is **Exponential Smoothing**.

Summary Key Findings and Insights

As we discussed above, Exponential Smoothing forecast was the best for monthly energy price forecasting. Its RMSE is the lowest - almost 23. Also, it took the shortest to train - Auto ARIMA tries many models to select the best one, and RNN takes some time depending on numbers of cells and epochs.

The plot below compares the forecasts obtained with all three prediction methods and actual 2020 monthly price values. As we can see, actual data is much lower than for predictions. However, Exponential Smoothing values tend to be the lowest (and thus the closest to the actual data), and they go down in the middle of year - the same as for the actual data. Other forecasting methods failed to predict this. So this is one more argument for Exponential Smoothing usage.

ARIMA forecast is pretty flat throughout 2020, and RNN forecast increases in the middle of year - the opposite direction than it should be.



However, none of the methods managed to predict price fall for 2020, compared to previous years. This might be due to a multi-year cyclical pattern that forecasting did not include, or simply because of prices changing randomly.

Next Steps

Some possible model improvement steps are described below.

We made an assumption about yearly seasonality, but this did not help to forecast a price downturn for 2020. Adding a cyclical component can be useful for long-term forecasting.

Exponential Smoothing model was trained based on certain assumptions, such as additive trend and seasonality. Experimenting with multiplicative trend and seasonality might be worthy as well.

ARIMA model was trained using auto.arima. It took some time, since many different models were tried to find the best one. However, some human judgement might also be useful to determine the final model. So relying on experimentation might help to come up with a different (simpler, or maybe even more efficient) model.

Simple RNN model could be modified changing a number of cell units and using more training epochs. Furthermore, since RNN model weights are initialized randomly, they depend on the selected random seed in our case. Trying a different random seed might yield different results, and probably a more effective model. Apart from that, we used default parameters (the activation function, initializers etc.) - trying to change these might be useful as well.

Also, instead of a Simple RNN model, an LSTM model can be used - experience shows that they are better at predicting more complex data patterns. However, it has more parameters, so takes more time to train.