

01_tsa_and_stationarity

September 29, 2021

1 Time Series Analysis and Stationarity

```
[1]: import warnings
warnings.filterwarnings('ignore')
```

```
[2]: %matplotlib inline

import pandas_datareader.data as web
import numpy as np

import statsmodels.tsa.api as tsa
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
from statsmodels.tsa.stattools import acf, q_stat, adfuller
from scipy.stats import probplot, moment

import matplotlib.pyplot as plt
import seaborn as sns
```

```
[3]: sns.set_style('whitegrid')
```

```
[4]: def plot_correlogram(x, lags=None, title=None):
    lags = min(10, int(len(x)/5)) if lags is None else lags
    with sns.axes_style('whitegrid'):
        fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(14, 8))
        x.plot(ax=axes[0][0], title='Residuals')
        x.rolling(21).mean().plot(ax=axes[0][0], c='k', lw=1)
        q_p = np.max(q_stat(acf(x, nlags=lags), len(x))[1])
        stats = f'Q-Stat: {np.max(q_p):>8.2f}\nADF: {adfuller(x)[1]:>11.2f}'
        axes[0][0].text(x=.02, y=.85, s=stats, transform=axes[0][0].transAxes)
        probplot(x, plot=axes[0][1])
        mean, var, skew, kurtosis = moment(x, moment=[1, 2, 3, 4])
        s = f'Mean: {mean:>12.2f}\nSD: {np.sqrt(var):>16.2f}\nSkew: {skew:12.
        ↳2f}\nKurtosis:{kurtosis:9.2f}'
        axes[0][1].text(x=.02, y=.75, s=s, transform=axes[0][1].transAxes)
        plot_acf(x=x, lags=lags, zero=False, ax=axes[1][0])
        plot_pacf(x, lags=lags, zero=False, ax=axes[1][1])
        axes[1][0].set_xlabel('Lag')
```

```

axes[1][1].set_xlabel('Lag')
fig.suptitle(title, fontsize=14)
sns.despine()
fig.tight_layout()
fig.subplots_adjust(top=.9)

```

1.1 Download Series

Load monthly industrial production and daily NASDAQ stock market index:

```

[5]: industrial_production = web.DataReader('IPGMFN', 'fred', '1988', '2017-12').
      ↪squeeze().dropna()
      nasdaq = web.DataReader('NASDAQCOM', 'fred', '1990', '2017-12-31').squeeze().
      ↪dropna()

```

1.2 Additive Decomposition

Time series data typically contains a mix of various patterns that can be decomposed into several components, each representing an underlying pattern category. In particular, time series often consist of the systematic components trend, seasonality and cycles, and unsystematic noise. These components can be combined in an additive, linear model, in particular when fluctuations do not depend on the level of the series, or in a non-linear, multiplicative model.

These components can be split up automatically. statsmodels includes a simple method to split the time series into a trend, seasonal, and residual component using moving averages. We can apply it to monthly data on industrial manufacturing production with both a strong trend and seasonality component, as follows:

```

[6]: components = tsa.seasonal_decompose(industrial_production, model='additive')

[7]: ts = (industrial_production.to_frame('Original')
      .assign(Trend=components.trend)
      .assign(Seasonality=components.seasonal)
      .assign(Residual=components.resid))
      with sns.axes_style('white'):
          ts.plot(subplots=True, figsize=(14, 8), title=['Original Series', 'Trend_
      ↪Component', 'Seasonal Component', 'Residuals'], legend=False)
          plt.suptitle('Seasonal Decomposition', fontsize=14)
          sns.despine()
          plt.tight_layout()
          plt.subplots_adjust(top=.91);

```



1.3 Time Series Stationarity

The statistical properties, such as the mean, variance, or autocorrelation, of a stationary time series are independent of the period, that is, they don't change over time. Hence, stationarity implies that a time series does not have a trend or seasonal effects and that descriptive statistics, such as the mean or the standard deviation, when computed for different rolling windows, are constant or do not change much over time. It reverts to its mean, and the deviations have constant amplitude, while short-term movements always look the same in the statistical sense.

More formally, strict stationarity requires the joint distribution of any subset of time series observations to be independent of time with respect to all moments. So, in addition to the mean and variance, higher moments such as skew and kurtosis, also need to be constant, irrespective of the lag between different observations. In most applications, we limit stationarity to first and second moments so that the time series is covariance stationary with constant mean, variance, and autocorrelation.

Note that we specifically allow for dependence between observations at different lags, just like we want the input data for linear regression to be correlated with the outcome. Stationarity implies that these relationships are stable, which facilitates prediction as the model can focus on learning systematic patterns that take place within stable statistical properties. It is important because classical statistical models assume that the time series input data is stationary.

To satisfy the stationarity assumption of linear time series models, we need to transform the original time series, often in several steps. Common transformations include the application of the (natural) logarithm to convert an exponential growth pattern into a linear trend and stabilize the variance. Deflation implies dividing a time series by another series that causes trending behavior, for example dividing a nominal series by a price index to convert it into a real measure.

1.3.1 Log Transformation

Double check for zero values

```
[8]: (nasdaq == 0).any(), (industrial_production==0).any()
```

```
[8]: (False, False)
```

```
[9]: nasdaq_log = np.log(nasdaq)
     industrial_production_log = np.log(industrial_production)
```

1.3.2 Differencing

In many cases, de-trending is not sufficient to make the series stationary. Instead, we need to transform the original data into a series of period-to-period and/or season-to-season differences. In other words, we use the result of subtracting neighboring data points or values at seasonal lags from each other. Note that when such differencing is applied to a log-transformed series, the results represent instantaneous growth rates or returns in a financial context.

If a univariate series becomes stationary after differencing d times, it is said to be integrated of the order of d , or simply integrated if $d=1$. This behavior is due to so-called unit roots.

Differencing of log series produces instantaneous returns.

```
[10]: nasdaq_log_diff = nasdaq_log.diff().dropna()

      # seasonal differencing => yoy instantanteous returns
      industrial_production_log_diff = industrial_production_log.diff(12).dropna()
```

1.3.3 Plot Series

The following chart shows time series for the NASDAQ stock index and industrial production for the 30 years through 2017 in original form, as well as the transformed versions after applying the logarithm and subsequently applying first and seasonal differences (at lag 12), respectively. The charts also display the ADF p-value, which allows us to reject the hypothesis of unit-root non-stationarity after all transformations in both cases:

```
[11]: with sns.axes_style('dark'):
      fig, axes = plt.subplots(nrows=3, ncols=2, figsize=(14, 8))

      nasdaq.plot(ax=axes[0][0],
                  title='NASDAQ Composite Index')
      axes[0][0].text(x=.03,
                      y=.85,
                      s=f'ADF: {tsa.adfuller(nasdaq.dropna())[1]:.4f}',
                      transform=axes[0][0].transAxes)
      axes[0][0].set_ylabel('Index')

      nasdaq_log.plot(ax=axes[1][0],
                      sharex=axes[0][0])
```

```

axes[1][0].text(x=.03, y=.85,
                s=f'ADF1: {tsa.adfuller(nasdaq_log.dropna())[1]:.4f}',
                transform=axes[1][0].transAxes)
axes[1][0].set_ylabel('Log')

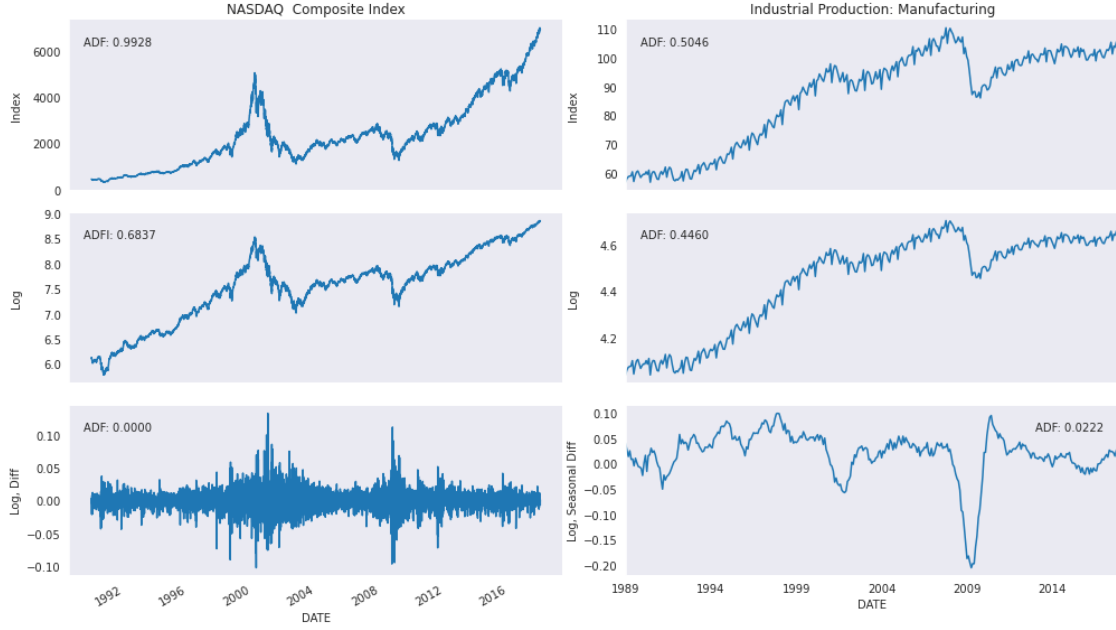
nasdaq_log_diff.plot(ax=axes[2][0],
                    sharex=axes[0][0])
axes[2][0].text(x=.03, y=.85,
                s=f'ADF: {tsa.adfuller(nasdaq_log_diff.dropna())[1]:.4f}',
                transform=axes[2][0].transAxes)
axes[2][0].set_ylabel('Log, Diff')

industrial_production.plot(ax=axes[0][1],
                           title='Industrial Production: Manufacturing')
axes[0][1].text(x=.03, y=.85,
                s=f'ADF: {tsa.adfuller(industrial_production)[1]:.4f}',
                transform=axes[0][1].transAxes)
axes[0][1].set_ylabel('Index')

industrial_production_log.plot(ax=axes[1][1],
                               sharex=axes[0][1])
axes[1][1].text(x=.03, y=.85,
                s=f'ADF: {tsa.adfuller(industrial_production_log.
→dropna())[1]:.4f}',
                transform=axes[1][1].transAxes)
axes[1][1].set_ylabel('Log')

industrial_production_log_diff.plot(ax=axes[2][1],
                                    sharex=axes[0][1])
axes[2][1].text(x=.83, y=.85,
                s=f'ADF: {tsa.adfuller(industrial_production_log_diff.
→dropna())[1]:.4f}',
                transform=axes[2][1].transAxes)
axes[2][1].set_ylabel('Log, Seasonal Diff')
sns.despine()
fig.tight_layout()
fig.align_ylabels(axes)

```



1.4 Correlogram

Autocorrelation (also called serial correlation) adapts the concept of correlation to the time series context: just as the correlation coefficient measures the strength of a linear relationship between two variables, the autocorrelation coefficient, ρ_k , measures the extent of a linear relationship between time series values separated by a given lag, k .

Hence, we can calculate one autocorrelation coefficient for each of the $T-1$ lags in a time series; T is the length of the series. The autocorrelation function (ACF) computes the correlation coefficients as a function of the lag. The autocorrelation for a lag larger than 1 (that is, between observations more than one time step apart) reflects both the direct correlation between these observations and the indirect influence of the intervening data points. The partial autocorrelation removes this influence and only measures the linear dependence between data points at the given lag distance. The partial autocorrelation function (PACF) provides all the correlations that result once the effects of a correlation at shorter lags have been removed.

There are algorithms that estimate the partial autocorrelation from the sample autocorrelation based on the exact theoretical relationship between the PACF and the ACF.

A correlogram is simply a plot of the ACF or PACF for sequential lags, $k=0,1,\dots,n$. It allows us to inspect the correlation structure across lags at one glance. The main usage of correlograms is to detect any autocorrelation after the removal of the effects of deterministic trend or seasonality. Both the ACF and the PACF are key diagnostic tools for the design of linear time series models and we will review examples of ACF and PACF plots in the following section on time series transformations.

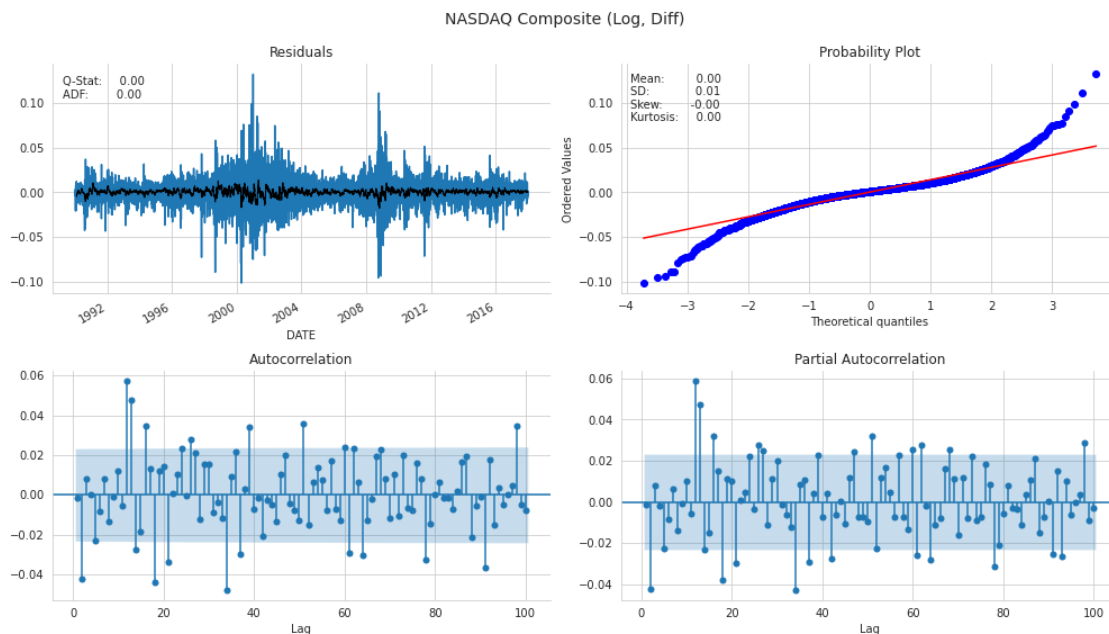
1.4.1 NASDAQ (log, diff)

We can further analyze the relevant time series characteristics for the transformed series using a Q-Q plot that compares the quantiles of the distribution of the time series observation to the

quantiles of the normal distribution and the correlograms based on the ACF and PACF.

For the NASDAQ plot, we notice that while there is no trend, the variance is not constant but rather shows clustered spikes around periods of market turmoil in the late 1980s, 2001, and 2008. The Q-Q plot highlights the fat tails of the distribution with extreme values more frequent than the normal distribution would suggest. The ACF and the PACF show similar patterns with autocorrelation at several lags appearing significant:

```
[12]: plot_correlogram(nasdaq_log_diff, lags=100, title='NASDAQ Composite (Log, Diff)')
```



1.4.2 Industrial Production (log, seasonl diff)

For the monthly time series on industrial manufacturing production, we notice a large negative outlier following the 2008 crisis as well as the corresponding skew in the Q-Q plot. The autocorrelation is much higher than for the NASDAQ returns and declines smoothly. The PACF shows distinct positive autocorrelation patterns at lag 1 and 13, and significant negative coefficients at lags 3 and 4:

```
[13]: plot_correlogram(industrial_production_log_diff, title='Industrial Production (Seasonal Diff)')
```

Industrial Production (Seasonal Diff)

