

## Autoregressiver integrierter gleitender Durchschnitt (ARIMA)

Der autoregressive integrierte gleitende Durchschnitt (ARIMA) ist eine statistische Methode, die ein Autoregressionsmodell für Zeitreihendaten verwendet, um auf der Grundlage historischer Daten zukünftige Prognosen abzuleiten. ARIMA kombiniert zwei Hauptmethoden: ein autoregressives Modell und ein gleitendes Durchschnittsmodell. Die Autoregressionsmethode lernt eine lineare Gleichung zwischen aktuellen und vergangenen Werten und verwendet sie, um zukünftige Prognosen der Variablen zu erstellen. Ein autoregressives Modell der Ordnung p AR(p) ist gegeben durch:

$$y_t = C + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + \varepsilon_t$$

Dabei ist  $\varepsilon_t$  das weiße Rauschen,  $\phi$  die Koeffizienten der linearen Gleichung. P ist die Reihenfolge einer Autoregression und stellt die Anzahl der vorhergehenden Werte der Variablen dar, die zur Erstellung zukünftiger Prognosen verwendet werden.

Andererseits verwendet ein gleitendes Durchschnittsmodell eine lineare Kombination von Prognosefehlern für seine Vorhersagen. Die folgende Gleichung erklärt ein gleitendes Durchschnittsmodell der Ordnung p:

$$y_t = C + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q}$$

Dabei ist  $\theta$  die Koeffizienten der linearen Gleichung.

Wenn die Differenzierung zwischen aufeinanderfolgenden Beobachtungen zusammen mit der Autoregression und einem gleitenden Durchschnittsmodell erfolgt, erhalten wir ein nicht-saisonales ARIMA-Modell. Das ARIMA-Modell kann mit der folgenden Gleichung dargestellt werden:

$$y'_t = C + \phi_1 y'_{t-1} + \dots + \phi_p y'_{t-p} + \theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q} + \varepsilon_t$$

wobei  $y'_t$  die differenzierte Reihe ist. Dieses Modell wird als ARIMA-Modell (p, d, q) bezeichnet, wobei p die Ordnung der autoregressiven Komponente, d die Ordnung der beteiligten Differenzierungen und q die Ordnung des gleitenden Durchschnitts sind. Die Differenzierung ist eine Transformationsmethode, mit der die zeitliche Abhängigkeit in der Zeitreihe beseitigt wird. Dazu gehört die Eliminierung von Komponenten wie Trends und Saisonalität, um Zeitreihen stationär zu machen. Die Häufigkeit, mit der die Differenzierung durchgeführt wird, wird als Differenzierungsreihenfolge bezeichnet. Die Differenzierung erster Ordnung ist also die Differenz zwischen den aktuellen und den vorherigen Zeitdatenpunkten ( $y_i - y_{i-1}$ ), während die Differenzierung zweiter Ordnung wie folgt beschrieben wird:  $(y_i - y_{i-1}) - (y_{i-1} - y_{i-2})$ . Obwohl die Differenzierung ein wichtiger Schritt bei der Vorbereitung von Daten für die Verwendung in einem ARIMA-Modell ist, muss sie sorgfältig durchgeführt werden, um eine übermäßige Differenzierung der Reihen zu vermeiden. Die richtige Differenzierungsreihenfolge ist die Mindestdifferenzierung, die erforderlich ist, um eine nahezu stationäre Reihe zu erhalten. Wenn die Serie beispielsweise für viele Verzögerungen (10 oder mehr) immer noch positive Autokorrelationen aufweist, muss die Serie weiter differenziert werden. Wenn andererseits die Autokorrelation bei Verzögerung 1 zu klein ist, müssen wir möglicherweise die Differenzierungsordnung verringern.

Das SARIMA-Modell (Seasonal ARIMA) ist eine Erweiterung des ARIMA-Modells. Dies wird durch die Addition einer linearen Kombination aus saisonalen Vergangenheitswerten und Prognosefehlern erreicht.

## ARIMA Übung

Für unsere Modelle werden wir [hier](#) tägliche Zeitreihen zur Stromproduktion verwenden. Zuerst müssen wir die Serie importieren und einen kurzen Blick darauf werfen:

```
In [1]: # lags: the number of lags to show the autocorrelation of data points at;
# the autocorrelation of values which are N/lags periods apart is.
# plot_acf(): statsmodels function to inspect which lags have a strong correlation to the forecasted value,
# this plot shows how much each previous lag influences the future lag.
# dataset : electric production

import pandas as pd
import matplotlib.pyplot as plt
from statsmodels.graphics.tsaplots import plot_acf
import numpy as np
df = pd.read_csv('Electric_Production.csv')
df['DATE']= pd.to_datetime(df['DATE'])
df
```

Out[1]:

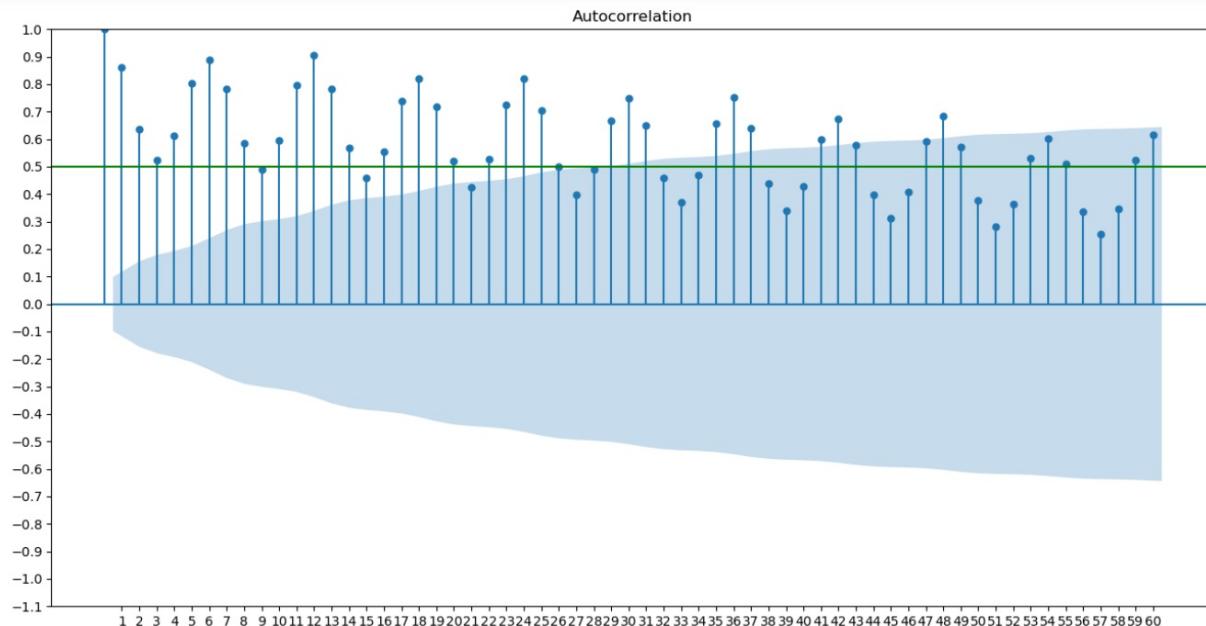
	DATE	Value
0	1985-01-01	72.5052
1	1985-02-01	70.6720
2	1985-03-01	62.4502
3	1985-04-01	57.4714
4	1985-05-01	55.3151
...	...	...
392	2017-09-01	98.6154
393	2017-10-01	93.6137
394	2017-11-01	97.3359

Jetzt überprüfen wir die Autokorrelation in der Zeitreihe mit der Funktion plot\_acf, die die Autokorrelationen bei n Verzögerungen darstellt. Hier setzen wir „Lags“ auf 60. Dieser Schritt ist wichtig, um die Anzahl der Lags zu bestimmen, die in das Modell einbezogen werden sollen.

Die Werte über dem Signifikanzniveau (0,5) weisen auf positive Autokorrelationen in der Zeitreihe mit der entsprechenden Verzögerung hin.  
Wir werden offensichtlich signifikante Autokorrelationen bis zur Verzögerung 60 bemerken:

```
In [4]: fig, ax = plt.subplots(figsize=(16,8))
# nlags= 10 *log10(N) ~ 60, N: dataset length; show autocorrelations at nlags
plot_acf(df['Value'], lags=60, ax=ax)
plt.ylim([-0.1,1])
plt.yticks(np.arange(-1.1, 1.1, 0.1))
plt.xticks(np.arange(1, 61, 1))
plt.axhline(y=0.5, color="green")
plt.show()

# high positive autocorelation (above the threshold :0.5 ) of data points till lag 60
```



Nachdem wir die Daten in Trainings- und Testsätze aufgeteilt haben, werden wir ein autoregressives Modell trainieren und bewerten. Die hier verwendete Bewertungsmetrik ist der mittlere durchschnittliche prozentuale Fehler (MAPE):

```
In [5]: # splitting dataset into train and test
train_data = df['Value'].iloc[:50]
test_data = df['Value'].iloc[-50:]
len(train_data)
from statsmodels.tsa.ar_model import AutoReg
# The number of lags to include in the model
model = AutoReg(train_data, lags=60).fit()

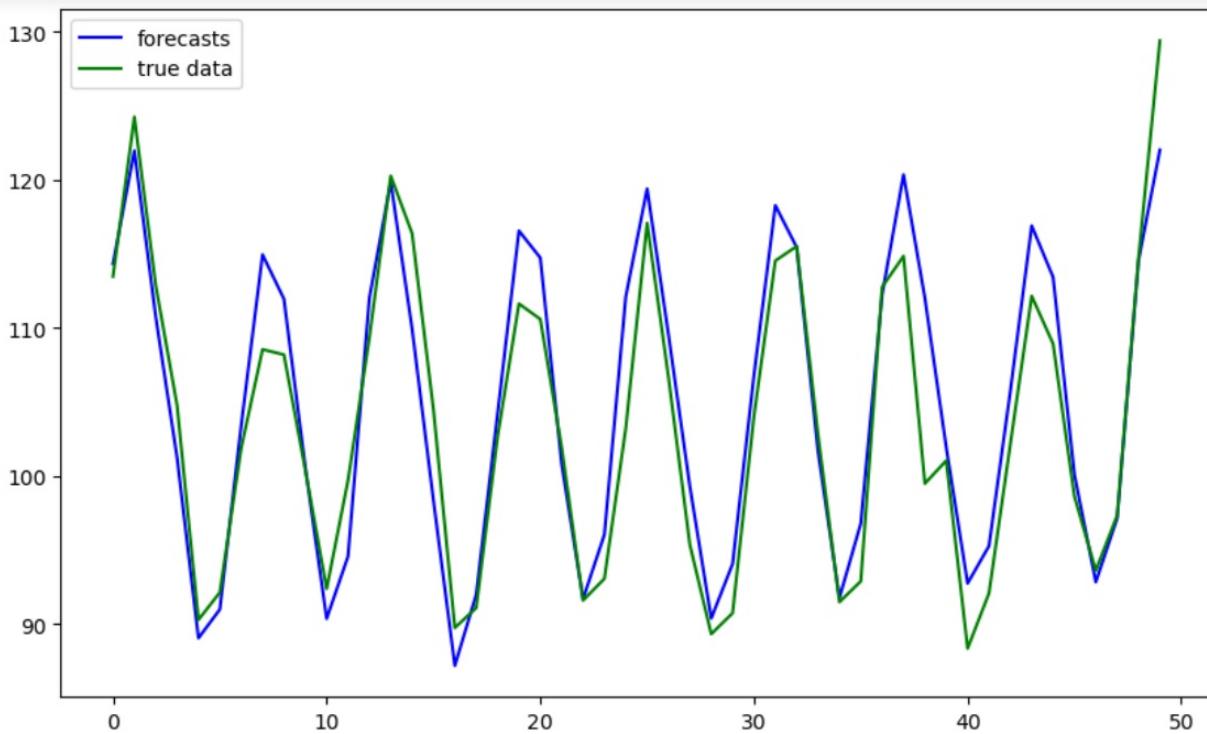
In [6]: # evaluating autoregressive model
forecasts = model.forecast(50).tolist()
test_values = test_data.tolist()

# evaluate the model using MAPE metric
from sklearn.metrics import mean_absolute_percentage_error
print(mean_absolute_percentage_error(test_values, forecasts))
# obviously, the error is about 3%
```

0.028449641966619997

Um eine visuelle Vorstellung davon zu bekommen, wie genau die Prognosen im Vergleich zu echten Daten sind, können wir sowohl die Prognosen als auch die Testdaten grafisch darstellen, um zu sehen, wie nahe die Kurven erscheinen:

```
In [7]: #To get a visual idea of how accurate the forecasts against true data are,
# we can plot both the forecasts and the test data to see how close the curves seem
fig = plt.subplots(figsize=(10,6))
plt.plot(forecasts, color="blue", label ='forecasts')
plt.plot(test_values,color="green", label ='true data')
plt.legend(loc="upper left")
plt.show()
# both curves seem to in agreement with each other
```

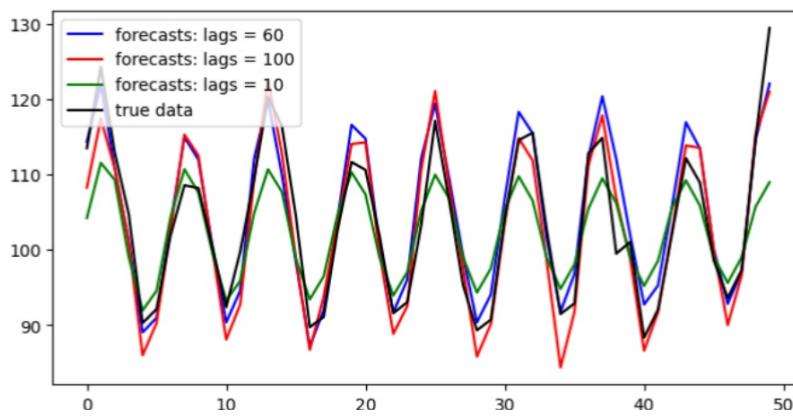


Jetzt werden wir das AR-Modell mit einer anderen Anzahl von Verzögerungen neu erstellen und die resultierenden Prognosen mit echten Daten vergleichen:

```
In [63]: # re-build model with different lags values:100, 10
model1 = AutoReg(train_data, lags=100).fit()
model2= AutoReg(train_data, lags=10).fit()
# evaluating autoregressive model
forecasts1 = model1.forecast(50).tolist()
forecasts2 = model2.forecast(50).tolist()
test_values = test_data.tolist()

# plot the curves against each other

fig = plt.subplots(figsize=(12,8))
plt.plot(forecasts, color="blue", label ='forecasts: lags = 60')
plt.plot(forecasts1, color="red", label ='forecasts: lags = 100')
plt.plot(forecasts2, color="green", label ='forecasts: lags = 10')
plt.plot(test_values,color="black", label ='true data')
plt.legend(loc="upper left")
plt.show()
```



Jetzt erstellen wir ein ARIMA-Modell unter Verwendung des Zeitreihendatensatzes [hier](#).

Zuerst prüfen wir mithilfe des Augmented Dickey Fuller-Tests (`adffuller()`) aus dem Statsmodels-Paket, ob die Daten stationär sind. Die Nullhypothese des ADF-Tests besagt, dass die Zeitreihe instationär ist. Wenn also der p-Wert des Tests unter dem Signifikanzniveau (0,05) liegt, lehnen Sie die Nullhypothese ab und folgern, dass die Zeitreihe tatsächlich stationär ist und umgekehrt:

```
In [28]: # import time series data
df2 = pd.read_csv('https://raw.githubusercontent.com/selva86/datasets/master/a10.csv', parse_dates=['date'])

df2.head(5)
#df2['date'] = pd.to_datetime(df2['date'])
df2 = pd.DataFrame(df2.set_index('date'))
df2.index = pd.DatetimeIndex(df2.index).to_period('M')
# 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.

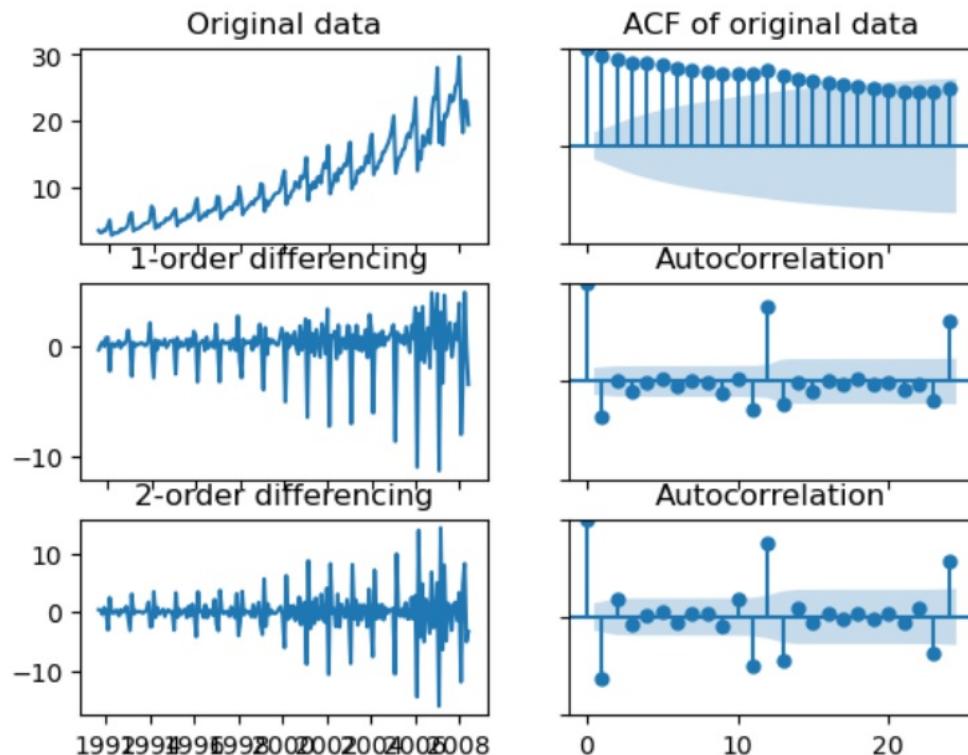
result = adfuller(df2.value)
print('ADF Statistic: %f' % result[0])
print('p-value: %f' % result[1])
# p-value is higher than significance level 0.05 > there is no reason to reject null_hypothesis and thus the data is non stationary
# with non stationary data, ARIMA can not be applied; use differencing
```

ADF Statistic: 3.145186  
p-value: 1.000000

Für ARIMA müssen drei Hauptargumente ermittelt werden: die Reihenfolge des AR-Terms (p), die Reihenfolge des MA-Terms (q) und die Reihenfolge der Differenzierung (d). Da die Zeitreihe nicht stationär ist, müssen wir Differenzierungen anwenden. Die Differenzierungsreihenfolge wird bestimmt, indem die Abnahme der Autokorrelation in den Daten nach jedem Differenzierungsschritt überprüft wird. Lassen Sie uns eine Differenzierung 1. und 2. Ordnung durchführen und die Ergebnisse überprüfen:

```
In [33]: # tuning p, d, q
# p is the order of the AR term
# q is the order of the MA term
# d is the order of differencing required to make the time series stationary
# time Series original acf plot

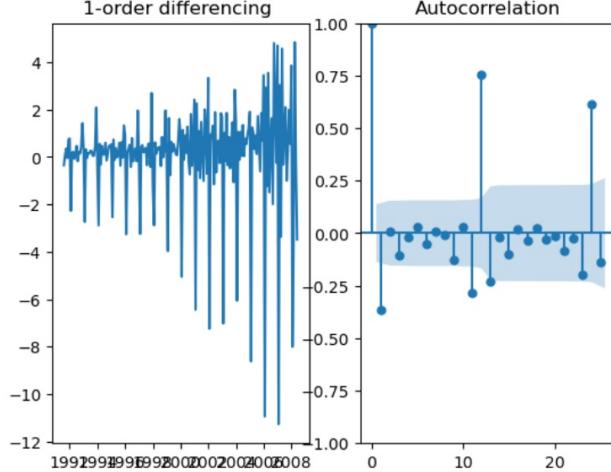
df1 = df2.to_timestamp()
figure, axis = plt.subplots(3, 2)
axis[0, 0].plot(df1)
axis[0, 0].set_title('Original data')
# default number of lags are considered
plot_acf(df2, ax=axis[0, 1])
axis[0, 1].set_title('ACF of original data')
axis[1, 0].plot(df1.diff())
axis[1, 0].set_title('1-order differencing')
axis[1, 1].plot(df2.diff().dropna(), ax=axis[1, 1])
axis[2, 0].plot(df1.diff().diff())
axis[2, 0].set_title('2-order differencing')
plot_acf(df2.diff().diff().dropna(), ax=axis[2, 1])
# Hide x Labels and tick labels for top plots and y ticks for right plots.
for ax in axis.flat:
    ax.label_outer()
# we will stop at first differencing since autocorrelation doesn't decrease after that.
```



ACF-Diagramme zeigen an, dass die Differenzierung 2. Ordnung keinen signifikanten Rückgang der Autokorrelation mit sich bringt, daher wird d auf 1 gesetzt. ACF-Diagramme geben auch den q-Wert an, den ARIMA annehmen sollte. Wenn wir uns das ACF-Diagramm der Differenzierung 1. Ordnung ansehen, können wir offensichtlich eine signifikante Autokorrelation bei Verzögerung 24 erkennen, daher ist q auf 24 gesetzt:

```
In [314]: fig, (ax1, ax2) = plt.subplots(1, 2)
ax1.plot(df1.diff().dropna())
ax1.set_title('1-order differencing ')
plot_acf(df2.diff().dropna(), lags=25, ax=ax2)
```

Out[314]:



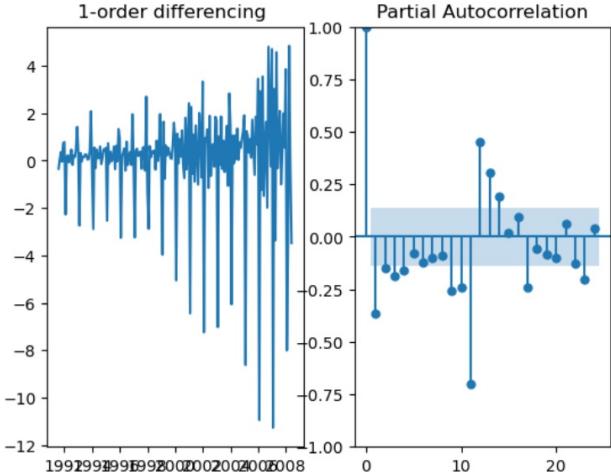
Jetzt müssen wir uns für den p-Wert entscheiden, der den AR-Term darstellt. Wir können die erforderliche Anzahl von AR-Termen herausfinden, indem wir das Diagramm der partiellen Autokorrelation (PACF) untersuchen. Partielle Autokorrelation kann man sich als die Korrelation zwischen der Reihe und ihrer Verzögerung vorstellen, nach Ausschluss der Beiträge aus den Zwischenverzögerungen. Lassen Sie uns die partielle Autokorrelation der differenzierten Daten erster Ordnung untersuchen:

```
In [315]: # AR term: p, You can find out the required number of AR terms
# by inspecting the Partial Autocorrelation (PACF) plot.
# Partial autocorrelation can be imagined as the correlation between the series and its lag,
# after excluding the contributions from the intermediate lags.

fig, (ax1, ax2) = plt.subplots(1, 2)
ax1.plot(df1.diff().dropna())
ax1.set_title('1-order differencing ')
plot_pacf(df2.diff().dropna(), ax=ax2)
# pacf starts from lag 1, so we should look at the first lag in plot
# the blue region is the significance threshold
# there is cutoff at lag 2, so we choose p=1

C:\Users\la2022\anaconda3\lib\site-packages\statsmodels\graphics\tsaplots.py:348: FutureWarning: The default method 'yw' can produce PACF values outside of the [-1,1] interval. After 0.13, the default will change to unadjusted Yule-Walker ('ywm'). You can use this method now by setting method='ywm'.
warnings.warn()
```

Out[315]:



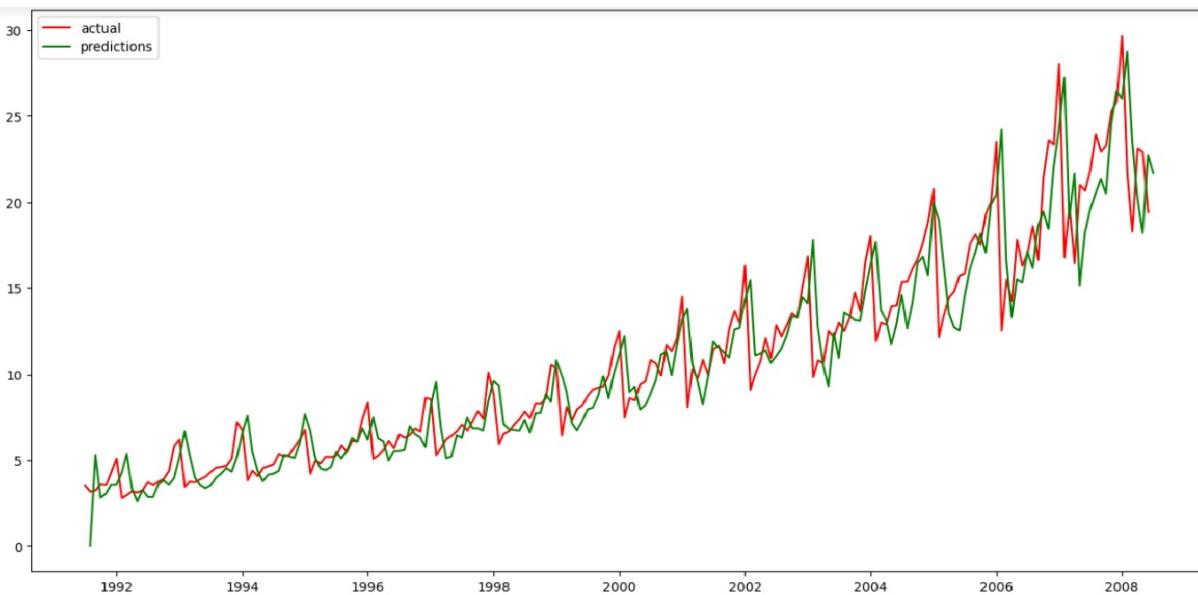
Jetzt erstellen wir ein ARIMA-Modell und verwenden es, um Vorhersagen innerhalb der Stichprobe zu treffen. Lassen Sie uns Vorhersagen anhand wahrer Daten darstellen und den resultierenden Fehlerwert überprüfen:

```
In [42]: # arima: order(p,d,q), p is the order (number of time lags) of the autoregressive model
from statsmodels.tsa.arima.model import ARIMA

# 1,1,24 ARIMA Model
model = ARIMA(df2, order=(1,1,24))
model_fit = model.fit()
print(model_fit.summary())
```

```
C:\Users\la2022\anaconda3\lib\site-packages\statsmodels\tsa\statespace\sarimax.py:978: UserWarning: Non-invertible starting MA parameters found. Using zeros as starting parameters.
warn('Non-invertible starting MA parameters found.'
```

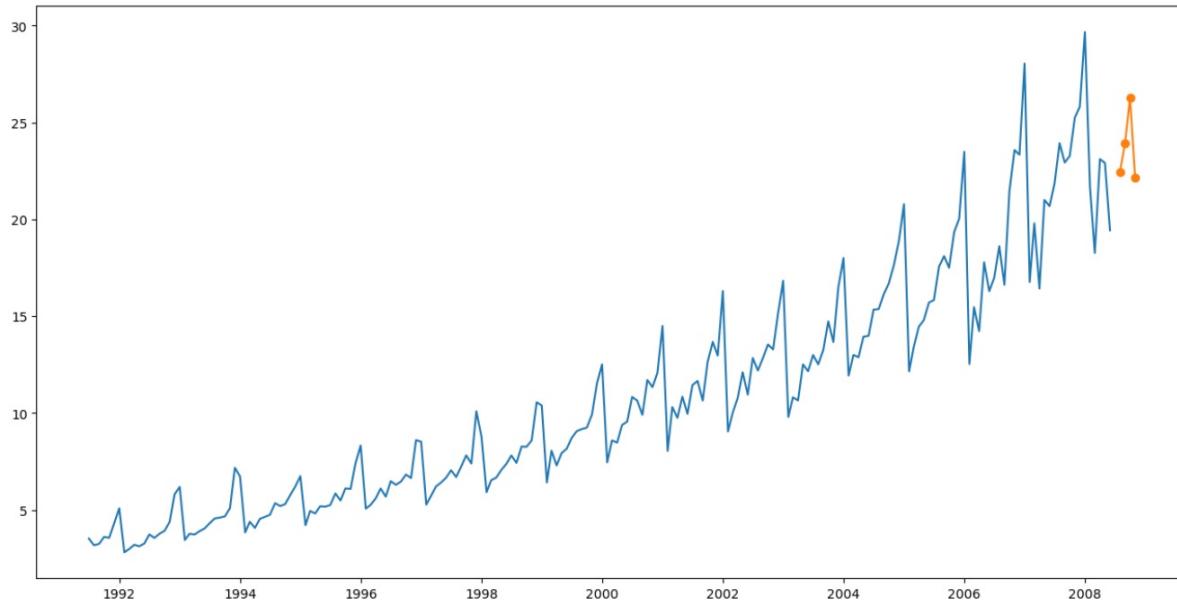
```
SARIMAX Results
=====
Dep. Variable:                      value
No. Observations:                  204
Model:             ARIMA(1, 1, 24)
Date:                Mon, 19 Jun 2023
Time:          16:13:16
Sample:        07-31-1991 - 06-30-2008
Covariance Type:            opg
```



```
In [44]: mape = mean_absolute_percentage_error(df2, pred)
mape
```

```
Out[44]: 0.08448952981466776
```

Jetzt erstellen wir ein ARIMA-Modell und verwenden es, um Vorhersagen innerhalb der Stichprobe zu treffen. Lassen Sie uns Vorhersagen anhand wahrer Daten darstellen und den resultierenden Fehlerwert überprüfen:



```

1 {
2 "cells": [
3 {
4   "cell_type": "code",
5   "execution_count": 1,
6   "id": "ac764d5d",
7   "metadata": {},
8   "outputs": [
9     {
10       "data": {
11         "text/html": [
12           "<div>\n",
13           "<style scoped>\n",
14           "  .dataframe tbody tr th:only-of-type {\n",
15             vertical-align: middle;\n",
16           }\n",
17           "\n",
18           ".dataframe tbody tr th {\n",
19             vertical-align: top;\n",
20           }\n",
21           "\n",
22           ".dataframe thead th {\n",
23             text-align: right;\n",
24           }\n",
25           "</style>\n",
26           "<table border=\"1\" class=\"dataframe\"\n",
27           "  <thead>\n",
28           "    <tr style=\"text-align: right;\"\n",
29           "      <th></th>\n",
30           "      <th>DATE</th>\n",
31           "      <th>Value</th>\n",

```

```

32     "</tr>\n",
33     "</thead>\n",
34     "<tbody>\n",
35     "    <tr>\n",
36     "        <th>0</th>\n",
37     "        <td>1985-01-01</td>\n",
38     "        <td>72.5052</td>\n",
39     "    </tr>\n",
40     "    <tr>\n",
41     "        <th>1</th>\n",
42     "        <td>1985-02-01</td>\n",
43     "        <td>70.6720</td>\n",
44     "    </tr>\n",
45     "    <tr>\n",
46     "        <th>2</th>\n",
47     "        <td>1985-03-01</td>\n",
48     "        <td>62.4502</td>\n",
49     "    </tr>\n",
50     "    <tr>\n",
51     "        <th>3</th>\n",
52     "        <td>1985-04-01</td>\n",
53     "        <td>57.4714</td>\n",
54     "    </tr>\n",
55     "    <tr>\n",
56     "        <th>4</th>\n",
57     "        <td>1985-05-01</td>\n",
58     "        <td>55.3151</td>\n",
59     "    </tr>\n",
60     "    <tr>\n",
61     "        <th>...</th>\n",
62     "        <td>...</td>\n",
63     "        <td>...</td>\n",
64     "    </tr>\n",
65     "    <tr>\n",
66     "        <th>392</th>\n",
67     "        <td>2017-09-01</td>\n",
68     "        <td>98.6154</td>\n",
69     "    </tr>\n",
70     "    <tr>\n",
71     "        <th>393</th>\n",
72     "        <td>2017-10-01</td>\n",
73     "        <td>93.6137</td>\n",
74     "    </tr>\n",
75     "    <tr>\n",
76     "        <th>394</th>\n",
77     "        <td>2017-11-01</td>\n",
78     "        <td>97.3359</td>\n",
79     "    </tr>\n",
80     "    <tr>\n",
81     "        <th>395</th>\n",
82     "        <td>2017-12-01</td>\n",
83     "        <td>114.7212</td>\n",
84     "    </tr>\n",
85     "    <tr>\n",
86     "        <th>396</th>\n",
87     "        <td>2018-01-01</td>\n",
88     "        <td>129.4048</td>\n",
89     "    </tr>\n",
90     "    </tbody>\n",
91     "</table>\n",
92     "<p>397 rows x 2 columns</p>\n",
93     "</div>\n",
94 ],
95 "text/plain": [
96     "      DATE      Value\n",
97     "0  1985-01-01    72.5052\n",
98     "1  1985-02-01    70.6720\n",
99     "2  1985-03-01    62.4502\n",
100    "3  1985-04-01    57.4714\n",
101    "4  1985-05-01    55.3151\n",
102    "...  ...  ...\n",
103    "392 2017-09-01    98.6154\n",
104    "393 2017-10-01    93.6137\n",
105    "394 2017-11-01    97.3359\n",
106    "395 2017-12-01   114.7212\n",
107    "396 2018-01-01   129.4048\n",
108    "\n",
109    "[397 rows x 2 columns]"
110  ]
111 },
112 "execution_count": 1,
113 "metadata": {},
114 "output_type": "execute_result"
115 }
116 ],
117 "source": [
118     "# lags: the number of lags to show the autocorrelation of data points at;\n",
119     "# the autocorrelation of values which are N/lags periods apart is.\n",
120     "# plot_acf() : statsmodels function to inspect which lags have a strong correlation to the forecasted value,\n",
121     "# this plot shows how much each previous lag influences the future lag.\n",
122     "# dataset : electric production\n",
123     "\n",
124     "import pandas as pd\n",
125     "import matplotlib.pyplot as plt\n",
126     "from statsmodels.graphics.tsaplots import plot_acf\n",
127     "import numpy as np\n",
128     "df = pd.read_csv('Electric_Production.csv')\n",
129     "df['DATE']= pd.to_datetime(df['DATE'])\n",
130     "df"
131  ]
132 },
133 {
134     "cell_type": "code",
135     "execution_count": 2,

```

```

136 "id": "a30dc660",
137 "metadata": {},
138 "outputs": [
139 {
140 "data": {
141 "image/png": "iVBORw0KGgoAAAANSUhEUgAAABRcAAAKoCAYAAADtfshrAAAAOXRFWHRTb2Z0d2FyZQBNYXRwbG90bGliIHZlcNpb24zLjcuMSwgaHR0",
142 "text/plain": [
143 "<Figure size 1600x800 with 1 Axes>",
144 ]
145 },
146 "metadata": {},
147 "output_type": "display_data"
148 }
149 ],
150 "source": [
151 "fig, ax = plt.subplots(figsize=(16,8))\n",
152 "# nlags= 10 * log10(N) ? 60, N: dataset length; show autocorrelations at nlags\n",
153 "plot_acf(df['Value'], lags=60, ax=ax)\n",
154 "plt.ylim([0,1])\n",
155 "plt.yticks(np.arange(-1.1, 1.1, 0.1))\n",
156 "plt.xticks(np.arange(1, 61, 1))\n",
157 "plt.axhline(y=0.5, color=\"green\")\n",
158 "plt.show()\n",
159 "\n",
160 "# high positive autocorelation (above the threshold :0.5 ) of data points till lag 60"
161 ]
162 },
163 {
164 "cell_type": "code",
165 "execution_count": 3,
166 "id": "04756836",
167 "metadata": {},
168 "outputs": [],
169 "source": [
170 "# splitting dataset into train and test\n",
171 "train_data = df['Value'].iloc[:50]\n",
172 "test_data = df['Value'].iloc[-50:]\n",
173 "len(train_data)\n",
174 "from statsmodels.tsa.ar_model import AutoReg\n",
175 "# The number of lags to include in the model \n",
176 "model = AutoReg(train_data, lags=60).fit()"
177 ]
178 },
179 {
180 "cell_type": "code",
181 "execution_count": 4,
182 "id": "ba2f555e",
183 "metadata": {},
184 "outputs": [
185 {
186 "name": "stdout",
187 "output_type": "stream",
188 "text": [
189 "0.028449641966619997\n"
190 ]
191 }
192 ],
193 "source": [
194 "# evaluating autoregressive model \n",
195 "forecasts = model.forecast(50).tolist()\n",
196 "test_values = test_data.tolist()\n",
197 "\n",
198 "# evaluate the model using MAPE metric\n",
199 "from sklearn.metrics import mean_absolute_percentage_error\n",
200 "print(mean_absolute_percentage_error(test_values, forecasts))\n",
201 "# obviously, the error is about 3%"
202 ]
203 },
204 {
205 "cell_type": "code",
206 "execution_count": null,
207 "id": "af933ddd",
208 "metadata": {},
209 "outputs": [],
210 "source": []
211 },
212 {
213 "cell_type": "code",
214 "execution_count": 5,
215 "id": "00c67b6e",
216 "metadata": {},
217 "outputs": [
218 {
219 "data": {
220 "image/png": "iVBORw0KGgoAAAANSUhEUgAAAz8AAAHH5CAYAACve4DDAAAAOXRFWHRTb2Z0d2FyZQBNYXRwbG90bGliIHZlcNpb24zLjcuMSwgaHR0",
221 "text/plain": [
222 "<Figure size 1000x600 with 1 Axes>"
223 ]
224 },
225 "metadata": {},
226 "output_type": "display_data"
227 }
228 ],
229 "source": [
230 "#To get a visual idea of how accurate the forecasts against true data are,\n",
231 "# we can plot both the forecasts and the test data to see how close the curves seem\n",
232 "fig = plt.subplots(figsize=(10,6))\n",
233 "plt.plot(forecasts, color=\"blue\", label ='forecasts')\n",
234 "plt.plot(test_values,color=\"green\", label ='true data')\n",
235 "plt.legend(loc=\"upper left\")\n",
236 "plt.show()\n",
237 "# both curves seem to in agreement with each other"
238 ]
239 },

```

```

240 {
241     "cell_type": "code",
242     "execution_count": 6,
243     "id": "01fdbc29",
244     "metadata": {},
245     "outputs": [],
246     "source": [
247         "# re-build model with different lags values:100, 10\n",
248         "model1 = AutoReg(train_data, lags=100).fit()\n",
249         "model2= AutoReg(train_data, lags=10).fit()\n",
250         "# evaluating autoregressive model \n",
251         "forecasts1 = model1.forecast(50).tolist()\n",
252         "forecasts2 = model2.forecast(50).tolist()\n",
253         "test_values = test_data.tolist()"
254     ]
255 },
256 {
257     "cell_type": "code",
258     "execution_count": 7,
259     "id": "fb3bc50b",
260     "metadata": {
261         "scrolled": true
262     },
263     "outputs": [
264     {
265         "data": {
266             "image/png": "iVBORw0KGgoAAAANSUhEUgAAAAqQAAFFCAYAAACY+97uAAAAOXRFWHRTb2Z0d2FyZQBNEYRwbG90bGliIH2lcnNpb24zLjcuMSwgaHR0",
267             "text/plain": [
268                 "<Figure size 800x400 with 1 Axes>"
269             ]
270         },
271         "metadata": {},
272         "output_type": "display_data"
273     }
274 ],
275 "source": [
276     "# plot the curves against each other \n",
277     "\n",
278     "fig = plt.subplots(figsize=(8,4))\n",
279     "plt.plot(forecasts, color='blue', label ='forecasts: lags = 60')\n",
280     "plt.plot(forecasts1, color='red', label ='forecasts: lags = 100')\n",
281     "plt.plot(forecasts2, color='green', label ='forecasts: lags = 10')\n",
282     "plt.plot(test_values,color='black', label ='true data')\n",
283     "plt.legend(loc='upper left')\n",
284     "plt.show()\n",
285     "
286 ],
287 },
288 {
289     "cell_type": "code",
290     "execution_count": 8,
291     "id": "fb54019c",
292     "metadata": {},
293     "outputs": [
294     {
295         "name": "stdout",
296         "output_type": "stream",
297         "text": [
298             "ADF Statistic: -1.773219\n",
299             "p-value: 0.39383653785330297\n"
300         ]
301     },
302     {
303         "name": "stderr",
304         "output_type": "stream",
305         "text": [
306             "C:\\\\Users\\\\la2022\\\\AppData\\\\Local\\\\Temp\\\\ipykernel_17244\\\\4186456414.py:26: FutureWarning: The default value of numer.
307             df2 = df2.groupby(level=0).sum()\n"
308         ]
309     },
310     {
311         "data": {
312             "text/html": [
313                 "<div>\n",
314                 "<style scoped>\n",
315                     ".dataframe tbody tr th:only-of-type {\n",
316                         vertical-align: middle;\n",
317                     }\n",
318                 "\n",
319                     .dataframe tbody tr th {\n",
320                         vertical-align: top;\n",
321                     }\n",
322                 "\n",
323                     .dataframe thead th {\n",
324                         text-align: right;\n",
325                     }\n",
326                 "</style>\n",
327                 "<table border=\"1\" class=\"dataframe\"\n",
328                     <thead>\n",
329                         <tr style=\"text-align: right;\"\n",
330                             <th></th>\n",
331                             <th>sales</th>\n",
332                         </tr>\n",
333                         <tr>\n",
334                             <th>date</th>\n",
335                             <th></th>\n",
336                         </tr>\n",
337                         </thead>\n",
338                         <tbody>\n",
339                             <tr>\n",
340                                 <th>2013-01-01</th>\n",
341                                 <td>2.0</td>\n",
342                             </tr>\n",
343                             <tr>\n",
344                         </tbody>\n"
345             ]
346         }
347     }
348 ]
349 }
```

```

344     "      <th>2013-01-02</th>\n",
345     "      <td>207.0</td>\n",
346     "      </tr>\n",
347     "      <tr>\n",
348     "        <th>2013-01-03</th>\n",
349     "        <td>125.0</td>\n",
350     "      </tr>\n",
351     "      <tr>\n",
352     "        <th>2013-01-04</th>\n",
353     "        <td>133.0</td>\n",
354     "      </tr>\n",
355     "      <tr>\n",
356     "        <th>2013-01-05</th>\n",
357     "        <td>191.0</td>\n",
358     "      </tr>\n",
359     "      <tr>\n",
360     "        <th>...</th>\n",
361     "        <td>...</td>\n",
362     "      </tr>\n",
363     "      <tr>\n",
364     "        <th>2017-08-11</th>\n",
365     "        <td>351.0</td>\n",
366     "      </tr>\n",
367     "      <tr>\n",
368     "        <th>2017-08-12</th>\n",
369     "        <td>369.0</td>\n",
370     "      </tr>\n",
371     "      <tr>\n",
372     "        <th>2017-08-13</th>\n",
373     "        <td>433.0</td>\n",
374     "      </tr>\n",
375     "      <tr>\n",
376     "        <th>2017-08-14</th>\n",
377     "        <td>337.0</td>\n",
378     "      </tr>\n",
379     "      <tr>\n",
380     "        <th>2017-08-15</th>\n",
381     "        <td>339.0</td>\n",
382     "      </tr>\n",
383     "    </tbody>\n",
384   "</table>\n",
385   "<p>1684 rows × 1 columns</p>\n",
386   "</div>"
```

387 ],  
388 "text/plain": [  
389 " sales\n",  
390 "date \n",  
391 "2013-01-01 2.0\n",  
392 "2013-01-02 207.0\n",  
393 "2013-01-03 125.0\n",  
394 "2013-01-04 133.0\n",  
395 "2013-01-05 191.0\n",  
396 "... ...\n",  
397 "2017-08-11 351.0\n",  
398 "2017-08-12 369.0\n",  
399 "2017-08-13 433.0\n",  
400 "2017-08-14 337.0\n",  
401 "2017-08-15 339.0\n",  
402 "\n",  
403 "[1684 rows × 1 columns]"  
404 ]  
405 },  
406 "execution\_count": 8,  
407 "metadata": {},  
408 "output\_type": "execute\_result"  
409 }  
410 ],  
411 "source": [  
412 "# AutoRegressive Integrated Moving Average(ARIMA) is a time series forecasting model that incorporates autocorrelation  
413 "# to model temporal structures within the time series data to predict future values.\n",  
414 "# The autoregression part of the model measures the dependency of a particular sample with a few past observations. \n",  
415 "# These differences are measured and integrated to make the data patterns stationary or\n",  
416 "# minimize the obvious correlation with past data (collinearity or linear dependence)\n",  
417 "# (since linear independence and no collinearity is one of the fundamental assumptions of the linear regression model).  
418 "# After this, a moving average helps condense and bring out significant features from the data.\n",  
419 "# stationary time series is Consequently, parameters such as mean and variance also do not change over time.\n",  
420 "# Thus, time series with trends, or with seasonality, are not stationary,\n",  
421 "# because trends and seasonality affect the time series at different time points.\n",  
422 "# ARIMA equation :y?t = c + ?1y?t?1 + ? + ?py?t?p + ?1?t?1 + ? + ?q?t?q + ?t\n",  
423 "\n",  
424 "# we will be using Store\_sale time series dataset: \n",  
425 "# https://www.kaggle.com/competitions/store-sales-time-series-forecasting/data\n",  
426 "# first we check if the data is stationary using the Augmented Dickey Fuller test (adfuller()),from the statsmodels pac:  
427 "from statsmodels.tsa.stattools import adfuller\n",  
428 "from numpy import log\n",  
429 "import pandas as pd\n",  
430 "df2 = pd.read\_csv('Store\_sales\_train.csv',parse\_dates=['date'])\n",  
431 "\n",  
432 "#df2['date'] = pd.to\_datetime(df2['date'])\n",  
433 "df2 = pd.DataFrame(df2, columns=['date','sales', 'family']).set\_index('date')\n",  
434 "df2.index = pd.DatetimeIndex(df2.index).to\_period('D')\n",  
435 "df2 = df2[df2.sales != 0]\n",  
436 "df2 = df2[df2.family == 'BEAUTY']\n",  
437 "df2 = df2.groupby(level=0).sum()\n",  
438 "## remove Family column\n",  
439 "df = pd.DataFrame(df2, columns=['sales'])\n",  
440 "\n",  
441 "result = adfuller(df)\n",  
442 "print('ADF Statistic: %f' % result[0])\n",  
443 "print('p-value: %s' % result[1])\n",  
444 "df\n"
445 ]
446 },
447 {

```

448     "cell_type": "code",
449     "execution_count": 9,
450     "id": "199c290b",
451     "metadata": {},
452     "outputs": [
453     {
454         "name": "stdout",
455         "output_type": "stream",
456         "text": [
457             "ADF Statistic: 3.145186\n",
458             "p-value: 1.000000\n"
459         ]
460     }
461 ],
462 "source": [
463     "# import time series data\n",
464     "df2 = pd.read_csv('https://raw.githubusercontent.com/selva86/datasets/master/a10.csv', parse_dates=['date'])\n",
465     "\n",
466     "df2.head(5)\n",
467     "#df2['date'] = pd.to_datetime(df2['date'])\n",
468     "df2 = pd.DataFrame(df2).set_index('date')\n",
469     "df2.index = pd.DatetimeIndex(df2.index).to_period('M')\n",
470     "# The null hypothesis of the ADF test is that the time series is non-stationary.\n",
471     "#So, if the p-value of the test is less than the significance level (0.05) then you reject the null hypothesis and\n",
472     "#infer that the time series is indeed stationary.\n",
473     "\n",
474     "result = adfuller(df2.value)\n",
475     "print('ADF Statistic: %f' % result[0])\n",
476     "print('p-value: %f' % result[1])\n",
477     "# p-value is higher than significance level 0.05 > there is no reason to reject null hypothesis and thus the data is non
478     "# with non stationary data, ARIMA can not be applied;use differencing\n"
479   ],
480 },
481 {
482     "cell_type": "code",
483     "execution_count": 17,
484     "id": "748cf318",
485     "metadata": {},
486     "outputs": [
487     {
488         "data": {
489             "image/png": "iVBORw0KGgoAAAANSUhEUgAABRcAAAKoCAYAAAAtfshrAAAAOXRFWHRB2Z0d2FyZQBNYXRwbG90bGliIHZlcnPb24zLjcuMSwgaHR0
490             "text/plain": [
491                 "<Figure size 1600x800 with 1 Axes>"
492             ]
493         },
494         "metadata": {},
495         "output_type": "display_data"
496     }
497 ],
498 "source": [
499     "import matplotlib.pyplot as plt\n",
500     "from statsmodels.graphics.tsaplots import plot_acf, plot_pacf\n",
501     "import numpy as np\n",
502     "fig, ax = plt.subplots(figsize=(16,8))\n",
503     "# nlags= 10 *log10(N) ? 32, N: dataset length; show autocorrelations at nlags\n",
504     "plot_acf(df2, lags=32, ax=ax)\n",
505     "plt.ylim([0,1])\n",
506     "plt.yticks(np.arange(-1.1, 1.1, 0.1))\n",
507     "plt.xticks(np.arange(1, 33, 1))\n",
508     "plt.axhline(y=0.5, color='green')\n",
509     "plt.show()\n",
510     "# we always ignore the long spike at lag 0, we look from lag 1 on .."
511   ],
512 },
513 {
514     "cell_type": "code",
515     "execution_count": 14,
516     "id": "181debe8",
517     "metadata": {},
518     "outputs": [
519     {
520         "data": {
521             "image/png": "iVBORw0KGgoAAAANSUhEUgAAAGxCAYAACju/aQAAAAXRFWHRB2Z0d2FyZQBNYXRwbG90bGliIHZlcnPb24zLjcuMSwgaHR0
522             "text/plain": [
523                 "<Figure size 640x480 with 6 Axes>"
524             ]
525         },
526         "metadata": {},
527         "output_type": "display_data"
528     }
529 ],
530 "source": [
531     "# tuning p, d, q \n",
532     "# p is the order of the AR term\n",
533     "# q is the order of the MA term\n",
534     "# d is the order of differencing required to make the time series stationary\n",
535     "# time Series original acf plot\n",
536     "\n",
537     "df1 = df2.to_timestamp()\n",
538     "figure, axis = plt.subplots(3, 2)\n",
539     "axis[0, 0].plot(df1)\n",
540     "axis[0, 0].set_title('Original data')\n",
541     "# default number of lags are considered\n",
542     "plot_acf(df2, ax=axis[0, 1])\n",
543     "axis[0, 1].set_title('ACF of original data')\n",
544     "axis[1, 0].plot(df1.diff())\n",
545     "axis[1, 0].set_title('1-order differencing ')\n",
546     "plot_acf(df2.diff().dropna(), ax=axis[1, 1])\n",
547     "\n",
548     "axis[2, 0].plot(df1.diff().diff())\n",
549     "axis[2, 0].set_title('2-order differencing')\n",
550     "plot_acf(df2.diff().diff().dropna(), ax=axis[2, 1])\n",
551     "# Hide x labels and tick labels for top plots and y ticks for right plots.\n",

```

```

552     "for ax in axis.flat:\n",
553     "    ax.label_outer()\n",
554     "# we will stop at first differencing since autocorrelation doesn't decrease afterthat.      "
555   },
556 },
557 {
558   "cell_type": "code",
559   "execution_count": 20,
560   "id": "506f6ea0",
561   "metadata": {},
562   "outputs": [
563     {
564       "data": {
565         "image/png": "iVBORw0KGgoAAAANSUhEUgAAAiIAAAGxCAYAAABfrtlAAAAOXRFWHRB2Z0d2FyZQBNYXRwbG90bGliIHZlcnNpb24zLjcuMSwgaHR0",
566         "text/plain": [
567           "<Figure size 640x480 with 4 Axes>"
568         ],
569       },
570       "metadata": {},
571       "output_type": "display_data"
572     },
573   ],
574   "source": [
575     "figure, (ax1, ax2, ax3, ax4) = plt.subplots(4)\n",
576     "plot_acf(df2, ax=ax1)\n",
577     "#add semicolon to suppress plotting returned output twice \n",
578     "plot_acf(df2.diff().dropna(), ax=ax2);\n",
579     "plot_acf(df2.diff().diff().dropna(), ax=ax3);\n",
580     "plot_acf(df2.diff().diff().diff().dropna(), ax=ax4);"
581   ],
582 },
583 {
584   "cell_type": "code",
585   "execution_count": 21,
586   "id": "c5bb284e",
587   "metadata": {},
588   "outputs": [
589     {
590       "data": {
591         "image/png": "iVBORw0KGgoAAAANSUhEUgAAIsAAAGxCAYAACju/aQAAAAXRFWHRB2Z0d2FyZQBNYXRwbG90bGliIHZlcnNpb24zLjcuMSwgaHR0",
592         "text/plain": [
593           "<Figure size 640x480 with 2 Axes>"
594         ],
595       },
596       "metadata": {},
597       "output_type": "display_data"
598     },
599   ],
600   "source": [
601     "fig, (ax1, ax2) = plt.subplots(1, 2)\n",
602     "ax1.plot(df1.diff().dropna())\n",
603     "ax1.set_title('1-order differencing ')\n",
604     "plot_acf(df2.diff().dropna(), lags=25, ax=ax2);"
605   ],
606 },
607 {
608   "cell_type": "code",
609   "execution_count": 23,
610   "id": "df86a764",
611   "metadata": {},
612   "outputs": [
613     {
614       "name": "stderr",
615       "output_type": "stream",
616       "text": [
617         "C:\\\\Users\\\\la2022\\\\anaconda3\\\\lib\\\\site-packages\\\\statsmodels\\\\graphics\\\\tsaplots.py:348: FutureWarning: The default :
618           warnings.warn(\n"
619       ],
620     },
621   ],
622   "data": {
623     "image/png": "iVBORw0KGgoAAAANSUhEUgAAIsAAAGxCAYAACju/aQAAAAXRFWHRB2Z0d2FyZQBNYXRwbG90bGliIHZlcnNpb24zLjcuMSwgaHR0",
624     "text/plain": [
625       "<Figure size 640x480 with 2 Axes>"
626     ],
627   },
628   "metadata": {},
629   "output_type": "display_data"
630 },
631 ],
632   "source": [
633     "# AR term: p, You can find out the required number of AR terms\n",
634     "# by inspecting the Partial Autocorrelation (PACF) plot.\n",
635     "# Partial autocorrelation can be imagined as the correlation between the series and its lag,\n",
636     "# after excluding the contributions from the intermediate lags.\n",
637     "\n",
638     "fig, (ax1, ax2) = plt.subplots(1, 2)\n",
639     "ax1.plot(df1.diff().dropna());\n",
640     "ax1.set_title('1-order differencing ')\n",
641     "plot_pacf(df2.diff().dropna(), ax=ax2);\n",
642     "# pacf starts from lag 1, so we should look at the first lag in plot\n",
643     "# the blue region is the significance threshold\n",
644     "# there is cutoff at lag 2, so we choose p=1"
645   ],
646 },
647 {
648   "cell_type": "code",
649   "execution_count": 24,
650   "id": "2d3eae28",
651   "metadata": {
652     "scrolled": true
653   },
654   "outputs": [
655     {

```

```

656     "data": {
657         "image/png": "iVBORw0KGgoAAAANSUhEUgAABSAAAkoCAYAAABnUVEqAAAAOXRFWHRtbZz0d2FyZQBNYXRwbG90bGliIHZlcNpb24zLjcuMSwgaHRO",
658         "text/plain": [
659             "<Figure size 1600x800 with 1 Axes>"
660         ],
661     },
662     "metadata": {},
663     "output_type": "display_data"
664   },
665 ],
666 "source": [
667     "# MA term q : from the ACF plot we decide the value of q.\n",
668     "# An MA term is technically, the error of the lagged forecast.\n",
669     "# The ACF tells how many MA terms (lags) are required to remove any autocorrelation in the stationarized series.\n",
670     "\n",
671     "fig, ax = plt.subplots(figsize=(16,8))\n",
672     "plot_acf(df2.diff().dropna(),lags=30, ax=ax);\n",
673     "# this plot indicates that after lag=24 the autocorrelation gets obviously insignificant"
674   ],
675 },
676 {
677     "cell_type": "code",
678     "execution_count": 25,
679     "id": "18e2d9df",
680     "metadata": {},
681     "outputs": [
682       {
683         "name": "stderr",
684         "output_type": "stream",
685         "text": [
686           "C:\\\\Users\\\\la2022\\\\anaconda3\\\\lib\\\\site-packages\\\\statsmodels\\\\tsa\\\\statespace\\\\sarimax.py:978: UserWarning: Non-inve.",
687           "    warn('Non-invertible starting MA parameters found.'\n"
688         ]
689       },
690       {
691         "name": "stdout",
692         "output_type": "stream",
693         "text": [
694             "                               SARIMAX Results\n",
695             "=====\\n",
696             "Dep. Variable:                      value    No. Observations:                  204\\n",
697             "Model:                          ARIMA(1, 1, 24)    Log Likelihood:                -310.453\\n",
698             "Date:                            Wed, 18 Oct 2023    AIC:                         672.907\\n",
699             "Time:                             17:15:51    BIC:                         759.050\\n",
700             "Sample:                           07-31-1991    HQIC:                        707.757\\n",
701             "                                    - 06-30-2008\n",
702             "Covariance Type:                   opg\n",
703             "=====\\n",
704             "      coef    std err        z     P>|z|    [0.025    0.975]\\n",
705             "-----\\n",
706             "ar.L1      -0.5384    0.138     -3.893    0.000    -0.809    -0.267\\n",
707             "ma.L1      -0.1004    0.672     -0.149    0.881    -1.417    1.216\\n",
708             "ma.L2      -0.3508    0.414     -0.847    0.397    -1.162    0.461\\n",
709             "ma.L3      -0.0110    0.709     -0.016    0.988    -1.402    1.380\\n",
710             "ma.L4      -0.1149    0.706     -0.163    0.871    -1.499    1.269\\n",
711             "ma.L5      -0.0072    0.591     -0.012    0.990    -1.166    1.152\\n",
712             "ma.L6      -0.3124    0.283     -1.103    0.270    -0.867    0.242\\n",
713             "ma.L7      0.0010    0.736     0.001    0.999    -1.441    1.443\\n",
714             "ma.L8      0.2245    0.712     0.315    0.752    -1.171    1.620\\n",
715             "ma.L9      -0.2578    0.782     -0.330    0.742    -1.790    1.275\\n",
716             "ma.L10     -0.2377    0.544     -0.437    0.662    -1.304    0.829\\n",
717             "ma.L11     0.0050    0.788     0.006    0.995    -1.540    1.550\\n",
718             "ma.L12     1.2720    0.700     1.817    0.069    -0.100    2.644\\n",
719             "ma.L13     -0.0417    0.690     -0.060    0.952    -1.395    1.311\\n",
720             "ma.L14     -0.4580    0.621     -0.737    0.461    -1.676    0.760\\n",
721             "ma.L15     0.0620    0.697     0.089    0.929    -1.304    1.428\\n",
722             "ma.L16     0.0705    0.659     0.107    0.915    -1.221    1.362\\n",
723             "ma.L17     -0.0183    0.571     -0.032    0.974    -1.137    1.100\\n",
724             "ma.L18     -0.2395    0.416     -0.575    0.565    -1.056    0.577\\n",
725             "ma.L19     -0.1094    0.608     -0.180    0.857    -1.302    1.083\\n",
726             "ma.L20     0.0460    0.654     0.070    0.944    -1.235    1.327\\n",
727             "ma.L21     -0.2108    0.656     -0.321    0.748    -1.496    1.075\\n",
728             "ma.L22     -0.3021    0.518     -0.583    0.560    -1.317    0.713\\n",
729             "ma.L23     0.1558    0.527     0.295    0.768    -0.878    1.189\\n",
730             "ma.L24     0.7804    0.666     1.172    0.241    -0.525    2.086\\n",
731             "sigma2     0.9935    0.788     1.261    0.207    -0.550    2.537\\n",
732             "=====\\n",
733             "Ljung-Box (L1) (Q):                 0.49    Jarque-Bera (JB):            30.32\\n",
734             "Prob(Q):                           0.48    Prob(JB):                  0.00\\n",
735             "Heteroskedasticity (H):              9.09    Skew:                      -0.28\\n",
736             "Prob(H) (two-sided):                0.00    Kurtosis:                  4.81\\n",
737             "=====\\n",
738             "\n",
739             "Warnings:\\n",
740             "[1] Covariance matrix calculated using the outer product of gradients (complex-step).\\n"
741       ],
742     },
743   },
744   {
745     "name": "stderr",
746     "output_type": "stream",
747     "text": [
748       "C:\\\\Users\\\\la2022\\\\anaconda3\\\\lib\\\\site-packages\\\\statsmodels\\\\base\\\\model.py:604: ConvergenceWarning: Maximum Likeli",
749       "    warnings.warn(\"Maximum Likelihood optimization failed to \"\n"
750     ]
751   },
752   "source": [
753     "# arima: order(p,d,q), p is the order (number of time lags) of the autoregressive model\\n",
754     "from statsmodels.tsa.arima.model import ARIMA\\n",
755     "\n",
756     "# 1,1,24 ARIMA Model\\n",
757     "model = ARIMA(df2, order=(1,1,24))\\n",
758     "model_fit = model.fit()\\n",
759     "print(model_fit.summary())"

```

```

760     ],
761   },
762   {
763     "cell_type": "code",
764     "execution_count": 26,
765     "id": "a6b2adf6",
766     "metadata": {},
767     "outputs": [
768       {
769         "data": {
770           "image/png": "iVBORw0KGgoAAAANSUhEUgAAABQcAAAKTCAYAACkbb7tAAAAOXRFWHRTb2Z0d2FyZQBNYXRwbG90bGliIHZlcNpb24zLjcuMSwgaHR0",
771           "text/plain": [
772             "<Figure size 1600x800 with 1 Axes>",
773           ],
774         },
775         "metadata": {},
776         "output_type": "display_data"
777       }
778     ],
779   },
780   "source": [
781     "# Actual vs Fitted\n",
782     "from sklearn.metrics import mean_absolute_percentage_error\n",
783     "fig, ax = plt.subplots(figsize=(16,8))\n",
784     "#pred= model_fit.forecast(steps= 100, dynamic=False)\n",
785     "# in sample predictions: predict() After a model has been fit predict returns the fitted values.\n",
786     "pred= model_fit.predict( dynamic=False)\n",
787     "\n",
788     "#df1 = df.to_timestamp()\n",
789     "plt.plot(df1, label ='actual',color='r')\n",
790     "plt.plot(pred,color='g',label= 'predictions')\n",
791     "plt.legend(loc=\"upper left\")\n",
792     "plt.show()"
793   ],
794 },
795   {
796     "cell_type": "code",
797     "execution_count": null,
798     "id": "48b76cbb",
799     "metadata": {},
800     "outputs": [],
801     "source": [
802       "mape = mean_absolute_percentage_error(df2, pred)\n",
803       "mape"
804     ],
805   },
806   {
807     "cell_type": "code",
808     "execution_count": 27,
809     "id": "61bc7451",
810     "metadata": {},
811     "outputs": [
812       {
813         "name": "stderr",
814         "output_type": "stream",
815         "text": [
816           "C:\\\\Users\\\\la2022\\\\AppData\\\\Local\\\\Temp\\\\ipykernel_17244\\\\3493184615.py:3: FutureWarning: The pandas.datetime class is\n           from pandas import datetime\n"
817         ]
818       },
819     ],
820     "data": {
821       "image/png": "iVBORw0KGgoAAAANSUhEUgAAABQcAAAKTCAYAACkbb7tAAAAOXRFWHRTb2Z0d2FyZQBNYXRwbG90bGliIHZlcNpb24zLjcuMSwgaHR0",
822       "text/plain": [
823         "<Figure size 1600x800 with 1 Axes>",
824       ],
825     },
826     "metadata": {},
827     "output_type": "display_data"
828   },
829 },
830   "source": [
831     "#out of samples forecasts\n",
832     "\n",
833     "from pandas import datetime\n",
834     "#start_index = datetime(2008, 7)\n",
835     "#end_index = datetime(2008, 12)\n",
836     "forecasts = model_fit.forecast(steps=4)\n",
837     "fig, ax = plt.subplots(figsize=(16,8))\n",
838     "# convert from period into timestamp\n",
839     "\n",
840     "plt.plot(df1, label ='true data')\n",
841     "plt.plot(forecasts, label='forecasts', marker='o')\n",
842     "plt.show()\n"
843   ],
844 },
845 },
846   "cell_type": "code",
847   "execution_count": null,
848   "id": "153762cb",
849   "metadata": {},
850   "outputs": [],
851   "source": []
852 },
853 },
854   "metadata": {
855     "kernelspec": {
856       "display_name": "Python 3 (ipykernel)",
857       "language": "python",
858       "name": "python3"
859     },
860     "language_info": {
861       "codemirror_mode": {
862         "name": "ipython",
863         "version": 3

```

```

864     },
865     "file_extension": ".py",
866     "mimetype": "text/x-python",
867     "name": "python",
868     "nbconvert_exporter": "python",
869     "pygments_lexer": "ipython3",
870     "version": "3.9.13"
871   }
872 },
873 "nbformat": 4,
874 "nbformat_minor": 5
875 }

```

## Trigonometric seasonality, Box-Cox transformation, ARMA errors, Trend and Seasonal components (TBATS)

Während die meisten gängigen Modelle wie ARIMA und exponentielle Glättung nur einzelnes Saisonalitätsmuster lernen können, wurde die TBATS-Prognosemethode entwickelt, um Zeitreihen zu modellieren, die komplexe und mehrere saisonale Muster umfassen, indem exponentielle Glättung verwendet wird, beispielsweise tägliche Daten, die ein wöchentliches und ein jährliches Muster aufweisen können.

TBATS ist die Abkürzung für:

- Trigonometrische Saisonalität
- Box-Cox-Transformation
- ARMA-Fehler
- Trend
- Saisonale Komponenten

In TBATS wird ein Box-Cox-Transformator verwendet, um Zeitreihen so zu transformieren, dass die Daten eine Normalverteilung aufweisen. Dies ist nützlich, da viele statistische Methoden eine Normalverteilung der Modellfehler annehmen. Anschließend lernt TBATS eine lineare Kombination aus einem exponentiell geglätteten Trend, einer saisonalen Komponente und ARMA-Fehlern. Die saisonale Komponente wird durch trigonometrische Funktionen berechnet. Dadurch ist TBATS in der Lage, genauere Langzeitprognosen zu erstellen und andere Prognosemodelle zu übertreffen. Ein Nachteil dieses Modells ist seine langsame Vorhersage- und Lernleistung bei der Modellierung langer Zeitreihen.

## TBATS Übung

Hier werden wir Zeitreihen für Ladenverkäufe [hier](#) mithilfe der TBATS-Methode modellieren. Dieser Datensatz enthält tägliche Daten zu Verkäufen verschiedener Artikel. Importieren wir zunächst die Daten:

```
In [1]: import pandas as pd
# https://www.kaggle.com/competitions/store-sales-time-series-forecasting/data
data = pd.read_csv('Store_sales_train.csv')
data
```

```
Out[1]:      id    date  store_nbr       family   sales  onpromotion
0         0 2013-01-01          1  AUTOMOTIVE  0.000          0
1         1 2013-01-01          1  BABY CARE  0.000          0
2         2 2013-01-01          1        BEAUTY  0.000          0
3         3 2013-01-01          1  BEVERAGES  0.000          0
4         4 2013-01-01          1        BOOKS  0.000          0
...
3000883  3000883 2017-08-15          9      POULTRY  438.133          0
3000884  3000884 2017-08-15          9  PREPARED FOODS  154.553          1
3000885  3000885 2017-08-15          9        PRODUCE  2419.729         148
3000886  3000886 2017-08-15          9  SCHOOL AND OFFICE SUPPLIES  121.000          8
3000887  3000887 2017-08-15          9        SEAFOOD  16.000          0
```

3000888 rows × 6 columns

Als nächstes extrahieren wir Zeitreihen für einen Artikel, die für das Modell verwendet werden sollen: Wir benötigen nur die Spalte „Verkäufe“, die mithilfe der Spalte „Datum“ indiziert wird. Extrahieren wir Zeitreihen von Geflügelverkäufen und sehen wir, wie die Saisonalität aussieht:

```
In [8]: # modelling poultry sales
#df = data[(data['store_nbr'] == 2) & (data['family'] == 'BEVERAGES') & (data['sales'] > 0) ]
df = data[(data['store_nbr'] == 2) & (data['family'] == 'POULTRY') & (data['sales'] > 0) ]
df = df.set_index('date')
train_data = df['sales']
train_data2 = train_data[:-100]
type(train_data.index)
print(len(train_data))

1679
```

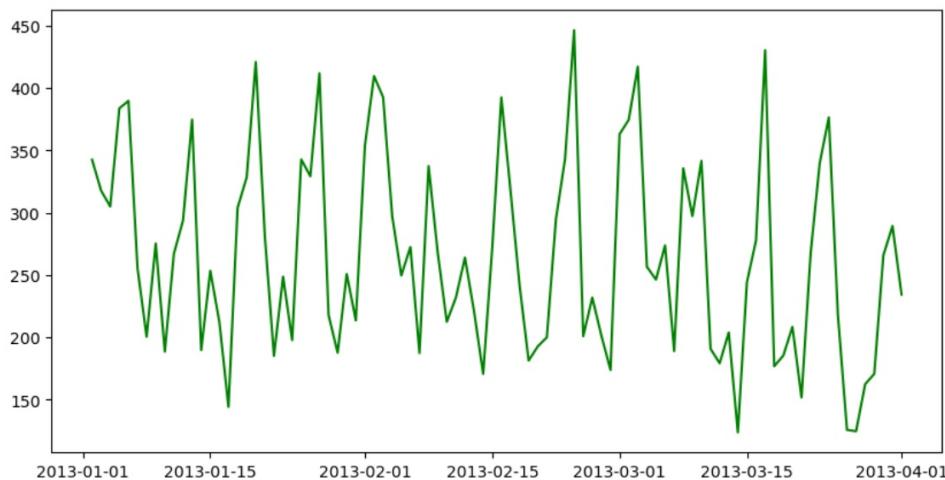
```
In [9]: #checking seasonality in the data using plots
import matplotlib.pyplot as plt
fig, ax = plt.subplots(figsize = (12,8))
data_to_plot = train_data
data_to_plot.index = pd.DatetimeIndex(data_to_plot.index).to_period('d')
data2 = data_to_plot.to_timestamp(freq='d')

plt.plot(data2[:90], label = 'actual data', color="green" )
```

```
Out[9]: [

```

```
Out[11]: [<matplotlib.lines.Line2D at 0x162403bbc10>]
```



Wenn wir uns das Datendiagramm ansehen, sehen wir eine Art wöchentliche Saisonalität.

Jetzt trainieren wir zwei TBATS-Modelle: mit und ohne COX\_box-Transformator. Anschließend werden wir die Leistung beider Modelle anhand von Testdaten vergleichen. TBATS wird wöchentliche und monatliche Saisonalitäten verwenden, um das Modell zu erstellen. Importieren wir zunächst die erforderlichen Python-Pakete:

```
In [21]: # !pip install tbats
#!pip install sktime
from sktime.forecasting.tbats import TBATS
from sktime.forecasting.model_selection import temporal_train_test_split
from sktime.forecasting.base import ForecastingHorizon

y_train, y_test = temporal_train_test_split(data2, test_size = 0.1)# 10% of the data for testing
y_train.index= pd.DatetimeIndex(y_train.index).to_period('d')
y_test.index= pd.DatetimeIndex(y_test.index).to_period('d')
fh = ForecastingHorizon(y_test.index, is_relative=False)

#box-cox transformation:a statistical technique that transforms your target variable
#so that your data closely resembles a normal distribution. This help making the residuals of the model as normal as possible
#and thus leads to better forecasting performance
# n_jobs: How many jobs to run in parallel when fitting TBATS model.
# monthly and weekly seasonalities
forecastor = TBATS(sp=[ 7, 30], use_box_cox=True, n_jobs=5)# with box_cox transformation
forecastor2 = TBATS(sp=[ 7, 30], use_box_cox=False, n_jobs=5)# box_cox transformation is disabled
fitted_model = forecastor.fit(y_train)
fitted_model2 = forecastor2.fit(y_train)
print(forecastor, forecastor2)

TBATS(n_jobs=5, sp=[7, 30], use_box_cox=True) TBATS(n_jobs=5, sp=[7, 30], use_box_cox=False)
```

```
In [262]: # prediction
y_pred = forecastor.predict(fh)
y_pred2 = forecastor2.predict(fh)

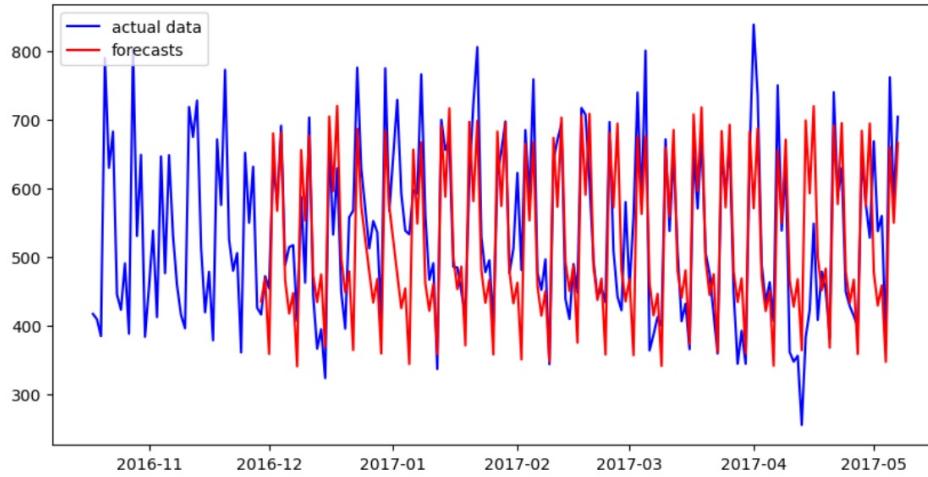
y_pred[:31]
```

```
Out[262]: date
2016-11-29    435.293772
2016-11-30    469.053532
2016-12-01    358.737733
2016-12-02    680.352287
2016-12-03    567.444919
2016-12-04    682.392106
2016-12-05    466.123288
2016-12-06    417.723508
2016-12-07    447.631427
2016-12-08    340.625047
2016-12-09    656.131723
2016-12-10    552.965208
2016-12-11    677.404144
2016-12-12    472.873148
```

Nach der Anpassung der TBATS-Modelle werden wir sie an 10 % der Zeitreihen testen. Lassen Sie uns Diagramme der tatsächlichen Werte mit den Prognosen beider Modelle vergleichen:

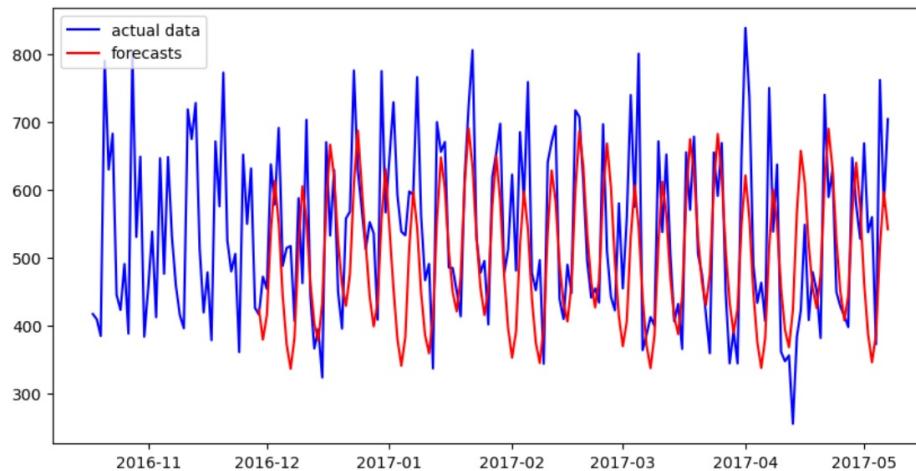
```
In [263]: df_forecast=pd.DataFrame(data= y_pred, index=y_test.index)

df_forecast = df_forecast.to_timestamp(freq='d')
fig, ax = plt.subplots(figsize = (10,5))
train_data2.index = pd.DatetimeIndex(train_data2.index).to_period('d')
train = train_data2.to_timestamp(freq='d')
plt.plot(train[-200:], label = 'actual data', color="blue" )
plt.plot(df_forecast, label= 'forecasts', color = 'red')
plt.legend(loc="upper left")
plt.show()
```



```
In [264]: df_forecast2=pd.DataFrame(data= y_pred2, index=y_test.index)
df_forecast2 = df_forecast2.to_timestamp(freq='d')
fig, ax = plt.subplots(figsize = (10,5))

plt.plot(train[-200:], label = 'actual data', color="blue" )
plt.plot(df_forecast2, label= 'forecasts', color = 'red')
plt.legend(loc="upper left")
plt.show()
```



Wenn wir abschließend die resultierenden Fehlerwerte beider Modelle überprüfen und vergleichen, sehen wir offensichtlich, dass die Einbeziehung des Cox-Box-Transformators zu einer besseren Leistung führt:

```
In [265]: # TBATS forecasting model evaluation:
from sklearn.metrics import mean_absolute_percentage_error
print('mean absolute percentage error on test data using box-cox transformer is:', mean_absolute_percentage_error(y_pred, y_test))
print('mean absolute percentage error on test data without using box-cox transformer is:', mean_absolute_percentage_error(y_pred2))
```

mean absolute percentage error on test data using box-cox transformer is: 0.12534766422768567  
mean absolute percentage error on test data without using box-cox transformer is: 0.1808882858413106

```
1 {
2   "cells": [
3     {
4       "cell_type": "code",
5       "execution_count": 8,
6       "id": "89de8ca2",
7       "metadata": {},
8       "outputs": [
9         {
10           "data": {
11             "text/html": [
12               "<div>\n",
13               "<style scoped>\n",
14               "  .dataframe tbody tr th:only-of-type {\n",
15               "    vertical-align: middle;\n",
16               "  }\n",
17               "\n",
18               "  .dataframe tbody tr th {\n",
19               "    vertical-align: top;\n",
20               "  }\n",
21               "\n",
22               "  .dataframe thead th {\n",
23               "    text-align: right;\n",
24               "  }\n",
25               "</style>\n",
26               "<table border=\"1\" class=\"dataframe\">\n",
```

```
27     "    <thead>\n",
28     "      <tr style=\"text-align: right;\">\n",
29     "        <th></th>\n",
30     "        <th>id</th>\n",
31     "        <th>date</th>\n",
32     "        <th>store_nbr</th>\n",
33     "        <th>family</th>\n",
34     "        <th>sales</th>\n",
35     "        <th>nonpromotion</th>\n",
36     "      </tr>\n",
37     "    </thead>\n",
38     "    <tbody>\n",
39     "      <tr>\n",
40     "        <th>0</th>\n",
41     "        <td>0</td>\n",
42     "        <td>2013-01-01</td>\n",
43     "        <td>1</td>\n",
44     "        <td>AUTOMOTIVE</td>\n",
45     "        <td>0.000</td>\n",
46     "        <td>0</td>\n",
47     "      </tr>\n",
48     "      <tr>\n",
49     "        <th>1</th>\n",
50     "        <td>1</td>\n",
51     "        <td>2013-01-01</td>\n",
52     "        <td>1</td>\n",
53     "        <td>BABY CARE</td>\n",
54     "        <td>0.000</td>\n",
55     "        <td>0</td>\n",
56     "      </tr>\n",
57     "      <tr>\n",
58     "        <th>2</th>\n",
59     "        <td>2</td>\n",
60     "        <td>2013-01-01</td>\n",
61     "        <td>1</td>\n",
62     "        <td>BEAUTY</td>\n",
63     "        <td>0.000</td>\n",
64     "        <td>0</td>\n",
65     "      </tr>\n",
66     "      <tr>\n",
67     "        <th>3</th>\n",
68     "        <td>3</td>\n",
69     "        <td>2013-01-01</td>\n",
70     "        <td>1</td>\n",
71     "        <td>BEVERAGES</td>\n",
72     "        <td>0.000</td>\n",
73     "        <td>0</td>\n",
74     "      </tr>\n",
75     "      <tr>\n",
76     "        <th>4</th>\n",
77     "        <td>4</td>\n",
78     "        <td>2013-01-01</td>\n",
79     "        <td>1</td>\n",
80     "        <td>BOOKS</td>\n",
81     "        <td>0.000</td>\n",
82     "        <td>0</td>\n",
83     "      </tr>\n",
84     "      <tr>\n",
85     "        <th>...</th>\n",
86     "        <td>...</td>\n",
87     "        <td>...</td>\n",
88     "        <td>...</td>\n",
89     "        <td>...</td>\n",
90     "        <td>...</td>\n",
91     "        <td>...</td>\n",
92     "      </tr>\n",
93     "      <tr>\n",
94     "        <th>3000883</th>\n",
95     "        <td>3000883</td>\n",
96     "        <td>2017-08-15</td>\n",
97     "        <td>9</td>\n",
98     "        <td>POULTRY</td>\n",
99     "        <td>438.133</td>\n",
100    "        <td>0</td>\n",
101   "      </tr>\n",
102   "      <tr>\n",
103   "        <th>3000884</th>\n",
104   "        <td>3000884</td>\n",
105   "        <td>2017-08-15</td>\n",
106   "        <td>9</td>\n",
107   "        <td>PREPARED FOODS</td>\n",
108   "        <td>154.553</td>\n",
109   "        <td>1</td>\n",
110   "      </tr>\n",
111   "      <tr>\n",
112   "        <th>3000885</th>\n",
113   "        <td>3000885</td>\n",
114   "        <td>2017-08-15</td>\n",
115   "        <td>9</td>\n",
116   "        <td>PRODUCE</td>\n",
117   "        <td>2419.729</td>\n",
118   "        <td>148</td>\n",
119   "      </tr>\n",
120   "      <tr>\n",
121   "        <th>3000886</th>\n",
122   "        <td>3000886</td>\n",
123   "        <td>2017-08-15</td>\n",
124   "        <td>9</td>\n",
125   "        <td>SCHOOL AND OFFICE SUPPLIES</td>\n",
126   "        <td>121.000</td>\n",
127   "        <td>8</td>\n",
128   "      </tr>\n",
129   "      <tr>\n",
130   "        <th>3000887</th>\n",
```

```

131     "<td>3000887</td>\n",
132     "<td>2017-08-15</td>\n",
133     "<td>9</td>\n",
134     "<td>SEAFOOD</td>\n",
135     "<td>16.000</td>\n",
136     "<td>0</td>\n",
137     "</tr>\n",
138     "</tbody>\n",
139     "</table>\n",
140     "<p>3000888 rows x 6 columns</p>\n",
141     "</div>\n",
142   ],
143   "text/plain": [
144     "      id      date  store_nbr          family    sales  \\\n",
145     "0      0  2013-01-01      1  AUTOMOTIVE  0.000  \n",
146     "1      1  2013-01-01      1  BABY CARE  0.000  \n",
147     "2      2  2013-01-01      1  BEAUTY    0.000  \n",
148     "3      3  2013-01-01      1  BEVERAGES 0.000  \n",
149     "4      4  2013-01-01      1  BOOKS    0.000  \n",
150     "...     ...     ...     ...     ...  \n",
151     "3000883  3000883  2017-08-15      9  POULTRY  438.133  \n",
152     "3000884  3000884  2017-08-15      9  PREPARED FOODS  154.553  \n",
153     "3000885  3000885  2017-08-15      9  PRODUCE  2419.729  \n",
154     "3000886  3000886  2017-08-15      9  SCHOOL AND OFFICE SUPPLIES  121.000  \n",
155     "3000887  3000887  2017-08-15      9  SEAFOOD  16.000  \n",
156     "\n",
157     "      onpromotion  \n",
158     "0      0  \n",
159     "1      0  \n",
160     "2      0  \n",
161     "3      0  \n",
162     "4      0  \n",
163     "...     .  \n",
164     "3000883      0  \n",
165     "3000884      1  \n",
166     "3000885      148 \n",
167     "3000886      8  \n",
168     "3000887      0  \n",
169     "\n",
170     "[3000888 rows x 6 columns]"
171   ]
172 },
173 "execution_count": 8,
174 "metadata": {},
175 "output_type": "execute_result"
176 }
177 ],
178 "source": [
179   "import pandas as pd\n",
180   "# https://www.kaggle.com/competitions/store-sales-time-series-forecasting/data\n",
181   "data = pd.read_csv('Store_sales_train.csv')\n",
182   "data"
183 ],
184 },
185 {
186   "cell_type": "code",
187   "execution_count": 18,
188   "id": "26e30043",
189   "metadata": {},
190   "outputs": [
191     {
192       "name": "stdout",
193       "output_type": "stream",
194       "text": [
195         "1679\n"
196       ]
197     }
198   ],
199   "source": [
200     "# modelling poultry sales\n",
201     "#df = data[(data['store_nbr'] == 2) & (data['family'] == 'BEVERAGES') & (data['sales'] > 0)]\n",
202     "#df = data[(data['store_nbr'] == 2) & (data['family'] == 'POULTRY') & (data['sales'] > 0)]\n",
203     "df = df.set_index('date')\n",
204     "train_data = df['sales']\n",
205     "train_data2 = train_data[-100]\n",
206     "type(train_data.index)\n",
207     "print(len(train_data))"
208   ],
209 },
210 {
211   "cell_type": "code",
212   "execution_count": 10,
213   "id": "4b1ae127",
214   "metadata": {},
215   "outputs": [
216     {
217       "data": {
218         "text/plain": [
219           "[<matplotlib.lines.Line2D at 0x2086fc84220>]"
220         ]
221       },
222       "execution_count": 10,
223       "metadata": {},
224       "output_type": "execute_result"
225     },
226     {
227       "data": {
228         "image/png": "iVBORw0KGgoAAAANSUhEUgAAA0QAAAGsCAYAAD9kWfoAAAAOXRFWHRTbZ20d2FyZQBNYXRwbG90bGliIHZlcNpb24zLjcuMSwgaHRO",
229         "text/plain": [
230           "<Figure size 1000x500 with 1 Axes>"
231         ]
232       },
233       "metadata": {},
234       "output_type": "display_data"
235     }
236   ]

```

```

235     }
236   ],
237   "source": [
238     "#checking seasonality in the data using plots\n",
239     "import matplotlib.pyplot as plt\n",
240     "fig, ax = plt.subplots(figsize = (10,5))\n",
241     "data_to_plot = train_data2\n",
242     "data_to_plot.index = pd.DatetimeIndex(data_to_plot.index).to_period('d')\n",
243     "data2 = data_to_plot.to_timestamp(freq='d')\n",
244     "\n",
245     "plt.plot(data2[:90], label = 'actual data', color=\"green\" )"
246   ],
247 },
248 {
249   "cell_type": "code",
250   "execution_count": 21,
251   "id": "05e90646",
252   "metadata": {},
253   "outputs": [
254     {
255       "name": "stdout",
256       "output_type": "stream",
257       "text": [
258         "TBATS(n_jobs=5, sp=[7, 30], use_box_cox=True) TBATS(n_jobs=5, sp=[7, 30], use_box_cox=False)\n"
259       ]
260     }
261   ],
262   "source": [
263     "# !pip install tbats\n",
264     "#!pip install sktime\n",
265     "from sktime.forecasting.tbats import TBATS\n",
266     "from sktime.forecasting.model_selection import temporal_train_test_split\n",
267     "from sktime.forecasting.base import ForecastingHorizon\n",
268     "\n",
269     "\n",
270     "y_train, y_test = temporal_train_test_split(data2, test_size = 0.1)# 10% of the data for testing \n",
271     "y_train.index= pd.DatetimeIndex(y_train.index).to_period('d')\n",
272     "y_test.index= pd.DatetimeIndex(y_test.index).to_period('d')\n",
273     "fh = ForecastingHorizon(y_test.index, is_relative=False)\n",
274     "\n",
275     "#box-cox transformation:a statistical technique that transforms your target variable\n",
276     "#so that your data closely resembles a normal distribution. This help making the residuals of the model as normal as po
277     # and thus leads to better forecasting performance\n",
278     "# n_jobs: How many jobs to run in parallel when fitting TBATS model.\n",
279     "# monthly and weekly seasonalities\n",
280     "forecaster = TBATS(sp=[ 7, 30], use_box_cox=True, n_jobs=5)# with box_cox transformation \n",
281     "forecaster2 = TBATS(sp=[ 7, 30], use_box_cox=False, n_jobs=5)# box_cox transformation is disabled\n",
282     "fitted_model = forecaster.fit(y_train)\n",
283     "fitted_model2 = forecaster2.fit(y_train)\n",
284     "print(forecaster, forecaster2)"
285   ],
286 },
287 {
288   "cell_type": "code",
289   "execution_count": 22,
290   "id": "bdaa6a22",
291   "metadata": {},
292   "outputs": [
293     {
294       "data": {
295         "text/plain": [
296           "date\n",
297           "2016-11-29    435.293771\n",
298           "2016-11-30    469.053531\n",
299           "2016-12-01    358.737733\n",
300           "2016-12-02    680.352288\n",
301           "2016-12-03    567.444920\n",
302           "2016-12-04    682.392107\n",
303           "2016-12-05    466.123289\n",
304           "2016-12-06    417.723508\n",
305           "2016-12-07    447.631427\n",
306           "2016-12-08    340.625047\n",
307           "2016-12-09    656.131724\n",
308           "2016-12-10    552.965209\n",
309           "2016-12-11    677.404144\n",
310           "2016-12-12    472.873147\n",
311           "2016-12-13    434.279893\n",
312           "2016-12-14    474.687027\n",
313           "2016-12-15    369.224240\n",
314           "2016-12-16    704.953463\n",
315           "2016-12-17    595.895730\n",
316           "2016-12-18    720.441934\n",
317           "2016-12-19    498.540929\n",
318           "2016-12-20    448.743349\n",
319           "2016-12-21    479.056565\n",
320           "2016-12-22    364.409213\n",
321           "2016-12-23    687.148571\n",
322           "2016-12-24    574.165854\n",
323           "2016-12-26    478.701776\n",
324           "2016-12-27    433.658950\n",
325           "2016-12-28    468.119896\n",
326           "2016-12-29    359.386117\n",
327           "2016-12-30    684.294173\n",
328           "Freq: D, Name: sales, dtype: float64"
329         ]
330       },
331       "execution_count": 22,
332       "metadata": {},
333       "output_type": "execute_result"
334     }
335   ],
336   "source": [
337     "# forecasting using the first and second models: with and without Cox_Box transformation\n",
338     "y_pred = forecaster.predict(fh)\n",

```

```

339 "y_pred2 = forecastor2.predict(fh)\n",
340 "\n",
341 "y_pred[:31]"
342 },
343 {
344 "cell_type": "code",
345 "execution_count": 23,
346 "id": "7e3c45a2",
347 "metadata": {},
348 "outputs": [
349 {
350 "data": {
351 "image/png": "iVBORw0KGgoAAAANSUhEUgAAz8AAAGsCAYAADzOBmHAAAAOXRFWHRTbZ0d2FyZQBNYXRwbG90bGliIHZlcnPb24zLjcuMSwgaHR0",
352 "text/plain": [
353 "<Figure size 1000x500 with 1 Axes>\n"
354 ],
355 },
356 },
357 "metadata": {},
358 "output_type": "display_data"
359 }
360 ],
361 "source": [
362 "# plotting true against forecasts using the first model.\n",
363 "df_forecast=pd.DataFrame(data= y_pred, index=y_test.index)\n",
364 "\n",
365 "df_forecast = df_forecast.to_timestamp(freq='d')\n",
366 "fig, ax = plt.subplots(figsize = (10,5))\n",
367 "train_data2.index = pd.DatetimeIndex(train_data2.index).to_period('d')\n",
368 "train = train_data2.to_timestamp(freq='d')\n",
369 "plt.plot(train[-200:], label = 'actual data', color=\\"blue\\" )\n",
370 "plt.plot(df_forecast, label= 'forecasts', color = 'red')\n",
371 "plt.legend(loc=\\"upper left\\")\n",
372 "plt.show()\n"
373 ],
374 },
375 {
376 "cell_type": "code",
377 "execution_count": 24,
378 "id": "7a3dae9b",
379 "metadata": {},
380 "outputs": [
381 {
382 "data": {
383 "image/png": "iVBORw0KGgoAAAANSUhEUgAAz8AAAGsCAYAADzOBmHAAAAOXRFWHRTbZ0d2FyZQBNYXRwbG90bGliIHZlcnPb24zLjcuMSwgaHR0",
384 "text/plain": [
385 "<Figure size 1000x500 with 1 Axes>\n"
386 ],
387 },
388 "metadata": {},
389 "output_type": "display_data"
390 }
391 ],
392 "source": [
393 "# plotting true against forecasts using the second model.\n",
394 "df_forecast2=pd.DataFrame(data= y_pred2, index=y_test.index)\n",
395 "\n",
396 "df_forecast2 = df_forecast2.to_timestamp(freq='d')\n",
397 "fig, ax = plt.subplots(figsize = (10,5))\n",
398 "\n",
399 "plt.plot(train[-200:], label = 'actual data', color=\\"blue\\" )\n",
400 "plt.plot(df_forecast2, label= 'forecasts', color = 'red')\n",
401 "plt.legend(loc=\\"upper left\\")\n",
402 "plt.show()\n"
403 ],
404 },
405 {
406 "cell_type": "code",
407 "execution_count": 25,
408 "id": "fbabbd7a",
409 "metadata": {},
410 "outputs": [
411 {
412 "name": "stdout",
413 "output_type": "stream",
414 "text": [
415 "mean absolute percentage error on test data using box-cox transformer is: 0.12534766373642992\n",
416 "mean absolute percentage error on test data without using box-cox transformer is: 0.18088828584131061\n"
417 ]
418 }
419 ],
420 "source": [
421 "# TBATS forecasting model evaluation: comparison\n",
422 "from sklearn.metrics import mean_absolute_percentage_error\n",
423 "print('mean absolute percentage error on test data using box-cox transformer is:', mean_absolute_percentage_error(y_pred2))\n",
424 "print('mean absolute percentage error on test data without using box-cox transformer is:', mean_absolute_percentage_error(y_pred))\n"
425 ],
426 },
427 {
428 "cell_type": "code",
429 "execution_count": null,
430 "id": "b4f011e8",
431 "metadata": {},
432 "outputs": [],
433 "source": []
434 }
435 ],
436 "metadata": {
437 "kernelspec": {
438 "display_name": "Python 3 (ipykernel)",
439 "language": "python",
440 "name": "python3"
441 },
442 "language_info": {

```

```

443     "codemirror_mode": {
444         "name": "ipython",
445         "version": 3
446     },
447     "file_extension": ".py",
448     "mimetype": "text/x-python",
449     "name": "python",
450     "nbconvert_exporter": "python",
451     "pygments_lexer": "ipython3",
452     "version": "3.9.13"
453 }
454 },
455 "nbformat": 4,
456 "nbformat_minor": 5
457 }

```

## Multivariate Zeitreihenvorhersage

Eine multivariate Zeitreihe besteht aus mehreren voneinander abhängigen Variablen, wobei jede Variable nicht nur eine Zeitabhängigkeit von ihren historischen Daten, sondern auch von den vergangenen Werten anderer Variablen aufweist. Multivariate Zeitreihenvorhersagen finden Anwendung in vielen Bereichen wie Energie, Aerologie, Meteorologie, Finanzen, Transportwesen usw. Eine der am häufigsten verwendeten Methoden für multivariate Zeitreihenvorhersagen ist die Vektorautoregression (VAR). In einem VAR-Algorithmus ist jeder Variablenwert eine lineare Beziehung der vergangenen Werte beider: der Variablen selbst und aller anderen Variablen, die als Vektoren dargestellt werden: Gewichtsvektor, Fehlervektor usw. Die VAR-Gleichung wird als beschrieben :

$$y_t = a + w_1 * y_{t-1} + \dots + w_p * y_{t-p} + e$$

wobei  $a$ ,  $w$ ,  $y$ ,  $e$  alle Vektoren sind und daher die vorherige Gleichung wie folgt ausgedrückt wird:

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} + \begin{bmatrix} w_{11} & \dots & w_{1n} \\ \vdots & \vdots & \vdots \\ w_{n1} & \dots & w_{nn} \end{bmatrix} \begin{bmatrix} y_1(t-1) \\ \vdots \\ y_n(t-1) \end{bmatrix} + \dots + \begin{bmatrix} w'_{11} & \dots & w'_{1n} \\ \vdots & \vdots & \vdots \\ w'_{n1} & \dots & w'_{nn} \end{bmatrix} \begin{bmatrix} y_1(t-p) \\ \vdots \\ y_n(t-p) \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

Dabei ist  $P$  die Anzahl der vergangenen Verzögerungen, die zur Berechnung der Prognosen verwendet wurden,  $n$  die Anzahl der voneinander abhängigen Variablen und  $\varepsilon$  der Vektor des weißen Rauschens.

## Multivariate Zeitreihenvorhersage Übung

Wir beginnen mit dem Import der Zeitreihen. Wir werden zwei Datensätze, Geld und Ausgaben, zu einer Zeitreihe mit täglichen Aufzeichnungen zusammenführen. Anschließend werden wir anhand dieser Zeitreihen eine Autoregression und ein multivariates LSTM-Zitreihenvorhersagemodell erstellen. Lasst uns die Datensätze importieren und einen Blick darauf werfen

```
In [3]: # downloading dataset
import pandas as pd
money = pd.read_csv('M2SLMoneyStock.csv')
spending = pd.read_csv('PCEPersonalSpending.csv')
df = pd.merge(money, spending, on='Date', how='outer')
# set index to Date
df = df.set_index('Date')
# convert index into datetime format
df.index = pd.DatetimeIndex(df.index).to_period('m')
df
```

Out[3]:

	Money	Spending
Date		
1995-01	3492.4	4851.2
1995-02	3489.9	4850.8
1995-03	3491.1	4885.4
1995-04	3499.2	4890.2
1995-05	3524.2	4933.1
...	...	...
2015-08	12096.8	12394.0
2015-09	12153.8	12392.8
2015-10	12187.7	12416.1
2015-11	12277.4	12450.1

Wir werden 90 % der Zeitreihen zum Training der Modelle verwenden. Wir werden ein AR-Modell unter Verwendung der letzten 5 Werte jeder Variablen erstellen. Lassen Sie uns die Zeitreihe aufteilen, die erforderlichen Python-Pakete importieren, das Modell an den Trainingssatz anpassen und schließlich das angepasste Modell testen.

```
In [5]: # splitting the data into train and validation sets: 90% for training, 10% for validation
data = df.dropna()
train = data[:int(0.9*(len(df)))]
test = data[int(0.9*(len(df))):]

In [6]: # using vector autoregression model
from statsmodels.tsa.vector_ar.var_model import VAR
# a pth-order VAR refers to a VAR model which includes lags for the last p time periods
# fit(maxlags): Maximum number of lags to check for order selection,
# defaults to 12 * (nobs/100.)**(1./4)
# we set order = 5
model = VAR(endog=train, freq='M')
fitted_var = model.fit(5)# fitted_var.summary(): to access model output

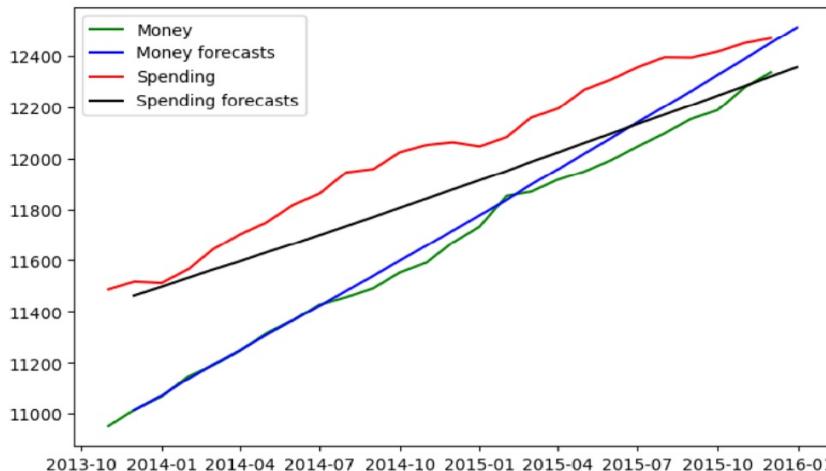
# make predictions and evaluate the model on validation set
# VARResults.forecast: Produce linear minimum MSE forecasts for desired number of steps ahead, using prior values y
lagged_Vals = train.values[-10:]# The VAR .forecast() requires passing in a lag order number of previous observations.
predictions = fitted_var.forecast(y=lagged_Vals, steps=len(test))
```

Jetzt werten wir das angepasste Autoregressionsmodell am Testsatz aus und zeichnen die tatsächlichen Werte im Vergleich zu den Vorhersagen für beide Variablen auf: „Money“ und „Spending“.

```
In [7]: #the predictions are made as an array, where each list represents the predictions of a row in the dataset.
# we need first to convert the array into a dataframe
from sklearn.metrics import mean_absolute_percentage_error
# Locate the predictions in a new dataframe with the same col names: without 'CO(GT)'
cols = df.columns
idx = test.index
pred = pd.DataFrame(index=idx, columns=[cols])
for j in range(0,2):
    for i in range(0, len(predictions)):
        pred.iloc[i][j] = predictions[i][j]

#calculate MAPE for all cols
for col in cols:
    print('MAPE value for', col, 'is : ', mean_absolute_percentage_error(pred[col], test[col]))
```

MAPE value for Money is : 0.008557826316923147  
MAPE value for Spending is : 0.009869428903370645



Jetzt verwenden wir den LSTM-Algorithmus für denselben Datensatz, um ein multivariates Zeitreihenmodell zu erstellen. Zuerst müssen wir die Zeitreihe mit dem StandardScaler-Transformator in eine Standardverteilung umwandeln. Wir werden die Variable „Geld“ als Zielvariable des LSTM-Modells verwenden:

```
In [15]: # using LSTM;
# first we need to scale input and output of LSTM using standard scaler
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
x_train_scaled = scaler.fit_transform(train)
x_train_scaled.shape
```

Out[15]: (226, 2)

```
In [319]: # to predict money variable
scaler2= StandardScaler()
y_train_scaled =scaler2.fit_transform(train[['Money']])
y_train_scaled.shape
```

Out[319]: (226, 1)

Ein wichtiger Schritt besteht darin, die Eingabe- und Ausgabesequenzen für LSTM vorzubereiten. Wir werden die gleiche Anzahl von Verzögerungen wie für AR verwenden, also 5 Verzögerungen:

```
In [19]: # for lag= 5 we need to create the sequences of the input array manually;reshaping input data
import numpy as np
x_train=[]
y_train = []
for i in range(5, 226):
    x_train.append(x_train_scaled[i-5:i])
    y_train.append(y_train_scaled[i][0])
x_train= np.array(x_train)
y_train = np.array(y_train)
print(x_train.shape)

(221, 5, 2)
```

Jetzt erstellen wir unser LSTM-Modell, das aus einer Eingabeschicht, einer verborgenen Schicht, zwei Dropout-Schichten und einer Ausgabeschicht besteht. Lassen Sie uns das Modell kompilieren und an Trainingsdaten anpassen:

```
In [21]: # build LSTM model
from tensorflow import keras
from tensorflow.keras.layers import Dropout, Dense, LSTM
from tensorflow.keras.models import Sequential

model = keras.Sequential()
model.add(LSTM(units=50, return_sequences=True, input_shape=(5,2)))
model.add(Dropout(0.2))
model.add(LSTM(units=50))
model.add(Dropout(0.2))
model.add(Dense(1))
```

```
In [22]: model.compile(loss='mean_squared_error',optimizer='adam')
history = model.fit(x_train, y_train, epochs=200, batch_size= 32, shuffle=False, validation_split=0.2)
6/6 [=====] - 0s 10ms/step - loss: 0.0078 - val_loss: 0.0078
Epoch 168/200
6/6 [=====] - 0s 10ms/step - loss: 0.0064 - val_loss: 0.0371
Epoch 169/200
6/6 [=====] - 0s 9ms/step - loss: 0.0053 - val_loss: 0.0170
Epoch 170/200
6/6 [=====] - 0s 9ms/step - loss: 0.0054 - val_loss: 0.0192
Epoch 171/200
6/6 [=====] - 0s 11ms/step - loss: 0.0058 - val_loss: 0.0295
Epoch 172/200
6/6 [=====] - 0s 11ms/step - loss: 0.0053 - val_loss: 0.0368
Epoch 173/200
6/6 [=====] - 0s 11ms/step - loss: 0.0066 - val_loss: 0.0092
Epoch 174/200
6/6 [=====] - 0s 9ms/step - loss: 0.0061 - val_loss: 0.0344
Epoch 175/200
6/6 [=====] - 0s 9ms/step - loss: 0.0074 - val_loss: 0.0133
Epoch 176/200
6/6 [=====] - 0s 9ms/step - loss: 0.0041 - val_loss: 0.0094
```

Um das Modell am Testsatz auszuwerten, müssen wir die letzten 5 Werte des Trainingssatzes anhängen, da das Modell Vorhersagen mit 5 Verzögerungen erzeugt. Die Testsequenz sollte skaliert und vorbereitet werden, bevor Vorhersagen getroffen werden, die wiederum ebenfalls umgekehrt skaliert werden sollten. `model.predict()` wird verwendet, um Vorhersagen von Testzeitreihen zu berechnen:

```
In [32]: # appending the last 10 of train set to the test set
train_last_10= train[-5:]
extended_test= pd.concat((train_last_10, test), axis=0)
# rescaling test data
scaled_test= scaler.fit_transform(extended_test)
scaled_test.shape

# rescaling y
y_test_scaled = scaler2.fit_transform(extended_test[['Money']])
y_test_scaled.shape
```

Out[32]: (31, 1)

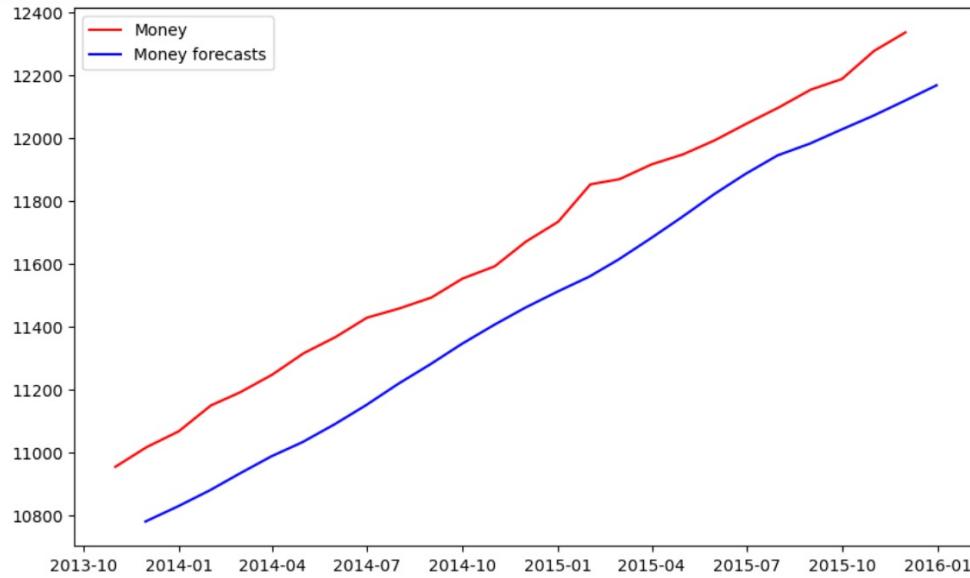
```
In [33]: # reshaping input data for testing the model
x_test=[]
for i in range(5, 31):
    x_test.append(scaled_test[i-5:i])
x_test= np.array(x_test)
print(x_test.shape)

(26, 5, 2)
```

```
In [34]: # make predictions
y_test=model.predict(x_test)
# inverse scaling of money forecasts
y= scaler2.inverse_transform(y_test)
y
```

Überprüfen wir die Darstellung der Vorhersagen anhand der tatsächlichen Geldwerte:

```
In [35]: # plotting predictions against actual money values
fig = plt.subplots(figsize=(12,8))
test2 = test.to_timestamp(freq='d')
df_forecast=pd.DataFrame(data=y, index=id)
plt.plot(test2['Money'],color="red", label ='Money')
plt.plot(df_forecast, color="blue", label ='Money forecasts')
plt.legend(loc="upper left")
plt.show()
```



Lassen Sie uns den Vorhersagefehlerwert für die Variable „money“ überprüfen:

```
In [36]: # MAPE between forecasts and actual money
print('MAPE value for Money is:', mean_absolute_percentage_error(y, test['Money']))
```

MAPE value for Money is: 0.014847821070793958

```

1 {
2     "cells": [
3         {
4             "cell_type": "code",
5             "execution_count": null,
6             "id": "a833112b",
7             "metadata": {},
8             "outputs": [],
9             "source": [
10                 "# downloading dataset\n",
11                 "import pandas as pd\n",
12                 "money = pd.read_csv('M2SLMoneyStock.csv')\n",
13                 "spending = pd.read_csv('PCEPersonalSpending.csv')\n",
14                 "df = pd.merge(money, spending, on='Date', how='outer')\n",
15                 "#df2 = pd.merge(passengers, restaurant, on='Date', how='outer')\n",
16                 "# set index to Date\n",
17                 "df = df.set_index('Date')\n",
18                 "# convert index into datetime format\n",
19                 "df.index = pd.DatetimeIndex(df.index).to_period('m')\n",
20                 "df"
21             ]
22         },
23     {
24         "cell_type": "code",
25         "execution_count": null,
26         "id": "85e39b7d",
27         "metadata": {},
28         "outputs": [],
29         "source": [
30             "#checking stationarity, \n",
31             "# Cointegration is a statistical method used to test the correlation between two or more non-stationary time series \n",
32             "# in the long run or for a specified period\n",
33             "from statsmodels.tsa.vector_ar.vecm import coint_johansen\n",
34             "#since the test works for only 12 variables, I have randomly dropped\n",
35             "#in the next iteration, I would drop another and check the eigenvalues\n",
36             "\n",
37             "#coint_johansen(endog, det_order, k_ar_diff); k_ar_diff Number of lagged differences in the model.\n",
38             "# det_order: determination order; -1 no-deterministic terms: the order of time polynomial in the null-hypothesis\n",
39             "data= df.dropna()\n",
40             "# Eigenvalues of coefficient matrix\n",
41             "result=coint_johansen(data,-1,1).eig\n",
42             "print(result)"
43         ]
44     },
45     {
46         "cell_type": "code",
47         "execution_count": null,
48         "id": "28clf9fb",
49         "metadata": {},
50         "outputs": [],
51         "source": [
52             "df"
53         ]
54     },
55     {
56         "cell_type": "code",
57         "execution_count": null,
58         "id": "647d01ba",
59         "metadata": {},
60         "outputs": [],
61         "source": [
62             "# splitting the data into train and validation sets: 90% for training, 10% for validation\n",

```

```

63     "data= df.dropna()\n",
64     "train = data[:int(0.9*(len(df)))]\n",
65     "test = data[int(0.9*(len(df))):]"  

66   },
67 },
68 {
69   "cell_type": "code",
70   "execution_count": null,
71   "id": "ed8eb81e",
72   "metadata": {},
73   "outputs": [],
74   "source": [
75     "# using vector autogression model \n",
76     "from statsmodels.tsa.vector_ar.var_model import VAR\n",
77     "# a pth-order VAR refers to a VAR model which includes lags for the last p time periods\n",
78     "# fit(maxlags): Maximum number of lags to check for order selection,\n",
79     "# defaults to 12 * (nobs/100.)**(1./4)\n",
80     "# we set order = 5 \n",
81     "model = VAR(endog=train, freq='M')\n",
82     "fitted_var = model.fit(5)# fitted_var.summary(): to access model output\n",
83     "\n",
84     "\n",
85     "# make predictions and evaluate the model on validation set\n",
86     "# VARResults.forecast: Produce linear minimum MSE forecasts for desired number of steps ahead, using prior values y\n",
87     "lagged_Vals = train.values[-10:]# The VAR .forecast() requires passing in a lag order number of previous observations.\n",
88     "predictions = fitted_var.forecast(y=lagged_Vals, steps=len(test))"  

89   ],
90 },
91 {
92   "cell_type": "code",
93   "execution_count": null,
94   "id": "b81c5df3",
95   "metadata": {},
96   "outputs": [],
97   "source": [
98     "# the predictions are made as an array, where each list represents the predictions of a row in the dataset.\n",
99     "# we need first to convert the array into a dataframe\n",
100    "from sklearn.metrics import mean_absolute_percentage_error\n",
101    "# locate the predictions in a new dataframe with the same col names: without 'CO(GT) '\n",
102    "cols =df.columns\n",
103    "idx = test.index\n",
104    "pred = pd.DataFrame(index=idx,columns=[cols])\n",
105    "for j in range(0,2):\n",
106      " for i in range(0, len(predictions)):\n",
107        " pred.iloc[i][j] = predictions[i][j]\n",
108      "\n",
109      "#calculate MAPE for all cols\n",
110      "for col in cols:\n",
111        " print('MAPE value for', col, 'is : ', mean_absolute_percentage_error(pred[col], test[col]))"  

112   ],
113 },
114 {
115   "cell_type": "code",
116   "execution_count": null,
117   "id": "26c20cfb",
118   "metadata": {},
119   "outputs": [],
120   "source": [
121     "#plot cols against predictions\n",
122     "import matplotlib.pyplot as plt\n",
123     "fig = plt.subplots(figsize=(8,5))\n",
124     "test2 = test.to_timestamp(freq='d')\n",
125     "\n",
126     "df_forecast=pd.DataFrame(data=pred, index=idx, columns=cols)\n",
127     "plt.plot(test2['Money'],color=\"green\", label ='Money')\n",
128     "plt.plot(pred['Money'], color=\"blue\", label ='Money forecasts')\n",
129     "plt.plot(test2['Spending'],color=\"red\", label ='Spending')\n",
130     "plt.plot(pred['Spending'], color=\"black\", label ='Spending forecasts')\n",
131     "plt.legend(loc="upper left")\n",
132     "plt.show()\n",
133     "# since each variable is a linear function of past lags of itself and past lags of the other variables, we get fitted 1.  

134   ],
135 },
136 {
137   "cell_type": "code",
138   "execution_count": null,
139   "id": "eedc9c9f",
140   "metadata": {},
141   "outputs": [],
142   "source": [
143     "# using LSTM;\n",
144     "# first we need to scale input and output of LSTM using standard scaler\n",
145     "from sklearn.preprocessing import StandardScaler\n",
146     "scaler = StandardScaler()\n",
147     "x_train_scaled =scaler.fit_transform(train)\n",
148     "x_train_scaled.shape "
149   ],
150 },
151 {
152   "cell_type": "code",
153   "execution_count": null,
154   "id": "514e95e2",
155   "metadata": {},
156   "outputs": [],
157   "source": [
158     "#using money variable as input time series\n",
159     "scaler2= StandardScaler()\n",
160     "y_train_scaled =scaler2.fit_transform(train[['Money']])\n",
161     "y_train_scaled.shape"
162   ],
163 },
164 {
165   "cell_type": "code",
166   "execution_count": null,

```

```

167     "id": "bb46576b",
168     "metadata": {},
169     "outputs": [],
170     "source": [
171         "# for lag= 5 we need to create the sequences of the input array manually;reshaping input data\n",
172         "import numpy as np\n",
173         "x_train=[]\n",
174         "y_train = []\n",
175         "for i in range(5, 226):\n",
176             x_train.append(x_train_scaled[i-5:i])\n",
177             y_train.append(y_train_scaled[i][0])\n",
178         "x_train= np.array(x_train)\n",
179         "y_train = np.array(y_train)\n",
180         "print(x_train.shape)"\n
181     ],
182 },
183 {
184     "cell_type": "code",
185     "execution_count": null,
186     "id": "dca9df45",
187     "metadata": {},
188     "outputs": [],
189     "source": [
190         "# build LSTM model \n",
191         "from tensorflow import keras\n",
192         "from tensorflow.keras.layers import Dropout, Dense, LSTM\n",
193         "from tensorflow.keras.models import Sequential\n",
194         "\n",
195         "model = keras.Sequential()\n",
196         "model.add(LSTM(units=50, return_sequences=True, input_shape=(5,2)))\n",
197         "model.add(Dropout(0.2))\n",
198         "model.add(LSTM(units=50))\n",
199         "model.add(Dropout(0.2))\n",
200         "model.add(Dense(1))"\n
201     ],
202 },
203 {
204     "cell_type": "code",
205     "execution_count": null,
206     "id": "16544270",
207     "metadata": {},
208     "outputs": [],
209     "source": [
210         "model.compile(loss='mean_squared_error',optimizer='adam')\n",
211         "history = model.fit(x_train, y_train, epochs=200, batch_size= 32, shuffle=False, validation_split=0.2)"\n
212     ],
213 },
214 {
215     "cell_type": "code",
216     "execution_count": null,
217     "id": "43f649a7",
218     "metadata": {},
219     "outputs": [],
220     "source": [
221         "# appending the last 5 of train set to the test set: to be used for forecasting, since we use 5 lags \n",
222         "train_last_5= train[-5:]",
223         "extended_test= pd.concat((train_last_5, test), axis=0)\n",
224         "# rescaling test data\n",
225         "scaled_test= scaler.fit_transform(extended_test)\n",
226         "scaled_test.shape\n",
227         "\n",
228         "# rescaling y \n",
229         "y_test_scaled = scaler2.fit_transform(extended_test[['Money']])\n",
230         "y_test_scaled.shape"
231     ],
232 },
233 {
234     "cell_type": "code",
235     "execution_count": null,
236     "id": "2a5c8162",
237     "metadata": {},
238     "outputs": [],
239     "source": [
240         "# reshaping input data for testing the model\n",
241         "x_test=[]\n",
242         "for i in range(5, 31):\n",
243             x_test.append(scaled_test[i-5:i])\n",
244         "x_test= np.array(x_test)\n",
245         "print(x_test.shape)"\n
246     ],
247 },
248 {
249     "cell_type": "code",
250     "execution_count": null,
251     "id": "ba2490c1",
252     "metadata": {},
253     "outputs": [],
254     "source": [
255         "# forecasting on test data\n",
256         "y_test=model.predict(x_test)\n",
257         "# inverse scaling of money forecasts\n",
258         "y= scaler2.inverse_transform(y_test)\n",
259         "y"
260     ],
261 },
262 {
263     "cell_type": "code",
264     "execution_count": null,
265     "id": "bfd59e08",
266     "metadata": {},
267     "outputs": [],
268     "source": [
269         "# plotting predictions against actual money values\n",
270         "fig = plt.subplots(figsize=(10,6))\n",

```

```

271 "test2 = test.to_timestamp(freq='d')\n",
272 "df_forecast=pd.DataFrame(data=y, index=idx)\n",
273 "plt.plot(test2['Money'],color=\"red\", label ='Money')\n",
274 "plt.plot(df_forecast, color=\"blue\", label ='Money forecasts')\n",
275 "plt.legend(loc=\"upper left\")\n",
276 "plt.show()\n"
277 ]
278 },
279 {
280 "cell_type": "code",
281 "execution_count": null,
282 "id": "1e7d73eb",
283 "metadata": {},
284 "outputs": [],
285 "source": [
286 "# MAPE between forecasts and actual money\n",
287 "print('MAPE value for Money is:', mean_absolute_percentage_error(y, test['Money']))\n"
288 ]
289 }
290 ],
291 "metadata": {
292 "kernelspec": {
293 "display_name": "Python 3 (ipykernel)",
294 "language": "python",
295 "name": "python3"
296 },
297 "language_info": {
298 "codemirror_mode": {
299 "name": "ipython",
300 "version": 3
301 },
302 "file_extension": ".py",
303 "mimetype": "text/x-python",
304 "name": "python",
305 "nbconvert_exporter": "python",
306 "pygments_lexer": "ipython3",
307 "version": "3.9.13"
308 },
309 },
310 "nbformat": 4,
311 "nbformat_minor": 5
312 }

```

Die resultierenden Fehlerwerte beider Modelle weisen auf eine höhere Leistung des Autoregressionsvektormodells als des LSTM-Modells hin.