

Einführung

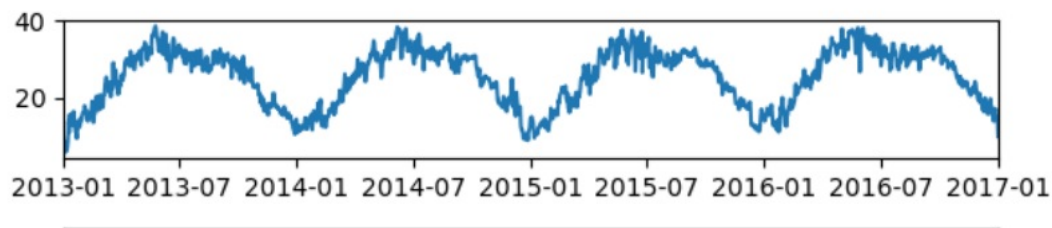
Unter Zeitreihenprognosen versteht man die Vorhersage zukünftiger Werte auf der Grundlage historischer, mit Zeitstempeln versehener Daten. Die Analyse historischer Daten kann uns helfen, zukünftige Trends zu verstehen und das Verhalten einer bestimmten Variablen in der Zukunft vorherzusagen. Im Vergleich zu Standardtechniken des maschinellen Lernens wie Random Forest und zeitverzögerten neuronalen Netzen verfügt die Zeitreihenprognosemethodik über die Fähigkeit, Muster außerhalb der Trainingsdaten zu extrapolieren und übertrifft damit die meisten Algorithmen des maschinellen Lernens, die nur Muster aus dem Trainingsdatensatz lernen können. Hier kommen Techniken zur Zeitreihenprognose ins Spiel.

Zeitreihenvorhersagen werden in vielen industriellen und wissenschaftlichen Bereichen eingesetzt, darunter Steuerungstechnik, Geschäftsplanung, Modellierung der Ausbreitung von Krankheiten, Signalverarbeitung und Wettervorhersage.

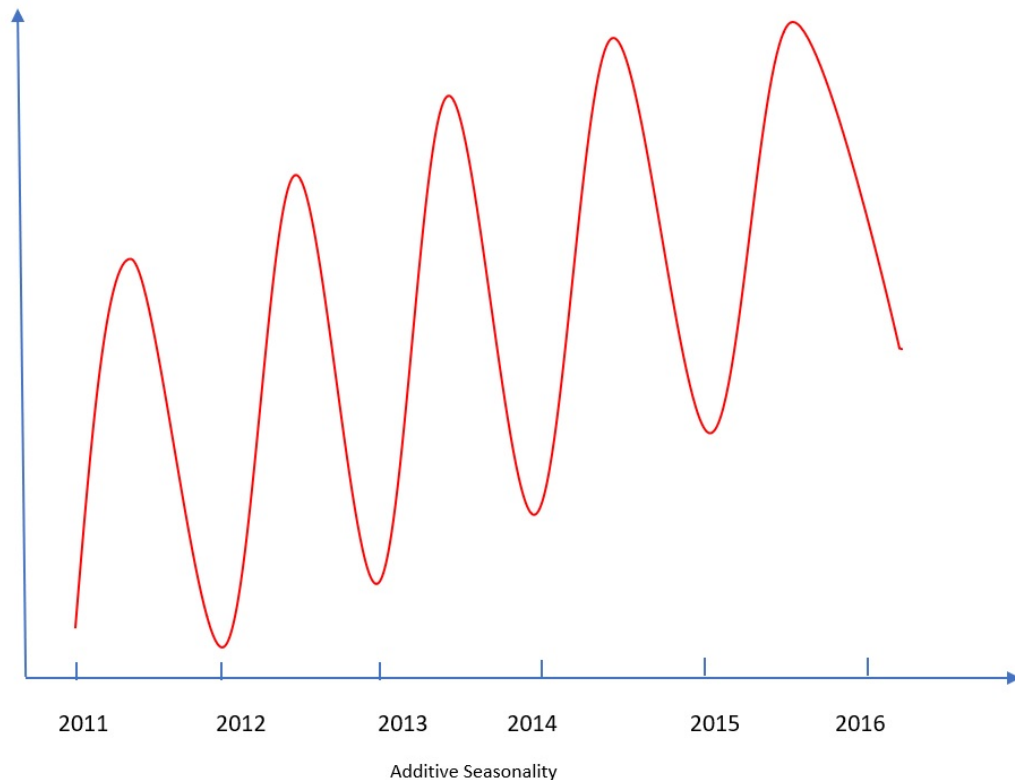
Für die Zeitreihenvorhersage werden drei Hauptbegriffe analysiert:

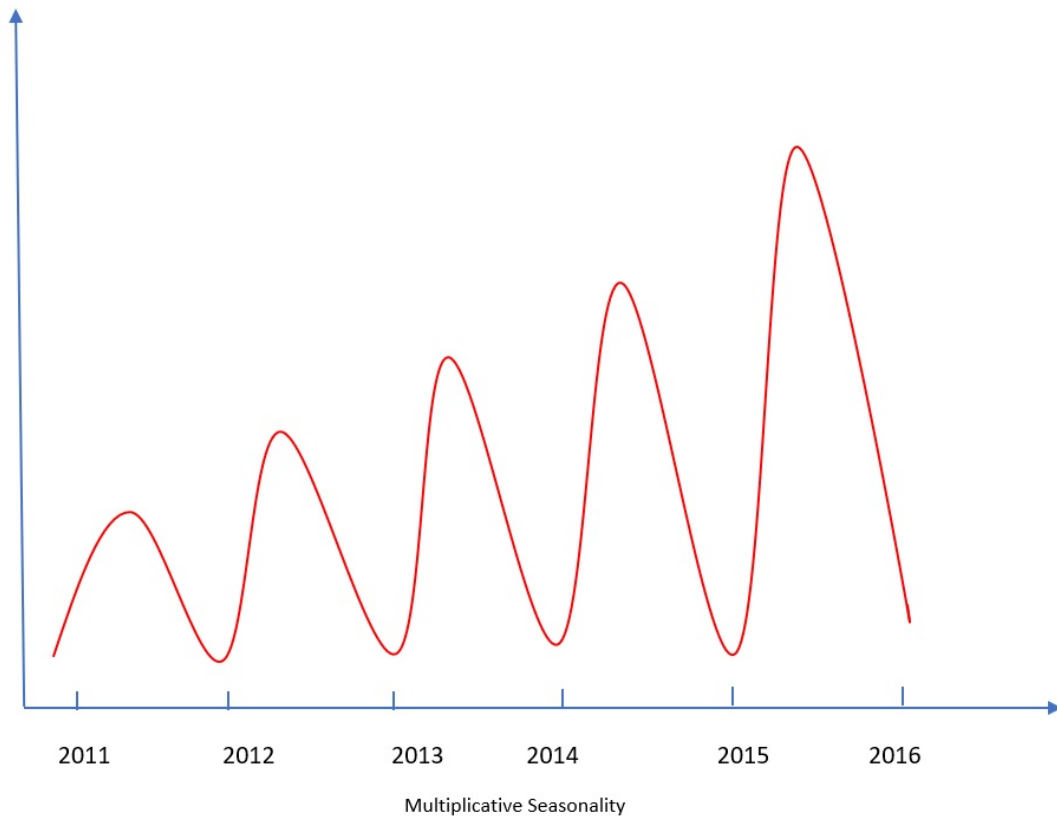
Saisonalität:

Unter Saisonalität versteht man die Wiederkehrenden Datenmuster zu bestimmten Zeitpunkten. Mit anderen Worten: Saisonalitäten sind Pattern, die sich regelmäßig in Zeitreihen über bestimmte Zeiträume wie Wochen oder Monate wiederholen. beispielsweise um Spitzenwerte in den Daten von Reiseunternehmen im November und Dezember aufgrund der Ferienzeit zu beobachten. Die Vorhersage von Zeitreihendaten erfordert die Analyse und Erfassung der Saisonalität. Die folgenden Abbildungen zeigen beispielhaft einen jährlichen Saisonzyklus:



Wir haben zwei Haupttypen der Saisonalität: die additive und die multiplikative Saisonalität. In additiven Saisonperioden zeigt der Trend Linearität, wobei die Amplitude der Saisonalität über die Zeit gleich bleibt, während wir bei multiplikativer Saisonalität Veränderungen in der Breite oder Höhe von Saisonperioden feststellen werden und der Trend über die Zeit nicht linear verläuft. Die folgenden Abbildungen zeigen den Unterschied zwischen den beiden Typen:





Trend:

Der Trend ist ein weiterer wichtiger Faktor bei der Zeitreihenprognose. Er beschreibt die Zunahme oder Abnahme einer Variablen in einem bestimmten Zeitraum. Wenn beispielsweise der Umsatz eines Unternehmens über einen bestimmten Zeitraum steigt, weist die Umsatzvariable einen steigenden Trend auf.

Zyklische Variation:

Die Werte der Daten weisen im Zeitverlauf Schwankungen auf, die keine feste Häufigkeit haben. Diese seltenen Änderungen in den Daten werden als Zyklen bezeichnet und treten häufig aufgrund wirtschaftlicher Bedingungen auf.

Unregelmäßige Variation:

Sie stellen die Schwankungen dar, die völlig zufällig wirken oder unregelmäßig und haben kein bestimmtes Muster. Die Schwankungen sind nicht vorhersehbar, müssen jedoch bei der Erstellung von Zeitreihenmodellen berücksichtigt werden.

Long Short Term Memory (LSTM)

Die Annahme einer linearen Beziehung zwischen Eingaben und Ausgaben trifft nicht unbedingt auf jede Zeitreihe zu und daher entstand der Bedarf an Algorithmen, die jede nichtlineare Funktion für Zeitreihenprognosen annähern können, ohne vorherige Kenntnisse über die Merkmale der Daten. Long Short-Term Memory (LSTM) ist eines der beliebtesten Zeitreihen-Prognosemodelle, das seine Fähigkeit unter Beweis gestellt hat, anhand historischer Daten zukünftige Prognosen zu erstellen. Dank des Langzeitgedächtnisses ist LSTM in der Lage, noch mehr Parameter und Langzeittrends in der Zeitreihe effizient zu lernen.

In diesem [Kapitel](#) wird ausführlich auf die Architektur neuronaler LSTM-Netze eingegangen.

LSTM Übung

Wir werden die [hier](https://www.kaggle.com/datasets/sumanthvrao/daily-climate-time-series-data) verfügbaren täglichen Klimazeitreihendaten verwenden: <https://www.kaggle.com/datasets/sumanthvrao/daily-climate-time-series-data>.

LSTM wird auf die Spalte „meantemp“ angewendet, um zukünftige Werte vorherzusagen. Werfen wir einen Blick auf den Datensatz:

```
In [37]: import scalecast
import tensorflow
import pandas as pd
import numpy as np
import pickle
import seaborn as sns
import matplotlib.pyplot as plt
from scalecast.Forecaster import Forecaster
# If a column or index cannot be represented as an array of datetimes,
# say because of an unparseable value or a mixture of timezones,
# the column or index will be returned unaltered as an object data type.

df = pd.read_csv('DailyDelhiClimateTrain.csv', parse_dates=['date'])
# take a look at the data
df.head()
#len(df)
```

```
Out[37]:
```

	date	meantemp	humidity	wind_speed	meanpressure
0	2013-01-01	10.000000	84.500000	0.000000	1015.666667
1	2013-01-02	7.400000	92.000000	2.980000	1017.800000
2	2013-01-03	7.166667	87.000000	4.633333	1018.666667
3	2013-01-04	8.666667	71.333333	1.233333	1017.166667
4	2013-01-05	6.000000	86.833333	3.700000	1016.500000

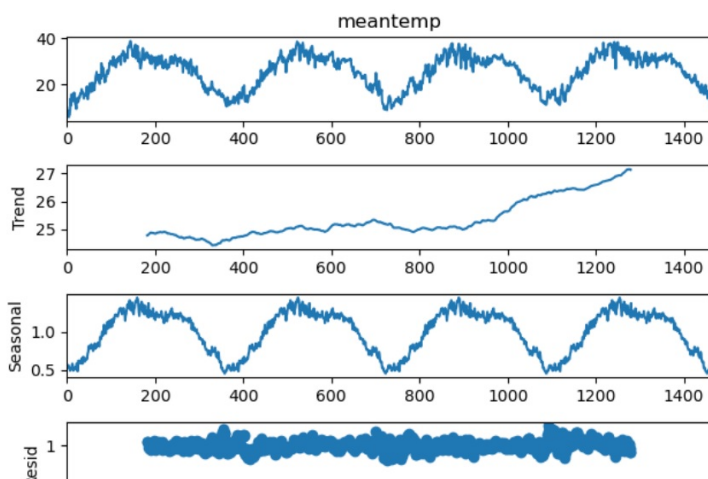
Es ist nicht immer einfach, ein leistungsstarkes LSTM-Modell zu entwerfen und zu implementieren. Aus diesem Grund stellen wir die Scalecast-Bibliothek vor, die ein Wrapper-Prognosestool für viele Zeitreihen-Prognosemodelle wie Arima und LSTM ist. Es verwendet TensorFlow LSTM und implementiert alle Datenverarbeitungsschritte, die zum Erstellen des Modells erforderlich sind, einschließlich Zerlegung der Zeitreihe in ihre Hauptbestandteile, Scaling, Unskalierung, Anpassung und Testen des Modells. Mit Scalecast Forecaster lässt sich LSTM viel einfacher anwenden und Zukunftsprognosen erstellen.

```
In [38]: #we must first call the Forecaster object with the y and current_dates parameters as 'meantemp' and 'date' variable specified
f = Forecaster(y=df['meantemp'], current_dates=df['date'])
f
```

```
Out[38]: Forecaster(
  DateStartActuals=2013-01-01T00:00:00.000000000
  DateEndActuals=2017-01-01T00:00:00.000000000
  Freq=D
  N_actuals=1462
  ForecastLength=0
  Xvars=[]
  TestLength=0
  ValidationLength=1
  ValidationMetric=rmse
  ForecastsEvaluated=[]
  CILevel=None
  CurrentEstimator=mlr
  GridsFile=Grids
)
```

Bevor wir ein LSTM-Modell erstellen, sollten wir die Hauptkomponenten der Zeitreihe wie Saisonalität und Trend überprüfen. Diese Komponenten helfen uns, ein Prognosemodell effizienter zu erstellen und können andere Merkmale wie Stationarität in Zeitreihen zeigen. Stationarität bedeutet, dass sich die Eigenschaften von Zeitreihen unabhängig von dem Zeitpunkt ändern, zu dem die Zeitreihe beobachtet wird, und dass eine stationäre Zeitreihe daher keine signifikante Saisonalität oder einen signifikanten Trend aufweist, da diese Komponenten die Zeitreihe zu bestimmten Zeitpunkten beeinflussen. Daher dient die Untersuchung des saisonalen Zyklus und der Trends dazu, Einblicke in die Stationarität von Zeitreihen zu gewinnen. Um die Zeitreihe in ihre Hauptkomponenten zu zerlegen, verwenden wir `seasonal_decompose` Funktion von Statsmodels Bibliothek wie im Folgenden:

```
In [17]: # Let's further decompose the series into its trend, seasonal, and residual parts:
from matplotlib import pyplot
from statsmodels.tsa.seasonal import seasonal_decompose
result = seasonal_decompose(df['meantemp'], model='multiplicative', period=365)# frequency of the observations is 1
result.plot()
pyplot.show()
#The figure obviously indicates yearly seasonality and increasing trend over time.
```



Die vorherige Abbildung zeigt offensichtlich eine erhebliche jährliche Saisonalität und einen zunehmenden Trend im Zeitverlauf.

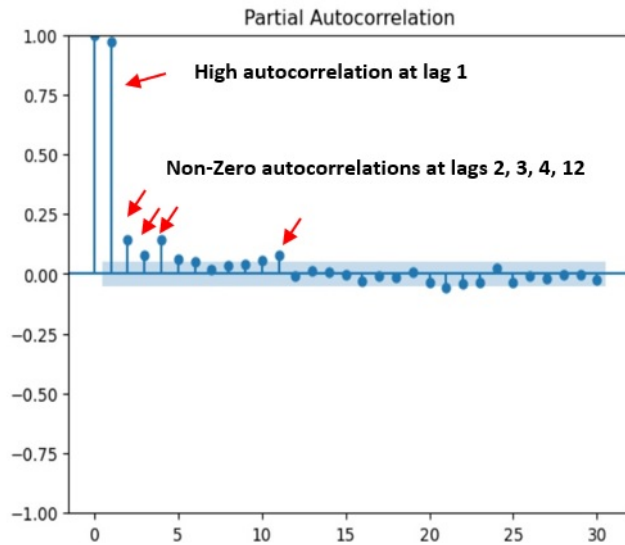
Lassen Sie uns nun die Autokorrelation der Durchschnittstemperatur mit ihren vergangenen Werten überprüfen. Zu diesem Zweck verwenden wir das Partial Auto Correlation Function (PACF)-Diagramm, das misst, wie stark die y-Variable (meantemp) statistisch mit früheren Werten von sich selbst korreliert. Wir prüfen die Autokorrelationen auf bis zu 30 Verzögerungen.

Der blaue Bereich des PACF-Diagramms zeigt den Signifikanzschwellenwert. Das bedeutet, dass Verzögerungen, die sich in diesem Bereich be-

finden, statistisch gesehen nahe Null liegen und daher keine signifikante Autokorrelation zwischen Datenpunkten besteht. PACF stellt intuitiv Korrelationen von 1 bei Verzögerung 0 dar, da dies die Korrelation der Zeitreihe mit sich selbst darstellt.

```
In [39]: f.plot_pacf(lags=30)
plt.show()
```

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()



Testen wir die Stationarität der Reihe mit dem Augmented Dickey-Fuller (ADF)-Test. null Die Hypothese dieses Tests ist, dass die Zeitreihe nicht stationär ist (es gibt eine Einheitswurzel). Wenn der resultierende p_value über einer kritischen Größe liegt (Standard ist 0,05), können wir nicht ausschließen, dass es eine Einheitswurzel gibt und die Reihe daher instationär ist, siehe [hier](#). Der resultierende p_Wert zeigt an, dass unsere Zeitsequenz nicht stationär ist.

```
In [8]: # Let's test the series' stationarity.
# If bool (full_res = False), returns whether the test suggests stationarity.
# null Hypothesis of Augmented Dickey-Fuller (ADF) test: time series is non_stationary ( there is a unit root,)
# If the pvalue is above a critical size (Default is 0.05), then we cannot reject that there is a unit root.
stat = f.adf_test(full_res=True)
print(stat)

# p_value is 0.28 > 0.05, then we cannot reject null hypothesis and thus time series is non stationary

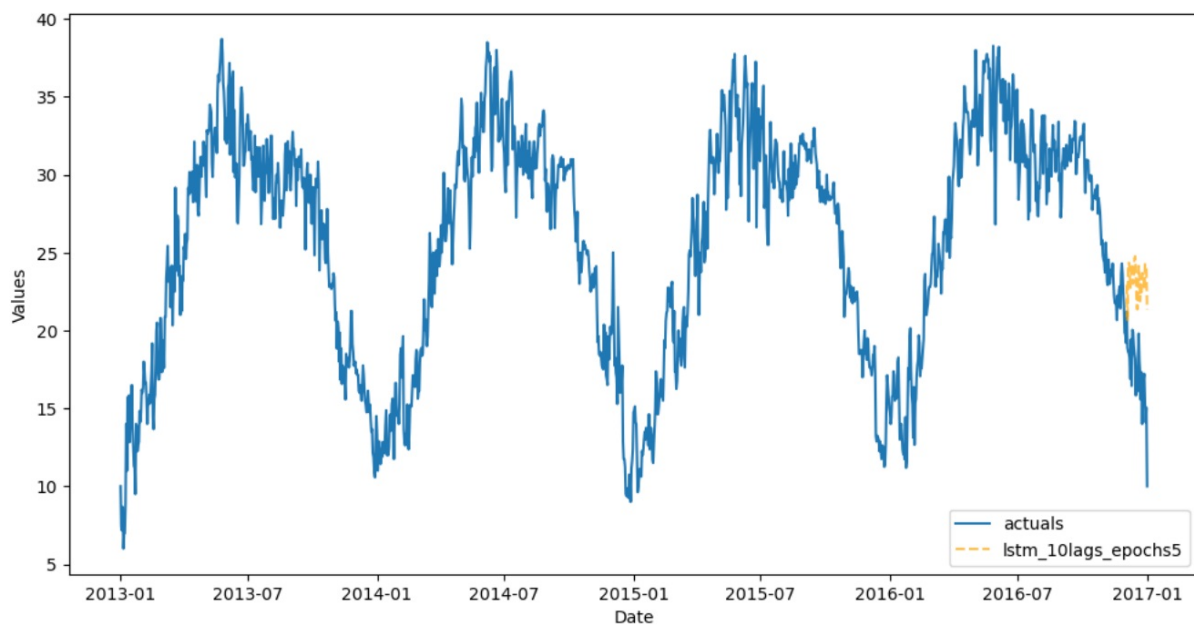
(-2.021069055920671, 0.2774121372301611, 10, 1451, {'1%': -3.4348647527922824, '5%': -2.863533960720434, '10%': -2.567831568508802}, 5423.895746470953)
```

Wir werden unser LSTM mit Scalecast Forecaster erstellen. Standardmäßig, Dieses Modell wird mit einer einzelnen Eingabeebene von 8 Einheiten, Adam-Optimierer, Tanh-Aktivierung, einer Lernrate von 0,001 und ohne Dropout ausgeführt. Wir müssen jedoch einige Parameter festlegen, z. B. Testsatz, Validierungsmetrik, Anzahl der Epochen und Anzahl der Verzögerungen, die für die Vorhersage verwendet werden. für den ersten Lauf werden wir verwenden 10 Verzögerungen und 5 Epochen. Schauen wir uns die resultierenden Diagramme an:

```
In [67]: # Now, to call an LSTM forecast. By default,
# this model will be run with a single input layer of 8 units, Adam optimizer, tanh activation,
# a Learning rate of 0.001, and no dropout.

# generate future dates: The number of dates you generate in this step will determine how long all models will be forecast out.
f.set_validation_metric('mape')
f.set_test_length(30) # 30 observations to test the results
f.generate_future_dates(30) # 30 future points to forecast
f.set_estimator('lstm') # LSTM neural network
f.manual_forecast(call_me='lstm_10lags_epochs5',lags=10, epochs=5)
f.plot_test_set(ci=True)

Epoch 1/5
44/44 [=====] - 1s 3ms/step - loss: 0.5967
Epoch 2/5
44/44 [=====] - 0s 3ms/step - loss: 0.4379
Epoch 3/5
44/44 [=====] - 0s 3ms/step - loss: 0.2333
Epoch 4/5
44/44 [=====] - 0s 3ms/step - loss: 0.1348
Epoch 5/5
44/44 [=====] - 0s 3ms/step - loss: 0.0981
1/1 [=====] - 0s 231ms/step
Epoch 1/5
45/45 [=====] - 1s 3ms/step - loss: 0.5467
Epoch 2/5
45/45 [=====] - 0s 4ms/step - loss: 0.3424
Epoch 3/5
45/45 [=====] - 0s 4ms/step - loss: 0.2424
```



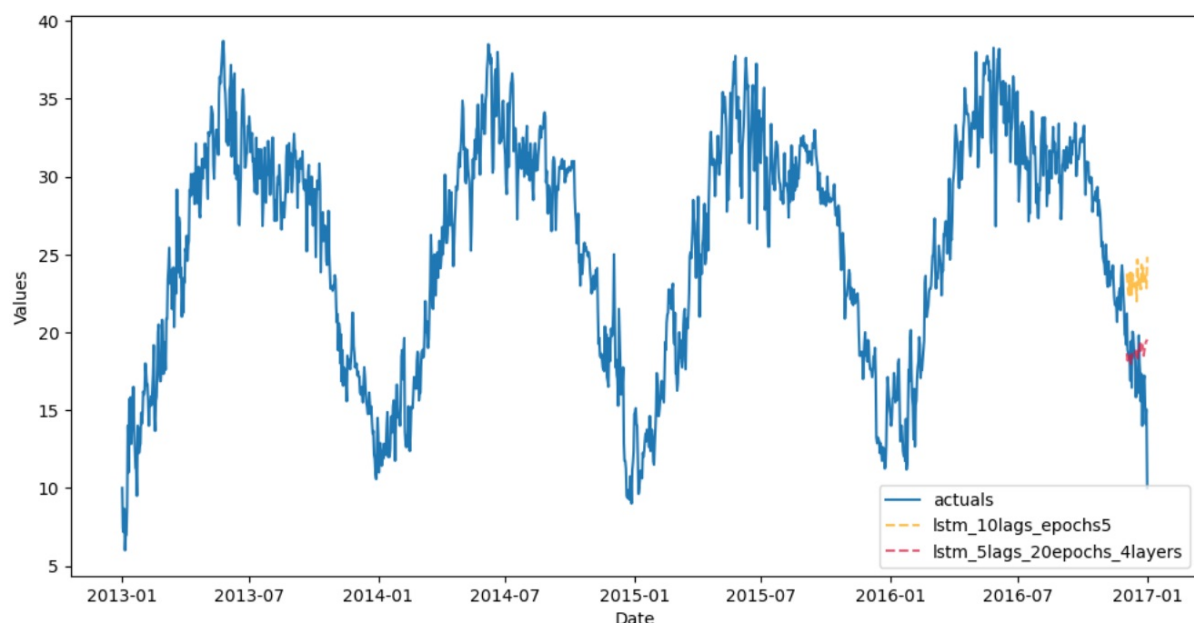
Die vorherige Darstellung zeigt Prognosen, die relativ weit von den tatsächlichen Daten entfernt liegen. Daher ist es notwendig, unsere Parameter abzustimmen und dann das Modell zu testen.

Jetzt ändern wir die Anzahl der Verzögerungen auf 5 und die Anzahl der Epochen auf 20. Außerdem werden wir die Anzahl der Ebenen erhöhen und weitere Parameter optimieren. Lassen Sie uns die Leistung und Plots überprüfen:

```
In [*]:
# All data is scaled going into the model with a min-max scaler and un-scaled coming out.
#Anything you can pass to the fit() method in TensorFlow,
# you can also pass to the scalecast manual_forecast() method.
#Plots all test-set predictions with the actuals.
#ci (bool) - Default False. Whether to display the confidence intervals.
# 5 lags, since we noticed 5 days autocorrelation

# Let's try increasing the number of layers in the network to 4,
#increasing epochs to 10, but monitoring the validation loss value and telling the model to quit after more
#than 5 iterations in which that doesn't improve. This is known as early stopping.
from tensorflow.keras.callbacks import EarlyStopping
f.manual_forecast(
    call_me='lstm_5lags_20epochs_4layers',
    lags=5,
    epochs=20,
    batch_size=16,
    activation='tanh',
    optimizer='Adam',
    shuffle=True,
    learning_rate=0.01,

    lstm_layer_sizes=(72,)*4, # 4 layers, each 72 units (size)
    dropout=(0,)*4, # dropout rate for each layer
    plot_loss=True
)
f.plot_test_set(ci=True)
```



Das zweite Modell zeigt offensichtlich eine bessere Leistung. Die MAPE-Ergebnisse für den Testsatz zeigen, dass das letztere Modell das erste übertrifft:

```
In [83]: # Lets have a look on the statistics of our models
res = f.export(dfs=['model_summaries'])
models = res['ModelNickname']
for m in models:
    print(m, res.loc[res['ModelNickname'] == m, 'LevelTestSetMAPE'])

lstm_10lags_epochs5 0    0.37005
Name: LevelTestSetMAPE, dtype: float64
lstm_5lags_earlystop_4layers 1    0.10882
Name: LevelTestSetMAPE, dtype: float64
```

```
1 {
2 "cells": [
3 {
4 "cell_type": "code",
5 "execution_count": null,
6 "id": "1a69b9d5",
7 "metadata": {},
8 "outputs": [],
9 "source": [
10 "# we will use climate time series data (daily):https://www.kaggle.com/datasets/sumanthvrao/daily-climate-time-series-da
11 "# install scalecast library\n",
12 "!pip install scalecast --upgrade"
13 ]
14 },
15 {
16 "cell_type": "code",
17 "execution_count": 1,
18 "id": "5fb69a3b",
19 "metadata": {},
20 "outputs": [
21 {
22 "data": {
23 "text/html": [
24 "<div>\n",
25 "<style scoped>\n",
26 " .dataframe tbody tr th:only-of-type {\n",
27 " vertical-align: middle;\n",
28 " }\n",
29 "\n",
30 " .dataframe tbody tr th {\n",
31 " vertical-align: top;\n",
32 " }\n",
33 "\n",
34 " .dataframe thead th {\n",
35 " text-align: right;\n",
36 " }\n",
37 "</style>\n",
38 "<table border=\"1\" class=\"dataframe\">\n",
39 " <thead>\n",
40 " <tr style=\"text-align: right;\">\n",
41 " <th></th>\n",
42 " <th>date</th>\n",
43 " <th>meantemp</th>\n",
44 " <th>humidity</th>\n",
45 " <th>wind_speed</th>\n",
46 " <th>meanpressure</th>\n",
47 " </tr>\n",
48 " </thead>\n",
49 " <tbody>\n",
50 " <tr>\n",
51 " <th>0</th>\n",
52 " <td>2013-01-01</td>\n",
53 " <td>10.000000</td>\n",
54 " <td>84.500000</td>\n",
55 " <td>0.000000</td>\n",
56 " <td>1015.666667</td>\n",
57 " </tr>\n",
58 " <tr>\n",
59 " <th>1</th>\n",
60 " <td>2013-01-02</td>\n",
61 " <td>7.400000</td>\n",
62 " <td>92.000000</td>\n",
63 " <td>2.980000</td>\n",
64 " <td>1017.800000</td>\n",
65 " </tr>\n",
66 " <tr>\n",
67 " <th>2</th>\n",
68 " <td>2013-01-03</td>\n",
69 " <td>7.166667</td>\n",
70 " <td>87.000000</td>\n",
71 " <td>4.633333</td>\n",
72 " <td>1018.666667</td>\n",
73 " </tr>\n",
74 " <tr>\n",
75 " <th>3</th>\n",
76 " <td>2013-01-04</td>\n",
77 " <td>8.666667</td>\n",
78 " <td>71.333333</td>\n",
79 " <td>1.233333</td>\n",
80 " <td>1017.166667</td>\n",
81 " </tr>\n",
82 " <tr>\n",
83 " <th>4</th>\n",
84 " <td>2013-01-05</td>\n",
85 " <td>6.000000</td>\n",
86 " <td>86.833333</td>\n",
87 " <td>3.700000</td>\n",
88 " <td>1016.500000</td>\n",
89 " </tr>\n",
```

```

90         "</tbody>\n",
91         "</table>\n",
92         "</div>"
93     ],
94     "text/plain": [
95         "      date    meantemp    humidity    wind speed    meanpressure\n",
96         "0 2013-01-01  10.000000   84.500000   0.000000   1015.666667\n",
97         "1 2013-01-02   7.400000   92.000000   2.980000   1017.800000\n",
98         "2 2013-01-03   7.166667   87.000000   4.633333   1018.666667\n",
99         "3 2013-01-04   8.666667   71.333333   1.233333   1017.166667\n",
100        "4 2013-01-05   6.000000   86.833333   3.700000   1016.500000"
101    ]
102 },
103 "execution_count": 1,
104 "metadata": {},
105 "output_type": "execute_result"
106 }
107 ],
108 "source": [
109     "import scalecast\n",
110     "import tensorflow\n",
111     "import pandas as pd\n",
112     "import numpy as np\n",
113     "import pickle\n",
114     "import seaborn as sns\n",
115     "import matplotlib.pyplot as plt\n",
116     "from scalecast.Forecaster import Forecaster\n",
117     "# If a column or index cannot be represented as an array of datetimes,\n",
118     "# say because of an unparsable value or a mixture of timezones,\n",
119     "# the column or index will be returned unaltered as an object data type.\n",
120     "\n",
121     "df = pd.read_csv('DailyDelhiClimateTrain.csv', parse_dates=['date'])\n",
122     "# take a look at the data\n",
123     "df.head()\n",
124     "# len(df)"
125 ]
126 },
127 {
128     "cell_type": "code",
129     "execution_count": 2,
130     "id": "7c9592f0",
131     "metadata": {},
132     "outputs": [
133         {
134             "data": {
135                 "text/plain": [
136                     "Forecaster\n",
137                     "      DateStartActuals=2013-01-01T00:00:00.000000000\n",
138                     "      DateEndActuals=2017-01-01T00:00:00.000000000\n",
139                     "      Freq=D\n",
140                     "      N_actuals=1462\n",
141                     "      ForecastLength=0\n",
142                     "      Xvars=[]\n",
143                     "      TestLength=0\n",
144                     "      ValidationLength=1\n",
145                     "      ValidationMetric=rmse\n",
146                     "      ForecastsEvaluated=[]\n",
147                     "      CILevel=None\n",
148                     "      CurrentEstimator=mlr\n",
149                     "      GridsFile=Grids\n",
150                     "]"
151                 ]
152             },
153             "execution_count": 2,
154             "metadata": {},
155             "output_type": "execute_result"
156         }
157     ],
158     "source": [
159         "# we must first call the Forecaster object with the y and current_dates parameters as 'meantemp' and 'date' variable spe\n",
160         "f = Forecaster(y=df['meantemp'], current_dates=df['date'])\n",
161         "f"
162     ]
163 },
164 {
165     "cell_type": "code",
166     "execution_count": 3,
167     "id": "8f8c814e",
168     "metadata": {},
169     "outputs": [
170         {
171             "name": "stderr",
172             "output_type": "stream",
173             "text": [
174                 "C:\\Users\\la2022\\anaconda3\\lib\\site-packages\\statsmodels\\graphics\\tsaplots.py:348: FutureWarning: The default i\n",
175                 "      warnings.warn(\n",
176             ]
177         },
178         {
179             "data": {
180                 "image/png": "iVBORw0KgAAAAANSUHEUgAAAJgAAAGxCAYAAABvIsx7AAAAOXRFWHRTb2Z0d2FyZQBNYXRwbG90bGliIHZlcnNpb24zLjcuMSwgSwgAHR0",
181                 "text/plain": [
182                     "<Figure size 640x480 with 1 Axes>"
183                 ]
184             },
185             "metadata": {},
186             "output_type": "display_data"
187         }
188     ],
189     "source": [
190         "# Let's decompose this time series by viewing the PACF (Partial Auto Correlation Function) plot,\n",
191         "# which measures how much the y variable (meantemp) is correlated to past values of itself.\n",
192         "# blue area PACF plots depicts the significance threshold.\n",
193         "# That means, lags that located within this area is statistically close to zero and thus insignificant autocorrelation\

```



```

194     "# between data points. \n",
195     "f.plot_pacf(lags=30)# up to 30 lags\n",
196     "plt.show()\n",
197     "# PACF will depicts intuitively correlations of 1 at lag 0,\n",
198     "# since this represents the correlation of the time series with itself.\n",
199     "#this plot indicate significant autocorrelation at lag 1 which means that adjacent points (have lag of 1) are highly co-
200     "# there are non zero autocorrelation at different lags as well"
201 ]
202 },
203 {
204     "cell_type": "code",
205     "execution_count": 5,
206     "id": "5c1c43fc",
207     "metadata": {},
208     "outputs": [
209         {
210             "data": {
211                 "image/png": "iVBORw0KGgoAAAANSUHEUgAAANyAAAHWCAYAAAD6oMSKAAAOXRFWHRTb2Z0d2FyZQBNYXRwbG90bGliIHZlcnNpb24zLjcuMSwgaHR0
212                 "text/plain": [
213                     "<Figure size 640x480 with 4 Axes>"
214                 ]
215             },
216             "metadata": {},
217             "output_type": "display_data"
218         }
219     ],
220     "source": [
221         "# Let's further decompose the series into its trend, seasonal, and residual parts:\n",
222         "from matplotlib import pyplot\n",
223         "from statsmodels.tsa.seasonal import seasonal_decompose\n",
224         "result = seasonal_decompose(df['meantemp'], model='multiplicative', period=365)# requency of the observations is 1\n",
225         "result.plot()\n",
226         "pyplot.show()\n",
227         "#The figure obviously indicates yearly seasonality and increasing trend over time."
228     ]
229 },
230 {
231     "cell_type": "code",
232     "execution_count": 6,
233     "id": "96cbfb55",
234     "metadata": {},
235     "outputs": [
236         {
237             "name": "stdout",
238             "output_type": "stream",
239             "text": [
240                 "(-2.021069055920671, 0.2774121372301611, 10, 1451, {'1%': -3.4348647527922824, '5%': -2.863533960720434, '10%': -2.56
241             ]
242         }
243     ],
244     "source": [
245         "# let's test the series? stationarity.\n",
246         "# If bool (full_res = False), returns whether the test suggests stationarity.\n",
247         "# null Hypothesis of Augmented Dickey-Fuller (ADF) test: time series is non_stationary ( there is a unit root,)\n",
248         "# If the pvalue is above a critical size (Default is 0.05), then we cannot reject that there is a unit root.\n",
249         "stat = f.adf_test(full_res=True)\n",
250         "print(stat)\n",
251         "\n",
252         "# p_value is 0.28 > 0.05, then we cannot reject null hypothesis and thus time series is non stationary"
253     ]
254 },
255 {
256     "cell_type": "code",
257     "execution_count": 7,
258     "id": "479c8afe",
259     "metadata": {},
260     "outputs": [],
261     "source": [
262         "# Now, to call an LSTM forecast. By default, \n",
263         "# this model will be run with a single input layer of 8 units, Adam optimizer, tanh activation,\n",
264         "# a learning rate of 0.001, and no dropout.\n",
265         "\n",
266         "# generate future dates: The number of dates you generate in this step will determine how long all models will be forec
267         "f.set_validation_metric('mape')\n",
268         "f.set_test_length(30) # 30 observations to test the results\n",
269         "f.generate_future_dates(30) # 30 future points to forecast\n",
270         "\n",
271         "# LSTM neural network\n",
272         "f.set_estimator('lstm')
273     ]
274 },
275 {
276     "cell_type": "code",
277     "execution_count": null,
278     "id": "c430afc2",
279     "metadata": {},
280     "outputs": [],
281     "source": []
282 },
283 {
284     "cell_type": "code",
285     "execution_count": 8,
286     "id": "ed39a4f1",
287     "metadata": {},
288     "outputs": [
289         {
290             "name": "stdout",
291             "output_type": "stream",
292             "text": [
293                 "Epoch 1/5\n",
294                 "44/44 [=====] - 2s 5ms/step - loss: 0.5603\n",
295                 "Epoch 2/5\n",
296                 "44/44 [=====] - 0s 5ms/step - loss: 0.3823\n",
297                 "Epoch 3/5\n",

```



```

298     "44/44 [=====] - 0s 5ms/step - loss: 0.2042\n",
299     "Epoch 4/5\n",
300     "44/44 [=====] - 0s 5ms/step - loss: 0.1186\n",
301     "Epoch 5/5\n",
302     "44/44 [=====] - 0s 5ms/step - loss: 0.0934\n",
303     "1/1 [=====] - 0s 444ms/step\n",
304     "Epoch 1/5\n",
305     "45/45 [=====] - 2s 5ms/step - loss: 0.5340\n",
306     "Epoch 2/5\n",
307     "45/45 [=====] - 0s 5ms/step - loss: 0.3052\n",
308     "Epoch 3/5\n",
309     "45/45 [=====] - 0s 4ms/step - loss: 0.1584\n",
310     "Epoch 4/5\n",
311     "45/45 [=====] - 0s 5ms/step - loss: 0.1089\n",
312     "Epoch 5/5\n",
313     "45/45 [=====] - 0s 5ms/step - loss: 0.0955\n",
314     "1/1 [=====] - 0s 426ms/step\n",
315     "45/45 [=====] - 0s 2ms/step\n"
316 ]
317 },
318 {
319     "name": "stderr",
320     "output_type": "stream",
321     "text": [
322         "C:\\Users\\la2022\\anaconda3\\lib\\site-packages\\scalecast\\_utils.py:55: Warning: Confidence intervals not found fo
323         " warnings.warn(\n"
324     ]
325 },
326 {
327     "data": {
328         "text/plain": [
329             "<Axes: xlabel='Date', ylabel='Values'>"
330         ]
331     },
332     "execution_count": 8,
333     "metadata": {},
334     "output_type": "execute_result"
335 },
336 {
337     "data": {
338         "image/png": "iVBORw0KGgoAAAANSUhEUgAAA+UAAAIOCAYAAADeN6x/AAAAOXRFWHRTb2Z0d2FyZQBNYXRwbG90bGliIHZlcnNpb24zLjcuMSwgaHR0
339         "text/plain": [
340             "<Figure size 1200x600 with 1 Axes>"
341         ]
342     },
343     "metadata": {},
344     "output_type": "display_data"
345 },
346 ],
347 "source": [
348     "f.manual_forecast(call_me='lstm_10lags_epochs5',lags=10, epochs=5)\n",
349     "##f.tf_model.save('lstm_10lags_epochs5_model.h5')\n",
350     "f.plot_test_set(ci=True)"
351 ]
352 },
353 {
354     "cell_type": "code",
355     "execution_count": 9,
356     "id": "e4fca805",
357     "metadata": {},
358     "outputs": [
359         {
360             "data": {
361                 "text/plain": [
362                     "<bound method Model.summary of <keras.engine.sequential.Sequential object at 0x000002205F08F040>>"
363                 ]
364             },
365             "execution_count": 9,
366             "metadata": {},
367             "output_type": "execute_result"
368         }
369     ],
370     "source": [
371         "#del f.tf_model\n",
372         "from keras.models import load_model\n",
373         "\n",
374         "model=load_model('lstm_10lags_epochs5_model.h5')\n",
375         "model.summary\n"
376     ]
377 },
378 {
379     "cell_type": "code",
380     "execution_count": 10,
381     "id": "ffa3a48a",
382     "metadata": {},
383     "outputs": [
384         {
385             "name": "stderr",
386             "output_type": "stream",
387             "text": [
388                 "C:\\Users\\la2022\\anaconda3\\lib\\site-packages\\scalecast\\_utils.py:55: Warning: Confidence intervals not found fo
389                 " warnings.warn(\n"
390             ]
391         },
392         {
393             "data": {
394                 "text/plain": [
395                     "<Axes: xlabel='Date', ylabel='Values'>"
396                 ]
397             },
398             "execution_count": 10,
399             "metadata": {},
400             "output_type": "execute_result"
401         }

```

```
402 {
403   "data": {
404     "image/png": "iVBORw0KGgoAAAANSUhEUgAAA+UAAAIOCAYAAADeN6x/AAAAOXRFWHRTb2Z0d2FyZQBNYXRwbG90bGliIHZlcnNpb24zLjcuMSwgHR0",
405     "text/plain": [
406       "<Figure size 1200x600 with 1 Axes>"
407     ]
408   },
409   "metadata": {},
410   "output_type": "display_data"
411 }
412 ],
413 "source": [
414   "f.plot_test_set(ci=True)"
415 ]
416 },
417 {
418   "cell_type": "code",
419   "execution_count": 18,
420   "id": "045bc465",
421   "metadata": {},
422   "outputs": [
423     {
424       "name": "stdout",
425       "output_type": "stream",
426       "text": [
427         "Epoch 1/20\n",
428         "88/88 [=====] - 8s 11ms/step - loss: 0.1365\n",
429         "Epoch 2/20\n",
430         "88/88 [=====] - 1s 11ms/step - loss: 0.0938\n",
431         "Epoch 3/20\n",
432         "88/88 [=====] - 1s 11ms/step - loss: 0.0934\n",
433         "Epoch 4/20\n",
434         "88/88 [=====] - 1s 11ms/step - loss: 0.0937\n",
435         "Epoch 5/20\n",
436         "88/88 [=====] - 1s 11ms/step - loss: 0.0886\n",
437         "Epoch 6/20\n",
438         "88/88 [=====] - 1s 11ms/step - loss: 0.0886\n",
439         "Epoch 7/20\n",
440         "88/88 [=====] - 1s 11ms/step - loss: 0.0884\n",
441         "Epoch 8/20\n",
442         "88/88 [=====] - 1s 11ms/step - loss: 0.0881\n",
443         "Epoch 9/20\n",
444         "88/88 [=====] - 1s 11ms/step - loss: 0.0865\n",
445         "Epoch 10/20\n",
446         "88/88 [=====] - 1s 11ms/step - loss: 0.0857\n",
447         "Epoch 11/20\n",
448         "88/88 [=====] - 1s 11ms/step - loss: 0.0892\n",
449         "Epoch 12/20\n",
450         "88/88 [=====] - 1s 11ms/step - loss: 0.0854\n",
451         "Epoch 13/20\n",
452         "88/88 [=====] - 1s 12ms/step - loss: 0.0846\n",
453         "Epoch 14/20\n",
454         "88/88 [=====] - 1s 11ms/step - loss: 0.0841\n",
455         "Epoch 15/20\n",
456         "88/88 [=====] - 1s 11ms/step - loss: 0.0866\n",
457         "Epoch 16/20\n",
458         "88/88 [=====] - 1s 11ms/step - loss: 0.0867\n",
459         "Epoch 17/20\n",
460         "88/88 [=====] - 1s 11ms/step - loss: 0.0843\n",
461         "Epoch 18/20\n",
462         "88/88 [=====] - 1s 11ms/step - loss: 0.0831\n",
463         "Epoch 19/20\n",
464         "88/88 [=====] - 1s 12ms/step - loss: 0.0833\n",
465         "Epoch 20/20\n",
466         "88/88 [=====] - 1s 11ms/step - loss: 0.0836\n",
467         "1/1 [=====] - 1s 1s/step\n",
468         "Epoch 1/20\n",
469         "90/90 [=====] - 8s 11ms/step - loss: 0.1291\n",
470         "Epoch 2/20\n",
471         "90/90 [=====] - 1s 12ms/step - loss: 0.0950\n",
472         "Epoch 3/20\n",
473         "90/90 [=====] - 1s 11ms/step - loss: 0.0947\n",
474         "Epoch 4/20\n",
475         "90/90 [=====] - 1s 10ms/step - loss: 0.0958\n",
476         "Epoch 5/20\n",
477         "90/90 [=====] - 1s 10ms/step - loss: 0.0913\n",
478         "Epoch 6/20\n",
479         "90/90 [=====] - 1s 10ms/step - loss: 0.0903\n",
480         "Epoch 7/20\n",
481         "90/90 [=====] - 1s 10ms/step - loss: 0.0892\n",
482         "Epoch 8/20\n",
483         "90/90 [=====] - 1s 10ms/step - loss: 0.0881\n",
484         "Epoch 9/20\n",
485         "90/90 [=====] - 1s 10ms/step - loss: 0.0855\n",
486         "Epoch 10/20\n",
487         "90/90 [=====] - 1s 10ms/step - loss: 0.0850\n",
488         "Epoch 11/20\n",
489         "90/90 [=====] - 1s 10ms/step - loss: 0.0849\n",
490         "Epoch 12/20\n",
491         "90/90 [=====] - 1s 10ms/step - loss: 0.0866\n",
492         "Epoch 13/20\n",
493         "90/90 [=====] - 1s 10ms/step - loss: 0.0854\n",
494         "Epoch 14/20\n",
495         "90/90 [=====] - 1s 11ms/step - loss: 0.0866\n",
496         "Epoch 15/20\n",
497         "90/90 [=====] - 1s 10ms/step - loss: 0.0870\n",
498         "Epoch 16/20\n",
499         "90/90 [=====] - 1s 10ms/step - loss: 0.0867\n",
500       ]
501     },
502   ],
503   "data": {
504     "image/png": "iVBORw0KGgoAAAANSUhEUgAAAwAAAHFCAYAAADR1KI/AAAAOXRFWHRTb2Z0d2FyZQBNYXRwbG90bGliIHZlcnNpb24zLjcuMSwgHR0",
505     "text/plain": [
```

```

506         "<Figure size 640x480 with 1 Axes>"
507     ]
508 },
509     "metadata": {},
510     "output_type": "display_data"
511 },
512 {
513     "name": "stdout",
514     "output_type": "stream",
515     "text": [
516         "1/1 [=====] - 1s 1s/step\n",
517         "45/45 [=====] - 0s 5ms/step\n"
518     ]
519 },
520 {
521     "name": "stderr",
522     "output_type": "stream",
523     "text": [
524         "C:\\Users\\la2022\\anaconda3\\lib\\site-packages\\scalecast\\_utils.py:55: Warning: Confidence intervals not found fo:
525         " warnings.warn(\n",
526         "C:\\Users\\la2022\\anaconda3\\lib\\site-packages\\scalecast\\_utils.py:55: Warning: Confidence intervals not found fo:
527         " warnings.warn(\n"
528     ]
529 },
530 {
531     "data": {
532         "text/plain": [
533             "<Axes: xlabel='Date', ylabel='Values'>"
534         ]
535     },
536     "execution_count": 18,
537     "metadata": {},
538     "output_type": "execute_result"
539 },
540 {
541     "data": {
542         "image/png": "iVBORw0KGgoAAAANSUHEUgAAA+UAAAIOCAYAAADeN6x/AAAAOXRFWHRTb2Z0d2FyZQBNYXRwbG90bGliIHZlcnNpb24zLjJlcuMSwgaHR0
543         "text/plain": [
544             "<Figure size 1200x600 with 1 Axes>"
545         ]
546     },
547     "metadata": {},
548     "output_type": "display_data"
549 },
550 ],
551 "source": [
552     "# All data is scaled going into the model with a min-max scaler and un-scaled coming out.\n",
553     "#Anything you can pass to the fit() method in TensorFlow,\n",
554     "# you can also pass to the scalecast manual_forecast() method.\n",
555     "#Plots all test-set predictions with the actuals.\n",
556     "#ci (bool) ? Default False. Whether to display the confidence intervals.\n",
557     "# 5 lags, since we noticed 5 days autocorrelation\n",
558     "\n",
559     "# let's try increasing the number of layers in the network to 4,\n",
560     "#increasing epochs to 10, but monitoring the validation loss value and telling the model to quit after more\n",
561     "#than 5 iterations in which that doesn't improve. This is known as early stopping.\n",
562     "from tensorflow.keras.callbacks import EarlyStopping\n",
563     "callback = EarlyStopping(monitor='loss', patience=5)\n",
564     "f.manual_forecast(\n",
565         "    call_me='lstm_5lags_20epochs_4layers',\n",
566         "    lags=5,\n",
567         "    epochs=20,\n",
568         "    batch_size=16,\n",
569         "    activation='tanh',\n",
570         "    optimizer='Adam',\n",
571         "    shuffle=True,\n",
572         "    learning_rate=0.01,\n",
573         "    callbacks=[callback],\n",
574         "    lstm_layer_sizes=(72,)*4, # 4 layers, each 72 units (size)\n",
575         "    dropout=(0,)*4, # dropout rate for each layer\n",
576         "    plot_loss=True\n",
577     )\n",
578     "\n",
579     "f.plot_test_set(ci=True)
580 ]
581 },
582 {
583     "cell_type": "code",
584     "execution_count": 12,
585     "id": "594568f8",
586     "metadata": {},
587     "outputs": [
588     {
589         "name": "stderr",
590         "output_type": "stream",
591         "text": [
592             "C:\\Users\\la2022\\anaconda3\\lib\\site-packages\\scalecast\\_utils.py:55: Warning: Confidence intervals not found fo:
593             " warnings.warn(\n",
594             "C:\\Users\\la2022\\anaconda3\\lib\\site-packages\\scalecast\\_utils.py:55: Warning: Confidence intervals not found fo:
595             " warnings.warn(\n"
596         ]
597     },
598     {
599         "data": {
600             "text/plain": [
601                 "<Axes: xlabel='Date', ylabel='Values'>"
602             ]
603         },
604         "execution_count": 12,
605         "metadata": {},
606         "output_type": "execute_result"
607     },
608     {
609         "data": {

```

```

610         "image/png": "iVBORw0KGgoAAAANSUhEUgAAA+UAAAIOCAYAAADeN6x/AAAAOXRFWHRTb2Z0d2FyZQBNYXRwbG90bGliIHZlcnNpb24zLjcuMSwgaHR0",
611         "text/plain": [
612             "<Figure size 1200x600 with 1 Axes>"
613         ],
614     },
615     "metadata": {},
616     "output_type": "display_data"
617 },
618 ],
619 "source": [
620     "f.plot_test_set(ci=True)"
621 ],
622 },
623 {
624     "cell_type": "code",
625     "execution_count": null,
626     "id": "060e5e65",
627     "metadata": {},
628     "outputs": [],
629     "source": [
630         "# plot the best 2 models based on MAPE metric\n",
631         "#f.plot_test_set(order_by='LevelTestSetMAPE',models='top_2',ci=True)## MAPE metric is used "
632     ]
633 },
634 {
635     "cell_type": "code",
636     "execution_count": 13,
637     "id": "27ec740f",
638     "metadata": {},
639     "outputs": [
640         {
641             "name": "stdout",
642             "output_type": "stream",
643             "text": [
644                 "lstm_10lags_epochs5 0      0.353709\n",
645                 "Name: LevelTestSetMAPE, dtype: float64\n",
646                 "lstm_5lags_20epochs_4layers 1      0.126901\n",
647                 "Name: LevelTestSetMAPE, dtype: float64\n"
648             ]
649         }
650     ],
651     "source": [
652         "# lets have a look on the statistics of our models\n",
653         "res = f.export(dfs=['model_summaries'])\n",
654         "models =res['ModelNickname']\n",
655         "for m in models:\n",
656             "    print(m, res.loc[res['ModelNickname'] == m, 'LevelTestSetMAPE'])\n",
657         " "
658     ]
659 },
660 {
661     "cell_type": "code",
662     "execution_count": 2,
663     "id": "balef624",
664     "metadata": {},
665     "outputs": [],
666     "source": []
667 },
668 {
669     "cell_type": "code",
670     "execution_count": null,
671     "id": "bb0a4723",
672     "metadata": {},
673     "outputs": [],
674     "source": []
675 },
676 ],
677 "metadata": {
678     "kernelspec": {
679         "display_name": "Python 3 (ipykernel)",
680         "language": "python",
681         "name": "python3"
682     },
683     "language_info": {
684         "codemirror_mode": {
685             "name": "ipython",
686             "version": 3
687         },
688         "file_extension": ".py",
689         "mimetype": "text/x-python",
690         "name": "python",
691         "nbconvert_exporter": "python",
692         "pygments_lexer": "ipython3",
693         "version": "3.9.13"
694     }
695 },
696 "nbformat": 4,
697 "nbformat_minor": 5
698 }

```

Exponentielle Glättung

Exponentielle Glättung ist eine Prognosemethode für univariate Zeitreihendaten. Im Gegensatz zur Methode des gleitenden Durchschnitts, bei der alle vergangenen Beobachtungen gleich gewichtet werden, wenn sie in das Fenster des gleitenden Durchschnitts fallen, verwendet die Methode der exponentiellen Glättung Gewichtungen, die mit zunehmendem Alter der Beobachtungen exponentiell abnehmen, wodurch Ausreißer oder Rauschen aus den Daten entfernt werden. Dies hilft dem Modell, deutliche und sich wiederholende Muster besser zu erkennen, die sonst im Rauschen verborgen wären, und kann erheblich zur Prognoseleistung beitragen, indem die Genauigkeit der Vorhersagen verbessert wird.

Zu den verschiedenen Arten der exponentiellen Glättung gehören die einfache exponentielle Glättung und die dreifache exponentielle Glättung (auch als Holt-Winters-Methode bekannt). Die einfache exponentielle Glättung befasst sich mit Datenreihen, die weder einen Trend noch eine saisonale Komponente aufweisen. Diese Methode verwendet nur einen Parameter namens Alpha (α) und aktualisiert die Pegelkomponente für jede Beobachtung. Da eine Komponente modelliert wird, verwendet es nur einen Gewichtsparameter, Alpha (α), der den Grad der Glättung bestimmt und Werte zwischen 0 und 1 annimmt. Niedrigere Werte führen zu glatteren angepassten Modellen, da den vergangenen Daten mehr Gewichtung gegeben wird und so die aktuellen Daten effizient mit den alten gemittelt werden. Umgekehrt können hohe Werte ein angepasstes Modell erzeugen, das schnell auf das zufällige Rauschen reagiert. Die einfache exponentielle Glättung wird in der folgenden mathematischen Formel dargestellt:

$$f_t = \alpha x_{t-1} + (1 - \alpha) f_{t-1}$$

Dabei ist x_{t-1} die Beobachtung im vorherigen Zeitschritt. Alpha ist der Prozentsatz, wie wichtig die jüngste Beobachtung im Vergleich zu historischen Daten im Modell ist. f_{t-1} die Vorhersage des Modells im vorherigen Zeitschritt. Die obige Gleichung kann in die folgende Formel übersetzt werden:

$$\begin{aligned} f_t &= \alpha x_{t-1} + (1 - \alpha) (\alpha x_{t-2} + (1 - \alpha) f_{t-2}) \\ &= \alpha x_{t-1} + \alpha (1 - \alpha) x_{t-2} + (1 - \alpha)^2 f_{t-2} \\ &= \alpha x_{t-1} + \alpha (1 - \alpha) x_{t-2} + \alpha (1 - \alpha)^2 x_{t-3} + (1 - \alpha)^3 f_{t-3} \end{aligned}$$

us dieser Gleichung können wir das Hauptkonzept der exponentiellen Glättung verstehen, bei der das Modell weiteren Beobachtungen exponentiell niedrigere Gewichte zuweist als den jüngsten.

Daher sollte der Alpha-Wert feinabgestimmt werden, um ein gut passendes Prognosemodell zu erhalten. Einfache exponentielle Glättung wird beispielsweise verwendet, um die Nachfrage nach einem Produkt zu modellieren.

Die dreifache exponentielle Glättung (Holt-Winters) wird für Zeitreihen verwendet, die Trend- und Saisonschwankungen aufweisen. Dies ist eine Erweiterung der einfachen exponentiellen Glättung, die die Einbeziehung von Saisonalitäts- und Trendkomponenten ermöglicht, um viele Werte über historische Daten zu berechnen und genauere aktuelle und zukünftige Prognosen zu erstellen.

Exponentielle Glättung Übung

Wir werden Covid-19 tägliche Zeitreihe verwenden. Diese Zeitreihe besteht aus den beiden Spalten „Date“ und „Confirmed“. Importieren Sie die Daten und werfen Sie einen Blick darauf:

```
In [6]: # downloading dataset
import pandas as pd
df = pd.read_csv('Covid19_ts.csv', parse_dates=['Date'])
df['Date'] = pd.to_datetime(df['Date'])
df = pd.DataFrame(df, columns=['Date', 'Confirmed']).set_index('Date')
train = df.iloc[:200, :]
train.index = pd.to_datetime(train.index) # convert index to datetime format
df
```

Out[6]:

	Date	Confirmed
0	2020-01-31	0.0
1	2020-02-01	0.0
2	2020-02-02	1.0
3	2020-02-03	1.0
4	2020-02-04	0.0
...
836	2022-05-16	324.0
837	2022-05-17	596.0
838	2022-05-18	501.0
839	2022-05-19	556.0
840	2022-05-20	558.0

841 rows × 2 columns

Wir werden ein einfaches exponentielles Glättungsmodell im Trainingsatz trainieren, indem wir die Python-Bibliothek statsmodels verwenden. Wir werden das Modell mit verschiedenen Alpha-Werten anpassen und die aus den Stichproben resultierenden Prognosen überprüfen:

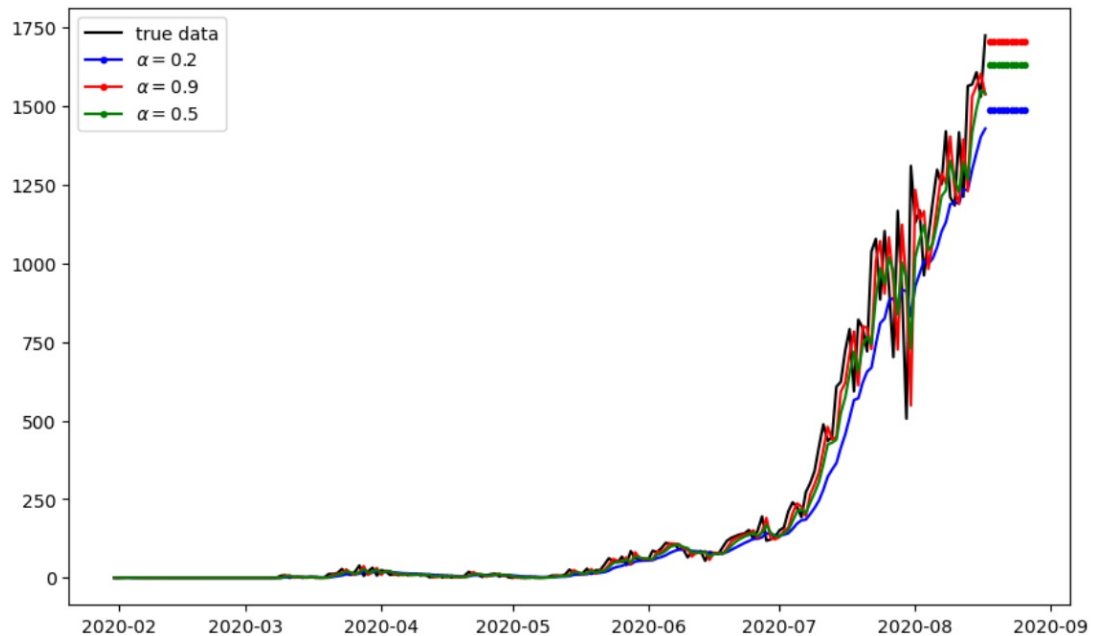
```
In [90]: #simple exponential smoothing
import matplotlib.pyplot as plt
import matplotlib.dates as mdates
from statsmodels.tsa.holtwinters import SimpleExpSmoothing, Holt
import numpy as np
train_data = pd.Series(train.Confirmed, train.index)

# building ES model
model = SimpleExpSmoothing(train_data, initialization_method="heuristic")

# The smoothing_level value of the simple exponential smoothing, if the value is set then this value will be used as the value.
fm1 = model.fit(smoothing_level=.2) # fitting the model
fcast_1 = fm1.forecast(9).rename(r"$\alpha=0.2$") # The number of out of sample forecasts from the end of the sample.
fm2 = model.fit(smoothing_level=.9)
fcast_2 = fm2.forecast(9).rename(r"$\alpha=0.9$")
fm3 = model.fit(smoothing_level=.5)
fcast_3 = fm3.forecast(9).rename(r"$\alpha=0.5$")
```

```
In [91]: plt.figure(figsize=(12, 8))
(line0,) = plt.plot(train_data, color="black")# plotting real temporal data
plt.plot(fm1.fittedvalues, color="blue")
(line1,) = plt.plot(fcast_1, marker=".", color="blue")
plt.plot(fm2.fittedvalues, color="red")
(line2,) = plt.plot(fcast_2, marker=".", color="red")
plt.plot(fm3.fittedvalues, color="green")
(line3,) = plt.plot(fcast_3, marker=".", color="green")
plt.legend([line1, line2, line3], [fcast_1.name, fcast_2.name, fcast_3.name])
plt.legend([line0, line1, line2, line3], ['true data', fcast_1.name, fcast_2.name, fcast_3.name])
# this will results in the forecasts with same value for all dates (9 out of the sample forecasts)
# because no trends are incorporated
```

Out[13]: <matplotlib.legend.Legend at 0x1c663aff190>



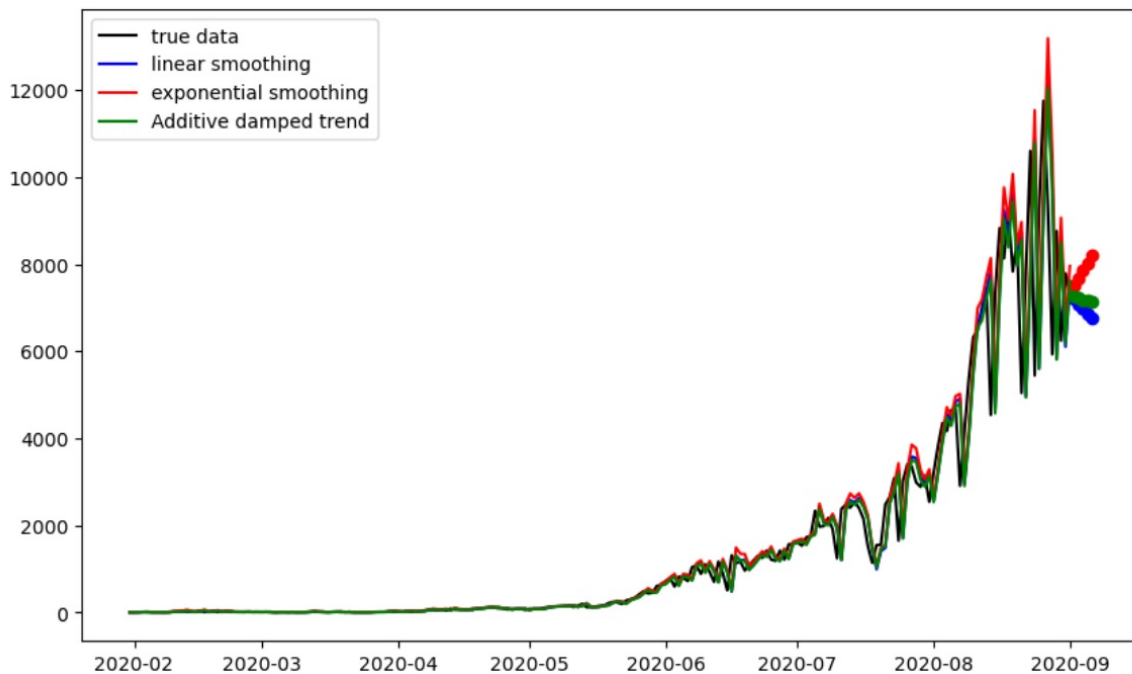
Wie in der Abbildung dargestellt, führt dieses Modell zu Prognosen mit demselben Wert für alle Daten (9 der Beispielpagnosen). Dies liegt daran, dass keine Trends in die Berechnung der Prognosen einbezogen werden. Jetzt werden wir die exponentielle Holt- α -Glättung so aufbauen, dass sie den Trend berücksichtigt. Anschließend werden wir das Modell mit verschiedenen Arten der Glättung anpassen: lineare Glättung, exponentielle Glättung und additiv gedämpfter Trend. Schauen wir uns die resultierenden Prognosen an:

```
In [94]: # Holt exponential smoothing
# preparing the data; removing rows with zero values
df = pd.read_csv('Covid19_ts.csv', parse_dates=['Date'])
# remove zeros from the dataset, this will result in disconnection in the time series,
# thats why, we generate the time series manually..
df1 = df[df['Confirmed'] > 0]
# get values without index
df1 = df1['Confirmed'].values
# convert from np array into list
df1 = df1.tolist()
# select the first 215 values
df1 = df1[:215]
# generate time sequence
index = pd.date_range(start='01/31/2020', end='9/01/2020')
len(index)
# join datetime into the data
train_data = pd.Series(df1, index)

# linear smoothing
Holt_fit1 = Holt(train_data, initialization_method="estimated").fit(
    smoothing_level=0.9, smoothing_trend=0.2, optimized=False)
fcast_holt_1 = Holt_fit1.forecast(5).rename("Holt's linear trend")
# exponential smoothing
Holt_fit2 = Holt(train_data, exponential=True, initialization_method="estimated").fit(
    smoothing_level=0.9, smoothing_trend=0.2, optimized=False)
fcast_holt_2 = Holt_fit2.forecast(5).rename("Exponential trend")

# Additive damped trend
Holt_fit3 = Holt(train_data, damped_trend=True, initialization_method="estimated").fit(
    smoothing_level=0.9, smoothing_trend=0.2)
fcast_holt_3 = Holt_fit3.forecast(5).rename("Additive damped trend")

plt.figure(figsize=(12, 8))
(line0,) = plt.plot(train_data, color="black", label='true data')
plt.plot(Holt_fit1.fittedvalues, color="blue", label='linear smoothing')
(line1,) = plt.plot(fcast_holt_1, marker="o", color="blue")
# plt.plot(fcast_holt_2, color="blue")
plt.plot(Holt_fit2.fittedvalues, color="red", label='exponential smoothing')
(line2,) = plt.plot(fcast_holt_2, marker="o", color="red")
plt.plot(Holt_fit3.fittedvalues, color="green", label='Additive damped trend')
(line3,) = plt.plot(fcast_holt_3, marker="o", color="green")
plt.legend(loc="upper left")
plt.show()
# forecasts show clearly the trend now..|
```



```

1 {
2   "cells": [
3     {
4       "cell_type": "code",
5       "execution_count": null,
6       "id": "d835e8e8",
7       "metadata": {},
8       "outputs": [],
9       "source": [
10        "# downloading dataset\n",
11        "import pandas as pd\n",
12        "df = pd.read_csv('Covid19_ts.csv', parse_dates=['Date'])\n",
13        "df['Date'] = pd.to_datetime(df['Date'])\n",
14        "df = pd.DataFrame(df, columns=['Date', 'Confirmed']).set_index('Date')\n",
15        "train = df.iloc[:200, :]\n",
16        "train.index = pd.to_datetime(train.index) # convert index to datetime format\n",
17        "df"
18      ]
19    },
20    {
21      "cell_type": "code",
22      "execution_count": null,
23      "id": "8891bebb",
24      "metadata": {},
25      "outputs": [],
26      "source": [
27        "# simple exponential smoothing\n",
28        "import matplotlib.pyplot as plt\n",
29        "import matplotlib.dates as mdates\n",
30        "from statsmodels.tsa.holtwinters import SimpleExpSmoothing, Holt\n",
31        "import numpy as np\n",
32        "train_data = pd.Series(train.Confirmed, train.index)\n",
33        "\n",
34        "# building ES model\n",
35        "model = SimpleExpSmoothing(train_data, initialization_method='heuristic')\n",
36        "\n",
37        "\n",
38        "# The smoothing_level value of the simple exponential smoothing, if the value is set then this value will be used as the\n",
39        "fm1 = model.fit(smoothing_level=.2) # fitting the model\n",
40        "fcst_1 = fm1.forecast(9).rename(r'${alpha=0.2$}') # The number of out of sample forecasts from the end of the sample.\n",
41        "fm2 = model.fit(smoothing_level=.9)\n",
42        "fcst_2 = fm2.forecast(9).rename(r'${alpha=0.9$}')\n",
43        "fm3 = model.fit(smoothing_level=.5)\n",
44        "fcst_3 = fm3.forecast(9).rename(r'${alpha=0.5$}')
45      ]
46    },
47    {
48      "cell_type": "code",
49      "execution_count": null,
50      "id": "77347cf3",
51      "metadata": {},
52      "outputs": [],
53      "source": [
54        "fcst_2"
55      ]
56    },
57    {
58      "cell_type": "code",
59      "execution_count": null,
60      "id": "8a7f3069",
61      "metadata": {},
62      "outputs": [],
63      "source": [
64        "plt.figure(figsize=(10, 6))\n",
65        "(line0,) = plt.plot(train_data, color='black') # plotting real temporal data\n",
66        "plt.plot(fm1.fittedvalues, color='blue')\n",
67        "(line1,) = plt.plot(fcst_1, marker='.', color='blue')\n",
68        "plt.plot(fm2.fittedvalues, color='red')

```



```

69     "(line2,) = plt.plot(fcast_2, marker=".", color="red")\n",
70     "plt.plot(fm3.fittedvalues, color="green")\n",
71     "(line3,) = plt.plot(fcast_3, marker=".", color="green")\n",
72     "plt.legend([line1, line2, line3], [fcast_1.name, fcast_2.name, fcast_3.name])\n",
73     "plt.legend([line0, line1, line2, line3], ['true data', fcast_1.name, fcast_2.name, fcast_3.name])\n",
74     "# this will results in the forecasts with same value for all dates (9 out of the sample forecasts)\n",
75     "# because no trends are incorporated "
76 ]
77 },
78 {
79     "cell_type": "code",
80     "execution_count": null,
81     "id": "cddc11a1",
82     "metadata": {},
83     "outputs": [],
84     "source": [
85         "# Holt exponential smoothing\n",
86         "# preparing the data; removing rows with zero values\n",
87         "df = pd.read_csv('Covid19_ts.csv', parse_dates=['Date'])\n",
88         "# remove zeros from the dataset, this will result in disconnection in the time series,\n",
89         "# thats why, we generate the time series manually..\n",
90         "df1 = df[df['Confirmed']> 0]\n",
91         "#get values without index\n",
92         "df1 = df1['Confirmed'].values\n",
93         "# convert from np array into list\n",
94         "df1 = df1.tolist()\n",
95         "# select the first 215 values\n",
96         "df1 = df1[:215]\n",
97         "# generate time sequence\n",
98         "index = pd.date_range(start='01/31/2020', end='9/01/2020')\n",
99         "len(index)\n",
100        "# join datetime into the data\n",
101        "train_data = pd.Series(df1, index)\n",
102        "\n",
103        "# linear smoothing\n",
104        "Holt_fit1 = Holt(train_data, initialization_method="estimated").fit(\n",
105            smoothing_level=0.9, smoothing_trend=0.2, optimized=False)\n",
106        "fcast_holt_1 = Holt_fit1.forecast(5).rename("Holt's linear trend")\n",
107        "#exponential smoothing\n",
108        "Holt_fit2 = Holt(train_data, exponential=True, initialization_method="estimated").fit(\n",
109            smoothing_level=0.9, smoothing_trend=0.2, optimized=False)\n",
110        "fcast_holt_2 = Holt_fit2.forecast(5).rename("Exponential trend")\n",
111        "\n",
112        "#Additive damped trend\n",
113        "Holt_fit3 = Holt(train_data, damped_trend=True, initialization_method="estimated").fit(\n",
114            smoothing_level=0.9, smoothing_trend=0.2)\n",
115        "fcast_holt_3 = Holt_fit3.forecast(5).rename("Additive damped trend")\n",
116        "\n",
117        "plt.figure(figsize=(10, 6))\n",
118        "(line0,) = plt.plot(train_data, color="black", label='true data')\n",
119        "plt.plot(Holt_fit1.fittedvalues, color="blue", label='linear smoothing')\n",
120        "(line1,) = plt.plot(fcast_holt_1, marker="o", color="blue")\n",
121        "#plt.plot(fcast_holt_2, color="blue")\n",
122        "plt.plot(Holt_fit2.fittedvalues, color="red", label='exponential smoothing')\n",
123        "(line2,) = plt.plot(fcast_holt_2, marker="o", color="red")\n",
124        "plt.plot(Holt_fit3.fittedvalues, color="green", label='Additive damped trend')\n",
125        "(line3,) = plt.plot(fcast_holt_3, marker="o", color="green")\n",
126        "plt.legend(loc="upper left")\n",
127        "plt.show()\n",
128        "# forecasts show clearly the trend now..\n",
129    ]
130 },
131 {
132     "cell_type": "code",
133     "execution_count": null,
134     "id": "1c7d91ee",
135     "metadata": {},
136     "outputs": [],
137     "source": []
138 }
139 ],
140 "metadata": {
141     "kernelspec": {
142         "display_name": "Python 3 (ipykernel)",
143         "language": "python",
144         "name": "python3"
145     },
146     "language_info": {
147         "codemirror_mode": {
148             "name": "ipython",
149             "version": 3
150         },
151         "file_extension": ".py",
152         "mimetype": "text/x-python",
153         "name": "python",
154         "nbconvert_exporter": "python",
155         "pygments_lexer": "ipython3",
156         "version": "3.9.13"
157     }
158 },
159 "nbformat": 4,
160 "nbformat_minor": 5
161 }

```