

```
In [ ]: version = "REPLACE_PACKAGE_VERSION"
```

---

## Assignment 2 Part 1: Time Series Patterns (50 pts)

In this assignment, we're going to practise some techniques that are useful for discerning patterns in a time series.

```
In [ ]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline

from pandas.plotting import register_matplotlib_converters
register_matplotlib_converters()

# Suppress all warnings
import warnings
warnings.filterwarnings("ignore")
```

## Question 1: Load Data (5 pts)

At the time of writing this assignment, August 2020, COVID-19 is still the most topical public-health crisis globally with nearly 300,000 new cases reported worldwide every day. **The number of daily new cases worldwide** is a time series that arises naturally from this topical event, and in this assignment we'll apply some of the techniques we learned in class to this very time series to discern any patterns it may contain.

You are provided with a csv file,

`assets/time_series_covid19_confirmed_global.csv`, which is part of the [Johns Hopkins University CSSE COVID-19 dataset \(https://github.com/CSSEGISandData/COVID-19/tree/master/csse\\_covid\\_19\\_data/csse\\_covid\\_19\\_time\\_series\)](https://github.com/CSSEGISandData/COVID-19/tree/master/csse_covid_19_data/csse_covid_19_time_series). As the name suggests, it contains the number of *cumulative* confirmed cases globally as of certain dates. However, we are interested in the number of *new* cases worldwide every day.

Create a function called `load_data` that reads in the csv file and produces a `pd.Series` that looks like:

```

2020-01-23      99.0
2020-01-24     287.0
2020-01-25     493.0
2020-01-26     684.0
2020-01-27     809.0
...
2020-08-17    209672.0
2020-08-18    255096.0
2020-08-19    274346.0
2020-08-20    267183.0
2020-08-21    270751.0
Length: 212, dtype: float64

```

where

- the index of the series is a `pd.DatetimeIndex`;
- the values of the series are daily *new* cases worldwide; and
- the series doesn't contain any `NaN` values.

**This function should return a `pd.Series` of length 212, whose index is a `pd.DatetimeIndex`.**

```
In [ ]: def load_data():
        daily_new_cases = None

        # YOUR CODE HERE
        raise NotImplementedError()

        return daily_new_cases
```

```
In [ ]: # Autograder tests

stu_ans = load_data()

assert isinstance(stu_ans, pd.Series), "Q1: Your function should return a pd.Series"
assert len(stu_ans) == 212, "Q1: The length of the series returned is incorrect"
assert isinstance(stu_ans.index, pd.DatetimeIndex), "Q1: The index of your series should be a DatetimeIndex"
assert (("2020-01-23" <= stu_ans.index) & (stu_ans.index <= "2020-08-21")), "Q1: Your series should cover the period from 2020-01-23 to 2020-08-21"
assert not stu_ans.isna().any(), "Q1: Your series contains NaN values."
assert np.issubdtype(stu_ans.dtype, np.floating), "Q1: Your series should be a floating point series"

# Some hidden tests

del stu_ans
```

```
In [ ]: # Let's plot and see the time series

fig, ax = plt.subplots(figsize=(10, 6))
ax.plot(load_data())
ax.set_xlabel("Day")
ax.set_ylabel("# Cases")
ax.set_title("Daily New COVID-19 Cases Worldwide")

del fig, ax
```

## Question 2: Perform a Seasonal Decomposition (5 pts)

With the time series ready, let's first perform a seasonal decomposition using tools from the `statsmodels` library to get a sense of what the possible patterns are hidden in the data. Complete the function below that takes a time series and an argument `model`, which indicates whether an additive or multiplicative seasonal decomposition should be performed, and that returns a `DecomposeResult` as produced by the `seasonal_decompose` function from the `statsmodels` library.

**This function should return a `statsmodels.tsa.seasonal.DecomposeResult`.**

```
In [ ]: from statsmodels.tsa.seasonal import seasonal_decompose, DecomposeResult

def sea_decomp(ser, model="additive"):
    """
    Takes in a series and a "model" parameter indicating which seasonal
    """
    result = None

    # YOUR CODE HERE
    raise NotImplementedError()

    return result
```

```
In [ ]: # Autograder tests

stu_ser = load_data()
stu_ans = sea_decomp(stu_ser, model="additive")

assert isinstance(stu_ans, DecomposeResult), "Q2: Your function should re

# Some hidden tests

del stu_ser, stu_ans
```

```
In [ ]: # Let's plot and see the seasonal decomposition

fig, axes = plt.subplots(4, 1, figsize=(10, 6), sharex=True)
res = sea_decomp(load_data(), model="additive")

axes[0].set_title("Additive Seasonal Decomposition")
axes[0].plot(res.observed)
axes[0].set_ylabel("Observed")

axes[1].plot(res.trend)
axes[1].set_ylabel("Trend")

axes[2].plot(res.seasonal)
axes[2].set_ylabel("Seasonal")

axes[3].plot(res.resid)
axes[3].set_ylabel("Residual")

axes[3].set_xlabel("Day")
fig.suptitle("Daily New COVID-19 Cases Worldwide", x=0.513, y=0.95)

del fig, axes, res
```

### Question 3: Fit a Trend Curve (15 pts)

The plot above suggests that there is a non-linear trend hidden in the time series. One approach to discover such a trend is to fit a regression model to the time series and ask the regression model to make predictions at each timestamp. When connected, these chronological predictions form a "trend curve". In the problem, we will explore how to fit a trend curve to our time series.

Complete the function below that fits an  $n$ -th order polynomial to the input time series and that returns the predictions as a `np.ndarray` of the same length. An  $n$ -th order polynomial regression model assumes that each dependent variable  $y_i$  is an  $n$ -th order polynomial function of the corresponding independent variable  $x_i$ :

$$y_i = c_0 + c_1 x_i + c_2 x_i^2 + \dots + c_n x_i^n$$

Now, the most interesting and important question to think about is, "**what are  $x_i$ 's and  $y_i$ 's in the problem?**". The  $y_i$ 's are the daily new cases worldwide at timestamps  $x_i$ 's, but **how should we represent the timestamps  $x_i$ 's in such a regression model?** There are many choices you may explore. In the function below, you are already given the code for training a polynomial regression model, but you have to figure out what `train_X` ( $x_i$ 's) and `train_y` ( $y_i$ 's) are. Since it's possible that everyone has a different design, this question is graded on the  $R^2$  score of your predictions. **For a 10-th order polynomial regression model, at least one choice of  $x_i$ 's leads to an  $R^2$  score  $\geq 0.95$ .**

**This function should return a `np.ndarray` of shape `(len(ser), )`, which represents the predictions of your polynomial regression model on the input time series. The predictions form the "trend curve" we are looking for.**

```
In [ ]: from sklearn.preprocessing import PolynomialFeatures
        from sklearn.linear_model import LinearRegression

        def fit_trend(ser, n):
            """
            Takes a series and fits an n-th order polynomial to the series.
            Returns the predictions.
            """

            trend_curve = None
            train_X, train_y = None, None # xi's and yi's

            # Create train_X and train_y
            # YOUR CODE HERE
            raise NotImplementedError()

            # Fit a polynomial regression model - code given to you
            train_X = PolynomialFeatures(n).fit_transform(train_X.reshape(-1, 1))
            lin_reg = LinearRegression().fit(train_X, train_y.reshape(-1))

            # Make predictions to create the trend curve
            # YOUR CODE HERE
            raise NotImplementedError()

            return trend_curve
```

```
In [ ]: # Autograder tests

        stu_ser = load_data()
        stu_ans = fit_trend(stu_ser, 10)

        assert isinstance(stu_ans, np.ndarray), "Q3: Your function should return
        assert stu_ans.shape == (len(stu_ser), ), "Q3: The shape of your np.ndar

        # Some hidden tests

        del stu_ser, stu_ans
```

```
In [ ]: # Let's plot and see your regression line

fig, ax = plt.subplots(figsize=(10, 6))
ser = load_data()
preds = fit_trend(ser, 10)
ax.plot(ser.index, ser.values, label="Original")
ax.plot(ser.index, preds, label="Fitted trend curve")
ax.set_xlabel("Day")
ax.set_ylabel("# Cases")
ax.set_title("Daily New COVID-19 Cases Worldwide")
ax.legend()

del fig, ax, ser, preds
```

It's worth mentioning that the `seaborn` library provides a function `regplot` (<https://seaborn.pydata.org/generated/seaborn.regplot.html>) that can plot both the data and the regression line in a few lines of code, thus saving you the trouble of fitting a regression model.

## Question 4: Calculate Weighted Moving Average (WMA) (15 pts)

The regression method seems to give a fairly accurate description of the trend hidden in the time series. In this problem and the next, we will explore an alternative method for discovering trends that is based on moving averages.

Recall from the lectures that a Weighted Moving Average (WMA) method applies the following transformation to each data point  $x_j$ :

$$x'_j = \frac{w_k x_j + w_{k-1} x_{j-1} + \dots + w_1 x_{j-k+1}}{w_k + w_{k-1} + \dots + w_1} \quad \text{if } j > k$$

$$x'_j = \frac{w_k x_j + w_{k-1} x_{j-1} + \dots + w_{k-j+1} x_1}{w_k + w_{k-1} + \dots + w_{k-j+1}} \quad \text{if } j \leq k$$

for a window of size  $k$ . Complete the function below that calculates the WMA for an input time series.

**This function should return a `np.ndarray` of shape `(len(ser), )` that represents the WMA values for the input time series.**

```
In [ ]: def calc_wma(ser, wd_size, weights=1):
        """
        Takes in a series and calculates the WMA with a window size of wd_size
        """
        wma = None
        if isinstance(weights, int):
            weights = np.full(wd_size, weights, dtype=float)

        assert len(weights) == wd_size, "Q4: The size of the weights must be"

        # YOUR CODE HERE
        raise NotImplementedError()

        return wma
```

```
In [ ]: # Autograder tests

wd_size = 7
weights = np.arange(1, wd_size + 1).astype(float) # linear weighting
stu_ser = load_data()
stu_ans = calc_wma(stu_ser, wd_size, weights)

assert isinstance(stu_ans, np.ndarray), "Q4: Your function should return"
assert stu_ans.shape == (len(stu_ser), ), "Q4: The np.ndarray returned is"
assert np.issubdtype(stu_ans.dtype, np.floating), "Q4: Your np.ndarray st"

# Some hidden tests

del wd_size, weights, stu_ser, stu_ans
```

```
In [ ]: # Let's plot and see your WMA

fig, ax = plt.subplots(figsize=(10, 6))
wd_size = 7
weights = np.arange(1, wd_size + 1)
ser = load_data()
wma = calc_wma(ser, wd_size, weights=weights)

ax.plot(ser.index, ser.values, label="Original")
ax.plot(ser.index, wma, label="WMA")
ax.set_xlabel("Day")
ax.set_ylabel("# Cases")
ax.set_title("Daily New COVID-19 Cases Worldwide")
ax.legend()

del fig, ax, wd_size, weights, ser, wma
```



## Question 5: Calculate "Time" Exponential Moving Average (EMA) (10 pts)

WMA usually works well if each data point is sampled at regular time intervals (which is the case for our time series). "Time" Exponential Moving Average (EMA), on the other hand, works well on both regular and irregular time series. Let's now explore how to apply EMA to our time series.

Recall from the lectures that a "time" EMA method applies the following transformation to each data point  $x_j$ :

$$x'_j = \frac{\sum_{i=1}^j \exp[-\lambda (t_j - t_i)] x_i}{\sum_{i=1}^j \exp[-\lambda (t_j - t_i)]}$$

where  $0 \leq \lambda \leq 1$  is the "decay rate". Also note that, when  $\lambda = 0$ , this is equivalent to a cumulative moving average (CMA). Complete the function below that calculates the "time" EMA for an input time series, **assuming the time intervals are days**.

**This function should return a `np.ndarray` of shape `(len(ser), )`, which represents the "time" EMA for the input time series.**

```
In [ ]: def calc_time_ema(ser, lmbd=0.0):
    """
    Takes in a series and calculates EMA with the lambda provided
    """

    time_ema = None

    # YOUR CODE HERE
    raise NotImplementedError()

    return time_ema
```

```
In [ ]: # Autograder tests

stu_ser = load_data()

# Sanity checks for a trivial case - CMA
stu_ans = calc_time_ema(stu_ser, lmbd=0.0)

assert isinstance(stu_ans, np.ndarray), "Q5: Your function should return
assert stu_ans.shape == (len(stu_ser), ), "Q5: The np.ndarray returned is
assert np.issubdtype(stu_ans.dtype, np.floating), "Q5: Your np.ndarray sl
assert np.isclose(stu_ans, np.cumsum(stu_ser) / np.arange(1, len(stu_ser)

# Some hidden tests

del stu_ser, stu_ans
```

```
In [ ]: # Let's plot and see your time EMA

fig, ax = plt.subplots(figsize=(10, 6))
ser = load_data()
ema = calc_time_ema(ser, lmbd=0.5)

ax.plot(ser.index, ser.to_numpy(), label="Original")
ax.plot(ser.index, ema, label="Time EMA")
ax.set_xlabel("Day")
ax.set_ylabel("# Cases")
ax.set_title("Daily New COVID-19 Cases Worldwide")
ax.legend()

del fig, ax, ser, ema
```

The `SimpleExpSmoothing`

([https://www.statsmodels.org/stable/examples/notebooks/generated/exponential\\_smoothing.html](https://www.statsmodels.org/stable/examples/notebooks/generated/exponential_smoothing.html))

class from the `statsmodels` library is a handy tool for EMA. See an example below.

```
In [ ]: from statsmodels.tsa.api import SimpleExpSmoothing

fig, ax = plt.subplots(figsize=(10, 6))

ser = load_data()
ema = SimpleExpSmoothing(ser, initialization_method=None).fit(smoothing_

ax.plot(ser, label="Original")
ax.plot(ema.fittedvalues, label="EMA")
ax.set_xlabel("Day")
ax.set_ylabel("# Cases")
ax.set_title("Daily New COVID-19 Cases Worldwide")
ax.legend()

del ser, ema, fig, ax
```