

Part 1: Cleaning Dataset

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import gamma
import math

In [ ]:
df = pd.read_csv('23.csv', header=0, names=['date', 'UT_conf', 'VA_conf', 'UT_death', 'VA_death'])
df[['UT_conf', 'VA_conf', 'UT_death', 'VA_death']] = df[['UT_conf', 'VA_conf', 'UT_death', 'VA_death']].diff().ff(
    print(df)

date UT_conf VA_conf UT_death VA_death
0 2020-01-22 0.0 0.0 0.0 0.0
1 2020-01-23 0.0 0.0 0.0 0.0
2 2020-01-24 0.0 0.0 0.0 0.0
3 2020-01-25 0.0 0.0 0.0 0.0
4 2020-01-26 0.0 0.0 0.0 0.0
...
432 2021-03-25 371.0 141.0 0.0 0.0
433 2021-03-26 371.0 141.0 0.0 0.0
434 2021-03-31 514.0 1057.0 4.0 8.0
435 2021-04-01 487.0 1047.0 2.0 11.0
436 2021-04-02 422.0 1523.0 5.0 11.0
437 2021-04-03 447.0 1523.0 1.0 6.0

[438 rows x 5 columns]

In [ ]:
print('Outliers')
labels = ['UT_conf', 'VA_conf', 'UT_death', 'VA_death']
for label in labels:
    print(label)
    Q1 = df[label].quantile(0.25)
    Q3 = df[label].quantile(0.75)
    IQR = Q3 - Q1
    lower_outlier = df.query('(< (Q1 - 1.5 * IQR)')
    large_outlier = df.query('(> (Q3 + 1.5 * IQR)')
    if len(smaller_outlier) > 0:
        print(smaller_outlier[['date', label]])
    if len(larger_outlier) > 0:
        print(larger_outlier[['date', label]])
    print()

Outliers
UT_conf
ABOVE Q3
date UT_conf
277 2020-10-25 3276.0
280 2020-11-23 2785.0
289 2020-11-06 2630.0
289 2020-11-07 3047.0
293 2020-11-10 2768.0
295 2020-11-12 3791.0
297 2020-11-14 5926.0
298 2020-11-17 3212.0
301 2020-11-18 3088.0
309 2020-11-19 3095.0
303 2020-11-20 4563.0
304 2020-11-21 3379.0
305 2020-11-22 3231.0
310 2020-11-27 5614.0
317 2020-12-02 4095.0
317 2020-12-04 3095.0
320 2020-12-05 4735.0
323 2020-12-10 3362.0
325 2020-12-12 3784.0
329 2020-12-16 2888.0
330 2020-12-17 3214.0
331 2020-12-08 3094.0
332 2020-12-09 4491.0
333 2020-12-10 3037.0
334 2020-12-11 3021.0
335 2020-12-12 4173.0
339 2020-12-16 3028.0
339 2020-12-17 3098.0
340 2020-12-18 3684.0
340 2020-12-19 3097.0
343 2020-12-30 2996.0
344 2020-12-31 4545.0
345 2021-01-06 3007.0
351 2021-01-07 4542.0
352 2021-01-08 4527.0
357 2021-01-13 5053.0
358 2021-01-14 2621.0

VA_conf
ABOVE Q3
date VA_conf
318 2020-12-05 3782.0
319 2020-12-06 3805.0
320 2020-12-07 3764.0
322 2020-12-08 3094.0
329 2020-12-09 4491.0
333 2020-12-10 3037.0
334 2020-12-11 3021.0
335 2020-12-12 4173.0
339 2020-12-16 3028.0
339 2020-12-17 3098.0
340 2020-12-18 3684.0
340 2020-12-19 3097.0
343 2020-12-30 2996.0
344 2020-12-31 4545.0
345 2021-01-06 3007.0
351 2021-01-07 4542.0
352 2021-01-08 4527.0
357 2021-01-13 5053.0
358 2021-01-14 2621.0

UT_death
ABOVE Q3
date UT_death
305 2020-11-25 26.0
308 2020-11-25 26.0
314 2020-12-01 26.0
321 2020-12-08 23.0
322 2020-12-09 23.0
323 2020-12-10 21.0
329 2020-12-16 21.0
330 2020-12-17 29.0
331 2020-12-18 30.0
333 2020-12-08 3094.0
334 2020-12-09 4491.0
335 2020-12-10 3037.0
336 2020-12-11 3021.0
337 2020-12-12 4173.0
339 2020-12-16 3028.0
339 2020-12-17 3098.0
340 2020-12-18 3684.0
340 2020-12-19 3097.0
343 2020-12-30 2996.0
344 2020-12-31 4545.0
345 2021-01-06 3007.0
351 2021-01-07 4542.0
352 2021-01-08 4527.0
357 2021-01-13 5053.0
358 2021-01-14 2621.0

VA_death
ABOVE Q3
date VA_death
237 2020-09-15 96.0
330 2020-12-17 93.0
353 2021-01-09 84.0
354 2021-01-10 84.0
357 2021-01-13 74.0
358 2021-01-14 74.0
365 2021-01-21 78.0
367 2021-01-23 76.0
370 2021-01-26 83.0
372 2021-01-28 86.0
373 2021-01-29 71.0
374 2021-01-30 70.0
379 2021-02-04 76.0
380 2021-02-05 82.0
384 2021-02-09 78.0
395 2021-02-20 95.0
397 2021-02-22 155.0
397 2021-02-23 171.0
399 2021-02-23 147.0
400 2021-02-25 154.0
401 2021-02-26 235.0
402 2021-02-27 182.0
403 2021-02-28 169.0
404 2021-03-01 231.0
405 2021-03-02 160.0
406 2021-03-03 304.0
408 2021-03-05 73.0
409 2021-03-06 90.0
410 2021-03-07 90.0
411 2021-03-08 87.0
412 2021-03-09 111.0
```

From the values above, we found that most of our dataset contains outliers. Initially, this was quite surprising to us and we thought something must have been done wrong. However, after looking at the COVID 19 dataset, it actually seems to be the case that these two states did indeed have outliers and both have unique ways of reporting their data.

We begin with Virginia. According to the data above, we see that January, February, and March has the majority of the death outliers in the dataset. The Virginia deaths are not really seemed to have been reported consistently, where some days there were dips (i.e. 0 deaths reported, 3 deaths reported, and then massive jumps. If you search up COVID deaths in Virginia on Google, there are a lot of articles saying Virginia had data glitches which resulted in death "data dumps" in February and March. A lot of the deaths reported in February and March included these dumps, so it makes sense that there are so many outliers in the Virginia death dataset.

Regarding the Virginia confirmed cases outliers, almost all of the outliers were in December and January. This corresponds with the winter surge in America. Interestingly enough, the data some days was also sporadic like the deaths data, where some days there were few cases and then a big spike. As the holiday season approached, there were also much lower confirmed cases on days like Christmas and New Years. This makes sense and happened not only in Virginia, but also across America. Finally, in January, there were massive spikes reported, 3 deaths reported in January 17th. With these spikes also came much lower cases on preceding on the next days, so it makes sense that we detected so many outliers for this month. Perhaps this could have something to do with the riots, election, and inauguration day since DC is nearby.

Finally, we look at the outliers for Utah. Regarding the deaths in Utah, Utah has a very jagged trend line if you search up the Utah deaths graph on Google. There are a lot of days where Utah reports few deaths, and then it jumps up. We see this in November around Thanksgiving, and also in a lot of December and January. Utah's death graph has a lot of gaps where some days there are 0 deaths and other days it jumps up, so it makes sense that so many outliers were detected for Utah.

Regarding cases in Utah, there is a similar trend too. For example, Utah reported 0 cases on Christmas and on January 1st, and then had resulting spikes the next couple of days. Perhaps this could be because Utah has a lot of Mormons and they closed all testing centers on Christmas and New Years. The Utah cases data seems to be in general smoother than their deaths data, but there are once again days where they go from lower reporting to higher reporting.

Note: since there are too many outliers that were found from Tukey's Rule, we do not remove the outliers from the dataset and use the original dataset for our required inferences

Part 2: Required Inferences

A: AR, EWMA, MSE, MAPE

EWMA

```
In [ ]:
import pandas as pd

df = pd.read_csv('23.csv', header=0, names=['date', 'UT_conf', 'VA_conf', 'UT_death', 'VA_death'])
df[['UT_conf', 'VA_conf', 'UT_death', 'VA_death']] = df[['UT_conf', 'VA_conf', 'UT_death', 'VA_death']].diff().ff(1)

for alpha in [0.5, 0.8]:
    print('alpha =', alpha)
    cols = ['UT_conf', 'VA_conf', 'UT_death', 'VA_death']
    ewma = {}
    mape = {}
    for col in cols:
        ewma[col] = training[col].tolist()[0]
        mape[col] = 0
        # Calculate MAPE for the first three weeks of august
        for i in range(1, 21):
            for col in cols:
                ewma[col] = alpha * df[col].tolist()[i] + (1 - alpha) * ewma[col]
                print('Day {} = {}'.format(i, ewma[col]))
            sample = df[col].tolist()[21 + i]
            # Calculate sum of squared errors
            mse[col] += (ewma[col] - sample) ** 2
            # Ignore 0 denominators
            if sample != 0:
                mape[col] += abs(ewma[col] - sample) / sample
            # Incorporate new sample into EWMA
            ewma[col] = alpha * df[col].tolist()[21 + i] + (1 - alpha) * ewma[col]
        print()

# Print accuracy results
for col in cols:
    print(col, 'MSE:', mse[col] / 7.0)
    print(col, 'MAPE:', mape[col] * 100.0 / 7.0)
    print()

alpha = 0.5
Day 22
UT_conf 449.93113362291504
VA_conf 994.4607778641348
UT_death 11.498978614307129
VA_death 11.467790603637695
Day 23
UT_conf 408.46556568185475
VA_conf 1558.4380832699223
UT_death 5.053885164254440
VA_death 4.748485307460565
UT_conf 976.0151944168461
VA_conf 2.374744653717022
UT_death 16.61694765999424
Day 25
UT_conf 293.3663914203644
VA_conf 829.0075972080231
UT_death 2.107373125808091
VA_death 10.30847826454712
Day 26
UT_conf 331.1831957301822
VA_conf 912.5037986040115
UT_death 5.053885164254440
VA_death 16.61694765999424
Day 27
UT_conf 375.091597850911
VA_conf 867.751893020958
UT_death 4.04684308172723
VA_death 16.61694765999424
Day 28
UT_conf 376.04578092734555
VA_conf 994.3759496510929
UT_death 3.5234215408503614
VA_death 15.41359528181839
UT_conf MSE: 6886.249511860623
UT_conf MAPE: 22.04778318019756
VA_conf MSE: 46542.94135325959
VA_conf MAPE: 20.6305287530803
UT_death MSE: 12.520215663308024
UT_death MAPE: 41.8512080340922
VA_death MSE: 97.3300723309794
VA_death MAPE: 86.64664928419849

alpha = 0.8
Day 22
UT_conf 406.5796172620181
VA_conf 950.815067498566
UT_death 0.316390674308505
VA_death 10.444182113439557
Day 23
UT_conf 391.71592255240364
VA_conf 1159.7631214981152
UT_death 6.46327813487371
VA_death 16.68836422686112
Day 24
UT_conf 296.7431845104887
VA_conf 947.1528242092227
UT_death 1.292655269747418
VA_death 20.737767284537224
Day 25
UT_conf 256.14863690209614
VA_conf 729.6385248091538
UT_death 3.131068240981508
VA_death 7.34755456907444
Day 26
UT_conf 346.4297273804192
VA_conf 948.1261849719889
UT_death 2.0524016092329
VA_death 19.8695160138140
Day 27
UT_conf 404.4859454760939
VA_conf 848.0252209043939
UT_death 2.0524016092329
VA_death 20.7739021382763
Day 28
UT_conf 382.4971809521675
VA_conf 1066.4680441988788
UT_death 3.131068240981508
VA_death 16.61694765999424
UT_conf MSE: 7272.562832934679
UT_conf MAPE: 22.376868252382328
VA_conf MSE: 56130.37144951825
VA_conf MAPE: 23.86206213686824
UT_death MSE: 18.958839270788146
UT_death MAPE: 52.52624016092329
VA_death MSE: 136.11454767014184
VA_death MAPE: 183.19426016019078

There were a few days in the last week of August where there were 0 new confirmed cases which had to be ignored in the MAPE calculation.

Each day in the last week was incorporated into the prediction for the next day, which is how this analysis would be done in real time.
```

AR

```
In [ ]:
'''AR(3) & AR(5)'''
import pandas as pd
import numpy as np
# Import matplotlib.pyplot as plt

from sklearn.linear_model import LinearRegression
from statsmodels.tsa.ar_model import Autoreg
from statsmodels.tsa.stattools import adfuller
from sklearn.metrics import mean_squared_error

File = 'https://raw.githubusercontent.com/michaelosbu/CSE-544-Datasets/main/States%20data/23.csv'
state = pd.read_csv(File)
state['date'] = pd.to_datetime(state['date'])

aug = (state['date'] >= '2020-08-01') & (state['date'] <= '2020-08-31')
state_aug = state[aug]
state_aug = state_aug.reset_index()

# mean absolute percent error
def mape(actual, Predicted):
    mape = np.mean(np.abs((actual - Predicted)/actual))*100
    return mape

# mean squared error
def msev(actual, Predicted):
    mse = np.square(np.subtract(actual, Predicted)).mean()
    return mse

# Autoregression model and p as the lags value
def ar(p, df, train_size, col):
    # Add f-p columns to dataframe
    for i in range(1, p+1):
        df[f'x_{i}'] = df[col].shift(i)
    # Breaking data into training and test set and remove the null value
    train_data = df.dropna(subset=[col])
    test = pd.DataFrame(df[train_size:].dropna(subset=[col]))

    # Separate the f-p columns as data, first column as labels
    train_data = train_data.dropna(subset=[col])
    train_data = train_data.reset_index(drop=True)
    test_data = test_data.dropna(subset=[col])
    test_data = test_data.reset_index(drop=True)

    # Running linear regression to generate the coefficients and intercepts
    lr = LinearRegression()
    lr.fit(train_data, train_data[col])

    # Apply the autoregression formula y(t) = intercept + i coefficient * y(t-i)
    pred = test_data.dot(lr.coef_) + lr.intercept_
    pred = pred.flatten()

    # Compute the mse and mape value
    mse = msev(test[col], pred)
    mape = mape(test[col], pred)

    # Report the accuracy
    print('AR(' + str(p) + ') mape: ', round(mape, 2), '%\tmse:', round(mse))

for col in state_aug.columns[2:]:
    data = state_aug[col]
    data = pd.DataFrame(state_aug[col])

    print(col)
    ar(3, data, 21, col)
    ar(5, data, 21, col)
    print()

AR(3) mape: 0.14 % mse: 6614
AR(5) mape: 0.33 % mse: 36549

VA confirmed mse: 143442
AR(3) mape: 0.28 % mse: 13
AR(5) mape: 0.48 % mse: 366682

UT deaths
AR(3) mape: 0.71 % mse: 9
AR(5) mape: 0.8 % mse: 13

VA deaths
AR(3) mape: 0.66 % mse: 351
AR(5) mape: 0.5 % mse: 244
```

B: Wald's test, Z test, and t-test

Wald's One Sample

```
In [ ]:
import pandas as pd
import numpy as np
state = pd.read_csv('23.csv')

state['date'] = pd.to_datetime(state['date'])
state_feb = (state['date'] >= '2021-02-01') & (state['date'] <= '2021-02-28')
state_mar = (state['date'] >= '2021-03-01') & (state['date'] <= '2021-03-31')
state_feb = state_feb[['date', 'VA_confirmed']]
state_mar = state_mar[['date', 'VA_confirmed']]
cases_count_VA_feb = cases_count_VA_mar[['date', 'VA_confirmed']]
deaths_count_VA_feb = deaths_count_VA_mar[['date', 'VA_confirmed']]
cases_count_UT_feb = cases_count_UT_mar[['date', 'VA_confirmed']]
deaths_count_UT_feb = deaths_count_UT_mar[['date', 'VA_confirmed']]
cases_count_VA_mar = cases_count_VA_mar[['date', 'VA_confirmed']]
deaths_count_VA_mar = deaths_count_VA_mar[['date', 'VA_confirmed']]
cases_count_UT_mar = cases_count_UT_mar[['date', 'VA_confirmed']]
deaths_count_UT_mar = deaths_count_UT_mar[['date', 'VA_confirmed']]

In [ ]:
def walds_one(feb, march):
    # Here we are using Poisson MLE.
    # The MLE for Poisson is simply the sample mean, so we just take the mean
    theta_hat = np.mean(feb)

    # We were told to use sample mean for theta_0
    theta_0 = np.mean(march)

    #se = np.sqrt(theta_hat / n)
    se = np.sqrt(theta_hat / len(march))

    #Compute Wald's statistic
    walds_statistic = (theta_hat - theta_0) / se

    #Critical threshold for alpha = 0.05 is 1.96
    if np.abs(walds_statistic) > 1.96:
        print('["I"] = "np.abs(walds_statistic), "> 1.96, therefore reject the Null Hypothesis")
    else:
        print('["I"] = "np.abs(walds_statistic), "<= 1.96, therefore accept the Null Hypothesis")

    walds_one(cases_count_UT_feb, cases_count_UT_mar)
    walds_one(cases_count_VA_feb, cases_count_VA_mar)
    walds_one(deaths_count_VA_feb, deaths_count_VA_mar)

["I"] = "np.abs(walds_statistic) > 1.96, therefore reject the Null Hypothesis"
["I"] = "np.abs(walds_statistic) > 1.96, therefore reject the Null Hypothesis"
["I"] = "np.abs(walds_statistic) > 1.96, therefore reject the Null Hypothesis"
["I"] = "np.abs(walds_statistic) > 1.96, therefore reject the Null Hypothesis"

In all of these cases the Wald's statistic W was much greater than the threshold value of 1.96. Based off of this, we can say the following:

For Virginia values, the mean number of cases AND the mean of daily deaths for Feb/21 is different from the corresponding mean of daily values for March/21

For Utah values, the mean number of cases AND the mean of daily deaths for Feb/21 is different from the corresponding mean of daily values for March/21

It seems that in this case, the 2 sample Wald's test is NOT applicable for this test. Wald's test assumes that the estimator theta_hat is AN, but in this case, the data is not AN.
```

Z test one sample

```
In [ ]:
def z_one(original, feb, march):
    #Z test requires true standard deviation, so we need to find that
    sigma = np.std(original)

    x_bar = np.mean(feb)
    mu_0 = np.mean(march)

    sigma_over_sqrt_n = sigma / np.sqrt(len(feb))

    z_statistic = (x_bar - mu_0) / sigma_over_sqrt_n

    if np.abs(z_statistic) > 1.96:
        print('["I"] = "np.abs(z_statistic), "> 1.96 therefore reject the null hypothesis")
    else:
        print('["I"] = "np.abs(z_statistic), "<= 1.96 therefore accept the null hypothesis")

    z_one(state['VA_confirmed'], cases_count_VA_feb, cases_count_VA_mar)
    z_one(state['VA deaths'], deaths_count_VA_feb, deaths_count_VA_mar)
    z_one(state['UT_confirmed'], cases_count_UT_feb, cases_count_UT_mar)
    z_one(state['UT deaths'], deaths_count_UT_feb, deaths_count_UT_mar)

["I"] = "np.abs(z_statistic) <= 1.96 therefore accept the null hypothesis"
["I"] = "np.abs(z_statistic) <= 1.96 therefore accept the null hypothesis"
["I"] = "np.abs(z_statistic) <= 1.96 therefore accept the null hypothesis"
["I"] = "np.abs(z_statistic) <= 1.96 therefore accept the null hypothesis"

In the case of mean monthly cases for Utah, and the case of mean monthly deaths for Utah and Virginia, here since |Z| <= 1.96, we accept the null hypothesis. This means that it is deemed true that the means are the same in February and March for these.
```

One exception to this is the mean monthly deaths in Virginia. Our Z value here was 5.11 which is greater than the critical threshold, so here we reject the null hypothesis and claim that the means are not the same.

We believe that the Z test is not applicable here. One of the criteria for the Z test is knowing the true standard deviation. However, as we see in the news and social media all the time, it seems that the true standard deviation of COVID cases can never really be known. This is because there is indeed a lot of asymptomatic cases which do not show up in the numbers, some false negatives, mistakes in testing like PCR thresholds, and finally people just simply refusing to get tested. These could all be contributing factors to now knowing the true standard deviation, so in this case, the Z test is not really applicable since it assumes we know the true standard deviation.

Wald's 2 sample

```
In [ ]:
def walds_two(feb, march):
    x_bar = np.mean(feb)
    y_bar = np.mean(march)
    delta_hat = x_bar - y_bar

    #Since we're using Poisson MLE, variance will be sample mean
    se_hat = np.sqrt(x_bar / len(feb)) + (y_bar / len(march))

    walds_statistic = delta_hat / se_hat

    if np.abs(walds_statistic) > 1.96:
        print('["I"] = "np.abs(walds_statistic), "> 1.96, therefore reject the Null Hypothesis")
    else:
        print('["I"] = "np.abs(walds_statistic), "<= 1.96, therefore accept the Null Hypothesis")

    walds_two(cases_count_UT_feb, cases_count_UT_mar)
    walds_two(deaths_count_UT_feb, deaths_count_UT_mar)
    walds_two(cases_count_VA_feb, cases_count_VA_mar)
    walds_two(deaths_count_VA_feb, deaths_count_VA_mar)

["I"] = "np.abs(walds_statistic) > 1.96, therefore reject the Null Hypothesis"
["I"] = "np.abs(walds_statistic) > 1.96, therefore reject the Null Hypothesis"
["I"] = "np.abs(walds_statistic) > 1.96, therefore reject the Null Hypothesis"
["I"] = "np.abs(walds_statistic) > 1.96, therefore reject the Null Hypothesis"

Similarly to the one sample Wald's test, in all of these cases the Wald's statistic W was much greater than the threshold value of 1.96. Based off of this, we can say the following:

For Virginia values, the mean number of cases AND the mean of daily deaths for Feb/21 is different from the corresponding mean of daily values for March/21

For Utah values, the mean number of cases AND the mean of daily deaths for Feb/21 is different from the corresponding mean of daily values for March/21

It seems that in this case, the 2 sample Wald's test is NOT applicable for this test. Wald's test assumes that the estimator theta_hat is AN, but in this case, the data is not AN.
```

t-test one sample

```
In [ ]:
def t_one(feb, march):
    x_bar = np.mean(feb)
    mu_0 = np.mean(march)
    ssd = np.std(feb)

    s_over_n = ssd / np.sqrt(len(feb))
    t = (x_bar - mu_0) / s_over_n
    #Critical threshold is 2.051831 found by table lookup
    if np.abs(t) > 2.051831:
        print('["I"] = "np.abs(t), "> 2.051831 therefore reject the null hypothesis")
    else:
        print('["I"] = "np.abs(t), "<= 2.051831 therefore accept the null hypothesis")

    t_one(cases_count_UT_feb, cases_count_UT_mar)
    t_one(deaths_count_UT_feb, deaths_count_UT_mar)
    t_one(cases_count_VA_feb, cases_count_VA_mar)
    t_one(deaths_count_VA_feb, deaths_count_VA_mar)

["I"] = "np.abs(t) > 2.051831 therefore reject the null hypothesis"
["I"] = "np.abs(t) > 2.051831 therefore reject the null hypothesis"
["I"] = "np.abs(t) > 2.051831 therefore reject the null hypothesis"
["I"] = "np.abs(t) > 2.051831 therefore reject the null hypothesis"

In all of these cases the Wald's T was greater than the threshold value of 2.051831. Based off of this, we can say the following:

For Virginia values, the mean number of cases AND the mean of daily deaths for Feb/21 is different from the corresponding mean of daily values for March/21

For Utah values, the mean number of cases AND the mean of daily deaths for Feb/21 is different from the corresponding mean of daily values for March/21

It seems that in this case, the T test is NOT applicable for this test. T test assumes that the data is normally distributed, and we do not know that.
```

Two sample unpaired t-test

```
In [ ]:
def t_two_unpaired(feb, march):
    x_bar = np.mean(feb)
    y_bar = np.mean(march)
    d_bar = x_bar - y_bar

    sx = np.var(feb) / len(feb)
    sy = np.var(march) / len(march)
    s_x_y = np.sqrt(sx + sy)

    t = (d_bar - 0) / s_x_y
    #Critical threshold is 2.02465 found by table lookup
    if np.abs(t) > 2.02465:
        print('["I"] = "np.abs(t), "> 2.02465 therefore reject the null hypothesis")
    else:
        print('["I"] = "np.abs(t), "<= 2.02465 therefore reject the null hypothesis")

    t_two_unpaired(cases_count_UT_feb, cases_count_UT_mar)
    t_two_unpaired(deaths_count_UT_feb, deaths_count_UT_mar)
    t_two_unpaired(cases_count_VA_feb, cases_count_VA_mar)
    t_two_unpaired(deaths_count_VA_feb, deaths_count_VA_mar)

["I"] = "np.abs(t) > 2.02465 therefore reject the null hypothesis"
["I"] = "np.abs(t) > 2.02465 therefore reject the null hypothesis"
["I"] = "np.abs(t) > 2.02465 therefore reject the null hypothesis"
["I"] = "np.abs(t) > 2.02465 therefore reject the null hypothesis"

In all of these cases the Wald's T was greater than the threshold value of 2.02465. Based off of this, we can say the following:

For Virginia values, the mean number of cases AND the mean of daily deaths for Feb/21 is different from the corresponding mean of daily values for March/21

For Utah values, the mean number of cases AND the mean of daily deaths for Feb/21 is different from the corresponding mean of daily values for March/21

It seems that in this case, the 2 sample unpaired T test is NOT applicable for this test. T test assumes that the data samples for X and Y are normally distributed, and we do not know that.
```

C: KS and Permutation

Two Sample KS test

```
In [ ]:
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import ks_2samp, ks_1samp, bisect_right
from scipy import warnings
warnings.simplefilter(action='ignore', category=SettingWithCopyWarning)

def plot_ecdf(f, labels):
    x = np.sort(f)
    yf = np.arange(len(x))/float(len(x))
    plt.step(x, yf, label=labels)
    return x, yf

def ks_2_sample(X, Y):
    x1, y1 = plot_ecdf(X, 'UT')
    x2, y2 = plot_ecdf(Y, 'VA')

    data_all = np.concatenate([x1, x2])
    #using searchsorted solves equal data problem
    idx1 = np.searchsorted(x1, data_all, side='right')
    idx2 = np.searchsorted(x2, data_all, side='right')
    cdf2 = idx2 / n2
    cdiffs = cdf1 - cdf2
    minS = np.min(cdiffs)
    maxS = np.max(cdiffs)
    idx = np.argmax(cdiffs)
    else:
        minS = maxS
        idx = np.argmax(cdiffs)
    maxdiff = x[idx]
    plt.xlabel('Num Cases')
    plt.ylabel('CDF')
    plt.show()

def run_ks(X, Y):
    print(stats.ks_2samp(X, Y))
    ks_2_sample(X, Y)

def main():
    df = pd.read_csv('23.csv', header=0, names=['date', 'UT_conf', 'VA_conf', 'UT_death', 'VA_death'])
    filtered_df = df[df['date'] == '2020-12-31'].index[0] + 1
    filtered_df = filtered_df[['date', 'UT_death', 'VA_death']]
    filtered_df = filtered_df[['date', 'UT_death', 'VA_death']]

    # For confirmed
    run_ks(X, Y)
    # For deaths
    run_ks(X, Y)
    # For confirmed
    run_ks(X, Y)
    # For deaths
    run_ks(X, Y)

ks_2sampresult(statistic=0.16304347826869857, pvalue=0.1736682169213413)

Deaths
ks_2sampresult(statistic=0.30695652173913043, pvalue=0.0189393748185584e-05)
```

Our KS hypothesis is that the distribution for UT and VA are the same for confirmed and deaths respectively for the two graphs above.

From the 2-sample KS test above, we see that for the confirmed cases, we see a p-value of 0.174 which is greater than 0.05. Thus, we cannot reject the null and we maintain that the UT and VA confirmed cases are distributed the same.

For deaths, we find that the p-value is very close to zero. Thus, we can reject the null, and the UT and VA deaths are distributed differently.

One Sample KS Test

