

Introduction

This study is to analyze mortality in the US based on causes across states. This study aims to answer some of questions such as are American facing increasing, decreasing or steady trend of death, what are the four leading causes of death, do individual states show the same 4 leading causes of death and if the year-by-year changes in the 4 leading causes of death nationwide.

There are 2 datasets used in this study and imported into the database. This study will be analyzed by Python programming.

Data Source

Publicly available files will be used.

- The first contains data on causes of death: NCHS_-_Leading_Causes_of_Death_United_States.csv
- The second contains population data: rst-es2018-01.xlsx

Both files have state-level information for multiple years.

0

Year	113 Cause Name	Cause Name	State	Deaths	Age-adjusted Death Rate	
0	2012	Nephritis, nephrotic syndrome and nephrosis (N...	Kidney disease	Vermont	21	2.6
1	2016	Nephritis, nephrotic syndrome and nephrosis (N...	Kidney disease	Vermont	30	3.7
2	2013	Nephritis, nephrotic syndrome and nephrosis (N...	Kidney disease	Vermont	30	3.8
3	2000	Intentional self-harm (suicide) (U03.X00-X04...	Suicide	District of Columbia	23	3.8
4	2014	Nephritis, nephrotic syndrome and nephrosis (N...	Kidney disease	Arizona	325	4.1

Getting to understand the table

Getting to understand the table

But here Column 113 Cause Name is not required in the project and hence will be dropped.

```
df=df.drop('113 Cause Name', axis=1)
df.head()
```

	Year	Cause Name	State	Deaths	Age-adjusted Death Rate
0	2012	Kidney disease	Vermont	21	2.6
1	2016	Kidney disease	Vermont	30	3.7
2	2013	Kidney disease	Vermont	30	3.8
3	2000	Suicide	District of Columbia	23	3.8
4	2014	Kidney disease	Arizona	325	4.1

1. Are Americans facing increasing, decreasing, or steady likelihood of death?

Total of deaths by year from 1999 to 2016

```
# not by All causes:
df1 = df[(df['State']=='United States') & (df['Cause Name'] != 'All causes')]

df1a = df1.groupby(['Year'])['Deaths'].sum().reset_index()
df1a
```

	Year	Deaths
0	1999	1905826

1. Are Americans facing increasing, decreasing, or steady likelihood of death?

Total of deaths by year from 1999 to 2016

```
df = df[df['State']=='United States'] & (df['Cause Name'] != 'All causes']
df1b = df.groupby(['Year'])['Deaths'].sum().reset_index()
df1b
```

Year	Deaths
0	1999 2291399
1	2000 2403951
2	2001 2416425
3	2002 2443387
4	2003 2448288
5	2004 2397615
6	2005 2448017
7	2006 2426264

Analysis:

- Americans are facing increasing trend likelihood of deaths (exclude All causes)

```
2011 2515458
13 2012 2543279
14 2013 2596993
15 2014 2626418
16 2015 2712630
17 2016 2724248
```

Analysis:

- Americans are facing increasing trend of deaths (only All causes)

2. What are the 4 leading causes of death for Americans?

Top 4 causes of deaths from 1999 to 2016

```
df2 = df[(df['State']=='United States') & (df['Cause Name'] != 'All causes')]
df2b = df2.groupby(['Cause Name'])['Deaths'].sum()
sorted_grouped = df2b.sort_values(ascending=False)
top_4 = sorted_grouped.nlargest(4)
top_4
```

Cause Name	
Heart disease	11575183
Cancer	10244936
Stroke	2580140

Analysis:

- Americans are facing increasing trend of deaths (only All causes)

2. What are the 4 leading causes of death for Americans?

Top 4 causes of deaths from 1999 to 2016

In [7]:

```
df2 = df[(df['State']=='United States') & (df['Cause Name'] != 'All causes')]
df2b = df2.groupby(['Cause Name'])['Deaths'].sum()
sorted_grouped = df2b.sort_values(ascending=False)
top_4 = sorted_grouped.nlargest(4)
top_4
```

Out [7]:

Cause Name	Deaths
Heart disease	11575183
Cancer	10244536
Stroke	2380140
CLRD	2434726
Name: Deaths, dtype: int64	

Analysis:

- Top 4 causes of death for Americans are: Heart disease, Cancer, Stroke and CLRD

3. Do individual states show the same four leading causes of death?

In [8]:

```
us_deaths_causes_df1 = df[(df['State'] != 'United States') & (df['Cause Name'] != 'All causes')]
rows = []

for state in us_deaths_causes_df1['State'].unique():
    state_cause_ydf = us_deaths_causes_df1[us_deaths_causes_df1['State'] == state].groupby(['Cause Name'])['Deaths'].sum()
    state_cause_df.sort_values(ascending=False, inplace=True)

    # Get the top 4 causes
    cause_list = list(state_cause_df.keys())[:4] + [None] * (4 - len(state_cause_df))

    rows.append(['State': state, '1st cause': cause_list[0], '2nd cause': cause_list[1],
                '3rd cause': cause_list[2], '4th cause': cause_list[3]])

df3 = pd.DataFrame(rows)
df3.sort_values(by='State', ascending=True, inplace=True)
display(df3)
```

Out [8]:

	State	1st cause	2nd cause	3rd cause	4th cause
45	Alabama	Heart disease	Cancer	Stroke	CLRD
12	Alaska	Cancer	Heart disease	Unintentional injuries	Stroke
2	Arizona	Heart disease	Cancer	Unintentional injuries	CLRD
50	Arkansas	Heart disease	Cancer	Stroke	CLRD
9	California	Heart disease	Cancer	Stroke	CLRD
24	Colorado	Cancer	Heart disease	Unintentional injuries	CLRD
16	Connecticut	Heart disease	Cancer	Stroke	CLRD
27	Delaware	Heart disease	Cancer	CLRD	Stroke
1	District of Columbia	Heart disease	Cancer	Stroke	Unintentional injuries
23	Florida	Heart disease	Cancer	CLRD	Stroke
36	Georgia	Heart disease	Cancer	Stroke	Unintentional injuries
21	Hawaii	Heart disease	Cancer	Stroke	Unintentional injuries
15	Idaho	Heart disease	Cancer	CLRD	Stroke
18	Illinois	Heart disease	Cancer	Stroke	CLRD
38	Indiana	Heart disease	Cancer	CLRD	Stroke
5	Iowa	Heart disease	Cancer	Stroke	CLRD
42	Kansas	Heart disease	Cancer	Stroke	CLRD
46	Kentucky	Heart disease	Cancer	CLRD	Unintentional injuries
37	Louisiana	Heart disease	Cancer	Unintentional injuries	Stroke
32	Maine	Cancer	Heart disease	CLRD	Stroke
22	Maryland	Heart disease	Cancer	Stroke	CLRD
7	Massachusetts	Cancer	Heart disease	Stroke	CLRD
31	Michigan	Heart disease	Cancer	Stroke	CLRD
19	Minnesota	Cancer	Heart disease	Stroke	Unintentional injuries
37	Mississippi	Heart disease	Cancer	Unintentional injuries	Stroke
47	Missouri	Heart disease	Cancer	Stroke	CLRD
19	Montana	Heart disease	Cancer	CLRD	Unintentional injuries
25	Nebraska	Heart disease	Cancer	CLRD	Stroke
28	Nevada	Heart disease	Cancer	CLRD	Unintentional injuries
20	New Hampshire	Cancer	Heart disease	CLRD	Stroke
8	New Jersey	Heart disease	Cancer	Stroke	CLRD
28	New Mexico	Heart disease	Cancer	Unintentional injuries	CLRD
6	New York	Heart disease	Cancer	CLRD	Stroke
44	North Carolina	Heart disease	Cancer	Stroke	CLRD
10	North Dakota	Heart disease	Cancer	Stroke	Alzheimer's disease
30	Ohio	Heart disease	Cancer	CLRD	Stroke
48	Oklahoma	Heart disease	Cancer	CLRD	Stroke
13	Oregon	Cancer	Heart disease	Stroke	CLRD
34	Pennsylvania	Heart disease	Cancer	Stroke	CLRD
11	Rhode Island	Heart disease	Cancer	CLRD	Stroke
40	South Carolina	Heart disease	Cancer	Stroke	Unintentional injuries
3	South Dakota	Heart disease	Cancer	Stroke	CLRD
33	Tennessee	Heart disease	Cancer	Stroke	CLRD
25	Texas	Heart disease	Cancer	Stroke	Unintentional injuries
26	Utah	Heart disease	Cancer	Unintentional injuries	Stroke
0	Vermont	Cancer	Heart disease	CLRD	Unintentional injuries
41	Virginia	Heart disease	Cancer	Stroke	CLRD
4	Washington	Cancer	Heart disease	Stroke	CLRD
49	West Virginia	Heart disease	Cancer	CLRD	Unintentional injuries
43	Wisconsin	Heart disease	Cancer	Stroke	Unintentional injuries
14	Wyoming	Heart disease	Cancer	CLRD	Unintentional injuries

Analysis:

- It can be seen that the common pattern of 4 leading causes of individual states are mainly: Heart disease, Cancer, Stroke and CLRD

4. Are there year-by-year changes in the four leading causes of death nationwide?

In [9]:

```
us_deaths_causes_dfyf = df[(df['State'] == 'United States') & (df['Cause Name'] != 'All causes')] rows = []  # Iterate through the years to find the top four causes of death for each year for yr in us_deaths_causes_dfyf['Year'].unique():     state_cause_ydf = us_deaths_causes_dfyf[us_deaths_causes_dfyf['Year'] == yr].groupby(['Cause Name'])['Deaths'].sum()     state_cause_ydf.sort_values(ascending=False, inplace=True)      # Get the top 4 causes     cause_list = list(state_cause_ydf.keys())[:4] + [None] * (4 - len(state_cause_ydf))      # Append the row to the list     rows.append(['State': yr, '1st cause': cause_list[0], '2nd cause': cause_list[1],                 '3rd cause': cause_list[2], '4th cause': cause_list[3]])  df4 = pd.DataFrame(rows) df4.sort_values(by='Year', ascending=True, inplace=True) display(df4)
```

Out [9]:

Year	1st cause	2nd cause	3rd cause	4th cause	
0	1999	Heart disease	Cancer	Stroke	CLRD
0	2000	Heart disease	Cancer	Stroke	CLRD
2	2001	Heart disease	Cancer	Stroke	CLRD
7	2002	Heart disease	Cancer	Stroke	CLRD
3	2003	Heart disease	Cancer	Stroke	CLRD
6	2004	Heart disease	Cancer	Stroke	CLRD
4	2005	Heart disease	Cancer	Stroke	CLRD
5	2006	Heart disease	Cancer	Stroke	CLRD
8	2007	Heart disease	Cancer	Stroke	CLRD
9	2008	Heart disease	Cancer	CLRD	Stroke
10	2009	Heart disease	Cancer	CLRD	Stroke
11	2010	Heart disease	Cancer	CLRD	Stroke
12	2011	Heart disease	Cancer	CLRD	Stroke
14	2012	Heart disease	Cancer	CLRD	Stroke
13	2013	Heart disease	Cancer	CLRD	Unintentional injuries
15	2014	Heart disease	Cancer	CLRD	Unintentional injuries
16	2015	Heart disease	Cancer	CLRD	Unintentional injuries
17	2016	Heart disease	Cancer	Unintentional injuries	CLRD

Analysis:

- Through out the year (1999-2016), Heart disease and Cancer remain top 1 and 2, respectively, while 3rd and 4th causes are normally Stroke, CLRD and Unintentional injuries.

Population

Import and explore Population data

```
file2 = "nwt-est2018-01.xlsx"
df_pop = pd.read_excel(path=file2, skiprows=3)
df_pop.head()
```

Unnamed: 0	Census	Estimates Base	2010	2011	2012	2013	2014	2015	2016	2017	2018	
0	United States	308745538.0	308758105.0	309320985.0	311580009.0	313874218.0	316057727.0	318386421.0	320742673.0	323071342.0	325147121.0	327167434.0
1	Northeast	55317240.0	55318430.0	55308045.0	55600532.0	55776729.0	55907823.0	56015884.0	56047587.0	56058789.0	56072676.0	56111079.0
2	Midwest	66927001.0	66929743.0	66974749.0	67152631.0	67336937.0	67564135.0	67752380.0	67869139.0	67969917.0	68105035.0	68306740.0
3	South	114555744.0	114563045.0	114867066.0	116039399.0	117217705.0	118393244.0	119657737.0	121037542.0	122401186.0	123596420.0	124733948.0
4	West	71345553.0	71946887.0	72193625.0	72787447.0	73489477.0	74192525.0	74960682.0	75789405.0	76614450.0	77319986.0	77933663.0

Cleaning the data

```
# Rename column
df_pop.rename(columns={'Unnamed: 0': 'State'}, inplace=True)

# Cleaning
df_clean['State'] = df_pop['State'].str.replace('\n', '', regex=True)
```

Analysis:

- Through out the year (1999-2016), Heart disease and Cancer remain top 1 and 2, respectively, while 3rd and 4th causes are normally Stroke, CLRD and Unintentional injuries.

Population

Import and explore Population data

8	Arkansas	39218718.0	29404077.0	2261092.0	2959549.0	2967726.0	2978407.0	2990410.0	30029970.0	3013825.0	
9	California	373203093.0	376418233.0	379607822.0	382968264.0	386251390.0	389531390.0	392971017.0	396393454.0	3957054.0	
10	Colorado	59481281.0	61217711.0	61937121.0	62704822.0	63512128.0	64521072.0	65490111.0	66519922.0	6695864.0	
11	Connecticut	35791225.0	35880223.0	35943951.0	35949151.0	35947883.0	3587509.0	35786174.0	35738800.0	3572665.0	
12	Delaware	899595.0	9073116.0	9151180.0	923638.0	932599.0	941143.0	9492166.0	9570789.0	9671171.0	
13	District of Columbia	605085.0	6196202.0	634725.0	6504310.0	6625133.0	675294.0	686575.0	6966110.0	702455.0	
14	Florida	18845785.0	19093352.0	19322230.0	19563168.0	19860300.0	20222449.0	20629992.0	20958112.0	21299923.0	
15	Georgia	17118110.0	9801578.0	9901496.0	9973320.0	10096011.0	10181111.0	10304763.0	10413055.0	10519475.0	
16	Hawaii	1363963.0	1378252.0	1384905.0	1414862.0	1424284.0	1426816.0	1428015.0	1424203.0	1420401.0	
17	Idaho	1570773.0	1583826.0	1595441.0	1611530.0	1631479.0	1651523.0	1662999.0	1718904.0	1742408.0	
18	Illinois	12840762.0	12867791.0	12884119.0	12898269.0	12888962.0	12884342.0	12826895.0	12789156.0	12740180.0	
19	Indiana	6490436.0	6516045.0	6537640.0	6568376.0	6593533.0	6608296.0	6633344.0	6660062.0	6691878.0	
20	Iowa	3050767.0	3060504.0	3070997.0	3083740.0	3109504.0	3123470.0	3137765.0	3143637.0	3156145.0	

Cleaning the data

24

Louisiana

454532.0

4675184.0

4600814.0

4624577.0

4644204.0

4664851.0

4676215.0

4670818.0

4659978.0

25

Maine

1327632.0

1328150.0

1327991.0

1328196.0

1330760.0

1330964.0

1329484.0

1331370.0

1335063.0

1334609.0

26

Maryland

5788642.0

5838991.0

5887072.0

5923704.0

5958116.0

5985173.0

6004692.0

6024981.0

6042718.0

27

Massachusetts

6566413.0

6613149.0

6663158.0

6713944.0

6763952.0

6795891.0

6826922.0

6863249.0

6902149.0

28

Michigan

987735.0

988152.0

989639.0

9913349.0

9930599.0

9932573.0

99518909.0

9976477.0

9999515.0

29

Minnesota

5310843.0

5340668.0

5376501.0

5413691.0

5451222.0

5482903.0

5523490.0

5568156.0

5611179.0

30

Mississippi

2970536.0

2978470.0

2983767.0

2988797.0

2996023.0

2998930.0

2998296.0

2999653.0

2996530.0

31

Missouri

5995976.0

6099641.0

6040693.0

6040658.0

6056239.0

6077145.0

60977023.0

61098612.0

6126452.0

32

Montana

990722.0

997221.0

1003754.0

1012564.0

1021891.0

1029503.0

1040863.0

1053090.0

1062395.0

33

Nebraska

182958.0

184058.0

185332.0

186544.0

1879522.0

1891507.0

1905924.0

1917575.0

193286.0

34

Nevada

2702464.0

2712799.0

2744566.0

2776972.0

2819012.0

2839440.0

2917727.0

2972405.0

3034392.0

35

New Hampshire

1361777.0

1319615.0

1323962.0

1326408.0

1332232.0

1336294.0

1342713.0

1349677.0

1356468.0

36

New Jersey

979624.0

9827783.0

9845483.0

9856362.0

9866709.0

9870898.0

9874550.0

9889854.0

9890520.0

37

New Mexico

2064588.0

2080395.0

2087549.0

2092792.0

2090342.0

2090211.0

2092789.0

2093358.0

2095428.0

38

New York

19400080.0

19489514.0

19574549.0

19628043.0

1966330.0

19664111.0

19641589.0

19652490.0

19642209.0

39

North Carolina

9547293.0

9656754.0

9749123.0

9843999.0

9933944.0

10033079.0

10156679.0

10270800.0

10336320.0

40

North Dakota

764273.0

765136.0

761116.0

721999.0

73782.0

754022.0

754393.0

755117.0

76077.0

41

Ohio

11539327.0

11543463.0

11548399.0

11576578.0

11620973.0

11617850.0

11653003.0

11664129.0

11698440.0

42

Oklahoma

3759632.0

378721.0

3818600.0

3832505.0

387837.0

3909831.0

3926769.0

3932640.0

3939679.0

43

Oregon

387352.0

3871728.0

3899118.0

3922908.0

3964106.0

4016918.0

4094044.0

4149942.0

4190713.0

44

Pennsylvania

1271158.0

12744583.0

12766827.0

12776612.0

1278911.0

12785579.0

12783588.0

12790447.0

12807680.0

45

Rhode Island

1063938.0

1053536.0

1046401.0

1055122.0

105601.0

1056713.0

1057063.0

1064466.0

1057315.0

46

South Carolina

4369566.0

4671422.0

471712.0

464153.0

482739.0

4892253.0

4958256.0

5021219.0

5084127.0

47

South Dakota

816105.0

823484.0

833496.0

842270.0

849088.0

853931.0

862890.0

873876.0

882235.0

48

Tennessee

6355301.0

6397104.0

6451281.0

6493432.0

6540826.0

6590089.0

6645011.0

6709047.0

6770014.0

49

Texas

25426729.0

2564227.0

26089920.0

2648944.0

26977142.0

27486814.0

27973422.0

28322771.0

28711815.0

50

Utah

2757334.0

28142216.0

2853467.0

2898727.0

2937399.0

2982497.0

3042613.0

31030134.0

3106465.0

51

Vermont

625990.0

626979.0

626963.0

626212.0

625218.0

625197.0

623644.0

642555.0

626299.0

52

Virginia

8032980.0

8104699.0

8185299.0

8253053.0

8312076.0

8379007.0

8410964.0

8456077.0

8517665.0

53

Washington

6742362.0

6821655.0

6892876.0

6962906.0

7056249.0

7163452.0

7284000.0

7425432.0

752591.0

54

West Virginia

1854214.0

1856074.0

1856764.0

1858378.0

1849467.0

1849467.0

18290929.0

1817046.0

1809832.0

55

Wisconsin

5690479.0

5704755.0

5719855.0

5736952.0

5751914.0

5761466.0

5772988.0

5792951.0

5813568.0

56

Wyoming

564483.0

567224.0

576270.0

582123.0

582548.0

585688.0

5842900.0

578934.0

577737.0

Veit the DataFrame to create a year column

```
df_melted = df_pop1.melt(id_vars=["State",
                             var_name="Year",
                             value_name="Population"])

# Convert the Year column to Integer type
df_melted["Year"] = df_melted["Year"].astype(int)
df_melted.head()
```