

Data Visualization Final Project

Module/Part 1:

My data set is from the U.S. Department of Health and Human Services; it can be found here: <https://catalog.data.gov/dataset/conditions-contributing-to-deaths-involving-coronavirus-disease-2019-covid-19-by-age-group>. The data set was published by the Centers For Disease Control and is titled “Conditions Contributing to COVID-19 Deaths, by State and Age, Provisional 2020-2023.” This dataset consists of all the COVID-19 deaths recorded in the United States and the health conditions (comorbidities) contributing to those deaths. The data is further grouped by age group and each state in the U.S. The data was collected between 2020 and 2023, and no new data points are being added to this data set. I want to pursue healthcare data science, and I was interested in looking at a data set related to public health. During the COVID-19 pandemic there was a lot of discussion around COVID comorbidities, which are health conditions that can increase an individual’s risk of having a severe COVID infection or increases an individual’s risk of dying from COVID. The data set includes various medical conditions mentioned in conjunction with COVID-19 deaths, these include influenza/pneumonia, chronic lower respiratory diseases, adult respiratory distress syndrome, respiratory failure, respiratory arrest, hypertensive diseases, ischemic heart disease, cardiac arrest, cardiac arrhythmia, heart failure, cerebrovascular diseases, sepsis, diabetes, obesity, Alzheimer’s disease, dementia, and renal failure. The attributes included in the data set are the U.S. states, the condition group, condition, ICD-10 code of the condition, age group, and the number of COVID-19 deaths associated with that state and age group; there are also totals for each age group and for the entire U.S. I will look at the total COVID-19 deaths in the United States, broken down by age group for 10 of the deadliest health conditions. The conditions I will look at are hypertensive diseases, diabetes, renal failure, ischemic heart disease, sepsis, chronic lower respiratory diseases, heart failure, cerebrovascular diseases, obesity, and Alzheimer’s disease. The age groups are broken down as follows: 0-24, 25-34, 35-44, 45-54, 55-64, 65-74, 75-84, and 85+.

The key questions I would want to answer include: What are the total COVID-19 deaths associated with each health condition? Which health conditions contributed to the highest number of COVID-19 deaths? Which health conditions most or least impacted each age group and are there differences in how each age group was affected? After doing some research on previous visualizations for COVID-19 death data I noticed that most graphs were either a variation of a bar chart or a scatter plot. Multiple bar graphs were often used to depict COVID-19 deaths for different age groups, but these usually only consisted of 2-3 different age groups. Having 8 bars per category (health condition) may be too many bars to efficiently look at the entire visual and may overwhelm the user. Many visuals also used random or rainbow colors or a few mutated colors. Since I have 8 categories for each age group, I would like to have 8 colors that are easy to distinguish from each other without implementing all the colors of the rainbow. Using hue for distinction would be appropriate since I’m using it for categorical encoding. One option would be to use a gradient for the different age groups but trying this the 8 colors were not distinctive enough from each other

to be effective in seeing each age group. I also saw several line graphs, but it was difficult to distinguish the colors of each line since lines can be very thin. Several of the bar graphs I saw had very narrow or short bars which made it difficult to compare each of the bars to each other.

Module/Part 2:

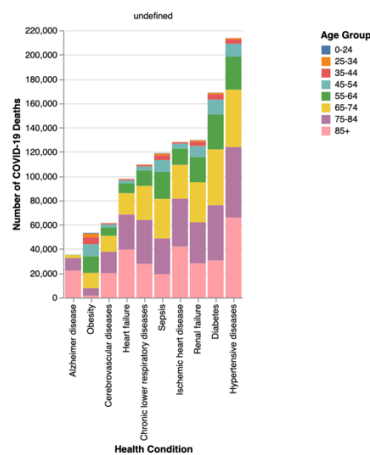
In module 2 I identified two tasks I would like to complete with my visual.

For the first task, the goal is: present/communicate information, the means is: navigate through the data, the characteristics are: identify patterns across each bar (how is the distribution of covid deaths by age group different for each health condition), the target data is: a relative reference frame, the workflow: user will look at the entire graph, distinguish each age group based on color using the legend, then compare the number of deaths there are for each age group within each bar, then compare the bars to each other to see how each condition compares, and the roles: the user has the goal of understanding how different age groups were impacted by different health conditions in contributing to COVID deaths. This gives a better understanding on the most common covid comorbidities for each age group. They also can determine that the majority of covid deaths occurred in people aged 45 and older.

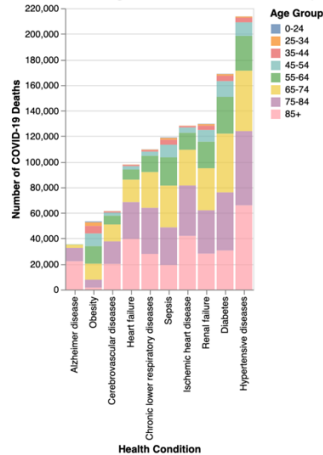
For the second task, the goal: present/communicate information, the means: build connections in the data, characteristics: identify the pattern in the distribution of covid-19 deaths (which conditions lead to the most covid-19 deaths), target data: relative reference frame, workflow: user will look at the entire graph, then use the x-axis to determine the total number of COVID-19 deaths that were associated with each comorbidity, and roles: The user will determine what the different covid comorbidities are, how many deaths were associated with each comorbidity, and determine which comorbidity was the most and least high risk for covid death.

These are some of my low fidelity prototypes that I created with Altair initially (I attached the code for the prototypes and final visual as well at the end):

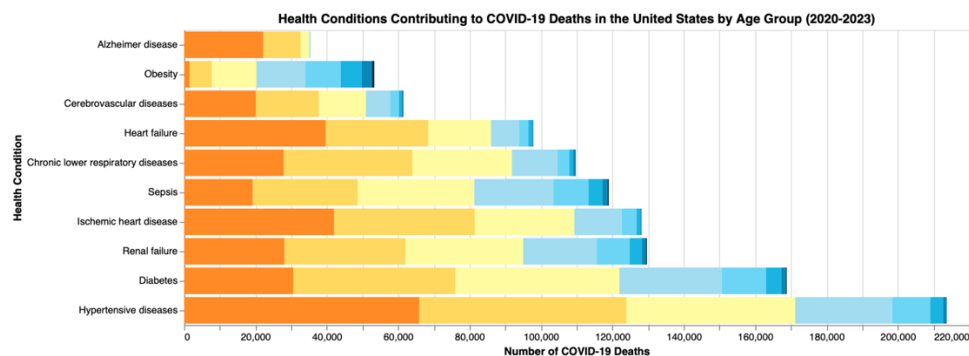
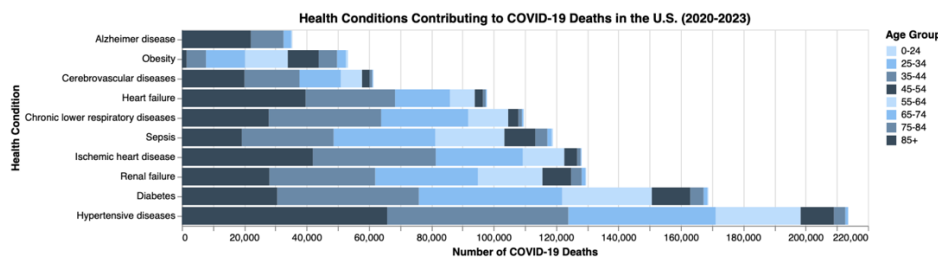
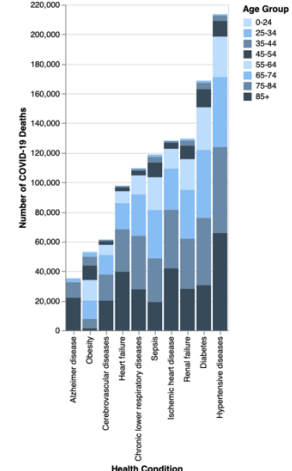
Health Conditions Contributing to COVID-19 Deaths in 2020-2023



Health Conditions Contributing to COVID-19 Deaths in the U.S. (2020-2023)



Health Conditions Contributing to COVID-19 Deaths in the U.S. (2020-2023)



Module/Part 3:

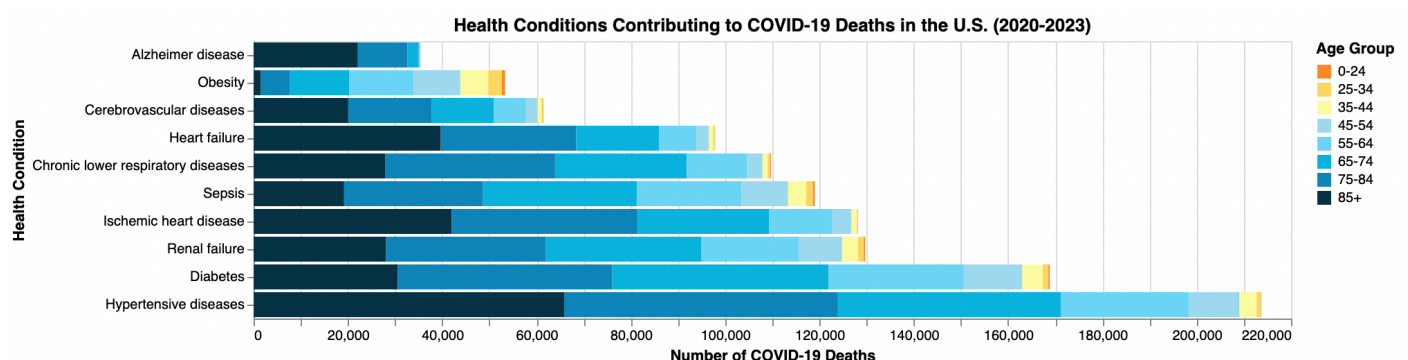
My preliminary design evaluation, would be an insight-based evaluation, and is as follows:

Target question I want to answer: Can people see the patterns they need to (and correctly) from my visual?

People recruited to answer that question: My family, colleagues/friends, and two acquaintances who are in doctoral programs who frequently review research and their corresponding graphs.

The kind of measures used and what the measures tell me about the core question: Determine if the people I've recruited can complete my two tasks as outlined in my Module 2 reflection. The measure would include insight depth and accuracy. I would want to know what insights my audience finds and if those are accurate with what is depicted in the graph.

The approach to answer that question: I would give each of my participants a survey and ask the following question: What types of insights can you discover or gain in this



Summary of the key elements of design:

I believe that the horizontal stacked bar graph was the best way to visualize the COVID-19 deaths for each health condition broken down by age group. This kept the graph as simple and easy to read as possible while still displaying all the data. I increased the length and width of each bar to make it easier to identify each age group within each bar. I kept the bars close to each other to make them easier to compare side by side. I made the stacked bar chart horizontal to make the names of each condition easy to read as it would've been difficult to read long names if they were vertical. Instead of a multiple bar chart I chose to use a stacked bar graph because a multiple bar graph would have shown a total of 80 bars, which could've been visually overwhelming. I ordered the bars (health conditions) from top to bottom in ascending order of the number of deaths to make it easier to see how the total number of deaths differed between each condition. I would have liked if all the colors used were a gradient of a hue to show the oldest age group as darkest and the youngest age group as lightest. However, given the number of categories (8 total age groups) it was too difficult to distinguish 8 shades of blue for example, especially since the younger age groups were the smallest. Given that each age group was categorical I made each color distinct and choose the colors for each group based on what would show up best for that category based on its size. Since there were a high number of categories (8) I chose colors that were distinct from each other but made several adjacent colors a gradient so that I could more easily be found on the legend. I used darker colors of blue for the older age groups and shades of orange/yellow for the younger age groups. This highlights that most COVID-19 deaths occurred in people aged 45 and older (as distinguished by the shades of blue). In addition, given the context of the topic, I tried to keep the colors engaging without being too vibrant or a rainbow of colors.

Discussion of my final evaluation approach:

My final evaluation consisted of administering my survey to four different people. I shared my survey via a google doc and my participants wrote in their responses. My four participants included two acquaintances who frequently read/review research, one friend, and one family member.

My survey asked the following questions: What types of insights can you discover or gain in this visual? Are the colors in the graph easy to distinguish and are they easy to pair with the legend? What are the total death counts associated with each health condition (comorbidity)? What types of trends do you notice in the data? Do you have any suggestions for improvement of this visual?

My evaluation results: After administering my surveys and receiving the responses I noticed a few things. One suggestion I had was to create a rainbow of colors to make each age group color distinct. Another suggestion was to only have two different age groups that split the data in half into a 0–54-year age group and a 55+ age group so that only two colors would be necessary. Another suggestion was to make a multiple bar graph with less age groups. Another suggestion I had was to make all 8 age group colors a gradient. Many of my suggestions seemed to be about the colors and the division of age groups. I decided

against most of these suggestions because I didn't think rainbow colors were meaningful for the data, I kept all 8 age groups because it highlighted important differences between the age groups, and opted to have a variation of a gradient where possible. All four of the participants were able to determine the total number of COVID-19 deaths per age group, although one participant initially thought that the stacked bars were cumulative. The trends that my participants noticed include: There were a higher portions of COVID-19 deaths for younger age groups associated with obesity and that obesity was the least impacting for the oldest age group (85+). The very youngest age group had the highest number of deaths associated with obesity and the oldest age group had the highest number of deaths associated with hypertensive disease. Overall, Alzheimer disease was associated with the least number of deaths and hypertensive disease was associated with the highest number of deaths. Most deaths occurred in people aged 45 and older.

A synthesis of your findings:

I think that my initial approach was to create low fidelity prototypes in Altair instead of by hand. This allowed me to very quickly try new things, like change the type of graph, color schemes, scales, axes, etc. I went through many changes before arriving at my final visual. I think that I would do my prototyping this way in the future after an initial draft sketch by hand to quickly eliminate the approaches that don't work well.

I think that conducting an insight-based evaluation using a survey was the most appropriate evaluation for this type of visual. I received a lot of feedback this way and was able to understand well what my participants were gaining from my visual. I think my biggest struggle was in choosing color, in the future I would aim to decrease the number of categories to 5-6 or less if it makes sense for the data. This would allow for more selective and intentional color choices. I would have like to have a gradient in my colors, but it would have been difficult for my participants to see every category with the stacked bar chart. I would also consider if using proportion instead of counts is appropriate for the data.

Overall, going through the steps of outlining my goals and tasks was very helpful, this kept my intentions for the visual aligned with its purpose. In the future I would implement the five design sheets process for a more complex visual and would hope to implement more interactive features like reconfiguration to show the age groups in a different order for example for easier comparison, or create an interaction that allowed users to see the proportion of deaths for each condition and age group.

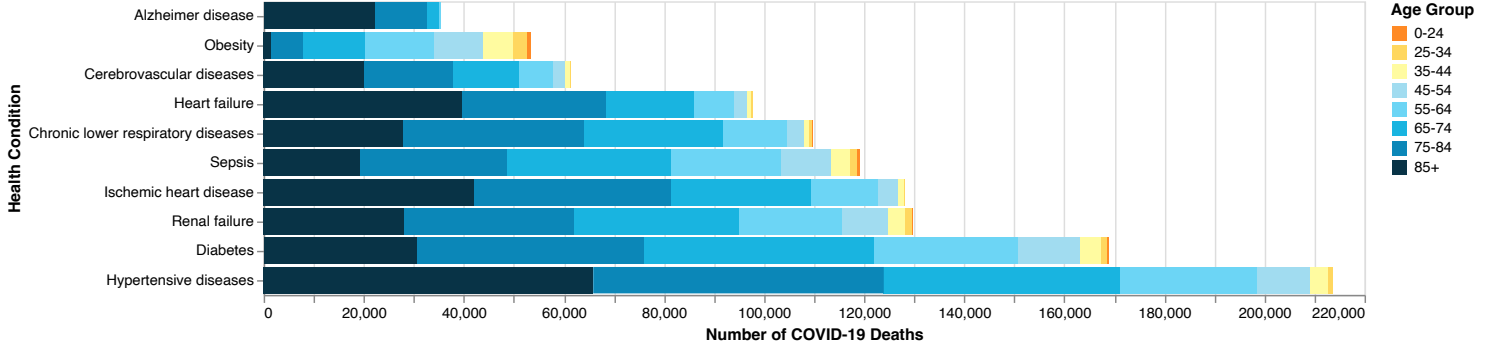
Health Condition

Age Group

- 0-24
- 25-34
- 35-44
- 45-54
- 55-64
- 65-74
- 75-84
- 85+

Number of COVID-19 Deaths

Health Condition	0-24	25-34	35-44	45-54	55-64	65-74	75-84	85+
Alzheimer disease	0	0	0	0	0	0	35,000	22,000
Obesity	0	0	0	0	0	0	10,000	8,000
Cerebrovascular diseases	0	0	0	0	0	0	45,000	20,000
Heart failure	0	0	0	0	0	0	95,000	10,000
Chronic lower respiratory diseases	0	0	0	0	0	0	110,000	10,000
Sepsis	0	0	0	0	0	0	82,000	18,000
Ischemic heart disease	0	0	0	0	0	0	82,000	20,000
Renal failure	0	0	0	0	0	0	95,000	10,000
Diabetes	0	0	0	0	0	0	75,000	90,000
Hypertensive diseases	0	0	0	0	0	0	170,000	20,000



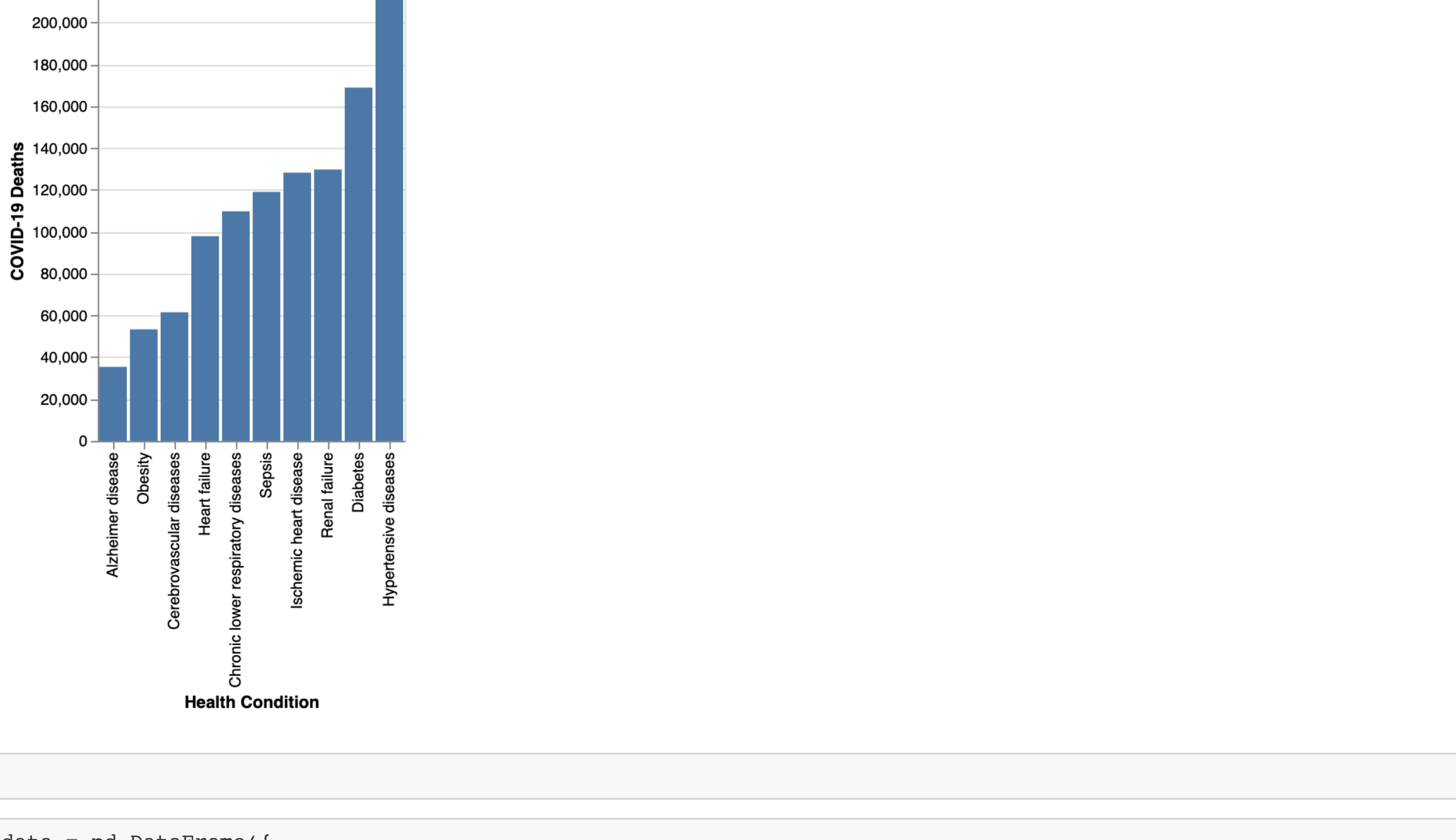

```
In [1]: import altair as alt
import pandas as pd
```

```
In [2]: data = pd.DataFrame({
    'Health Condition': ['Hypertensive diseases', 'Diabetes', 'Renal failure', 'Ischemic heart disease', 'Sepsis', 'Chronic lower respiratory diseases', 'Heart failure', 'Cerebrovascular diseases', 'Obesity', 'Alzheimer disease'],
    'COVID-19 Deaths': [213709, 168830, 129674, 128178, 119000, 109697, 97821, 61455, 53321, 35377]
})

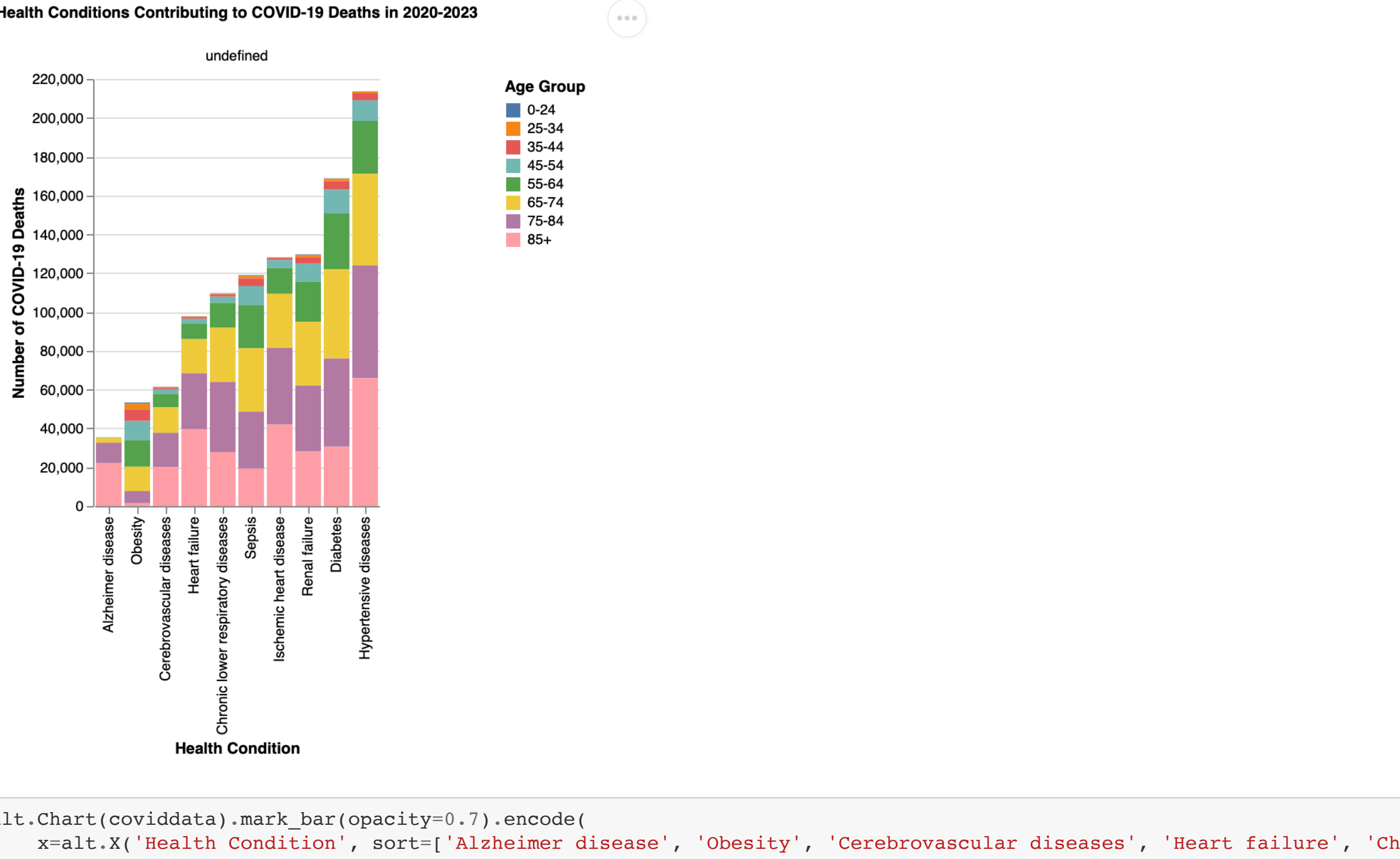
alt.Chart(data).mark_bar().encode(
    x=alt.X('Health Condition', sort=['Alzheimer disease', 'Obesity', 'Cerebrovascular diseases', 'Heart failure', 'Chronic lower respiratory diseases', 'Sepsis', 'Ischemic heart disease', 'Renal failure', 'Diabetes', 'Hypertensive diseases']),
    y='COVID-19 Deaths',
    properties(title="Health Conditions Contributing to COVID-19 Deaths in 2020-2023")
)
```

Out[2]: Health Conditions Contributing to COVID-19 Deaths in 2020-2023

Age group	Number of people
15-24	210,000
25+	180,000

[illegible][illegible]

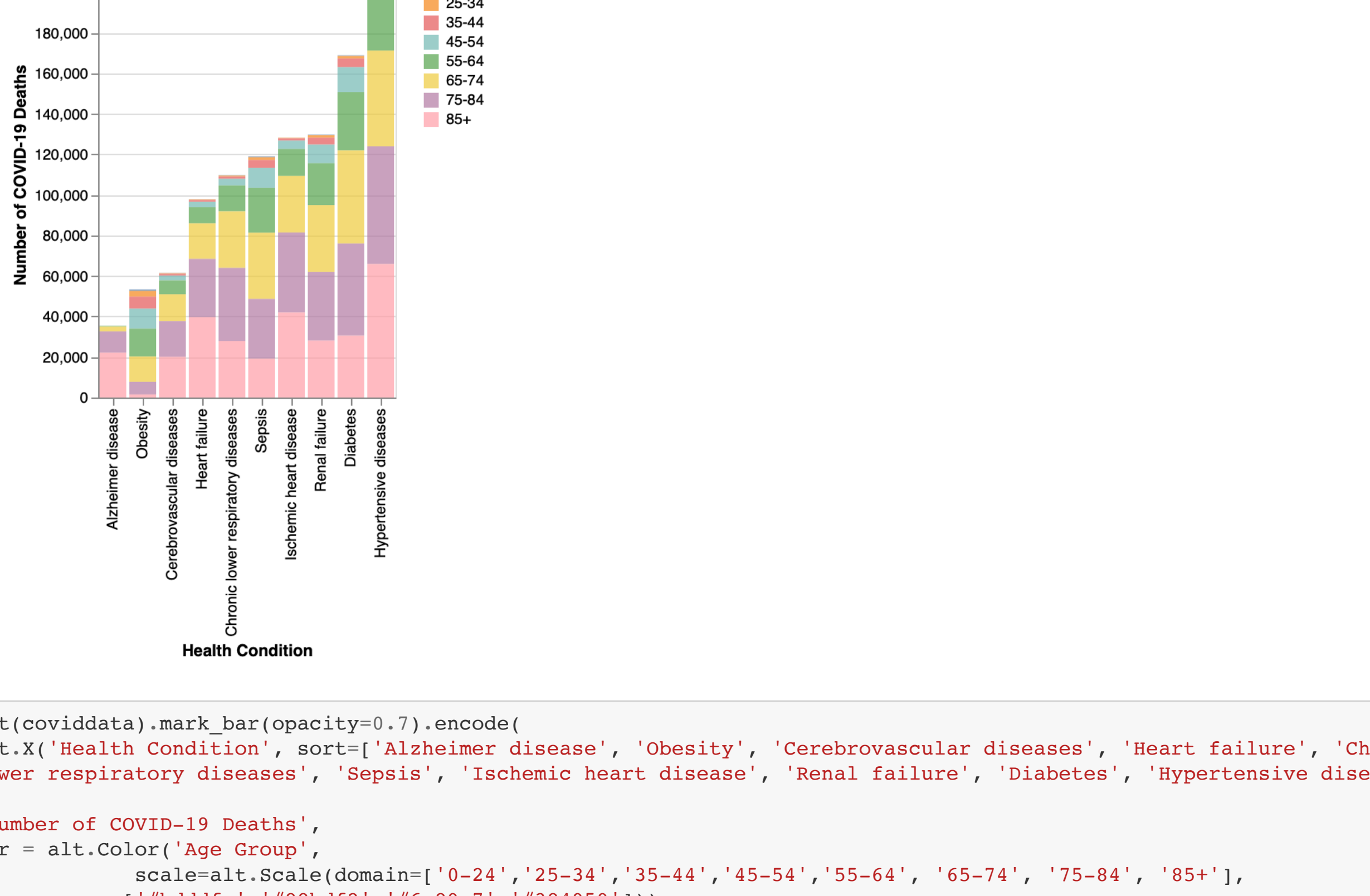
```
alt.Chart(coviddata).mark_bar().encode(
    x=alt.X('Health Condition', sort=['Alzheimer disease', 'Obesity', 'Cerebrovascular diseases', 'Heart failure', 'Chronic lower respiratory diseases', 'Sepsis', 'Ischemic heart disease', 'Renal failure', 'Diabetes', 'Hypertensive diseases']),
    y='Number of COVID-19 Deaths',
    color='Age Group',
    tooltip=[ 'primary_type:N',
              'year:O',
              alt.Tooltip('sum(Number_of_Incidents):Q',
                          title='Number of incidents')
            ]
).facet(
    column=alt.Column('year:O',
                      header=alt.Header(title="Health Conditions Contributing to COVID-19 Deaths in 2020-2023")),
).resolve_scale(
    x='independent'
).configure_view(
    stroke='transparent'
```

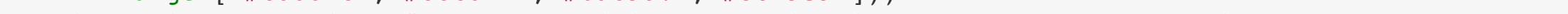


```
y='Number of COVID-19 Deaths',
color='Age Group')
).properties(title="Health Conditions Contributing to COVID-19 Deaths in the U.S. (2020-2023)")
```

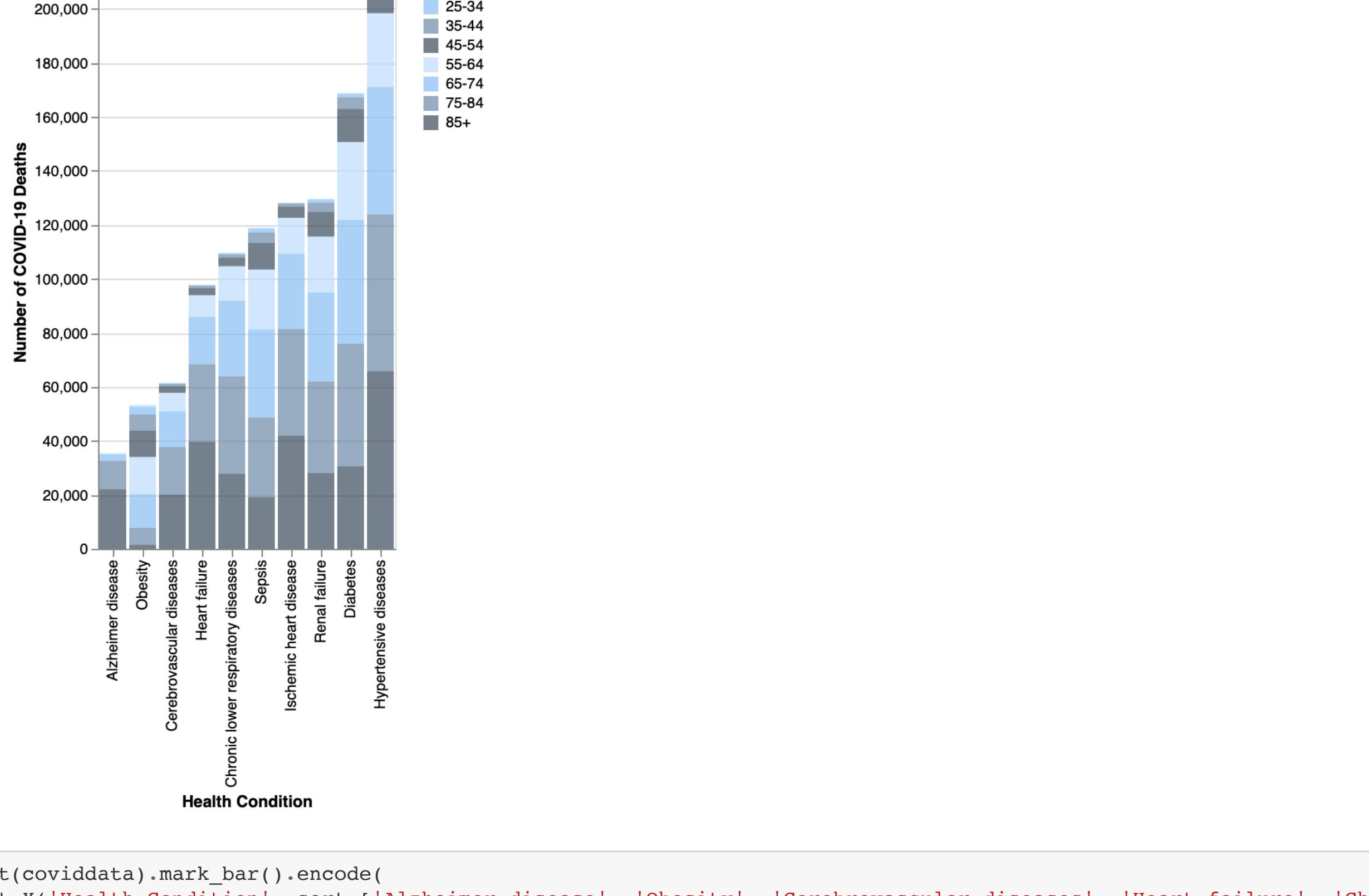
Out[7]: Health Conditions Contributing to COVID-19 Deaths in the U.S. (2020-2023)

Age Group	Number of People
0-24	~180,000
25-64	~120,000
65+	~80,000



Out[16]: 

Age Group	Percentage of Respondents
0-24	~10%
25-34	~15%
35-44	~20%
45-54	~25%
55-64	~30%
65-74	~35%
75+	~40%

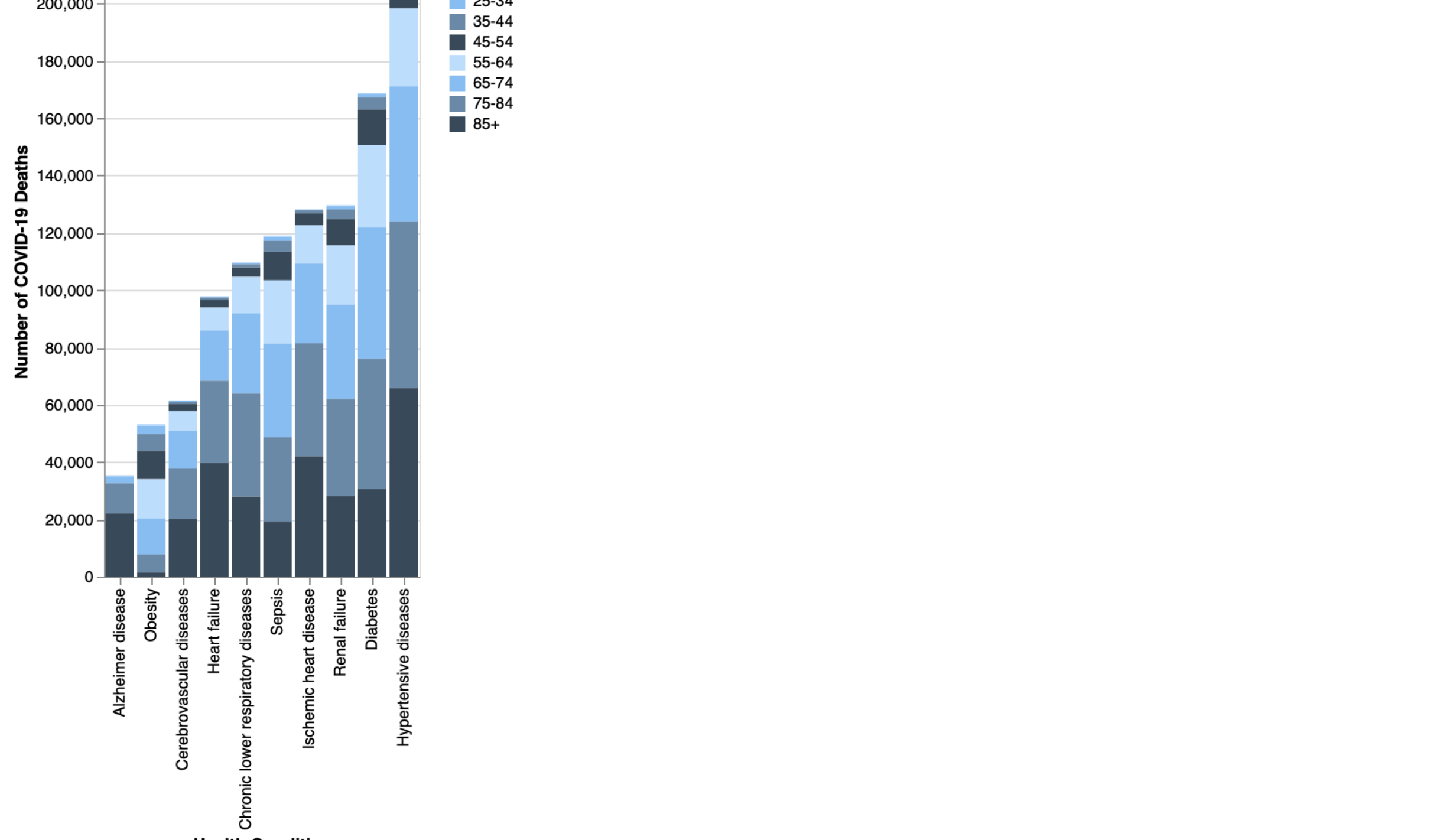


```

        ronic lower respiratory diseases', 'Sepsis', 'Ischemic heart disease', 'Renal failure', 'Diabetes', 'Hypertensive dise
        ases'})),
        y='Number of COVID-19 Deaths',
        color=alt.Color('Age Group'),
        scale=alt.Scale(domain=['0-24', '25-34', '35-44', '45-54', '55-64', '65-74', '75-84', '85+'],
        range=['#b8dfe2', '#8dbdf2', '#6a89a7', '#384959'])),
        .properties(height=400, title="Health Conditions Contributing to COVID-19 Deaths in the U.S. (2020-2023)")
    Out[17]: Health Conditions Contributing to COVID-19 Deaths in the U.S. (2020-2023)

```

0.24
0.25-0.34

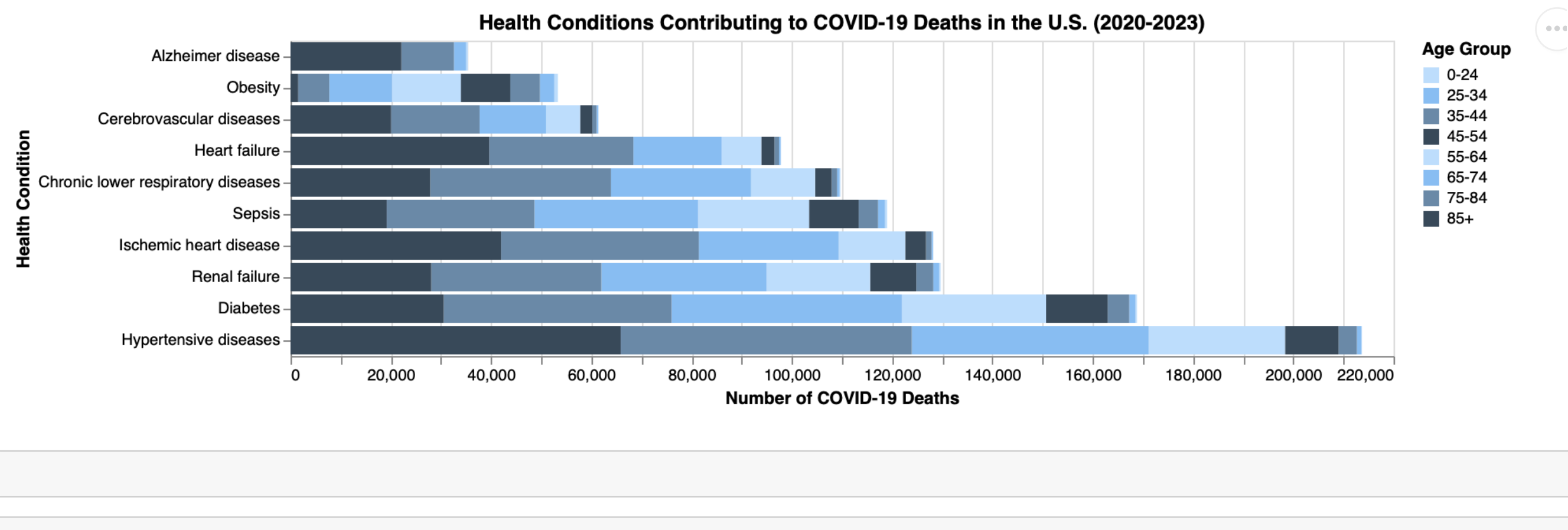


```
In [18]: alt.Chart(coviddata).mark_bar().encode(
    y=alt.Y('Health Condition', sort=['Alzheimer disease', 'Obesity', 'Cerebrovascular diseases', 'Heart failure', 'Chronic lower respiratory diseases', 'Sepsis', 'Ischemic heart disease', 'Renal failure', 'Diabetes', 'Hypertensive diseases']),
    x='Number of COVID-19 Deaths',
    color=alt.Color('Age Group',
        scale=alt.Scale(domain=['0-24', '25-34', '35-44', '45-54', '55-64', '65-74', '75-84', '85+'],
            range=['#b8dfdc', '#88bdf2', '#6a89a7', '#384959'])))
```

```

)properties(height: 200,width: 100,color: red,stroke: black)

```



[13]:

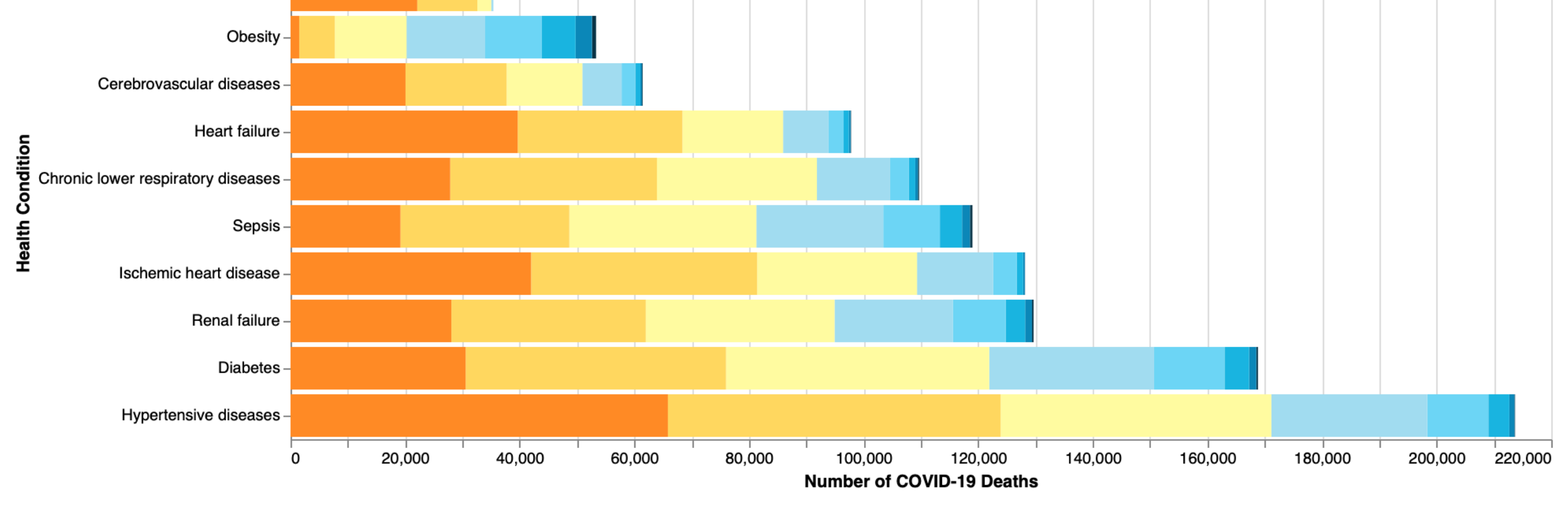
```

        ronic lower respiratory diseases", 'Sepsis', 'Ischemic heart disease', 'Renal failure', 'Diabetes', 'Hypertensive disease')],
        x= 'Number of COVID-19 Deaths',
        color= alt.Color('Age group'),
        scale=alt.Scale(domain=['0-24', '25-34', '35-44', '45-54', '55-64', '65-74', '75-84', '85+'],
            range=['#083B46', '#0B87B8', '#19B4E0', '#6CD5F5', '#A1D9F0', '#FFB400', '#FF7F50', '#FF8A25'])),
        ).properties(height=300,width=800,title="Health Conditions Contributing to COVID-19 Deaths in the United States by Age Group (2020-2023)")
    
```

Out[13]:

Health Conditions Contributing to COVID-19 Deaths in the United States by Age Group (2020-2023)

Alzheimer disease –

[illegible]

`'emic heart disease', 'Renal failure', 'Renal failure', 'Renal failure', 'Renal failure', 'Renal failure', 'Renal failure', 'Renal failure', 'Renal failure', 'Diabetes', 'Diabetes', 'Diabetes', 'Diabetes', 'Diabetes', 'Diabetes', 'Diabetes', 'Diabetes', 'Diabetes', 'Diabetes', 'Hypertensive diseases', 'Hypertensive diseases', 'Hypertensive diseases', 'Hypertensive diseases', 'Hypertensive diseases', 'Hypertensive diseases', 'Hypertensive diseases', 'Hypertensive diseases', 'Hypertensive diseases', 'Hypertensive diseases']`

`#Number of COVID-19 Deaths': [0, 2, 3, 36, 321, 2425, 10476, 22114, 707, 2865, 5910, 9878, 13715, 12516, 6208, 152]`

[illegible][illegible]

```
    alt.Chart(coviddata).mark_bar().encode(
      y=alt.Y('Health Condition', sort=['Alzheimer disease', 'Obesity', 'Cerebrovascular diseases', 'Heart failure', 'Ch
```

```

        'nonfatal lower respiratory diseases', 'sepsis', 'ischemic heart disease', 'renal failure', 'diabetes', 'hypertensive dise
ases']]),
        x = 'Number of COVID-19 Deaths',
        color = alt.Color('Age Group'),
        facet = alt.Facet('domain', ['0-24', '25-34', '35-44', '45-54', '55-64', '65-74', '75-84', '85+'])

```

```

    range='#FF8A25',#FFD75F',#ffffba',#a1dcf0',#6cd5f5',#19b4e0',#08b788',#033466'))
    ).properties(height=200,width=750,title="Health Conditions Contributing to COVID-19 Deaths in the U.S. (2020-2023)")

```

Out[20]:

Health Conditions Contributing to COVID-19 Deaths in the U.S. (2020-2023)

Disease	0-24	25-34	35-44	45-54	55-64
Alzheimer disease	0.00	0.00	0.00	0.00	0.00
Obesity	0.00	0.15	0.10	0.05	0.00
Cerebrovascular diseases	0.00	0.00	0.05	0.10	0.85
Heart failure	0.00	0.00	0.00	0.10	0.90

Health Condition	75-84 (%)	85+ (%)
Chronic lower respiratory diseases	~75	~25
Sepsis	~65	~35
Ischemic heart disease	~60	~40
Renal failure	~55	~45

Age Group	No comorbidity	Diabetes	Hypertensive diseases
18-24	~100	~100	~100
25-34	~100	~100	~100
35-44	~100	~100	~100
45-54	~100	~100	~100
55-64	~100	~100	~100
65-74	~100	~100	~100
75-84	~100	~100	~100
85+	~100	~100	~100