

SI649-21-Fall Lab 4 -> Altair Transformation

Overview

We will work with a dataset from the article "[Joining The Avengers Is As Deadly As Jumping Off A Four-Story Building](https://fivethirtyeight.com/features/avengers-death-comics-age-of-ultron/)" (<https://fivethirtyeight.com/features/avengers-death-comics-age-of-ultron/>).

We'll focus on **transformation**. In this lab, you will practice:

1. transform_calculation
2. transform_aggregate
3. transform_join_aggregate
4. transform_window
5. transform_filter
6. exporting to HTML

Data

The source data is available at this [link](https://github.com/fivethirtyeight/data/tree/master/avengers) (<https://github.com/fivethirtyeight/data/tree/master/avengers>). You can also find a description of each column on the website. We preprocessed the dataset for you, and you do not need to write any data processing code.

For this lab, please write Altair code to answer the questions. In many situations, you could also solve the problem using Pandas. However, we want code that can be deployed without using pandas, so it's better to practice to just do as much as we can in Altair. You can complete the entire lab without writing any pandas transformation.

It's fine if your visualization looks slightly different from the example (e.g., getting 1.1 instead of 1.0, use orange instead of red)

Lab Instructions (read the full version on the handout of the previous lab)

- Save, rename, and submit the ipynb file (use your username in the name).
- Run every cell (do Runtime -> Restart and run all to make sure you have a clean working version), print to pdf, submit the pdf file.
- For each visualization, we will ask you to write down a "Grammar of Graphics" plan first (basically a description of what you'll code).
- If you end up stuck, show us your work by including links (URLs) that you have searched for. You'll get partial credit for showing your work in progress.

```
In [78]: # imports we will use
import altair as alt
import pandas as pd
```

```

In [79]: #obtain the data and process the data
# you do not need to edit this block, just run it
# and use avenger_df as your starting point
avenger_link="https://raw.githubusercontent.com/fivethirtyeight/data/master
avenger_return=pd.read_csv(avenger_link,encoding="latin-1")
articleLink="https://fivethirtyeight.com/features/avengers-death-comics-age
#process the dataset
convertColumnList=[]
modified=[]
for i in range(5):
    convertColumnList.append("Death"+str(i+1))
    convertColumnList.append("Return"+str(i+1))

for c in convertColumnList:
    avenger_return[c]=avenger_return[c].astype(str)

for j in range(len(avenger_return)):
    modified.append([])
    storage=modified[-1]
    for c in convertColumnList:
        item=avenger_return[c][j]
        v = 1 if item=="YES" else 0
        storage.append(v)

organizedDR=pd.DataFrame(modified,columns=convertColumnList)
shortlist=avenger_return.drop(columns=convertColumnList)
avenger_df=shortlist.join(organizedDR)

```

```

In [80]: # you can set to whatever theme you prefer.
alt.themes.enable('fivethirtyeight')
avenger_df.sample(2)

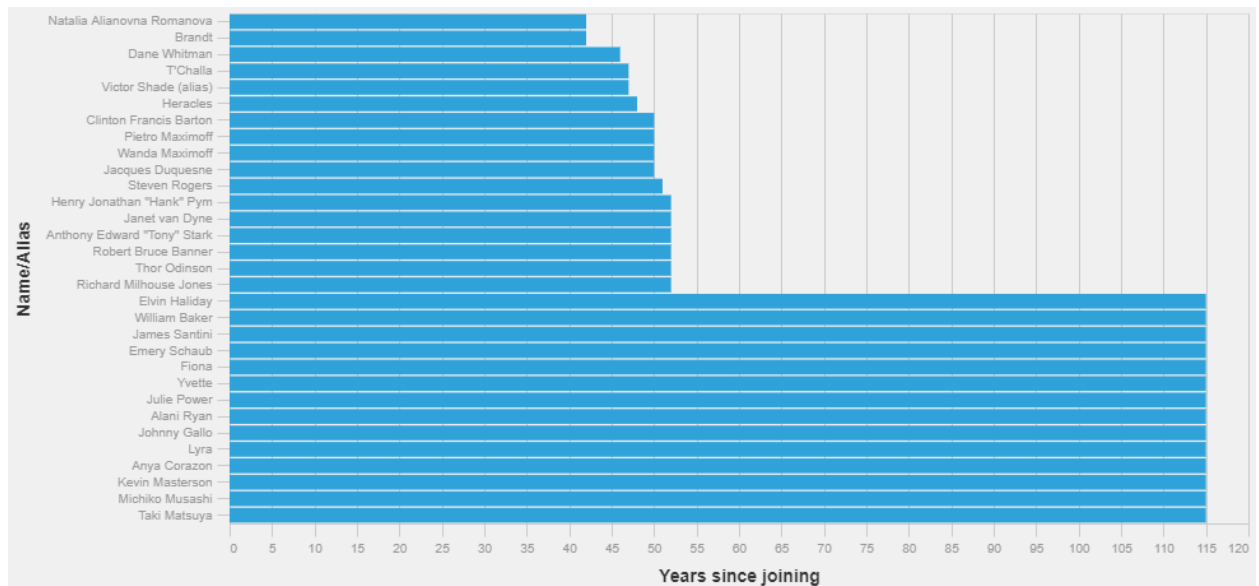
```

Out[80]:

	URL	Name/Alias	Appearances	Current?	Gender	Pi
154	http://marvel.wikia.com/Abyss_(Ex_Nihilo%27s)_...	NaN	25	YES	FEMALE	
130	http://marvel.wikia.com/Ken_Mack_(Earth-616)#	Ken Mack	59	NO	MALE	

2 rows x 21 columns

Visualization 1: Top 30 avengers who have been in the series for the longest time



Description of the visualization:

We want to find avengers who have been in the series for the longest time. We'd like to build a bar chart with the top 30 avengers ranked by their years since join.

- The number of years that they have been in the series (column "Years since joined") is plotted on the X axis
- The avenger's name is displayed on the y axis
- Only the top 30 avengers are included in the chart

Step 1: Write down your plan for the visualization (edit this cell)

- Data Name: `avenger_df`
- mark type: `TODO: mark_bar()`
- Encoding Specification:

- `x:TODO: Years since joining, Q`

- `y: TODO: Name/Alias, N`

- Transformation Plan:

- `step 1:TODO: I will produce the rank of 'Year since joining by using transform_window)`

- `step 2:TODO: I will filter the top 30 avengers by using transfor_filter`

Hint

- You need to rank avengers by the number of years they have joined. Review our lab demo and the altair documentation. Which transformation can generate the rank?
- Because we only want to include the top 30 avengers, which transformation can we use?
- Adding tooltips can be a helpful way to check if your calculation is correct.

Step 2: Create your chart.

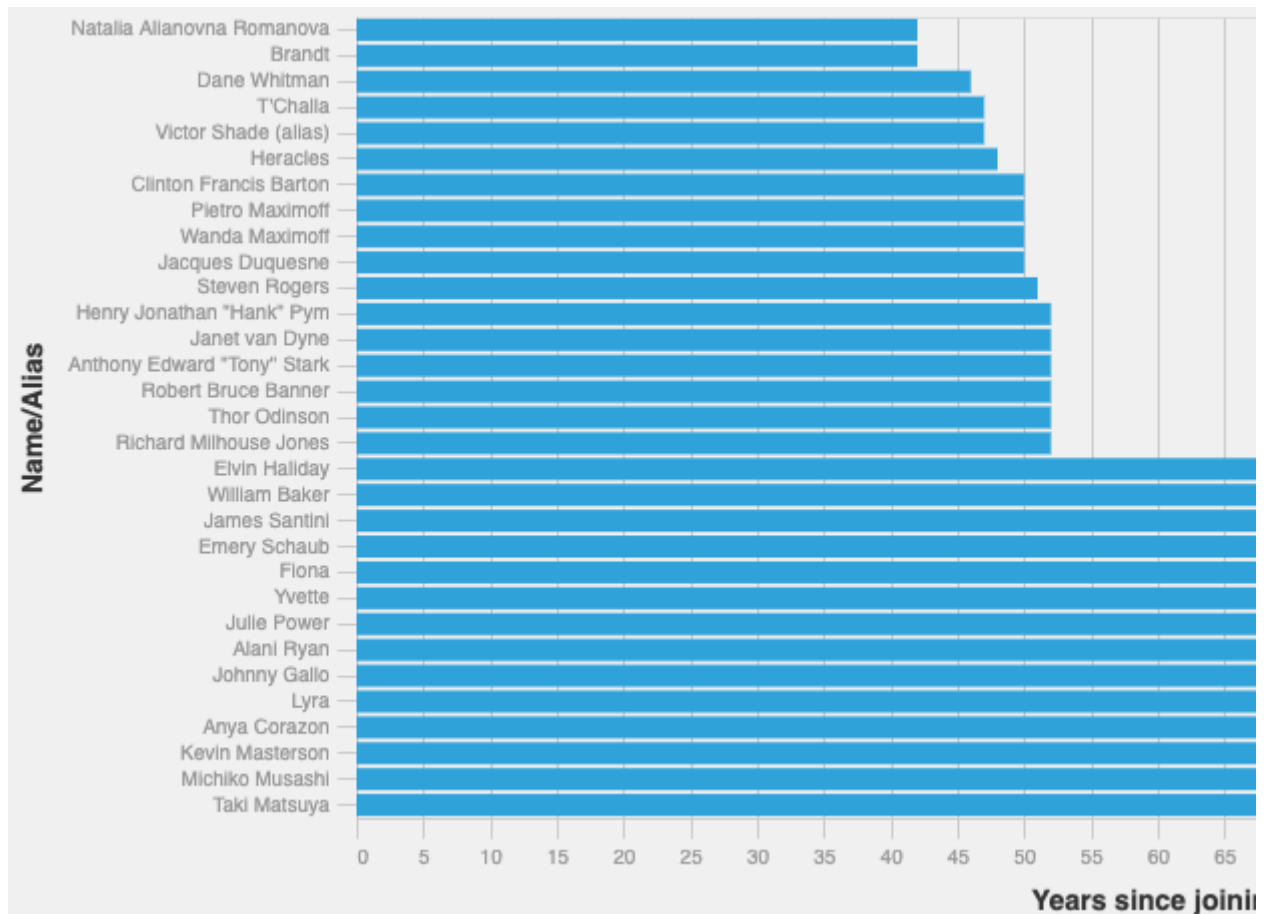
Please paste your FINAL answer to the cell immediately below this block (it will allow us to grade).
You can search for the keyword "TODO" to locate cells that need your edits

```
In [81]: # TODO: (vis1) top 30 avengers who have been in the series for the longest

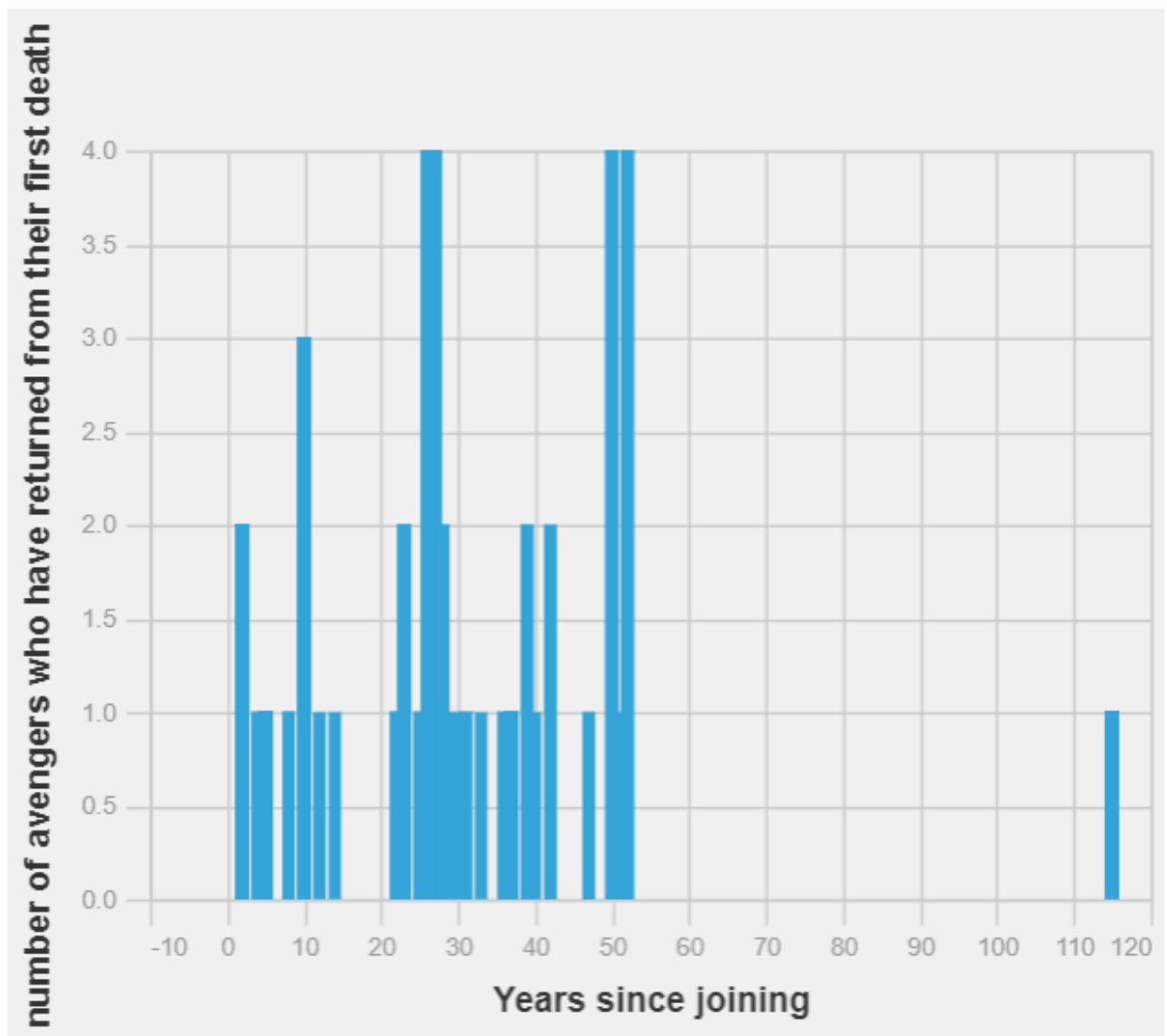
chart1 = alt.Chart(avenger_df, width = 800, height = 400).transform_window(
    sort = [alt.SortField('Years since joining', order = 'descending')],
    y_s_j = 'rank(*)'
).transform_filter(
    alt.datum.y_s_j < 31
).mark_bar(size = 11).encode(
    alt.X('Years since joining:Q', title = 'Years since joining'),
    alt.Y('Name/Alias:N', sort = alt.EncodingSortField(
        field = 'y_s_j', order = 'descending')),
    tooltip = 'Years since joining:Q'
)

chart1
```

Out[81]:



Visualization 2: Count of avengers who return from their first death



Description of the data: We want to count the number of avengers who have returned from their first death. We have processed the data so that for every Death/Return column:

- 0 indicate that they did not die or did not return
- 1 indicate that they died or returned
 - e.g., if an avenger died at least once and did not return, their data looks like this: (Death1=1, Return1=0)
 - e.g., if an avenger did not die twice (Death2==0), then they could not have returned twice (Return2==0)

Description of the visualization:

- a bar chart
- y-axis plot the total number of avengers who have returned within an "age group"
- x-axis plot the "age group": number of years since they have joined.
- y-axis uses the label: *number of avengers who have returned from their first death*

Step 1: Write down your plan for the visualization (edit this cell)

- Data Name: avenger_df

- mark type: *TODO:mar_bar()*
- Encoding Specification:

- *x:TODO: Years since joining, Q*

- *y: TODO: sum of Return1, Q*

- Transformation Plan:

- *step 1:TODO:I will produce the number of avenger that returned from their first death by using transform_aggregate to group them in ages and count the numbers that returned*

Hint

1. We want to plot a count of avengers who have returned from their first death. For these avengers, their Return1 value would be equal to 1. Do you need to look at other columns (such as Death1)?
2. Do you want to use joinaggregate or aggregate?
3. How can you produce a calculated/aggregated value within a group?

Step 2: Create your chart.

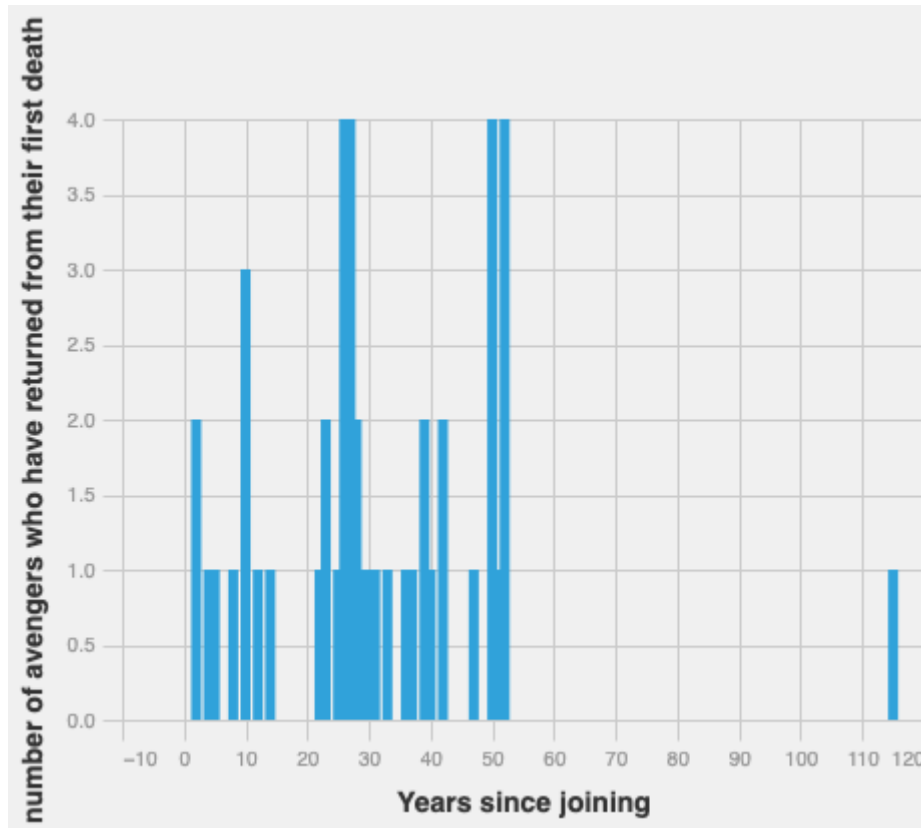
Please paste your FINAL answer to the cell immediately below this block (it will allow us to grade). You can search for the keyword "TODO" to locate cells that need your edits

```
In [82]: #TODO: vis2
# number of avengers who returned from the 1st death, breakdown by years si

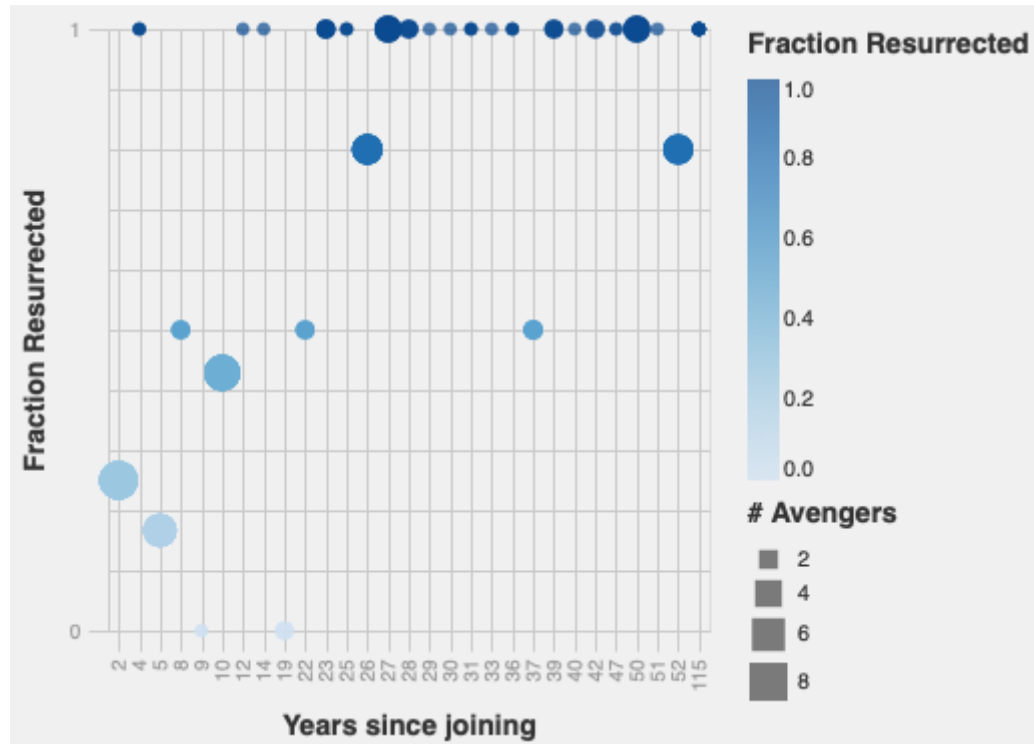
chart2 = alt.Chart(avenger_df).transform_aggregate(
    groupby = ['Years since joining'],
    return1_count = 'sum(Return1)'
).mark_bar().encode(
    alt.X('Years since joining:Q'),
    alt.Y('return1_count:Q', title = 'number of avengers who have returned
')
)

chart2
```

Out[82]:



Visualization 3: Chance of resurrection increases with years in service



The previous chart displays the number of time that avengers have returned from their first death, broken down by their years since joining. We also want to see how likely it is for them to be resurrected. Therefore, we are going to make a scatter plot with the fraction of avengers resurrected, the years since they've joined, and how many avengers died.

Description of the visualization:

- a scatterplot
- x-axis plot the number of years that they have joined
- y-axis plot the fraction of avengers that were resurrected, which is double encoded with color
 - y-axis uses the label: *Fraction Resurrected*
- circle size encodes the number of avengers who died at least once.
- (optional) Tooltips that display the encoded data

Hint:

- Just like the previous chart, you want to produce a count of avengers who have returned after their first death. Do you also need a count of avengers who have died at least one time?
- Should you use aggregate or join aggregate?
- How can you calculate the fraction of avengers that were resurrected? What is the numerator of the fraction and what is the denominator?

Step 1: Write down your plan for the visualization (edit this cell)

- Data Name: avenger_df
- mark type: *TODO:mark_circle()*
- Encoding Specification:

- *x:TODO: Years since joining, Q*

- *y: TODO: sum(Return1)/sum(Death1), Q*

- Transformation Plan:

- *step 1:TODO: I will produce sum(Return1) and sum(Death1) by using transform_aggregate*

- *step 2:TODO: I will produce fraction resurrected by using transform_calculate*

Step 2: Create your chart.

Please paste your FINAL answer to the cell immediately below this block (it will allow us to grade).
You can search for the keyword "TODO" to locate cells that need your edits

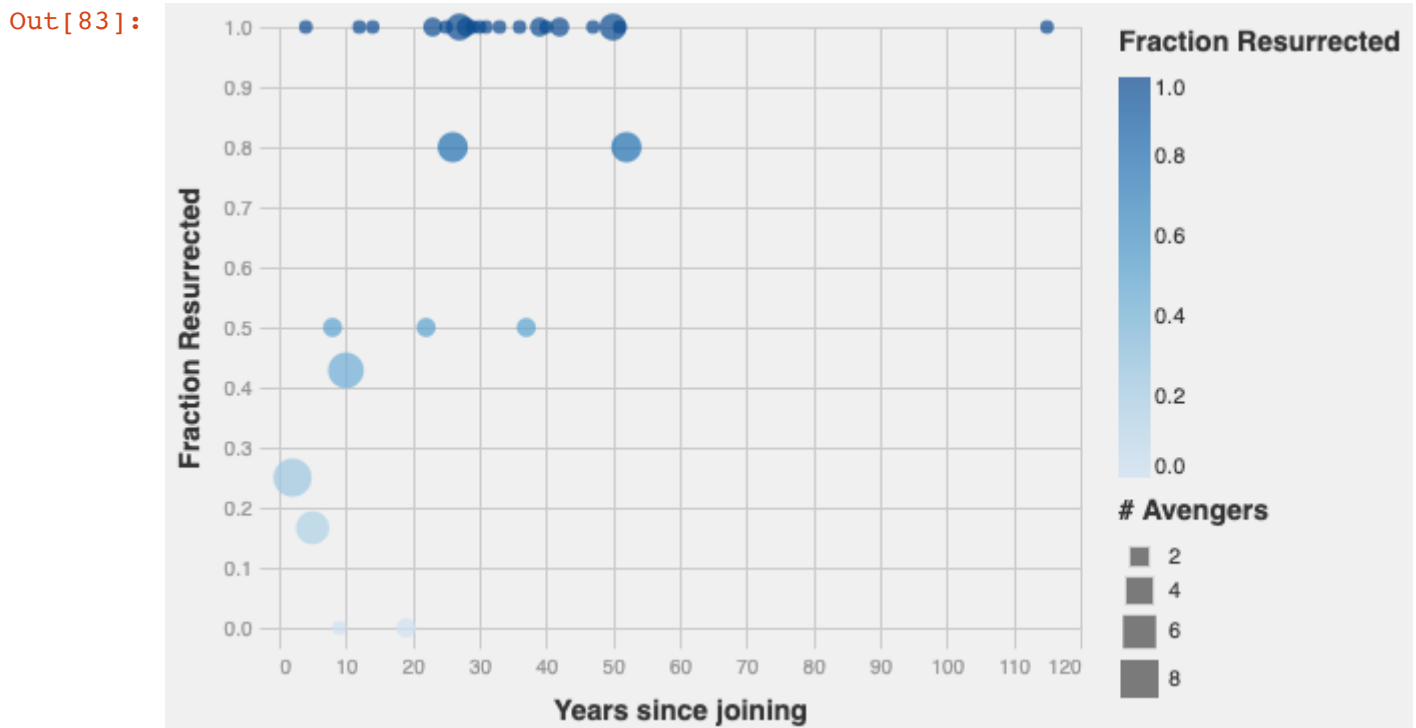
```

In [83]: #TODO vis 3
#chance of resurrection and years since joining

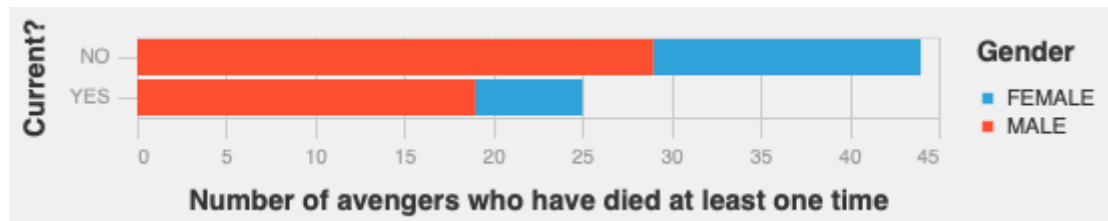
chart3 = alt.Chart(avenger_df).transform_aggregate(
    groupby = ['Years since joining'],
    return1_count = 'sum(Return1)',
    death1_count = 'sum(Death1)'
).transform_calculate(
    fraction_resurrected = alt.datum.return1_count / alt.datum.death1_count
).mark_circle().encode(
    alt.X('Years since joining:Q'),
    alt.Y('fraction_resurrected:Q', title = 'Fraction Resurrected'),
    alt.Color('fraction_resurrected:Q', title = 'Fraction Resurrected'),
    alt.Size('death1_count:Q', title = '# Avengers'),
    tooltip = ['Years since joining:Q', 'fraction_resurrected:Q', 'death1_c
)

chart3

```



Vis 4: Number of avengers who have died at least once, divided by current membership and gender



Instead of looking at years since joined, we want to investigate in the relationship between count of avengers who have died at least once and their other properties such as whether they are a current member or not, and gender.

Description of the visualization:

- a bar chart
- x-axis displays the total number of avengers who have **died at least one time** within an "age group"
 - x-axis uses the label: *number of avengers who have died at least one time*
- y-axis displays whether they are a current member or not
- use color to display the gender information

Hint:

- Do you want to use `joinaggregate` or `aggregate`?

Step 1: Write down your plan for the visualization (edit this cell)

- Data Name: `avenger_df`
- mark type: `TODO:mar_bar()`
- Encoding Specification:

- `x:TODO: sum(Death1), Q`

- `y: TODO: Current?, N`

- Transformation Plan:

- `step 1:TODO: I will produce the number of avengers who have died at least once by using transform_aggregate`

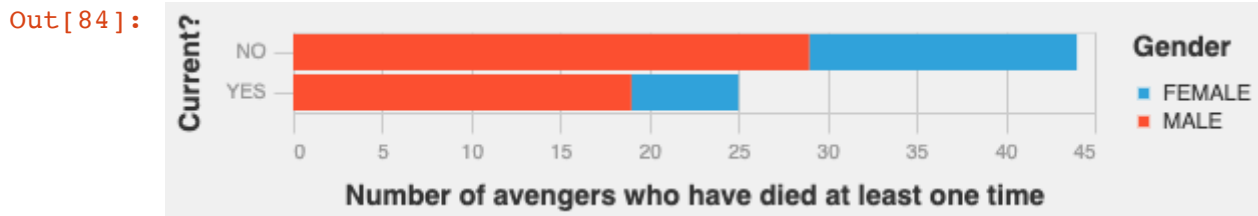
Step 2: Create your chart.

Please paste your FINAL answer to the cell immediately below this block (it will allow us to grade). You can search for the keyword "TODO" to locate cells that need your edits

```
In [84]: #TODO Vis 4
         #number of avenger who has died at least once,
         #divided by current membership and gender

chart4 = alt.Chart(avenger_df).transform_aggregate(
    groupby = ['Current?', 'Gender'],
    death1_count = 'sum(Death1)'
).mark_bar().encode(
    alt.X('death1_count:Q', title = 'Number of avengers who have died at le
    alt.Y('Current?:N'),
    alt.Color('Gender:N', title = 'Gender')
)

chart4
```



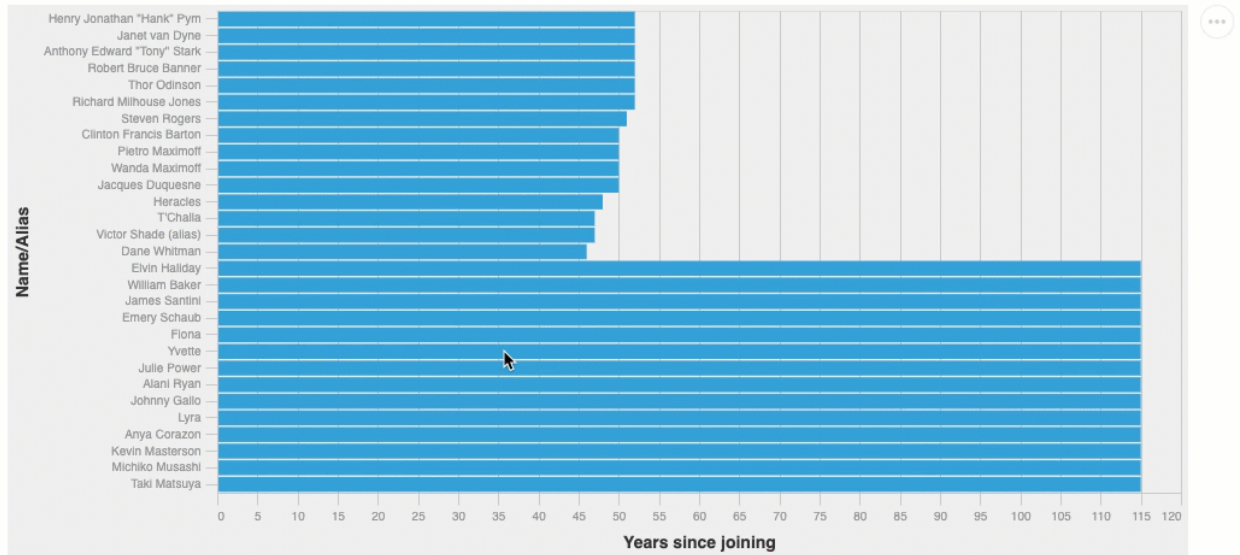
Submission

Once you are done with all the altair-related coding, we want you to embed these charts in an HTML file:

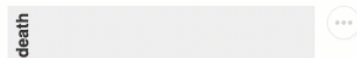
1. export these chart specifications to JSON
2. embed them into an HTML file.
3. If you use multiple files (e.g., HTML + JSON + javascript instead of pasting JSON into the HTML file), make sure you zip the entire folder.

The final website will look similar to this:

Vis 1



Vis 2



We also want to have a copy of your python script. Please run all cells (Runtime->Run all), and

1. save to PDF (File->Print->Save PDF)
2. save to ipynb

Rename all files with your username: e.g.,

1. username.pdf
2. username.ipynb
3. username.zip (html)

Upload all files to canvas. *End of Lab 4 Homework*

```
In [85]: chart1.save('chart1.json')
```

```
In [86]: chart2.save('chart2.json')
```

```
In [87]: chart3.save('chart3.json')
```

```
In [88]: chart4.save('chart4.json')
```