

Education Analysis

Work by: Nicholas Konovalenko
Team 3

Current Work

For this project I focused on looking at the educational attainment for Tracts 64, 102, 105, and 110. I found that within tracts 64 & 105, there is a much smaller percentage of people that achieves to top levels of education compared to tracts 102 & 110. I further break this down by demographic (Race & Gender) which show some stark differences. We saw the the percentage graduation rate for White-Alone Males and Females was significantly higher than Black-Alone Males and Females. In fact, in Census Tract 64 the high school graduation rate for Black-Alone Females was barely above 50%. Additionally, I looked at median earnings. For 4-year degrees, the median earnings across the tracts was relatively even. However, for the category 'Some college or associate's degree', the median earnings in CT 102/110 were 3 times higher than CT 64/105 ($80,000 vs 30,000$ per year). Finally, I took a look at access to computers, and found no visible correlation between computer access and education level attained.

Challenges Faced

The current biggest challenge faced, was the lack of all the tracts in all the years. This made it more difficult to do year-to-year analysis.

Next Steps

I think that the differences between schools must be analyzed much more closely. Why is the graduation rate for Census Tract 64 so much lower than in Census Tract 102? Why are the median earnings in CT 64 so much lower for the same level of education?

How to run

All the data is loaded to the 'Education Analysis/data/' directory, and all the plots get saved to the 'Education Analysis/Plots' directory. To run this, use either a jupyter notebook, or google colab and run all cells.

Importing the data

In [1]:

```
1 import numpy as np
2 import glob
3 import pandas as pd
4 import matplotlib.pyplot as plt
5
6 PATH = 'data/'
7
8 files = [
9     'ACSSST5Y2010.S1501.csv',
10    'ACSSST5Y2011.S1501.csv',
11    'ACSSST5Y2012.S1501.csv',
12    'ACSSST5Y2013.S1501.csv',
13    'ACSSST5Y2014.S1501.csv',
14    'ACSSST5Y2015.S1501.csv',
15    'ACSSST5Y2016.S1501.csv',
16    'ACSSST5Y2017.S1501.csv',
17    'ACSSST5Y2018.S1501.csv',
18    'ACSSST5Y2019.S1501.csv',
19    'ACSSST5Y2020.S1501.csv',
20    'ACSSST5Y2021.S1501.csv'
21 ]
22
23 df = pd.DataFrame()
24
25 ## Load files into df
26
27 # this currently only reads the 2021 CSV, as tracts are missing from the others.
28 for file in files[-1:]:
29     fpath = PATH + file
30
31     # Load each file into a temporary DataFrame
32     temp_df = pd.read_csv(fpath)
33
34     # Append the temporary DataFrame to the main DataFrame
35     df = df.append(temp_df, ignore_index=True)
36
```

C:\Users\konic\AppData\Local\Temp\ipykernel_15036\2308872588.py:35: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.

```
df = df.append(temp_df, ignore_index=True)
```

In [2]:

```
1 # Rename the columns using the .str.replace() method
2 df.columns = df.columns.str.replace('Census Tract ([\d\.]+), District of Columbia, District of Columbia!!([^\!]+)!!([^\!]+)', r'Tract \1, \2 \3'
```

C:\Users\konic\AppData\Local\Temp\ipykernel_15036\3899270774.py:3: FutureWarning: The default value of regex will change from True to False in a future version.

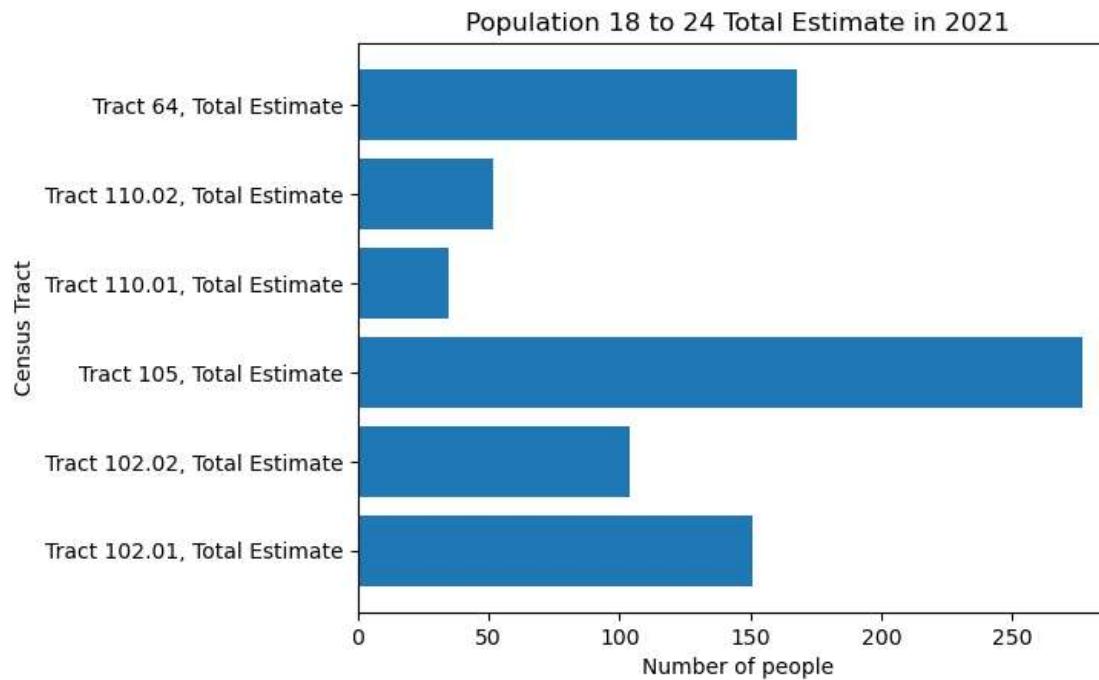
```
df.columns = df.columns.str.replace('Census Tract ([\d\.]+), District of Columbia, District of Columbia!!([^\!]+)!!([^\!]+)', r'Tract \1, \2 \3')
```

Population 18 to 24 years old

Let's start by looking at the overall number of student estimates.

In [3]:

```
1 # Extract the relevant data from the DataFrame
2 tract_cols = [col for col in df.columns if 'Tract' in col and 'Total' in col and 'Estimate' in col]
3 tract_data = df.loc[1, tract_cols]
4 tract_data = tract_data.sort_index()
5
6 # Create a figure and axis object
7 fig, ax = plt.subplots()
8
9 # Plot the data as a horizontal bar chart
10 ax.barh(tract_data.index, tract_data.values.astype(int))
11
12 # Set the axis labels and title
13 ax.set_xlabel('Number of people')
14 ax.set_ylabel('Census Tract')
15 ax.set_title('Population 18 to 24 Total Estimate in 2021')
16
17 # Display the plot
18 plt.show()
```



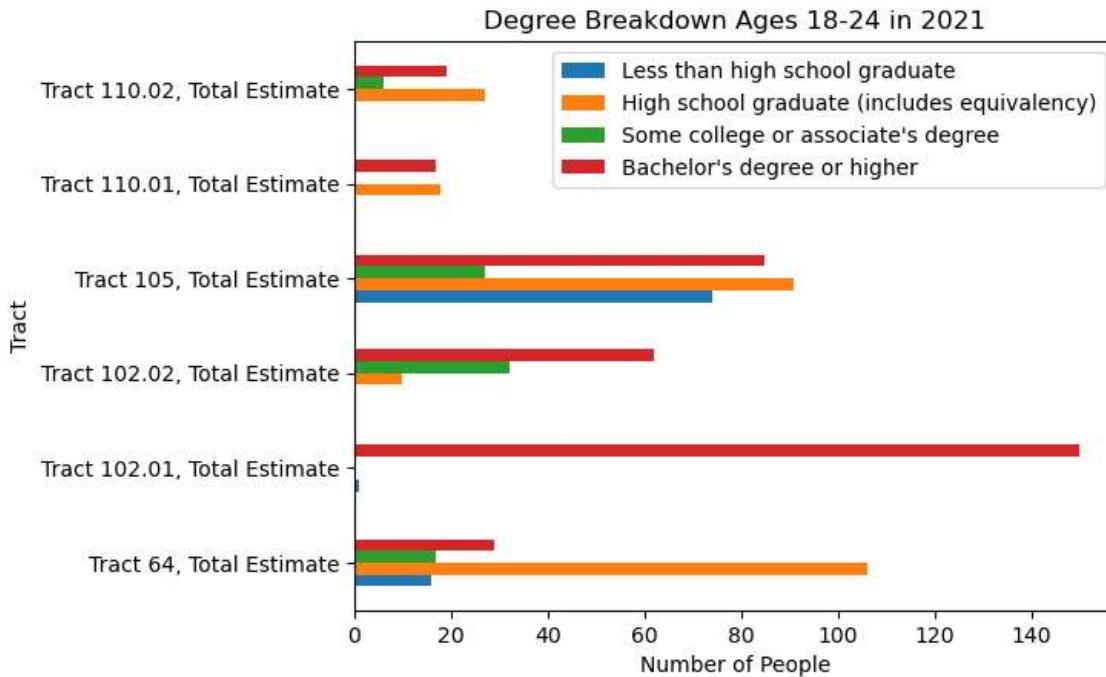
```
In [8]: # This is the overall breakdown, let's try breaking down by degree
# Extract the relevant data from the DataFrame
tract_cols = [col for col in df.columns if 'Tract' in col and 'Total' in col and 'Estimate' in col]
tract_data = df.loc[2:5, tract_cols]
tract_data = tract_data.sort_index()

#####
# Rename the rows
labels = ['Less than high school graduate', 'High school graduate (includes equivalency)',
          'Some college or associate\'s degree', 'Bachelor\'s degree or higher']
for i in range(len(labels)):
    tract_data = tract_data.rename(index={2+i: labels[i]})

#####
# Convert the columns to int
for col in tract_data.columns:
    tract_data[col] = tract_data[col].astype(int)
tract_data = tract_data.T

bar_plot = tract_data.plot(kind='barh', title='Degree Breakdown Ages 18-24 in 2021')
bar_plot.set_xlabel('Number of People')
bar_plot.set_ylabel('Tract')
plt.savefig(f'Plots/tract_18_24_breakdown.png')
bar_plot
```

Out[8]: <AxesSubplot:title={'center':'Degree Breakdown Ages 18-24 in 2021'}, xlabel='Number of People', ylabel='Tract'>



What are the takeaways from this graph?

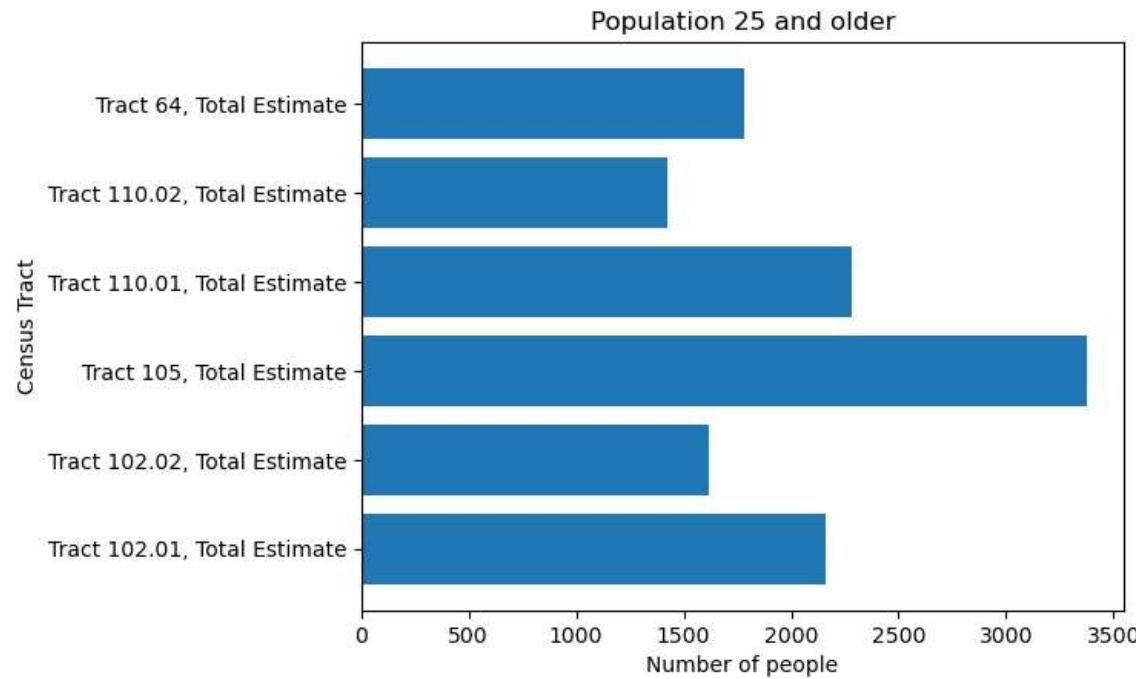
We can see that Tract 102 has the highest number of Bachelor's degree or higher graduates. Infact, Tract 102.01 is composed of essentially only 4-year degree graduates (for those between the ages of 18-24). Tract 110 is also heavily made up of college graduates. On the other hand, a large chunk of Tract 64 contains high school graduates, and Tract 105 has double the number of High school graduate + High school dropout compared to the number of 4-year graduates it has.

This shows that 18 to 24 year olds in Tracts 102 and 110 continue to get more advanced education, while those in Tracts 64 and 105 struggle to do the same.

Population 25+

In [10]:

```
1 # Extract the relevant data from the DataFrame
2 tract_cols = [col for col in df.columns if 'Tract' in col and 'Total' in col and 'Estimate' in col]
3 tract_data = df.loc[6, tract_cols]
4 tract_data = tract_data.str.replace(',', '').astype(int)
5 tract_data = tract_data.sort_index()
6
7 # Create a figure and axis object
8 fig, ax = plt.subplots()
9
10 # Plot the data as a horizontal bar chart
11 ax.bah(tract_data.index, tract_data.values.astype(int))
12
13 # Set the axis labels and title
14 ax.set_xlabel('Number of people')
15 ax.set_ylabel('Census Tract')
16 ax.set_title('Population 25 and older')
17
18 # Display the plot
19 plt.show()
```



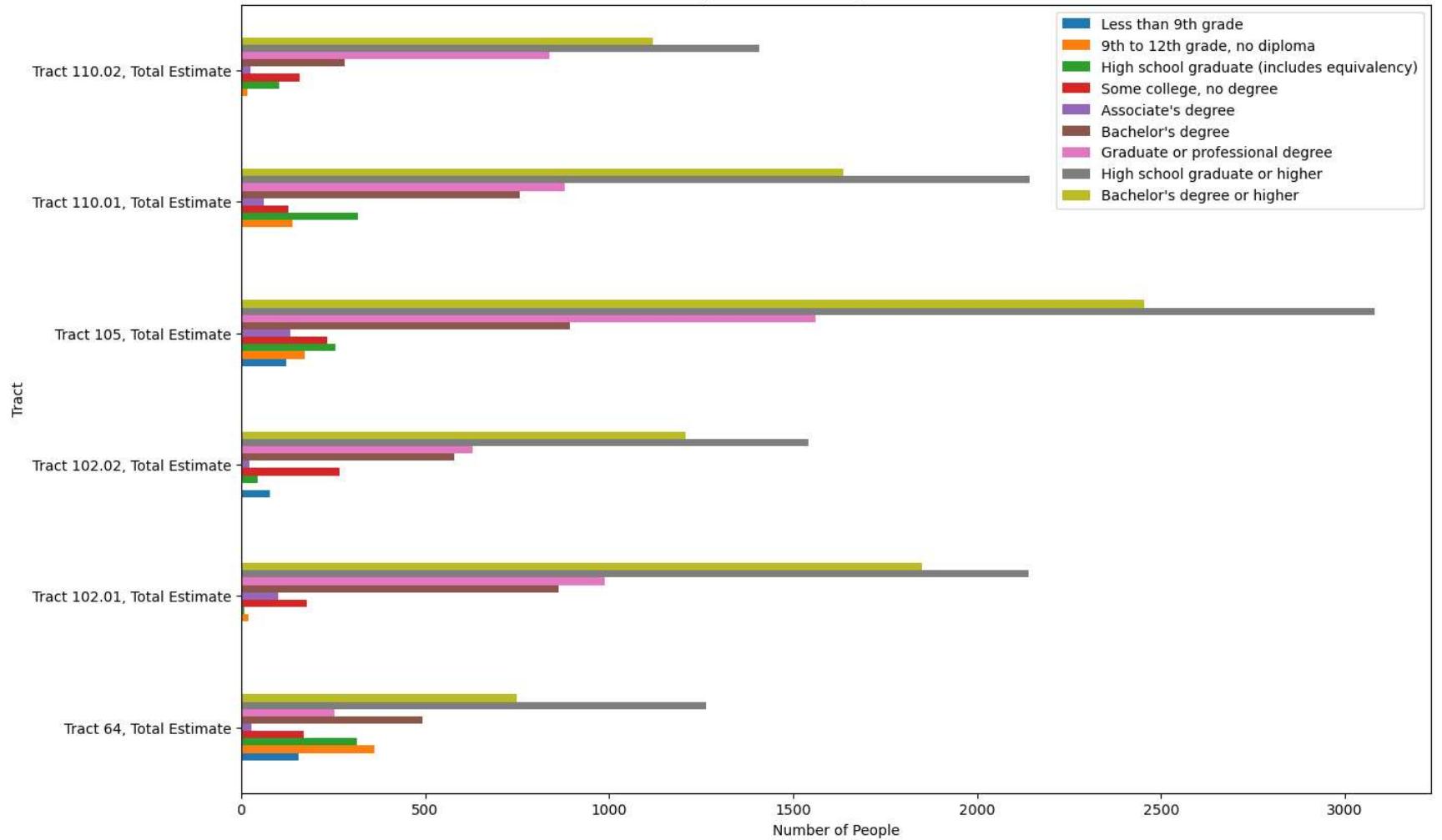
```
In [11]: # Extract the relevant data from the DataFrame
2 tract_cols = [col for col in df.columns if ('Tract' in col and 'Total' in col and 'Estimate' in col) or 'Label (Grouping)' in col]
3 tract_data = df.loc[7:15, tract_cols]
4 tract_data = tract_data.sort_index()
5
6 ##### Rename the rows
7 labels = tract_data['Label (Grouping)'].values
8
9 # Get rid of the '\xa0'
10 for i, label in enumerate(labels):
11     labels[i] = label.replace('\xa0', '')
12
13
14 for i in range(len(labels)):
15     tract_data = tract_data.rename(index={7+i: labels[i]})
```

16

```
17 ##### Drop the grouping column
18 tract_data = tract_data.drop(tract_data.columns[0], axis=1)
19
20 ##### Convert the columns to int
21 for col in tract_data.columns:
22     tract_data[col] = tract_data[col].str.replace(',', '').astype(int)
23
24 #####
25 tract_data = tract_data.T
26 bar_plot = tract_data.plot(kind='barh', title='Degree Breakdown Ages 25+ in 2021', figsize=(15, 10))
27 bar_plot.set_xlabel('Number of People')
28 bar_plot.set_ylabel('Tract')
29 plt.savefig(f'Plots/tract_25plus_degree_breakdown.png')
30 bar_plot
```

```
Out[11]: <AxesSubplot:title={'center':'Degree Breakdown Ages 25+ in 2021'}, xlabel='Number of People', ylabel='Tract'>
```

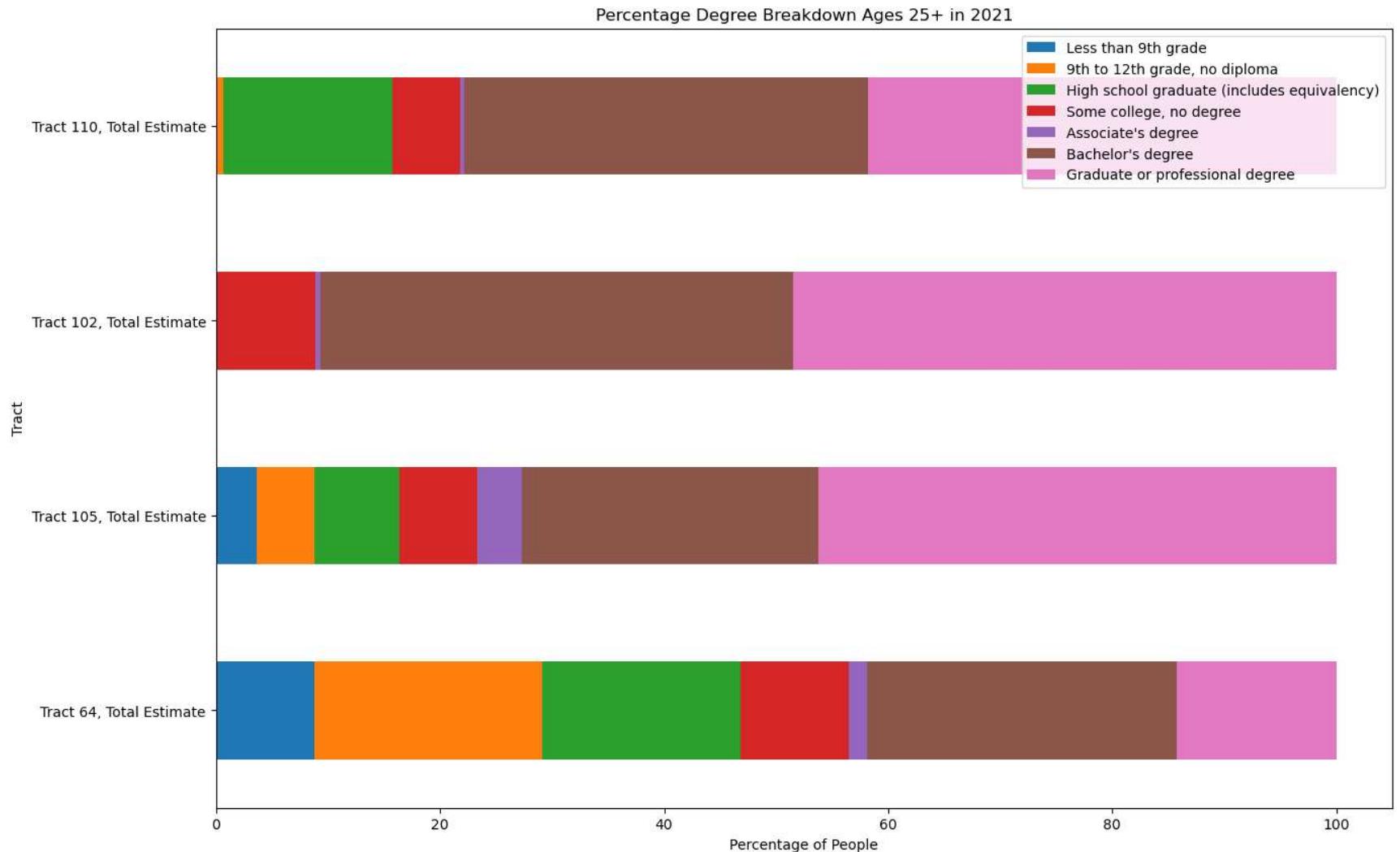
Degree Breakdown Ages 25+ in 2021



```
In [13]: # Now we want to see the percentages of each degree type
2
3 tract_cols = [col for col in df.columns if ('Tract' in col and 'Total' in col and 'Estimate' in col) or 'Label (Grouping)' in col]
4 tract_data = df.loc[7:15, tract_cols]
5 tract_data = tract_data.sort_index()
6
7 ##### Rename the rows
8 labels = tract_data['Label (Grouping)'].values
9
10 # Get rid of the '\xa0'
11 for i, label in enumerate(labels):
12     labels[i] = label.replace('\xa0', '')
13
14
15 for i in range(len(labels)):
16     tract_data = tract_data.rename(index={7+i: labels[i]})
17
18 ##### Drop the grouping column
19 tract_data = tract_data.drop(tract_data.columns[0], axis=1)
20
21 # Drop the "degree or higher rows"
22 tract_data = tract_data.drop(tract_data.index.values[-1])
23 tract_data = tract_data.drop(tract_data.index.values[-1])
24
25 ##### Combine decimal tracts
26 tract_data['Tract 102, Total Estimate'] = tract_data['Tract 102.01, Total Estimate'] + tract_data['Tract 102.02, Total Estimate']
27 tract_data['Tract 110, Total Estimate'] = tract_data['Tract 110.01, Total Estimate'] + tract_data['Tract 110.02, Total Estimate']
28
29 decimal_tracts = ['Tract 102.01, Total Estimate', 'Tract 102.02, Total Estimate', 'Tract 110.01, Total Estimate', 'Tract 110.02, Total Estimate']
30 for tract in decimal_tracts:
31     tract_data = tract_data.drop(tract, axis=1)
32 ##### Convert the columns to int
33 for col in tract_data.columns:
34     tract_data[col] = tract_data[col].str.replace(',', '').astype(int)
35
36 #####
37
38 tract_data = tract_data.T
39
40 # Turn into percentages
41 tracts = tract_data.T.columns.values
42 sums = tract_data.T.sum()
43
44 tract_data = tract_data.T
45
46 for label in tracts:
47     tract_data[label] = tract_data[label] / sums.T[label] * 100
48
49 for col in tract_data.columns:
50     plot_label = col.replace('Total Estimate', '')
51     sorted_data = tract_data[[col]].sort_values(by=col, ascending=False)
52     pie_plot = sorted_data.plot(kind='pie', y=col, title=f'Percentage {plot_label} Degree Breakdown Ages 25+ in 2021', figsize=(15, 10))
53     pie_plot.set_xlabel('')
54     pie_plot.set_ylabel('')
55     plt.savefig(f'Plots/{col}_pie_chart.png')
56     plt.clf()
```

```
57 | tract_data = tract_data.T
58 | bar_plot = tract_data.plot(kind='barh', stacked=True, title='Percentage Degree Breakdown Ages 25+ in 2021', figsize=(15, 10),)
59 | bar_plot.set_xlabel('Percentage of People')
60 | bar_plot.set_ylabel('Tract')
61 |
62 | plt.savefig(f'Plots/tract_25plus_percentage_breakdown.png')
63 | bar_plot
```

```
Out[13]: <AxesSubplot:title={'center':'Percentage Degree Breakdown Ages 25+ in 2021'}, xlabel='Percentage of People', ylabel='Tract'>
<Figure size 1500x1000 with 0 Axes>
```



What are the takeaways from this graph?

At a first glance, the degree breakdown is a lot to take in, and quite hard to read. However, by playing around with the data, we can reformat it to see the percentage level for each degree type for each tract.

What we can see here, is that Tracts 110 and 102 are primarily made up of advanced degrees. For Tract 110, roughly 40% have a bachelor's and 40% have a graduate/professional degree. In Tract 102, that combination makes up almost 90% of those above 25.

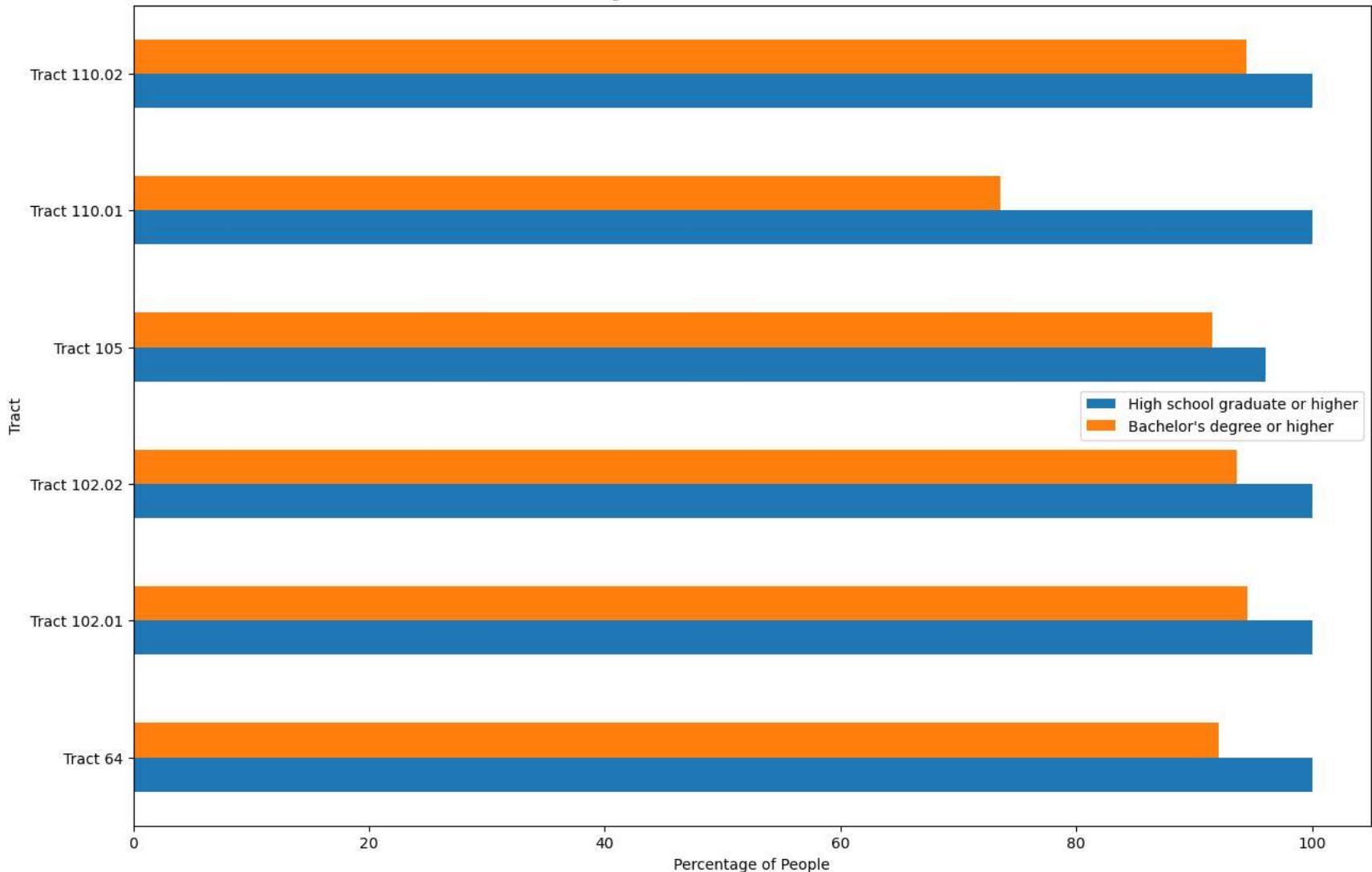
On the other hand, looking at Tracts 105 and 64, we start to see a big shake up. In Tract 64, roughly half of those older than 25 have **at most** have a high school education. Roughly 10% didn't go to high school at all, and another 20% didn't graduate from high school.

Demographic Information

```
In [16]: # Extract the relevant data from the DataFrame
1 tract_cols = [col for col in df.columns if ('Tract' in col and 'Percent' in col and \
3                                         ('Male Estimate' in col)) or 'Label (Grouping)' in col]
4
5 rows = [30, 31]
6 tract_data = df.loc[rows, tract_cols]
7
8 ##### Rename the rows
9 labels = tract_data['Label (Grouping)'].values
10
11 # Get rid of the '\xa0'
12
13 for i, label in enumerate(labels):
14     labels[i] = label.replace('\xa0', '')
15
16 for i in range(len(labels)):
17     tract_data = tract_data.rename(index={30+i: labels[i]})
18
19 tract_data = tract_data.T
20 index = tract_data.index.values
21 for i in range(len(index)):
22     new = index[i].replace(', Percent Male Estimate', '')
23     tract_data = tract_data.rename(index={index[i]: new})
24 tract_data = tract_data.T
25
26 ##### Drop the grouping column
27 tract_data = tract_data.drop(tract_data.columns[0], axis=1)
28 tract_data = tract_data.T
29
30 ##### Convert the columns to int
31 for col in tract_data.columns:
32     tract_data[col] = tract_data[col].str.replace('%', '').astype(float)
33
34 #####
35
36 bar_plot = tract_data.plot(kind='barh', stacked=False, title='Percentage White Alone Male Graduation in 2021', figsize=(15, 10),)
37 bar_plot.set_xlabel('Percentage of People')
38 bar_plot.set_ylabel('Tract')
39 plt.savefig(f'Plots/white_alone_male_graduation.png')
40 bar_plot
```

Out[16]: <AxesSubplot:title={'center':'Percentage White Alone Male Graduation in 2021'}, xlabel='Percentage of People', ylabel='Tract'>

Percentage White Alone Male Graduation in 2021



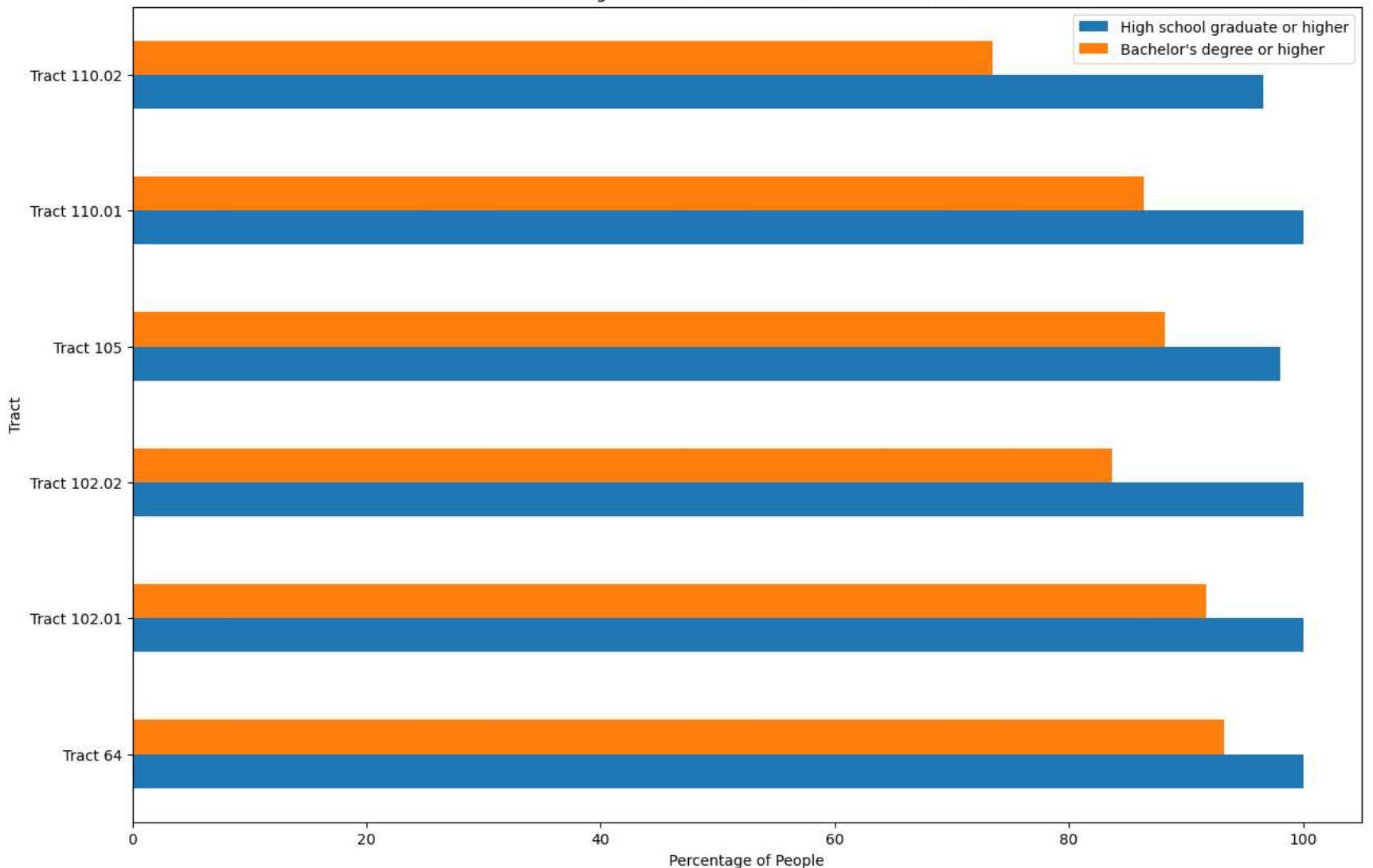
```
In [18]: # Extract the relevant data from the DataFrame
1 tract_cols = [col for col in df.columns if ('Tract' in col and 'Percent' in col and \
2                                         ('Female Estimate' in col)) or 'Label (Grouping)' in col]
3
4
5 rows = [30, 31]
6 tract_data = df.loc[rows, tract_cols]
7
8 ##### Rename the rows
9 labels = tract_data['Label (Grouping)'].values
10
11 # Get rid of the '\xa0'
12
13 for i, label in enumerate(labels):
14     labels[i] = label.replace('\xa0', '')
15
16 for i in range(len(labels)):
17     tract_data = tract_data.rename(index={30+i: labels[i]})
```

18

```
19 tract_data = tract_data.T
20 index = tract_data.index.values
21 for i in range(len(index)):
22     new = index[i].replace(' ', 'Percent Female Estimate', '')
23     tract_data = tract_data.rename(index={index[i]: new})
24 tract_data = tract_data.T
25
26 ##### Drop the grouping column
27 tract_data = tract_data.drop(tract_data.columns[0], axis=1)
28 tract_data = tract_data.T
29
30 ##### Convert the columns to int
31 for col in tract_data.columns:
32     tract_data[col] = tract_data[col].str.replace('%', '').astype(float)
33
34 #####
35
36 bar_plot = tract_data.plot(kind='barh', stacked=False, title='Percentage White Alone Female Graduation in 2021',\
37                             figsize=(15, 10),)
38 bar_plot.set_xlabel('Percentage of People')
39 bar_plot.set_ylabel('Tract')
40 plt.savefig(f'Plots/white_alone_female_graduation.png')
41 bar_plot
```

```
Out[18]: <AxesSubplot:title={'center':'Percentage White Alone Female Graduation in 2021'}, xlabel='Percentage of People', ylabel='Tract'>
```

Percentage White Alone Female Graduation in 2021



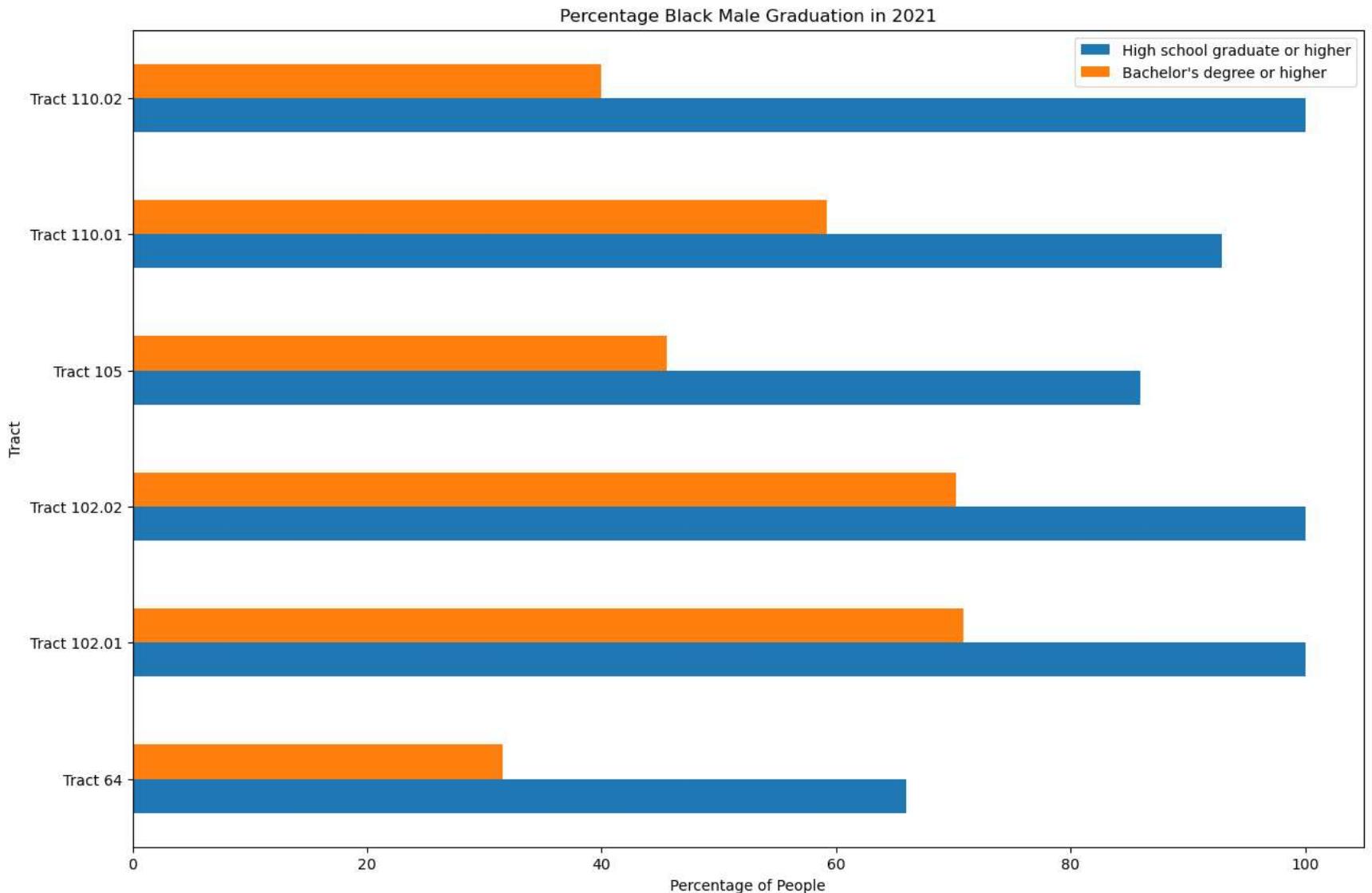
In [19]:

```
1 # Extract the relevant data from the DataFrame
2 tract_cols = [col for col in df.columns if ('Tract' in col and 'Percent' in col and \
3                                                 ('Male Estimate' in col)) or 'Label (Grouping)' in col]
4
5 rows = [36, 37]
6 tract_data = df.loc[rows, tract_cols]
7
8 ##### Rename the rows
9 labels = tract_data['Label (Grouping)'].values
10
11 # Get rid of the '\xa0'
12
13 for i, label in enumerate(labels):
14     labels[i] = label.replace('\xa0', '')
15
16 for i in range(len(labels)):
17     tract_data = tract_data.rename(index={36+i: labels[i]})
```

18

```
19 tract_data = tract_data.T
20 index = tract_data.index.values
21 for i in range(len(index)):
22     new = index[i].replace(' ', 'Percent Male Estimate', '')
23     tract_data = tract_data.rename(index={index[i]: new})
24 tract_data = tract_data.T
25
26 ##### Drop the grouping column
27 tract_data = tract_data.drop(tract_data.columns[0], axis=1)
28 tract_data = tract_data.T
29
30 ##### Convert the columns to int
31 for col in tract_data.columns:
32     tract_data[col] = tract_data[col].str.replace('%', '').astype(float)
33
34 #####
35
36 bar_plot = tract_data.plot(kind='barh', stacked=False, title='Percentage Black Male Graduation in 2021',\
37                             figsize=(15, 10),)
38 bar_plot.set_xlabel('Percentage of People')
39 bar_plot.set_ylabel('Tract')
40 plt.savefig(f'Plots/black_alone_male_graduation.png')
41 bar_plot
```

Out[19]: <AxesSubplot:title={'center':'Percentage Black Male Graduation in 2021'}, xlabel='Percentage of People', ylabel='Tract'>



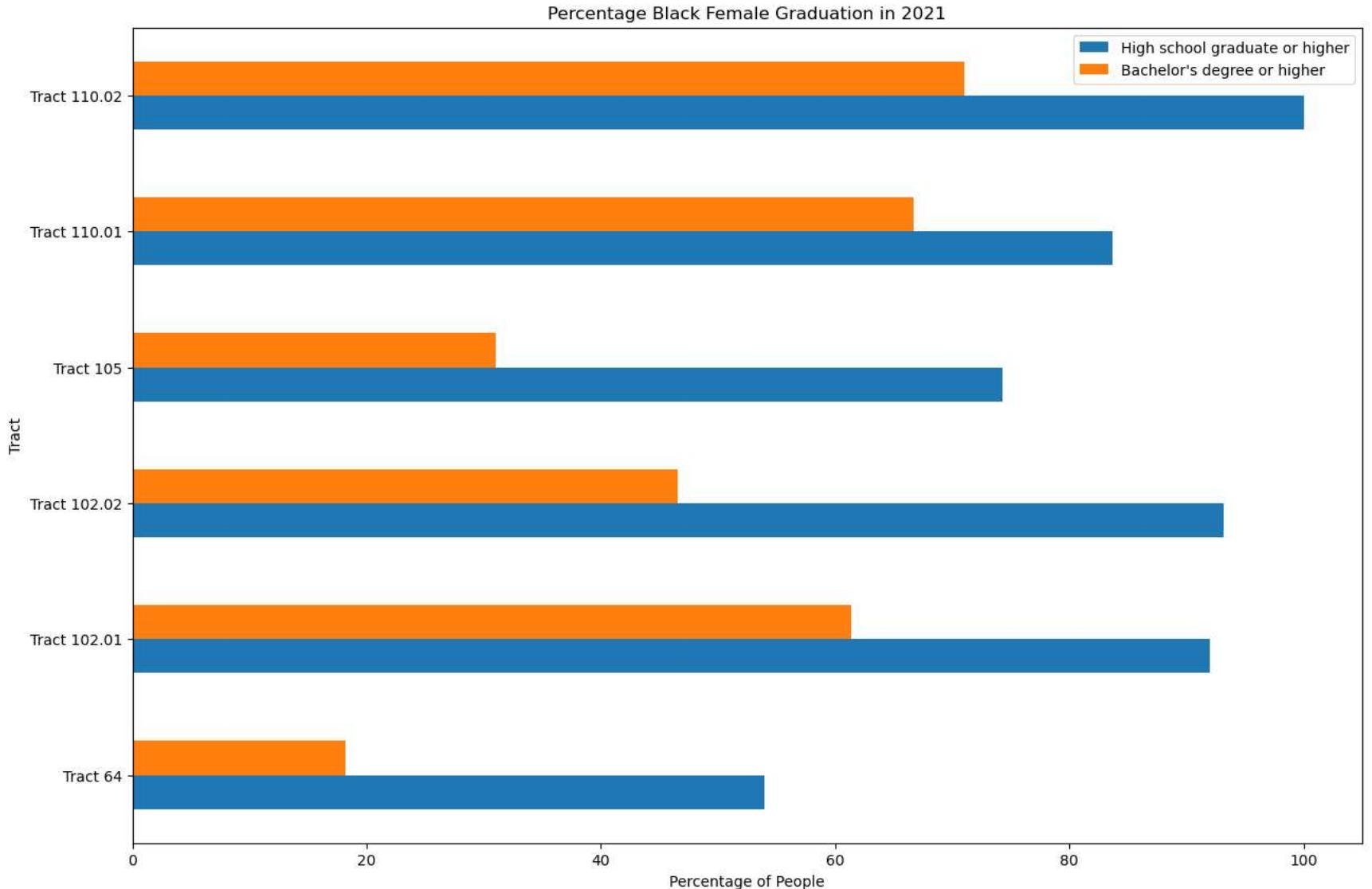
In [20]:

```
1 # Extract the relevant data from the DataFrame
2 tract_cols = [col for col in df.columns if ('Tract' in col and 'Percent' in col and \
3                                                 ('Female Estimate' in col)) or 'Label (Grouping)' in col]
4
5 rows = [36, 37]
6 tract_data = df.loc[rows, tract_cols]
7
8 ##### Rename the rows
9 labels = tract_data['Label (Grouping)'].values
10
11 # Get rid of the '\xa0'
12
13 for i, label in enumerate(labels):
14     labels[i] = label.replace('\xa0', '')
15
16 for i in range(len(labels)):
17     tract_data = tract_data.rename(index={36+i: labels[i]})
```

18

```
19 tract_data = tract_data.T
20 index = tract_data.index.values
21 for i in range(len(index)):
22     new = index[i].replace(' ', 'Percent Female Estimate', '')
23     tract_data = tract_data.rename(index={index[i]: new})
24 tract_data = tract_data.T
25
26 ##### Drop the grouping column
27 tract_data = tract_data.drop(tract_data.columns[0], axis=1)
28 tract_data = tract_data.T
29
30 ##### Convert the columns to int
31 for col in tract_data.columns:
32     tract_data[col] = tract_data[col].str.replace('%', '').astype(float)
33
34 #####
35
36 bar_plot = tract_data.plot(kind='barh', stacked=False, title='Percentage Black Female Graduation in 2021',\
37                             figsize=(15, 10),)
38 bar_plot.set_xlabel('Percentage of People')
39 bar_plot.set_ylabel('Tract')
40 plt.savefig(f'Plots/black_alone_female_graduation.png')
41 bar_plot
```

Out[20]: <AxesSubplot:title={'center':'Percentage Black Female Graduation in 2021'}, xlabel='Percentage of People', ylabel='Tract'>



What are the takeaways from these graphs?

What we can see is a white-alone males and females have high graduation rates in all the selected tracts, while black-alone males and females have significantly lower graduation rates, especially in Tracts 64 and 105. With these percentages, it will be important to look at how many students of each demographic there are per the tracts; however even so, it offers valuable insights that some sort of help could help reduce these disparities.

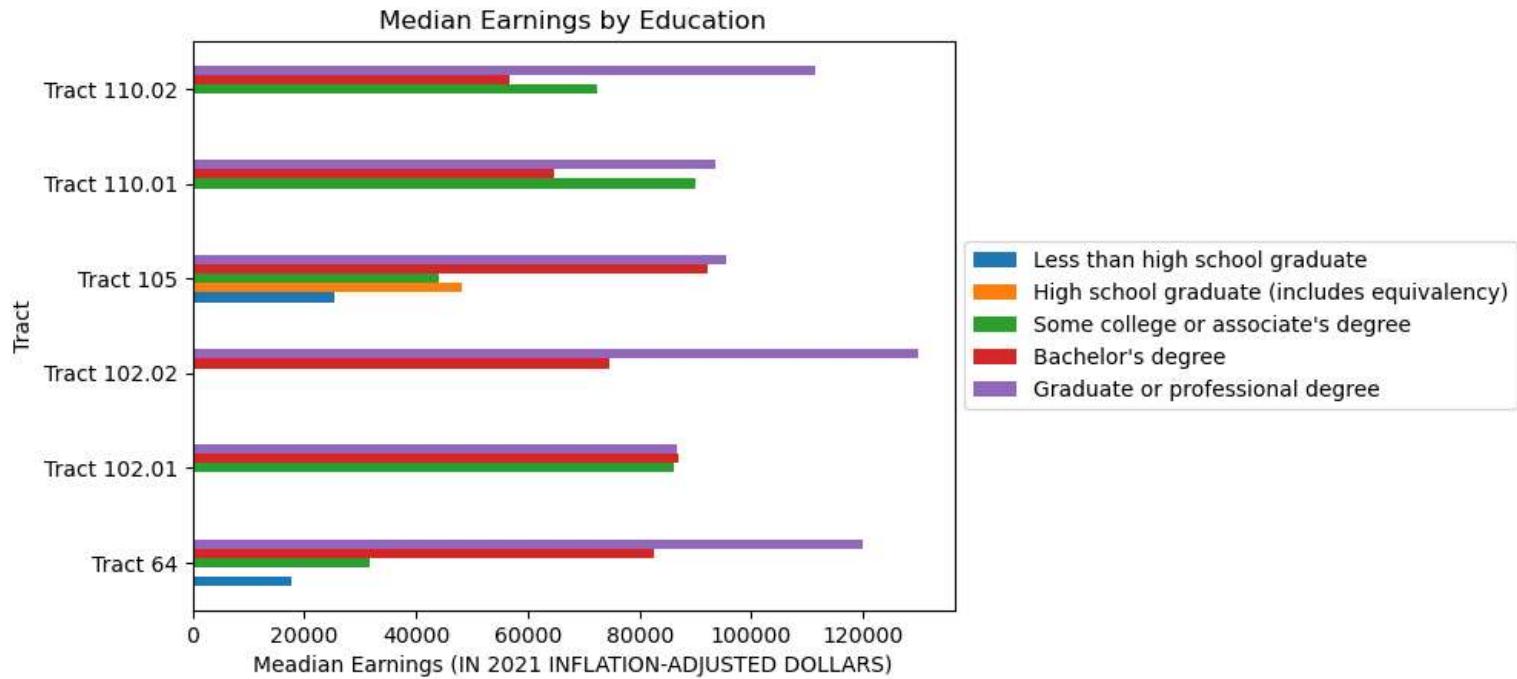
Education v Median Income

In [34]:

```
1 # This is the overall breakdown, let's try breaking down by degree
2
3 # Extract the relevant data from the DataFrame
4 tract_cols = [col for col in df.columns if 'Tract' in col and 'Total' in col and 'Estimate' in col]
5 tract_data = df.loc[63:67, tract_cols]
6 tract_data = tract_data.sort_index()
7
8 # Use a lambda function to remove the ', Total Estimate' suffix from each column name
9 tract_data.columns = tract_data.columns.map(lambda x: x.replace(', Total Estimate', ''))
10
11 ##### Rename the rows
12 labels = ['Less than high school graduate', 'High school graduate (includes equivalency)',
13           'Some college or associate\'s degree', 'Bachelor\'s degree', 'Graduate or professional degree']
14
15 for i in range(len(labels)):
16     tract_data = tract_data.rename(index={63+i: labels[i]})
```

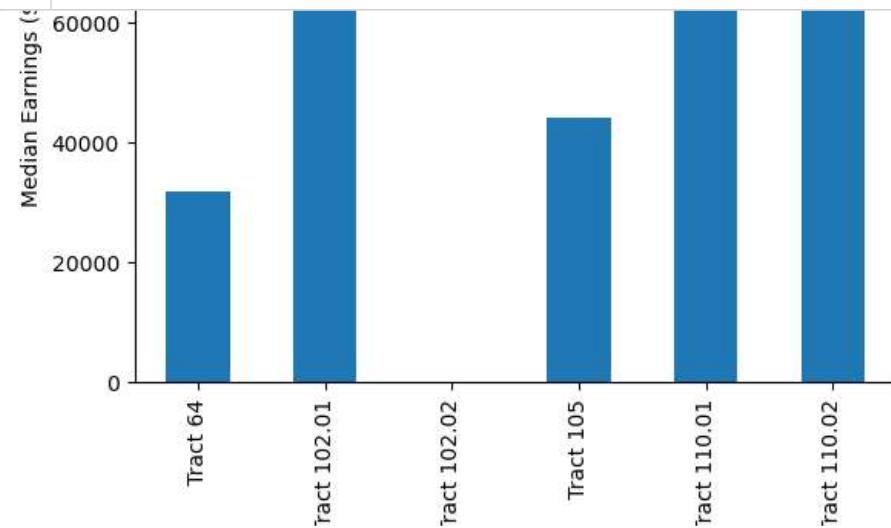
18

```
19 ##### Convert the columns to int
20
21 tract_data = tract_data.replace('-', '0')
22
23 for col in tract_data.columns:
24     tract_data[col] = tract_data[col].str.replace(',', '').astype(int)
25
26 #####
27 tract_data = tract_data.T
28
29
30 fig, ax = plt.subplots()
31 bar_plot = tract_data.plot(kind='barh', title='Median Earnings by Education', ax=ax)
32 bar_plot.set_xlabel('Median Earnings (IN 2021 INFLATION-ADJUSTED DOLLARS)')
33 bar_plot.set_ylabel('Tract')
34 # Move the Legend to the right of the plot and expand the plot
35 ax.legend(loc='center left', bbox_to_anchor=(1.0, 0.5))
36 plt.savefig(f'Plots/tract_18_24_breakdown.png')
```



In [35]:

```
1 # Graph with Less than high school graduate, tract on y axis, earnings x-axis
2
3 for label in labels:
4     education_level = tract_data[label]
5     title = f'Median Earnings for {label} per Tract'
6     bar_plot = education_level.plot(kind='bar', title=title)
7     bar_plot.set_ylabel('Median Earnings ($)')
8     bar_plot.set_xlabel('Census Tract')
9     plt.savefig(f'Plots/{title}.png')
10    plt.show()
```



What does this graph mean?

We can see that Tract 64 and 105 are the only ones that report some sort of earnings for those with a maximum education of high school. Their median is less than \$50,000 with half of those in tract 64 with less than a high school education earning less than \$17,500.

We also see a big difference in median earning for those with either some college education or an associate's degree. With the Tract 102 and 110 medians being 2x higher than tract 105, and 3x higher than tract 64. This seems to suggest that even with some college education, the people in tracts 64 & 105 aren't able to get jobs that match their education level.

However when we look at a bachelor's degree or higher, we see very similar median earnings, suggesting that those with a 4-year degree or higher have very similar earning opportunities.

The big thing to note though, is in the previous charts we see that tracts 102 & 110 are very heavily populated with 4-year degree holders or higher, while tracts 64 and 105 contain significantly more people that have at most a high school level education. This tracks with the big differences between how many people live below the poverty line between the two tracts. To gain a better understanding of why this may be happening, we would have to obtain data specific to the junior high and high schools in the area.

Computer Access Analysis

Here I take a look at B28006: EDUCATIONAL ATTAINMENT BY PRESENCE OF A COMPUTER AND TYPES OF INTERNET SUBSCRIPTION IN HOUSEHOLD broken down by block group for census tracts 64, 102, 105, 110.

Data Loading

```
In [3]: # Read in the file
      2 computer_table_name = 'ACSDT5Y2021.B28006.csv'
      3 computer_df = pd.read_csv(PATH + computer_table_name)
```

Data cleaning

```
In [23]: print('Initial column name format: ', computer_df.columns[1])
      2
      3 # Filter out columns with 'Margin of Error' in their names
      4 computer_df = computer_df.filter(regex='^(?!.*Margin of Error)')
      5
      6
      7 # Rename the columns using the .str.replace() method and regex
      8 # This removes ', District of Columbia, District of Columbia!!Estimate' from each column
      9 pattern = 'Block Group ([\d.]+), Census Tract ([\d.]+), District of Columbia, District of Columbia!!([^\!]+)'
     10 replacement = r'Block Group \g<1>, Tract \g<2>'
     11 computer_df.columns = computer_df.columns.str.replace(pattern, replacement)
     12
     13 print('Updated column name format: ', computer_df.columns[1])
```

Initial column name format: Block Group 1, Tract 64

Updated column name format: Block Group 1, Tract 64

C:\Users\konic\AppData\Local\Temp\ipykernel_15036\3982640668.py:11: FutureWarning: The default value of regex will change from True to False in a future version.

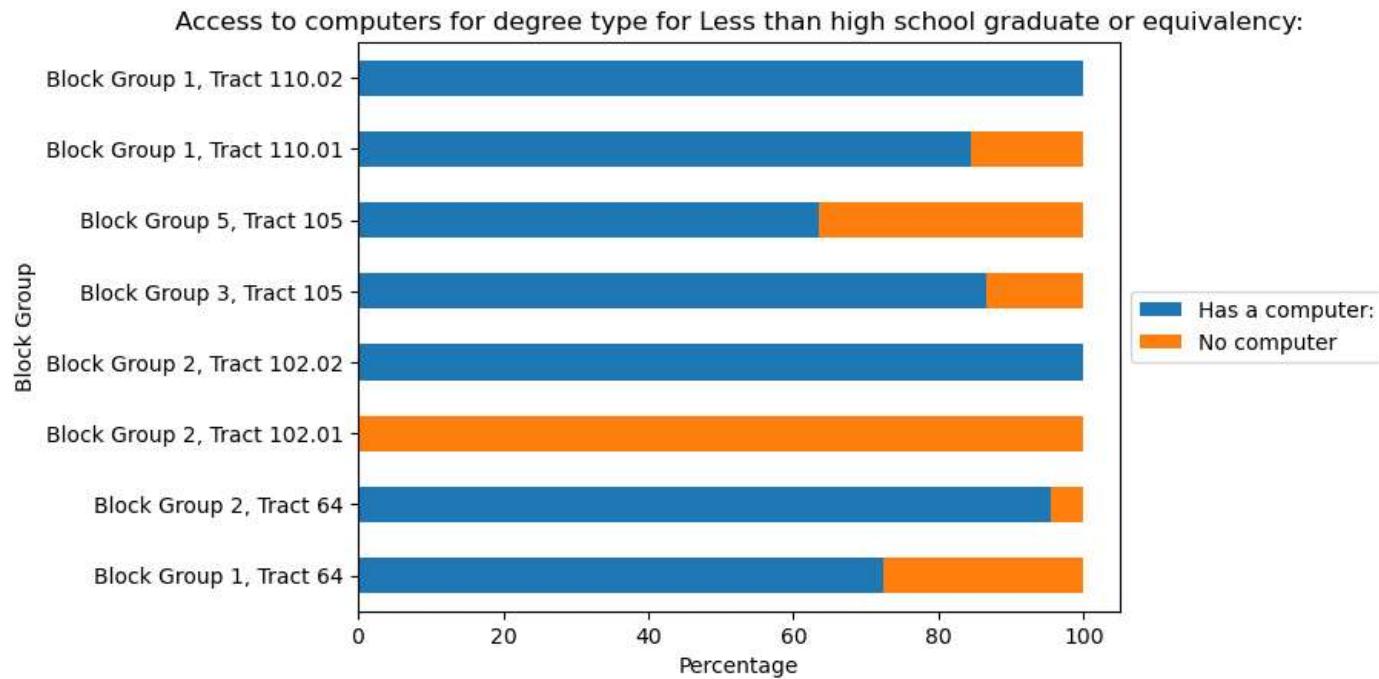
```
computer_df.columns = computer_df.columns.str.replace(pattern, replacement)
```

Computer Access Analysis

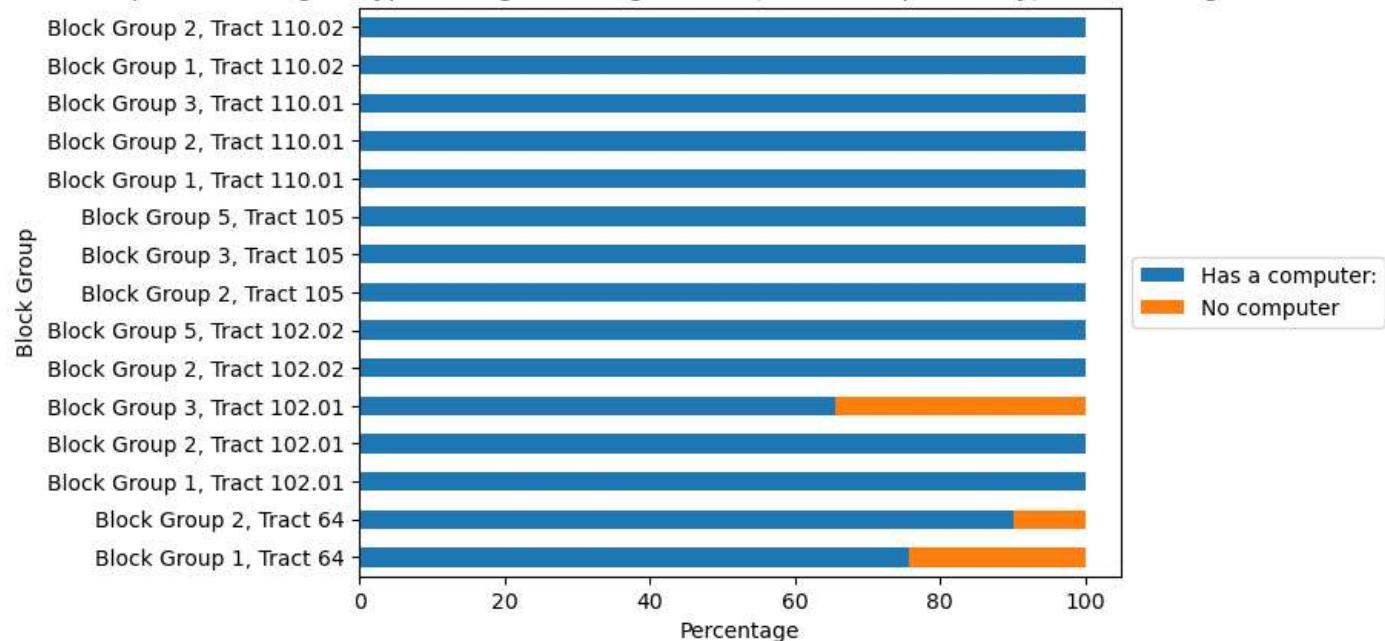
In [32]:

```
1 # Dataframe uses multiple dtypes, convert everything to string so we can remove symbols
2 computer_df = computer_df.astype(str)
3
4 # For each degree type
5 for i in range(1,14,6):
6     # Get the data for that degree
7     tract_data = computer_df[i:i+6]
8
9     labels = tract_data['Label (Grouping)'].values
10
11    # Get rid of the '\xa0'
12    for j, label in enumerate(labels):
13        labels[j] = label.replace('\xa0', '')
14
15    ##### Convert the columns to int
16    tract_data = tract_data.T[1:]
17    for col in tract_data.columns:
18        tract_data[col] = tract_data[col].str.replace(',', '').astype(float)
19
20    for j in range(len(labels)):
21        tract_data = tract_data.rename(columns={i+j: labels[j]})
22
23    ##### Remove any rows with values of 0
24    # Select rows where all columns have 0
25    zero_rows = tract_data[(tract_data == 0).all(axis=1)]
26
27    # Output the resulting rows
28    tract_data = tract_data.drop(zero_rows.index)
29
30    ##### Keep only 'Has a computer' and 'No computer' columns
31    degree_type = tract_data.columns[0]
32
33    # Total number of people of that degree per block group
34    normalizing_values = tract_data.loc[:, degree_type]
35
36    # Keep only these two columns
37    keep_columns = ['Has a computer:', 'No computer']
38    tract_data = tract_data.loc[:, keep_columns]
39
40    # Normalize & turn into percents
41    tract_data = tract_data.div(normalizing_values, axis=0) * 100
42
43    # Plot the chart
44    fig, ax = plt.subplots()
45    plot = tract_data.plot(kind='barh', stacked=True, ax=ax,
46                           title=f'Access to computers for degree type for {degree_type}')
47    ax.set_xlabel('Percentage')
48    ax.set_ylabel('Block Group')
49    ax.legend(loc='center left', bbox_to_anchor=(1.0, 0.5))
50
51    # Save the plot
52    filename = "Plots/comp_access"
53    if i == 1:
54        filename += "_less_than_highschool.png"
55    elif i == 7:
56        filename += "_highschool_associates.png"
```

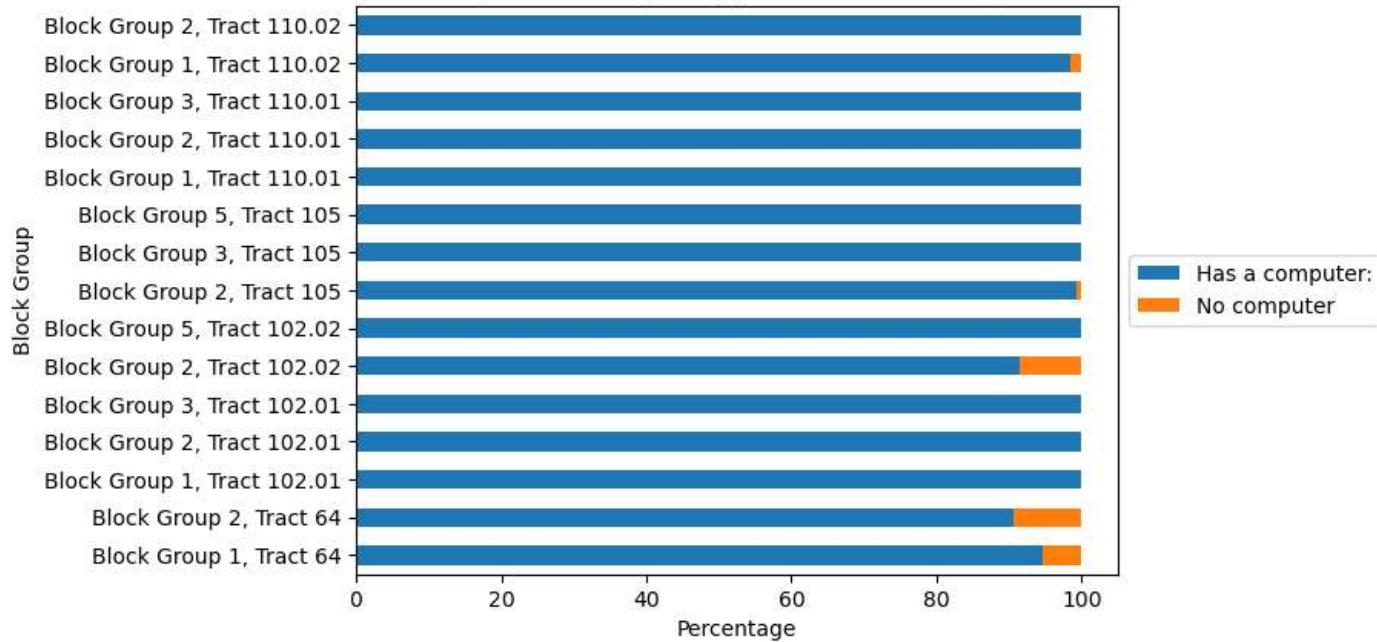
```
57 elif i == 13:  
58     filename += "_bachelors.png"  
59 plt.savefig(filename, bbox_inches='tight')  
60
```



Access to computers for degree type for High school graduate (includes equivalency) , some college or associate's degree :



Access to computers for degree type for Bachelor's degree or higher:



What does this graph mean?

Here we can see that most people do have computer access regardless of degree type, therefore, it's unlikely that computer access is the reason that there are such drastic differences in educational attainment between tracts.

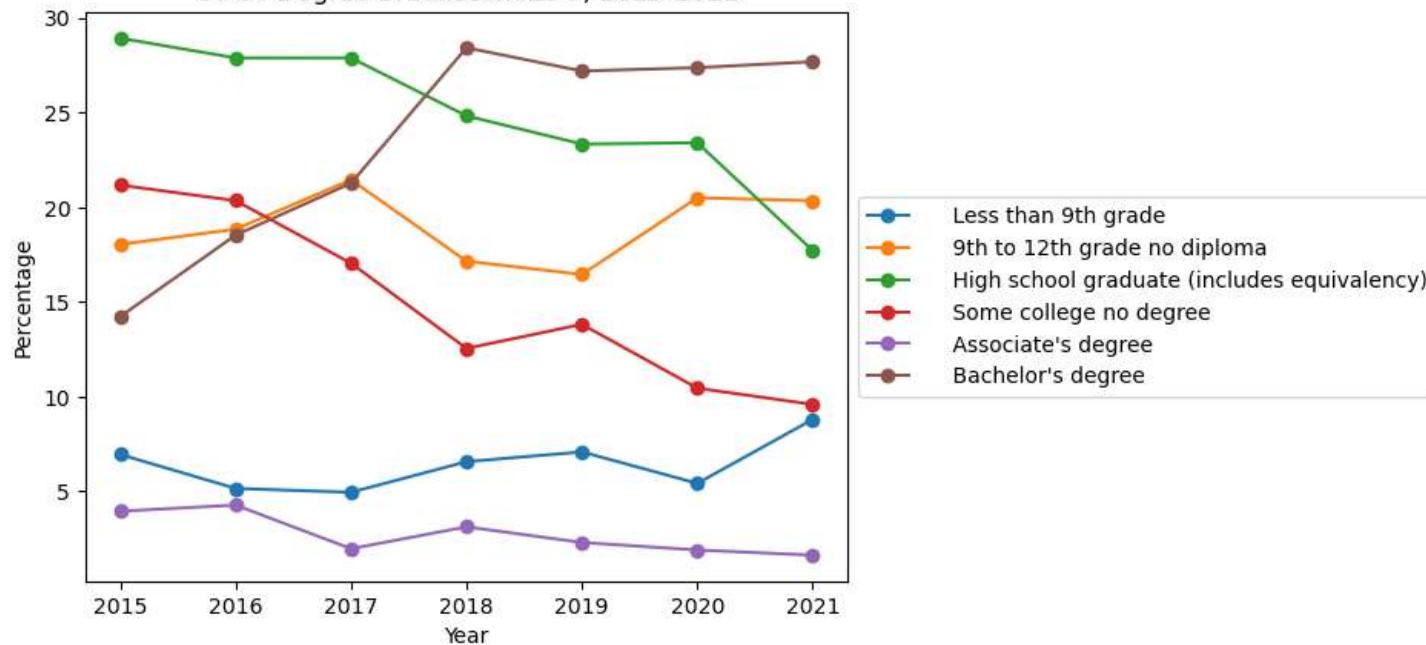
Trend over time analysis

This is the beginning of the trend over time analysis. However, this is where I ran into the challenge of not all the tracts being available.

In [183]:

```
1 ## Load files into df
2 dfs = []
3
4 # Loading in files from 2015-2021
5 # I don't Load the pre-2015 files because they store the data differently
6 for file in files[5:]:
7     year = file.split('.')[0][-4:]
8
9     temp_pd = pd.read_csv(PATH+file)
10    temp_pd.columns = temp_pd.columns.str.replace('Census Tract ([\d\.\.]+), District of Columbia, District of Columbia!!([!]+)!!([!]+)', r'Tract \1, DC, DC!!([!]+)!!([!]+)', re.IGNORECASE)
11
12    # Get just the 25+
13    temp_pd = temp_pd.iloc[6:13]
14    temp_pd = temp_pd.replace(',', ' ', regex=True)
15
16    if len(dfs) == 0:
17        temp_pd = temp_pd[['Label (Grouping)', 'Tract 64, Total Estimate']]
18        temp_pd.iloc[:, 1:] = temp_pd.iloc[:, 1:]._astype(int)
19        temp_pd.loc[7:13, 'Tract 64, Total Estimate'] /= temp_pd.loc[6, 'Tract 64, Total Estimate'] / 100
20
21    else:
22        temp_pd = temp_pd[['Tract 64, Total Estimate']]._astype(int)
23        col_name = 'Tract 64, Total Estimate'
24        divisor = temp_pd.loc[6, col_name]
25        temp_pd.loc[7:12, col_name] = 100 * temp_pd.loc[7:12, col_name] / divisor
26
27    temp_pd = temp_pd.rename(columns={'Tract 64, Total Estimate': year})
28    temp_pd = temp_pd.drop(index=6)
29    dfs.append(temp_pd)
30
31 df_concat = pd.concat(dfs, axis=1)
32
33 df_concat = df_concat.set_index('Label (Grouping)')
34 df_concat = df_concat.T
35
36 fig, ax = plt.subplots()
37 plot = df_concat.plot(kind='line', marker='o', ax=ax, title='CT 64 Degree Breakdown 25+, 2015-2021')
38 ax.legend(loc='center left', bbox_to_anchor=(1.0, 0.5))
39 ax.set_xlabel('Year')
40 ax.set_ylabel('Percentage')
41 plt.savefig('Plots/ct64_degrees_2015_2021.png', bbox_inches='tight')
```

CT 64 Degree Breakdown 25+, 2015-2021



In [182]:

```
1 ## Load files into df
2 dfs = []
3
4 # Loading in files from 2015-2021
5 # I don't Load the pre-2015 files because they store the data differently
6 for file in files[5:]:
7     year = file.split('.')[0][-4:]
8
9     temp_pd = pd.read_csv(PATH+file)
10    temp_pd.columns = temp_pd.columns.str.replace('Census Tract ([\d\.\.]+), District of Columbia, District of Columbia!!([!]+)!!([!]+)', r'Tract \1, DC, DC!!([!]+)!!([!]+)', re.IGNORECASE)
11
12    # Get just the 25+
13    temp_pd = temp_pd.iloc[1:6]
14    temp_pd = temp_pd.replace(',', '', regex=True)
15
16    if len(dfs) == 0:
17        temp_pd = temp_pd[['Label (Grouping)', 'Tract 64, Total Estimate']]
18        temp_pd.iloc[:, 1:] = temp_pd.iloc[:, 1:]._astype(int)
19        temp_pd.loc[1:6, 'Tract 64, Total Estimate'] /= temp_pd.loc[1, 'Tract 64, Total Estimate'] / 100
20
21    else:
22        temp_pd = temp_pd[['Tract 64, Total Estimate']]._astype(int)
23        col_name = 'Tract 64, Total Estimate'
24        divisor = temp_pd.loc[1, col_name]
25        temp_pd.loc[1:6, col_name] = 100 * temp_pd.loc[1:6, col_name] / divisor
26
27    temp_pd = temp_pd.rename(columns={'Tract 64, Total Estimate': year})
28    temp_pd = temp_pd.drop(index=1)
29    dfs.append(temp_pd)
30
31 df_concat = pd.concat(dfs, axis=1)
32
33 df_concat = df_concat.set_index('Label (Grouping)')
34 df_concat = df_concat.T
35
36 fig, ax = plt.subplots()
37 title = "CT 64 Degree Breakdown 18-24, 2015-2021"
38 plot = df_concat.plot(kind='line', marker='o', ax=ax, title=title)
39 ax.legend(loc='center left', bbox_to_anchor=(1.0, 0.5))
40 ax.set_xlabel('Year')
41 ax.set_ylabel('Percentage')
42 plt.savefig('Plots/ct64_degrees_18_2015_2021.png', bbox_inches='tight')
```

CT 64 Degree Breakdown 18-24, 2015-2021

