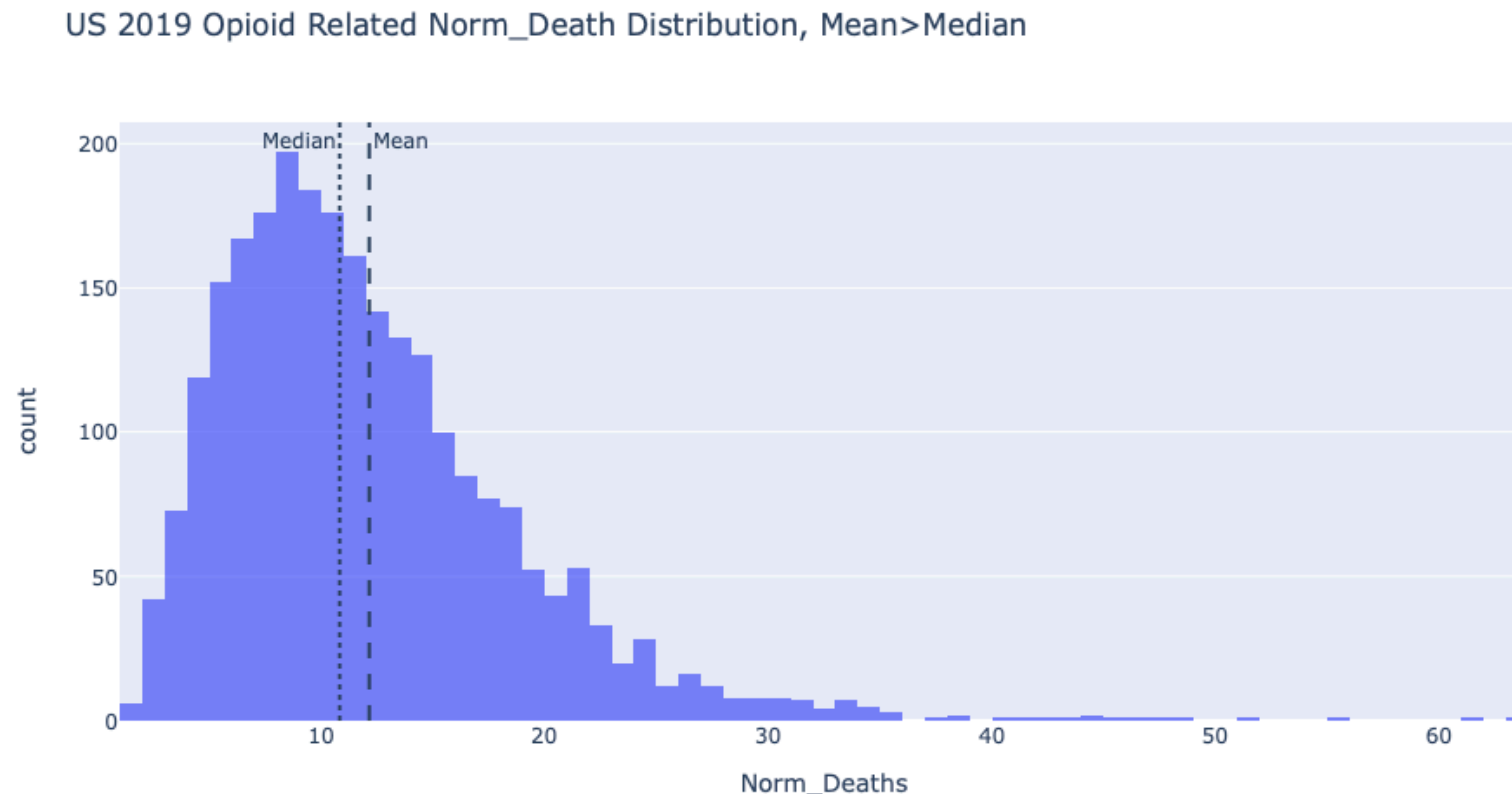


# Understand the distribution and measures of center for US opioid mortality for 2019

- For this task we need to import the data and then create a histogram.
- Adding the vertical lines of mean and median.
- If the histogram is skewed right, the mean is greater than the median (A few large values drive the mean upward).
- But they are really close to each other. I think it is better to consider it symmetric.



## Codify the column Normalized Deaths and store them as a “label” column

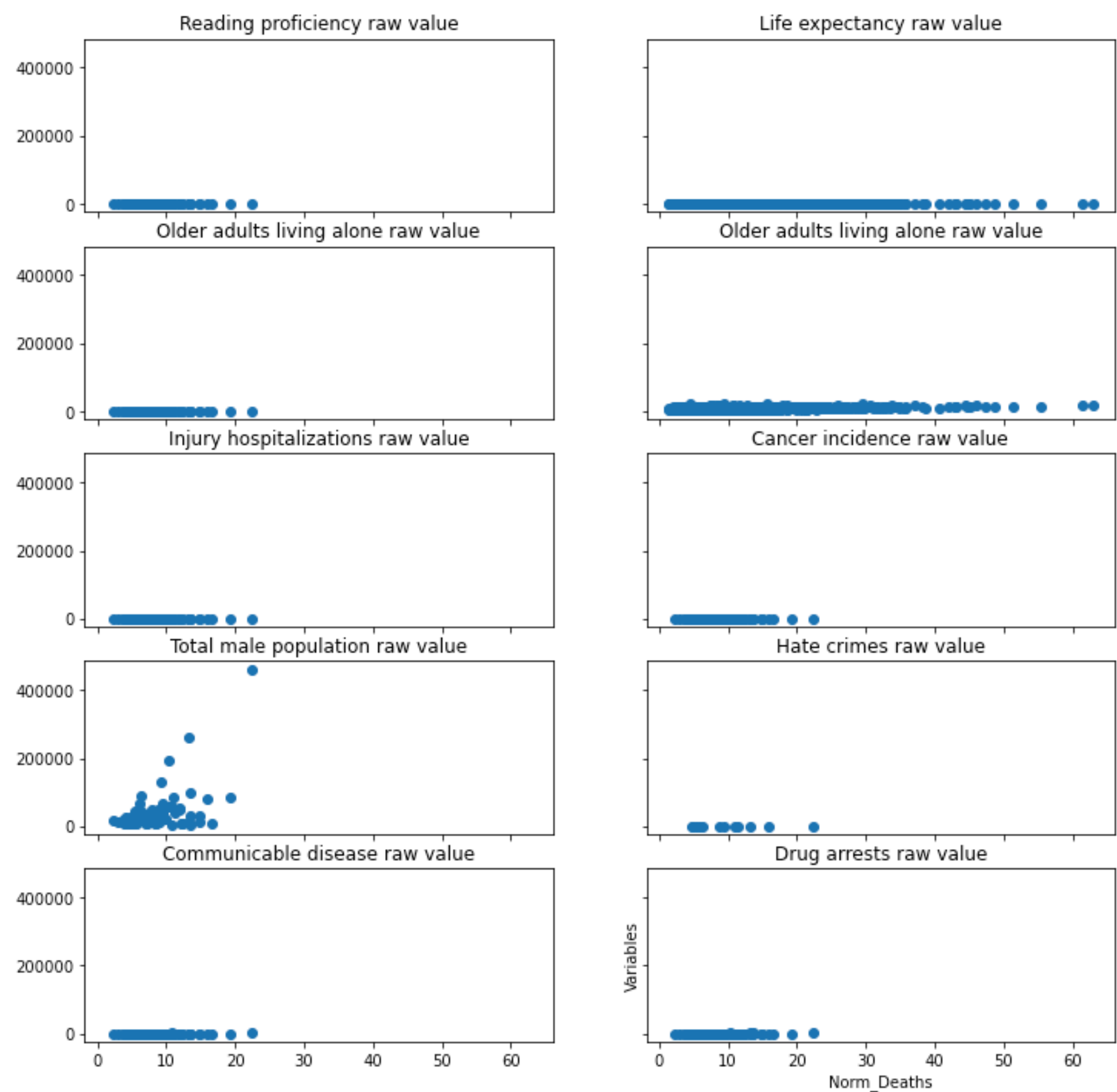
- Minimum Norm\_Deaths in my data frame is less than 1.999. So, I'll change the 1.99 to 0.
- I use pandas cut() method to categorize values in Norm\_Deaths columns into four categories.
- cut() method takes an input series that need to be grouped into various bins.
- cut() method takes a list of values that begins with the smallest value for the first bin and ends with the maximum value for the last bin. Break points between the bins also exist in the list. Here, I assumed the lower band to be inclusive.
- cut() method also takes a list of bins name. This list have one element less than bin list.

```
In [40]: super_df['label'] = pd.cut(super_df['Norm_Deaths'], bins = [0,8,11,16,64], labels = ['v_low','low','high','v_high'])
super_df['label'].unique()
```

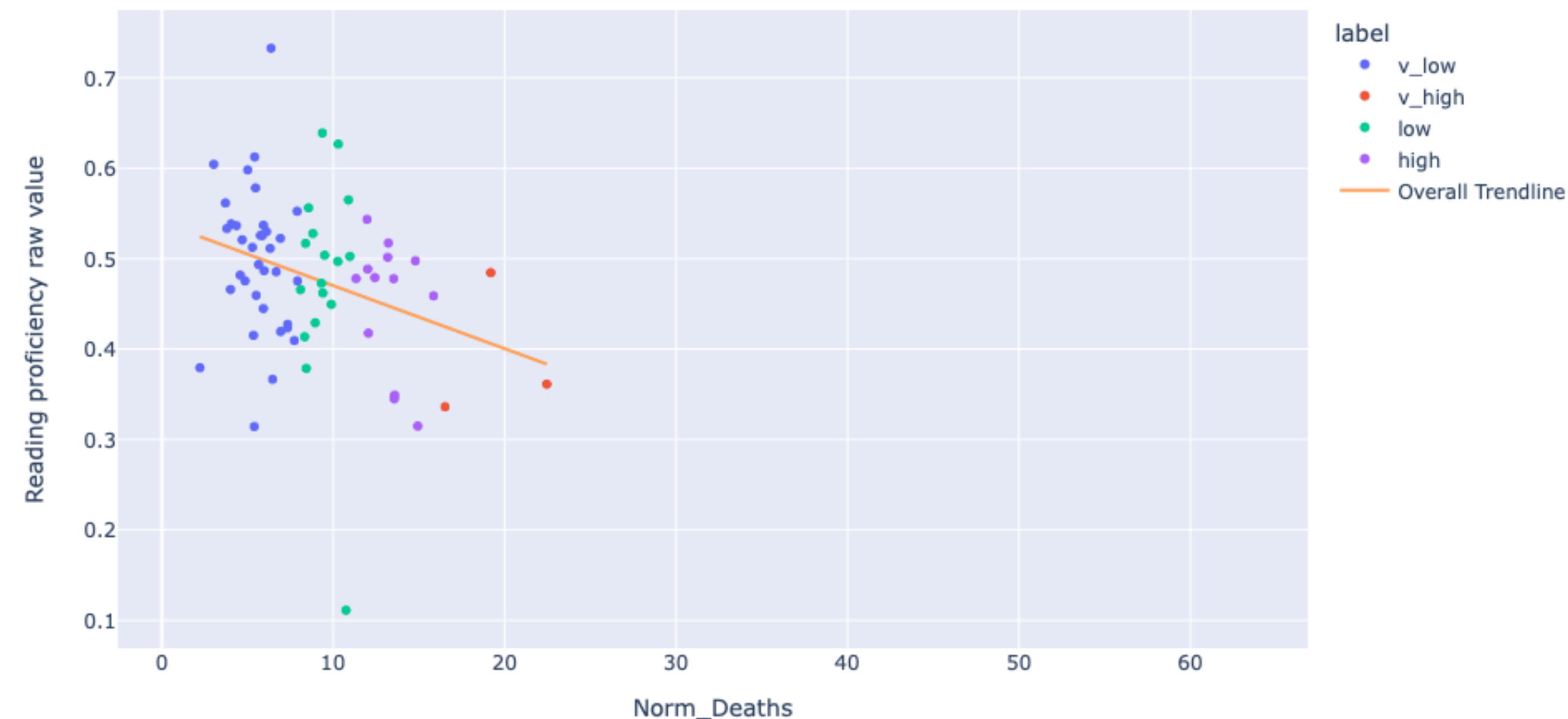
```
Out[40]: ['v_low', 'v_high', 'low', 'high']
Categories (4, object): ['v_low' < 'low' < 'high' < 'v_high']
```

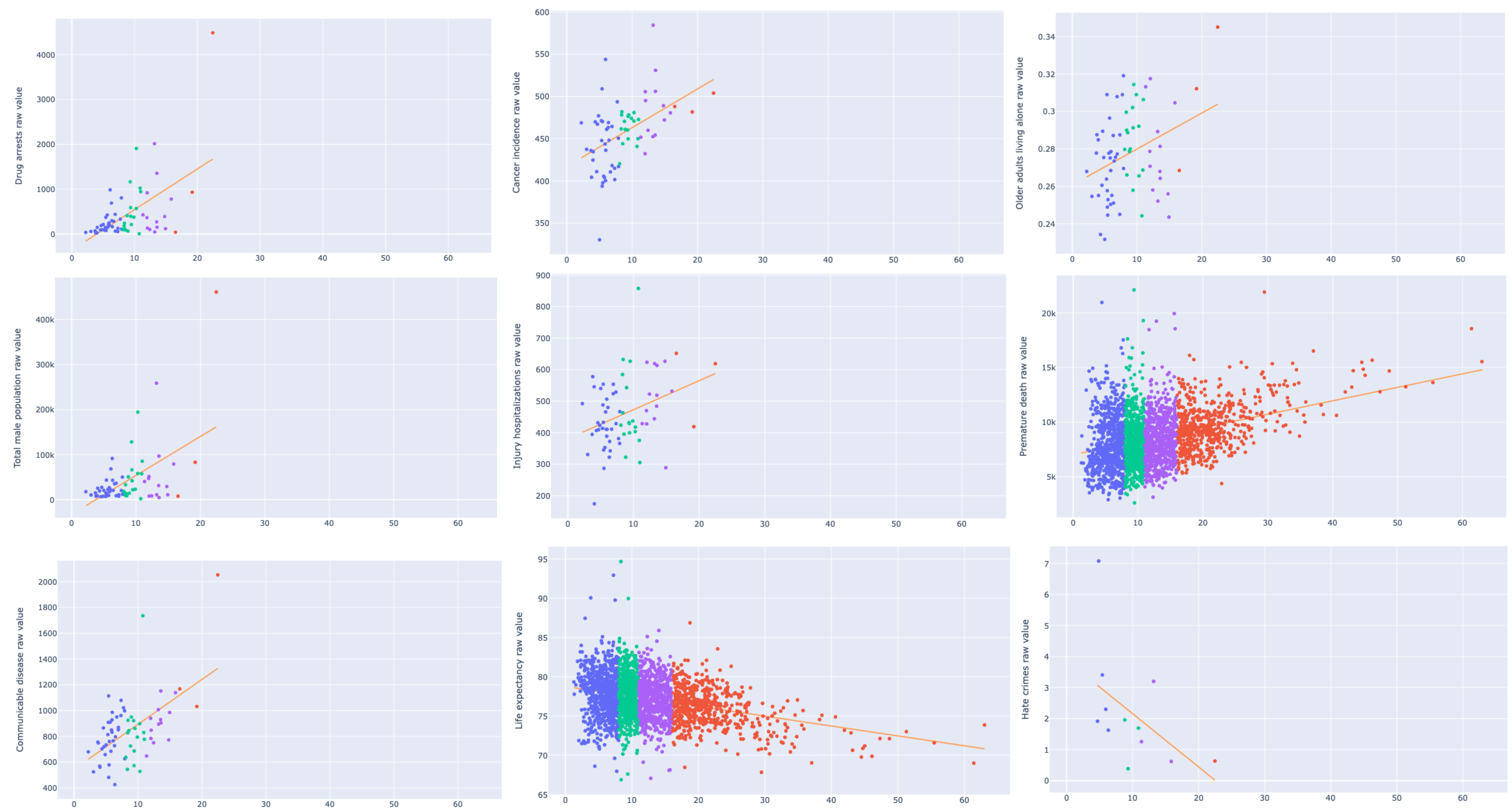
Now try the variables you have identified in Stage I and plot them as a second variable to Normalized Mortality in a scatter plot to observe any trends.

- For this task, first I tries to plot all of them in a subplot. But because their values are not in the same period the result was not satisfying.



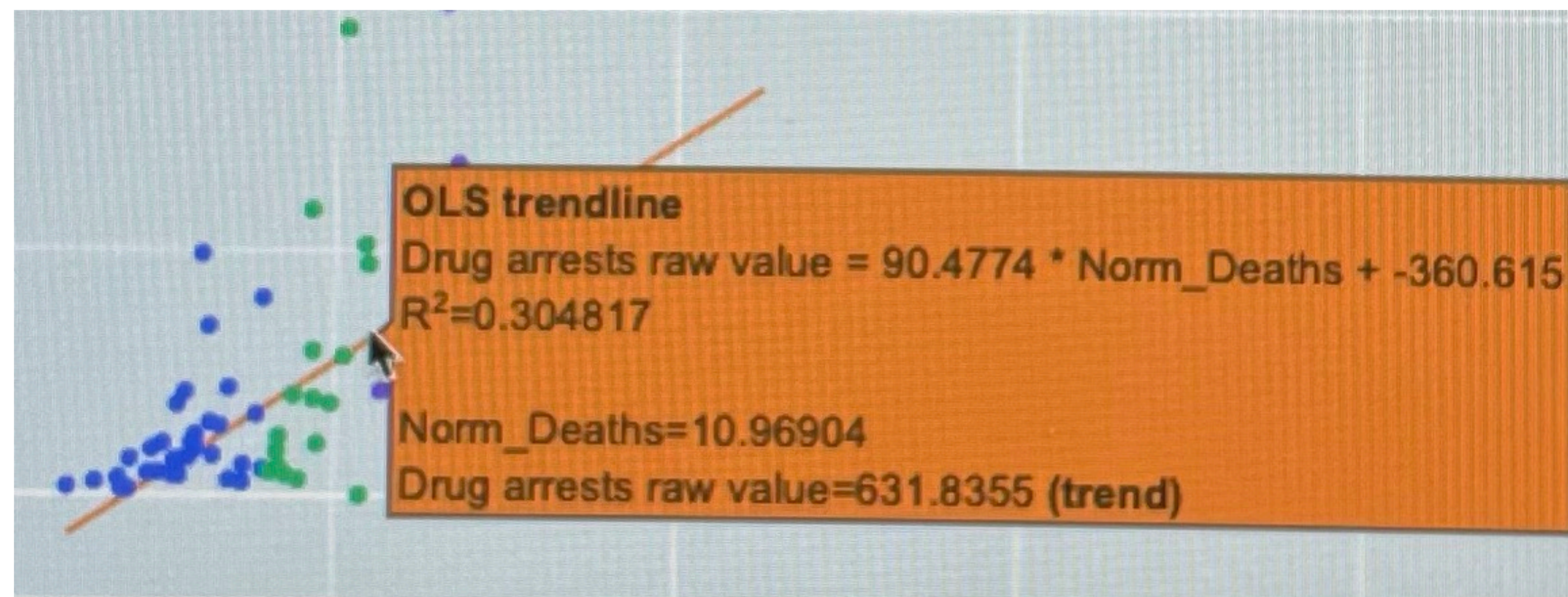
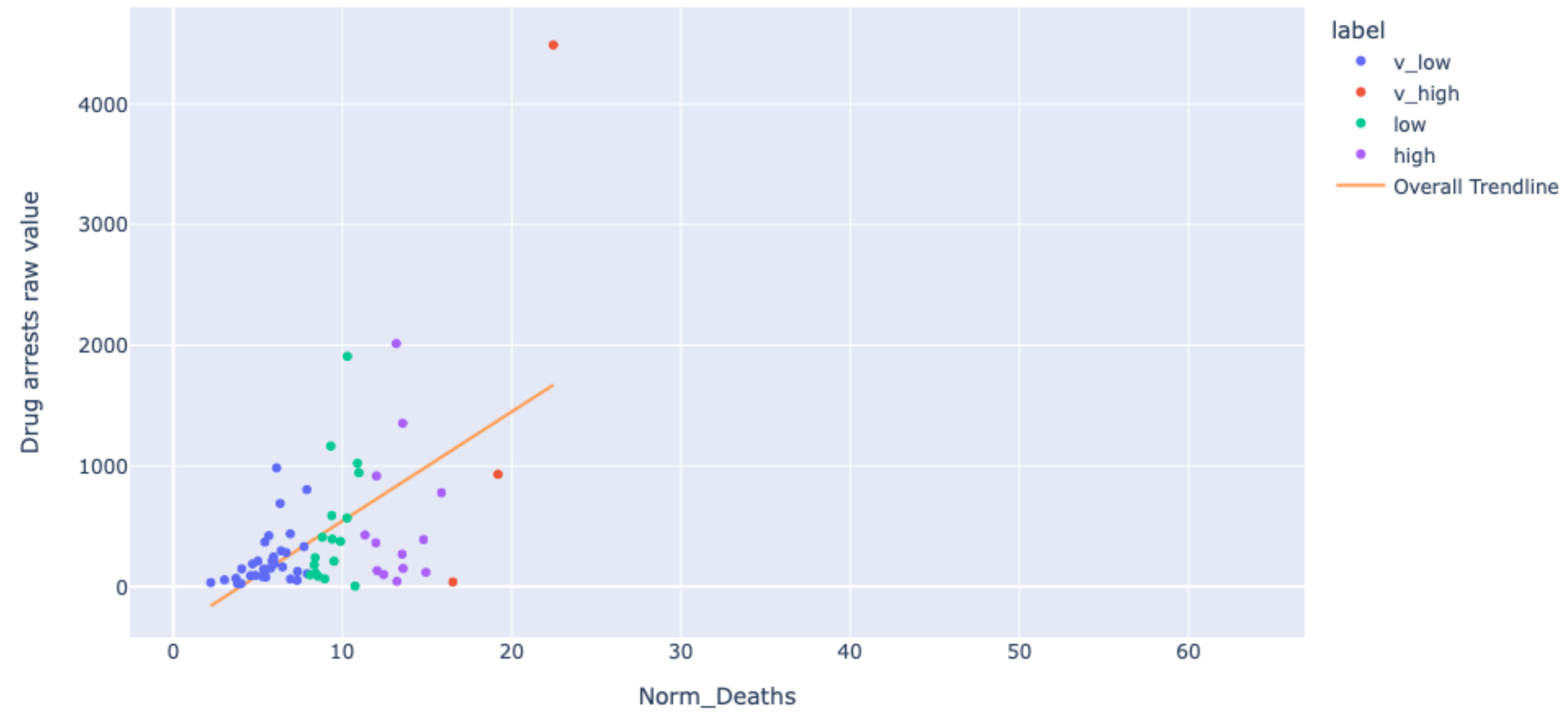
- So, I create separate scatter plots for each of them.
- In my first trial, I selected 10 columns that I believed correlated well with Norm\_Death column.
- When creating the scatter plot, I thought it'll be useful to add a trend\_line to my visual and color code observation based on the categories defined above. For trend\_line I used OLS method or ordinary least square approach. However, I am not sure if this is a right method or if there are other methods available.
- Only five out of ten variables that I have initially selected show strong visual correlation with the Norm\_Deaths variable. These are 'Drug arrest raw values', 'Communicable disease raw value', 'Total male population raw value', 'Hate crimes raw value' and 'Cancer incidence raw value' variables. I think I need to take a look at most columns and come up with a better selection!







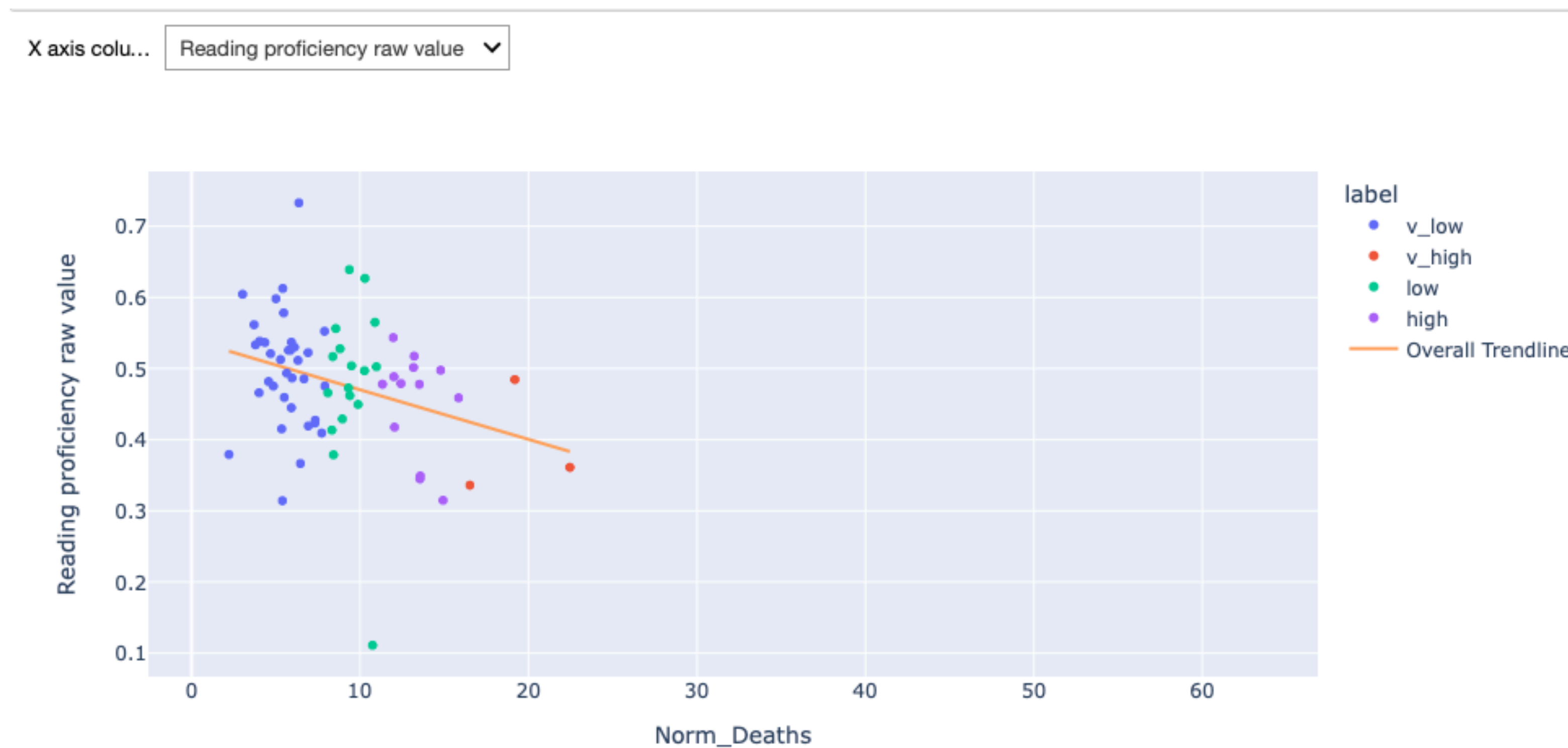
```
In [108]: px.scatter(super_df, y='Drug arrests raw value', x='Norm_Deaths', trendline='ols', color='label', trendline_scope='over
```



- The R square for the variables are listed in the table.

Variable	R square
Reading Proficiency	0.09
Life Expectancy	0.09
Premature Death	0.11
Older Adults	0.10
Injury Hospitalization	0.11
Cancer Incidence	0.21
Hate Crimes	0.25
Total Male Population	0.26
Communicable disease	0.31
Drug Arrest	0.30

- I used "Plotly widget" to develop a dashboard to visually see how these variables interact with Norm\_Death variable.
- So, I first created a list of these ten variables and named it as independent variables. Then, I created the widgets dropdown that can only takes values in this list. Then I followed the guid lines in provided link and added one line of code to show the scatter plot as well as the dropdown.





```

import plotly.graph_objects as go
from ipywidgets import widgets

independent_variables = ['Reading proficiency raw value', 'Life expectancy raw value', 'Older adults living alone raw',
, 'Injury hospitalizations raw value', 'Cancer incidence raw value', 'Hate crimes raw value', 'Total male population raw

X_Columns = widgets.Dropdown(
    options= independent_variables,
    description='X axis column',
)

g = go.FigureWidget(px.scatter(super_df, x=X_Columns.value, y="Norm_Deaths",trendline="ols",color='label',trendline_
    layout=go.Layout()))

def validate():
    if X_Columns.value in independent_variables:
        return True
    else:
        return False

def response(change):
    if validate():
        tmp_df = super_df
        with g.batch_update():
            g.data[0].x = tmp_df[X_Columns.value]
            g.layout.xaxis.title = X_Columns.value

X_Columns.observe(response, names="value")
container = widgets.HBox([X_Columns])
widgets.VBox([container,g])

```