

GT K CYBER

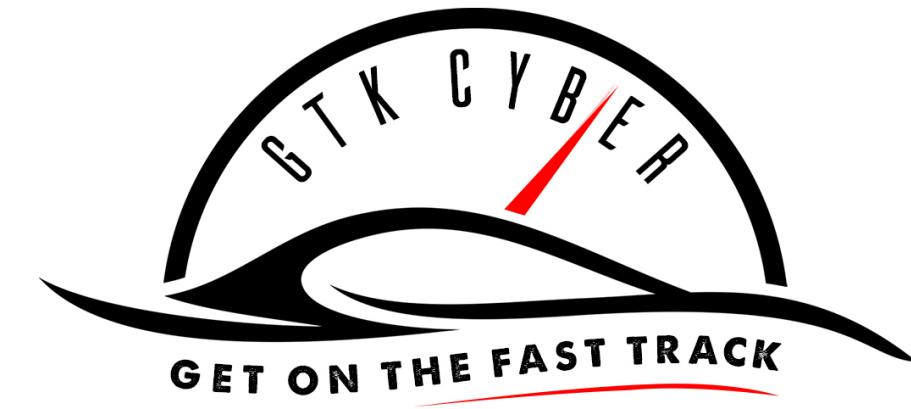
GET ON THE FAST TRACK

Data Science for Security Professionals - Day 2

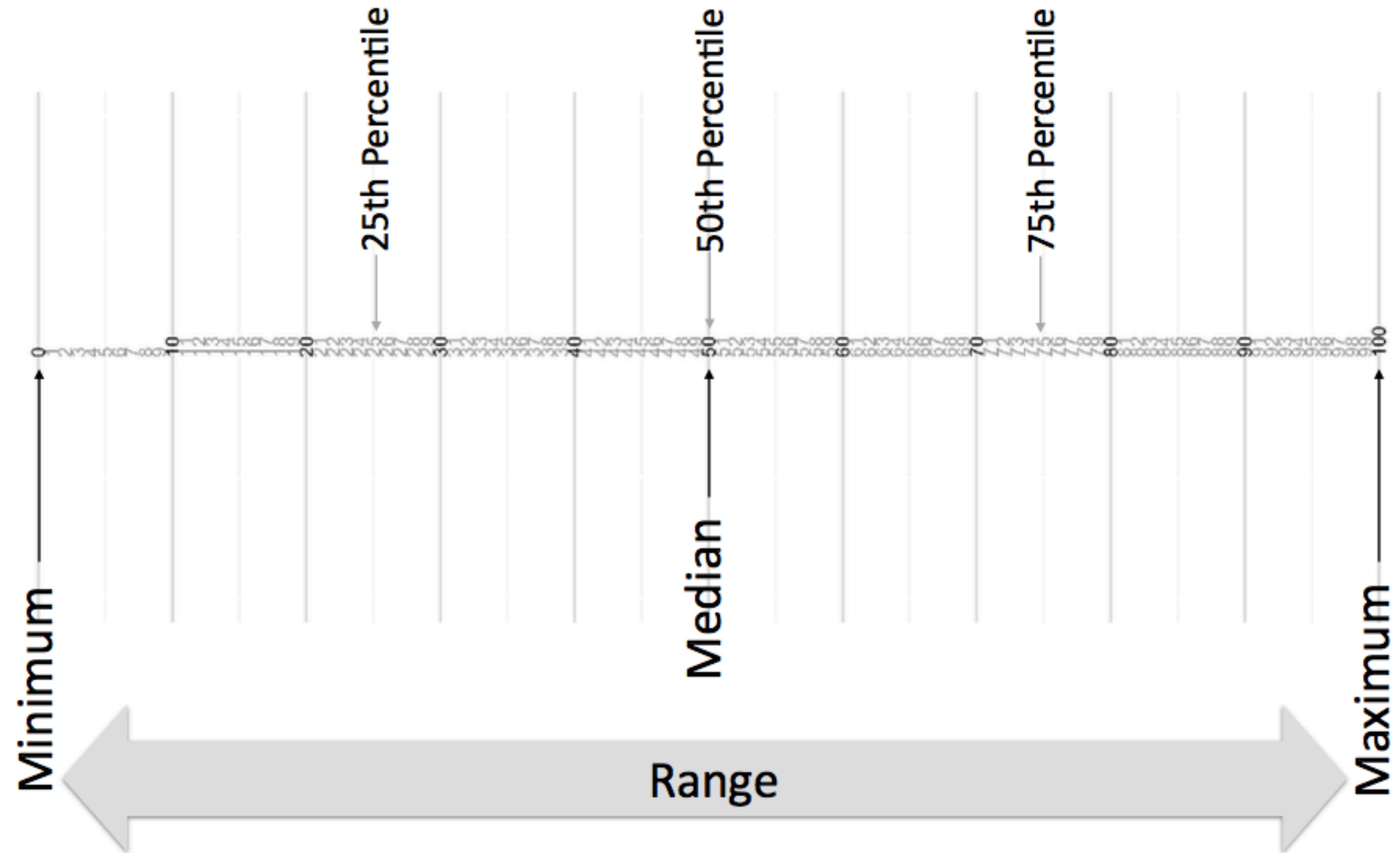
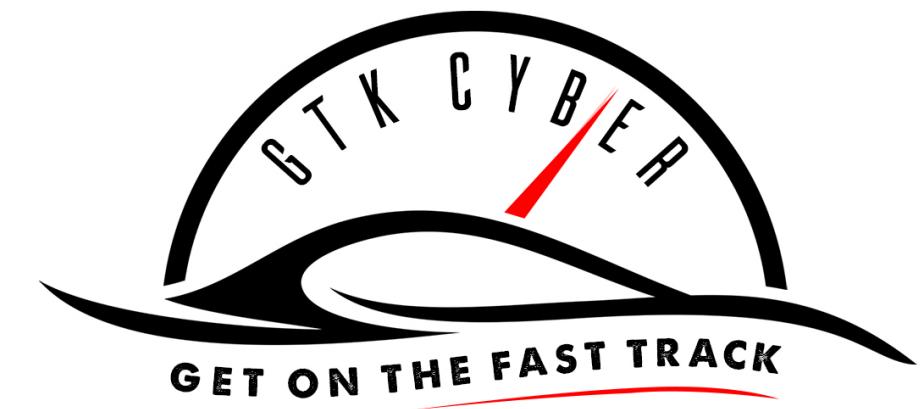
Statistical Summaries and Data Visualization

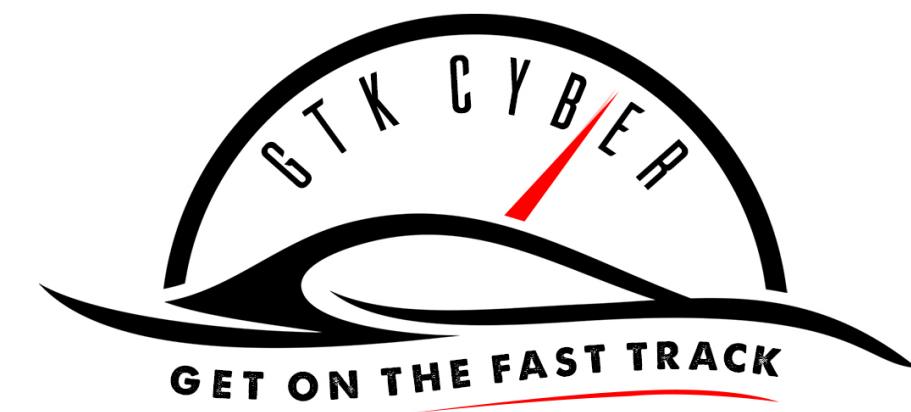
Explore Your Data



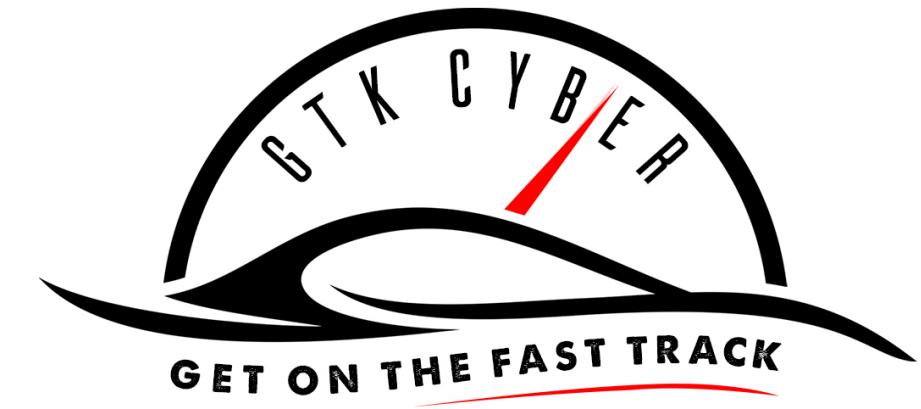


- **Mean:** The mean is the average value of a set of numbers
- **Median:** The median is the middle value of an ordered set of numbers
- **Mode:** The mode is the value from a set which is repeated the most frequently.
- **Range:** Difference between minimum and maximum
- **Standard Deviation:** Measures dispersion in a data set. Close to zero indicates little dispersion.





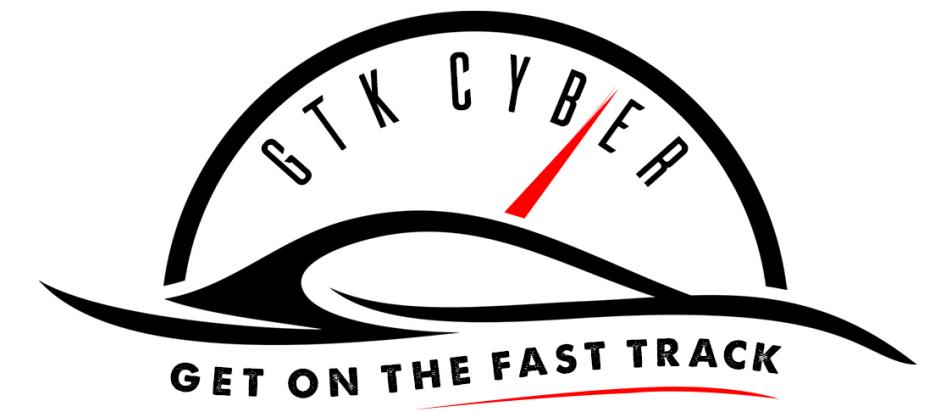
Series.abs()	Absolute Value of the Series
Series.count()	Returns number of non-empty values in the series
Series.max()	Returns maximum value in the Series
Series.mean()	Returns the mean of a Series
Series.median()	Returns the median of a Series
Series.min()	Returns the minimum value in a Series
Series.mode()	Take a guess..
Series.quantile([q])	Returns the quantiles of a Series
Series.sum	Returns the sum of a series
Series.std	Returns the standard deviation of a Series



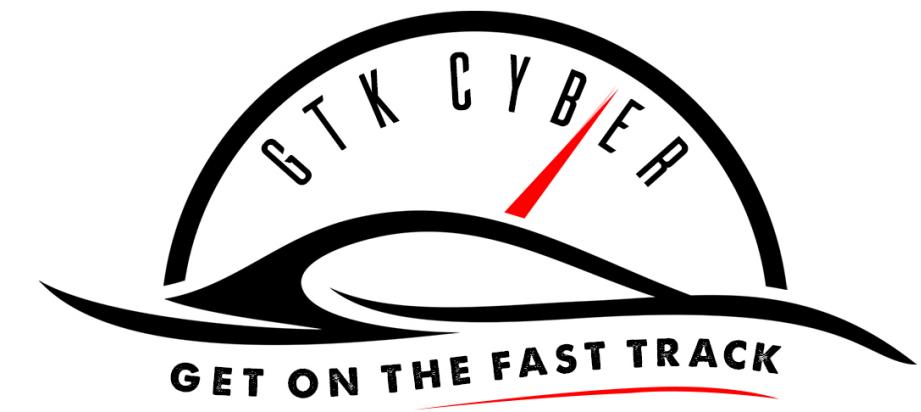
Tukey 5 Number Summary

- **Minimum:** The smallest value in the dataset
- **Lower Quartile:** Smallest 25% of the dataset
- **The Median:** The middle value of the dataset
- **Upper Quartile:** The largest 25% of the dataset
- **Maximum:** The largest value in the dataset



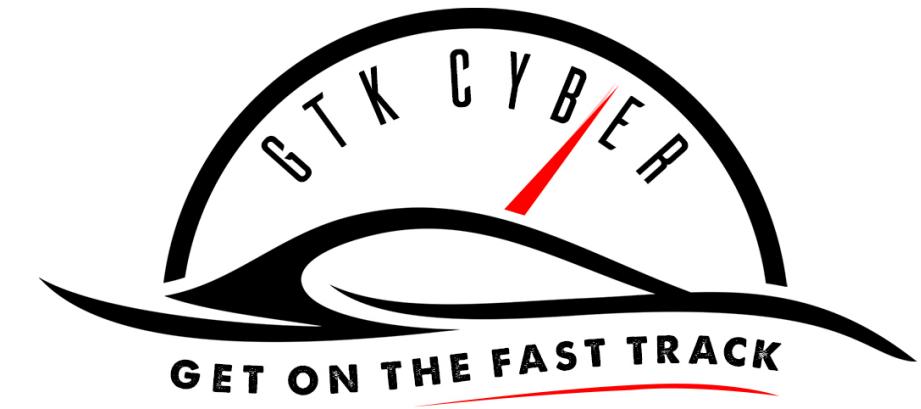


Series.describe()



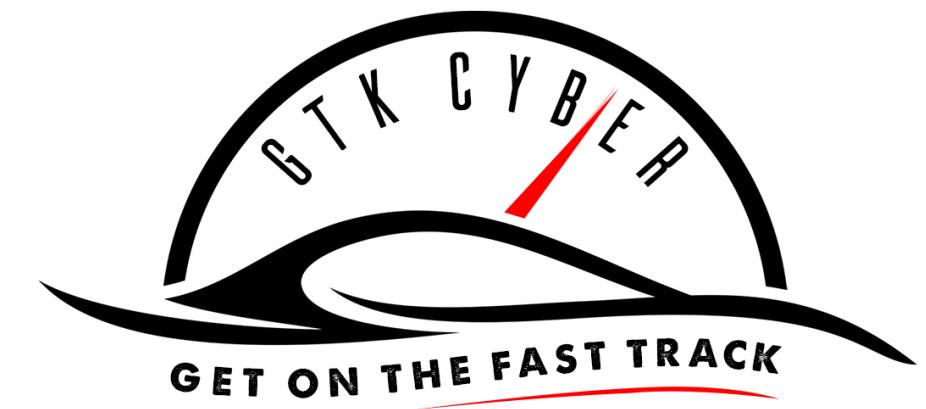
Series.describe()

```
>>> random_numbers.describe()
count      50.000000
mean      50.620000
std       30.102471
min       1.000000
25%      25.500000
50%      54.000000
75%      73.000000
max      99.000000
dtype: float64
```



```
names = pd.Series(  
    [ 'Jim', 'Bob',  
    'Bob', 'Steve', 'Jim', 'Jane', 'Steph', 'Emma', 'Rachel  
    ' ] )
```

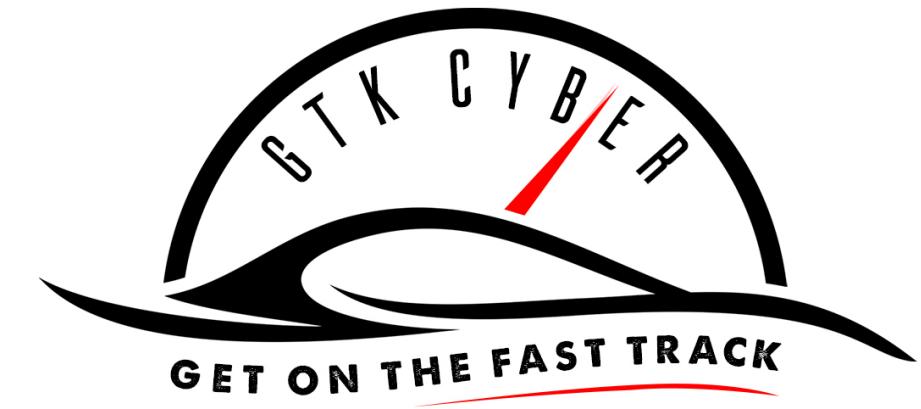
names.**describe()**



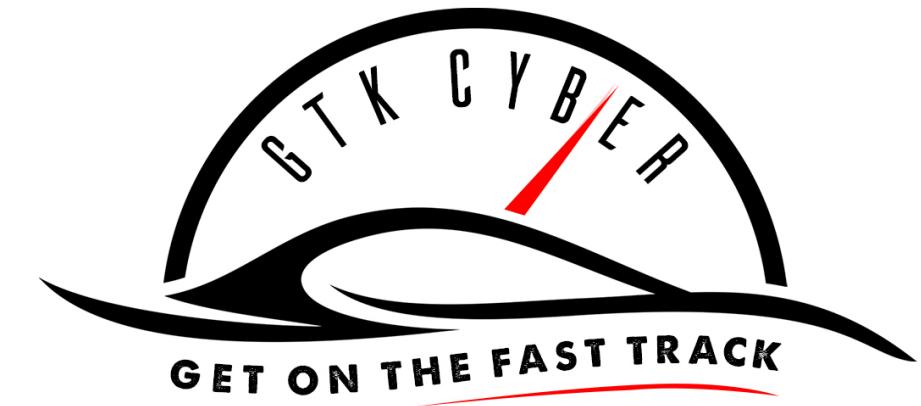
```
names = pd.Series(  
    [ 'Jim', 'Bob',  
    'Bob', 'Steve', 'Jim', 'Jane', 'Steph', 'Emma', 'Rachel  
    ' ] )
```

```
names.describe()
```

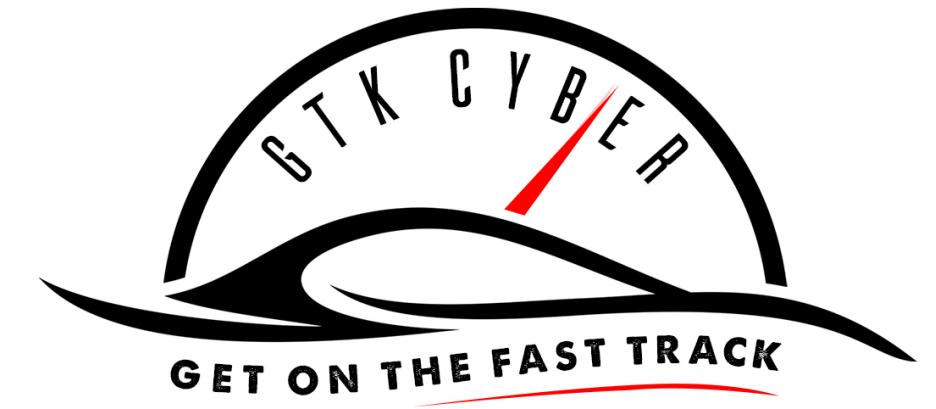
```
count          9  
unique         7  
top           Jim  
freq           2  
dtype: object
```



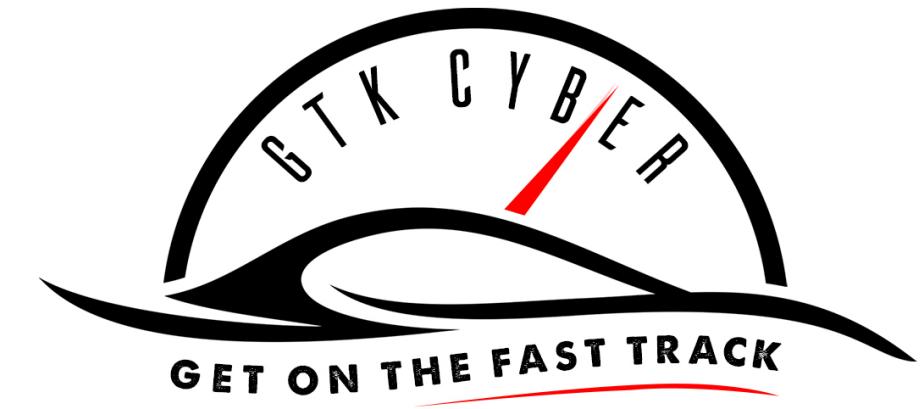
Finding Unique Values



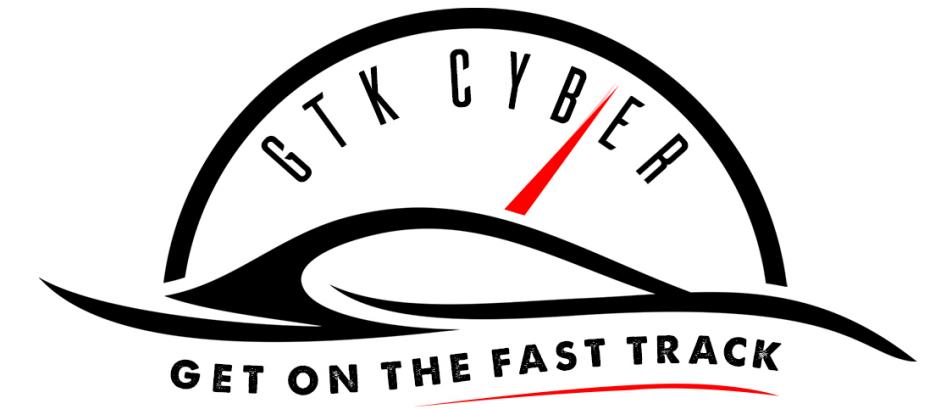
```
unique_nums = []
for key, value in random_numbers.iteritems():
    if value not in unique_nums:
        unique_nums.append(value)
```



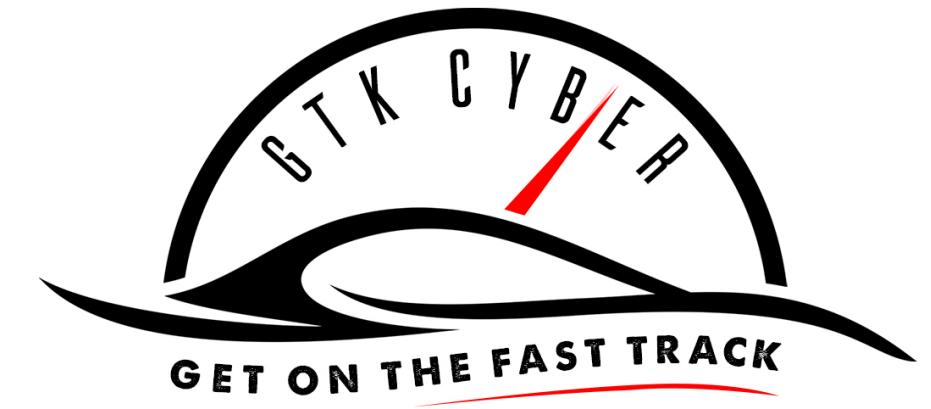
NO!!!



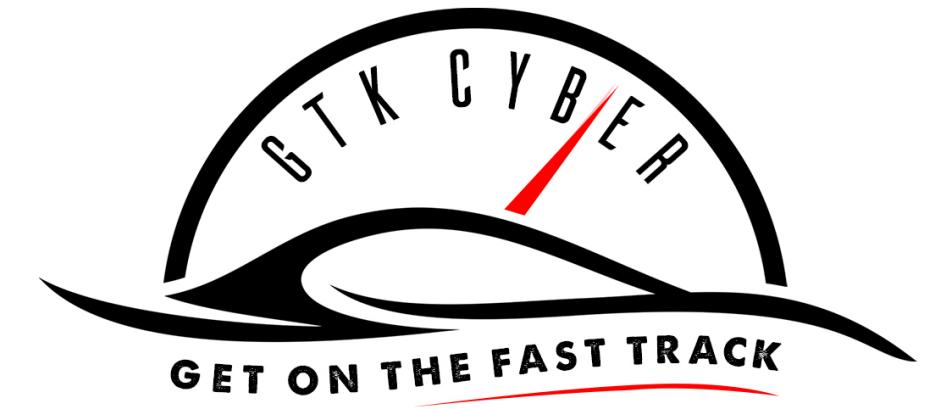
`series.drop_duplicates()`
`series.unique()`



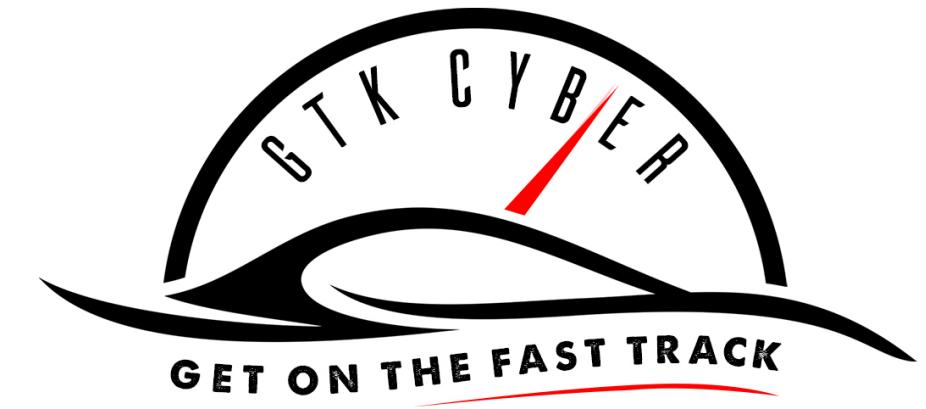
Finding Unique Values



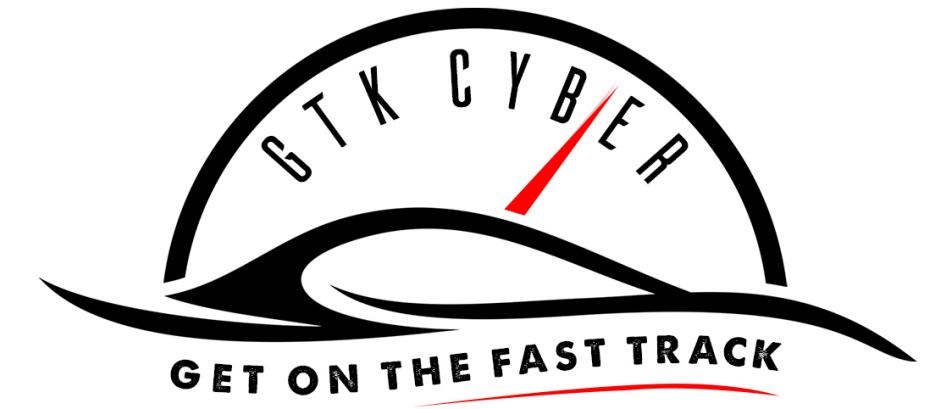
`series.value_counts()`



`series.value_counts(bins=5)`



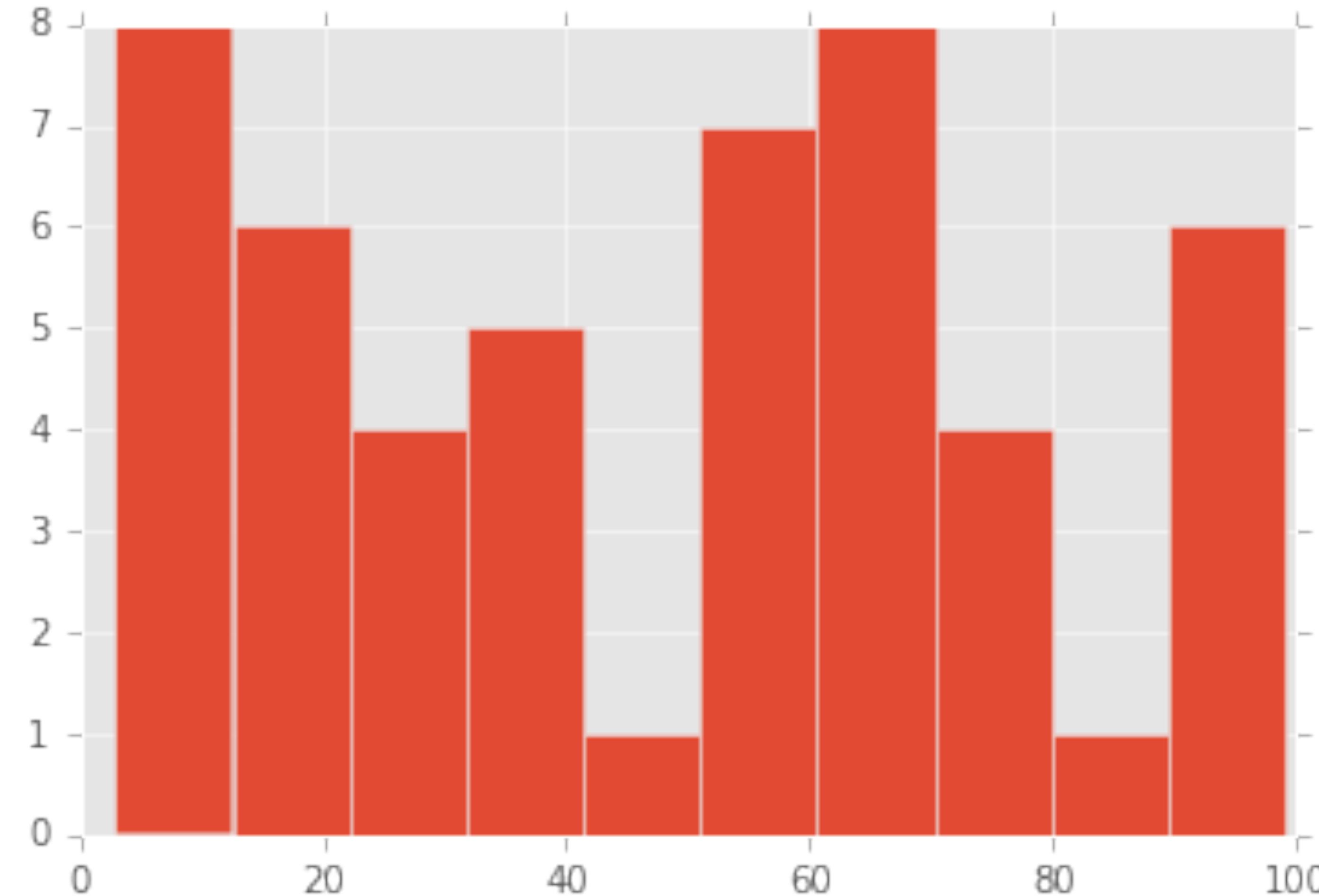
```
series.value_counts(bins=5,  
normalize=True)
```

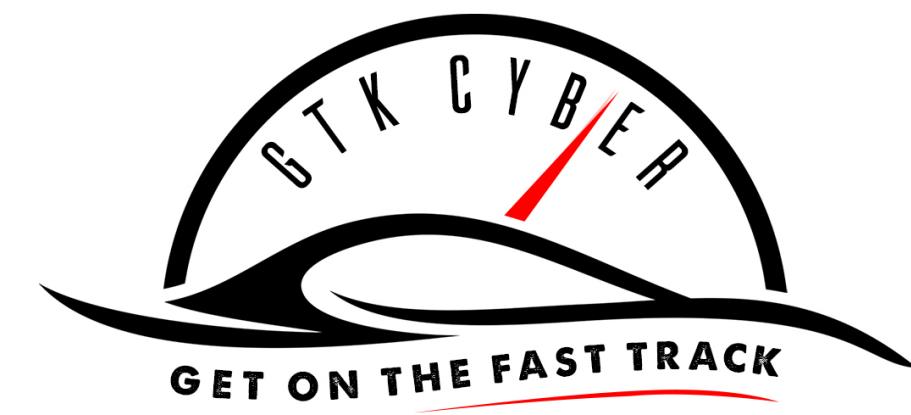


series.hist()

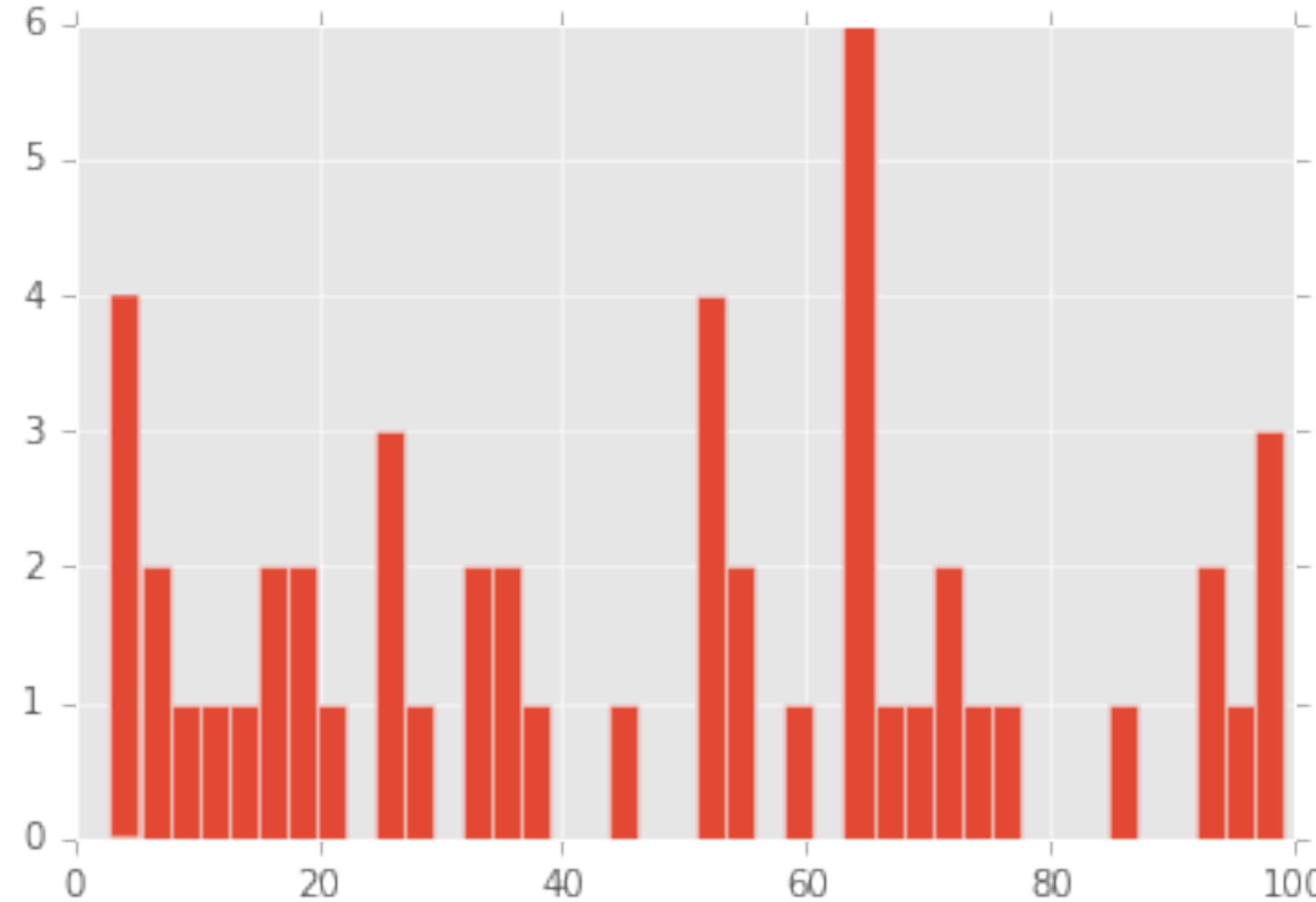


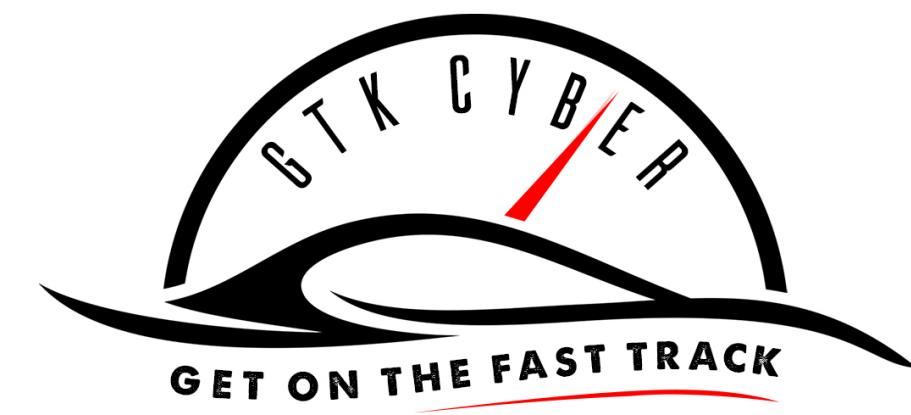
series.hist()



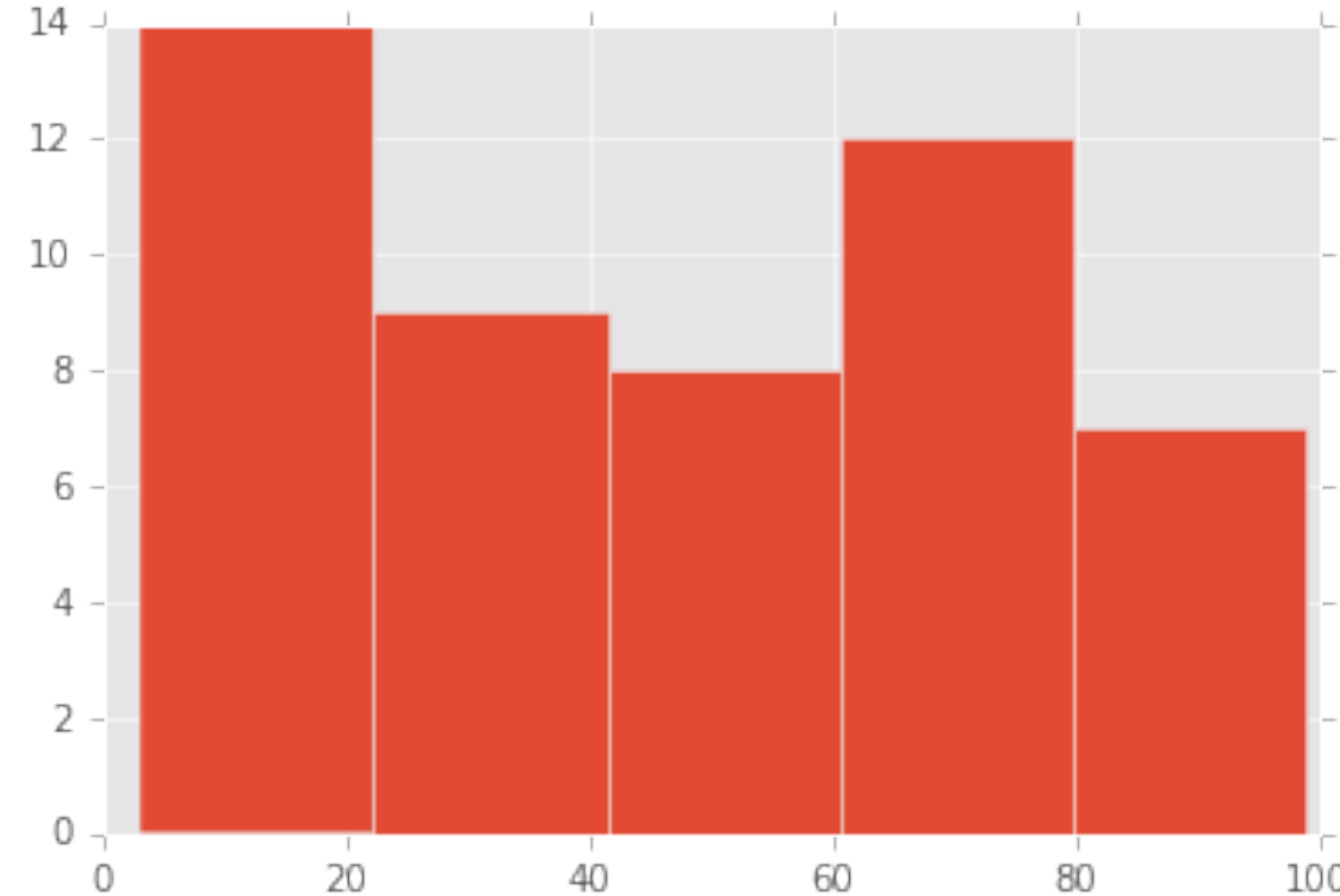


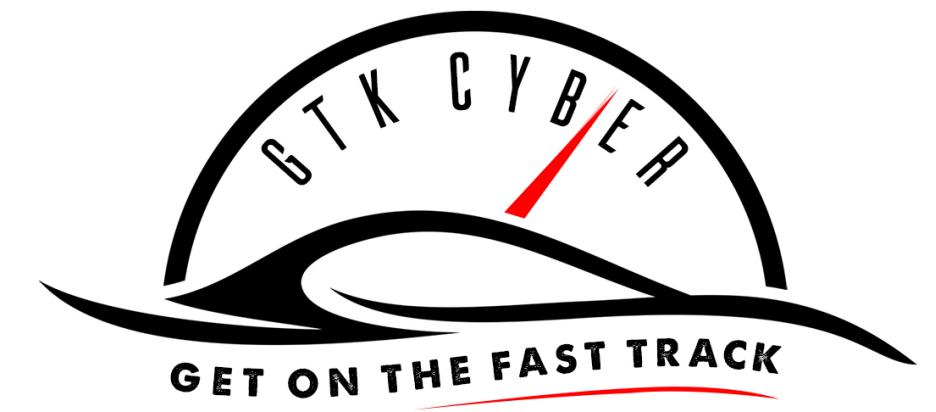
series.hist(bins=40)



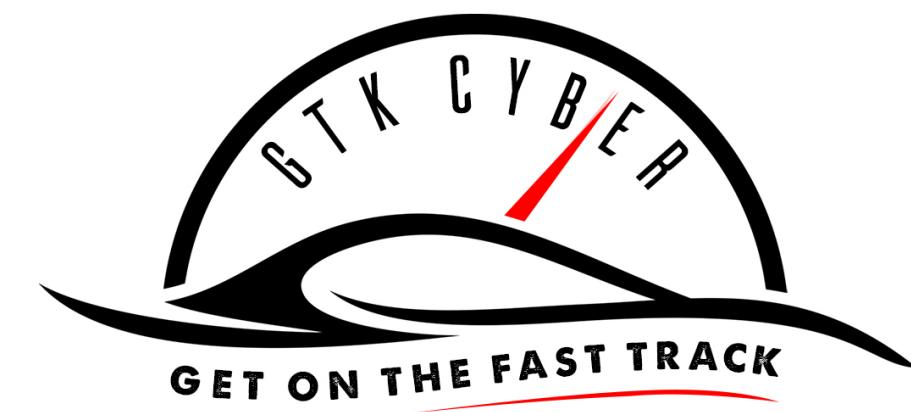


series.hist(bins=5)

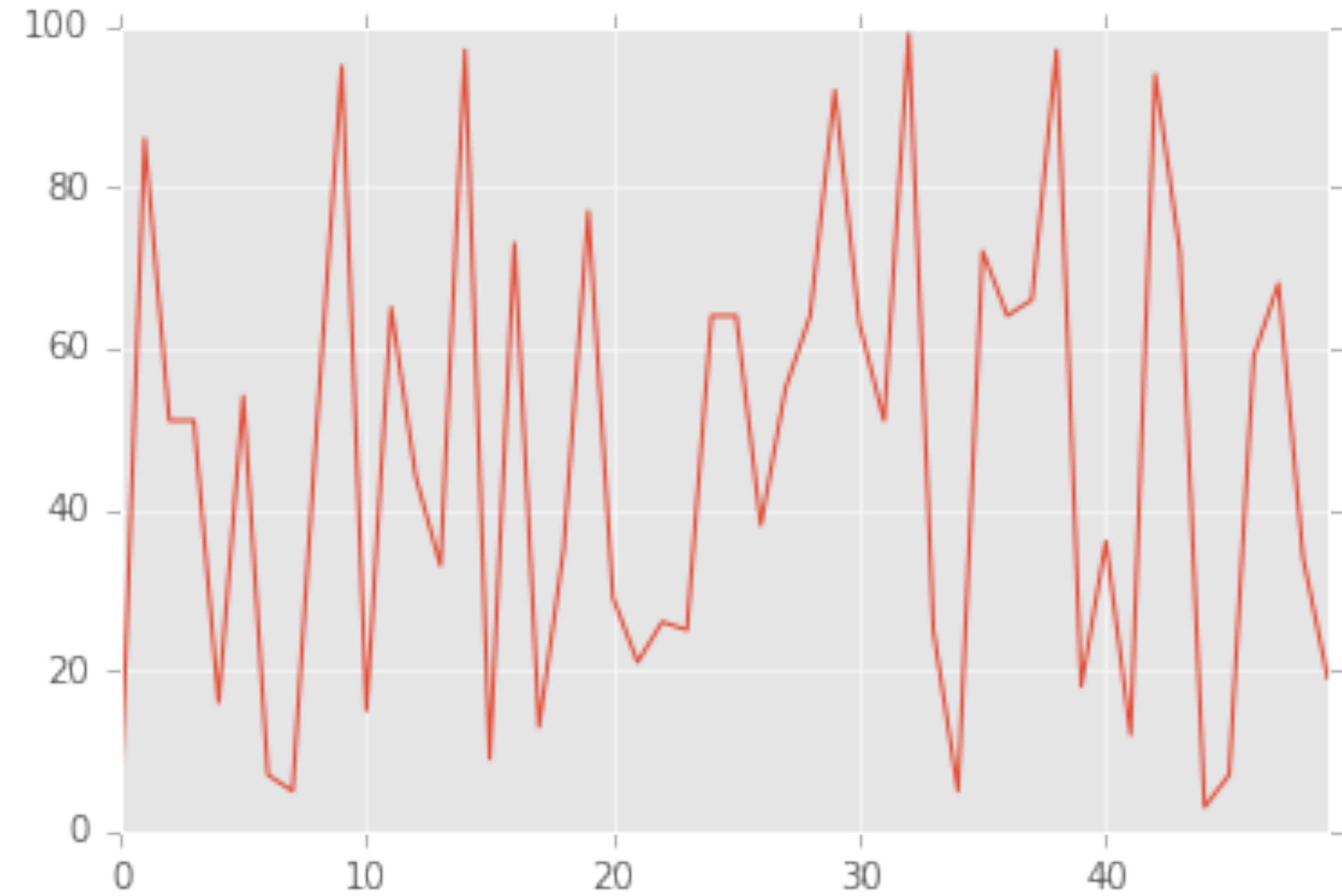




series.plot()

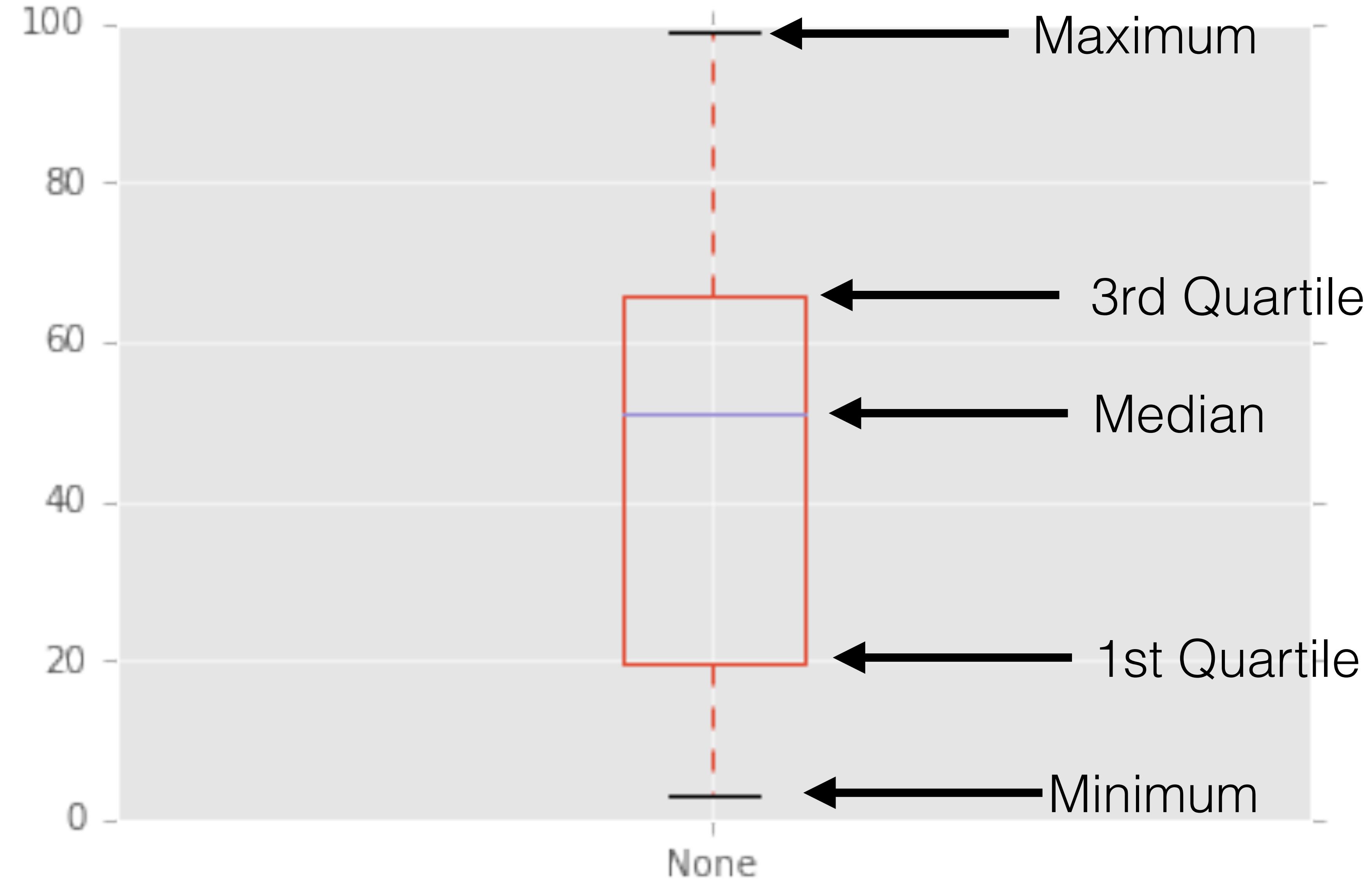


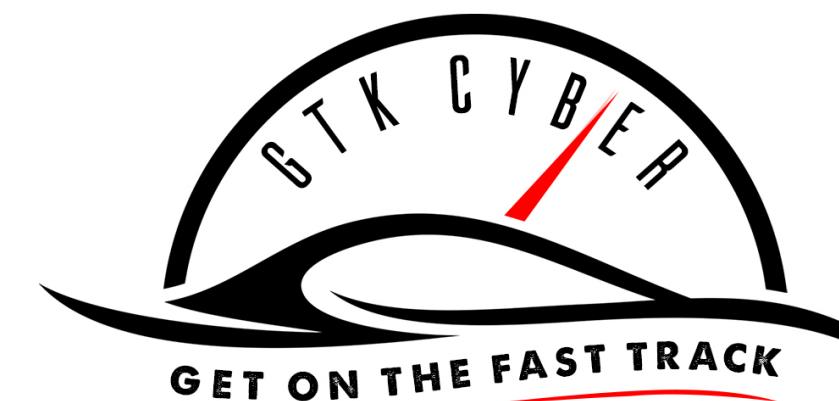
series.plot()





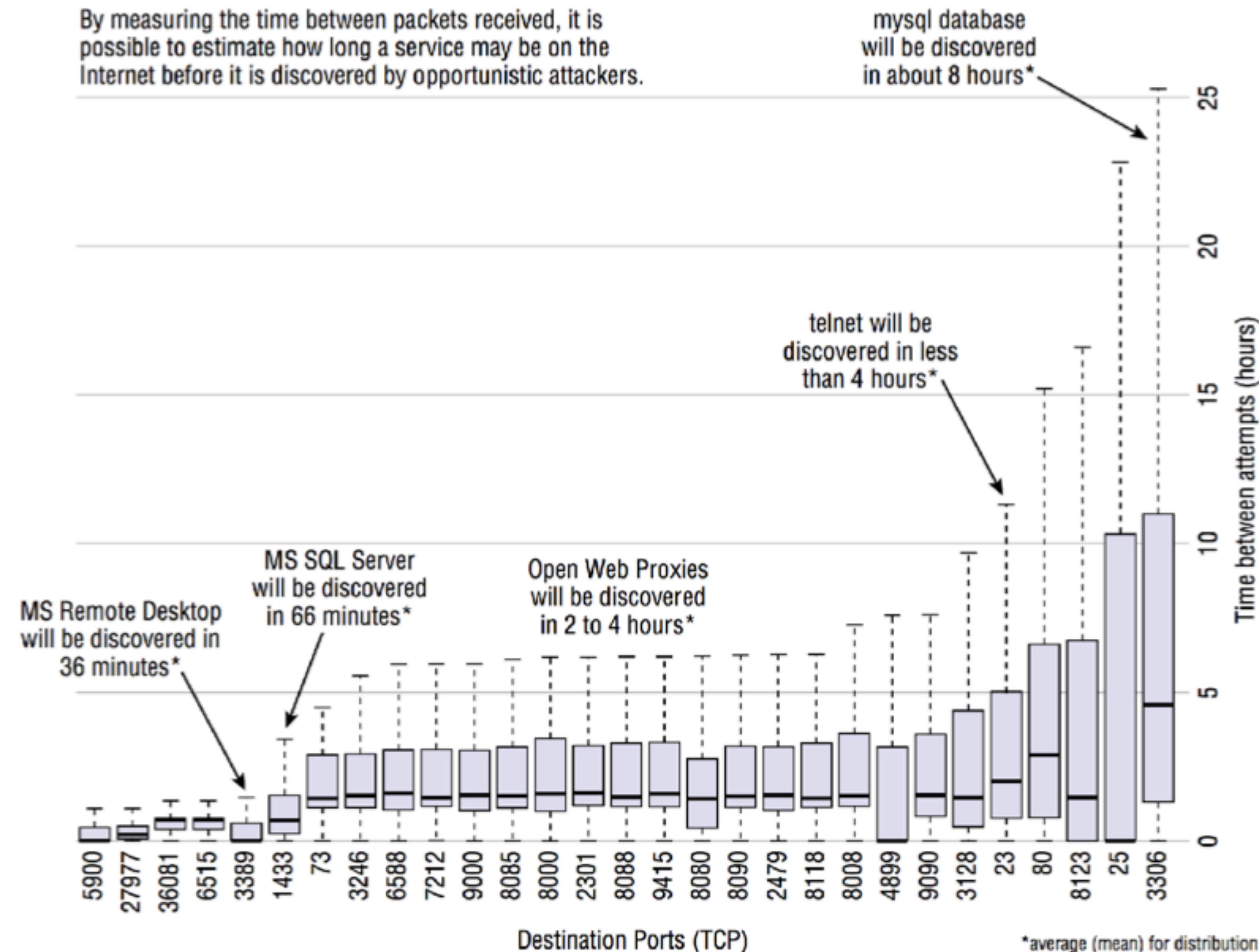
series.plot(kind='box')

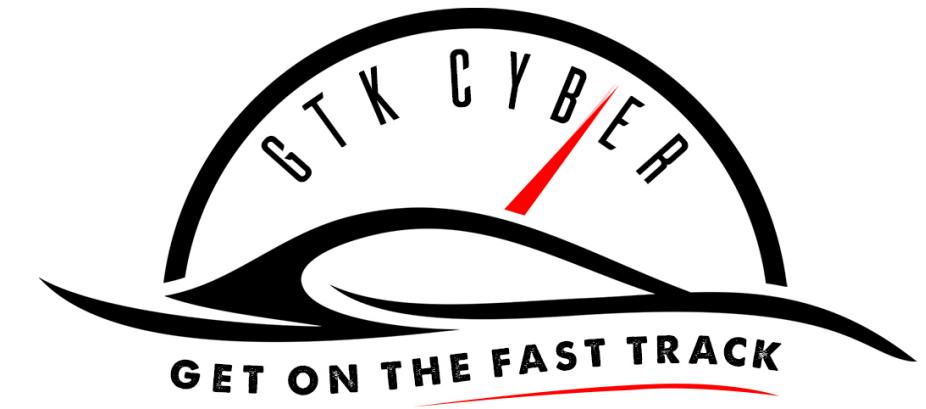




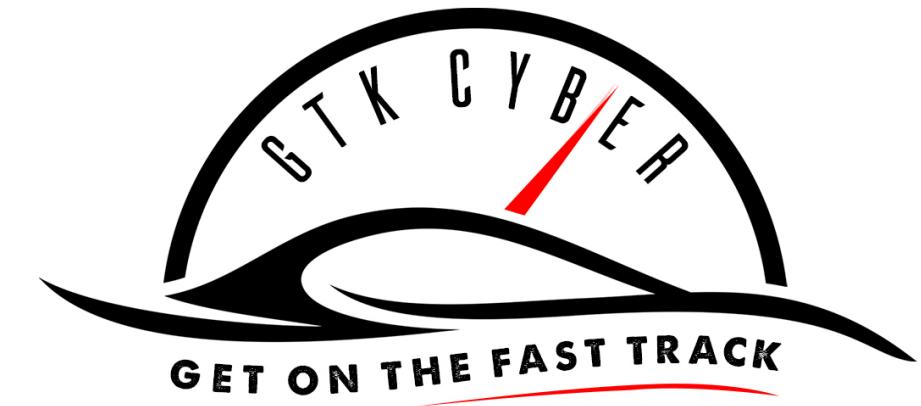
How long will a service go undiscovered by opportunistic attackers?

By measuring the time between packets received, it is possible to estimate how long a service may be on the Internet before it is discovered by opportunistic attackers.





Correlations between Datasets

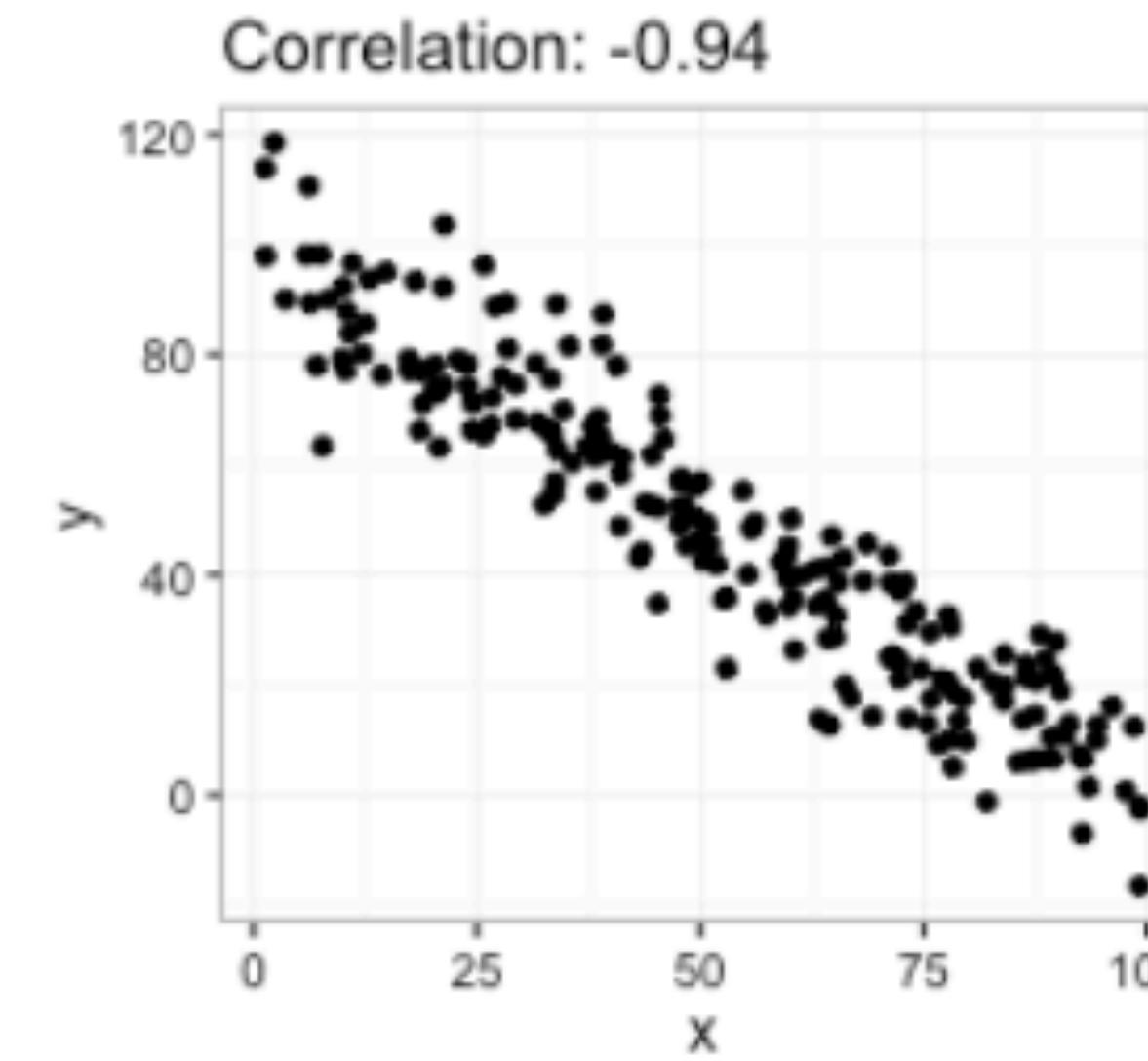
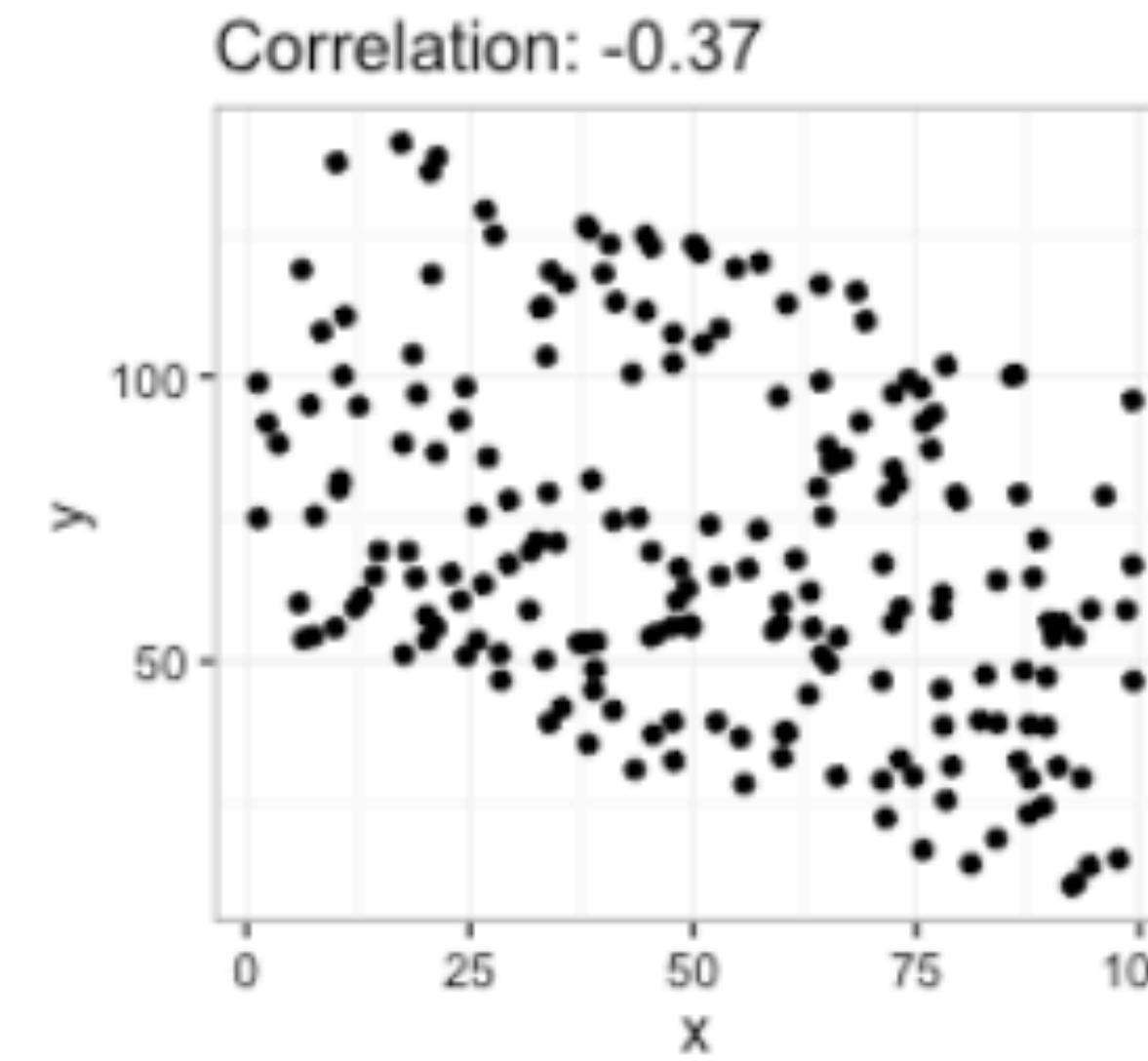
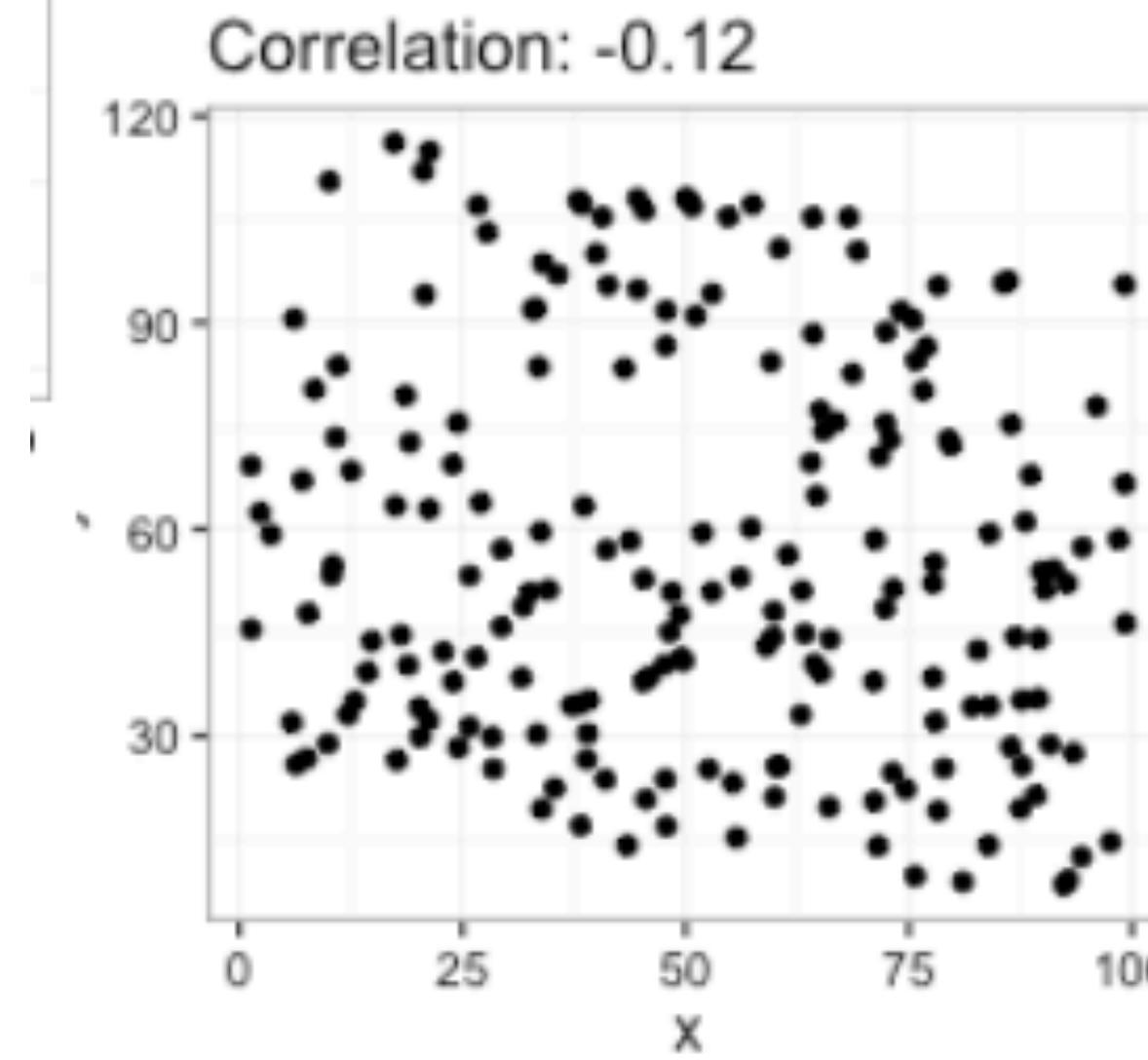
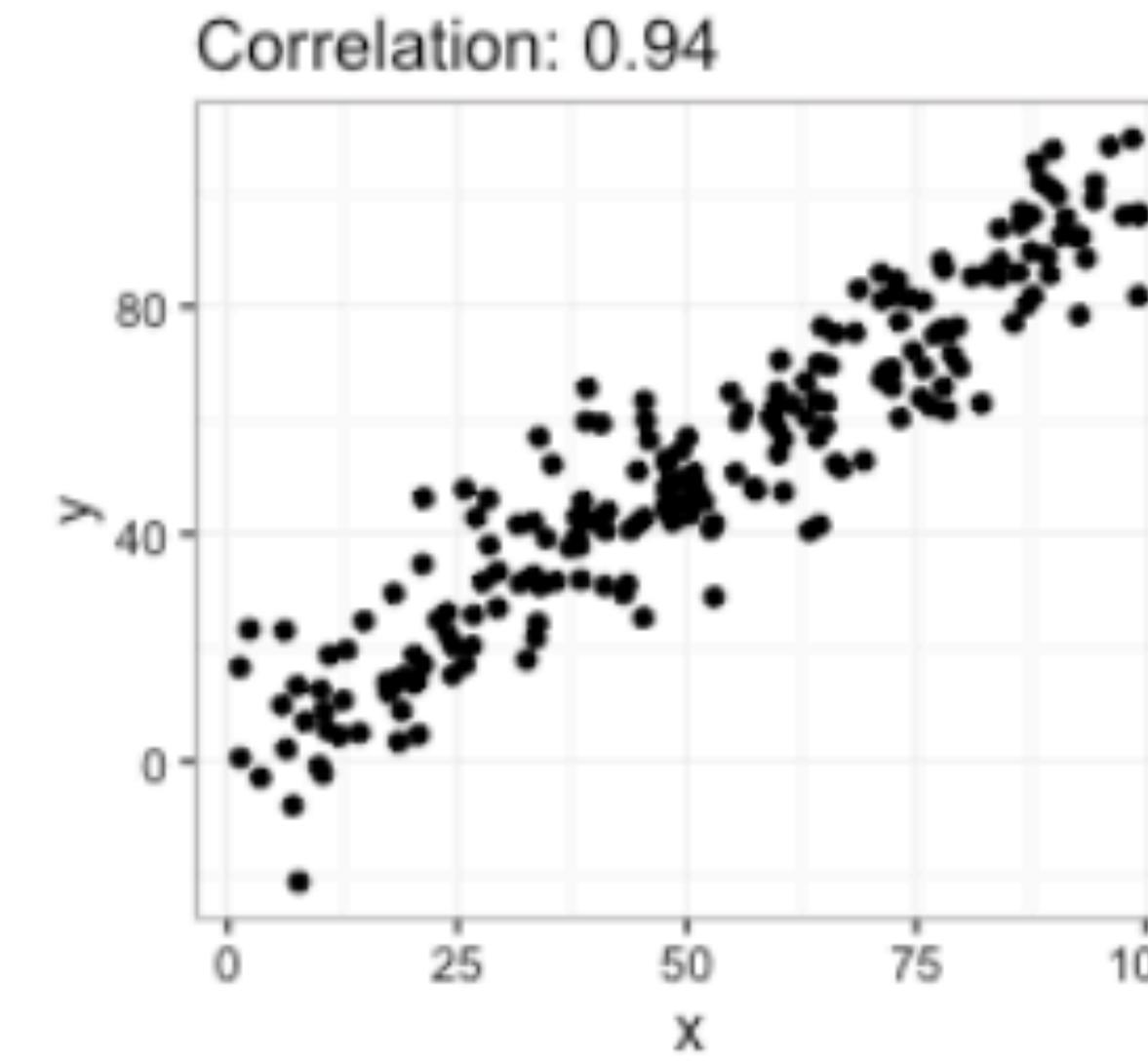
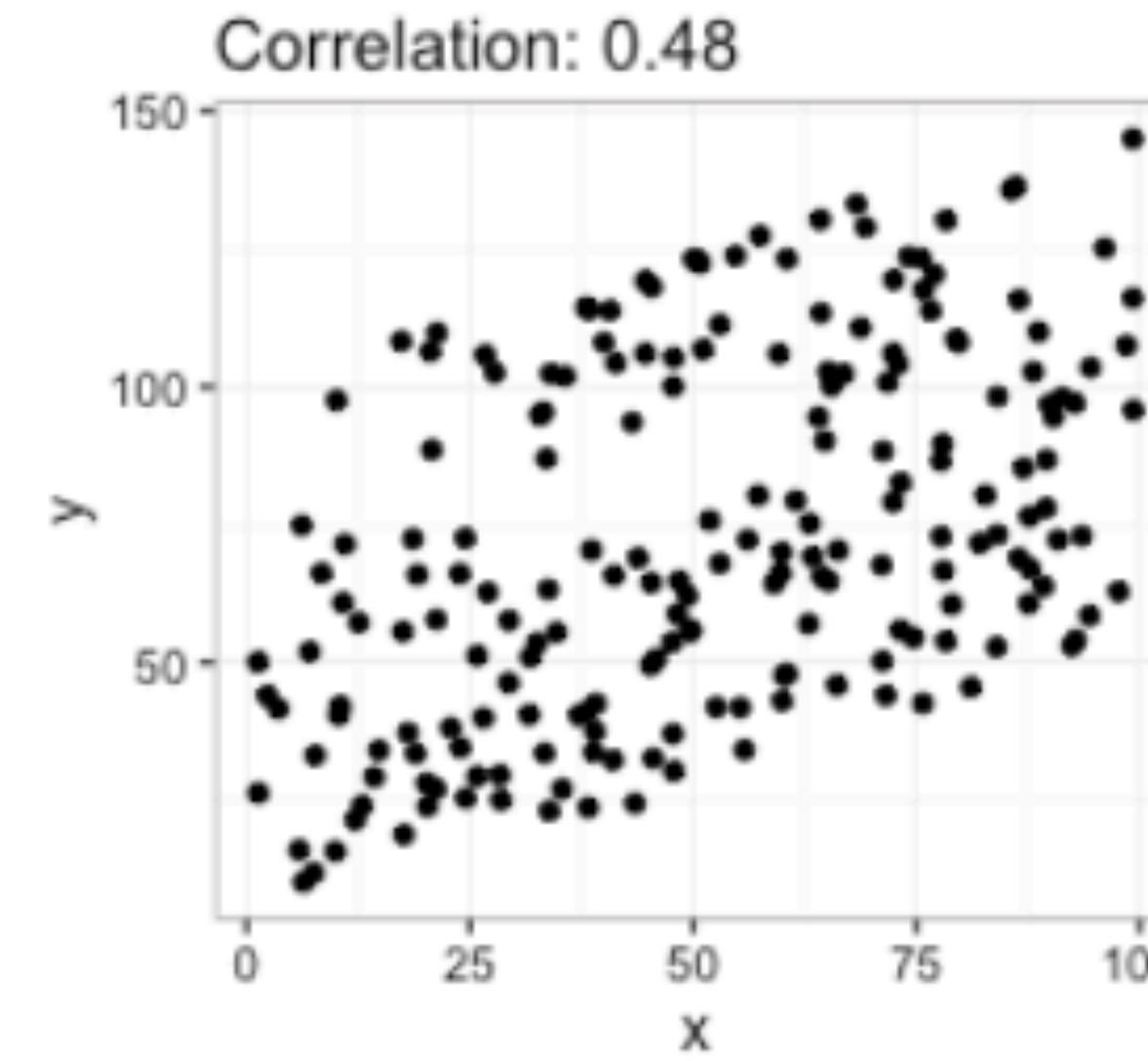
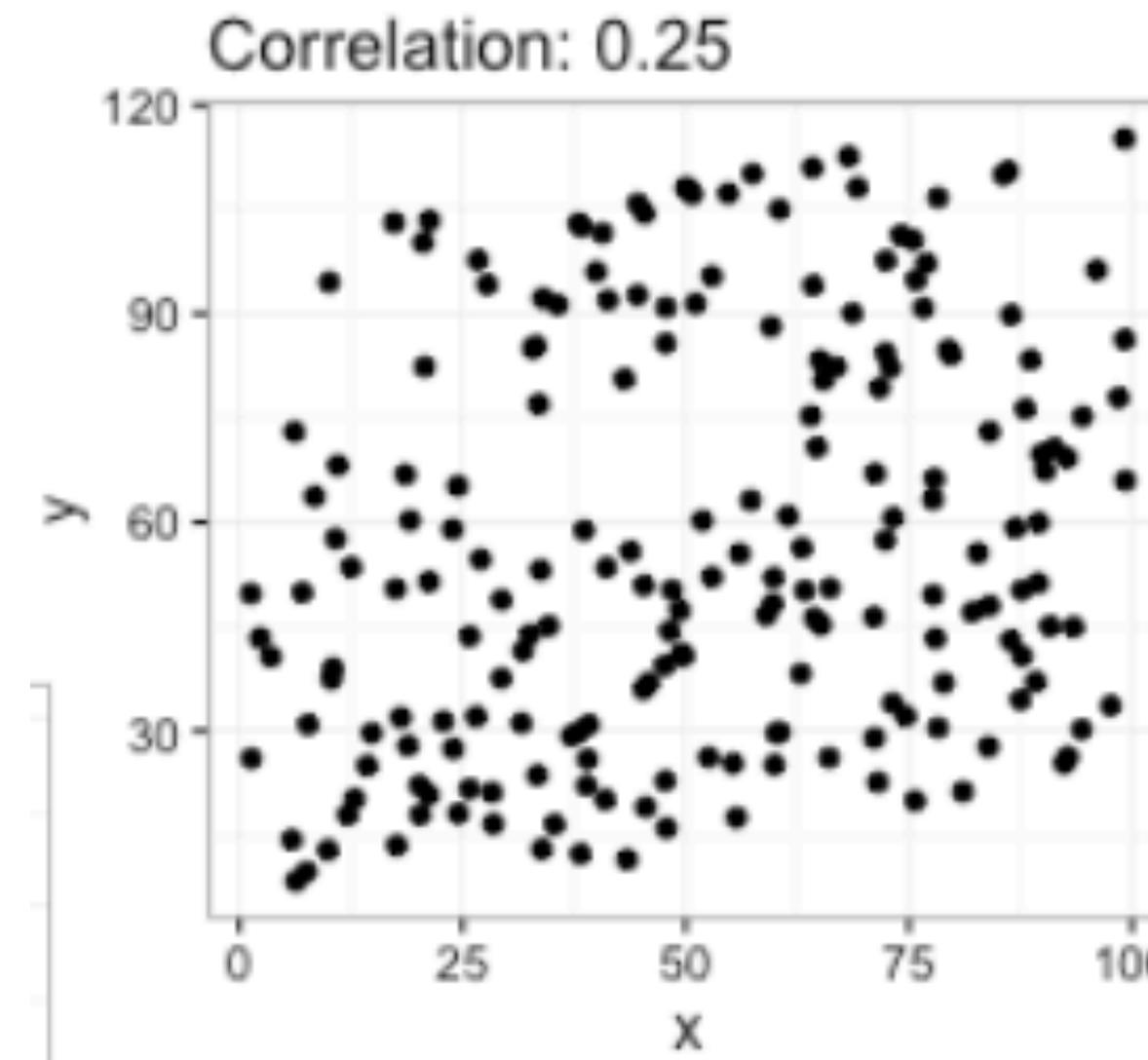


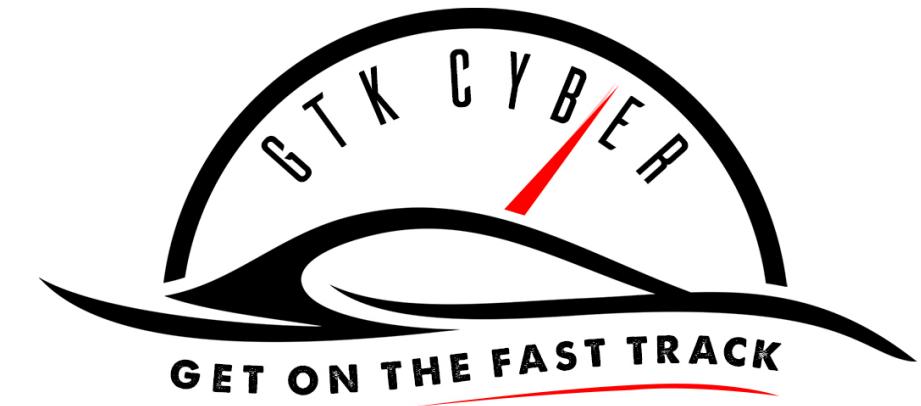
Correlations

- Measurement of the relationship between **two continuous variables**
- Output is between -1 and 1, with 1 being perfect correlation, -1 is perfect negative correlation, 0 is no correlation.
- Correlation measurement is referred to as r .



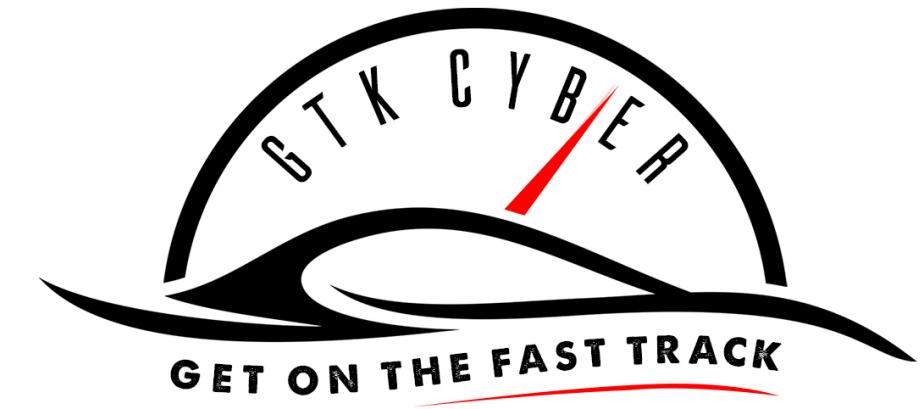
Correlations





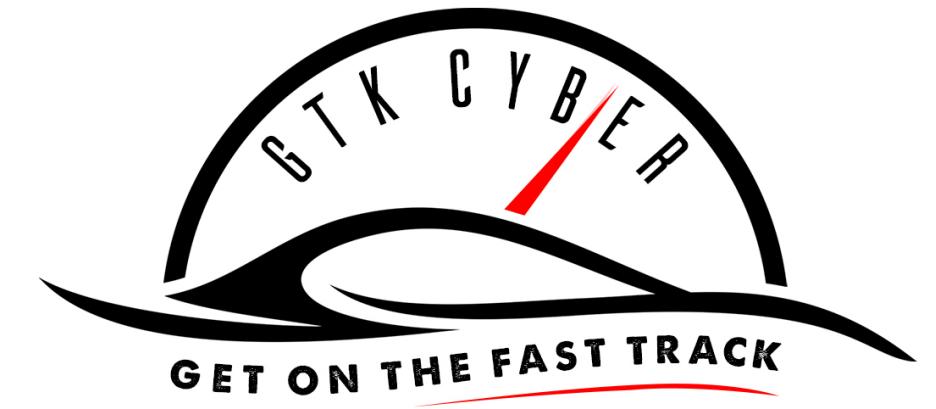
Correlations

```
series1.corr( series2 )
```

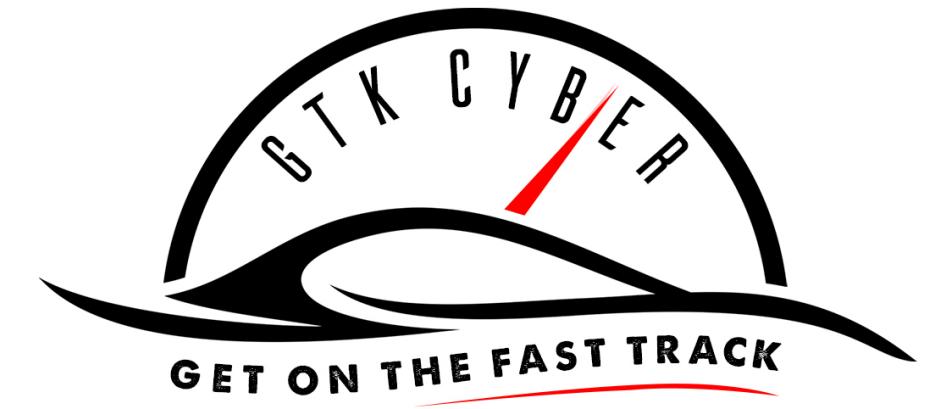


In Class Exercise:

Please complete Worksheet 3: EDA Worksheet



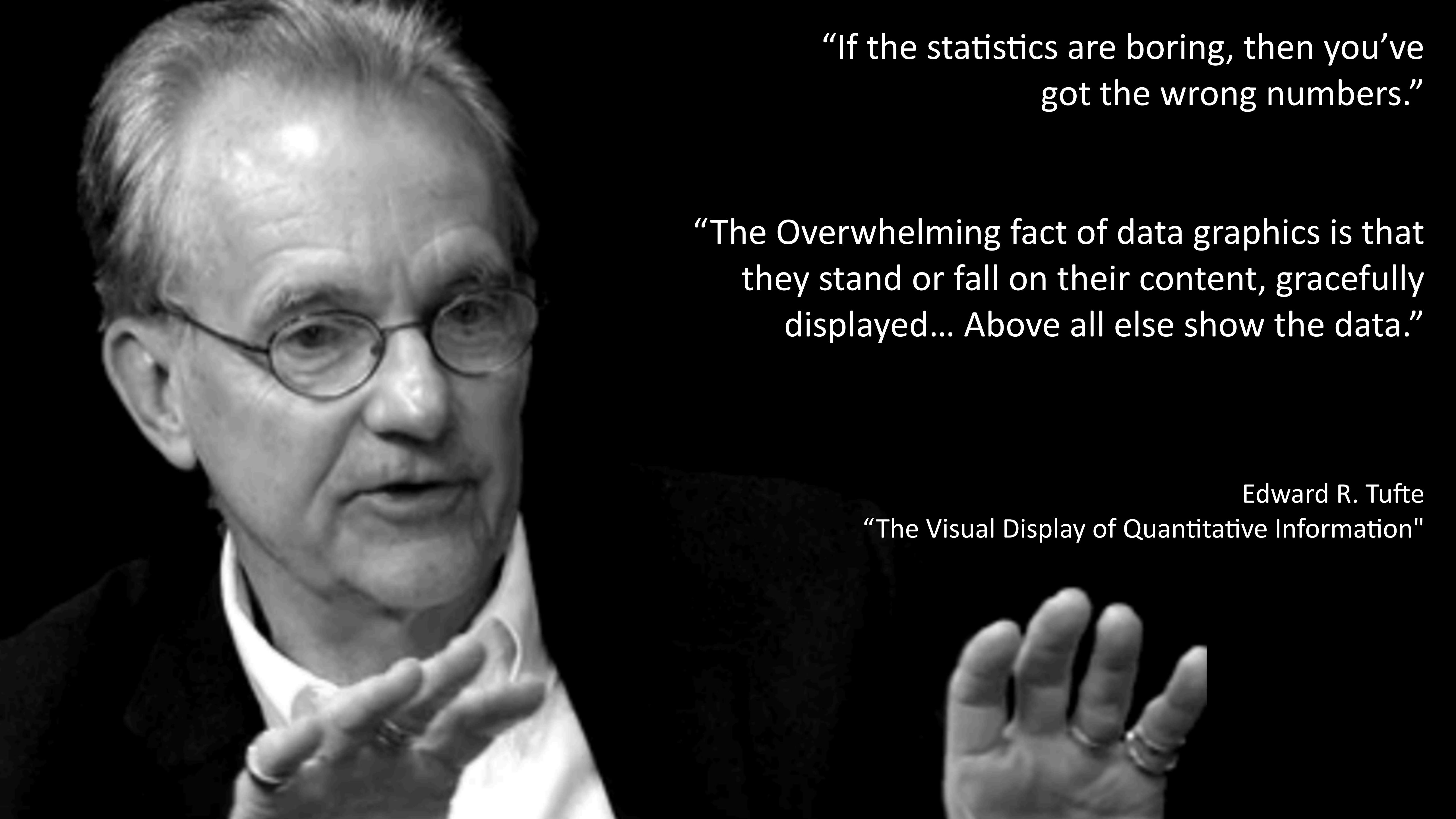
Questions?



Exploring Data Through Visualization

Visualization Goals

- Analyze
- Explore
- Assess
- Determine
- Decide
- Communicate
- Explain
- Present
- Prove
- Persuade

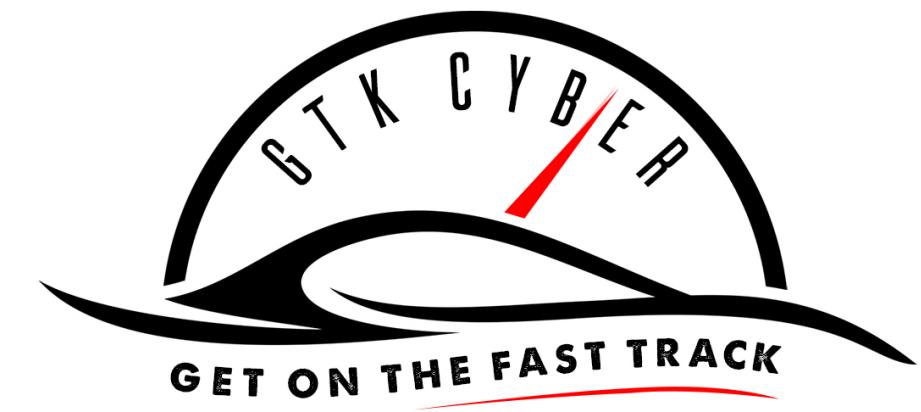


“If the statistics are boring, then you’ve got the wrong numbers.”

“The Overwhelming fact of data graphics is that they stand or fall on their content, gracefully displayed... Above all else show the data.”

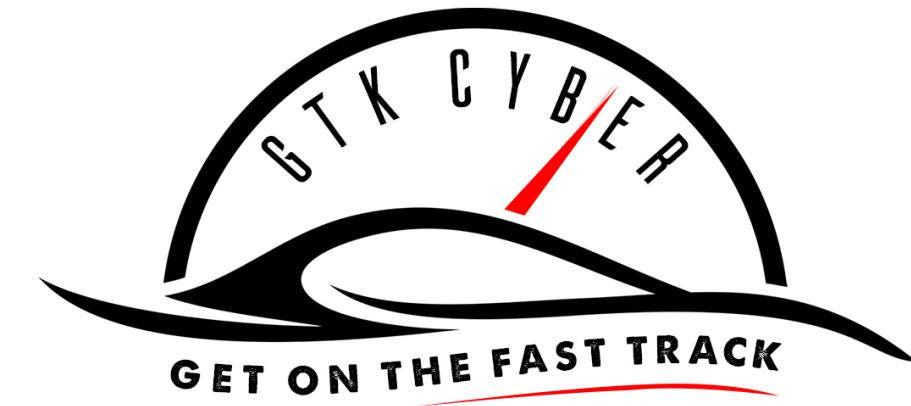
Edward R. Tufte

“The Visual Display of Quantitative Information”



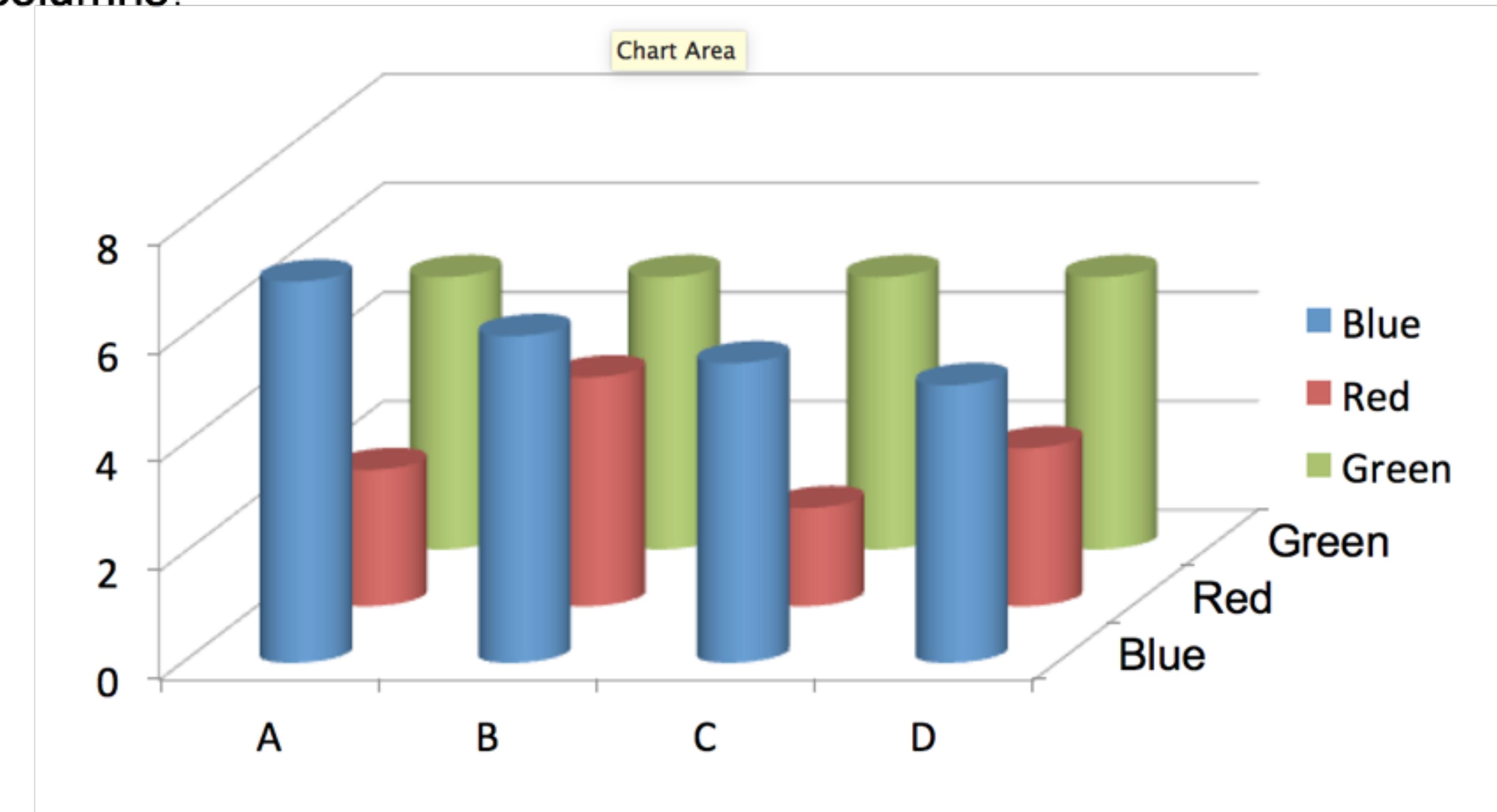
Elements of Good Visualizations

1. Graphical Integrity
2. Simple
3. Proper Display
4. Proper Color
5. Tells a story



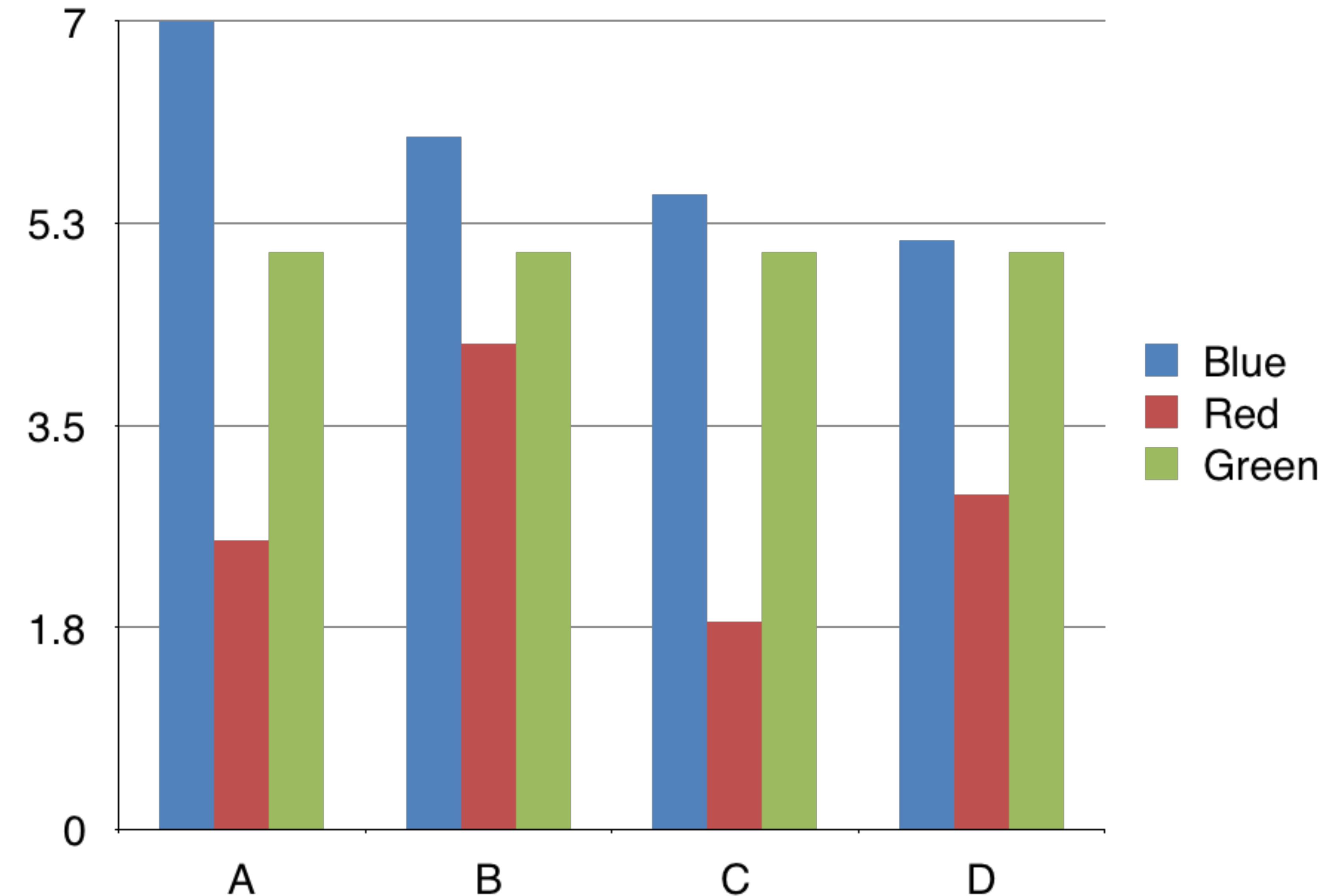
Proper Display

Questions: What is the height of the green columns? For which categories (A,B,C,D) are the blue columns taller than the green columns?



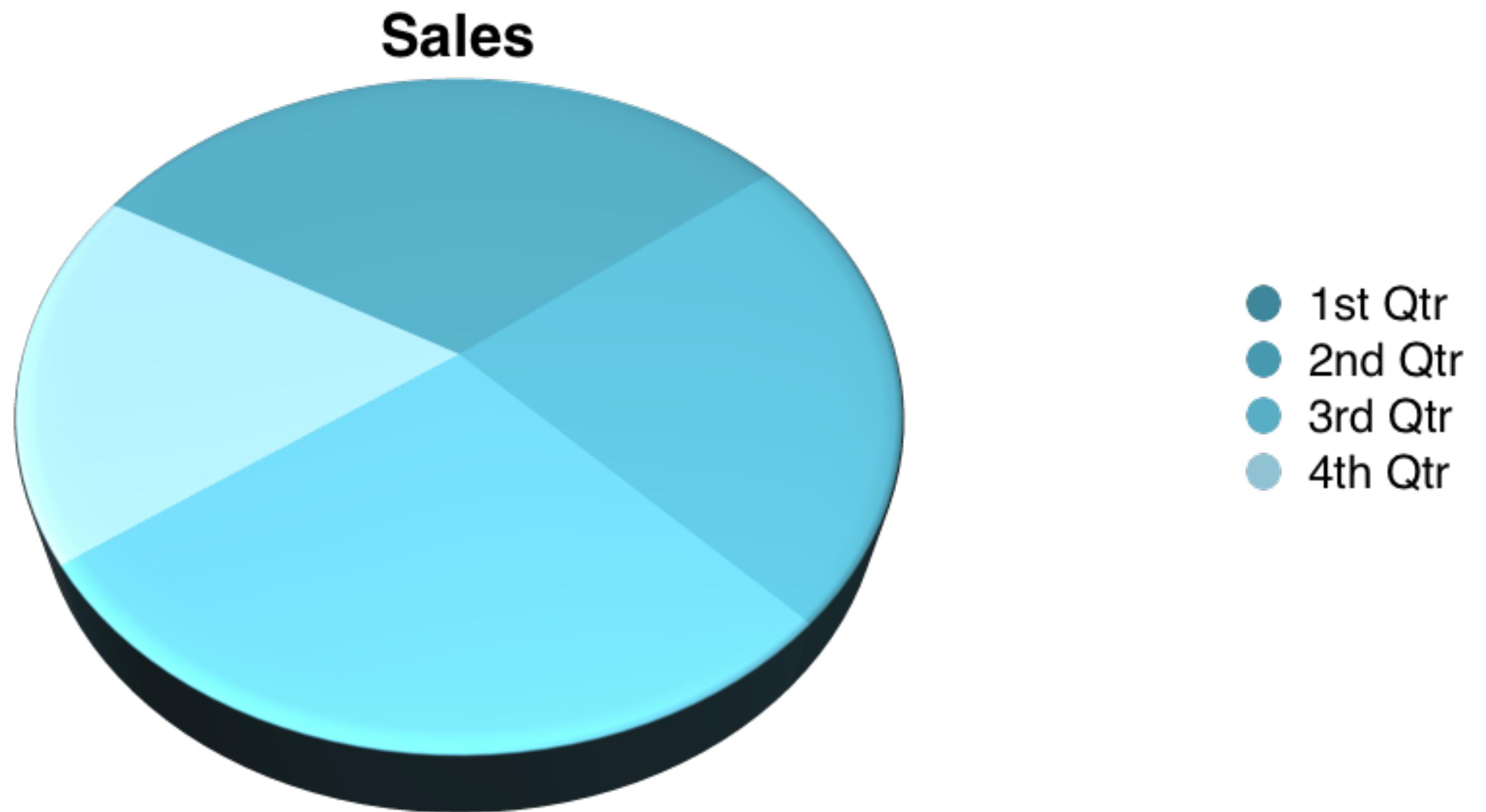


Proper Display



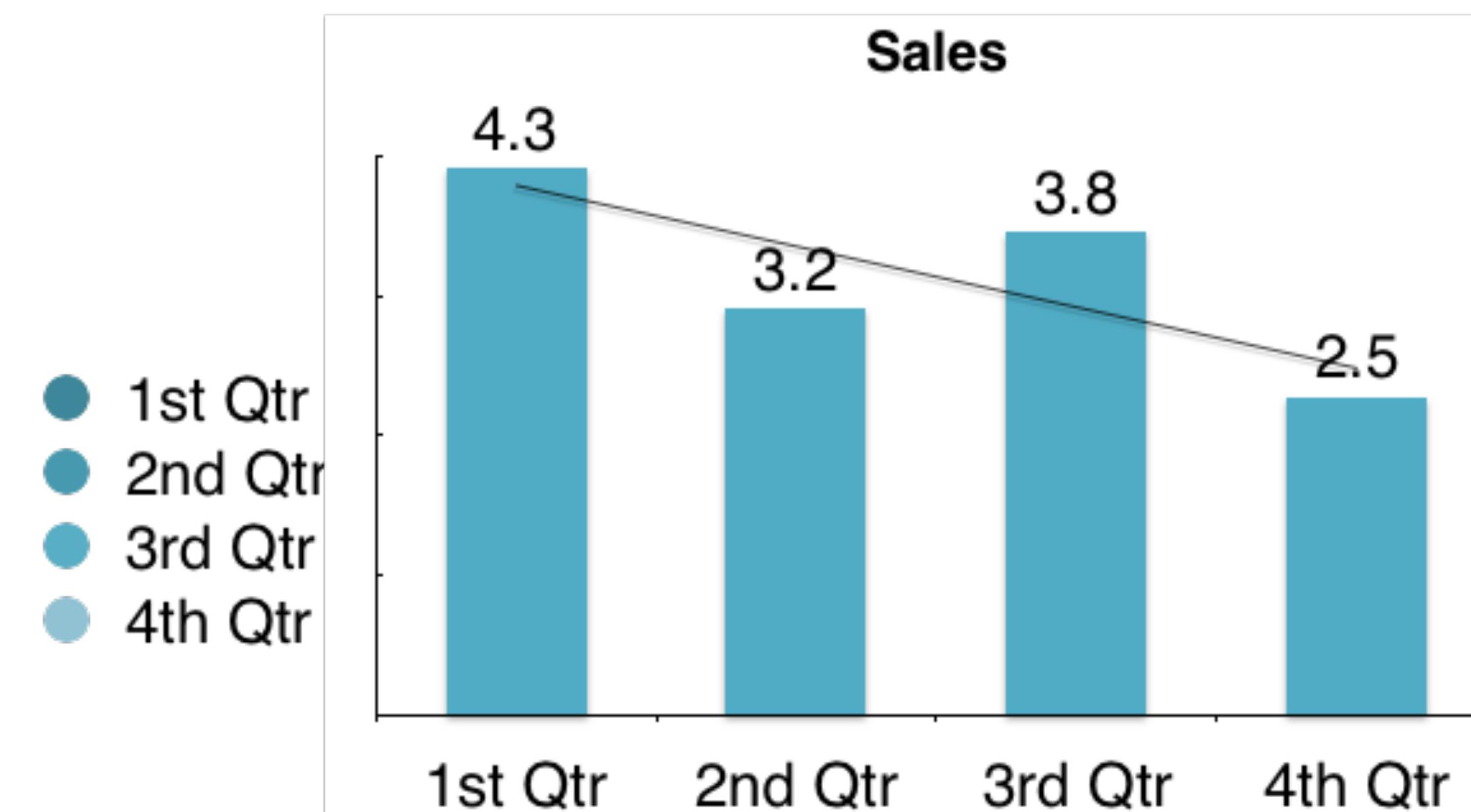
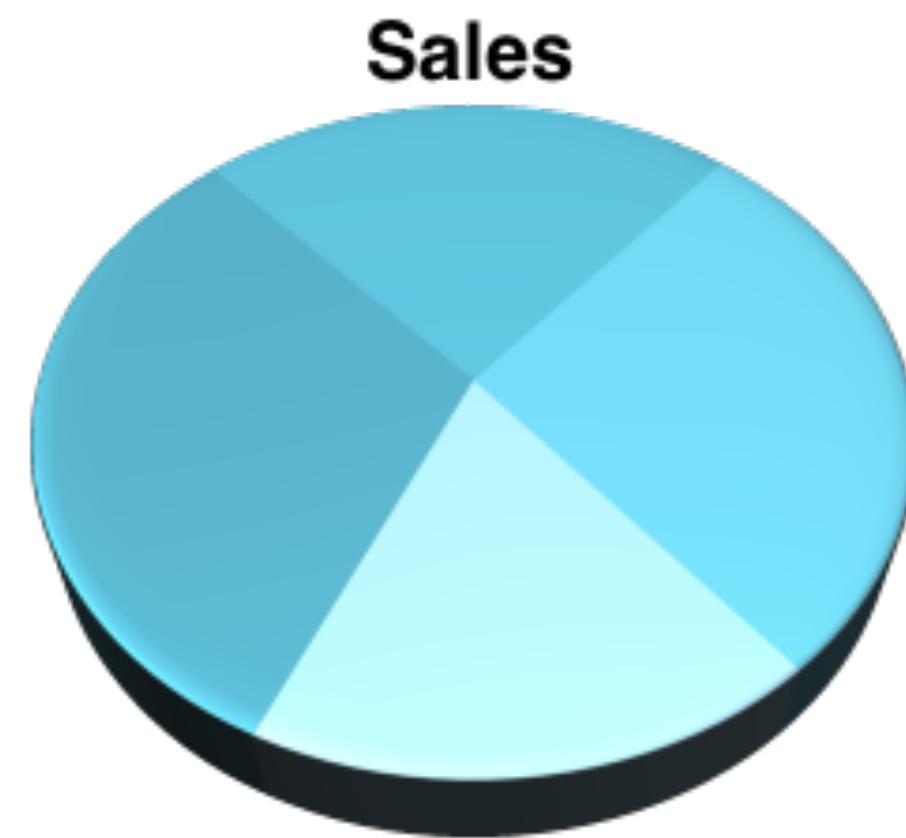


Proper Display



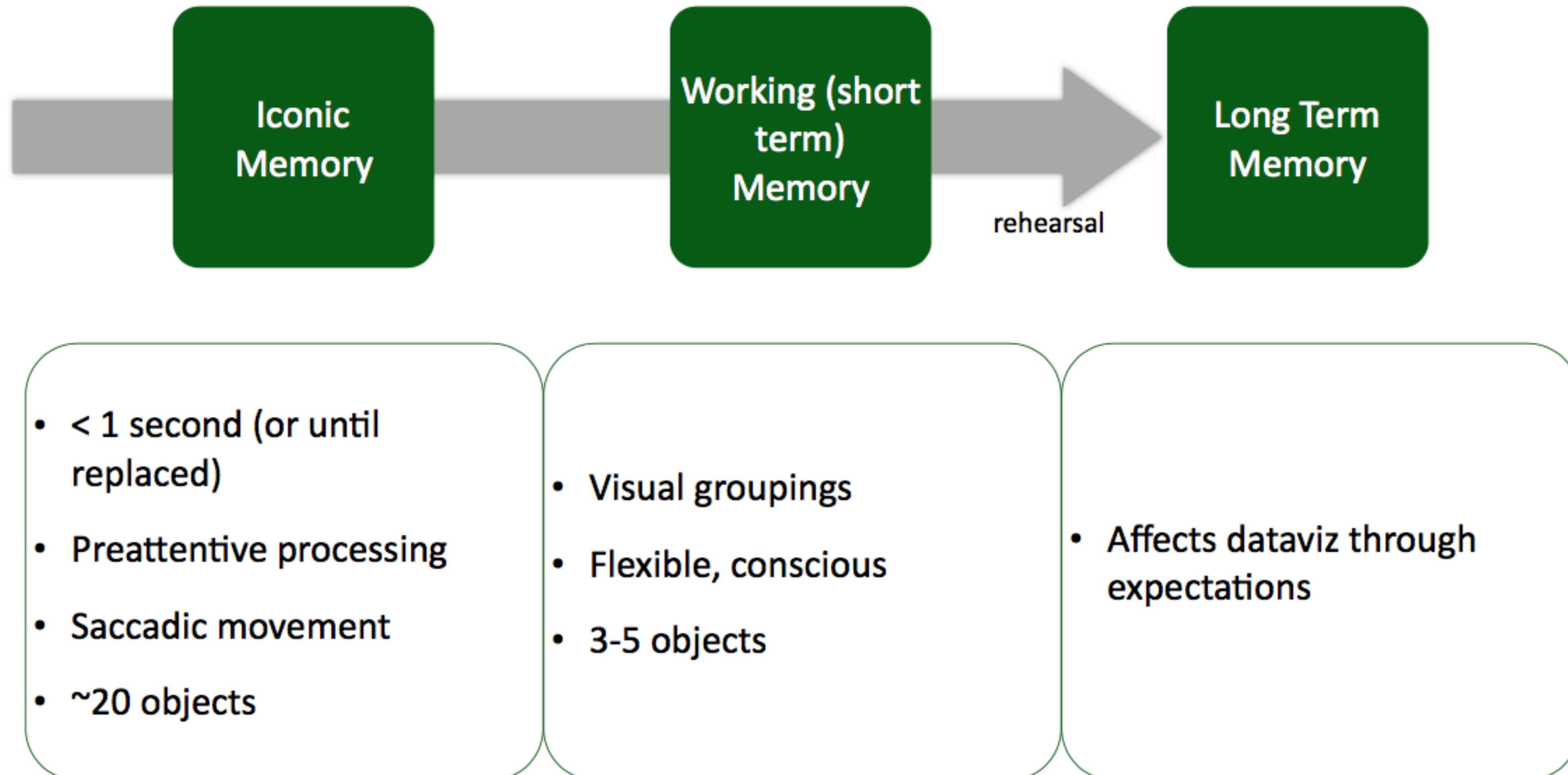


Proper Display





Visual Processing System

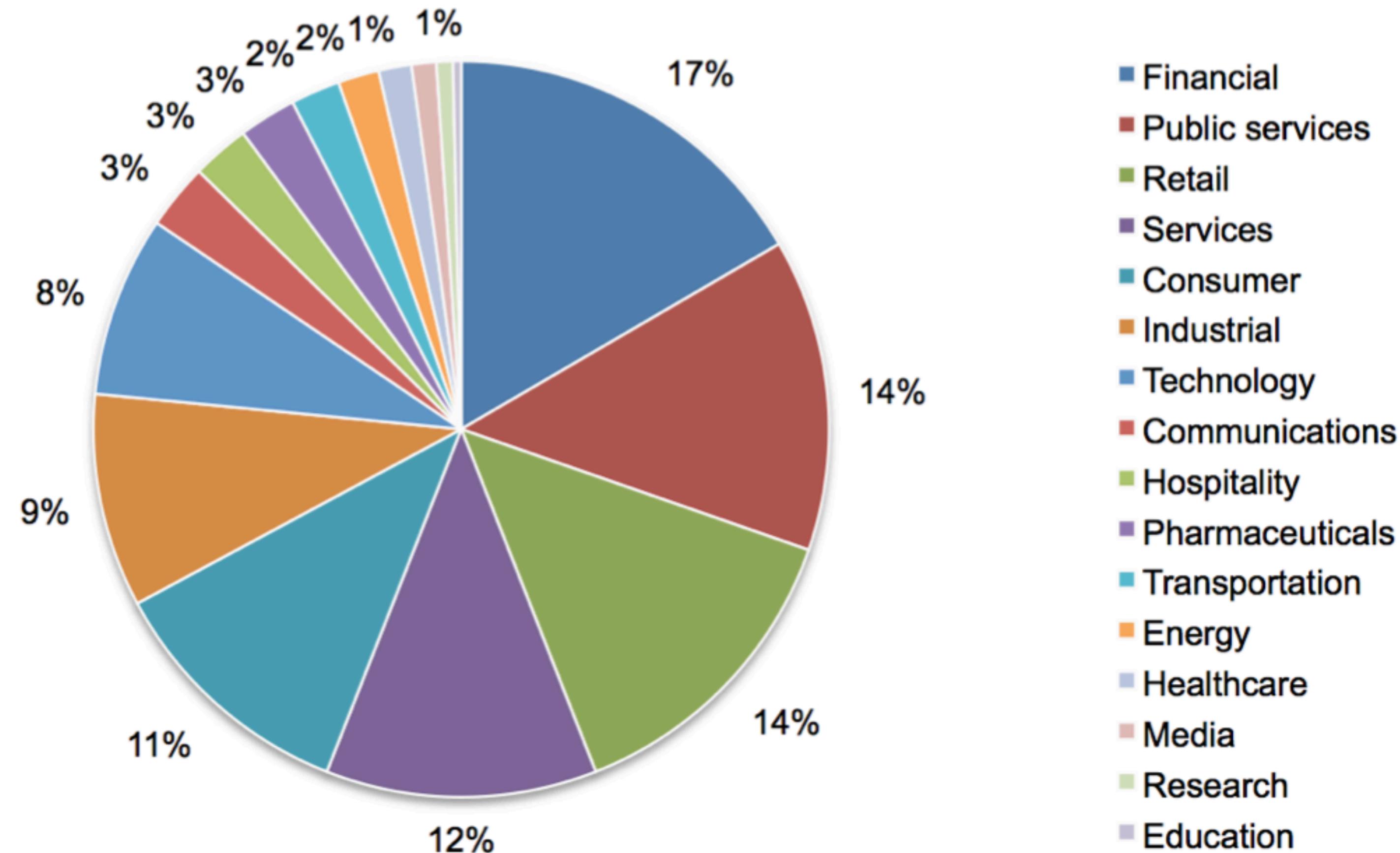


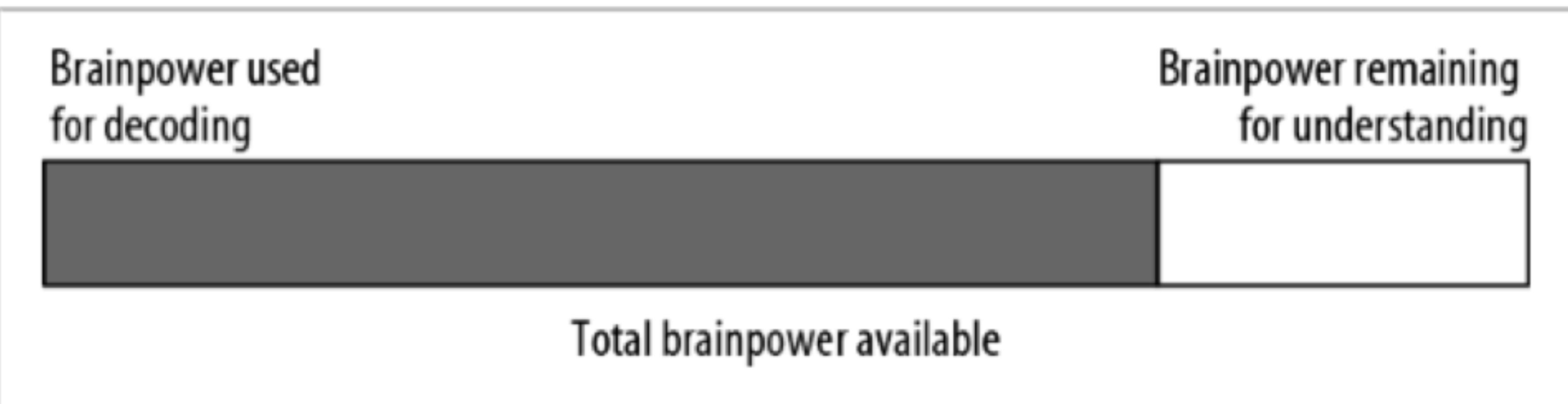
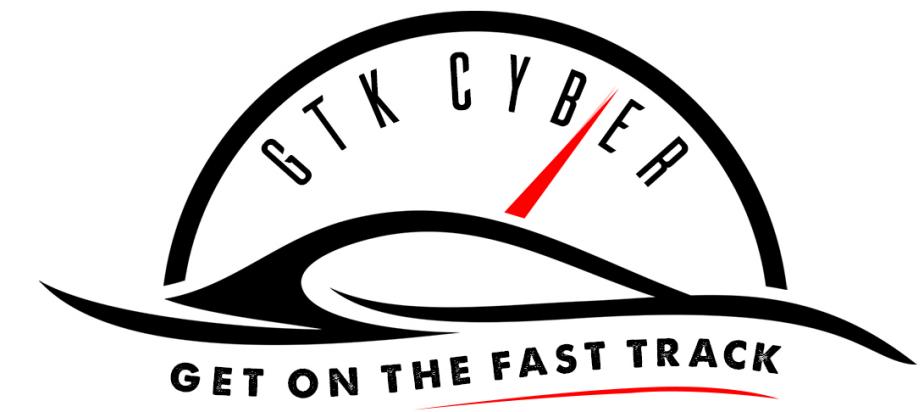


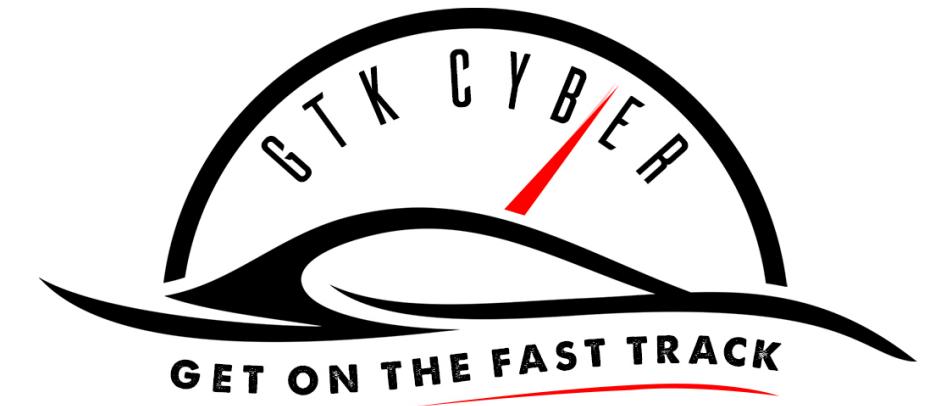
Overworking Visual Memory

Figure 20. Distribution of the benchmark sample by industry segment

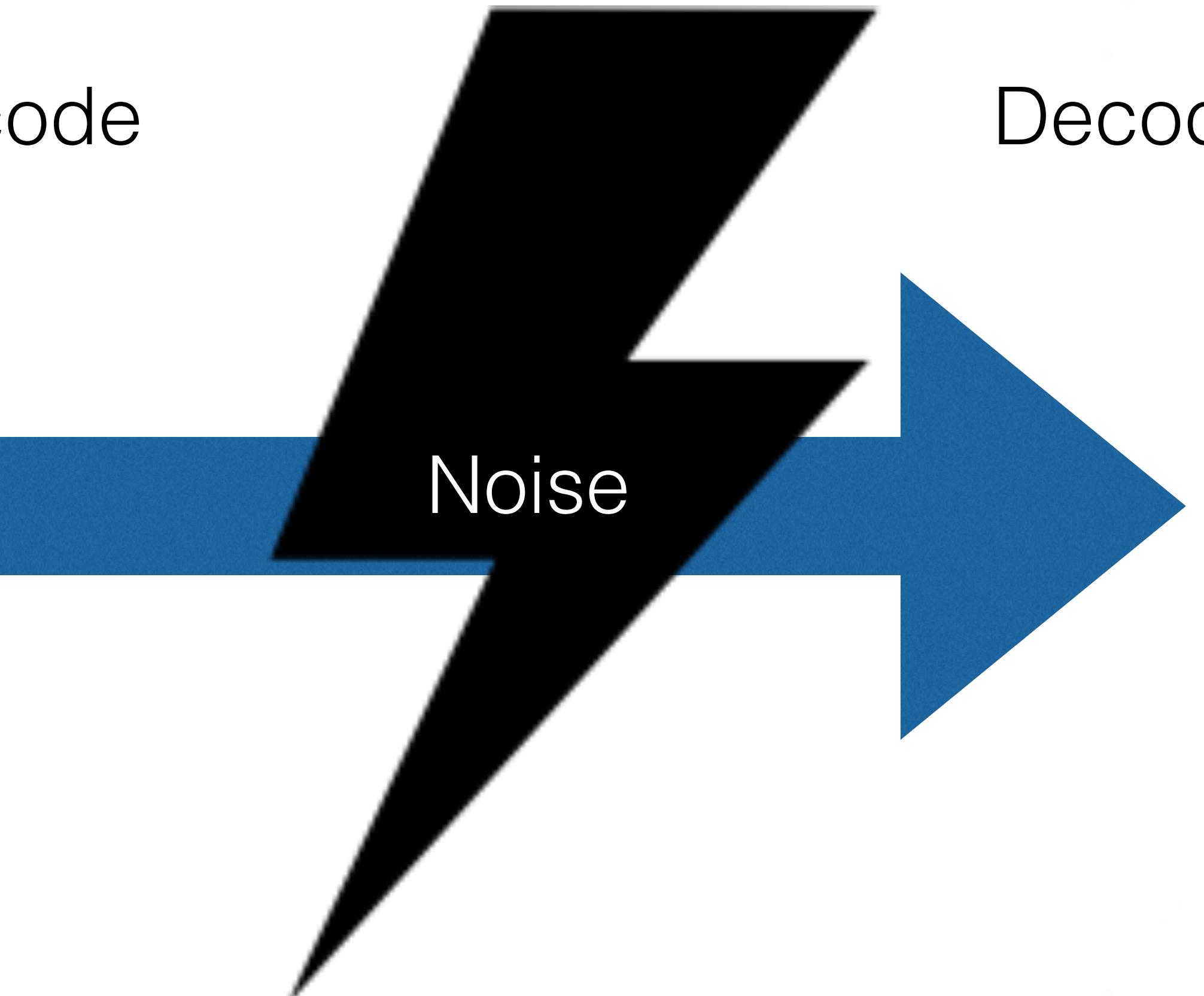
Consolidated (n = 277 organizations)

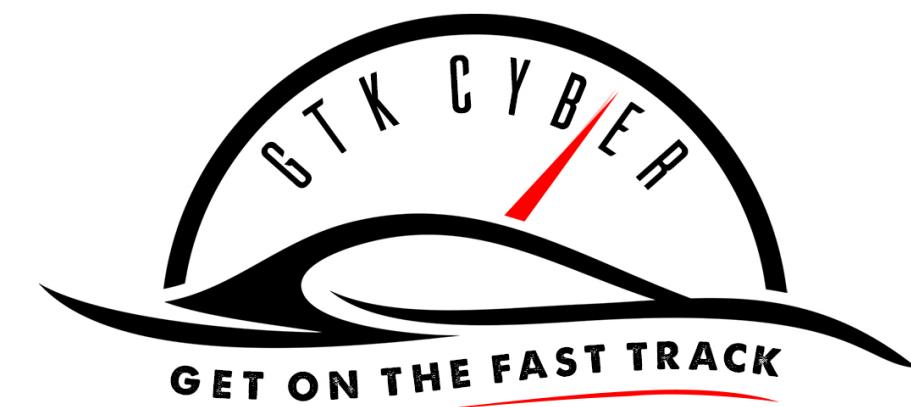




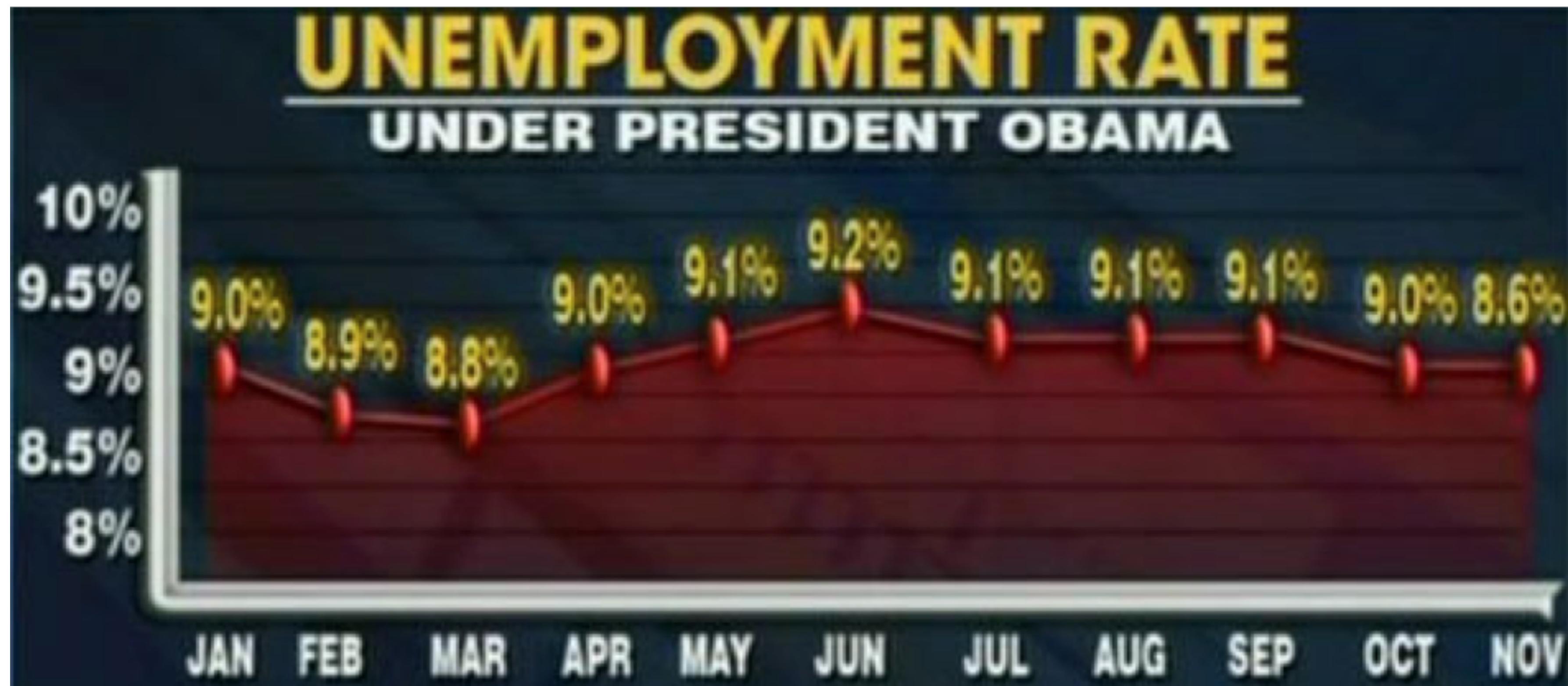


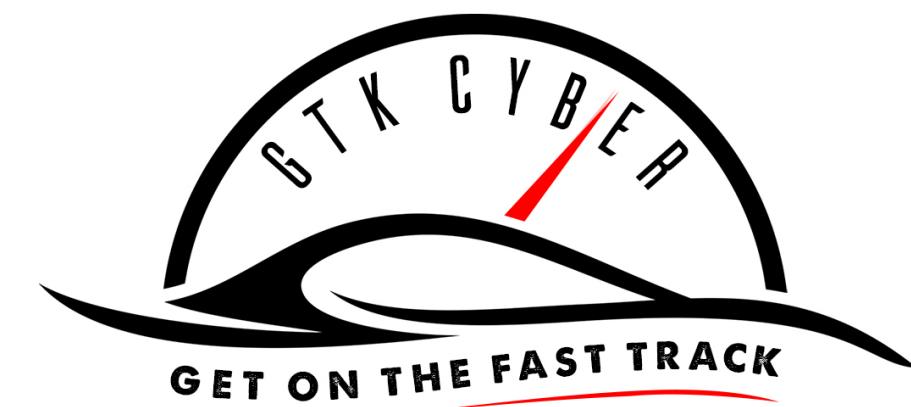
Encode



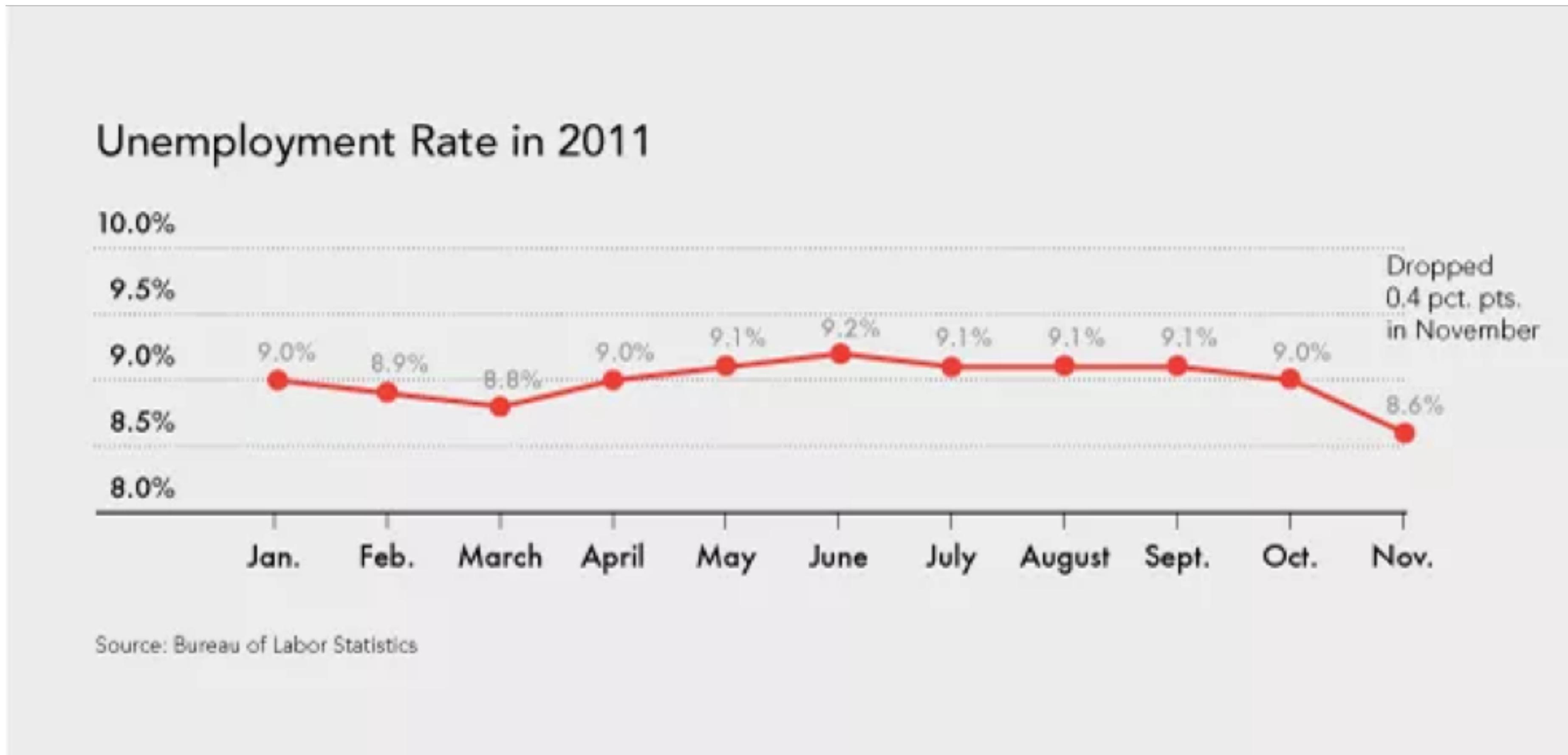


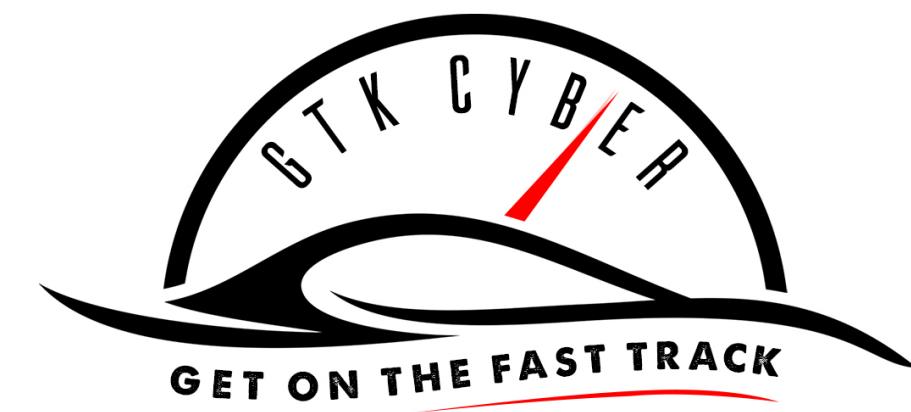
Graphical Integrity



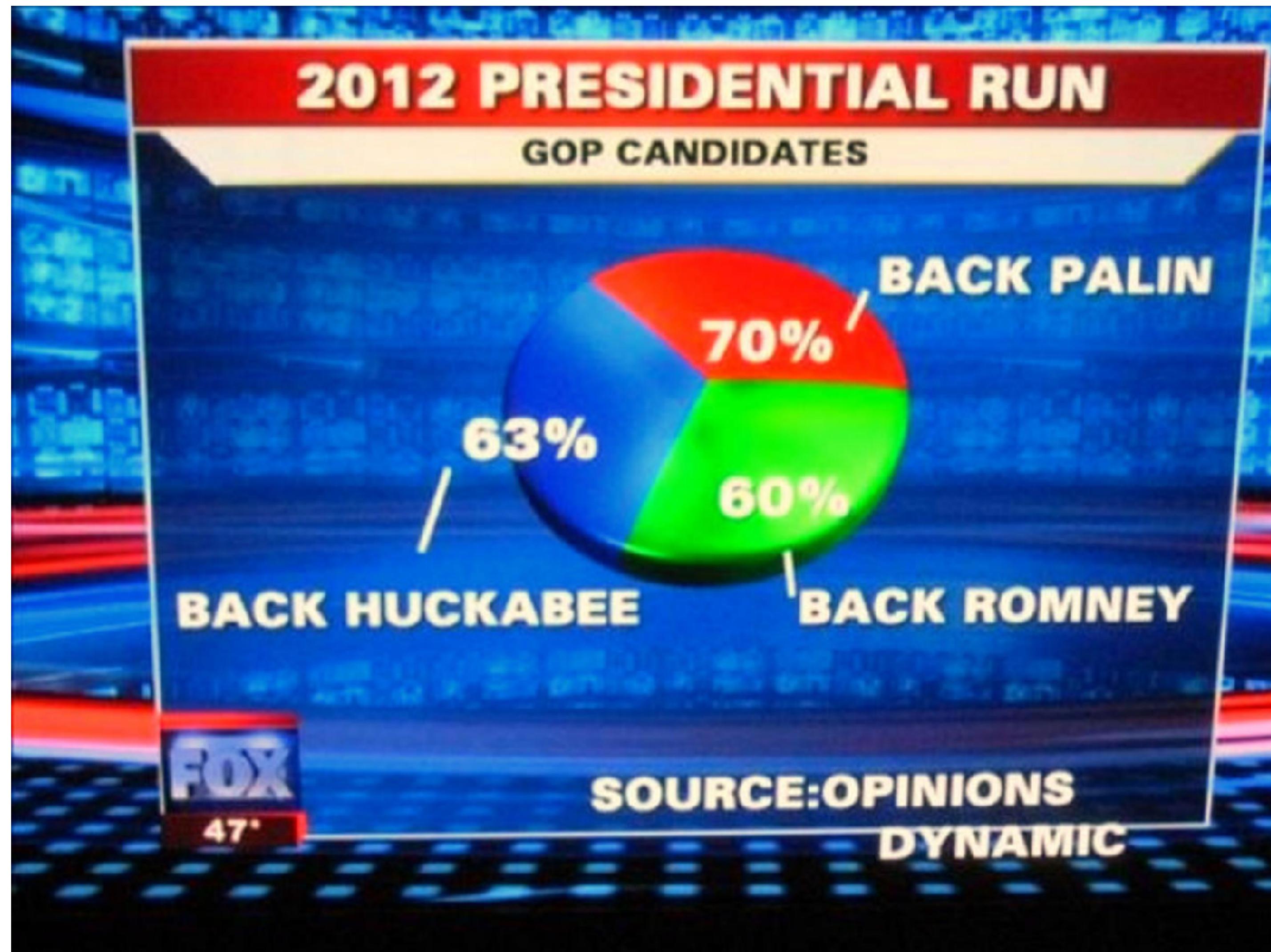


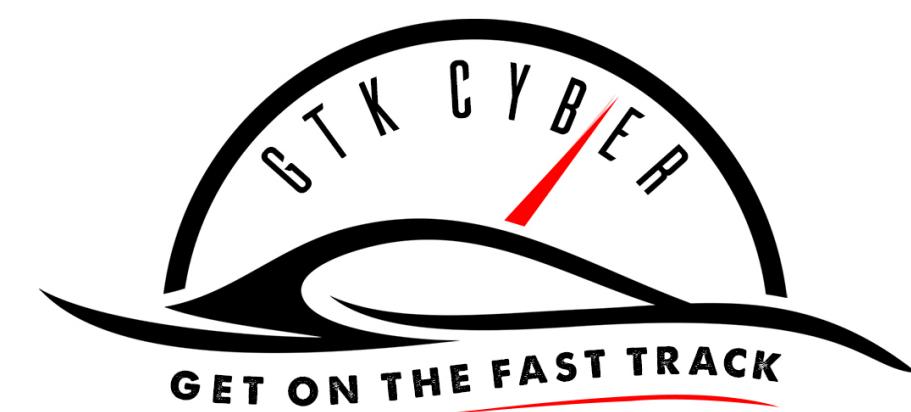
Graphical Integrity



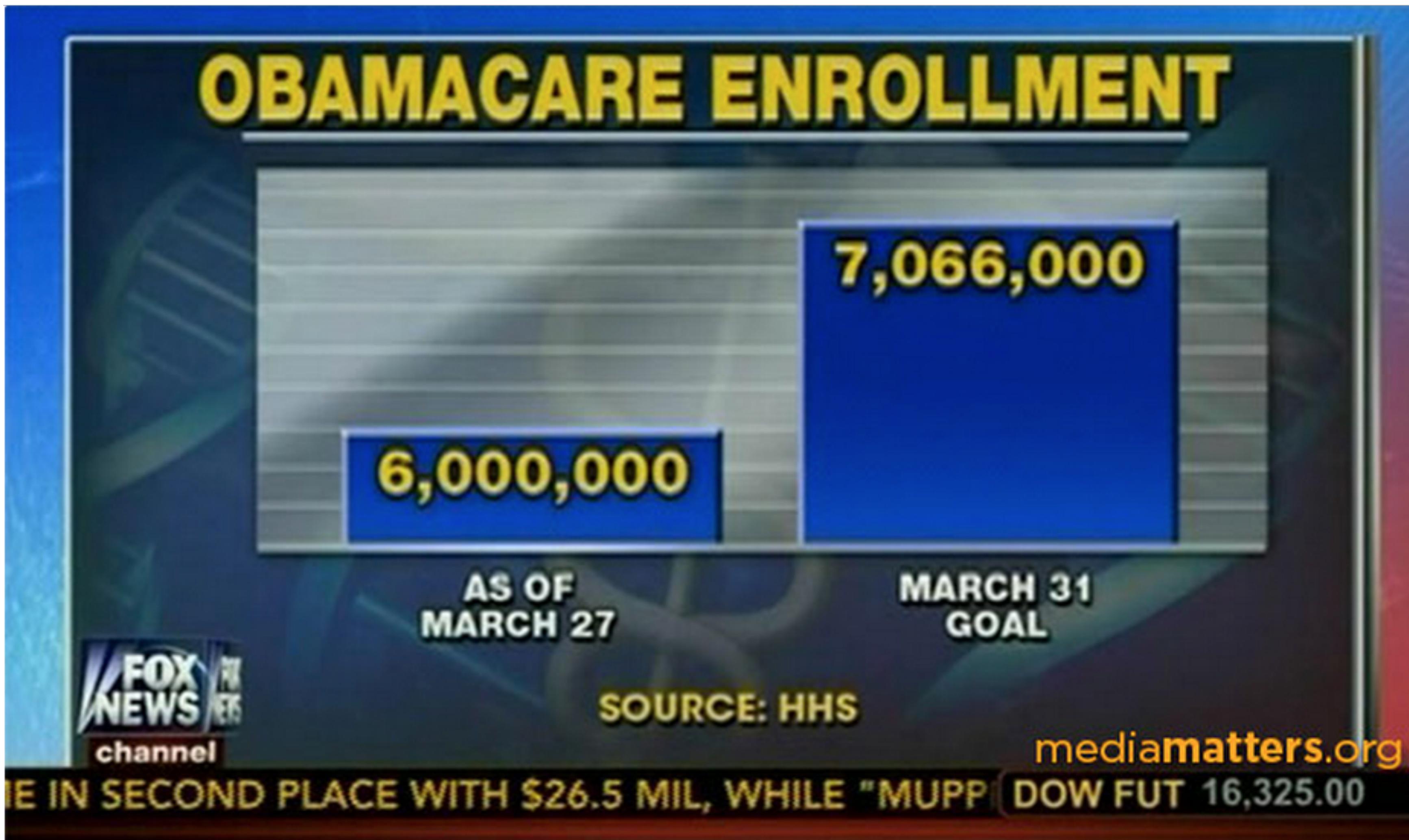


Graphical Integrity?



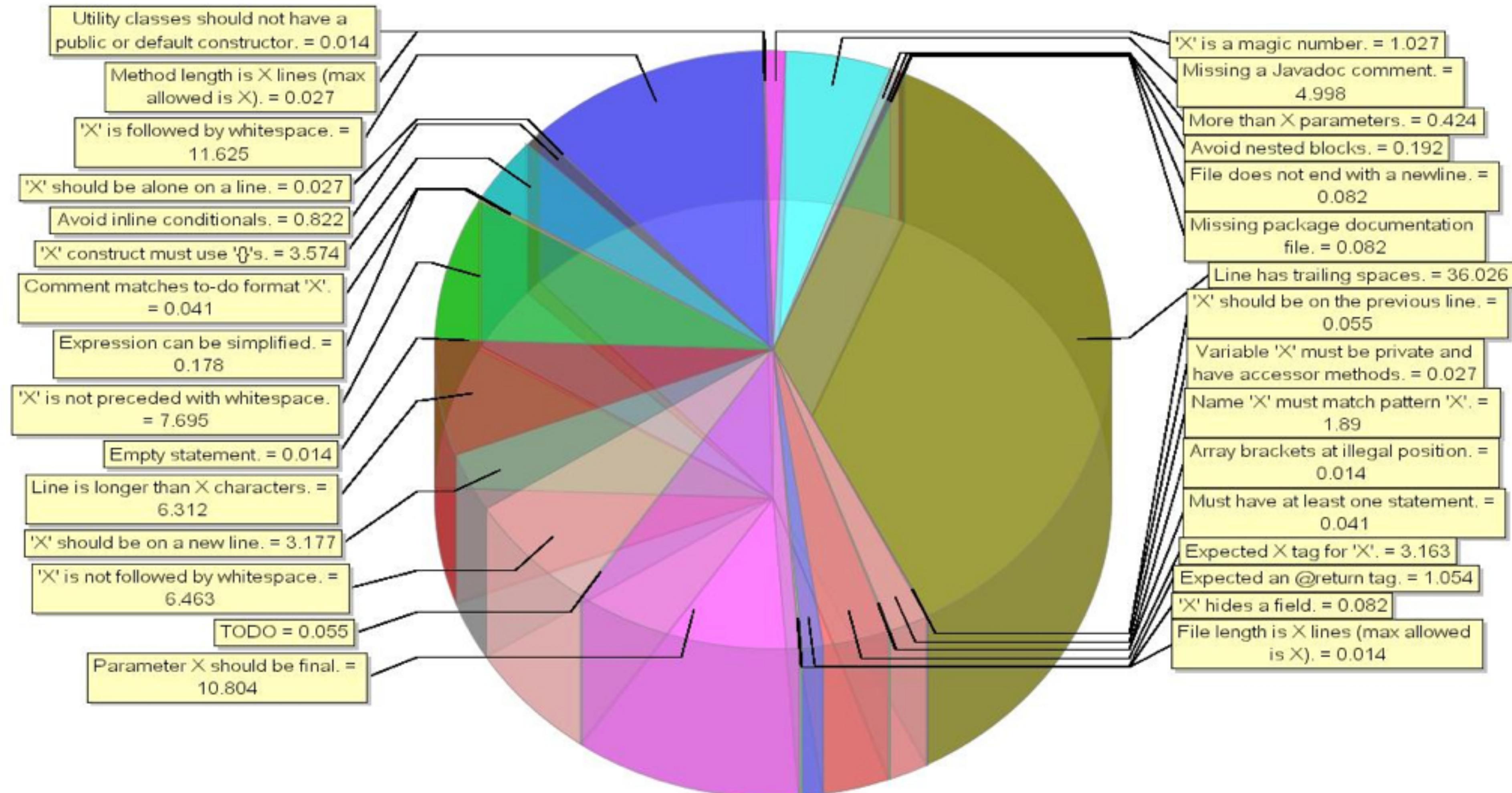


Graphical Integrity



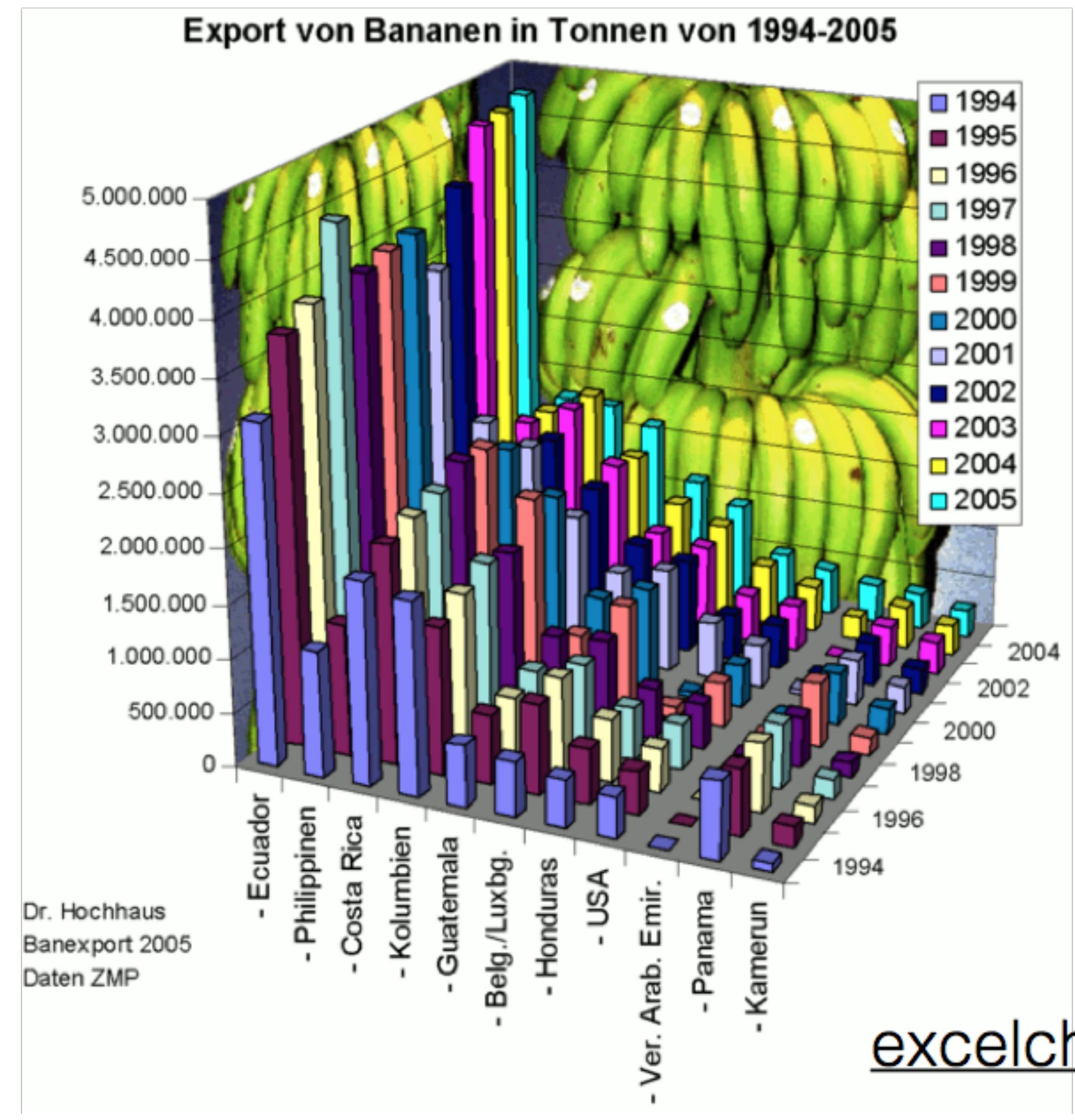


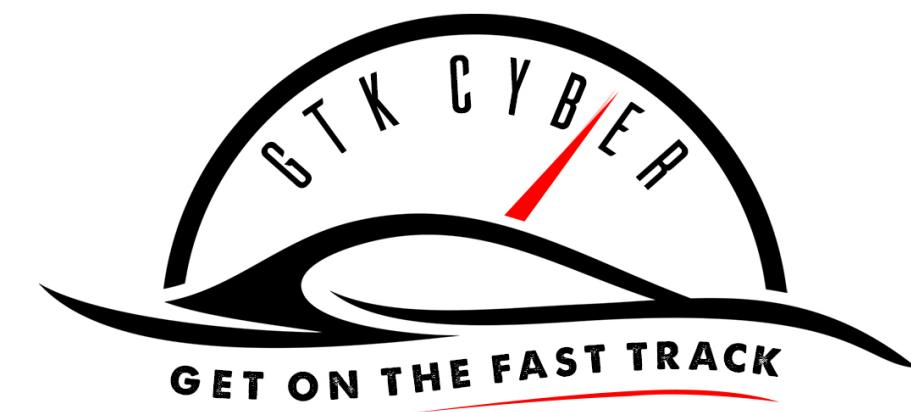
Simple



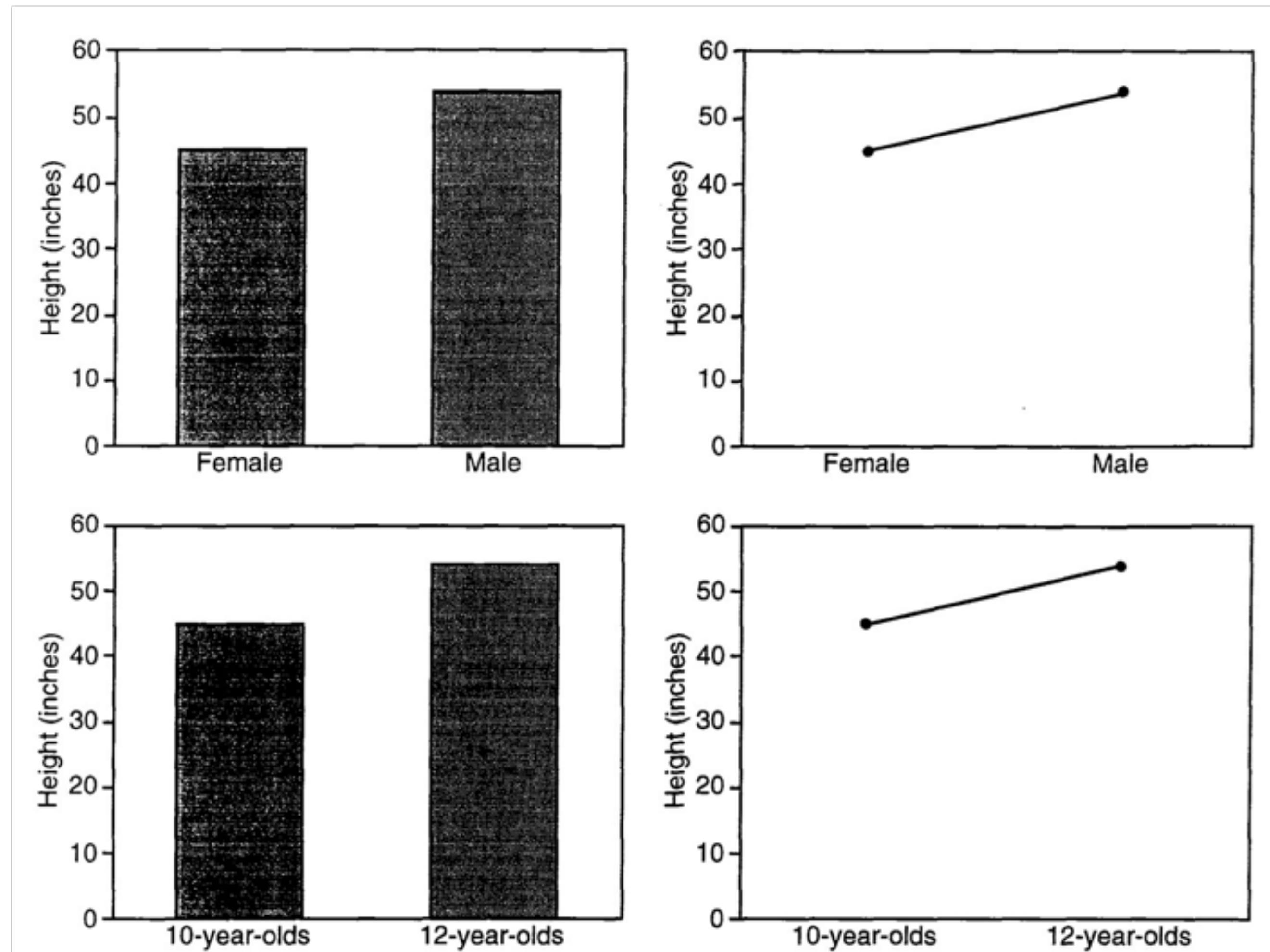


Simple





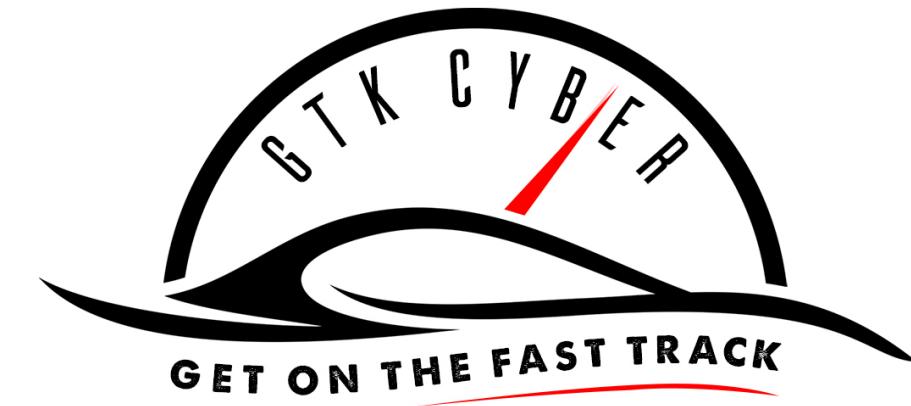
Proper Display



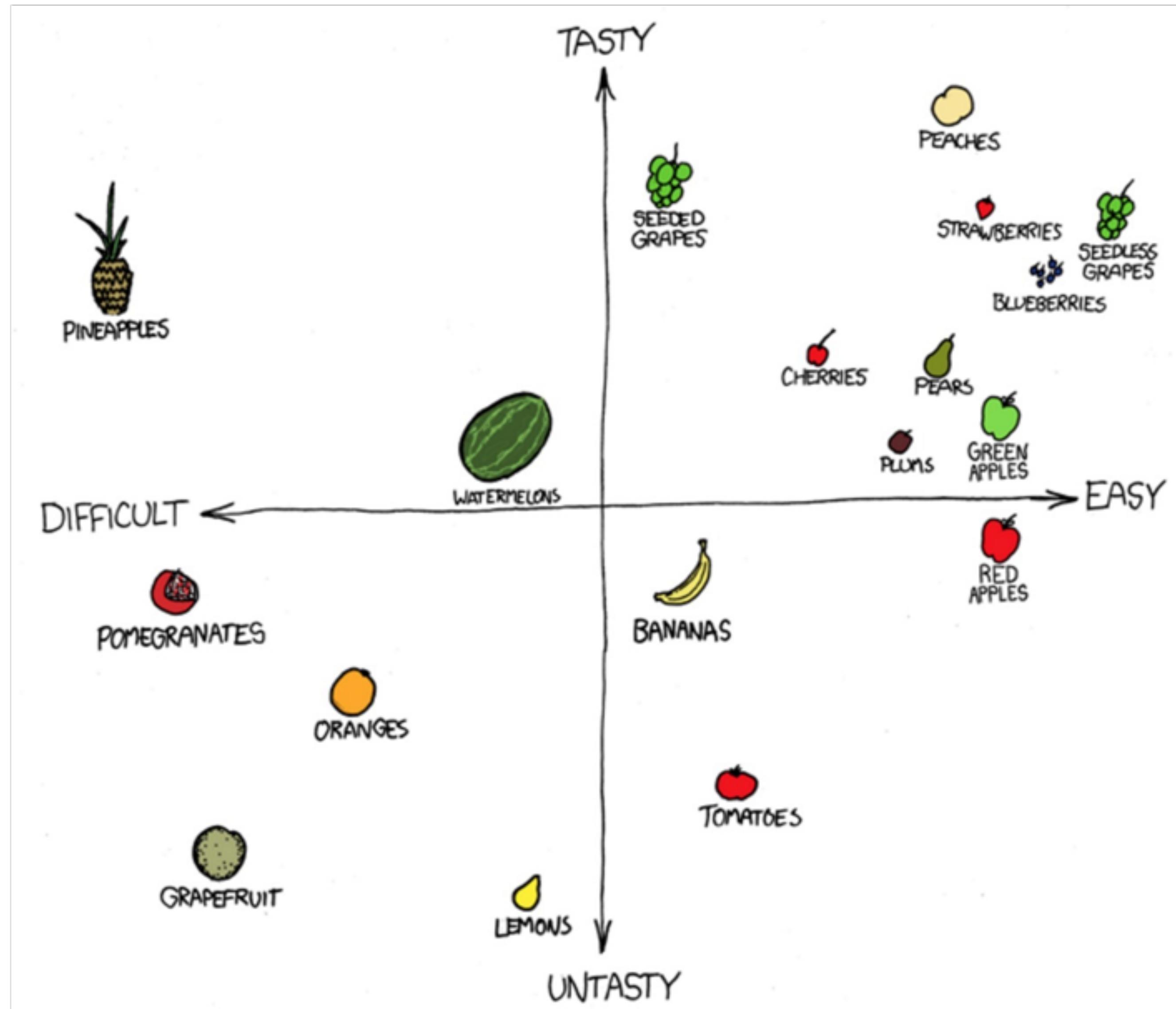


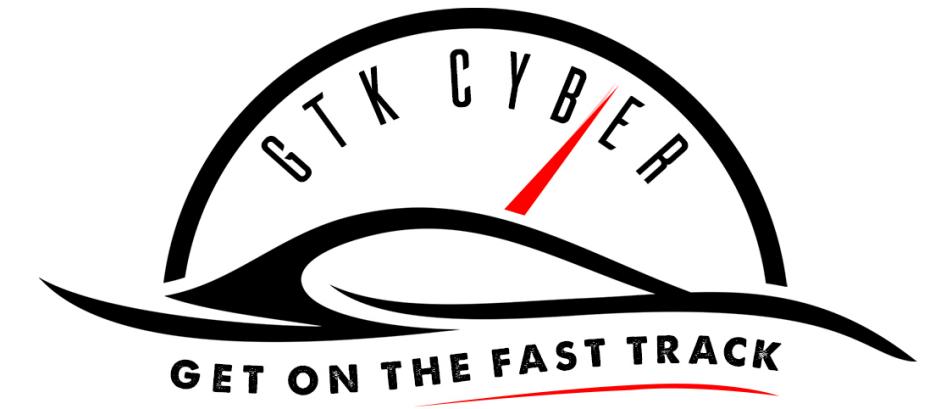
Proper Display



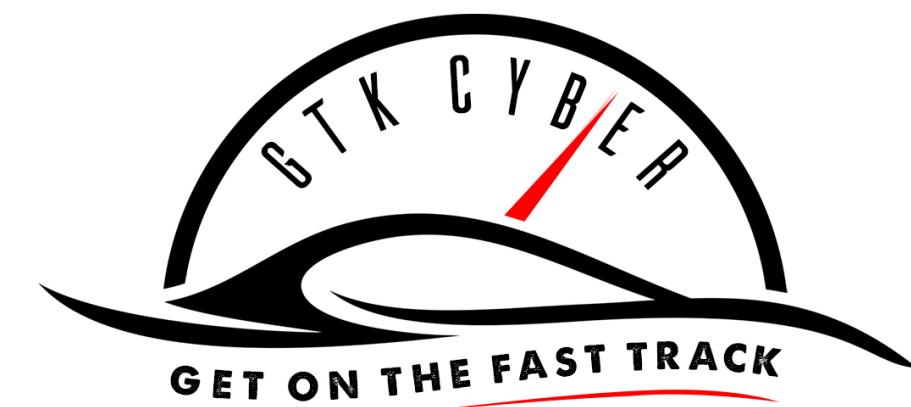


Proper Display

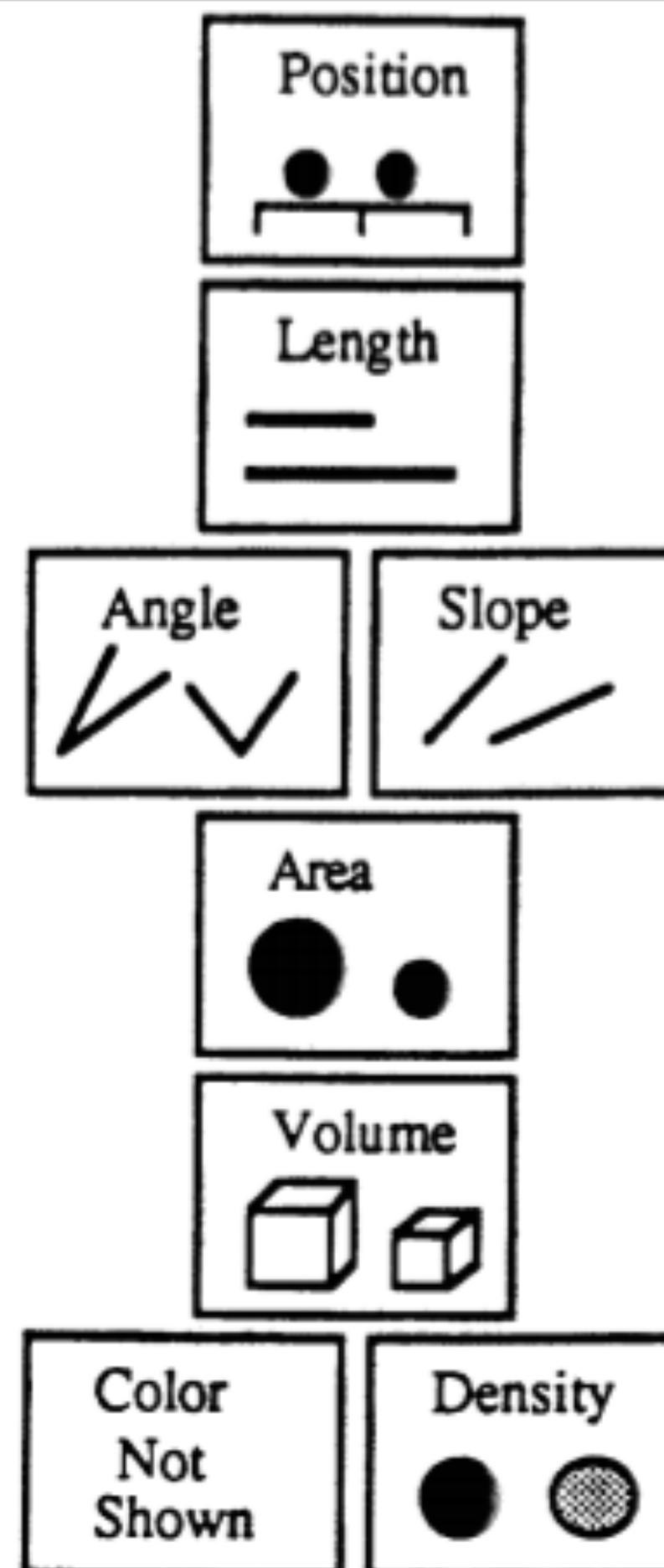




Visual Encoding



More accurate

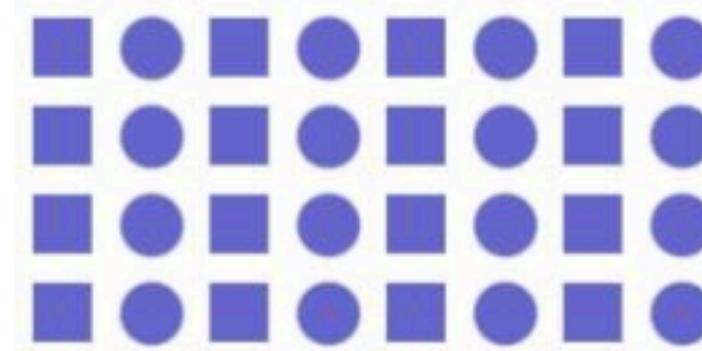


Less accurate





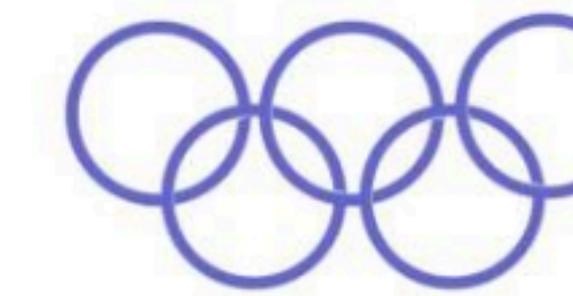
Gestalt Laws



Law of Similarity:

Items that are similar tend to be grouped together.

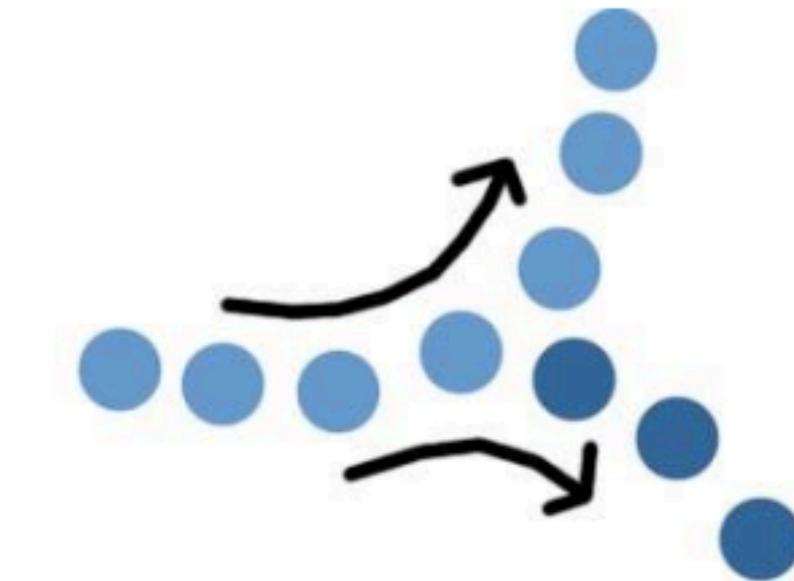
In the image above, most people see vertical columns of circles and squares.



Law of Pragnanz:

Reality is organized or reduced to the simplest form possible.

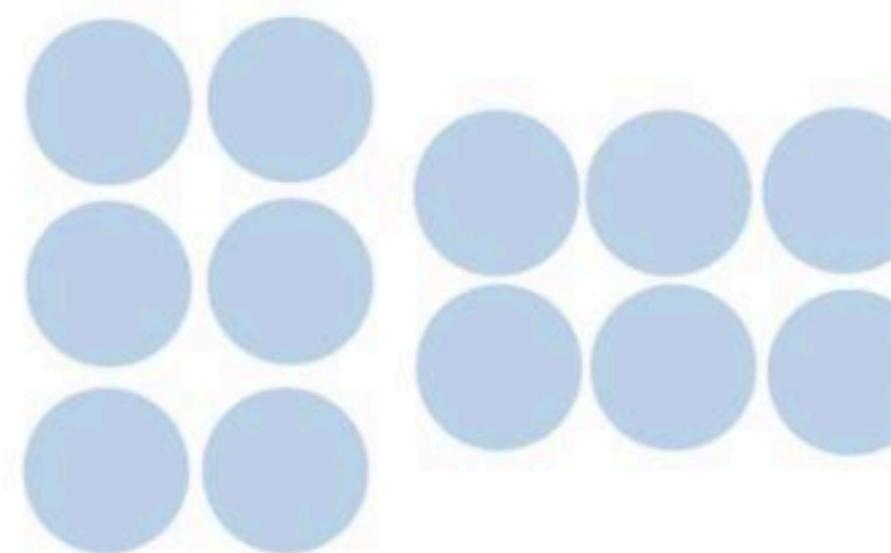
For example, we see the image above as a series of circles rather than as many much more complicated shapes.



Law of Continuity:

Lines are seen as following the smoothest path.

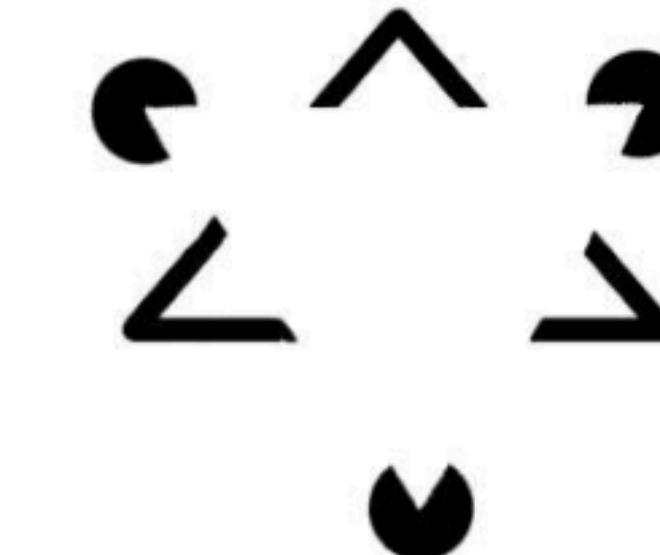
In the image above, the top branch is seen as continuing the first segment of the line. This allows us to see things as flowing smoothly without breaking lines up into multiple parts.



Law of Proximity:

Objects near each other tend to be grouped together.

The circles on the left appear to be grouped in vertical columns, while those on the right appear to be grouped in horizontal rows.



Law of Closure:

Objects grouped together are seen as a whole.

We tend to ignore gaps and complete contour lines. In the image above, there are no triangles or circles, but our minds fill in the missing information to create familiar shapes and images.



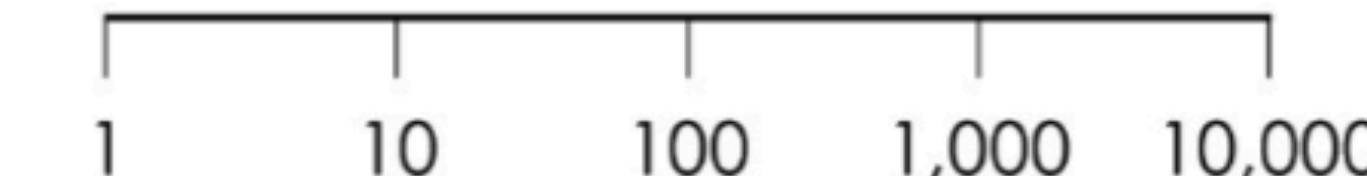
Linear

Values are evenly spaced



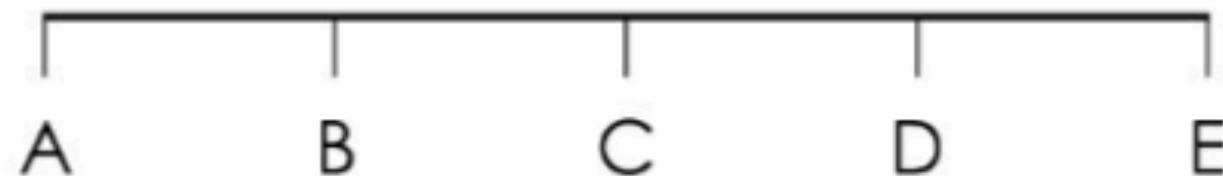
Logarithmic

Focus on percent change



Categorical

Discrete placement in bins



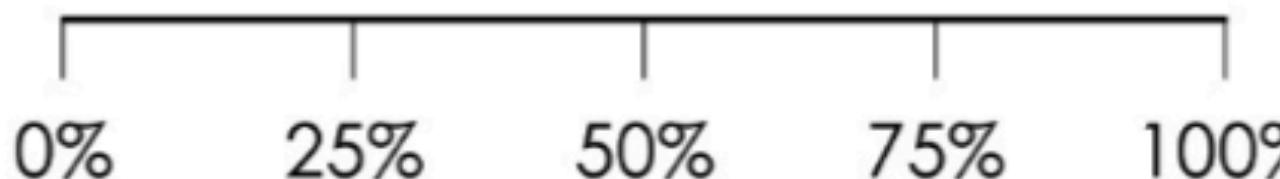
Ordinal

Categories where order matters



Percent

Representing parts of a whole

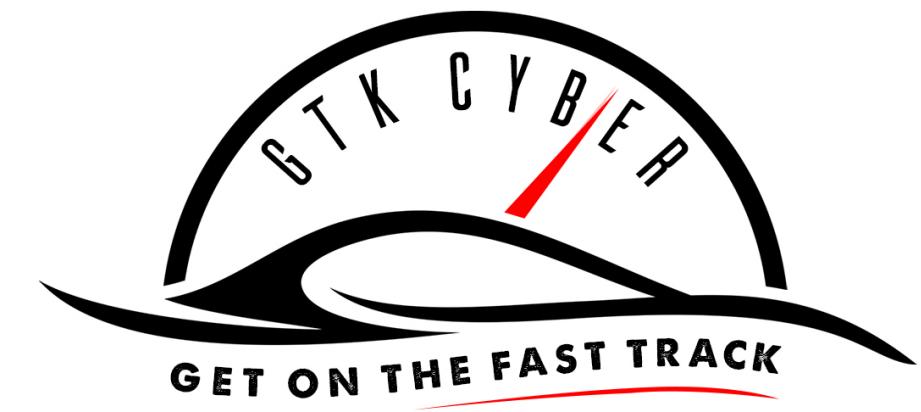


Time

Units of months, days, or hours



Source: Nathan Yau, Data Points

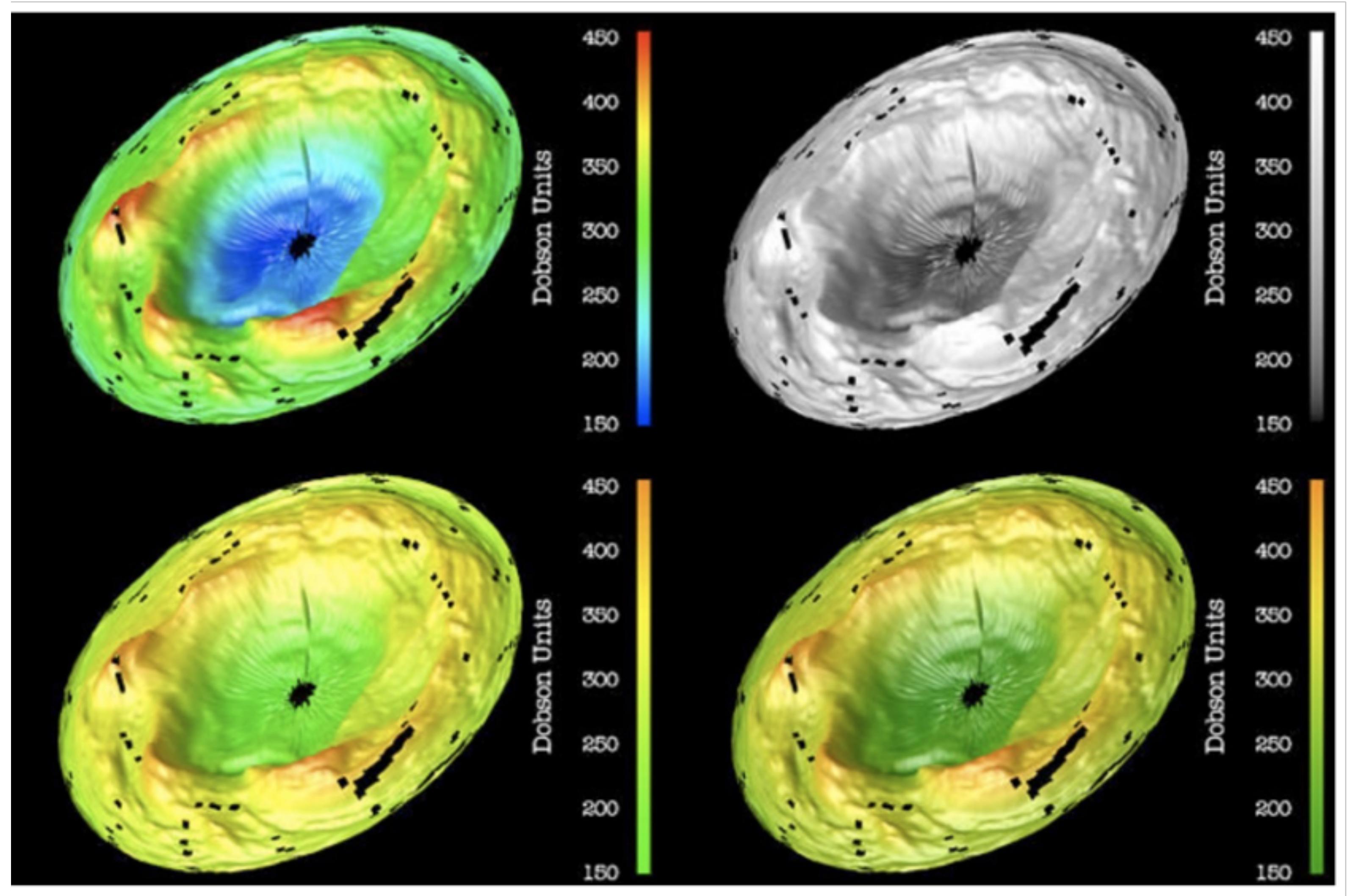


Color

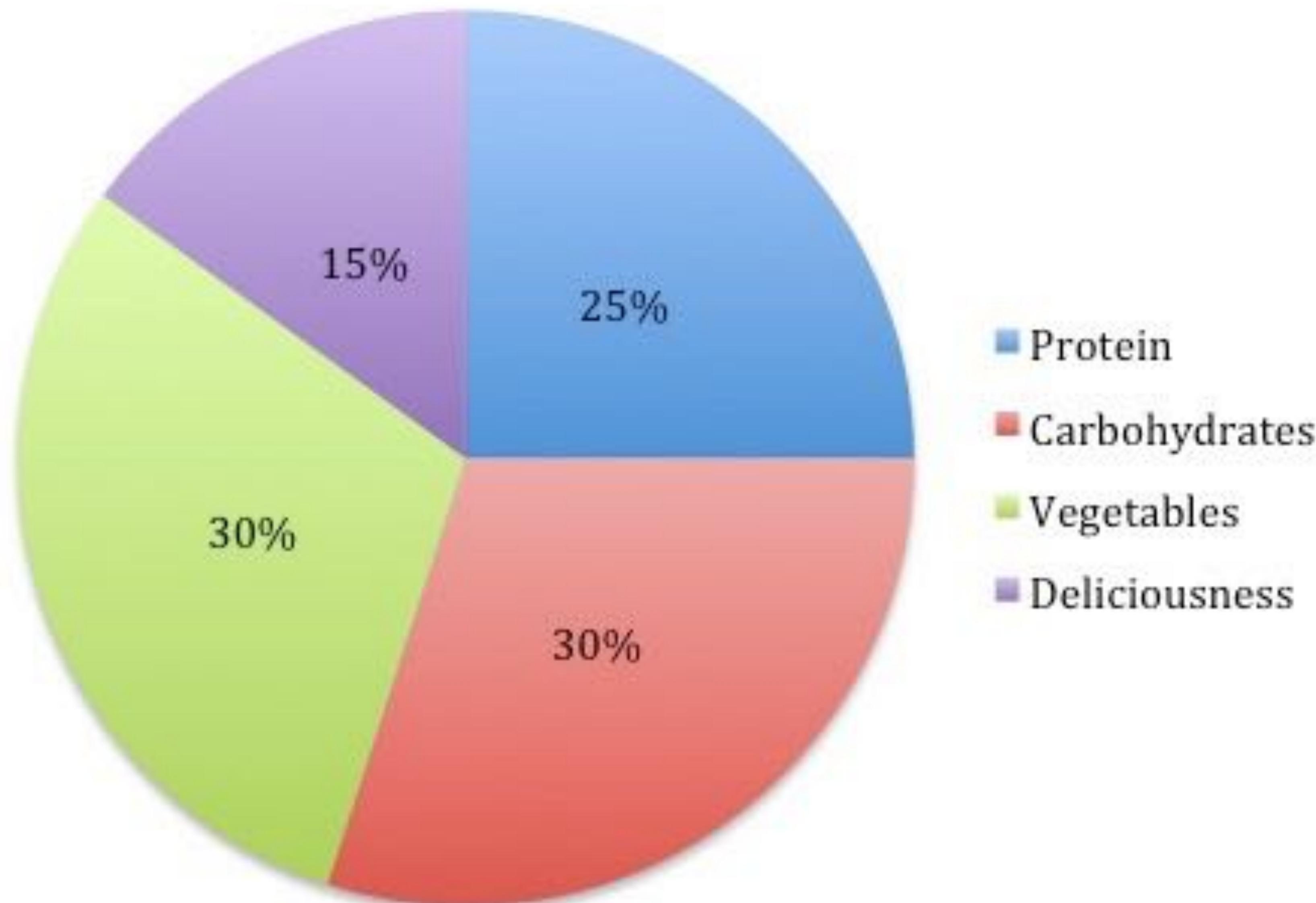




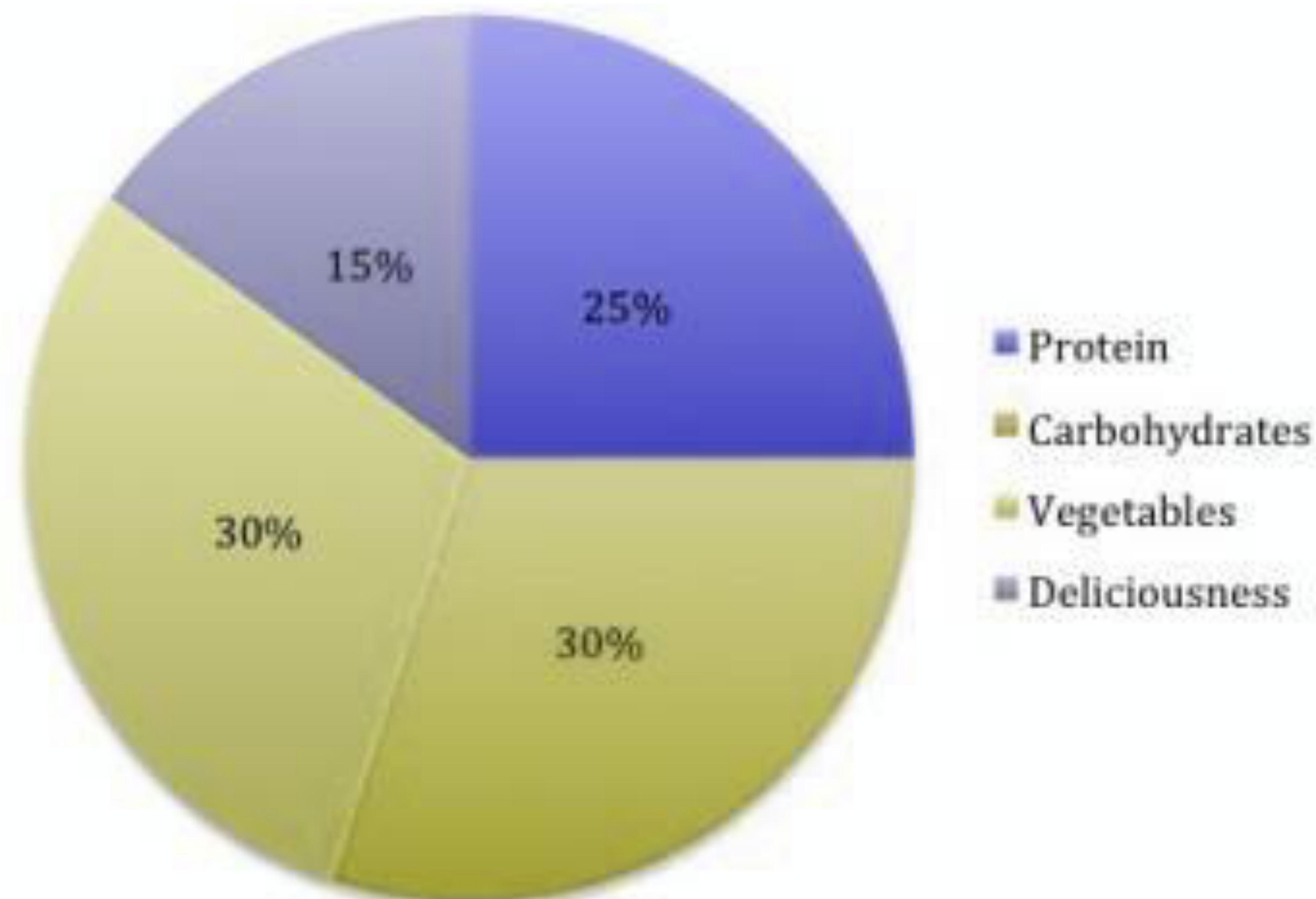
Color



A Healthy Meal



A Healthy Meal





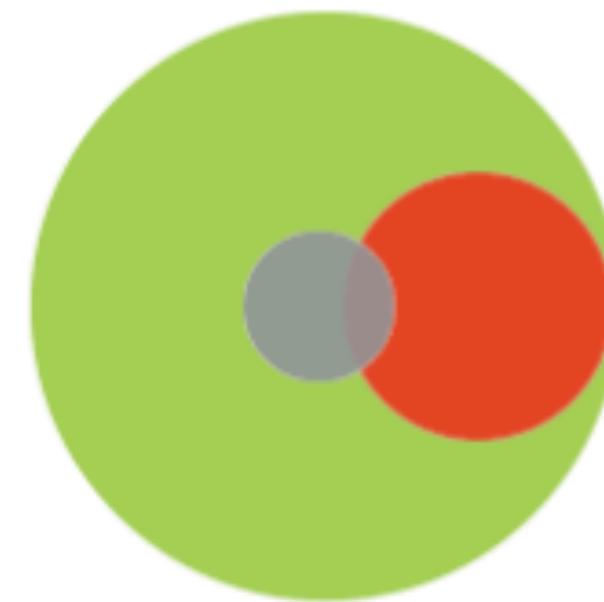
elastica



Dashboard

Users

Total (192)



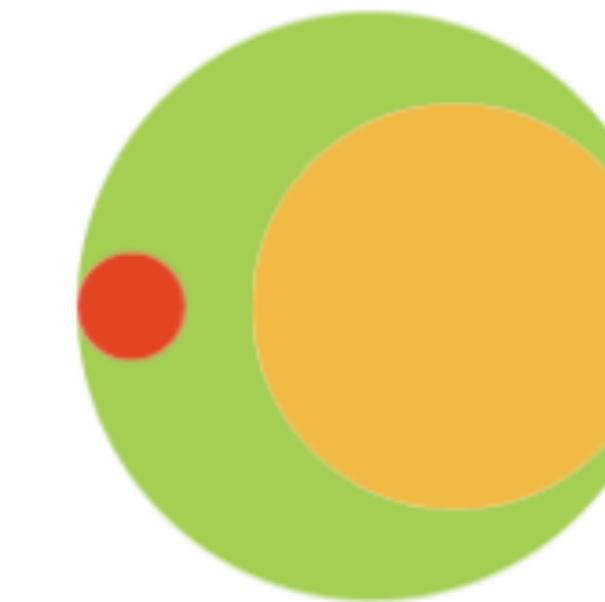
28
High Risk

0
Med Risk

4
Blocked

Policies

Total (231)



3
Blocking

142
Alerting

Policy Alerts

Alerting



Blocked



Rest



Threat Alerts

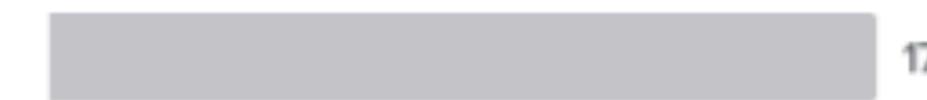
High Risk



Med Risk



Low Risk



Audited Services

by Users ▾

High Risk (736)

Medium Risk (3k)

Low Risk (3k)



3k
Users

494.3 GB
Traffic

963k
Sessions

243
Destinations

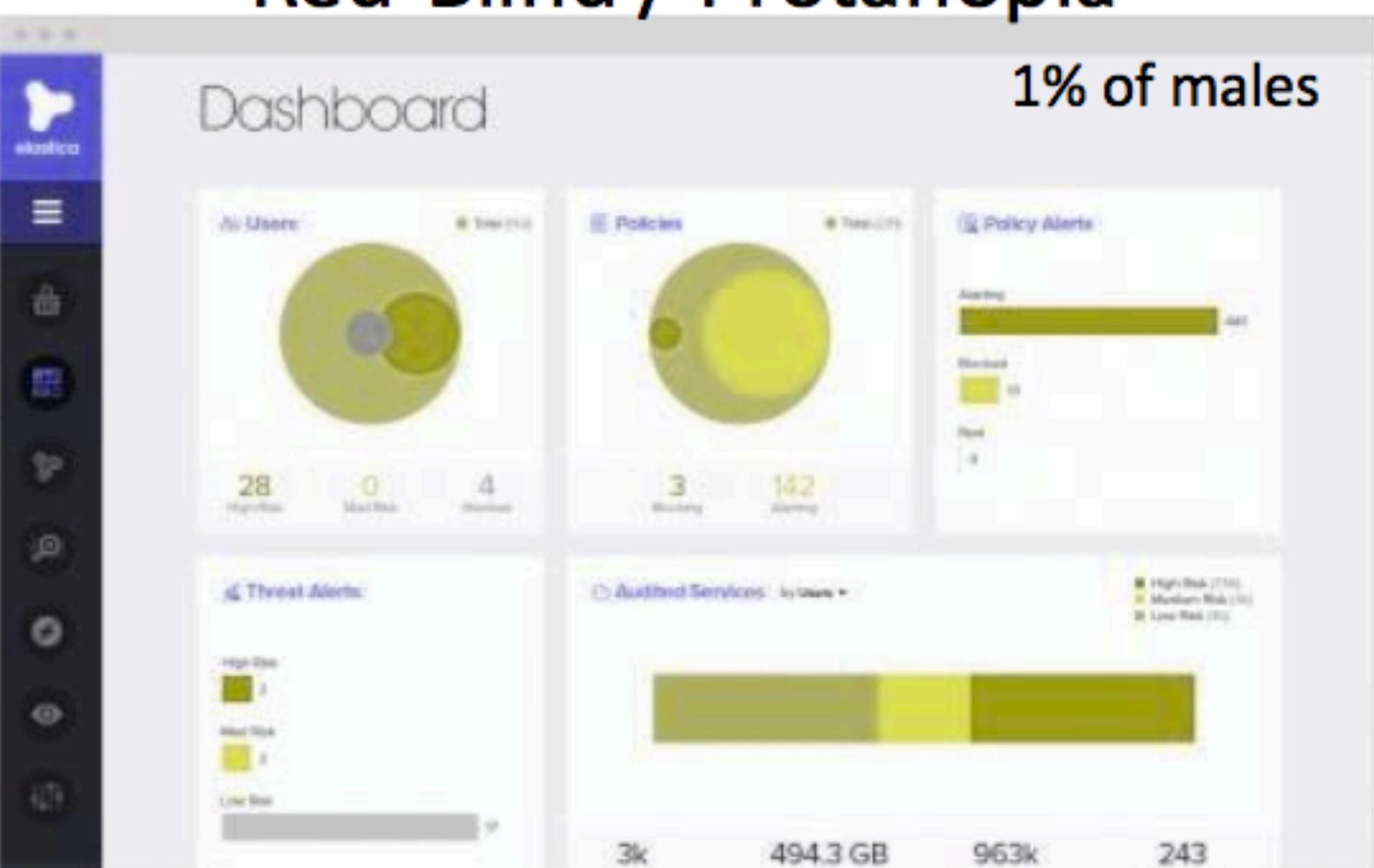
Original



Green-Blind / Deutanopia

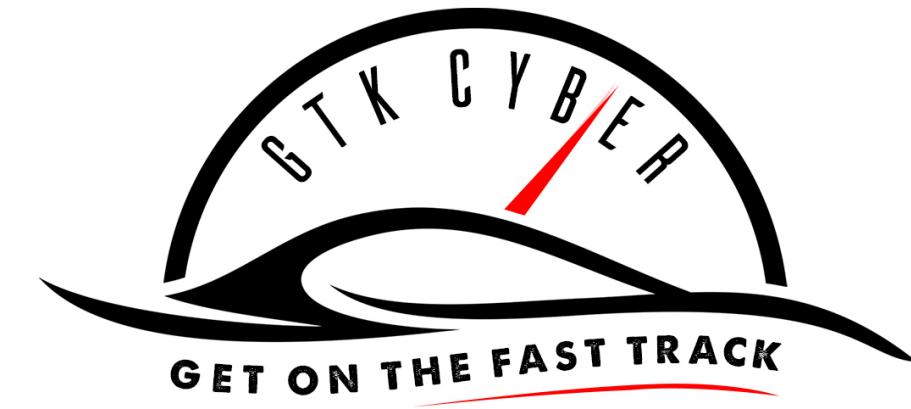


Red-Blind / Protanopia

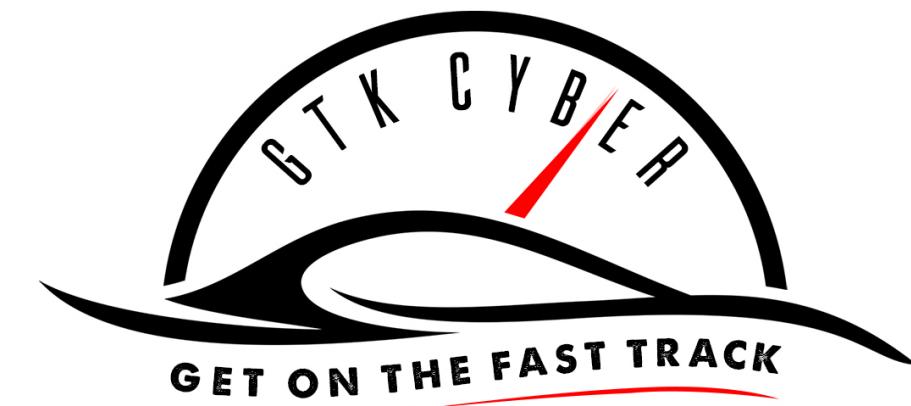


Blue-Blind / Tritanopia





<http://www.color-blindness.com/coblis-color-blindness-simulator/>



Color

Number of data classes: 3

Nature of your data:
 sequential diverging qualitative

Pick a color scheme:
Multi-hue: Single hue:

Only show:
 colorblind safe
 print friendly
 photocopy safe

Context:
 roads
 cities
 borders

Background:
 solid color terrain
 color transparency

how to use | updates | downloads | credits

COLORBREWER 2.0
color advice for cartography

3-class BuGn

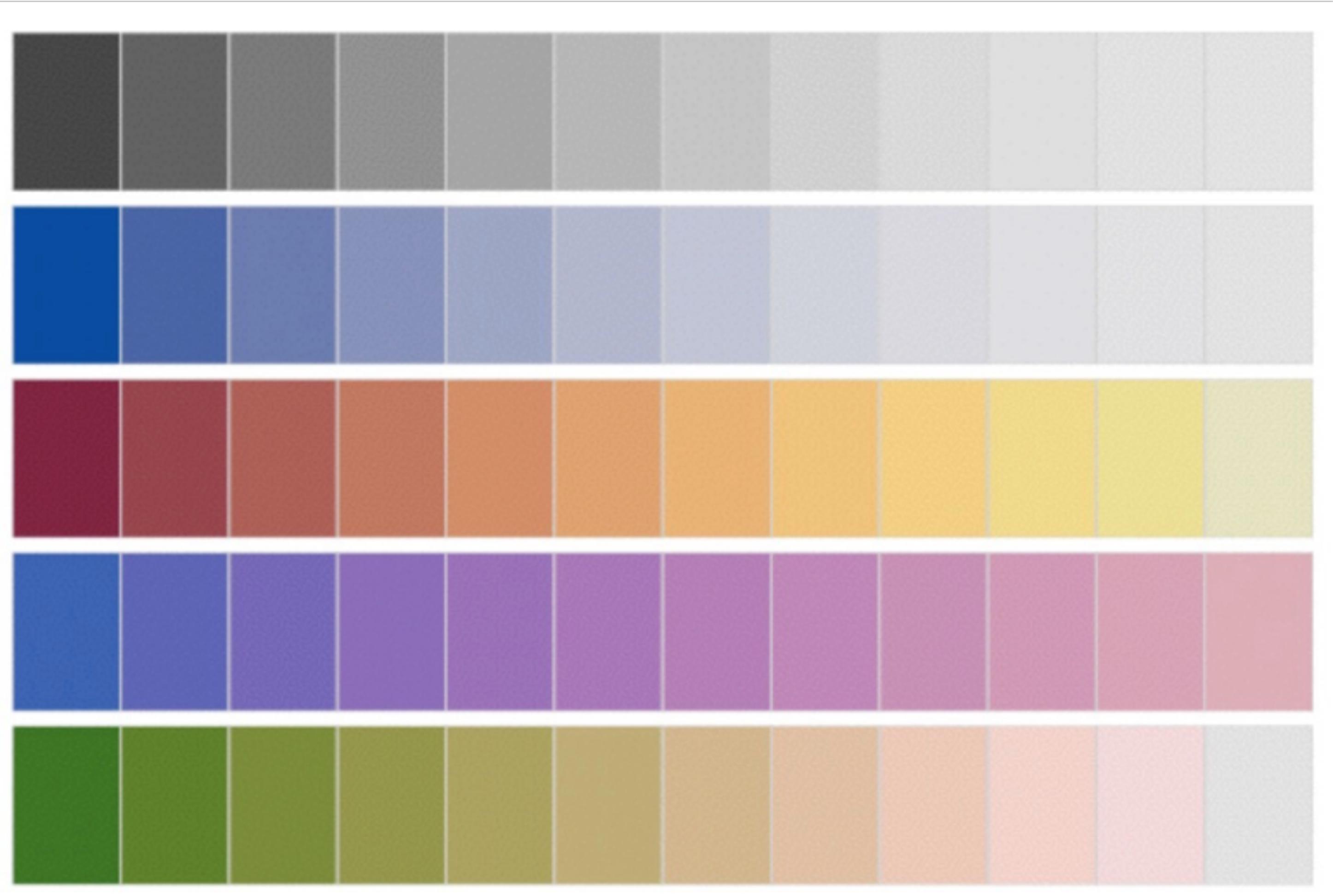
EXPORT

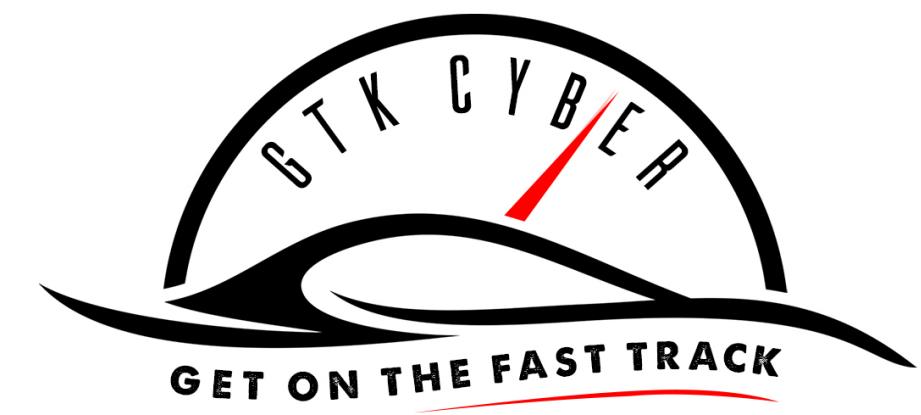
HEX

#e5f5f9
#99d8c9
#2ca25f

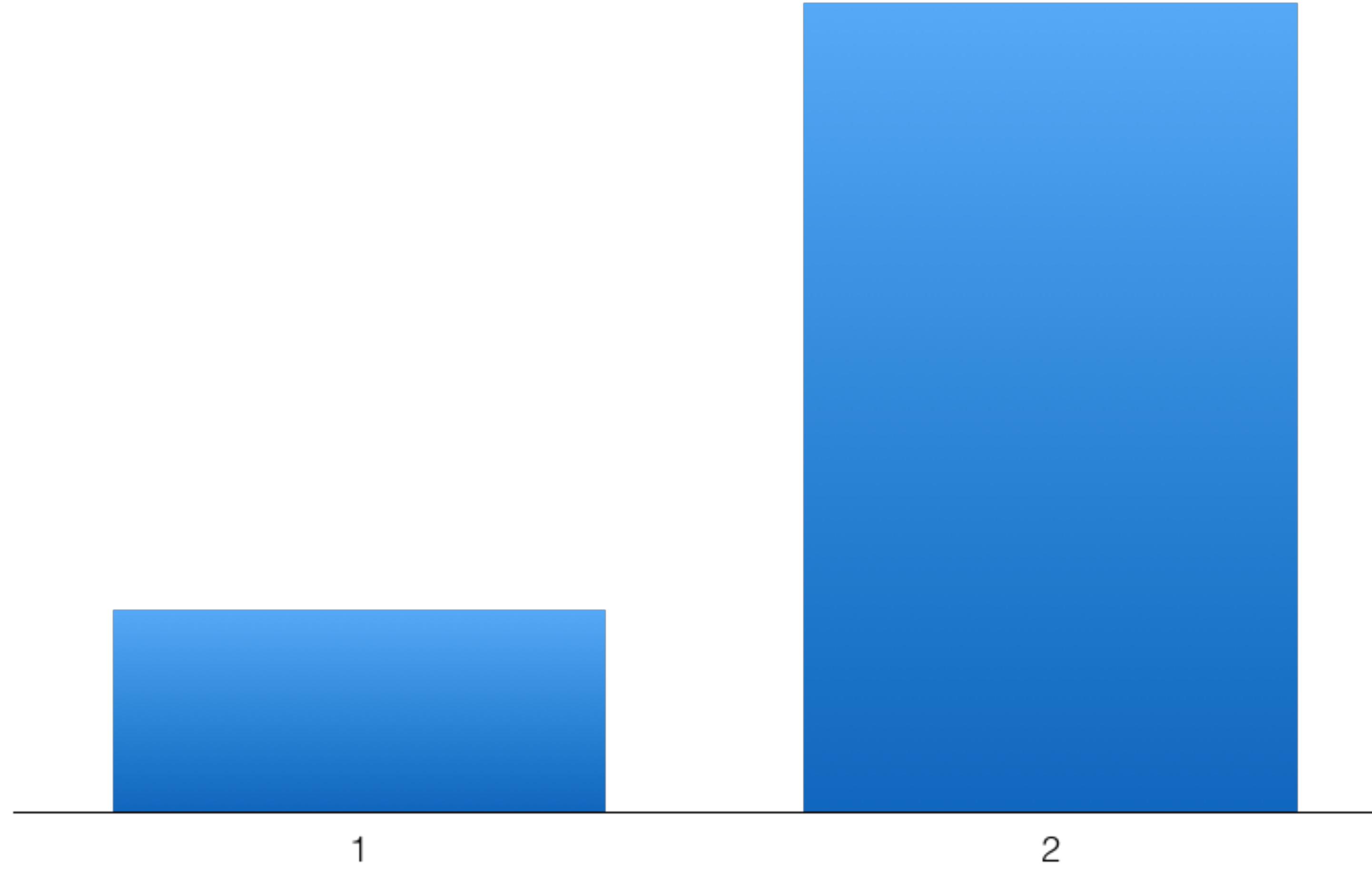


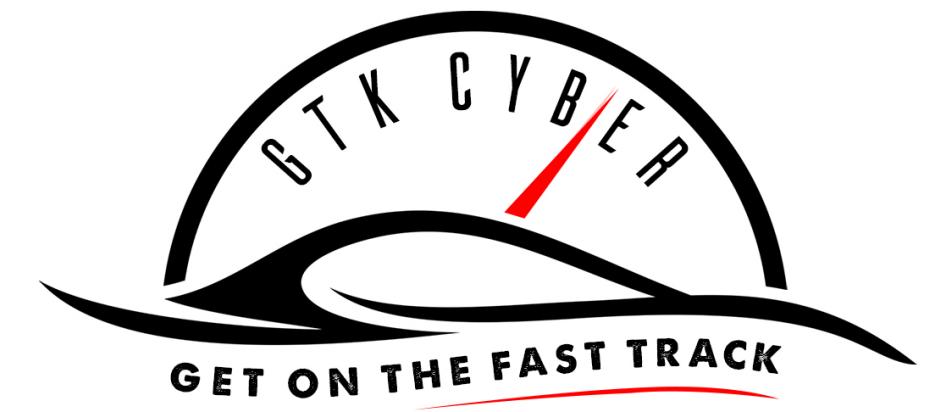
Color



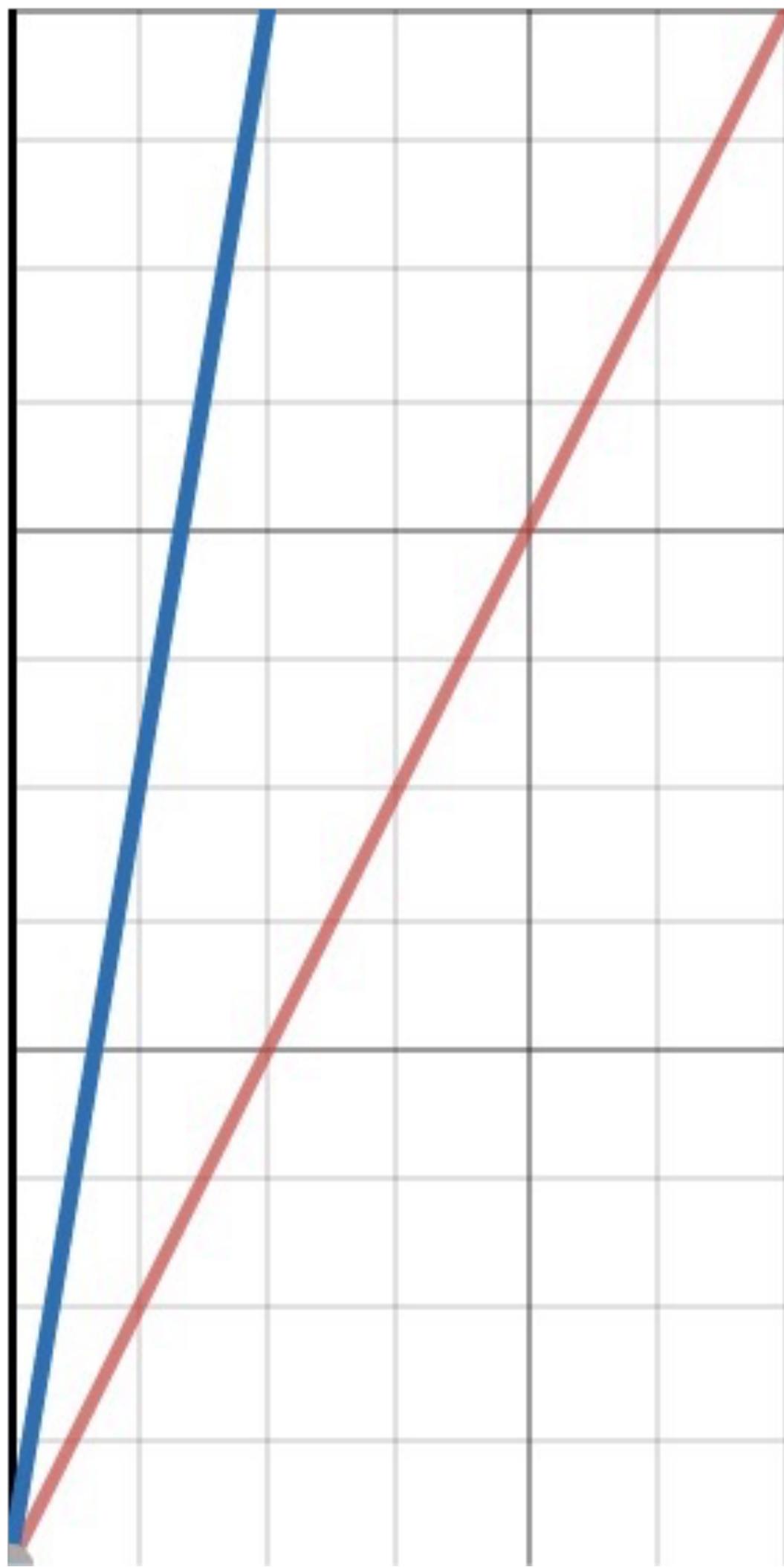


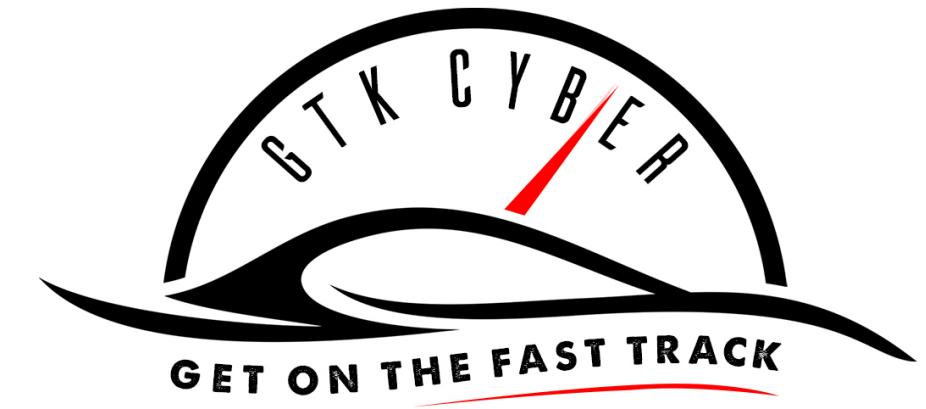
Height Difference?



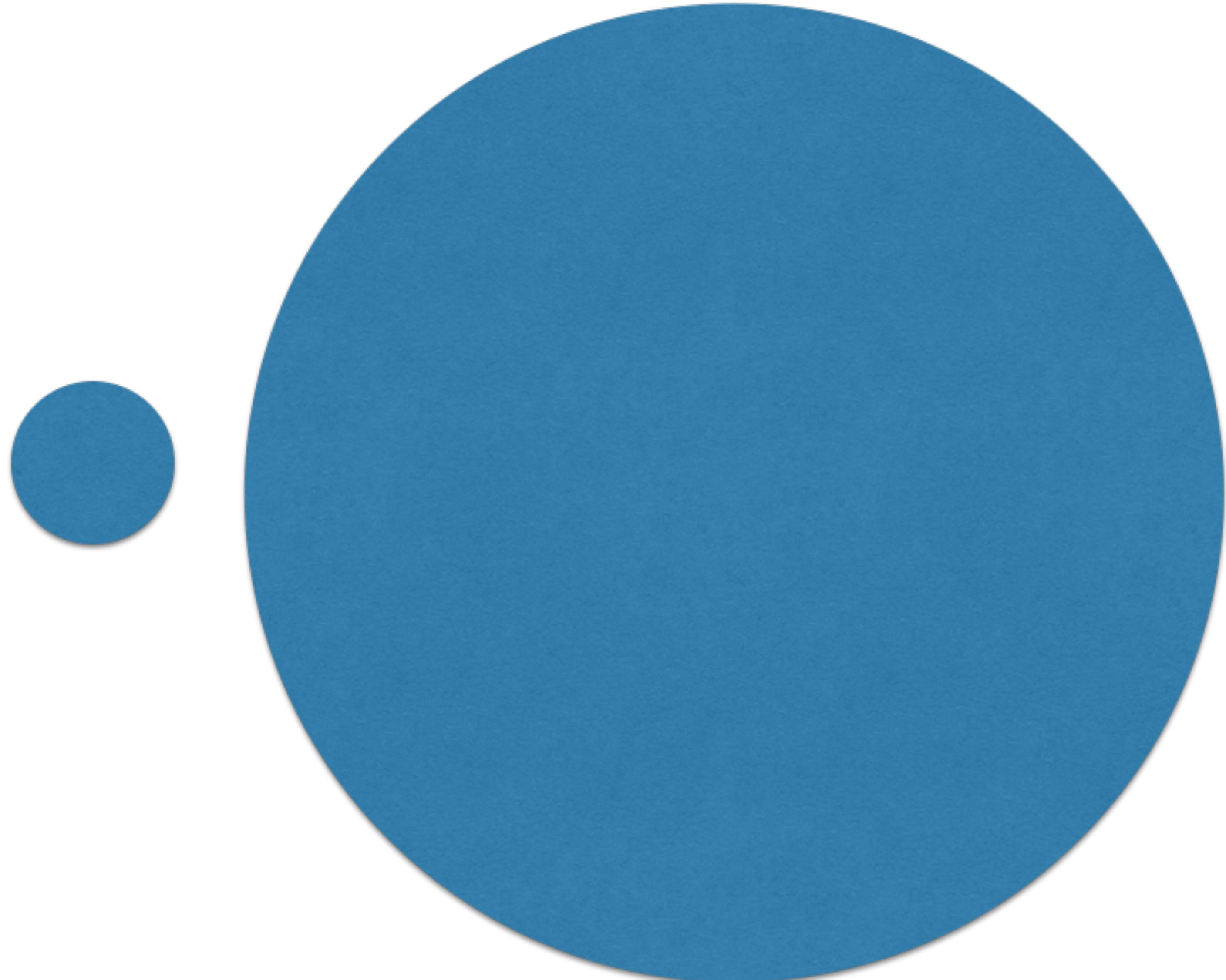


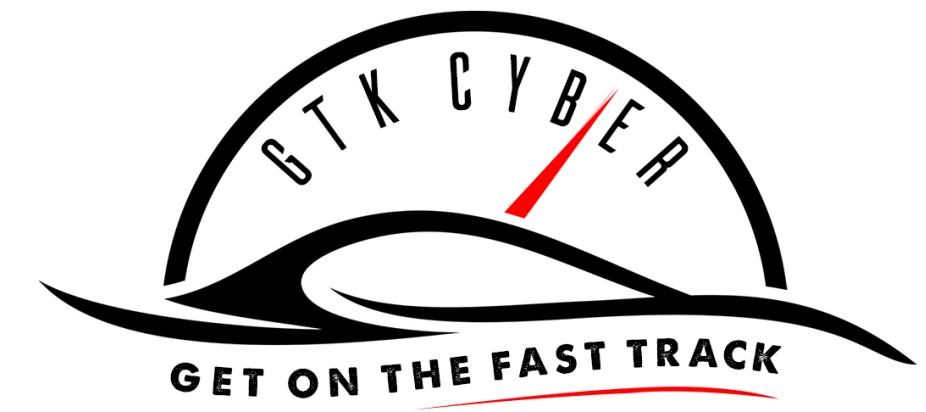
Slope Difference?



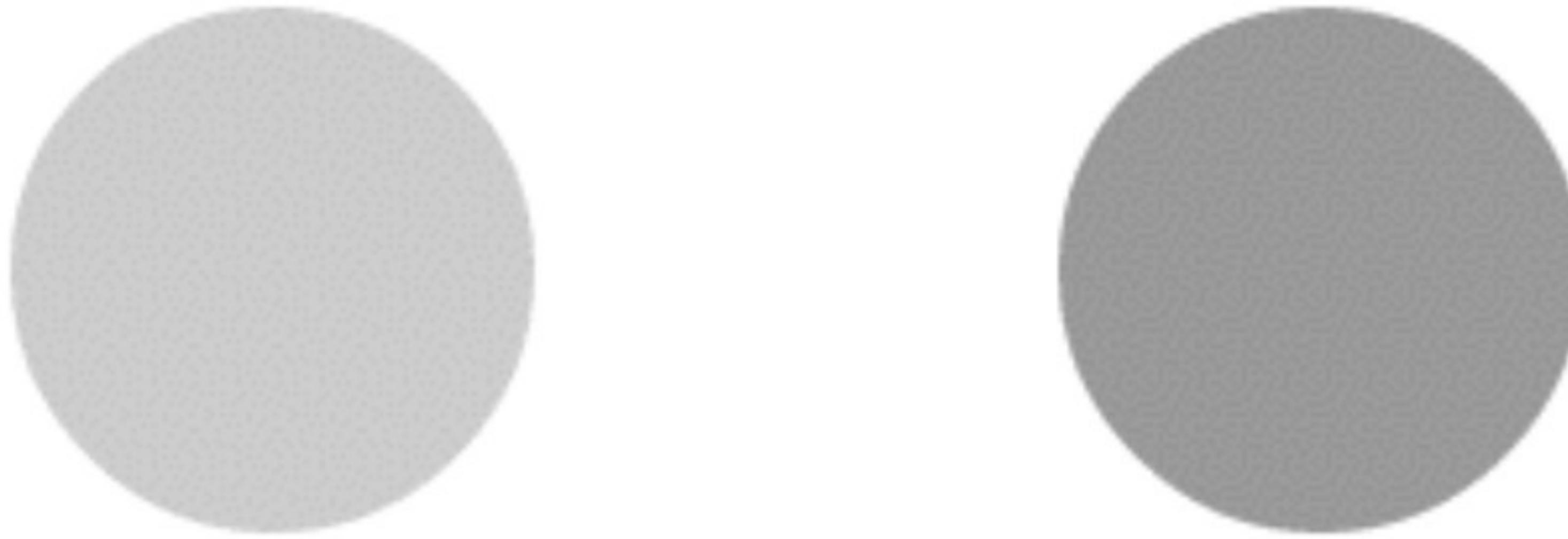


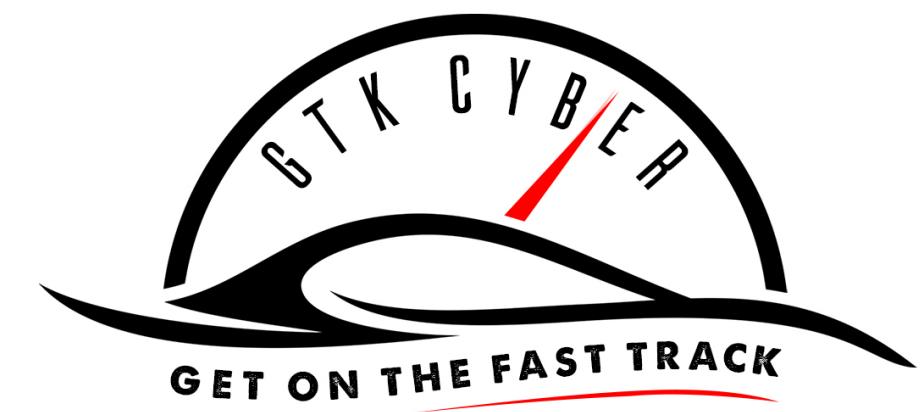
Area Difference?



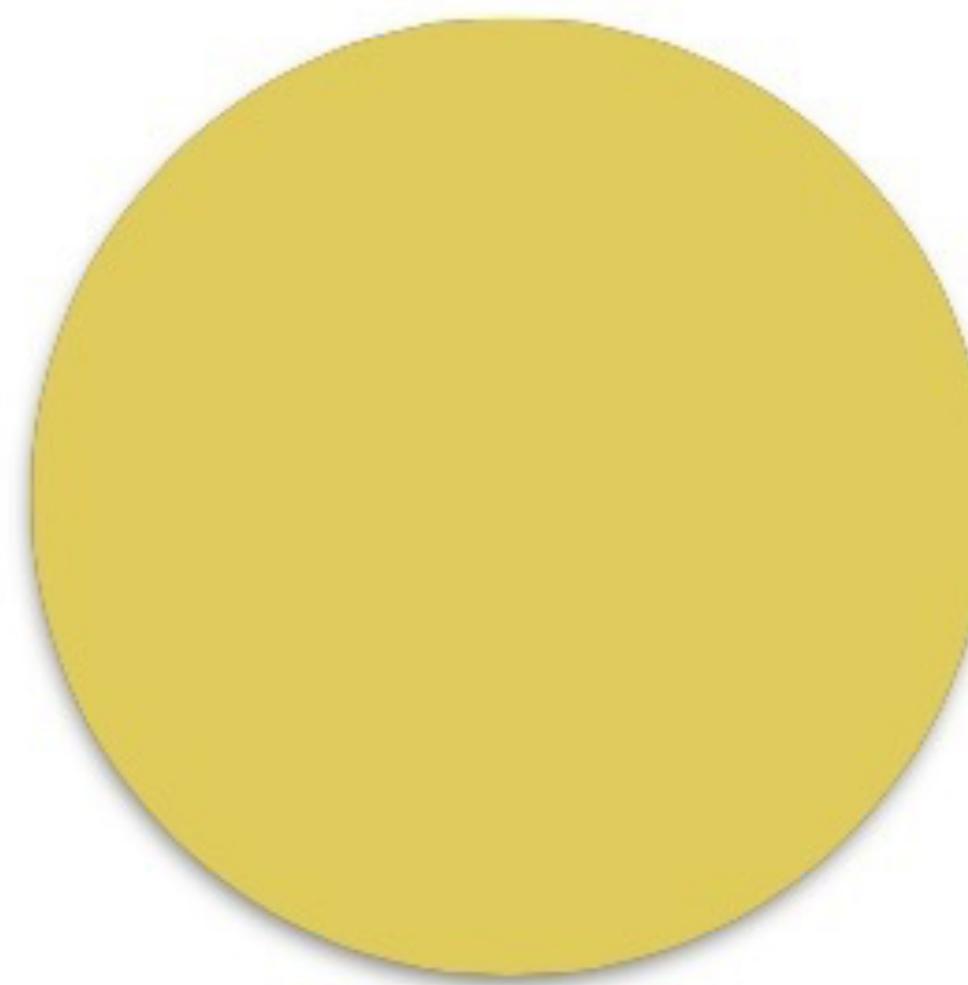
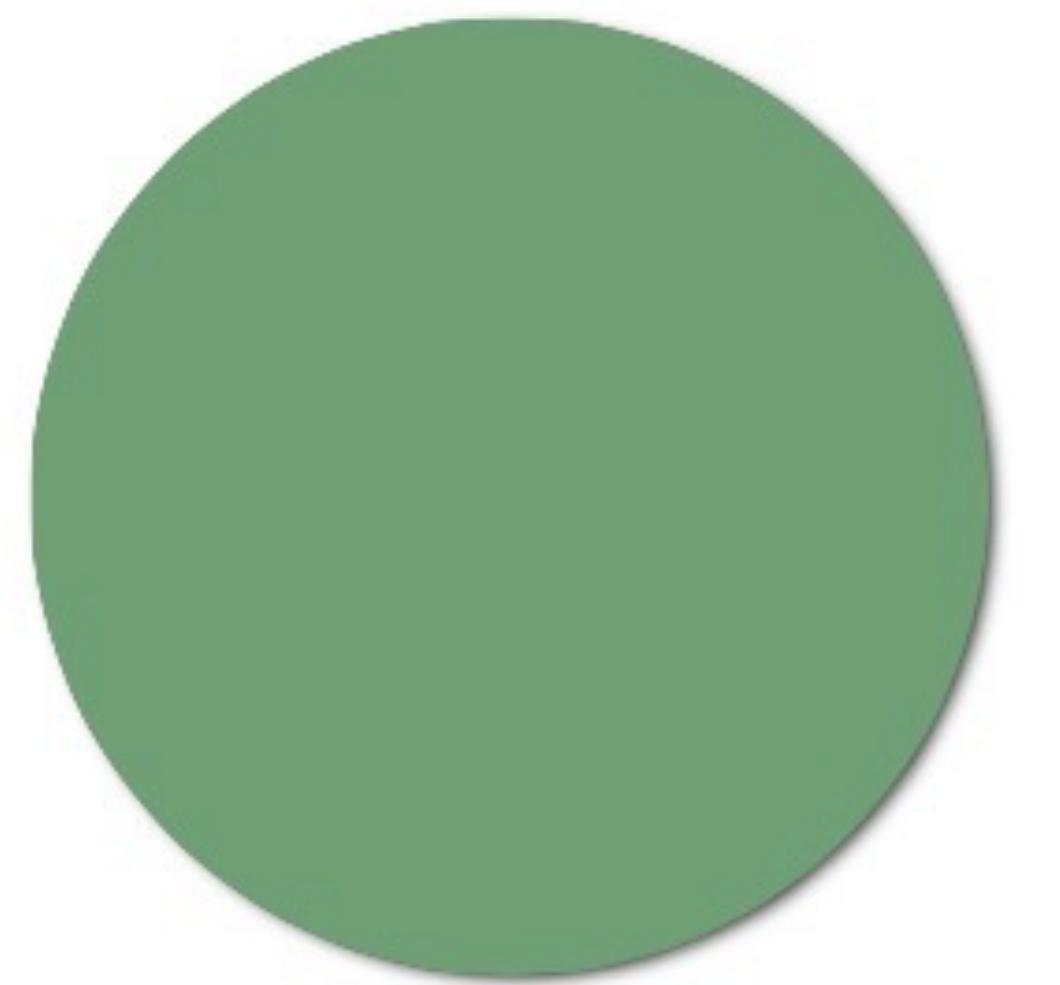


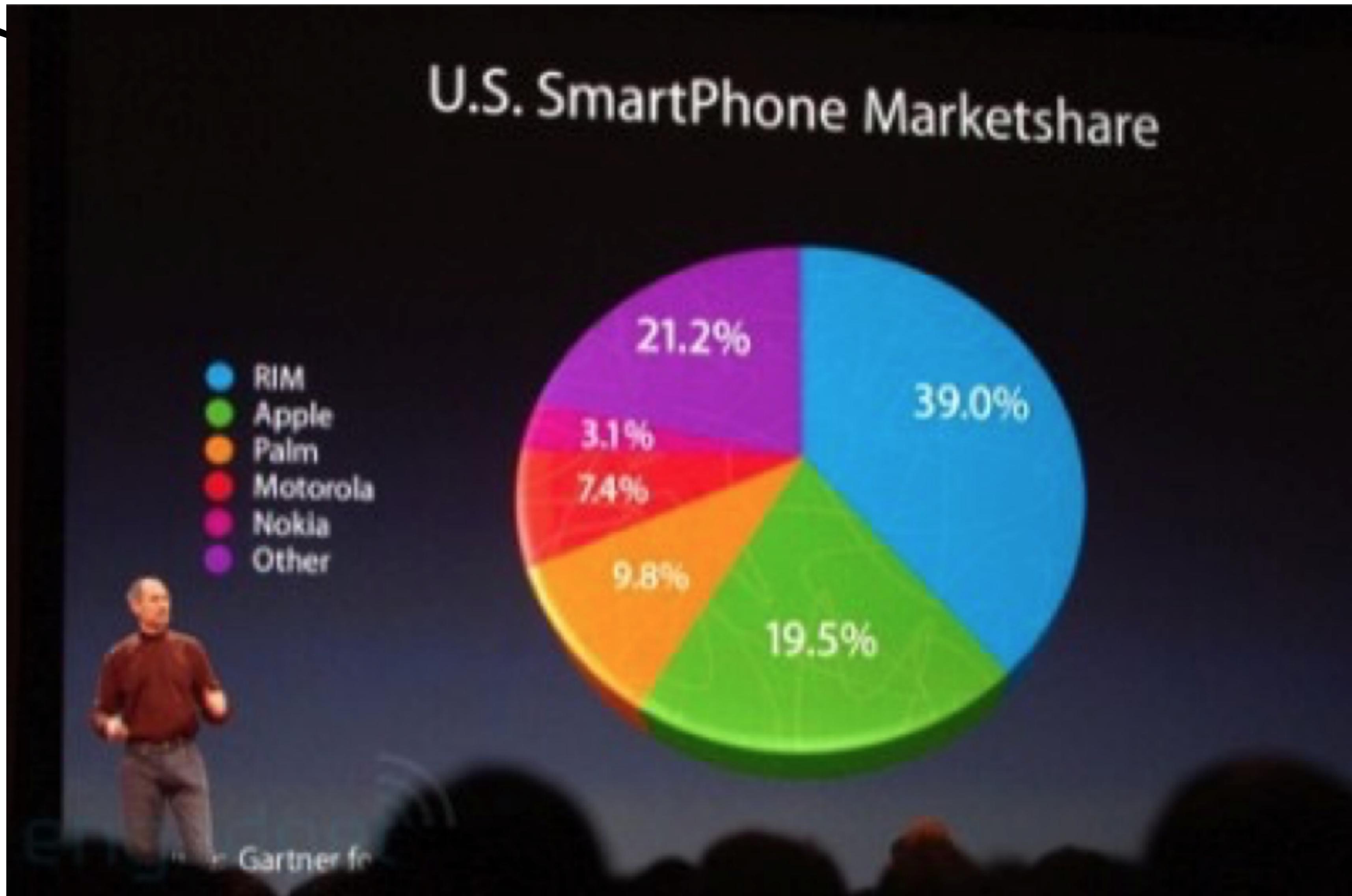
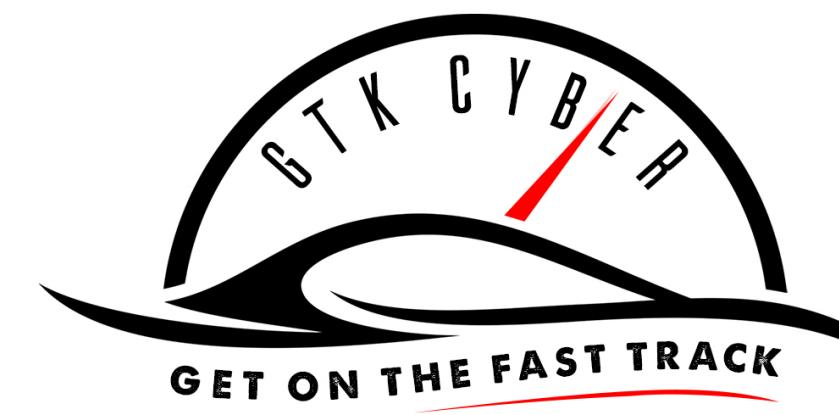
Saturation Difference?





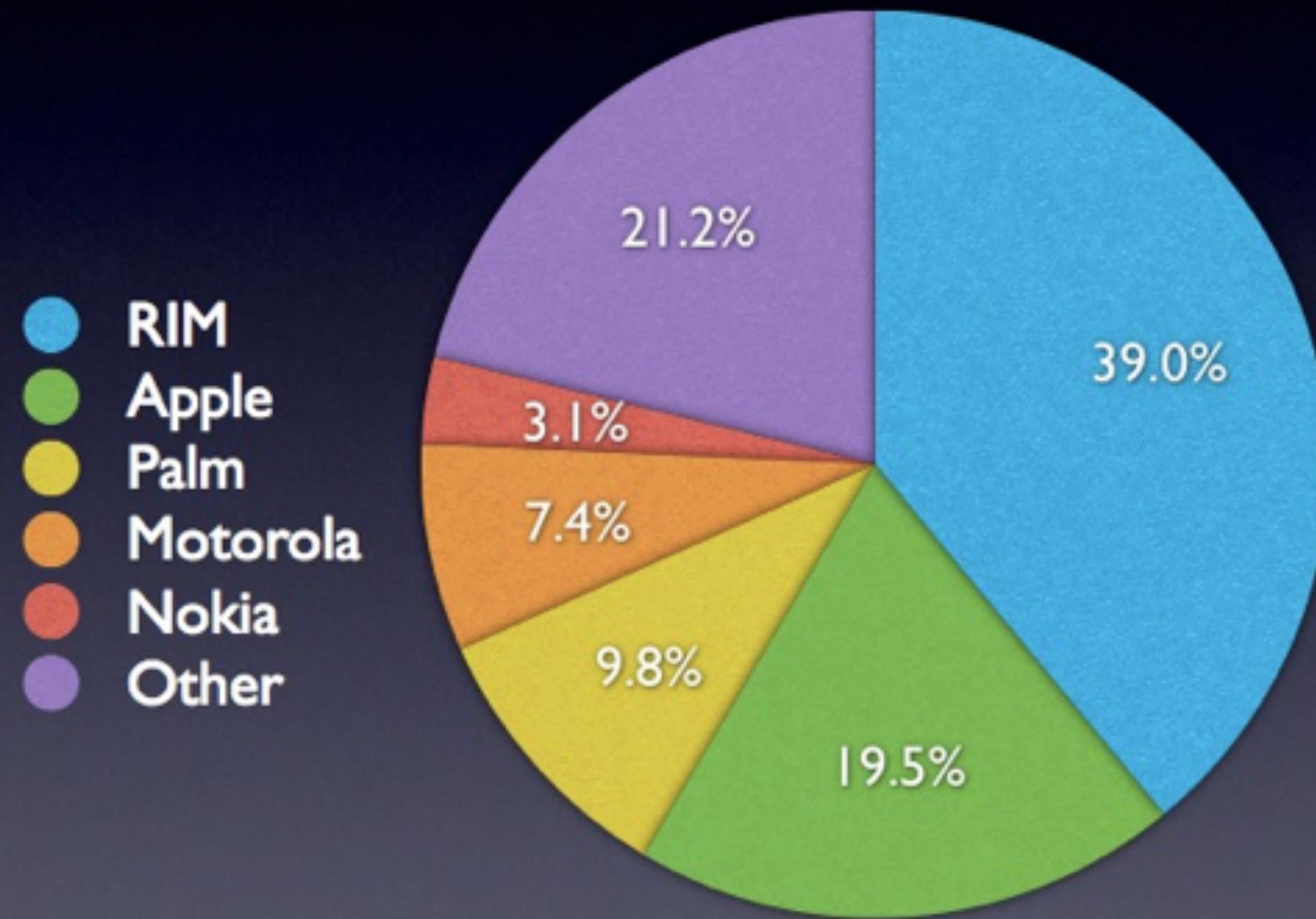
Mapping Difference?

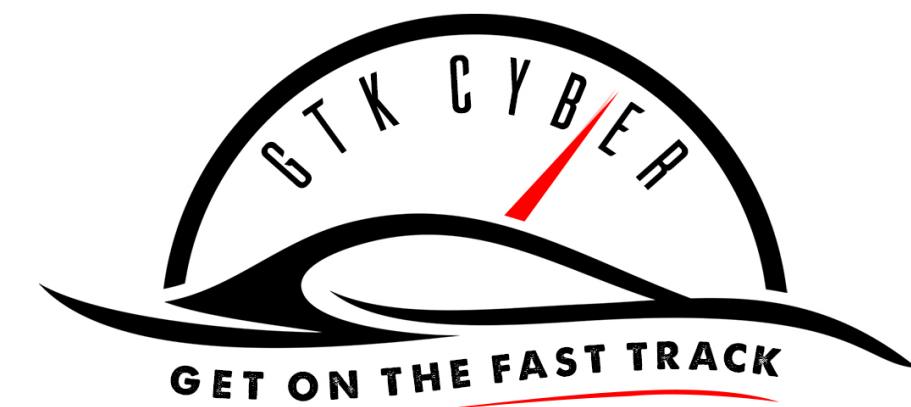






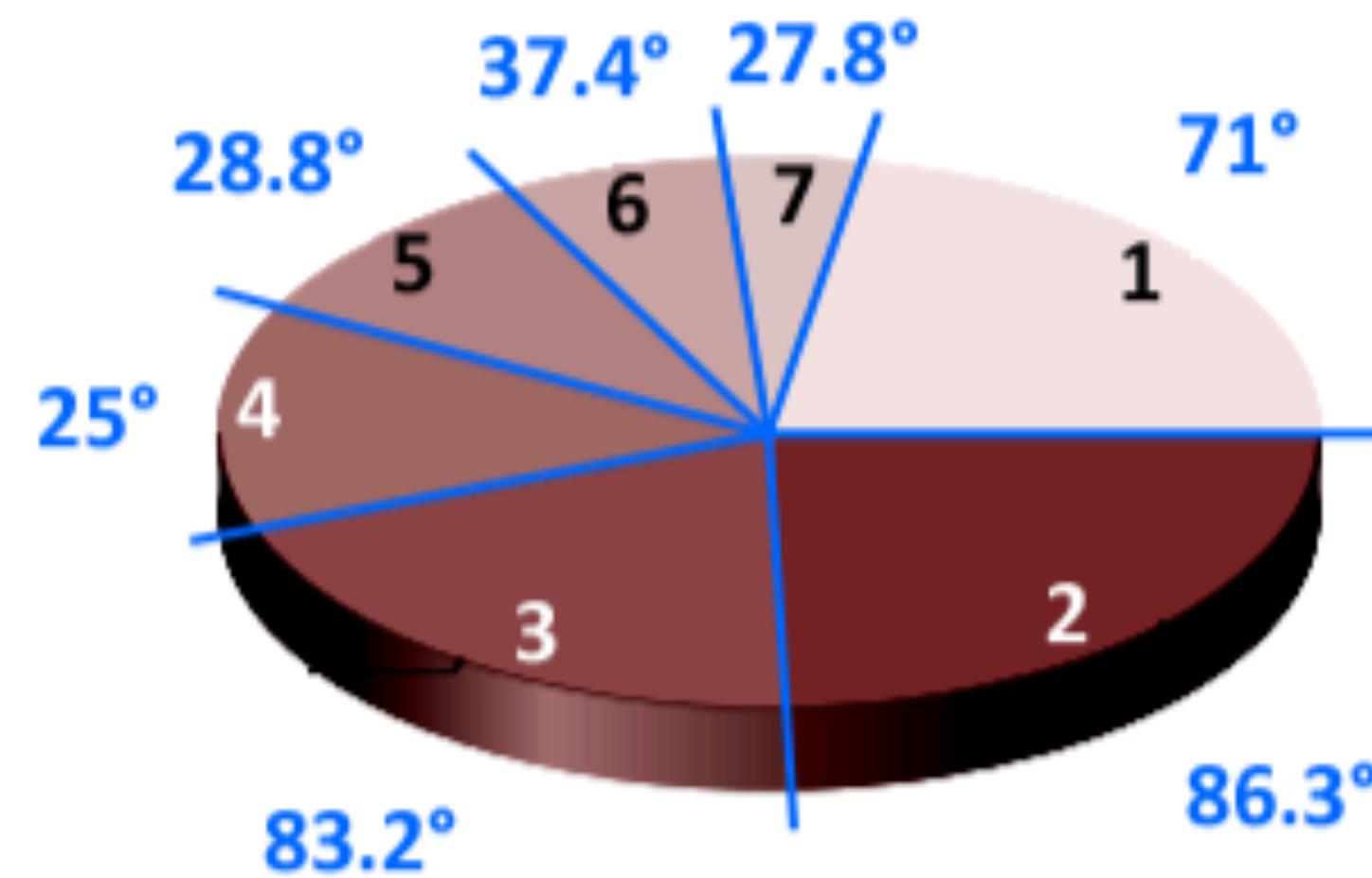
U.S. SmartPhone Marketshare



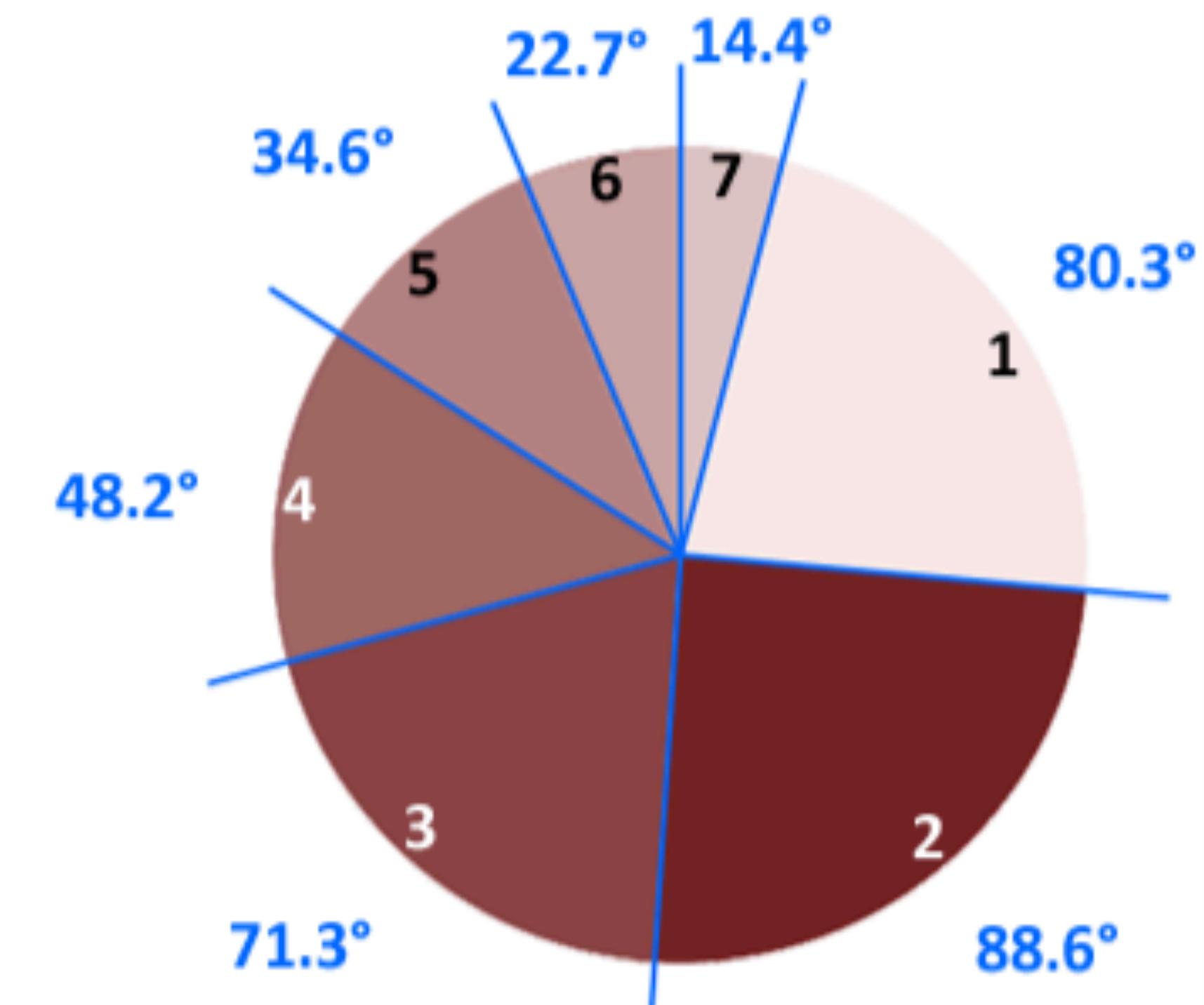


3D Pie Charts are BAD!

3D Pie Charts Distort Angles



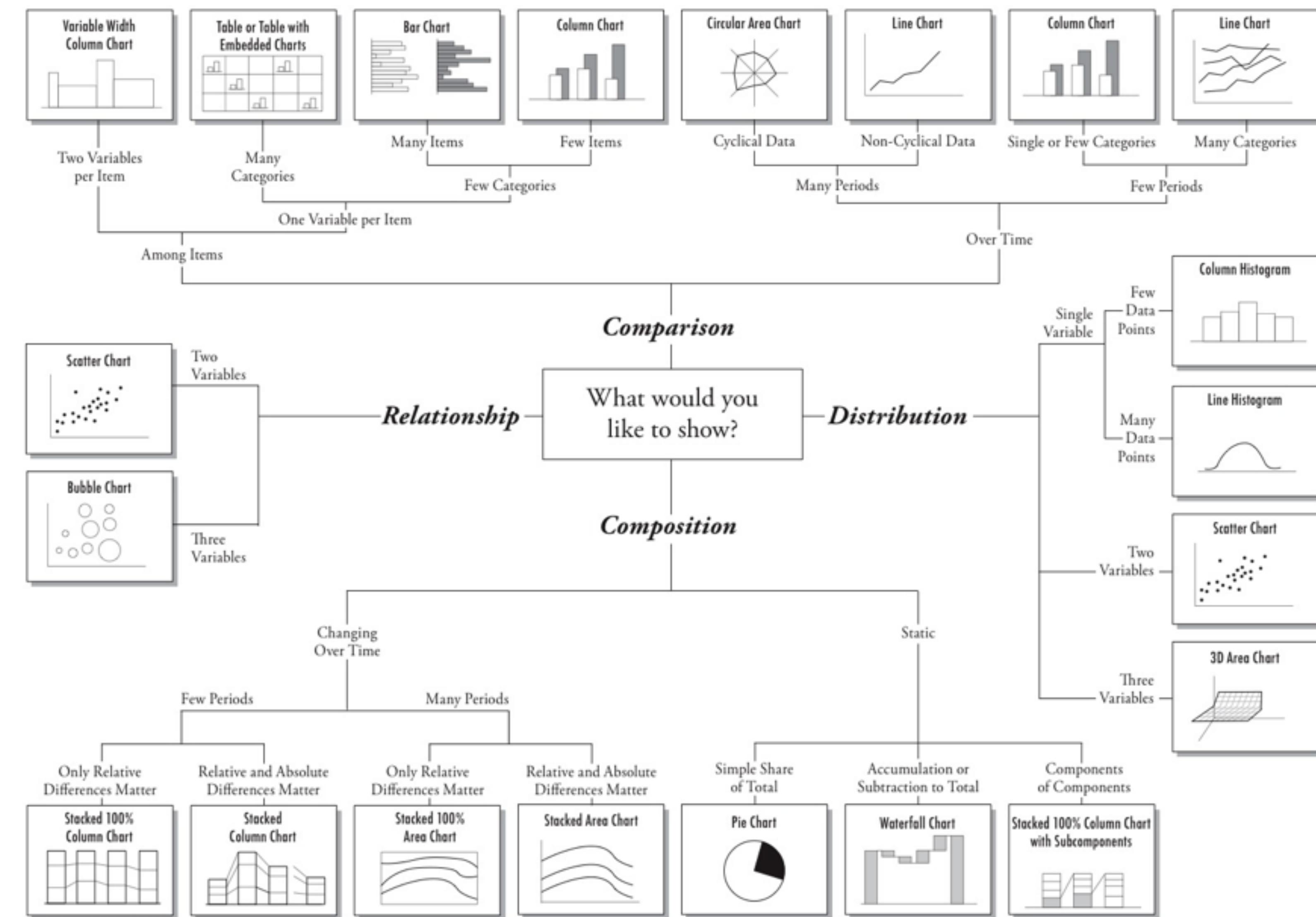
Angles on the original pie chart

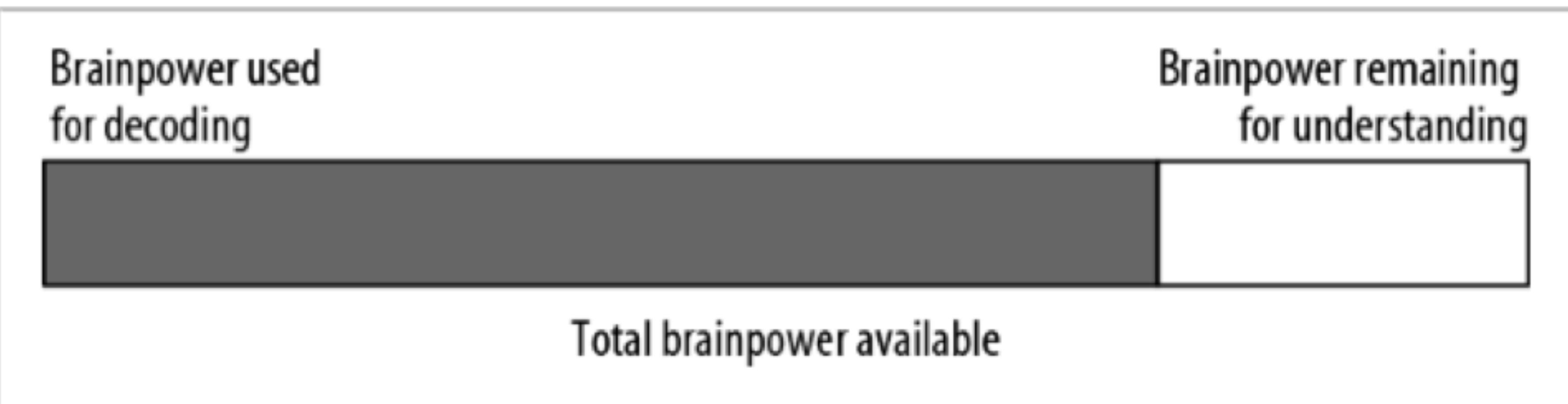
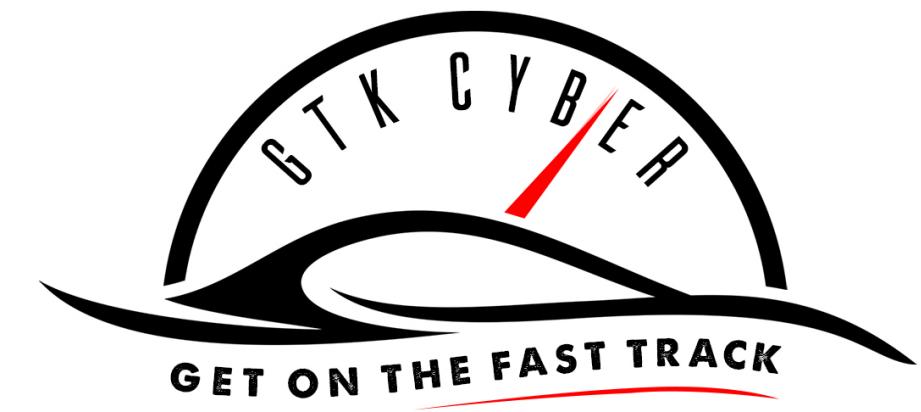


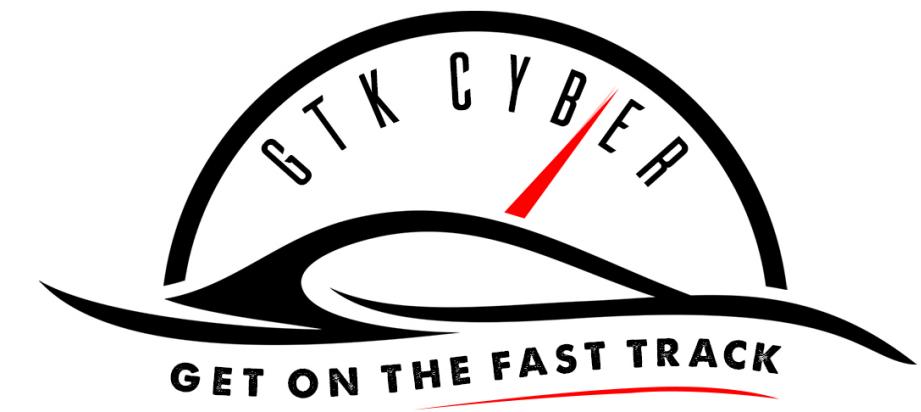
Angles on a non-3D pie chart



Chart Suggestions—A Thought-Starter







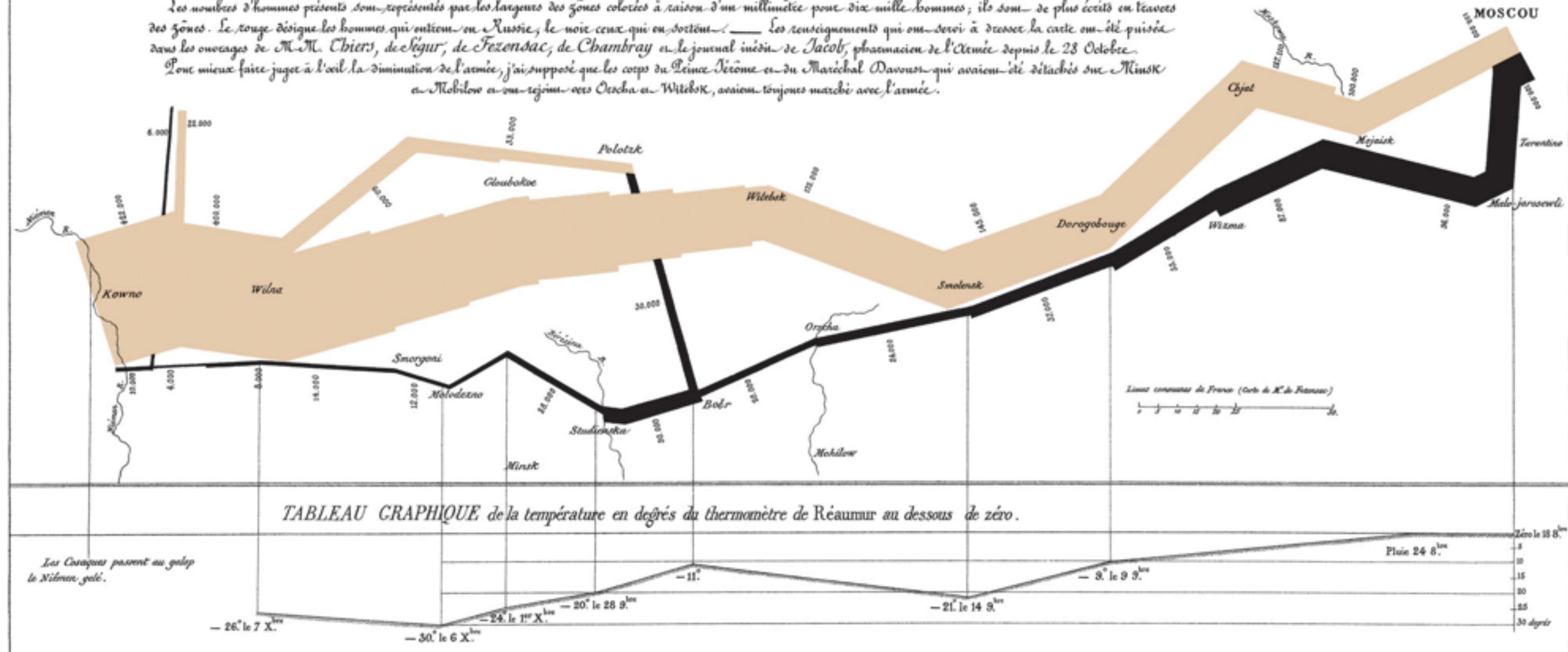


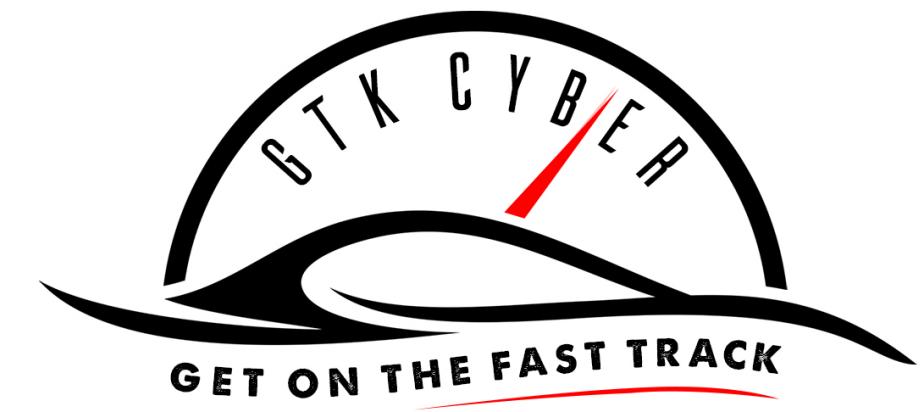
Carte Figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813
 Dressée par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite.

Précisé par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite Paris, le 20 Novembre 1869

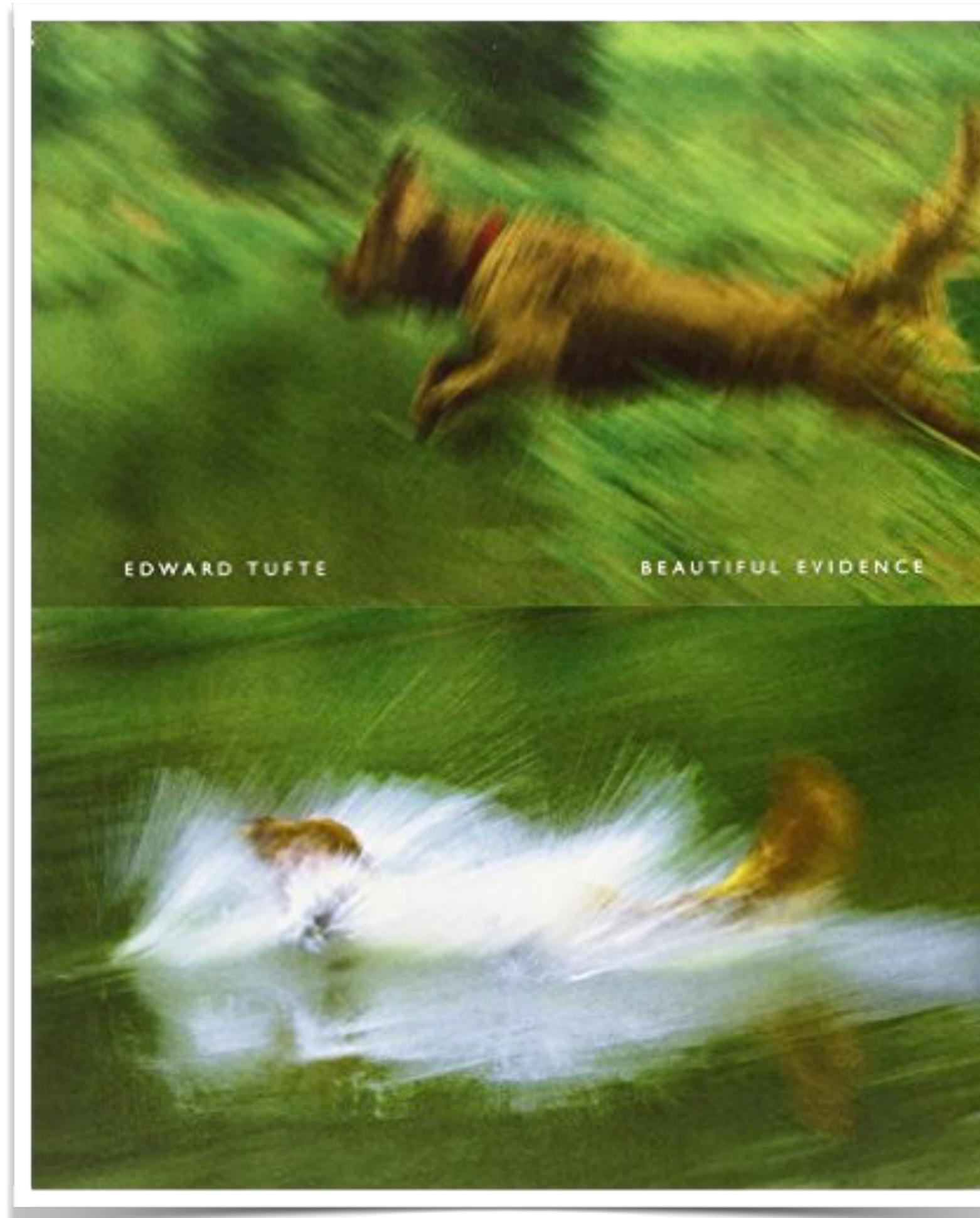
Les nombres d'hommes présents sont représentés par les largesur des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en lettres des zones. Le rouge désigne les hommes qui entrent en Russie, le noir ceux qui en sortent. — Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M. M. Chiers, de Ségur, de Fezensac, de Chambray et le journal intérieur de Jacob, pharmacien de l'Académie depuis le 28 Octobre.

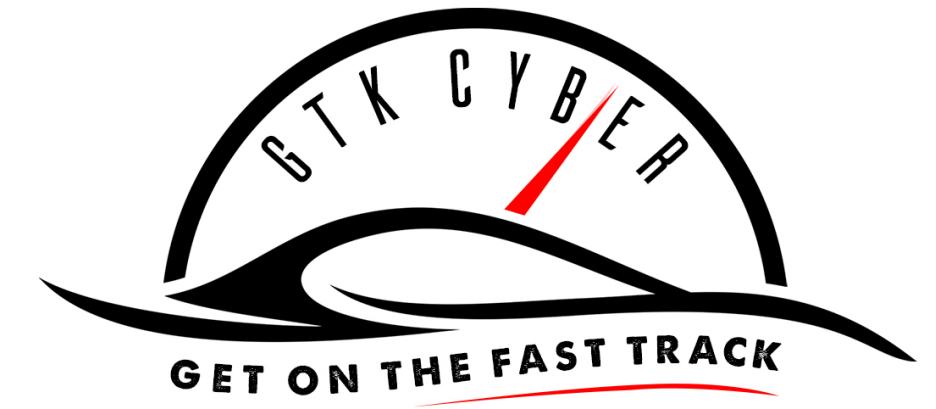
Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps du Prince Jérôme et du Maréchal Davout, qui avaient été détachés sur Minsk et Mogilow et qui rejoignaient Osscha et Witlobk, avaient toujours marché avec l'armée.



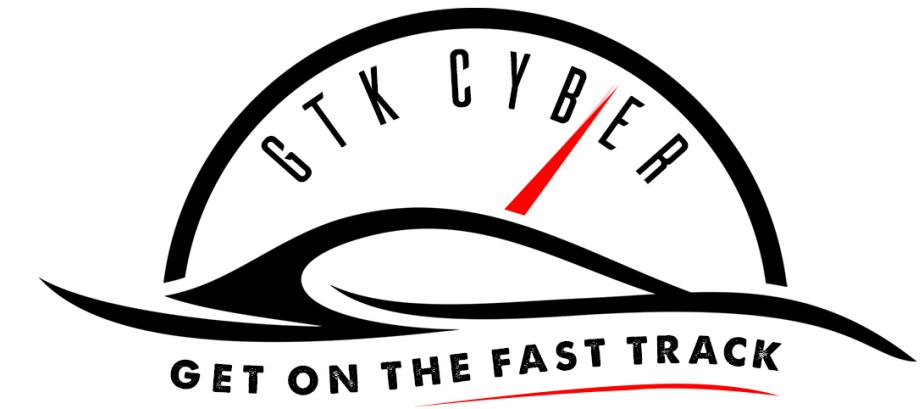


Recommended Reading





Storytelling



Storytelling

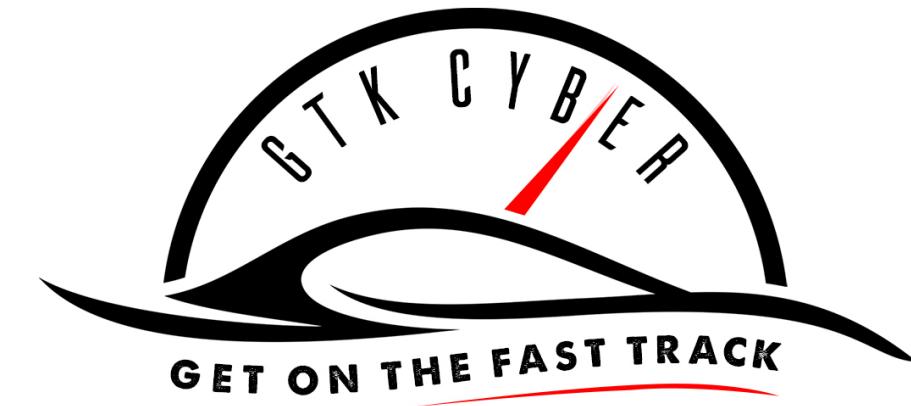
What is the Question?

Who Cares?



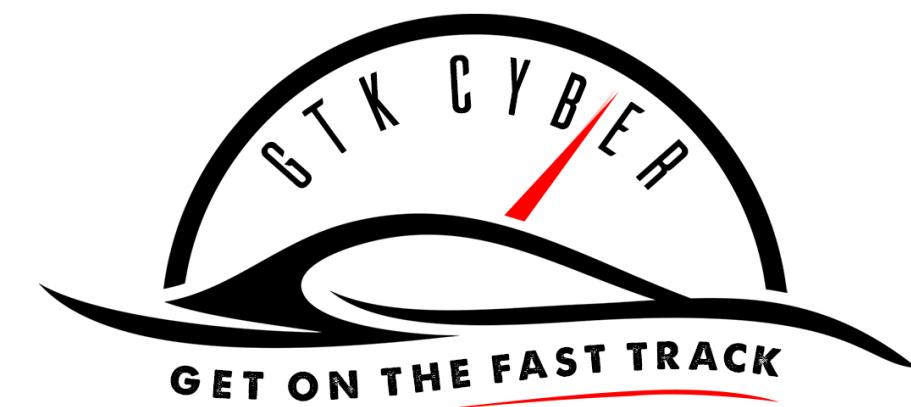
Storytelling

1. Know Your Audience
2. Tell a Story
3. Don't Confuse
4. Go with the Flow

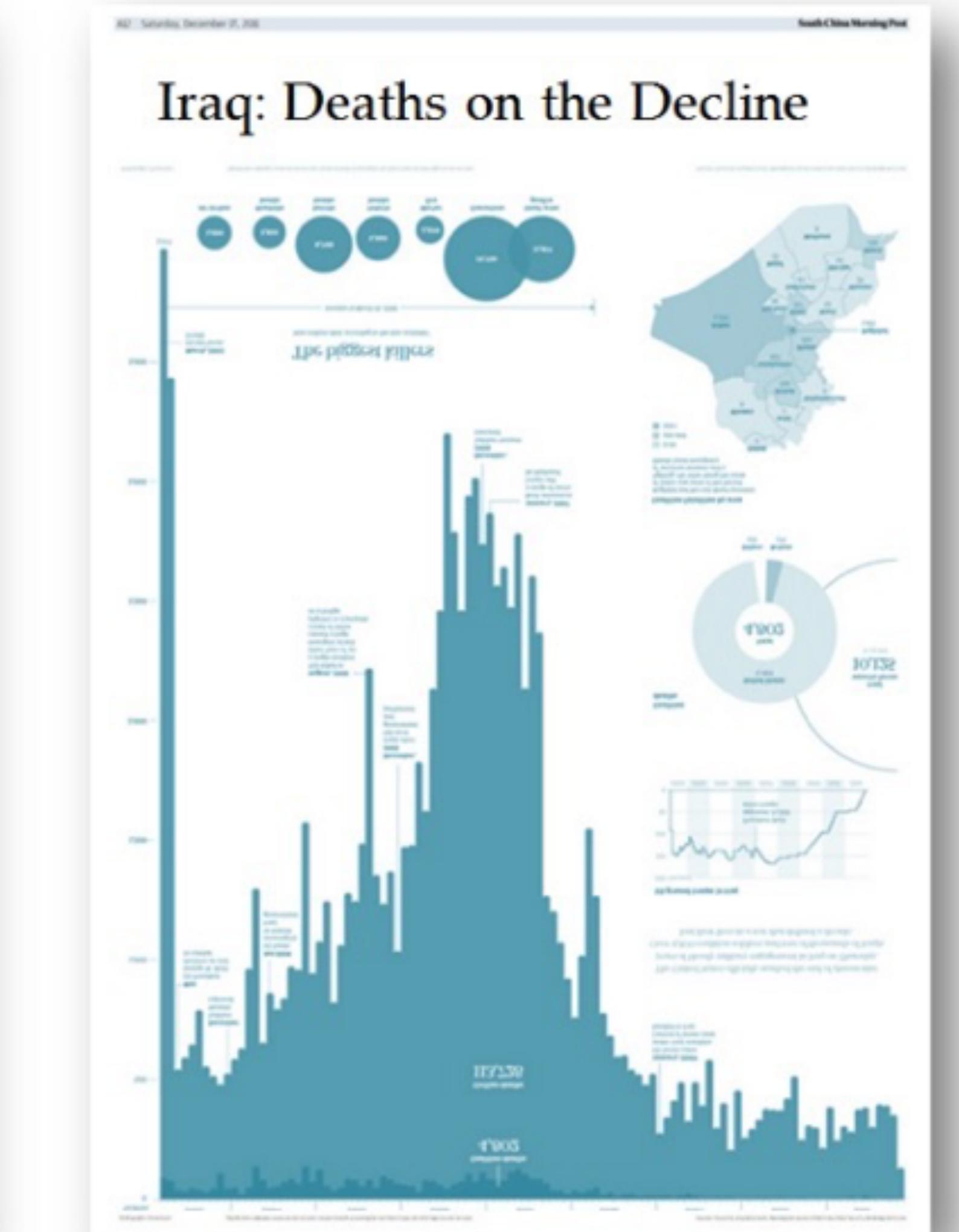
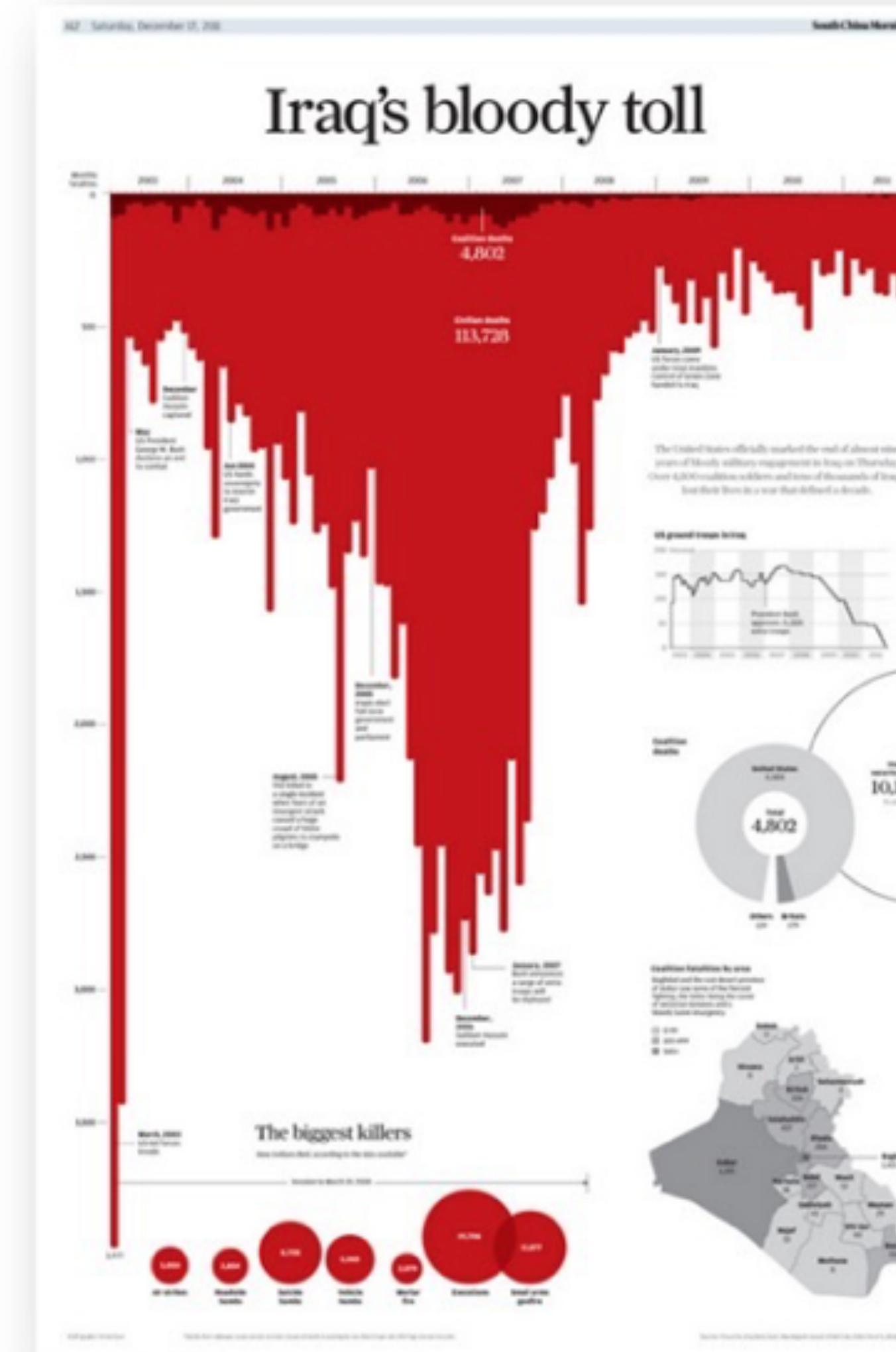


Elements of a Story

1. Characters
2. Conflict
3. Resolution
4. Sequel?



Influential Storytelling

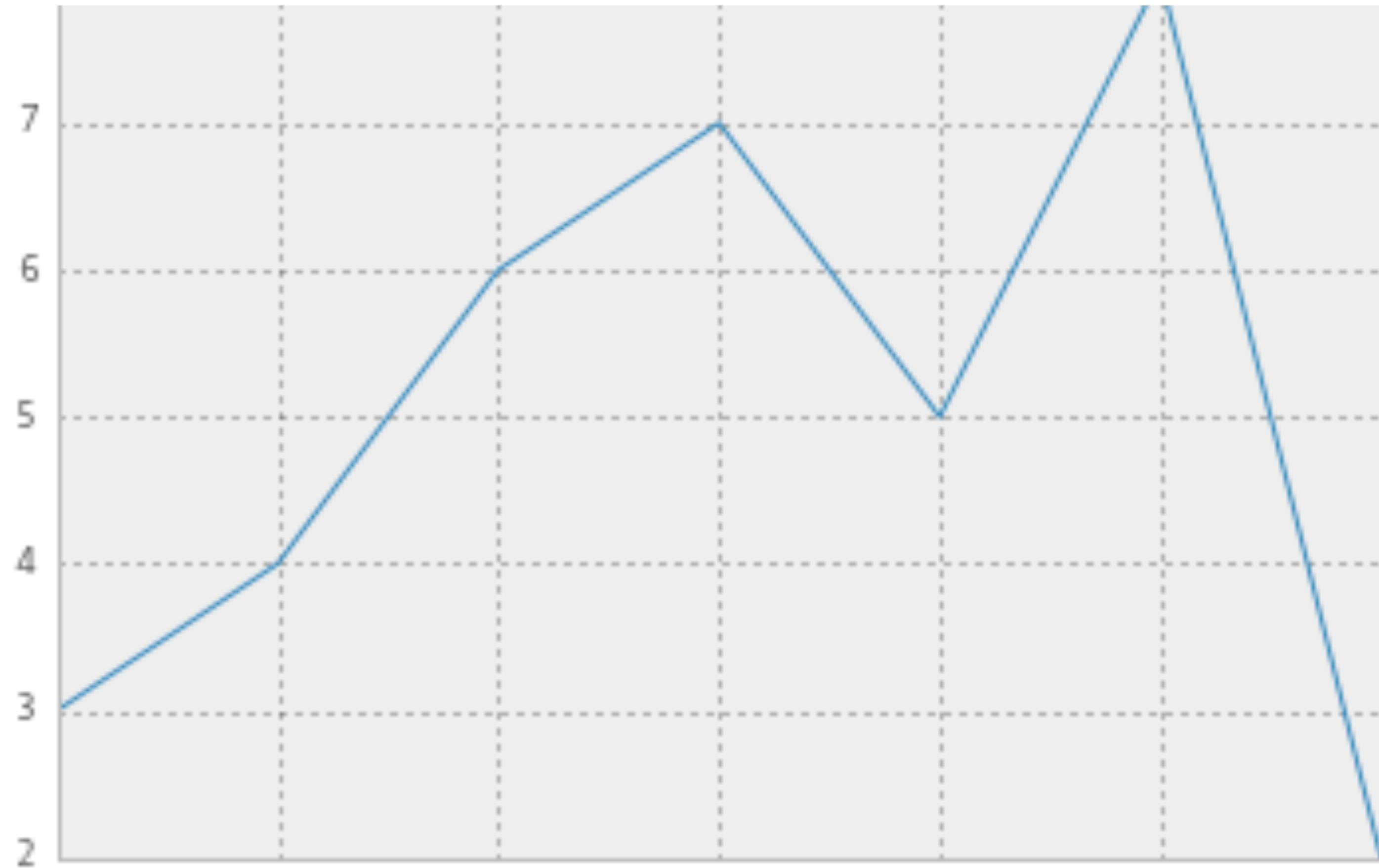
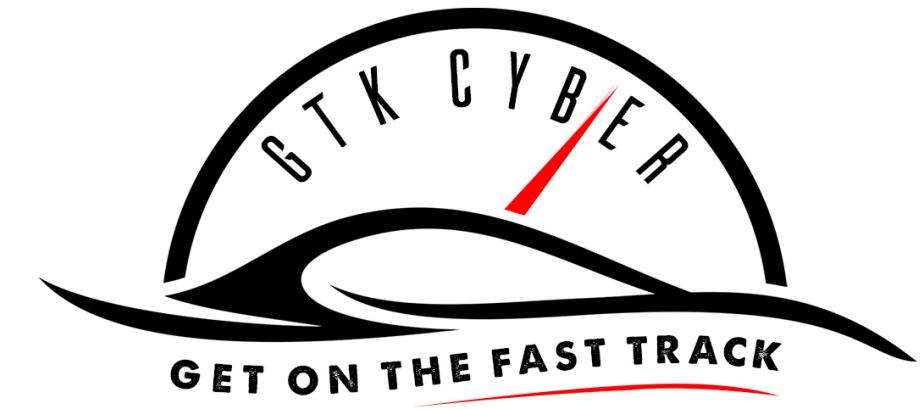


Andy Cotgreave

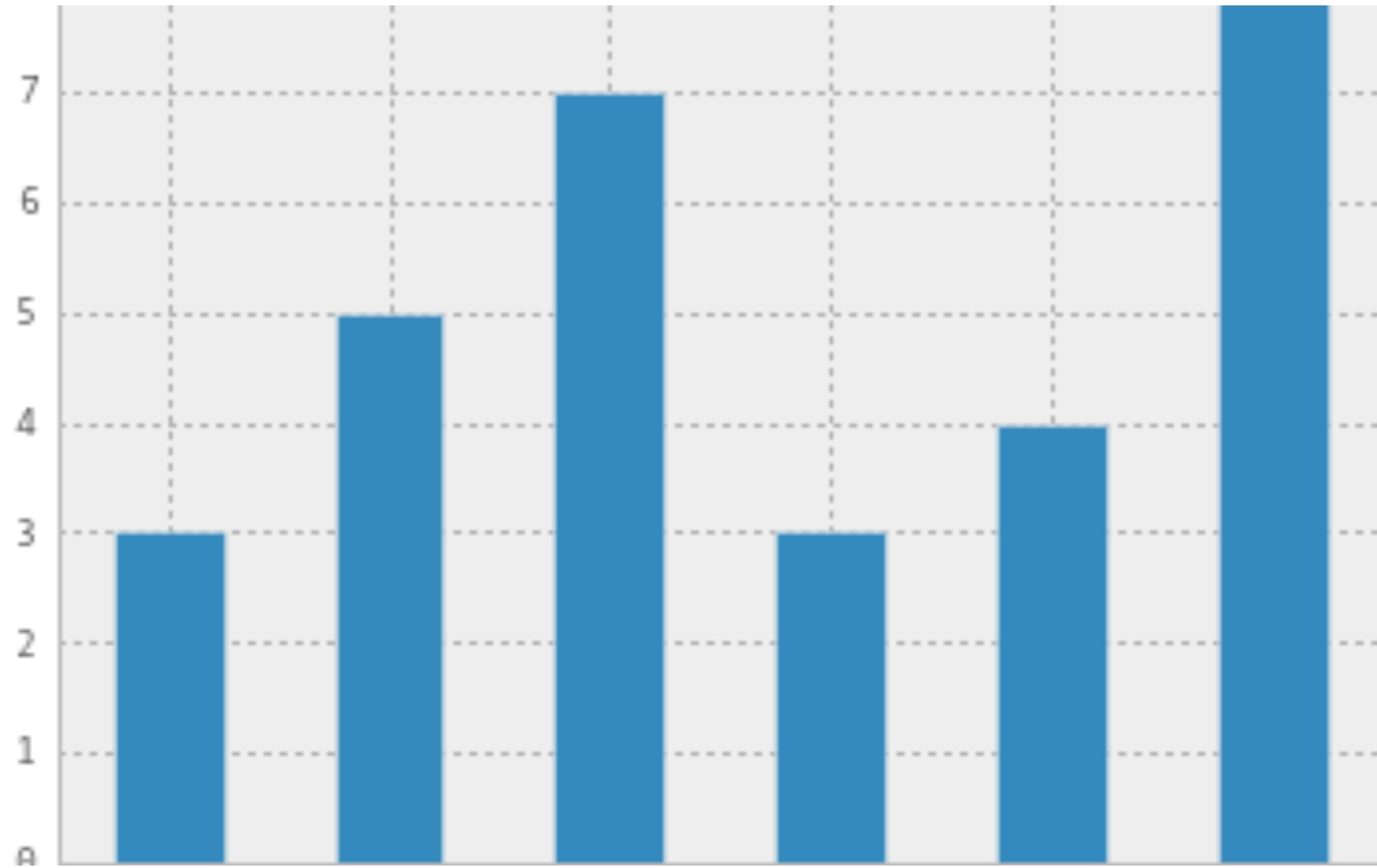
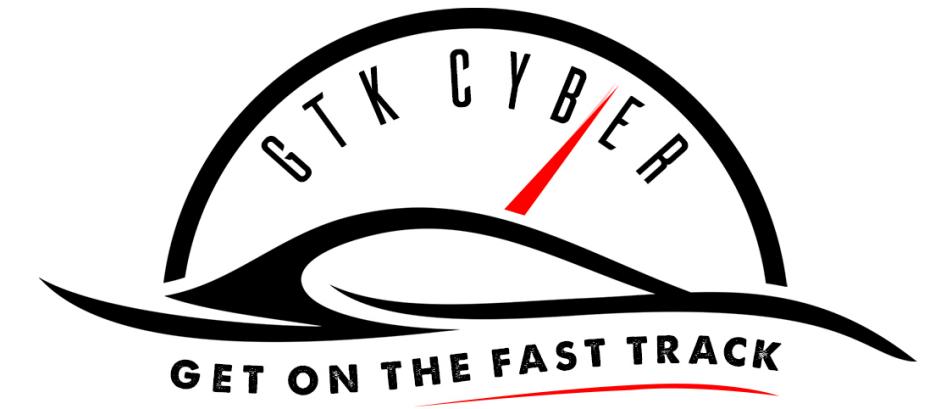
Data Visualization in Python

GET ON THE FAST TRACK

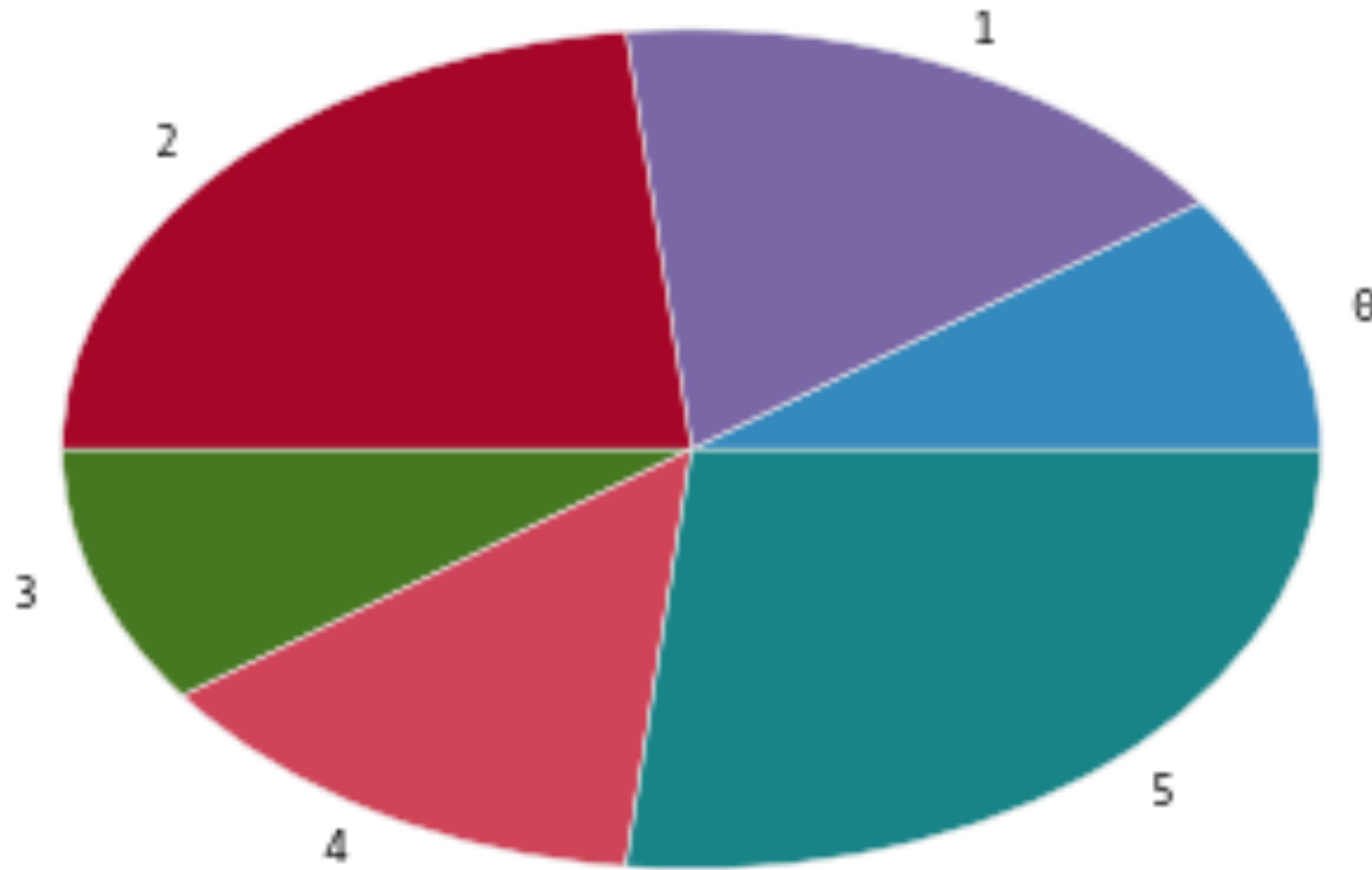
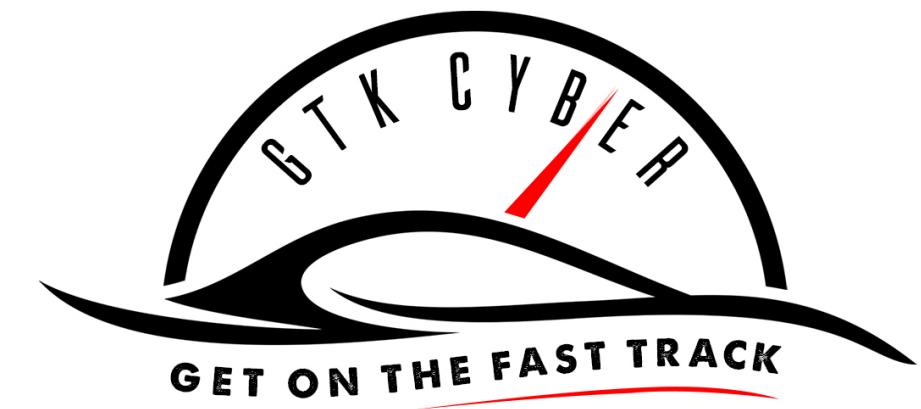
```
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
%matplotlib inline
pd.options.display.mpl_style = 'default'
```



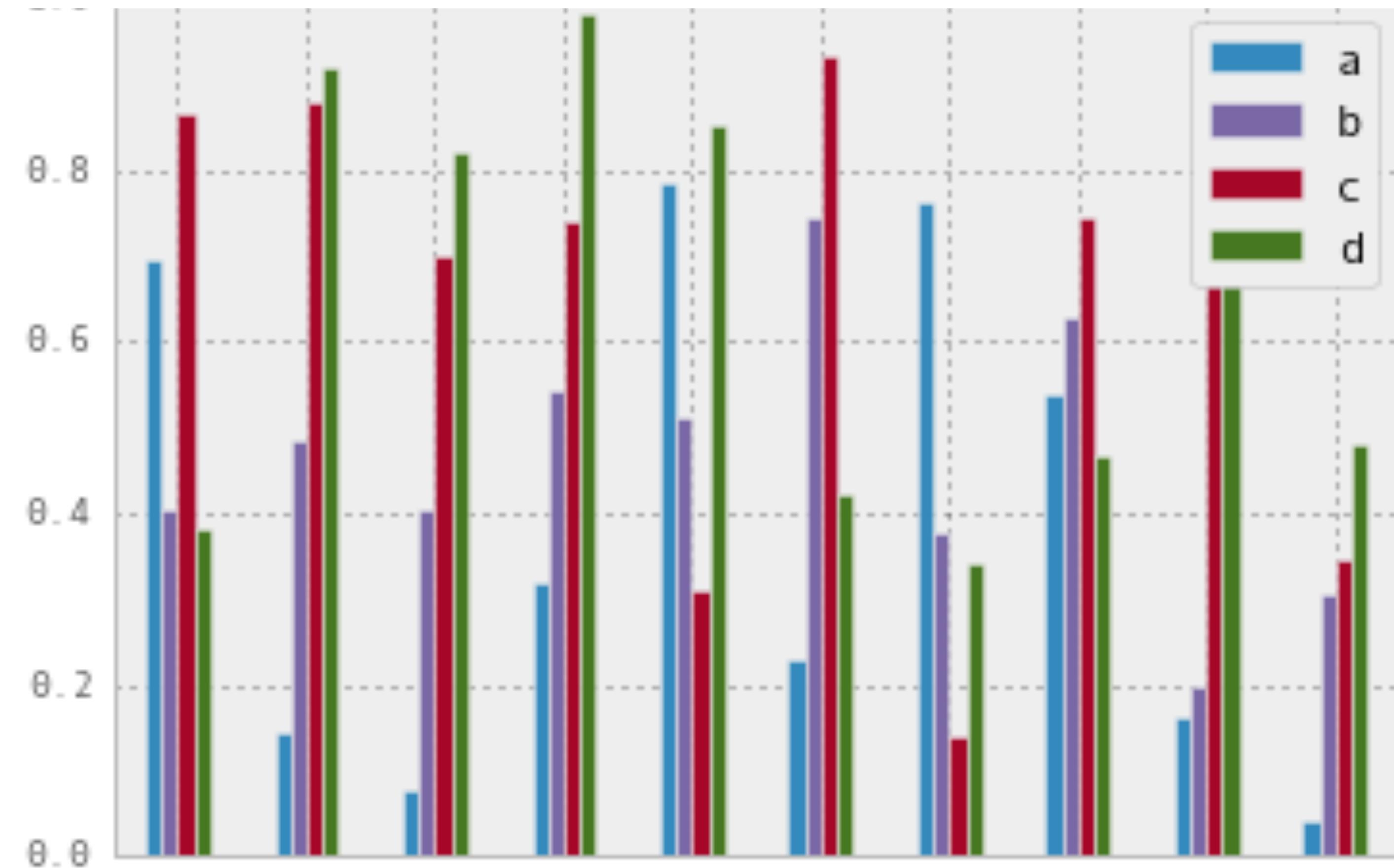
```
data = pd.Series( [ 3, 4, 6, 7, 5, 8, 2 ] )  
graph = data.plot()
```



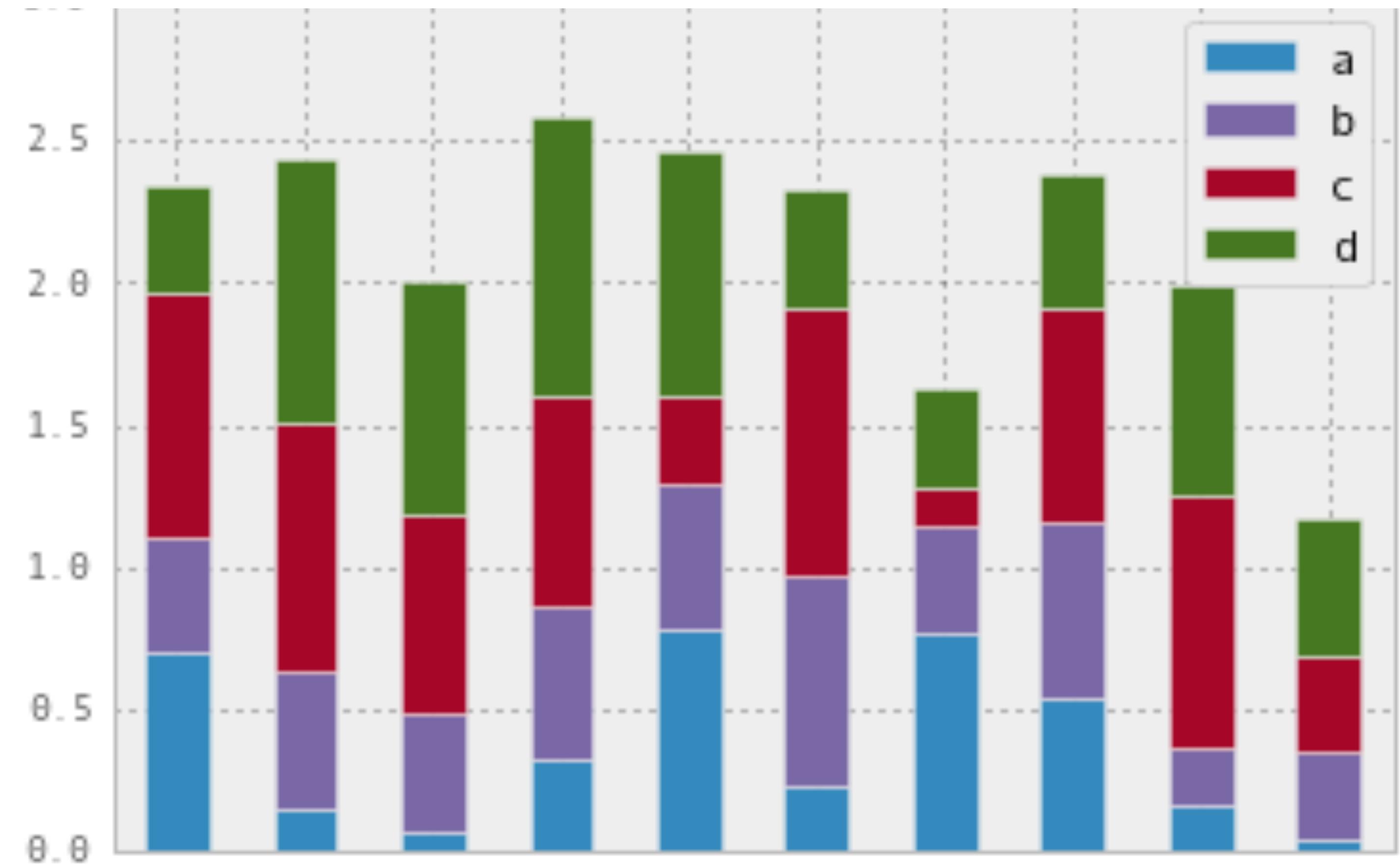
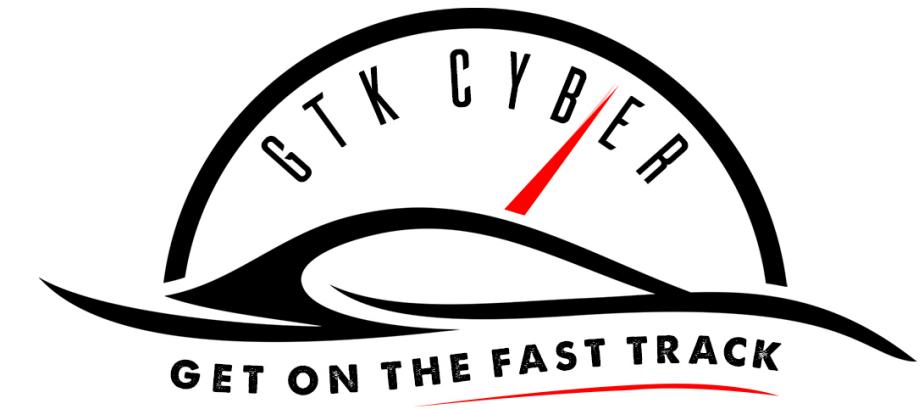
```
barchart = data.plot( kind="bar" )
```



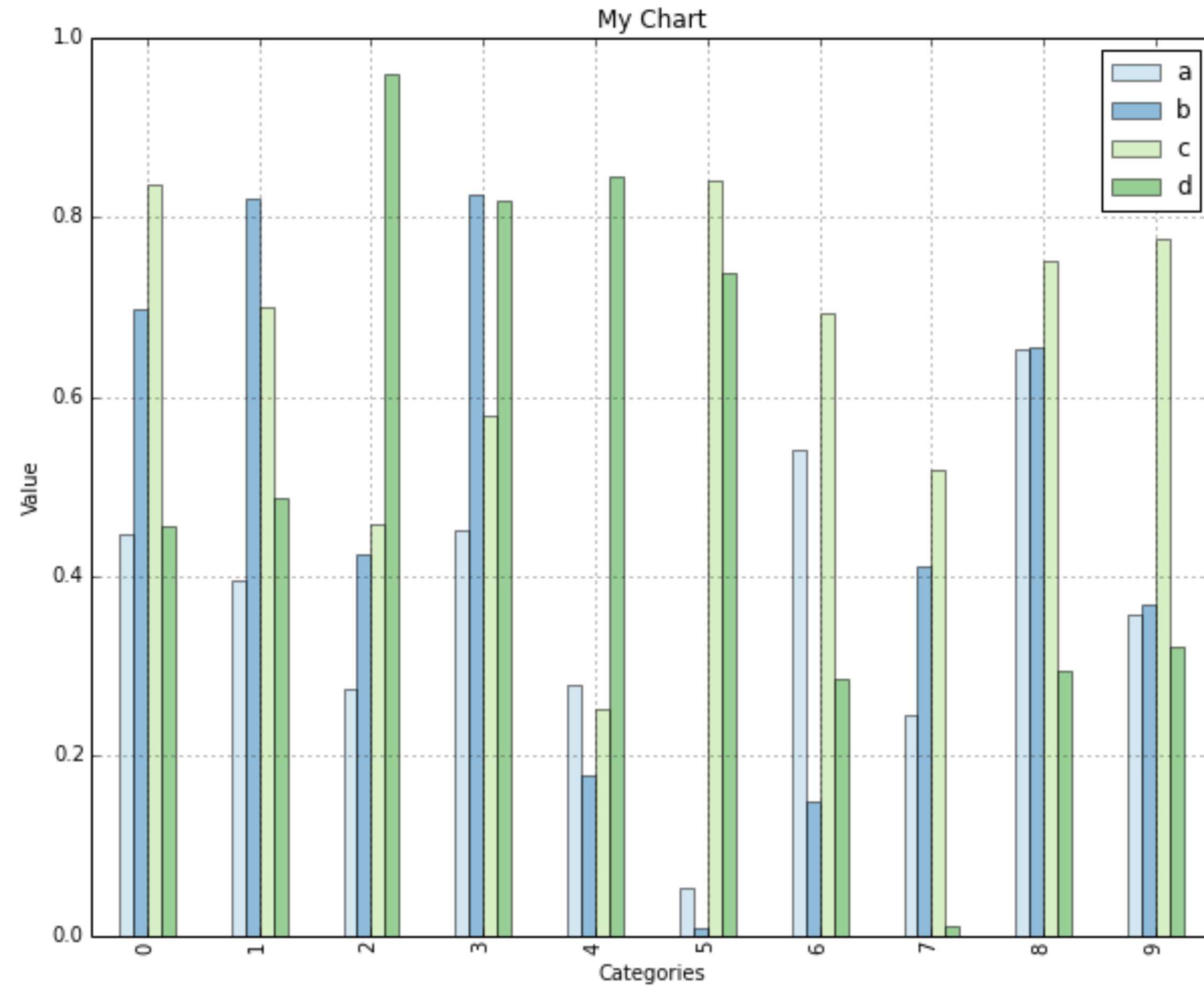
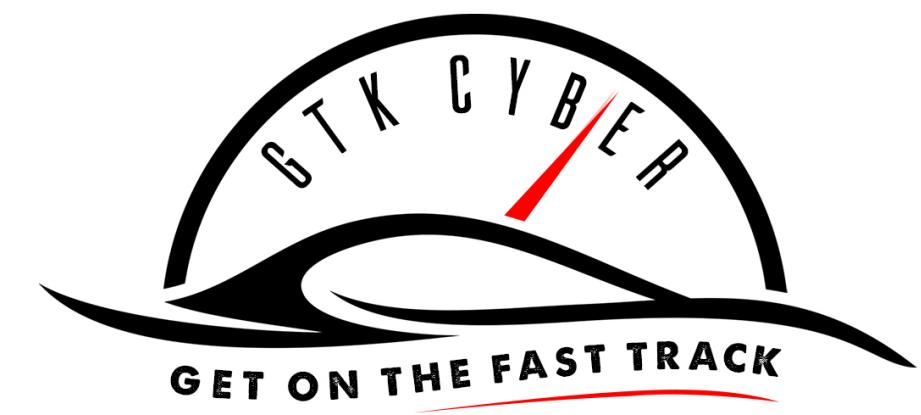
```
piechart = data.plot( kind="pie" )
```



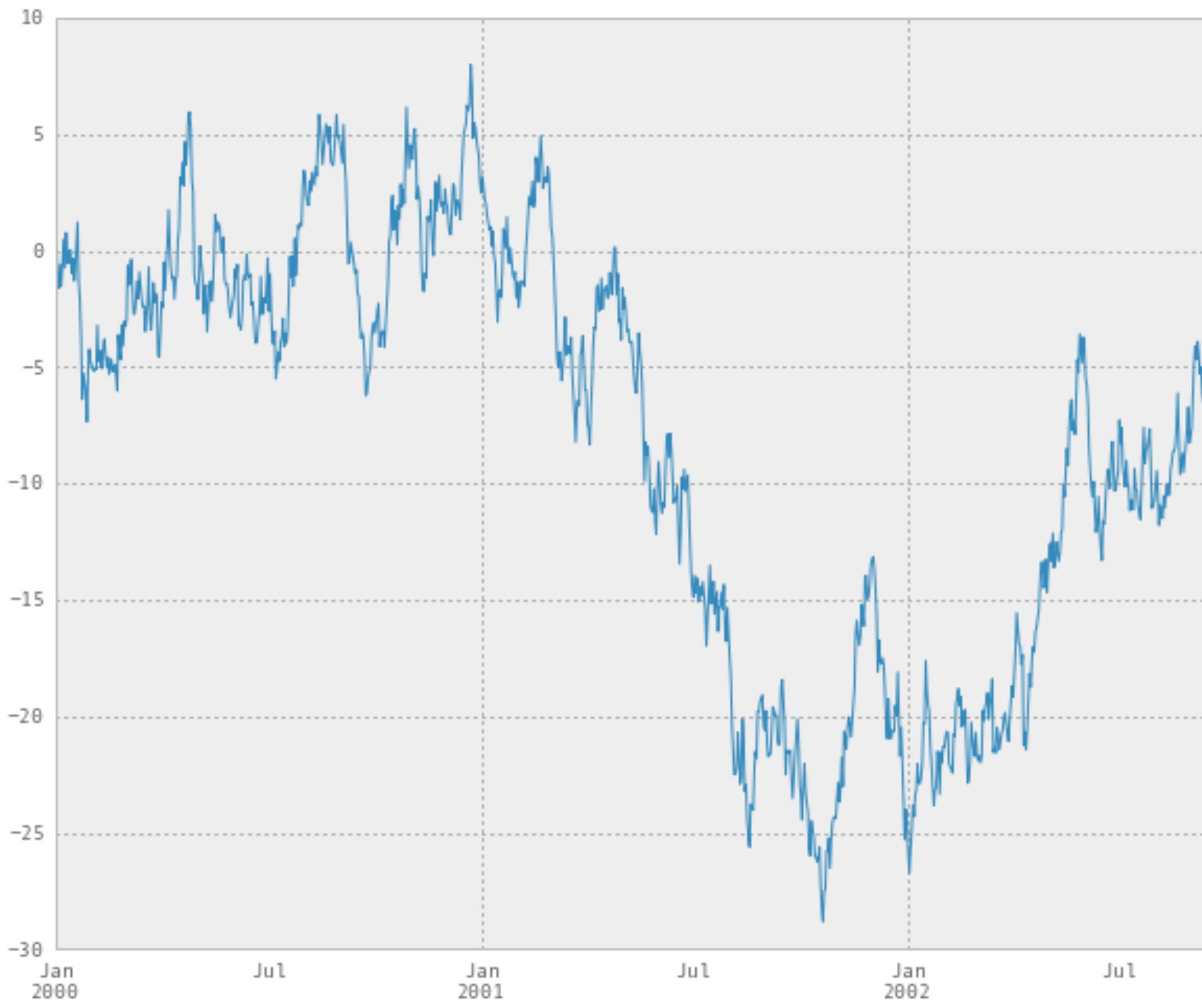
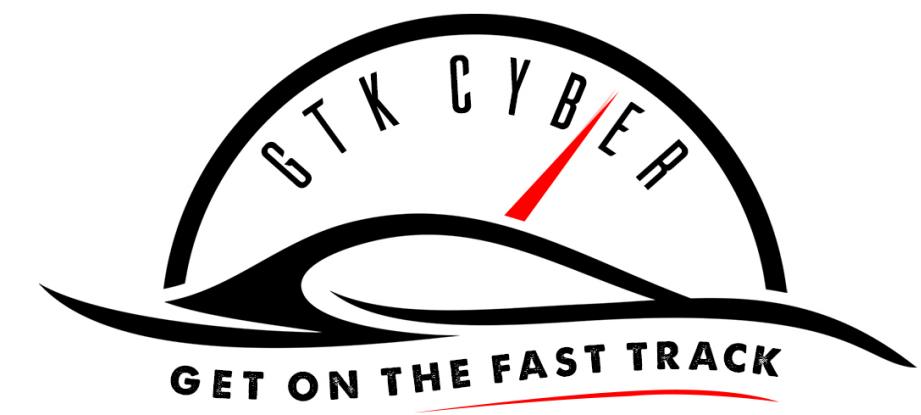
```
df2 = pd.DataFrame(np.random.rand(10, 4),  
columns=[ 'a' , 'b' , 'c' , 'd' ] )  
df2.plot( kind='bar' )
```



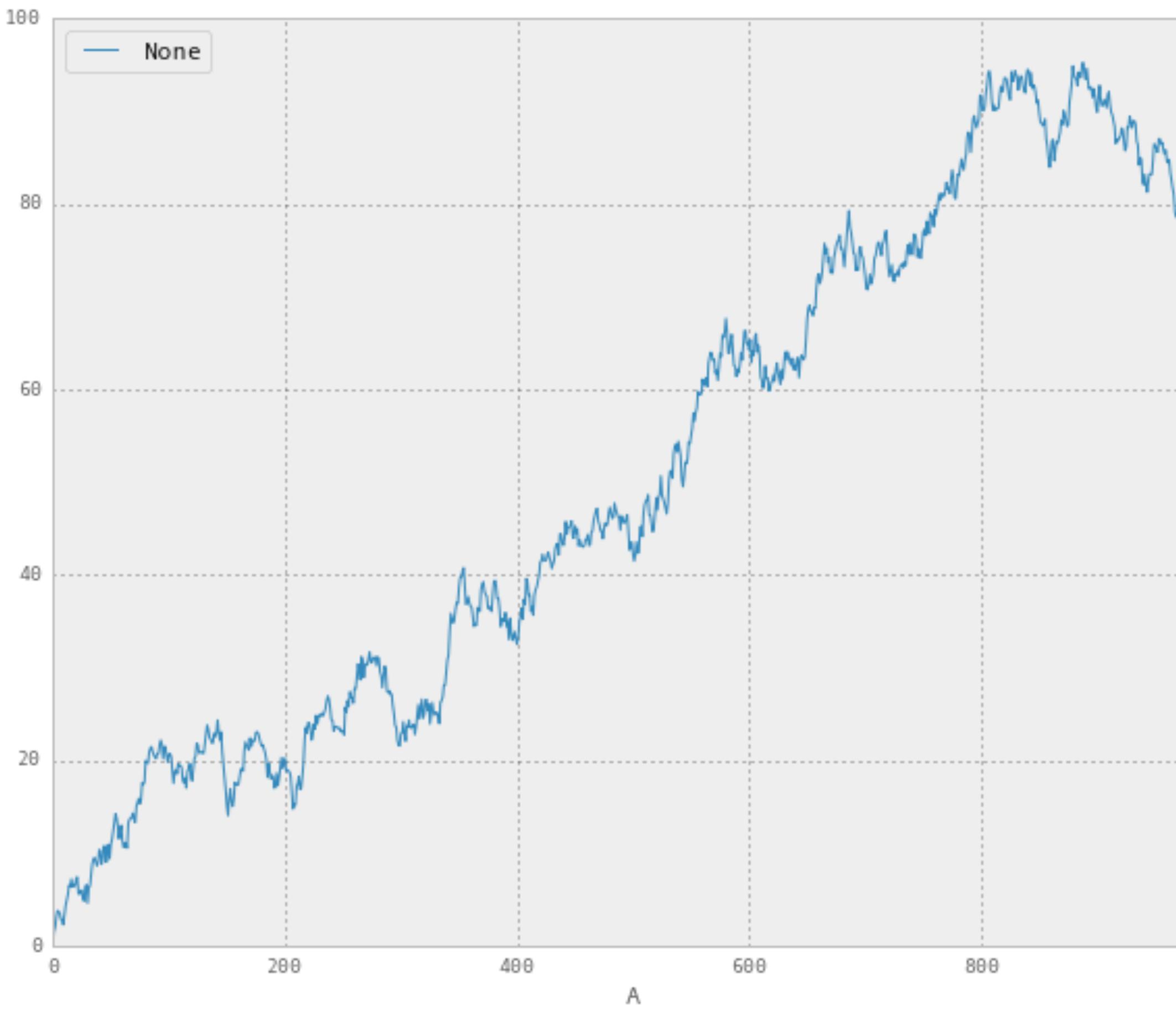
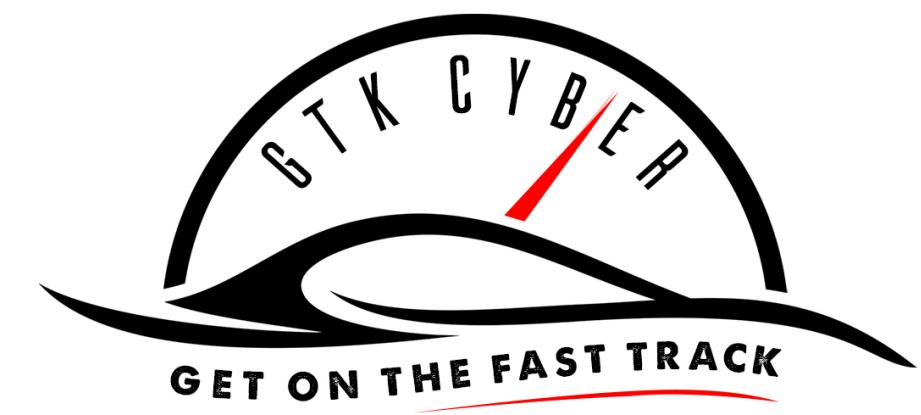
```
df2.plot( kind='bar', stacked=True )
```



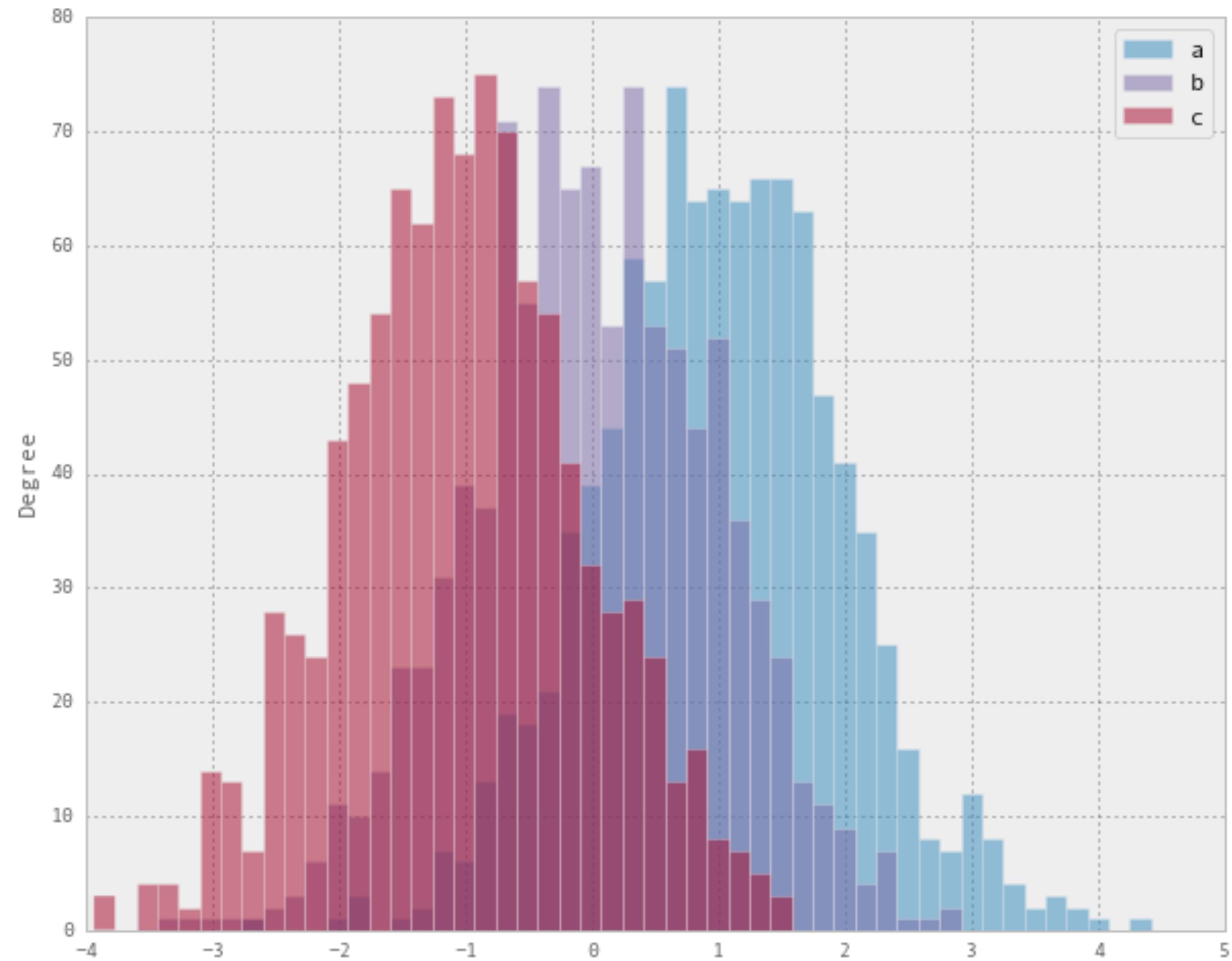
```
df2.plot( kind='bar',
          color=('#a6cee3', '#1f78b4', '#b2df8a', '#33a02c' ),
          alpha=0.5,
          width=0.5,
          figsize=(10, 8))
plt.title( "My Chart" )
plt.xlabel( "Categories" )
plt.ylabel( "Value" )
```



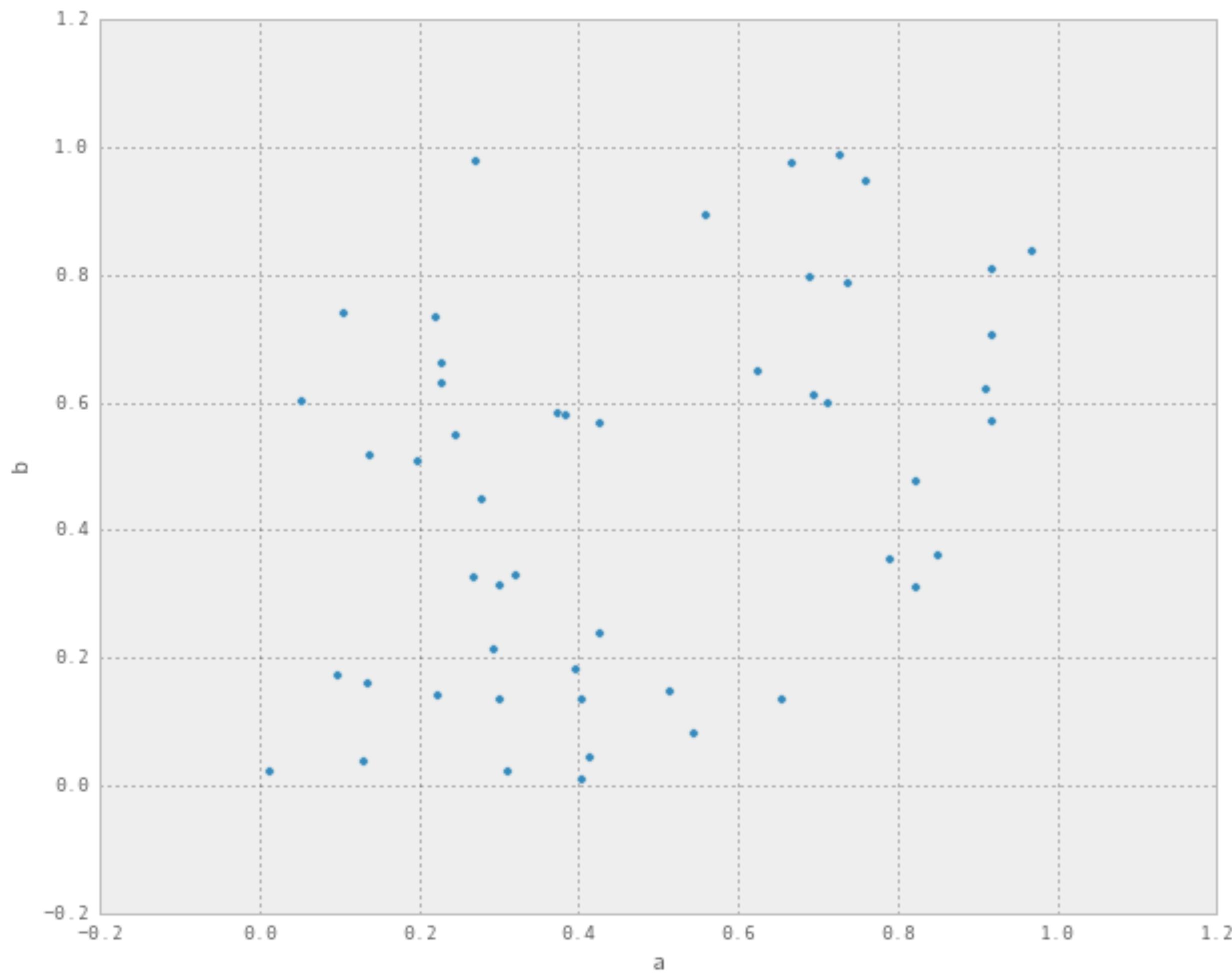
```
ts = pd.Series(np.random.randn( 1000 ),  
index=pd.date_range('1/1/2000', periods=1000))  
ts = ts.cumsum()  
timeseriesChart = ts.plot( figsize=(10, 8) )
```



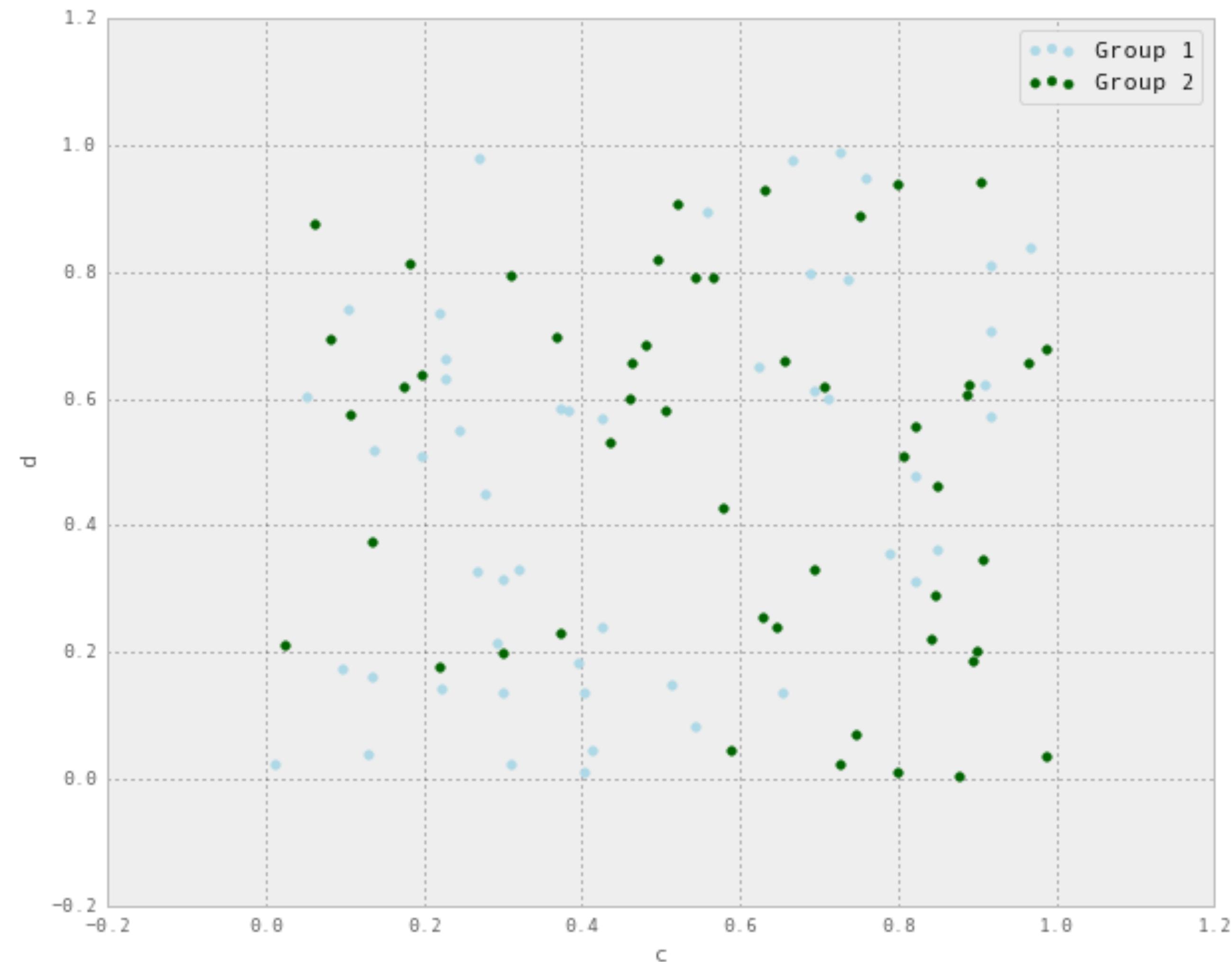
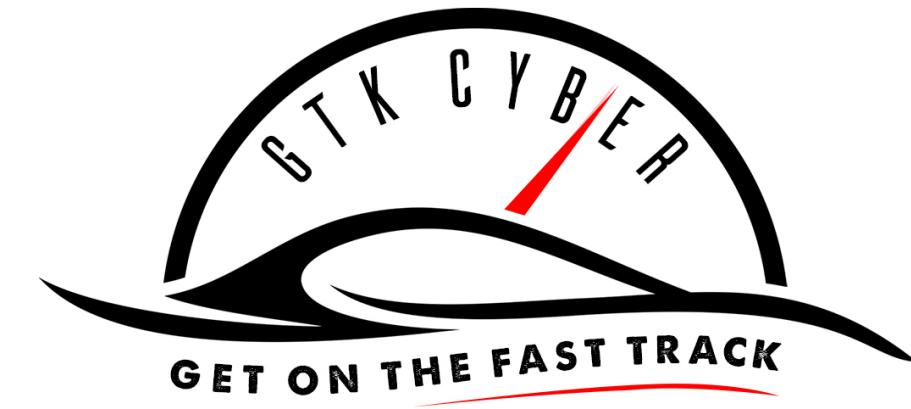
```
df3 = pd.DataFrame(np.random.randn(1000, 2),  
columns=[ 'B' , 'C' ]).cumsum()  
df3[ 'A' ] = pd.Series(list(range(len(df3))))  
  
df3.plot( x='A' , y='B' )
```



```
df4.plot(kind='hist',  
         alpha=0.5,  
         bins=50 )
```

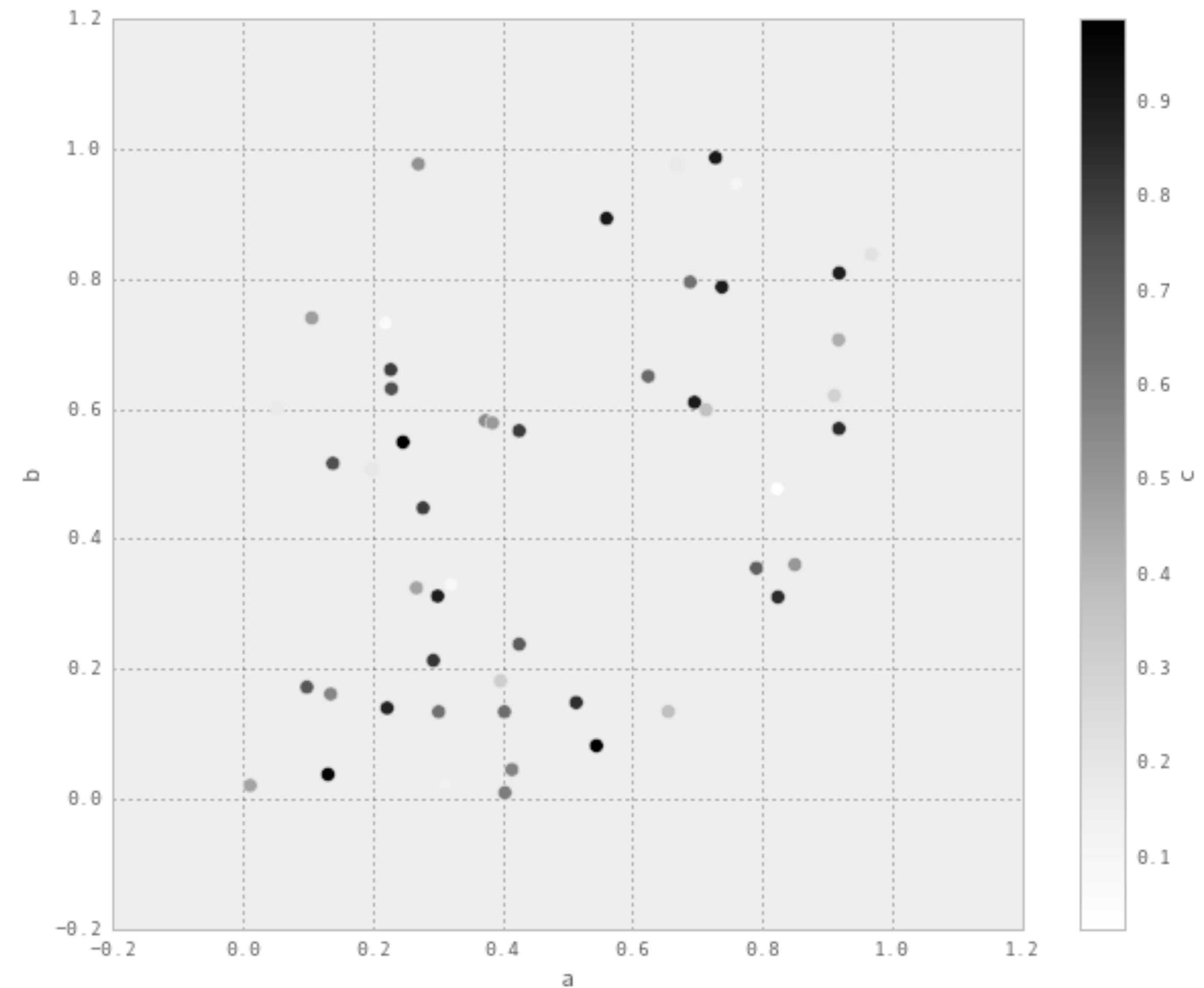
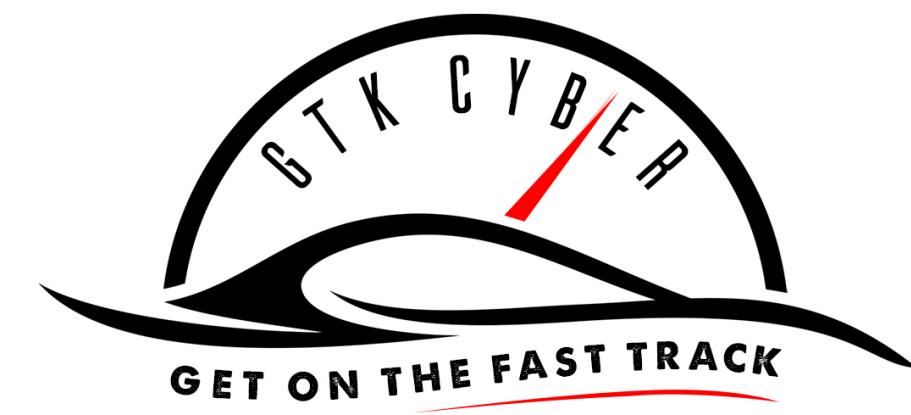


```
df5 = pd.DataFrame(np.random.rand(50, 4),  
columns=[ 'a' , 'b' , 'c' , 'd' ] )  
df5.plot(kind='scatter', x='a' , y='b' )
```

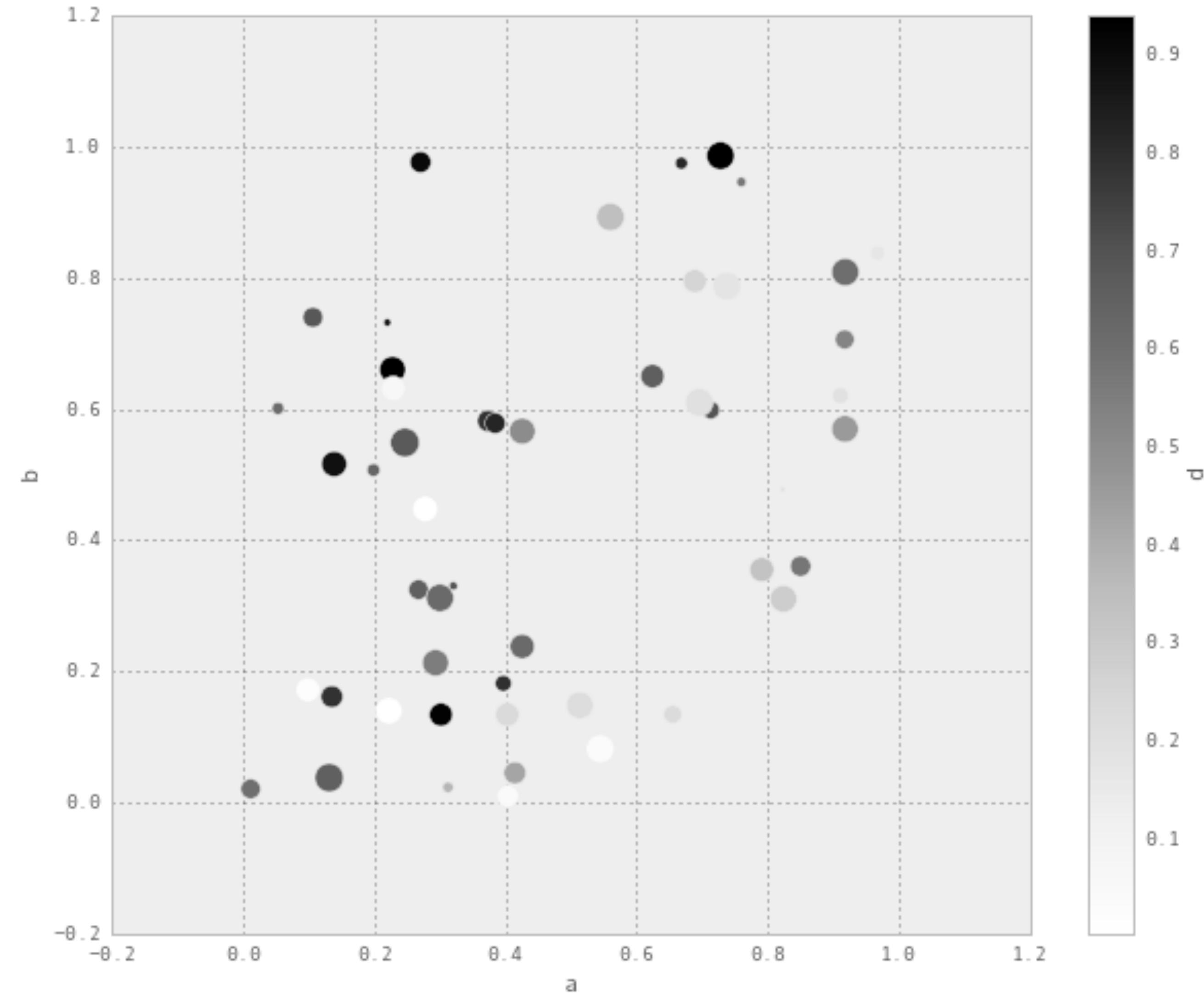
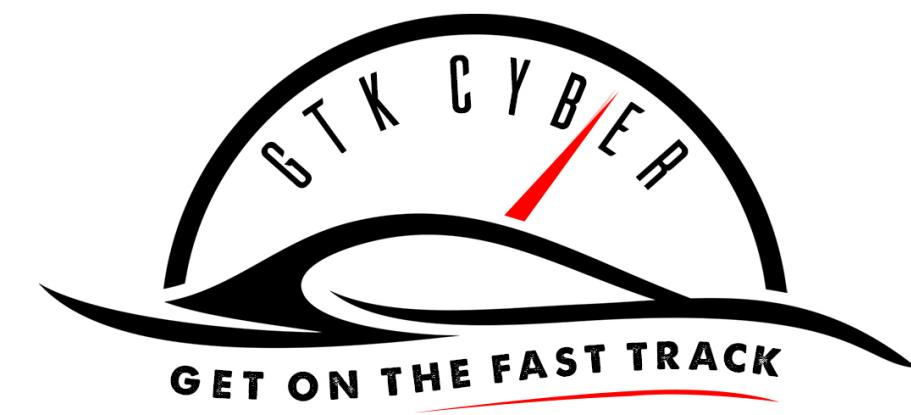


```
ax = df5.plot(kind='scatter', x='a', y='b',
               color='LightBlue',
               label='Group 1',
               figsize=(10, 8) )
```

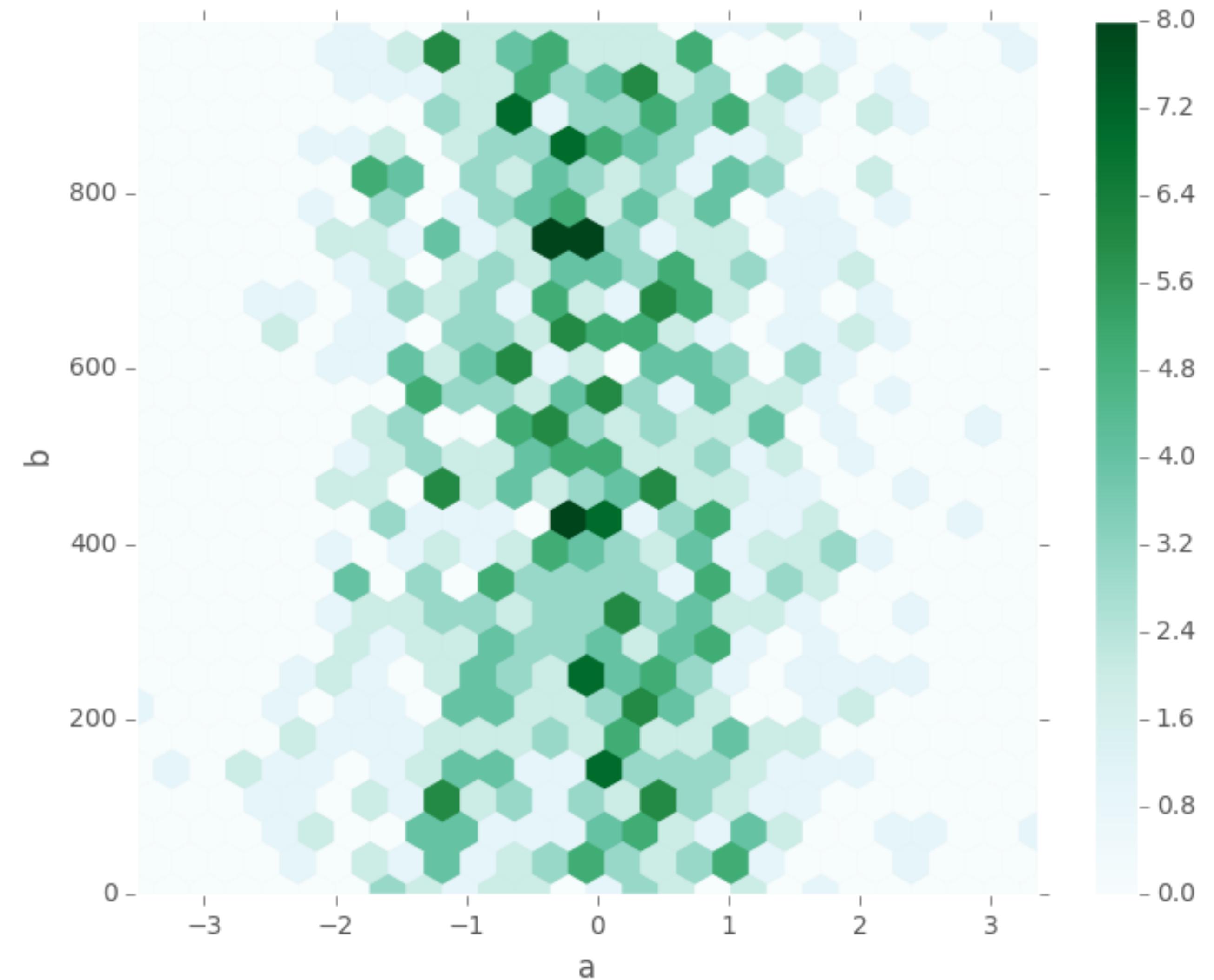
```
df5.plot(kind='scatter', x='c', y='d',
          color='DarkGreen',
          label='Group 2',
          ax=ax)
```



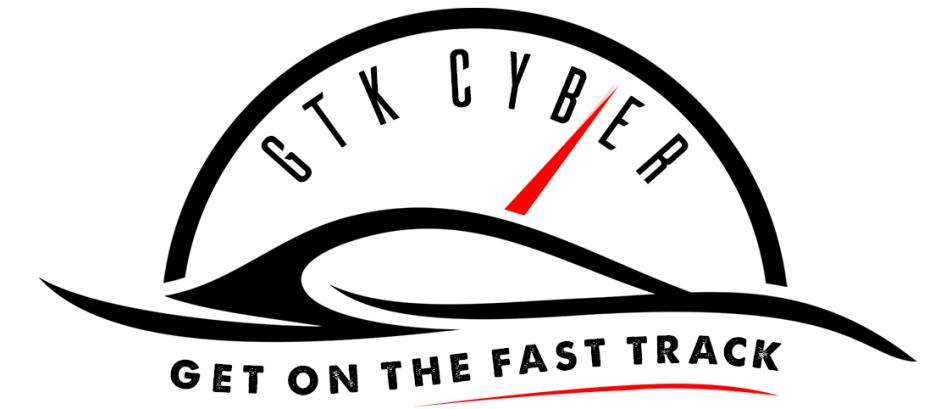
```
df5.plot(kind='scatter', x='a', y='b',  
c='c', s=50 )
```



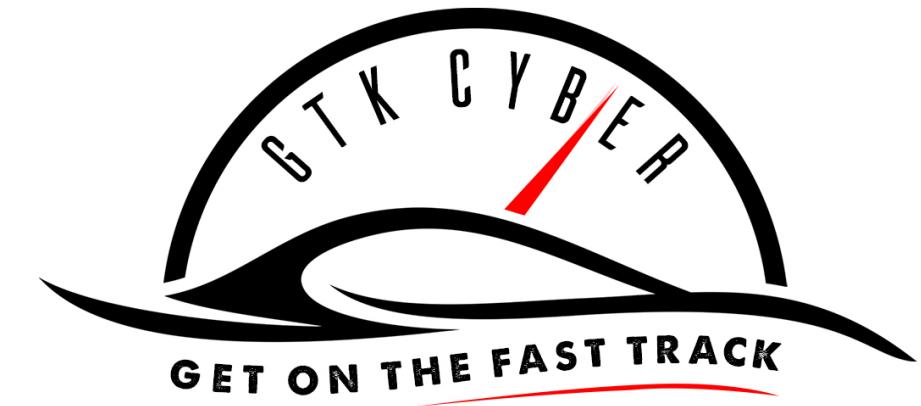
```
scatter = df5.plot(kind='scatter',
                    x='a',
                    y='b',
                    s=df5[ 'c' ]*200,
                    c='d')
```



```
df.plot(kind='hexbin', x='a', y='b',  
gridsize=25)
```

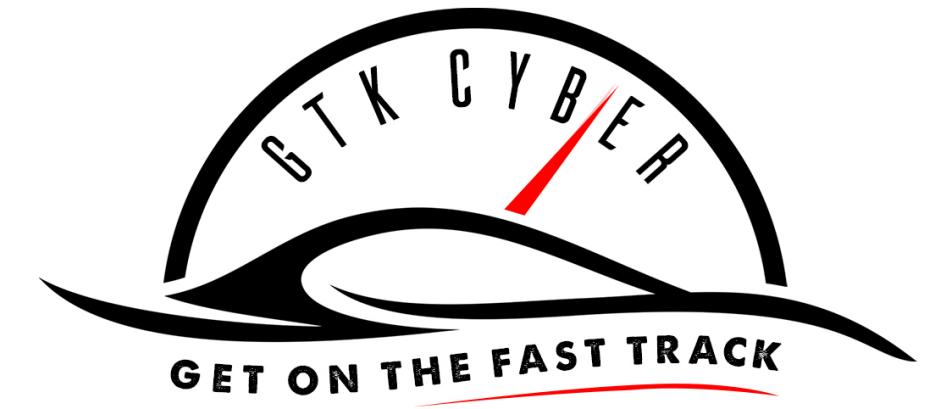


```
scatterPlot.get_figure().savefig( "scatterPlot.png" )
```

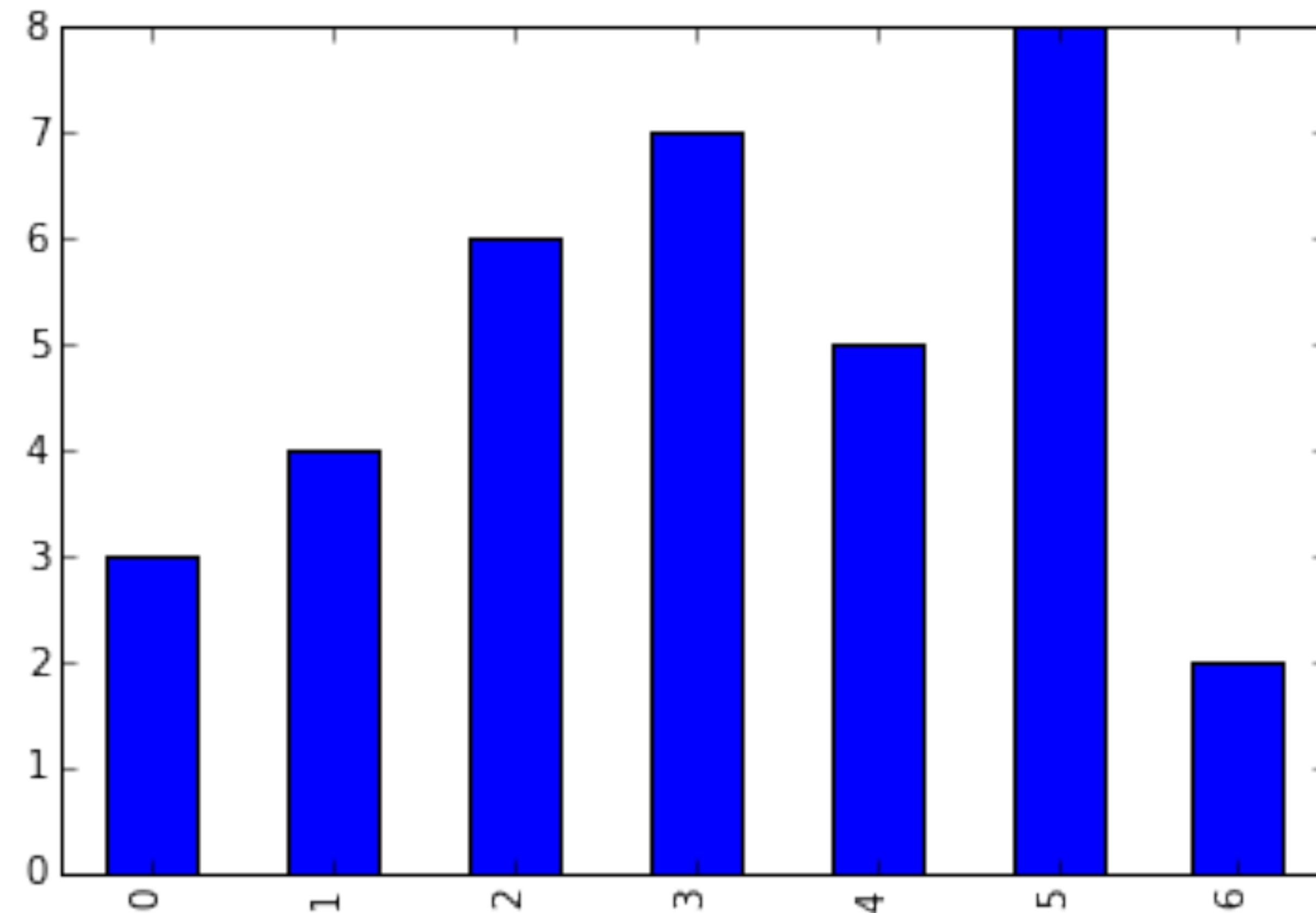


```
print(plt.style.available)

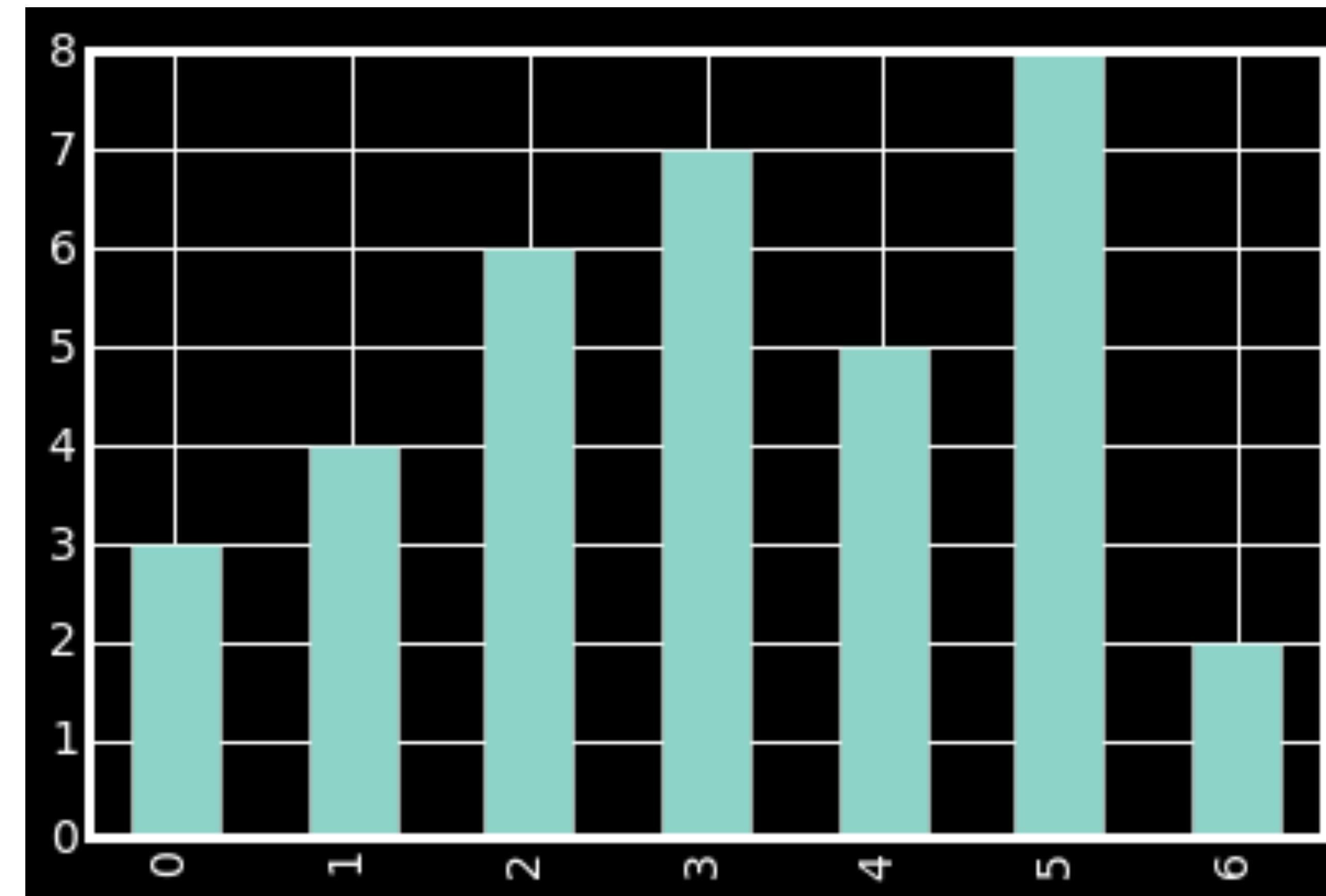
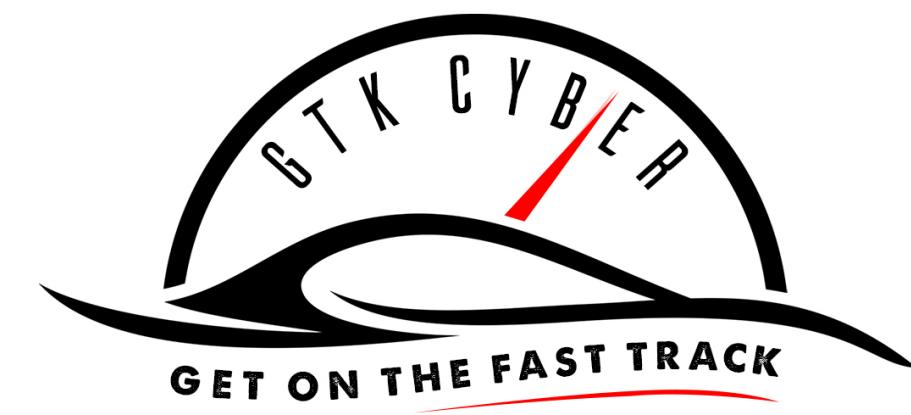
['dark_background', 'grayscale', 'ggplot',
 'bmh', 'fivethirtyeight']
```



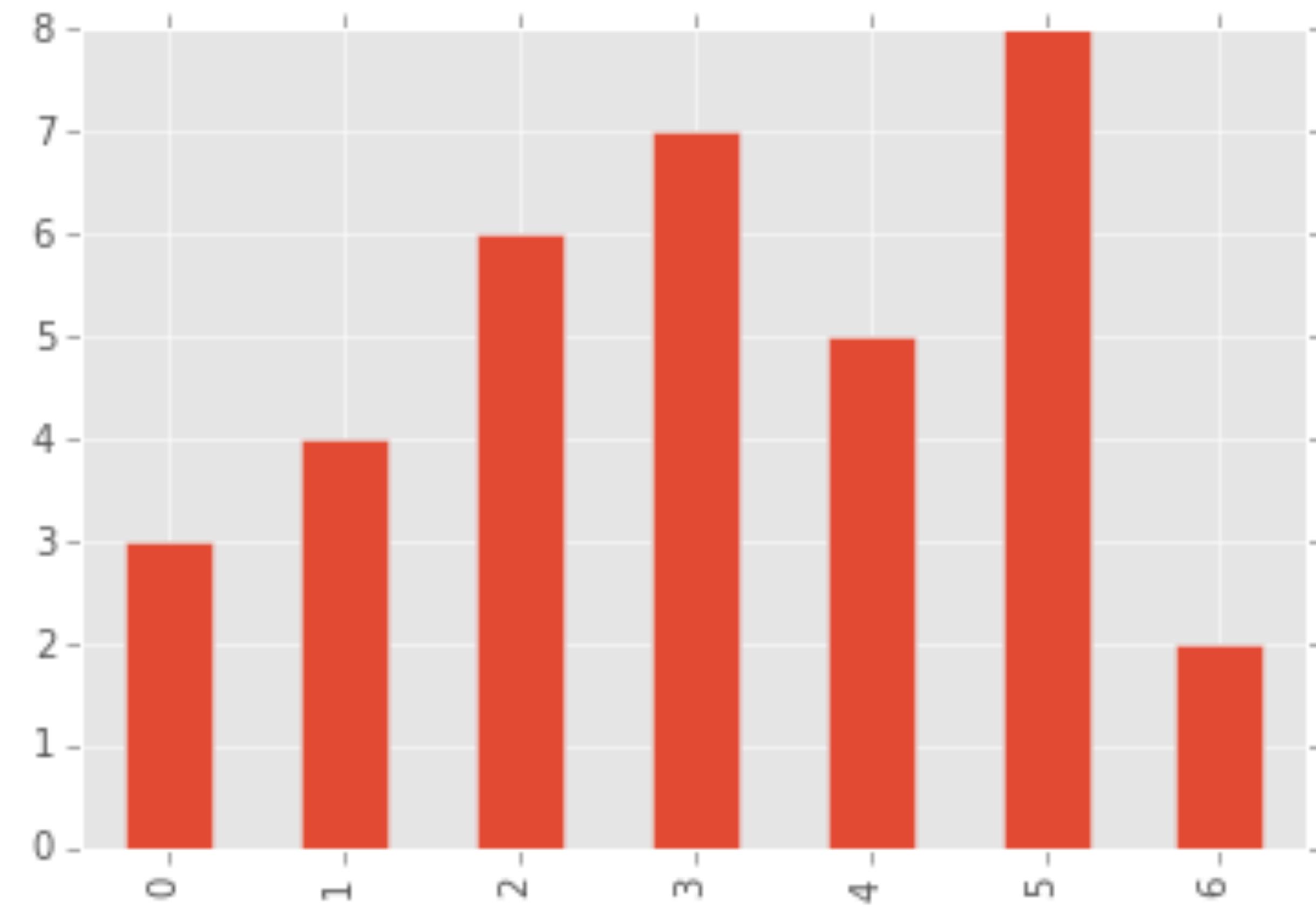
UGLY



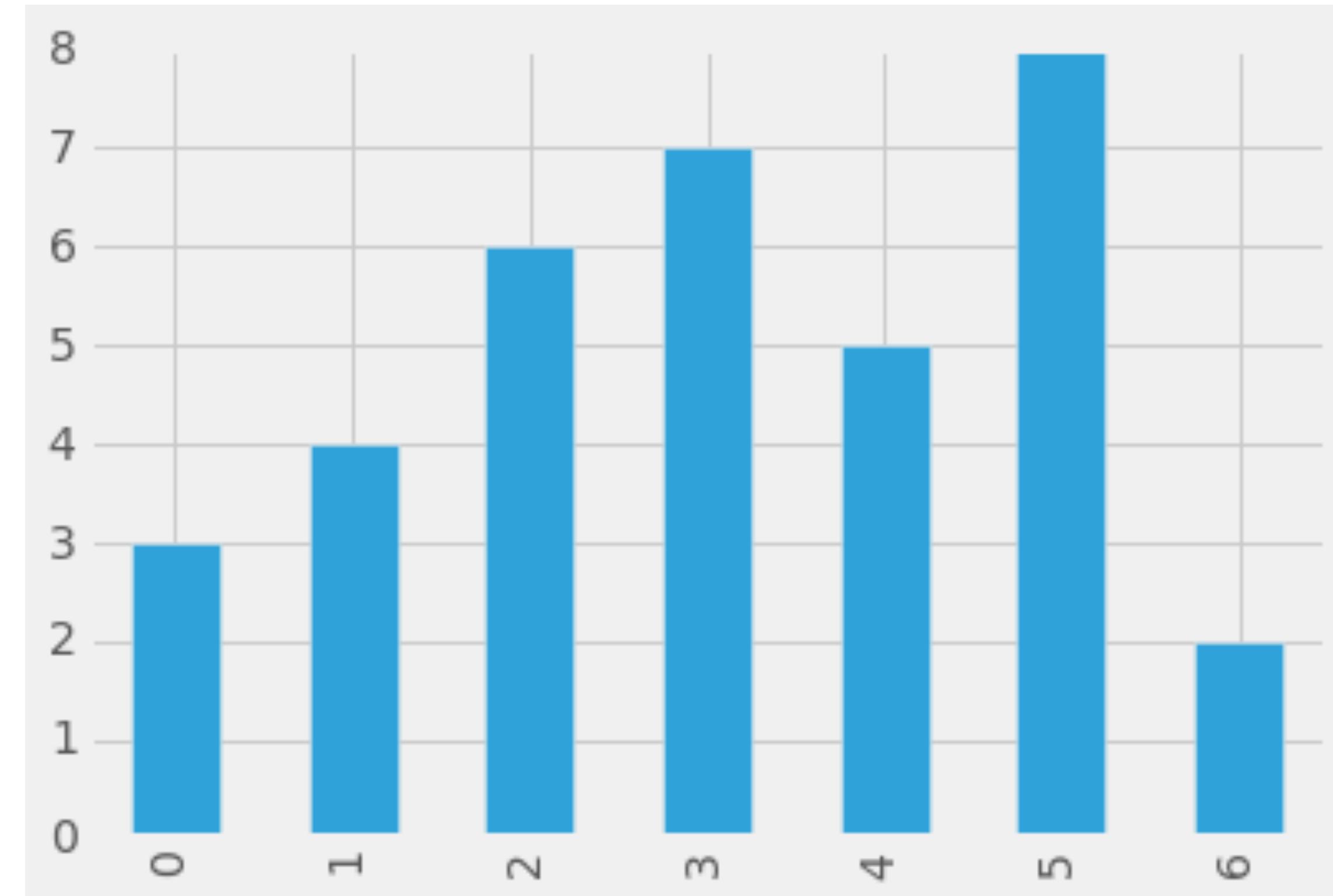
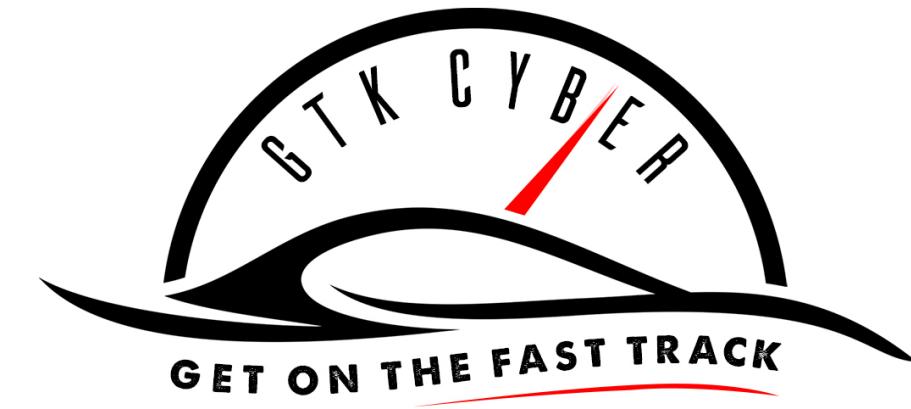
```
data = pd.Series( [3,4,6,7,5,8,2] )
barchart = data.plot( kind="bar" )
```



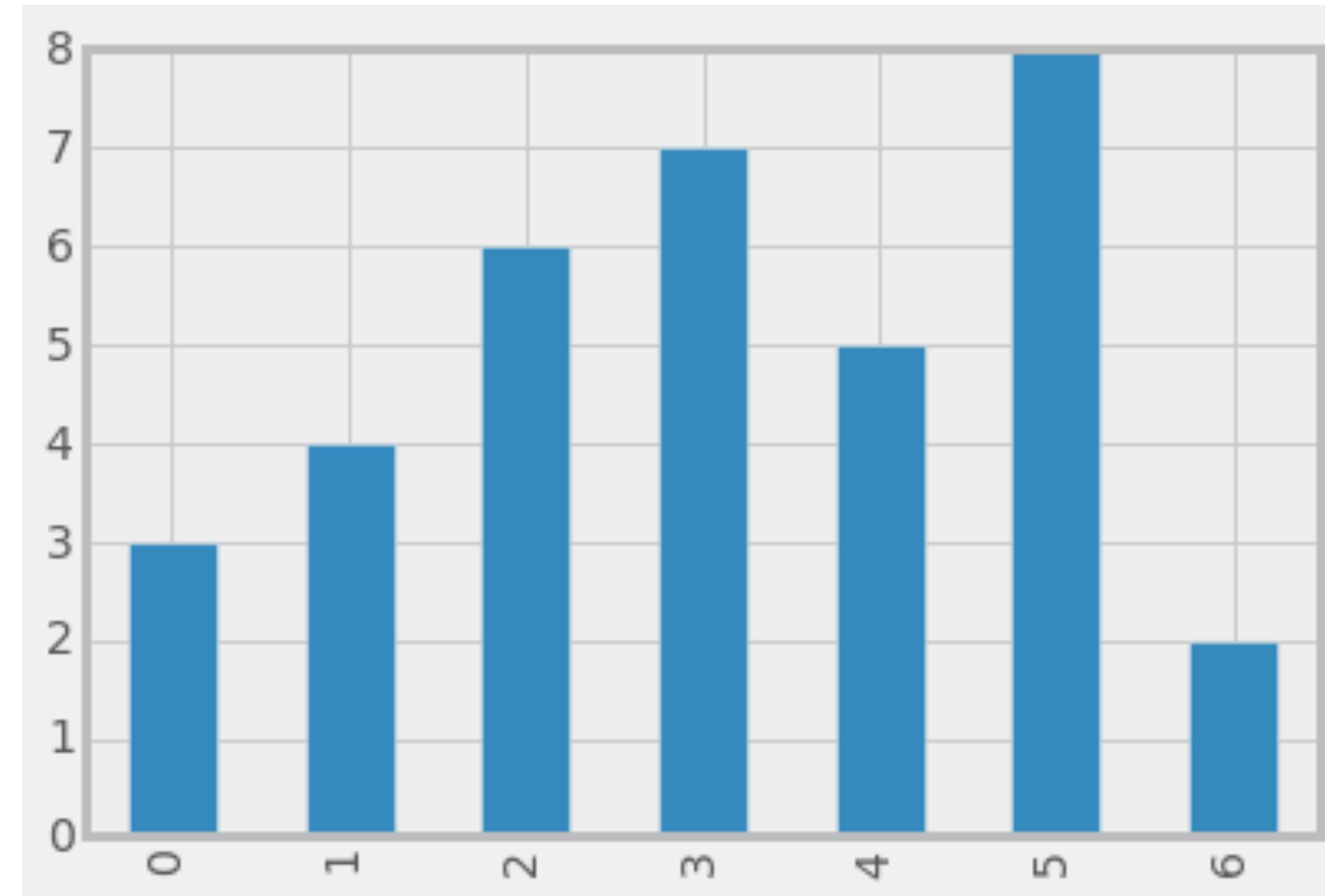
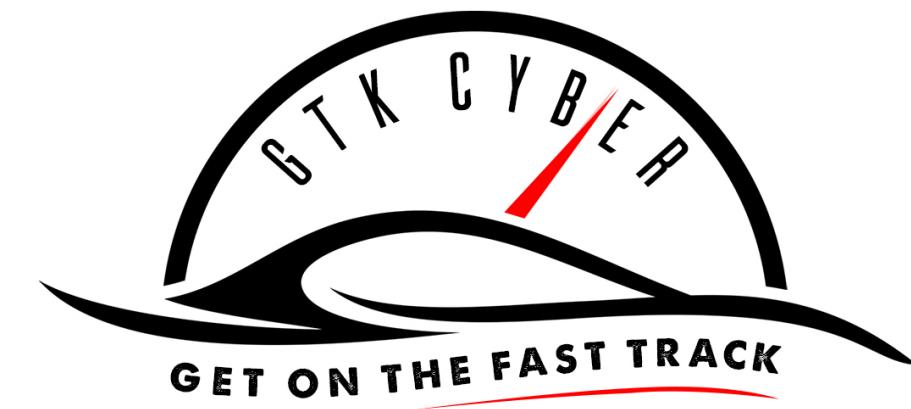
```
plt.style.use('dark_background')
barchart = data.plot( kind="bar" )
```



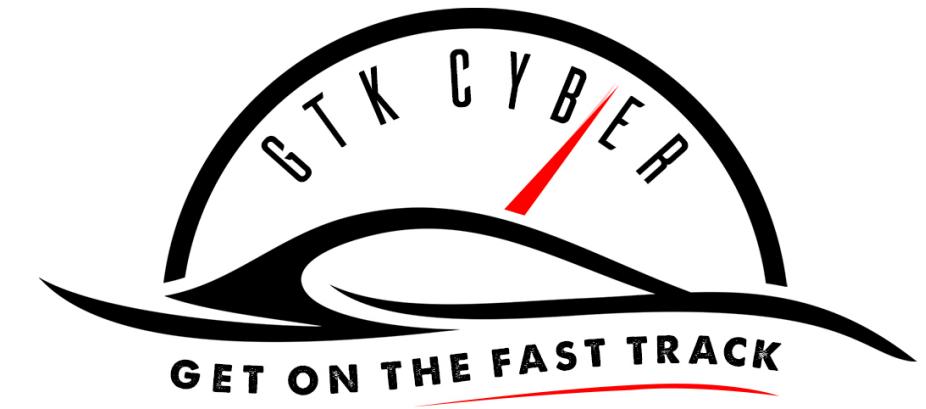
```
plt.style.use('ggplot')
barchart = data.plot( kind="bar" )
```



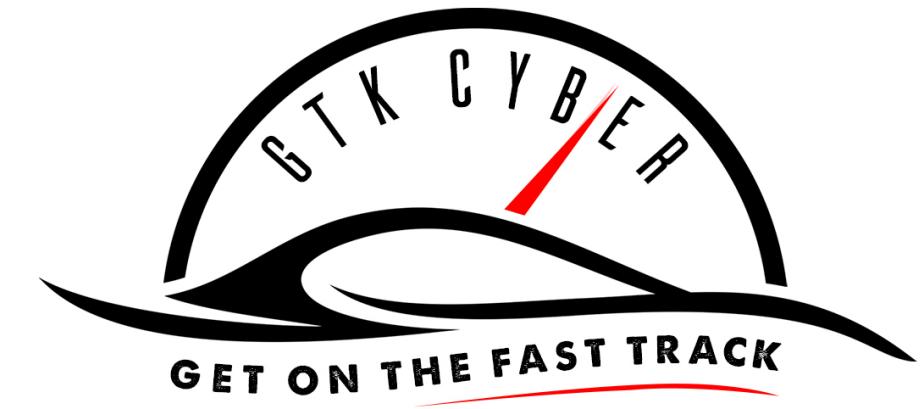
```
plt.style.use('fivethirtyeight')
barchart = data.plot( kind="bar" )
```



```
plt.style.use('bmh')
barchart = data.plot( kind="bar" )
```



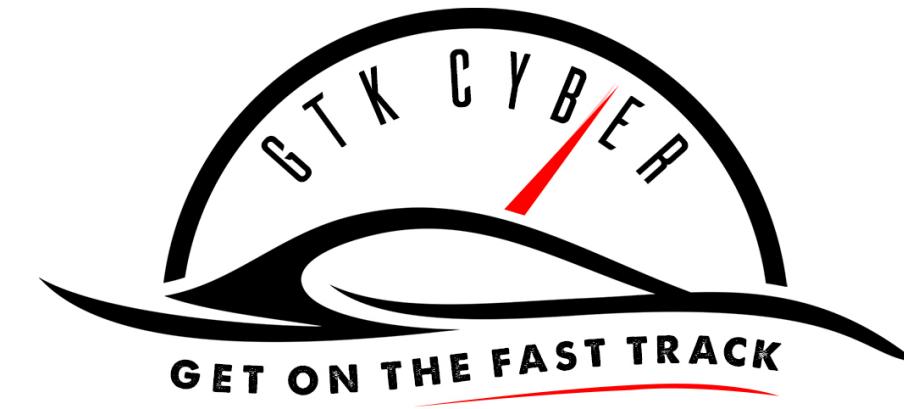
Interactive Visualizations



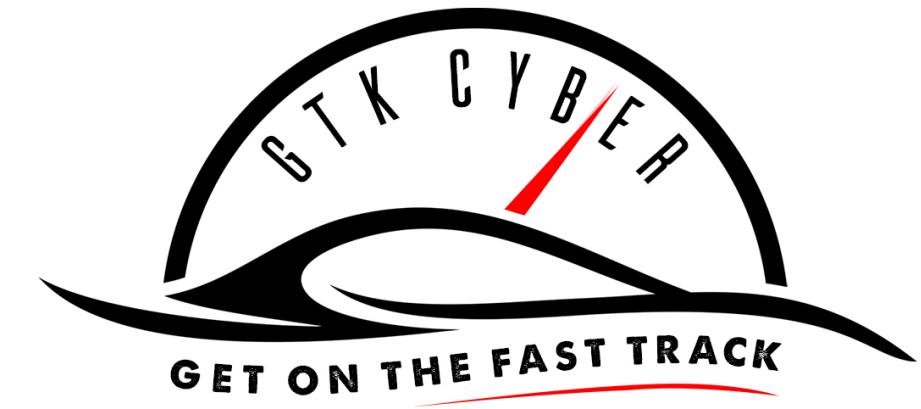
QlikView







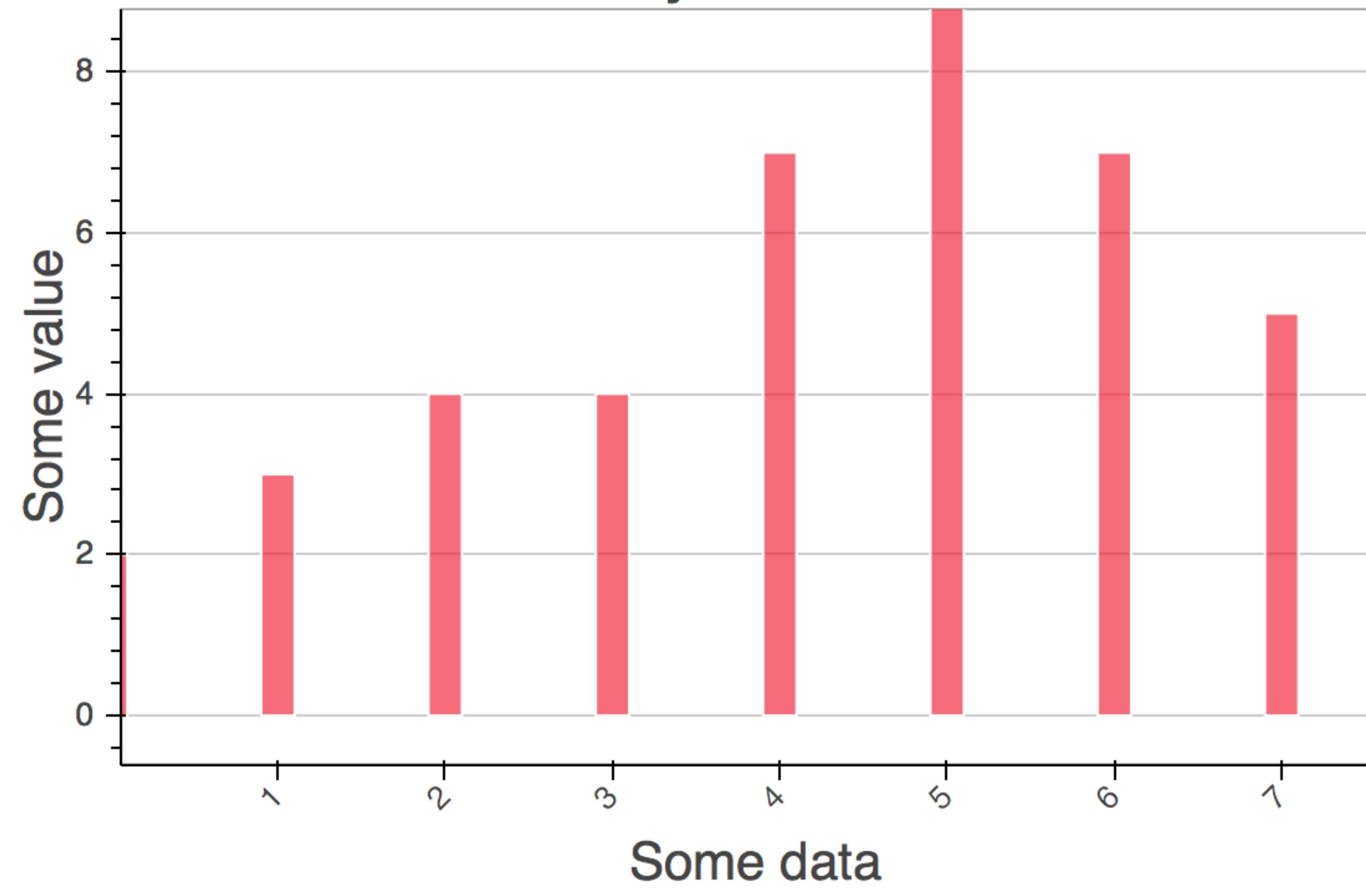
Easy to use... if you know R



Introducing Bokeh

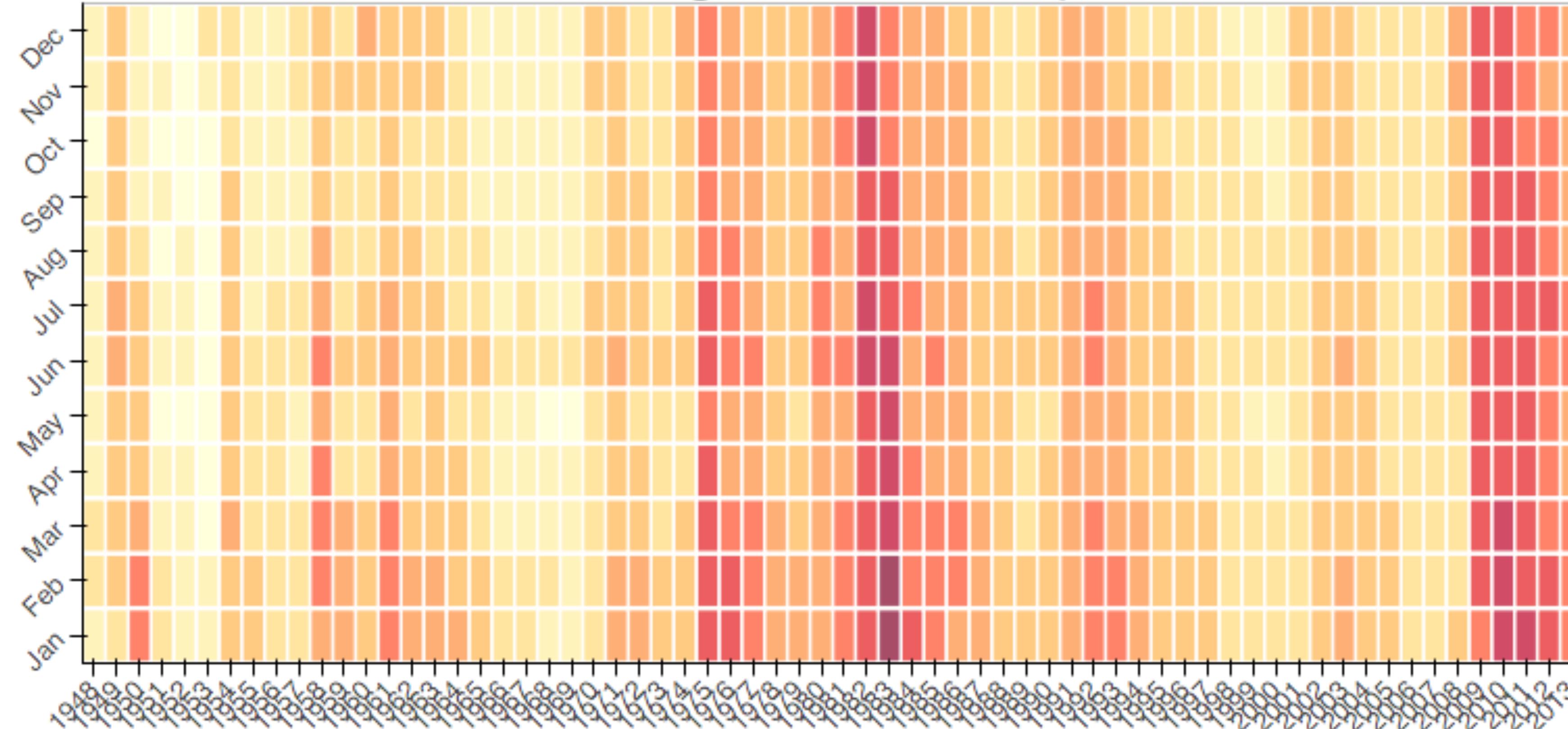


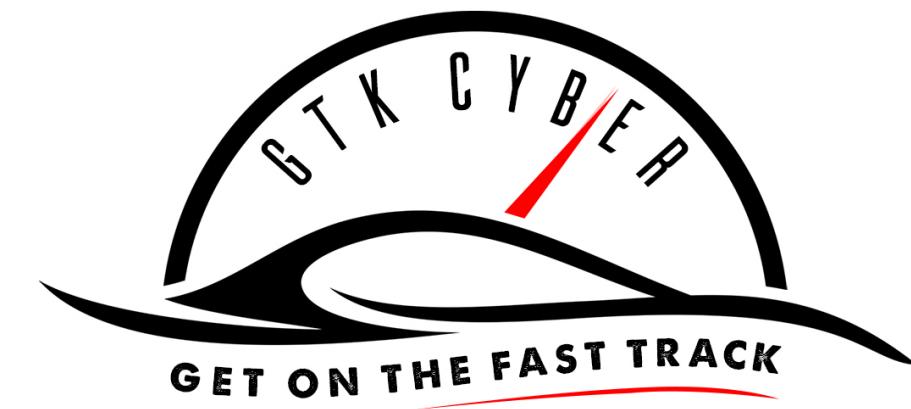
My Chart





categorical heatmap





```
from bokeh.charts import HeatMap, output_file, show
from bokeh.palettes import YlOrRd9 as palette
from bokeh.sampledata.unemployment1948 import data

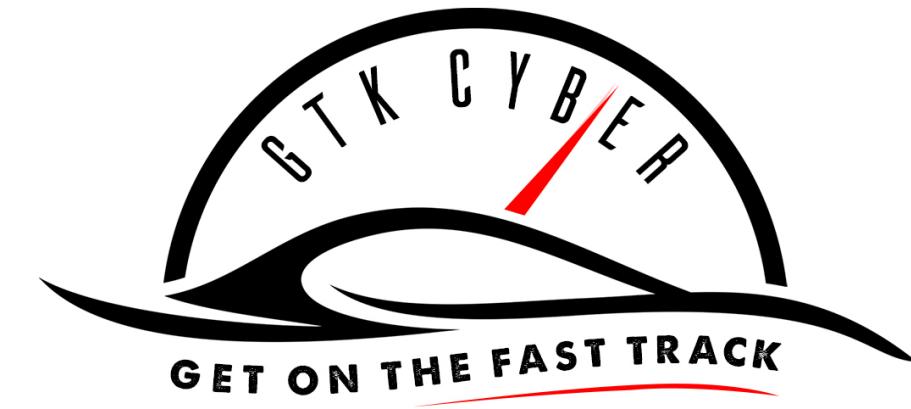
# pandas magic
df = data[data.columns[:-1]]
df2 = df.set_index(df[df.columns[0]].astype(str))
df2.drop(df.columns[0], axis=1, inplace=True)
df3 = df2.transpose()

output_file("cat_heatmap.html")

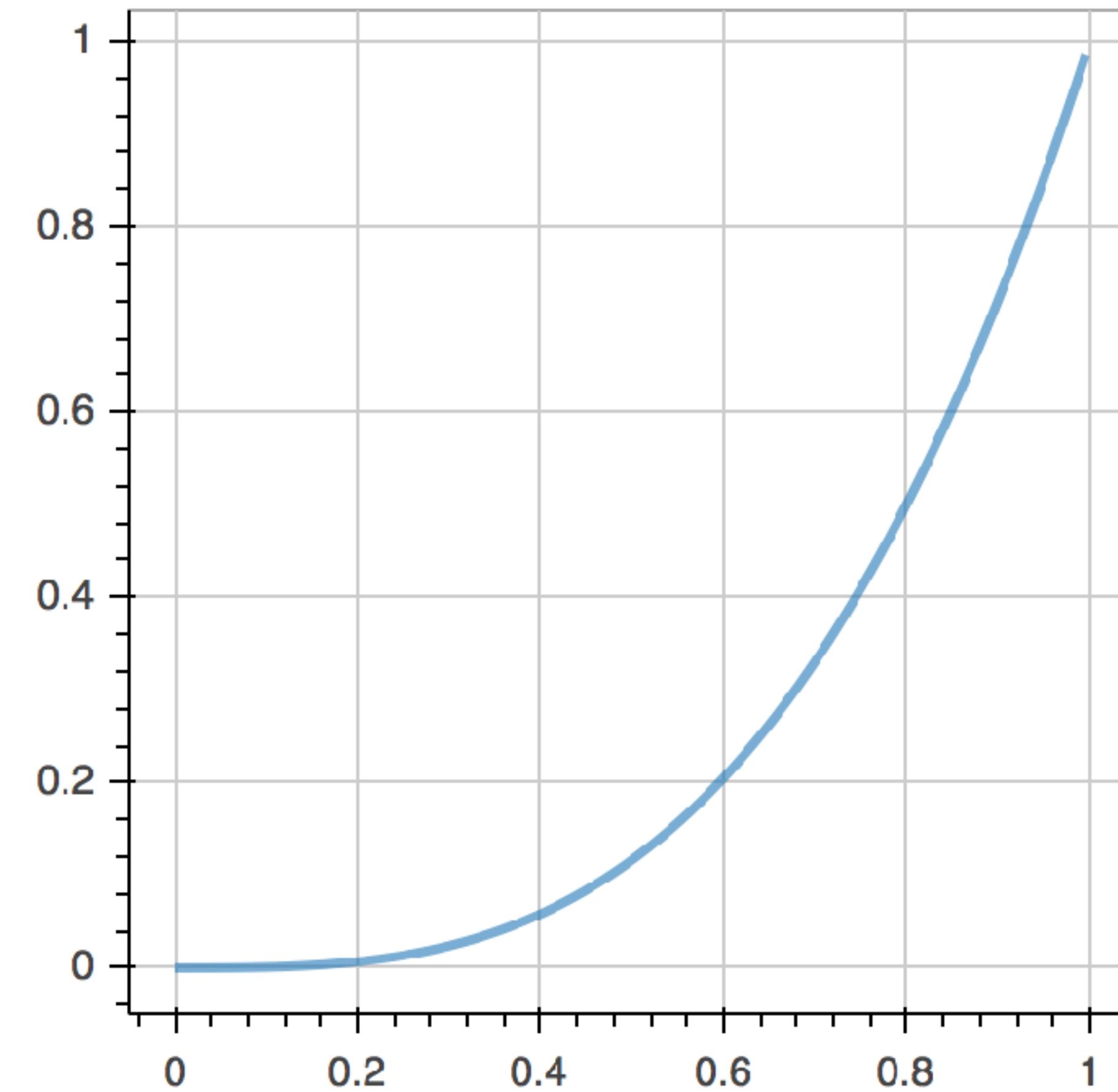
palette = palette[::-1] # Reverse the color order so dark red is highest un
hm = HeatMap(df3, title="categorical heatmap", width=800, palette=palette)

show(hm)
```

http://bokeh.pydata.org/en/latest/docs/gallery/cat_heatmap_chart.html

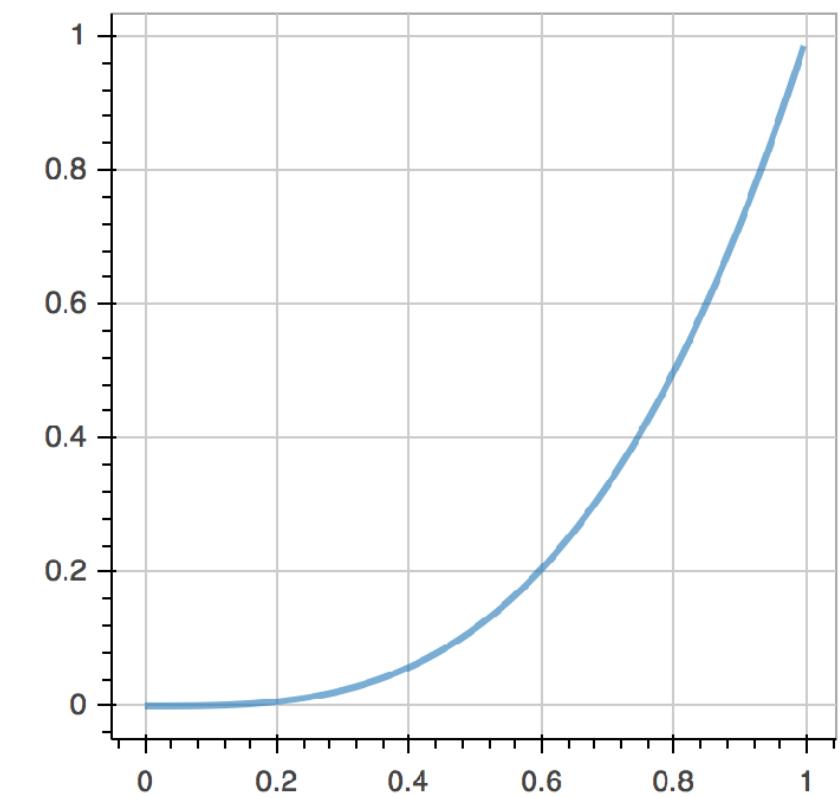


power: 3.1





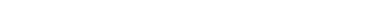
power: 3

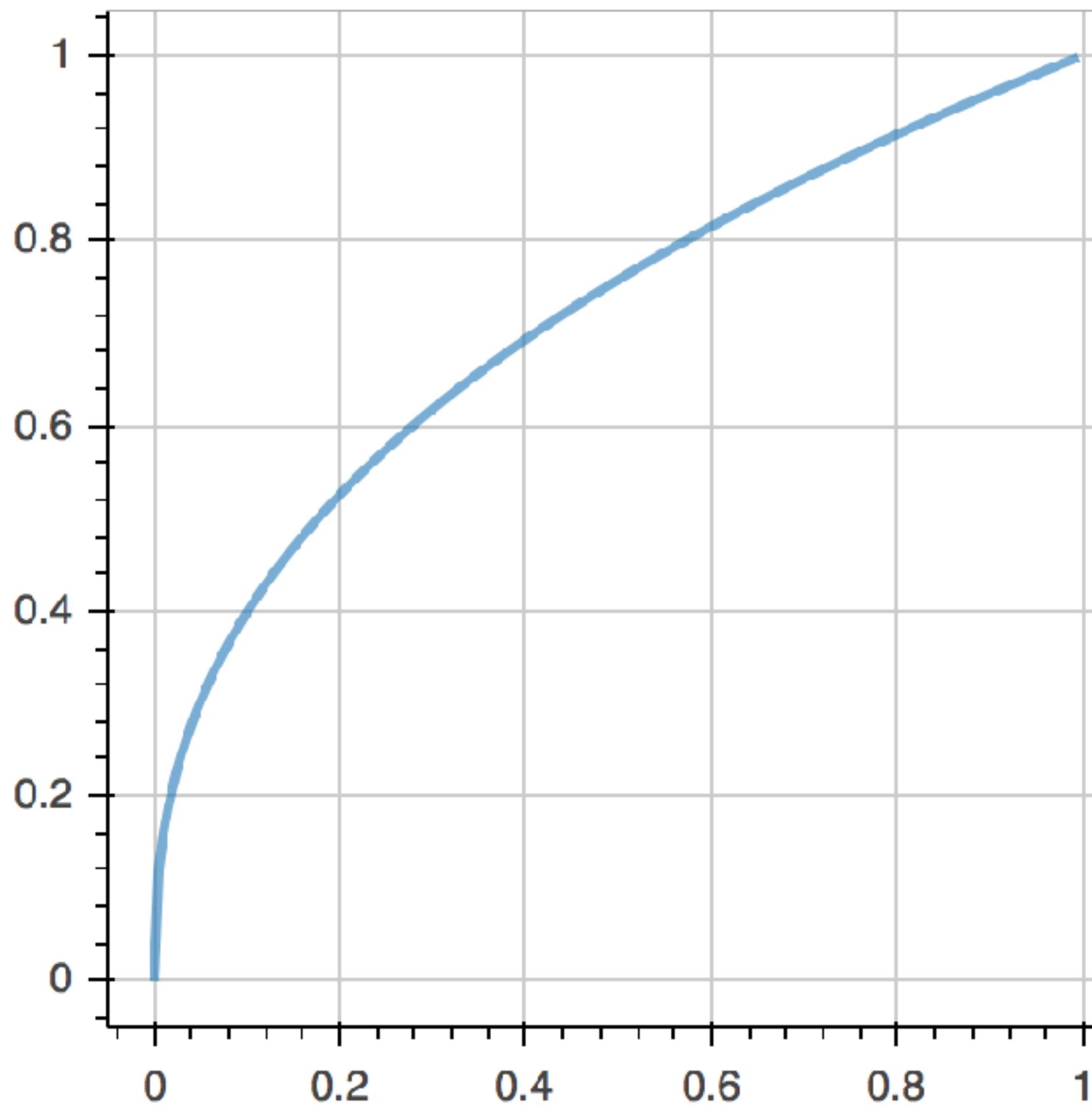


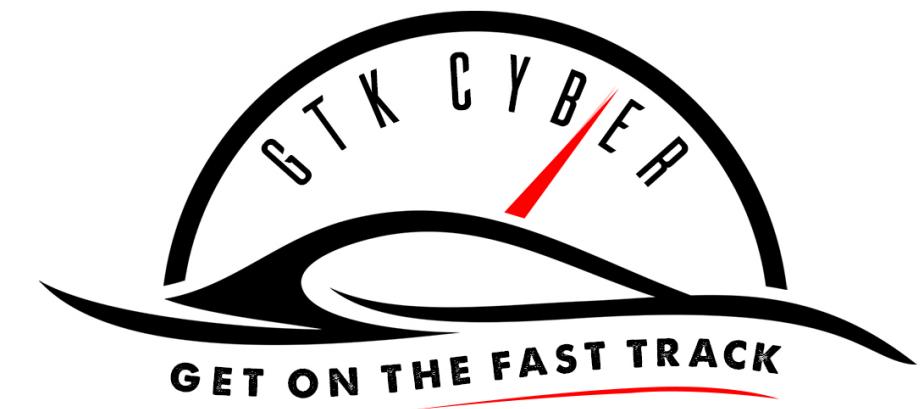
power: 0.4

1









```
from bokeh.io import vform
from bokeh.models import Callback, ColumnDataSource, Slider
from bokeh.plotting import figure, output_file, show

x = [x*0.005 for x in range(0, 200)]
y = x

source = ColumnDataSource(data=dict(x=x, y=y))

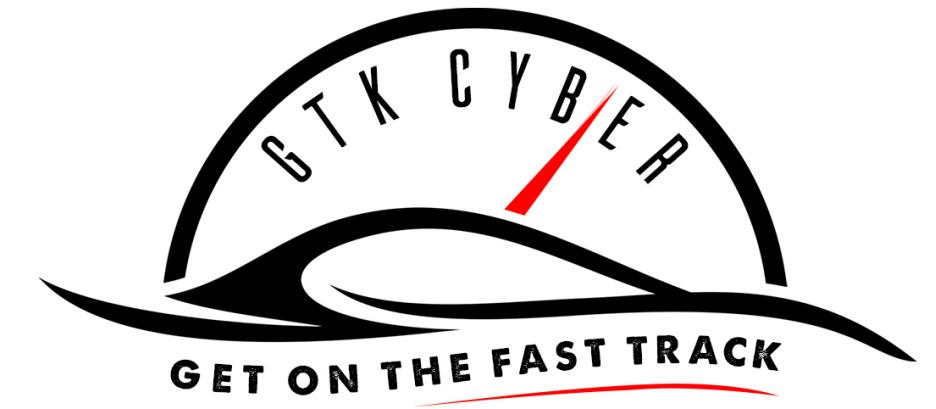
plot = figure(plot_width=400, plot_height=400)
plot.line('x', 'y', source=source, line_width=3, line_alpha=0.6)

callback = Callback(args=dict(source=source), code="""
    var data = source.get('data');
    var f = cb_obj.get('value')
    x = data['x']
    y = data['y']
    for (i = 0; i < x.length; i++) {
        y[i] = Math.pow(x[i], f)
    }
    source.trigger('change');
""")

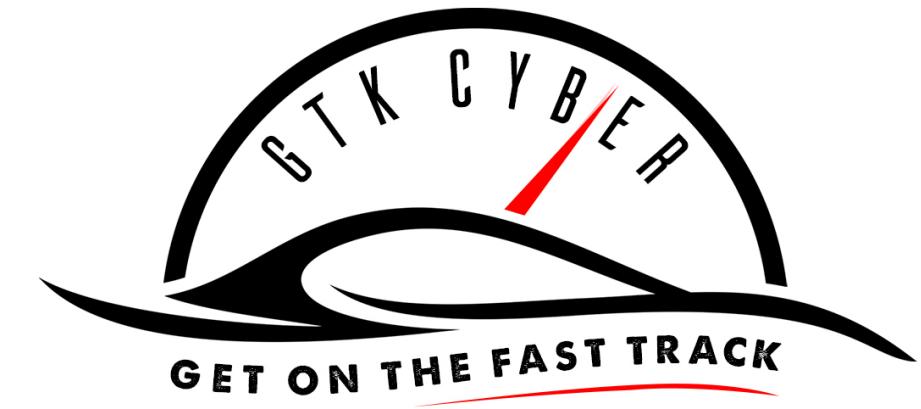
slider = Slider(start=0.1, end=4, value=1, step=.1, title="power", callback=callback)

layout = vform(slider, plot)

show(layout)
```



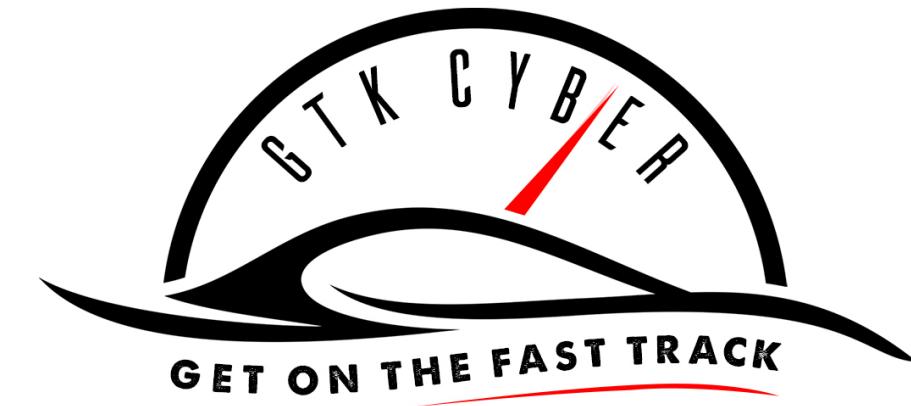
Using Bokeh



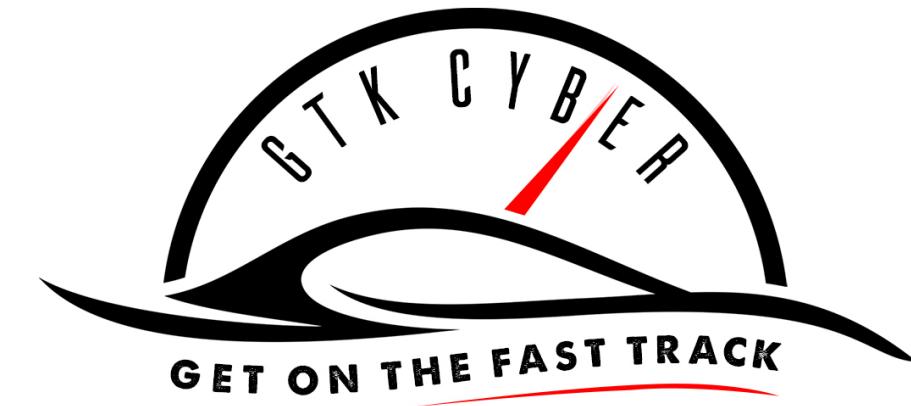
```
from bokeh.plotting import output_notebook  
output_notebook()
```



- Area (Overlapped & Stacked)
- Bar (Grouped & Stacked)
- BoxPlot
- Donut (ICK!)
- HeatMap
- Line
- Scatter
- Step
- Timeseries



```
from bokeh.charts import <chartname>
from bokeh.io import show
```



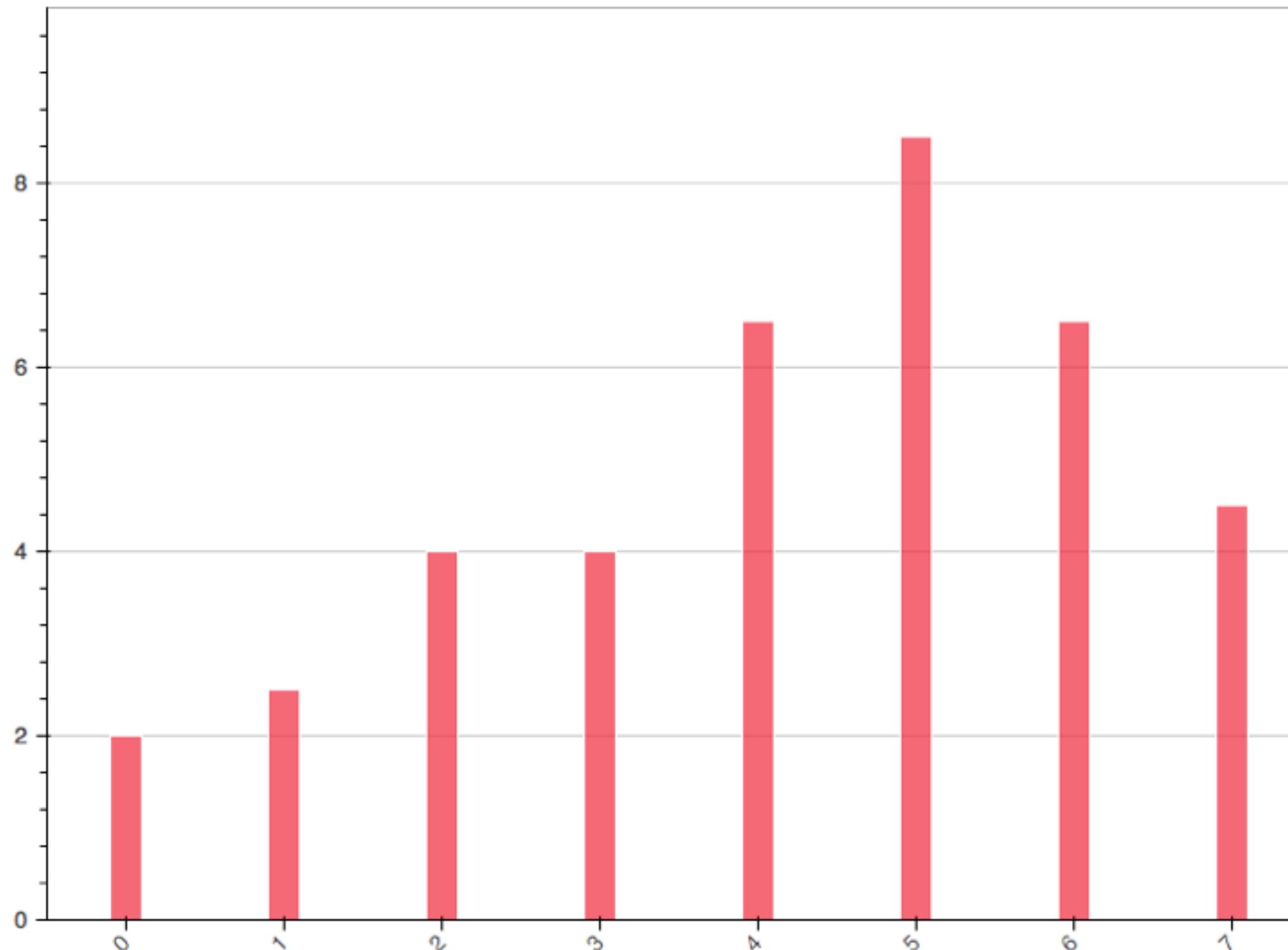
```
from bokeh.charts import Bar, show
data = [2,3,4,4,7,9,7,5]

barchart = Bar( data, notebook=True )
show()
```



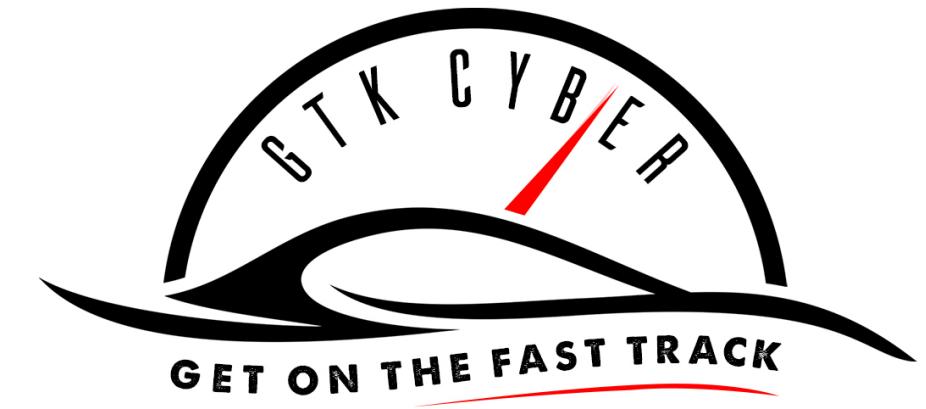
```
from bokeh.charts import Bar
from bokeh.io import show
data = [2,3,4,4,7,9,7,5]

barchart = Bar( data, notebook=True )
show( bar chart )
```



Method Chaining





```
barchart =  
    Bar( data, title="My Chart", xlabel="Categories", ylabel="Value",  
notebook=True )  
  
show( barchart )
```



```
barchart =  
    Bar( data, title="My Chart", xlabel="Categories", ylabel="Value",  
notebook=True )
```

```
barchart.show()
```

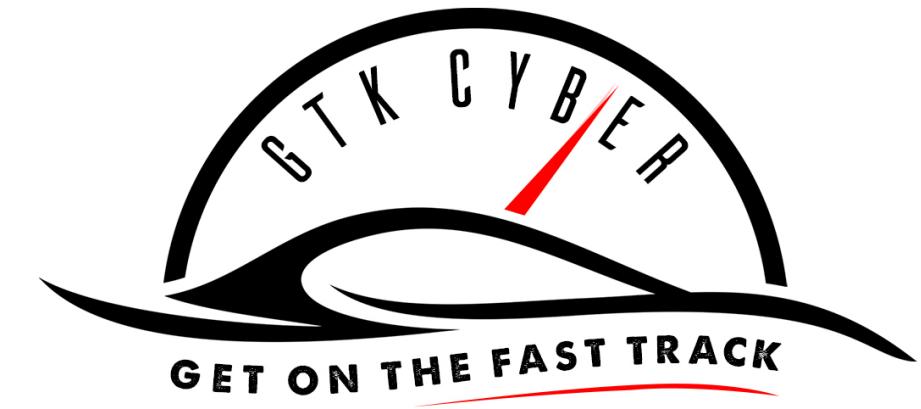
```
barchart3 = Bar(data)  
    .title( "My Chained Chart" )  
    .notebook( True )  
    .legend( True )  
    .xlabel( "Categories" )  
    .ylabel( "Value" )
```

```
show( barchart3 )
```



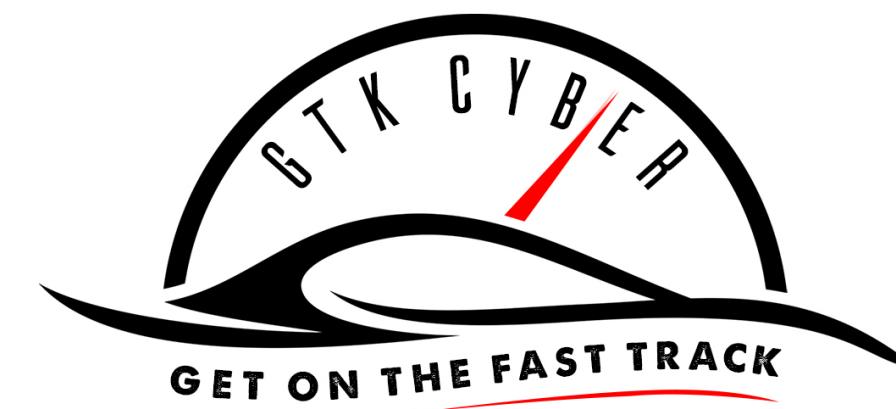
Bokeh Config Options

- **title** (str): the title of your plot.
- **xlabel** (str): the x-axis label of your plot.
- **ylabel** (str): the y-axis label of your plot.
- **legend** (str, bool): the legend of your plot.
- **xscale** (str): the x-axis type scale of your plot.
- **yscale** (str): the y-axis type scale of your plot.
- **width** (int): the width of your plot in pixels.
- **height** (int): the height of your plot in pixels.
- **tools** (bool): to enable or disable the tools in your plot.
- **filename** (str or bool): the name of the file where your plot will be written.
- **server** (str or bool): the name of your plot in the server.
- **notebook** (bool): if you want to output (or not) your plot into the IPython notebook.



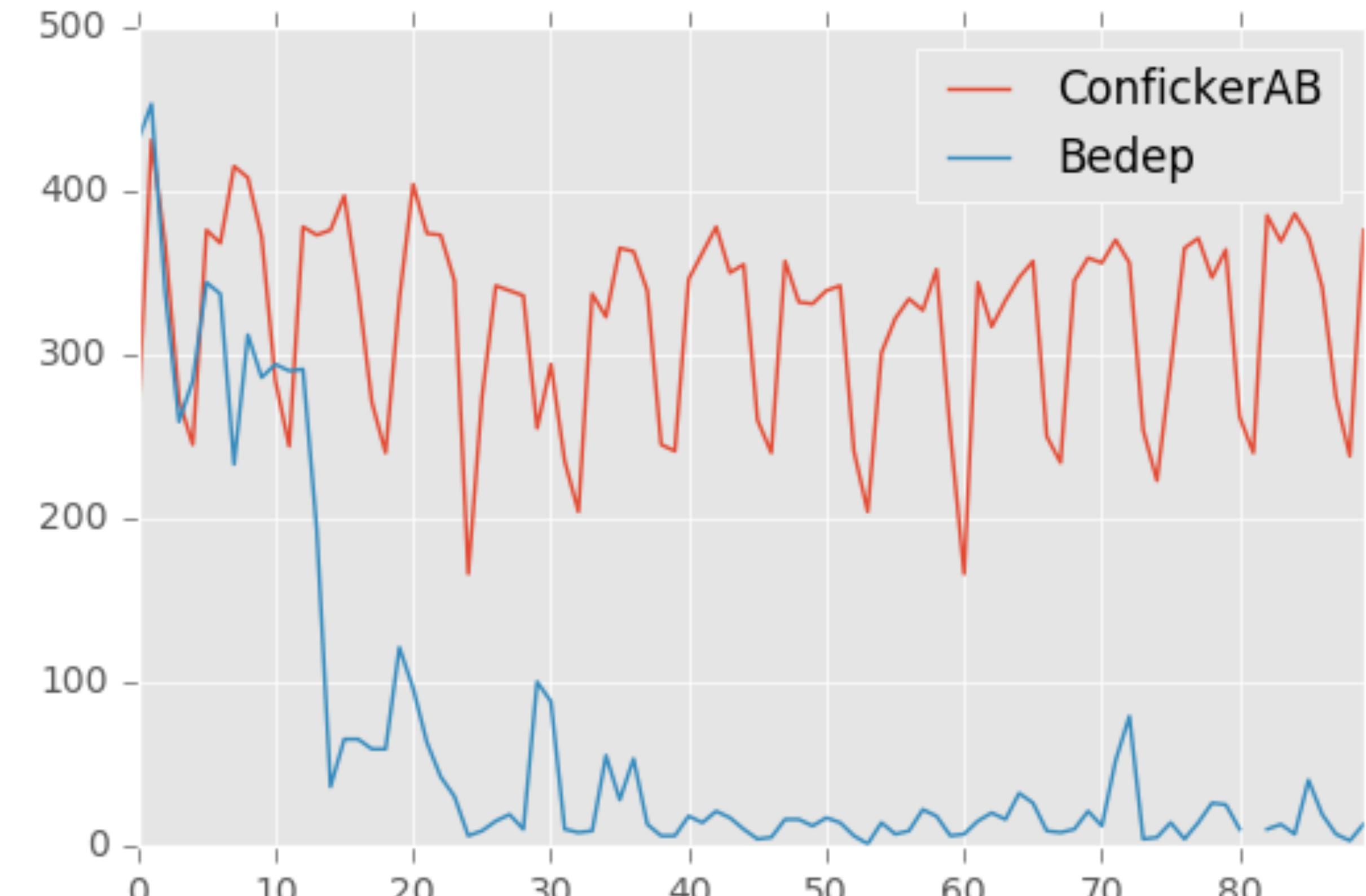
In Class Exercise

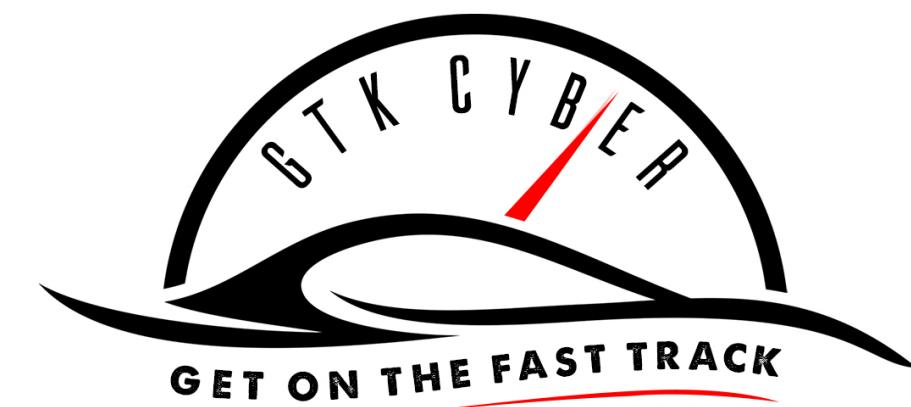
Please complete Worksheet 5: Data Visualization



Exercise 1

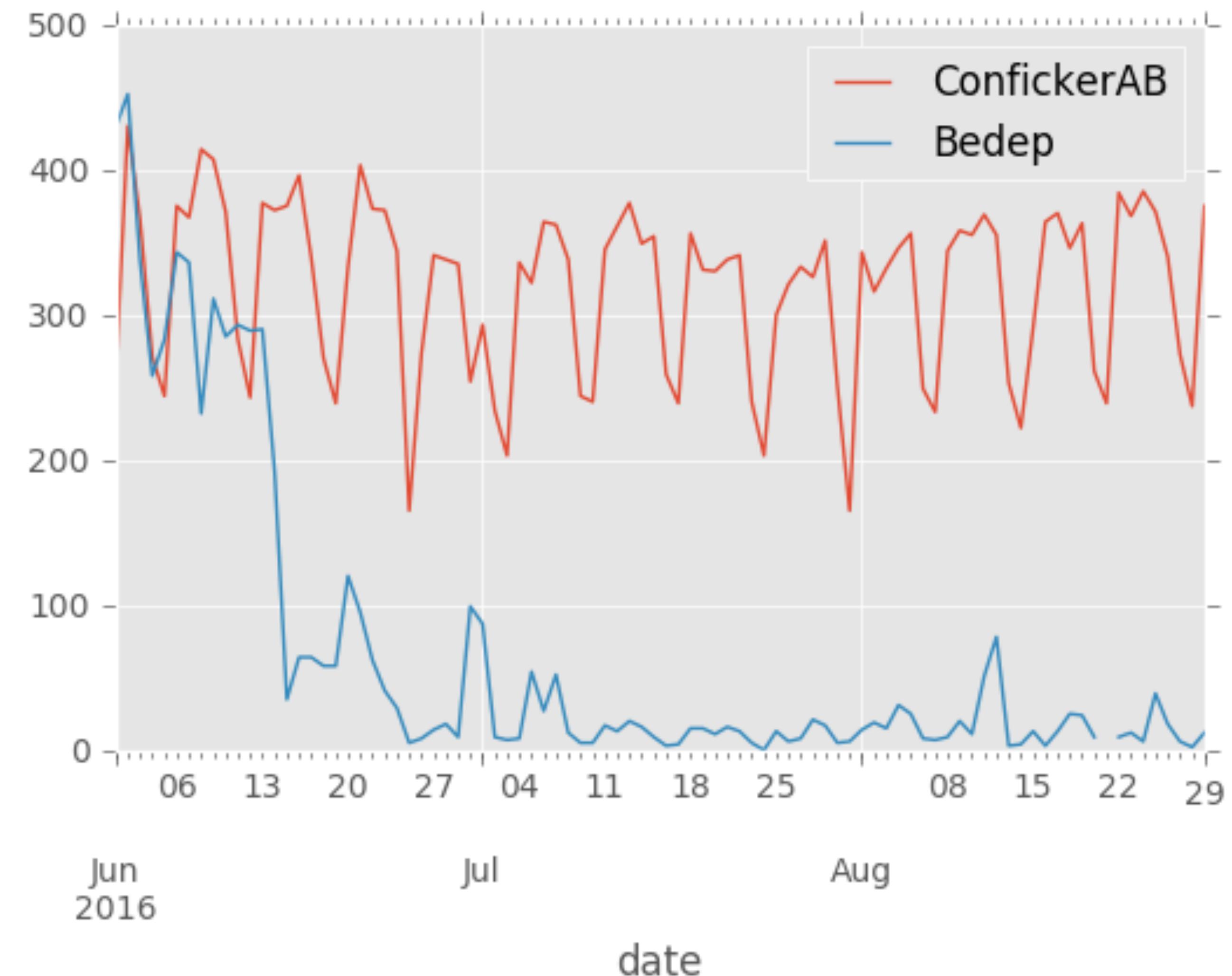
```
data = pd.read_csv('..../data/dailybots.csv')
filteredData = data[data['industry'] == "Government/Politics"]
filteredData2 = filteredData[filteredData['botfam'] == 'ConfickerAB'][['date', 'hosts']]
filteredData2.columns = ['date', 'ConfickerAB']
filteredData2.date = pd.to_datetime(filteredData2.date)
finalData.plot(kind='line')
```





Exercise 1

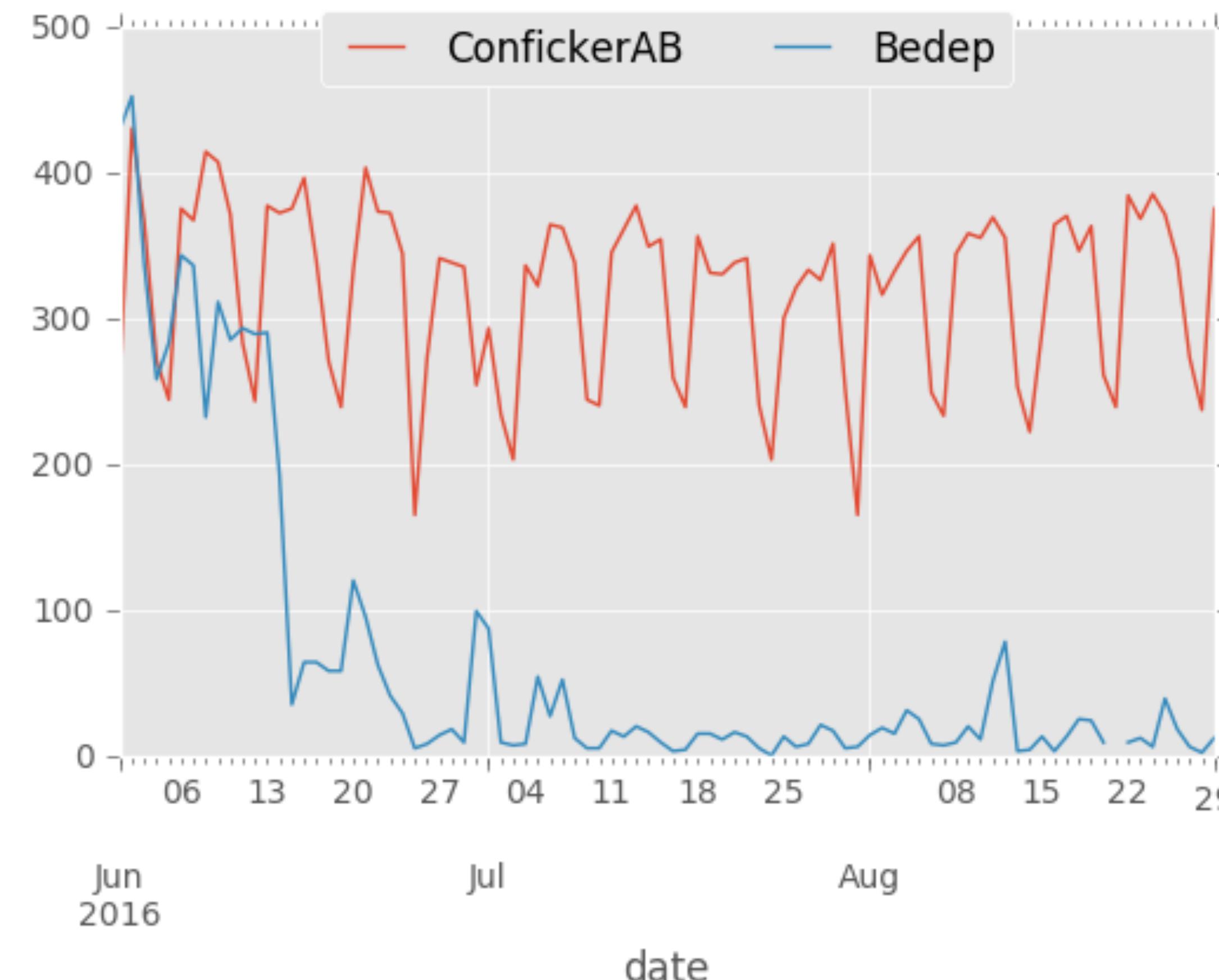
```
finalData.set_index('date', inplace=True)  
finalData.plot(kind='line')
```





Exercise 1

```
nicePlot = finalData.plot( kind="line")
nicePlot.legend(loc='upper center', bbox_to_anchor=(0.5, 1.05),
                 ncol=3, fancybox=True, shadow=False)
```

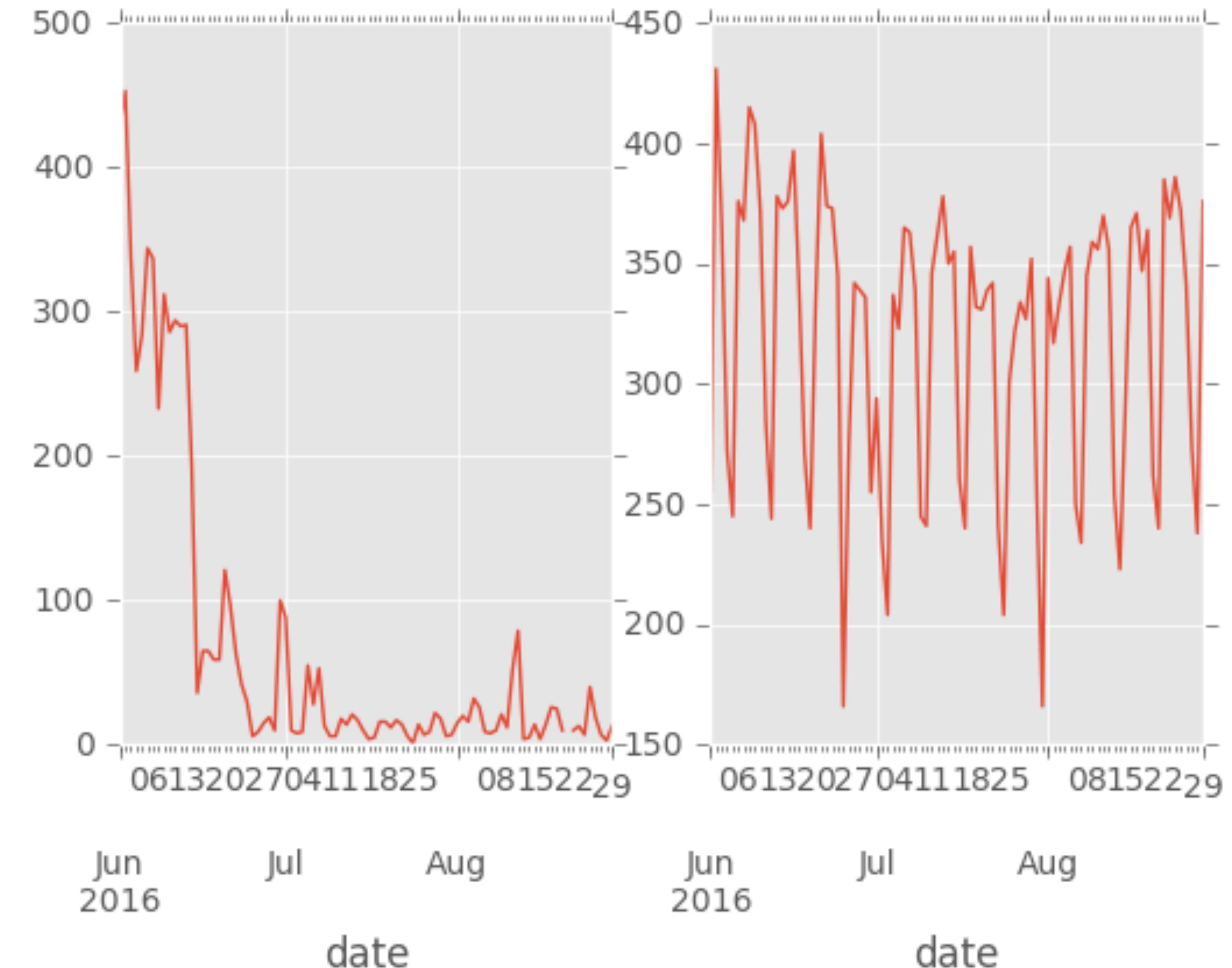


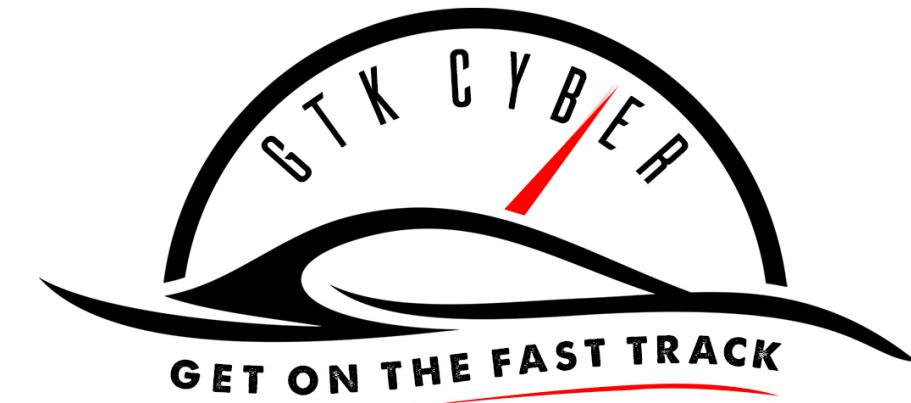


Exercise 1

```
fig, axes = plt.subplots(nrows=1, ncols=2)
```

```
finalData[ 'ConfickerAB' ].plot()  
finalData[ 'Bedep' ].plot(ax=axes[ 0 ])
```





Exercise 2

```
from bokeh.plotting import output_notebook
output_notebook()
from bokeh.charts import TimeSeries
from bokeh.io import show

linechart = TimeSeries( data=finalData,
                        title="ConfickerAB Hosts",
                        legend="top_left" )
show( linechart )
```



Exercise 2

