**Jericho McLeod**

**CSI-777 Assignment 2**

For provided glass dataset, plot histograms for numeric attributes. For provided iris data, create discretized intervals with consistent distributions.

First, we need to import appropriate libaries and the datasets.

```
In [1]:  import csv
         import matplotlib.pyplot as plt
         import numpy as np
```

```
In [2]:  glass = 'Glass&IrisDATA/glass.csv'
         iris = 'Glass&IrisDATA/iris.csv'

         def import_func(filename):
             data = []
             with open(filename) as csvfile:
                 csv_r = csv.reader(csvfile,delimiter = ',')
                 for i in csv_r:
                     data.append(i)
             return(data)

         glass_data = import_func(glass)
         iris_data = import_func(iris)
```

Now we need to examine the glass dataset to build a basic understanding of its contents.

```
In [3]:  for i in range(len(glass_data)):
             if i < 11:
                 print(glass_data[i])
```

```
['refractive_index', 'Sodium', 'Magnesium', 'Aluminum', 'Silicon', 'Pot
assium', 'Calcium', 'Barium', 'Iron', 'Type']
['1.51793', '12.79', '3.5', '1.12', '73.03', '0.64', '8.77', '0', '0',
"'build wind float'"]
['1.51643', '12.16', '3.52', '1.35', '72.89', '0.57', '8.53', '0', '0',
"'vehic wind float'"]
['1.51793', '13.21', '3.48', '1.41', '72.64', '0.59', '8.43', '0', '0',
"'build wind float'"]
['1.51299', '14.4', '1.74', '1.54', '74.55', '0', '7.59', '0', '0', 'ta
bleware']
['1.53393', '12.3', '0', '1', '70.16', '0.12', '16.19', '0', '0.24',
"'build wind non-float'"]
['1.51655', '12.75', '2.85', '1.44', '73.27', '0.57', '8.79', '0.11',
'0.22', "'build wind non-float'"]
['1.51779', '13.64', '3.65', '0.65', '73', '0.06', '8.93', '0', '0',
"'vehic wind float'"]
['1.51837', '13.14', '2.84', '1.28', '72.85', '0.55', '9.07', '0', '0',
"'build wind float'"]
['1.51545', '14.14', '0', '2.68', '73.39', '0.08', '9.07', '0.61', '0.0
5', 'headlamps']
['1.51789', '13.19', '3.9', '1.3', '72.33', '0.55', '8.44', '0', '0.2
8', "'build wind non-float'"]
```

The data is structured with headers in the first row and instances on subsequent rows. Since we will be dealing with attributes independently, converting it to a dictionary will simplify tasks.

We will also check input lengths here.

```
In [4]: headers = glass_data[0] # extract the attribute names
        glass_data_dict = {}
        glass_hist_dict = {}
        for i in headers:
            glass_data_dict[i] = []
            glass_hist_dict[i] = {}
        for i in range(len(glass_data)):
            if i > 0: # to skip the attribute names in the loop
                for j in range(len(glass_data[i])):
                    try: # try adding data as floats, on exception as as strings
                        glass_data_dict[headers[j]].append(float(glass_data[i][j
]))
                    except:
                        glass_data_dict[headers[j]].append(glass_data[i][j])

        for k,v in glass_data_dict.items():
            print(k,'contains',len(v),'values')
```

```
refractive_index contains 214 values
Sodium contains 214 values
Magnesium contains 214 values
Aluminum contains 214 values
Silicon contains 214 values
Potassium contains 214 values
Calcium contains 214 values
Barium contains 214 values
Iron contains 214 values
Type contains 214 values
```

Since it is easy to spot-check the first instance, we will use it as the determining 'type' to identify which attributes contain numeric vs. textual data.

To get histogram counts for an attribute, we need to know the range across which it is distributed, then divide that by some amount and count the instances which fall into the defined subset. Below, 'n_categories' specifies the histogram width.

```
In [5]:  n_categories = 20

         for i in headers:
             glass_hist_dict[i] = {}

         for k,v in glass_data_dict.items():
             if isinstance(v[0],float) or isinstance(v[0],int):
                 min_v   = min(v)
                 max_v   = max(v)
                 range_v = max_v-min_v
                 print(k,min_v,max_v)
                 step = range_v/n_categories
                 for i in range(n_categories):
                     lower = min_v + step*i
                     upper = lower + step
                     name = '>= '+str(lower)+' and < '+str(upper)
                     glass_hist_dict[k][name] = 0
                     for j in v:
                         if j >= lower and j < upper:
                             glass_hist_dict[k][name] +=1
```
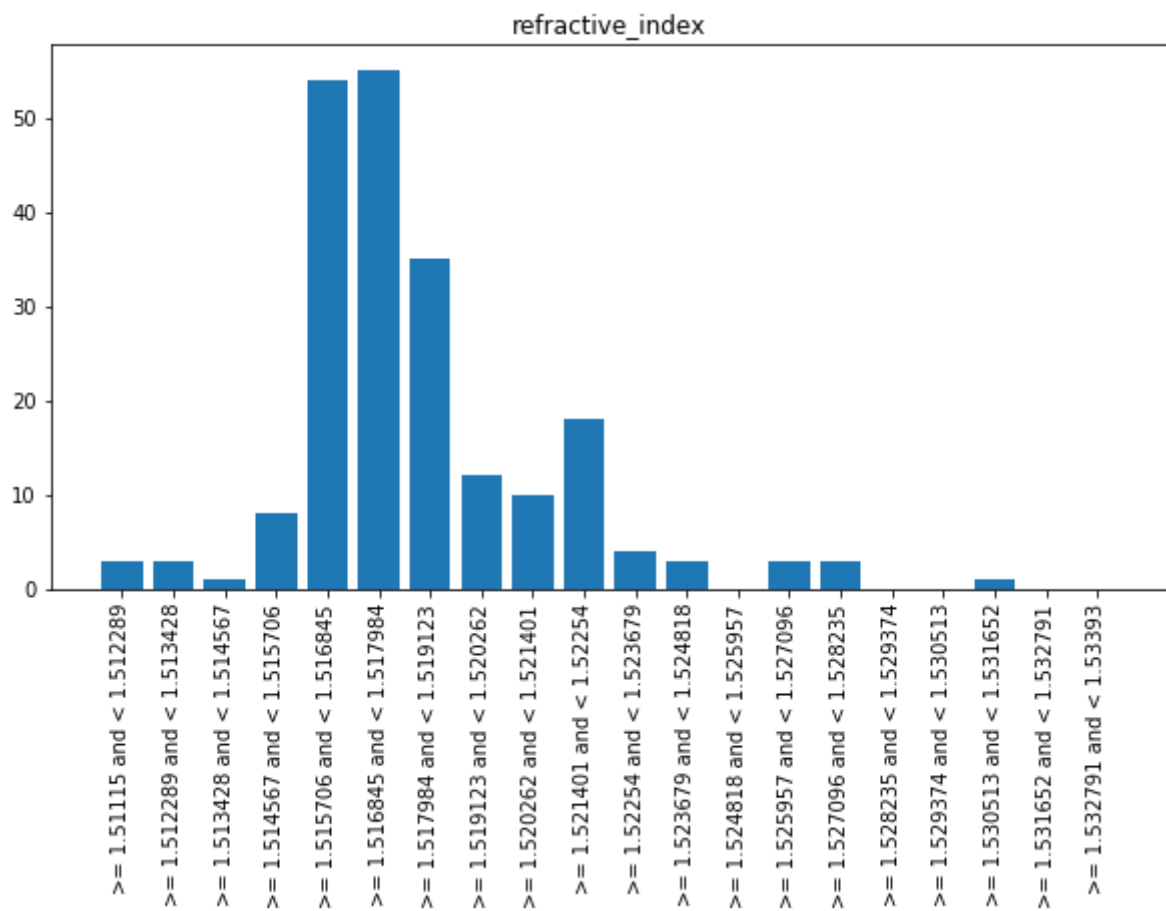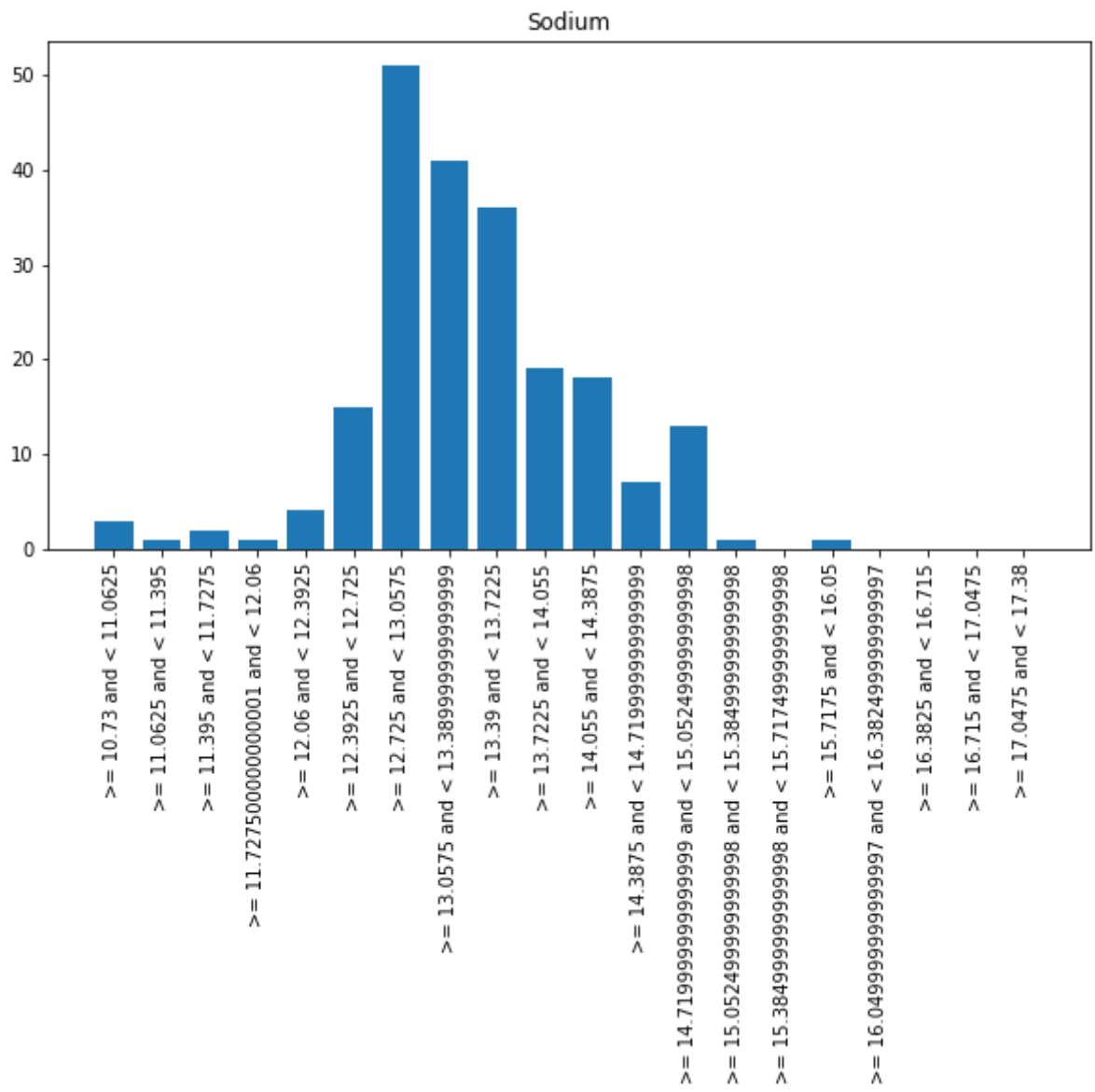
```
refractive_index 1.51115 1.53393
Sodium 10.73 17.38
Magnesium 0.0 4.49
Aluminum 0.29 3.5
Silicon 69.81 75.41
Potassium 0.0 6.21
Calcium 5.43 16.19
Barium 0.0 3.15
Iron 0.0 0.51
```
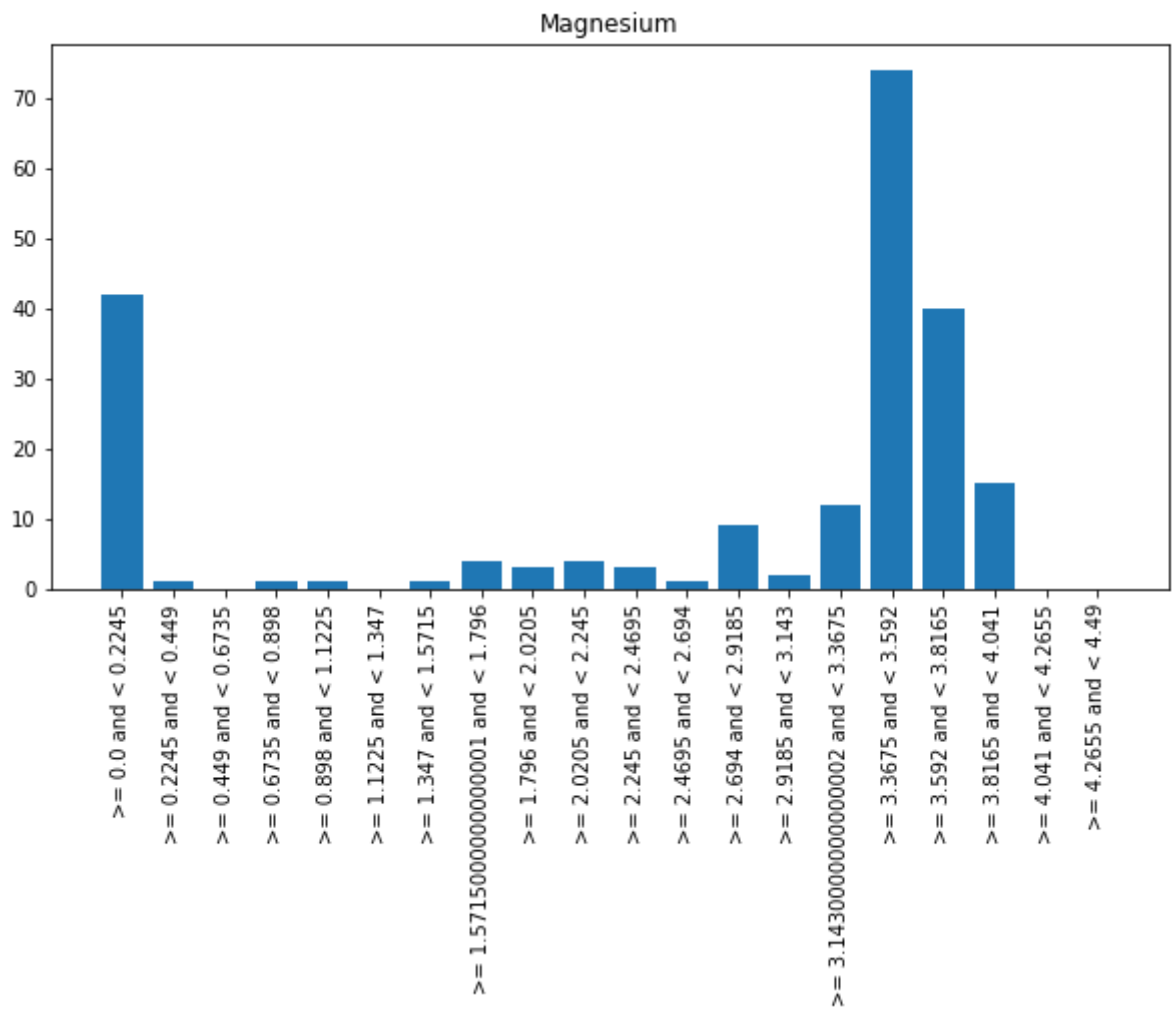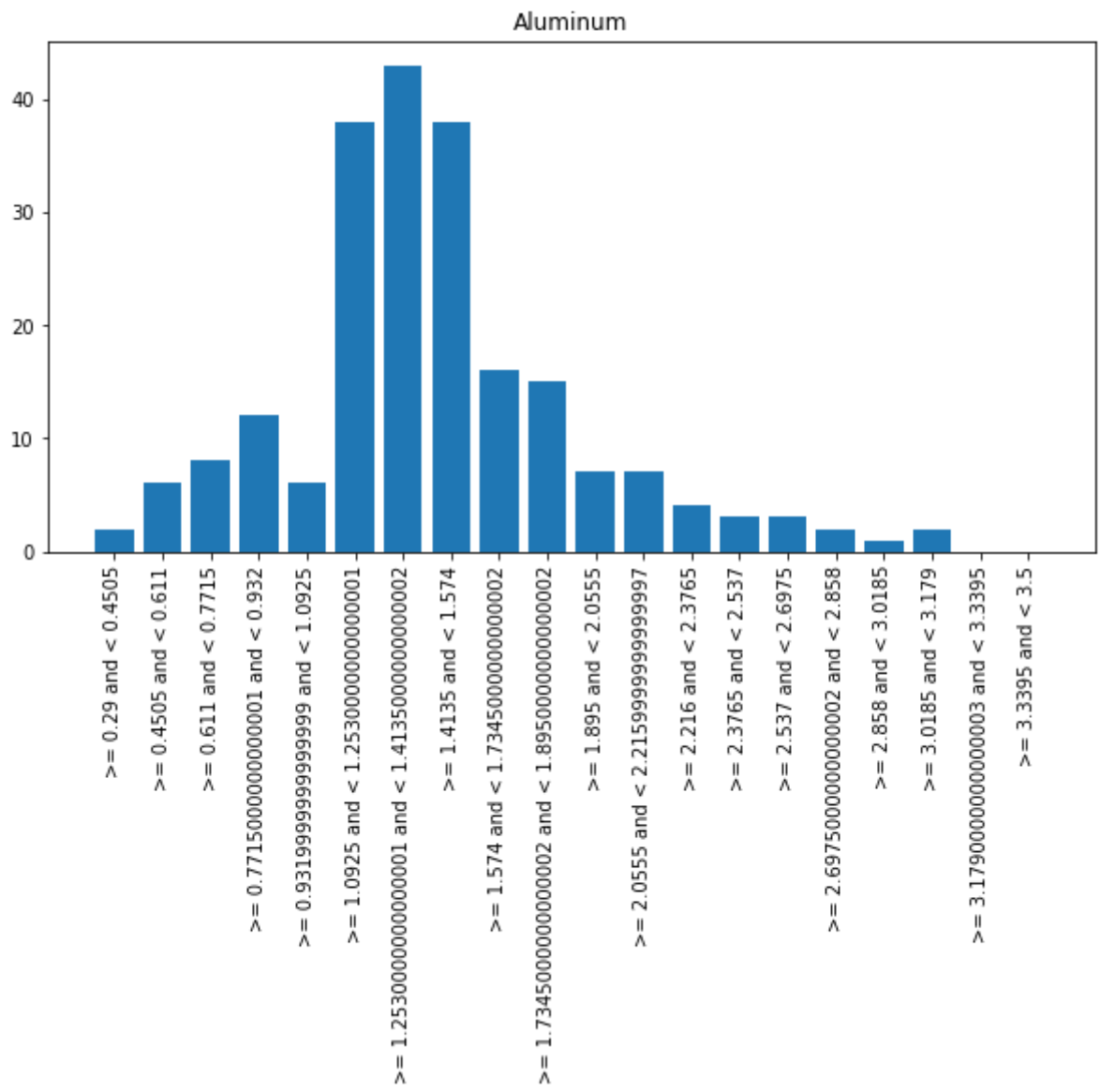
All that remains is to show plots for the actual histograms.

```python
In [6]: for k,v in glass_hist_dict.items():
            if len(v)>0:
                x = []
                x_labels = []
                for k2,v2 in v.items():
                    x.append(v2)
                    x_labels.append(k2)
                plt.figure(figsize=(10,5))
                plt.bar(x_labels,x)
                plt.xticks(x_labels,rotation='vertical')
                plt.title(k)
                plt.show()
```
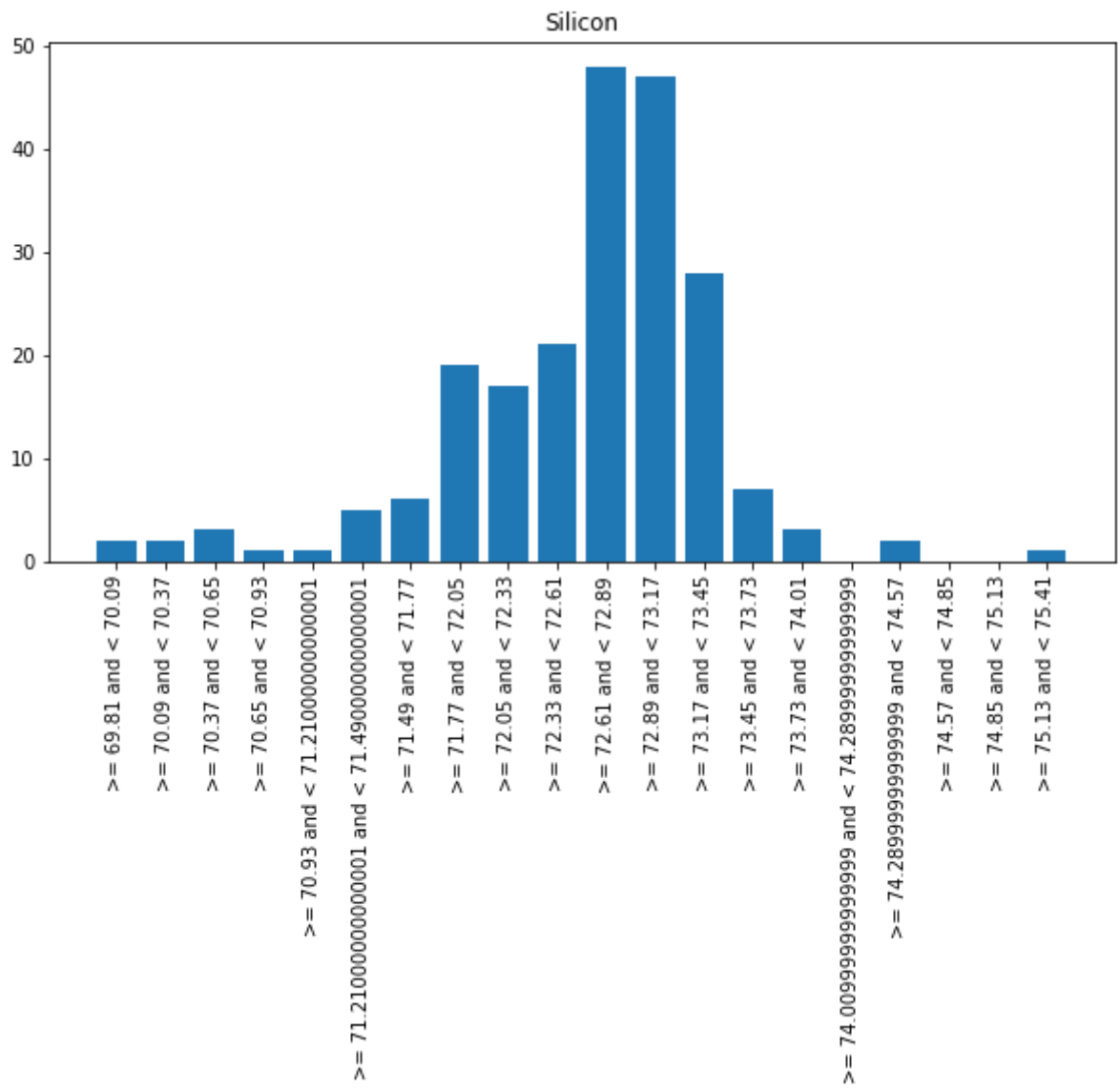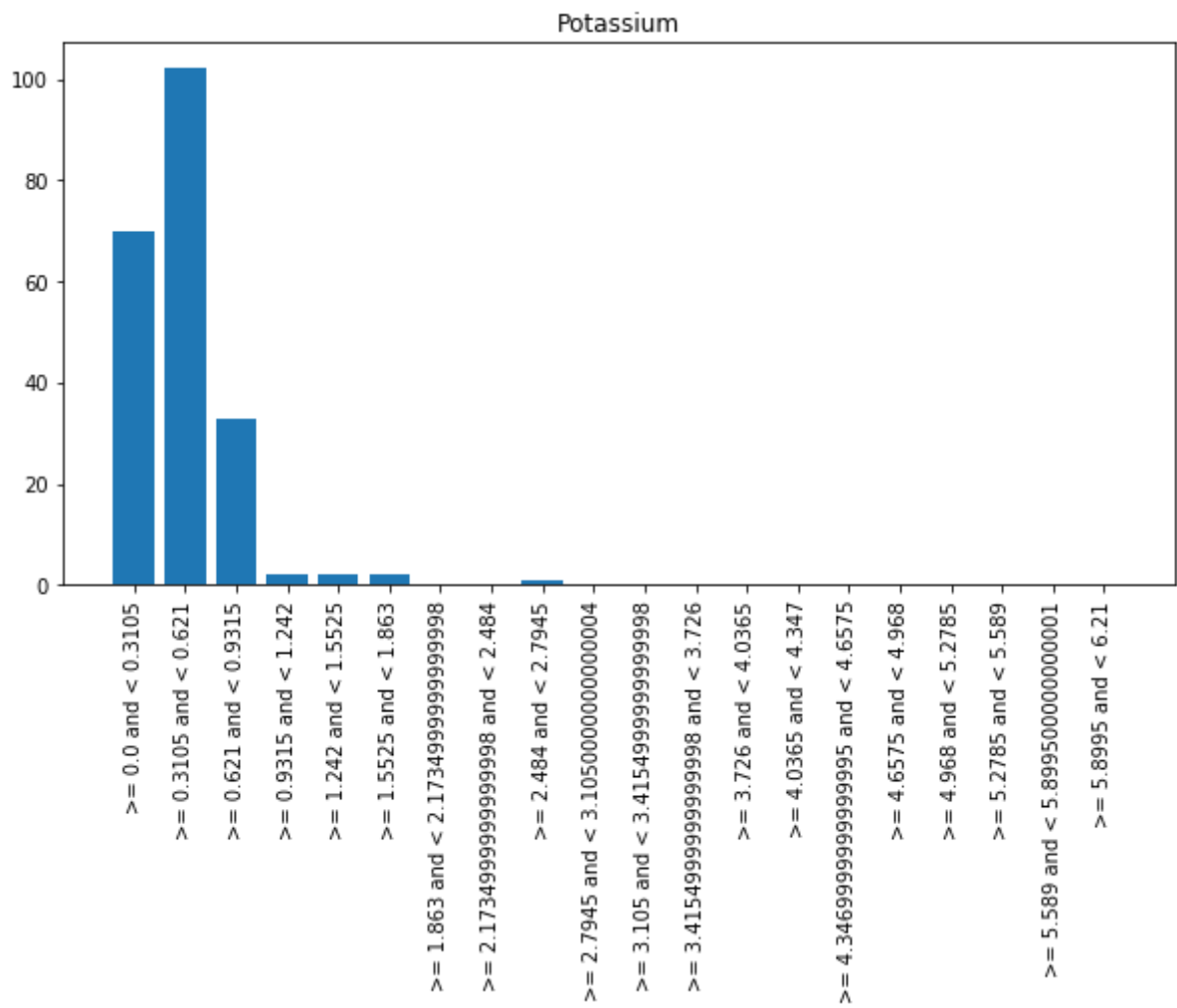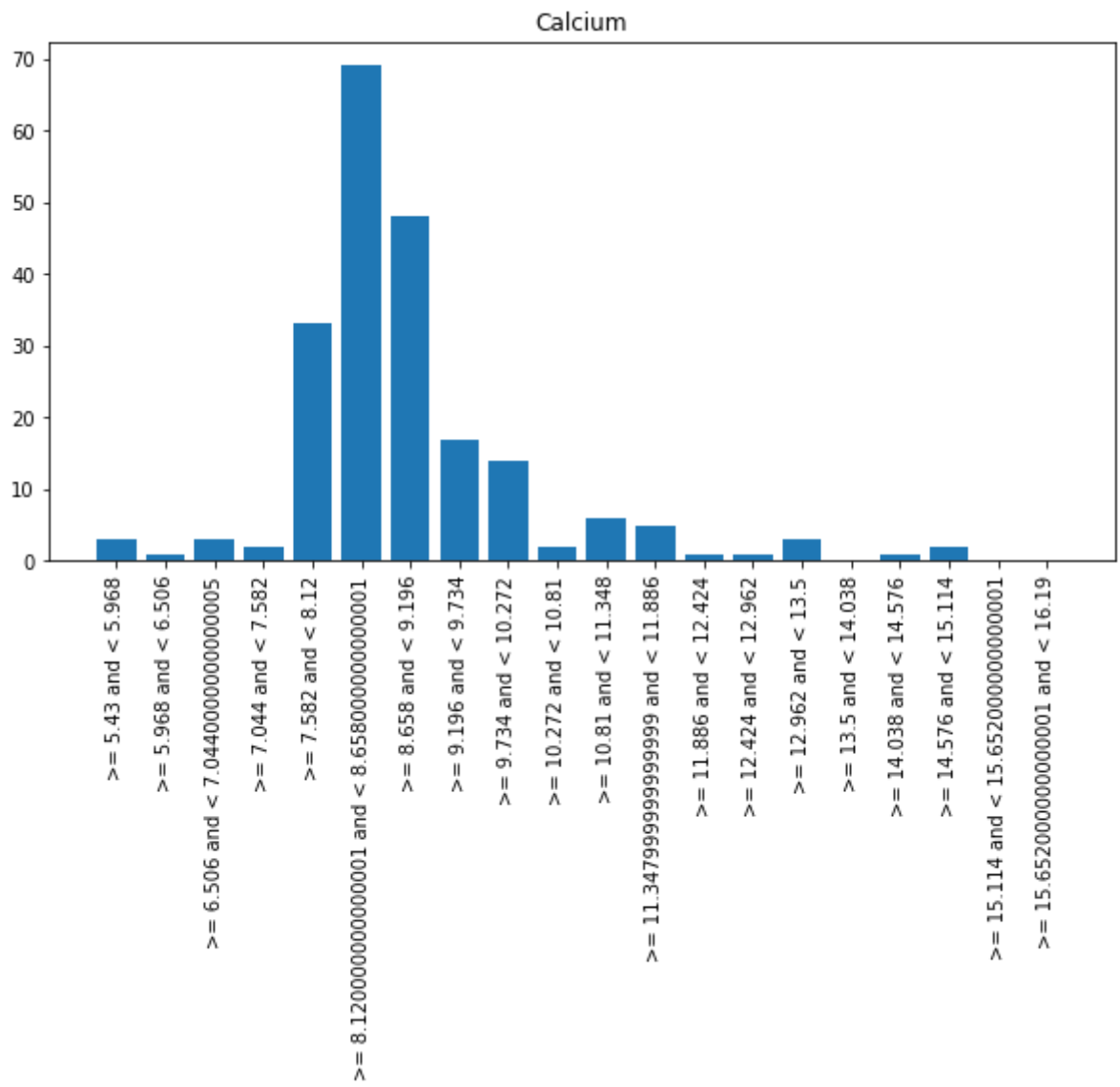
refractive_index

Sodium

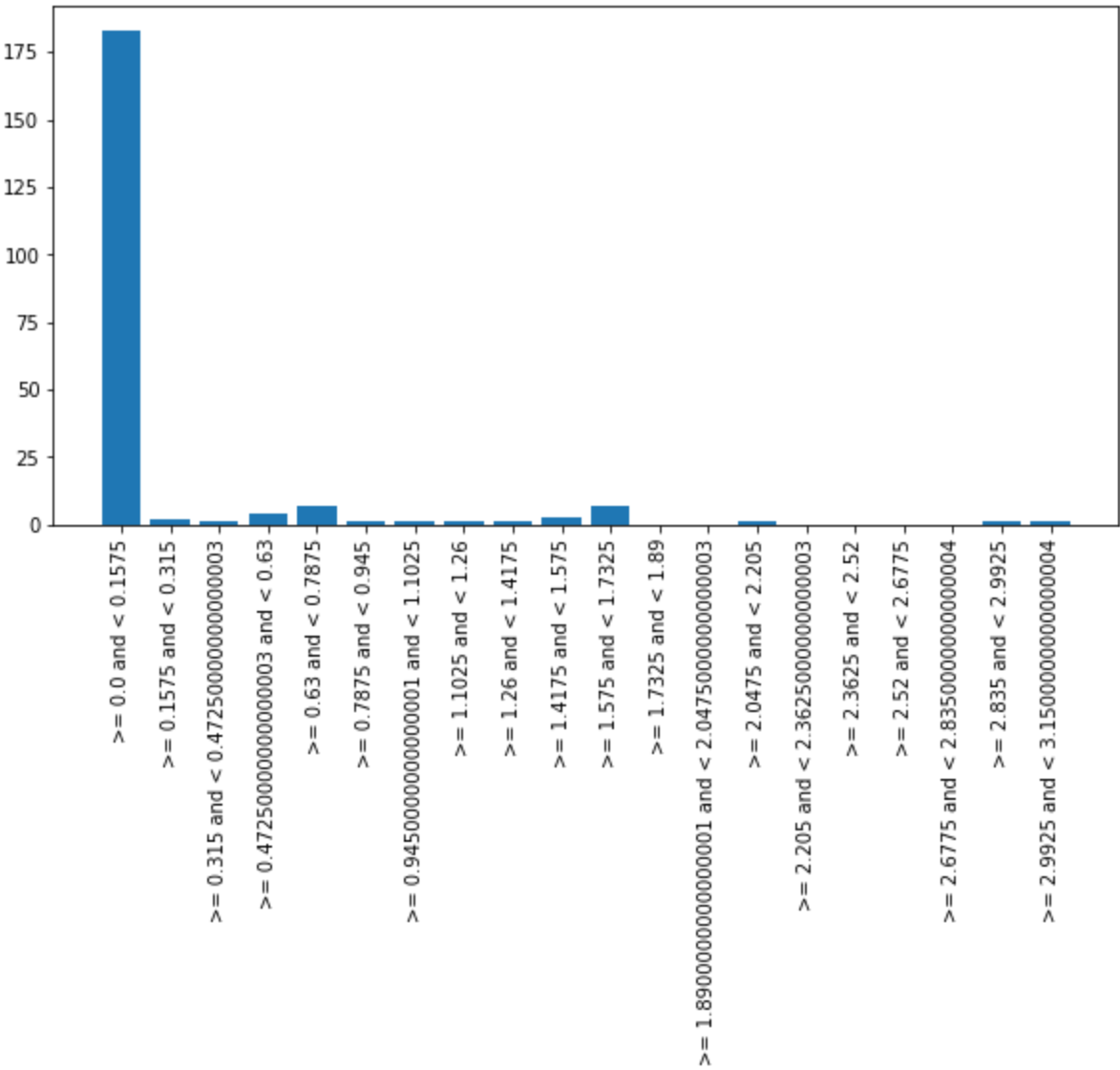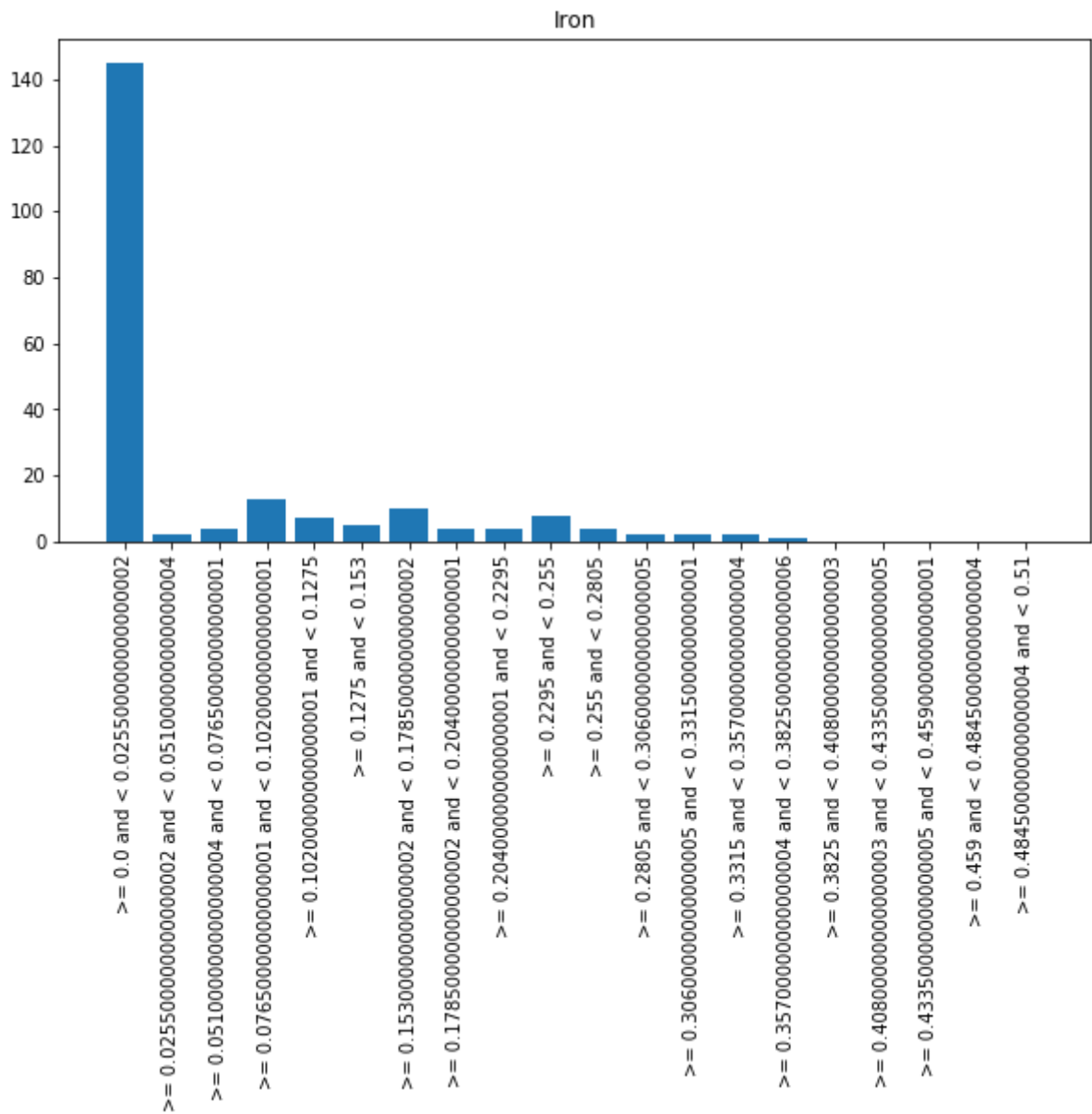Magnesium

Aluminum

Silicon

Potassium

Calcium

Barium

**Iron**

Moving on to the iris dataset, where we need to create discretized intervals with consistent distributions, unlike we processed for the Glass dataset.

```
In [7]: for i in range(len(iris_data)):
            if i < 11:
                print(iris_data[i])

        print('\nTotal instances',len(iris_data)-1)

        ['petallength', 'petalwidth', ' class']
        ['1.4', '0.2', 'Iris-setosa']
        ['1.4', '0.2', 'Iris-setosa']
        ['1.3', '0.2', 'Iris-setosa']
        ['1.5', '0.2', 'Iris-setosa']
        ['1.4', '0.2', 'Iris-setosa']
        ['1.7', '0.4', 'Iris-setosa']
        ['1.4', '0.3', 'Iris-setosa']
        ['1.5', '0.2', 'Iris-setosa']
        ['1.4', '0.2', 'Iris-setosa']
        ['1.5', '0.1', 'Iris-setosa']

        Total instances 150
```

Since there are just two attributes, we will store them in individual variables.

```
In [8]: petallength,petalwidth = [],[]
        for i in range(len(iris_data)):
            if i>0:
                petallength.append(iris_data[i][0])
                petalwidth.append(iris_data[i][1])
```

To discretize the data into $n$ categories, we need only sort the data, determine the number of items per category as $length/n$, and find the values at indexes of multiples of that value.

However, given that the iris data is relatively short in terms of significant digits, it is entirely possible that a chosen value could be in the middle of a collection of identical values, which will somewhat disrupt distributions.

```
In [9]: def create_indexes(variable,n):
            length = len(variable)
            cat_size = length/n
            indexes = []
            for i in range(n):
                min_val = cat_size * i
                max_val = (cat_size * (i+1))-1
                indexes.append((min_val,max_val))
            return(indexes)


        n = 10
        index_vals = create_indexes(petallength,n)
        for i in index_vals:
            print(i)
```

```
(0.0, 14.0)
(15.0, 29.0)
(30.0, 44.0)
(45.0, 59.0)
(60.0, 74.0)
(75.0, 89.0)
(90.0, 104.0)
(105.0, 119.0)
(120.0, 134.0)
(135.0, 149.0)
```

Next we need to convert the indexes to actual values in the dataset. If we skip this step the resulting categories will be completely equal distributions, but we may also see identical values in different categories arbitrarily.

```
In [10]: def find_values_by_index(variable,indexes):
             values = []
             for i in indexes:
                 low = int(i[0])
                 high = int(i[1])
                 sorted_var = sorted(variable)
                 values.append((sorted_var[low],sorted_var[high]))
             return(values)

         plength_cats = find_values_by_index(petallength,index_vals)
         pwidth_cats = find_values_by_index(petalwidth,index_vals)

         print(plength_cats)
         print(pwidth_cats)
```

```
[('1', '1.4'), ('1.4', '1.5'), ('1.5', '1.7'), ('1.7', '3.9'), ('3.9',
'4.3'), ('4.4', '4.6'), ('4.7', '5'), ('5', '5.3'), ('5.4', '5.8'),
('5.8', '6.9')]
[('0.1', '0.2'), ('0.2', '0.2'), ('0.2', '0.4'), ('0.4', '1.1'), ('1.
2', '1.3'), ('1.3', '1.5'), ('1.5', '1.8'), ('1.8', '1.9'), ('1.9', '2.
2'), ('2.2', '2.5')]
```

It is readily apparent that the distributions will be unequal. This is an artifact of having data points which do not divide evenly over our number of categories. For datasets that contain no duplicate values this will not be the case.
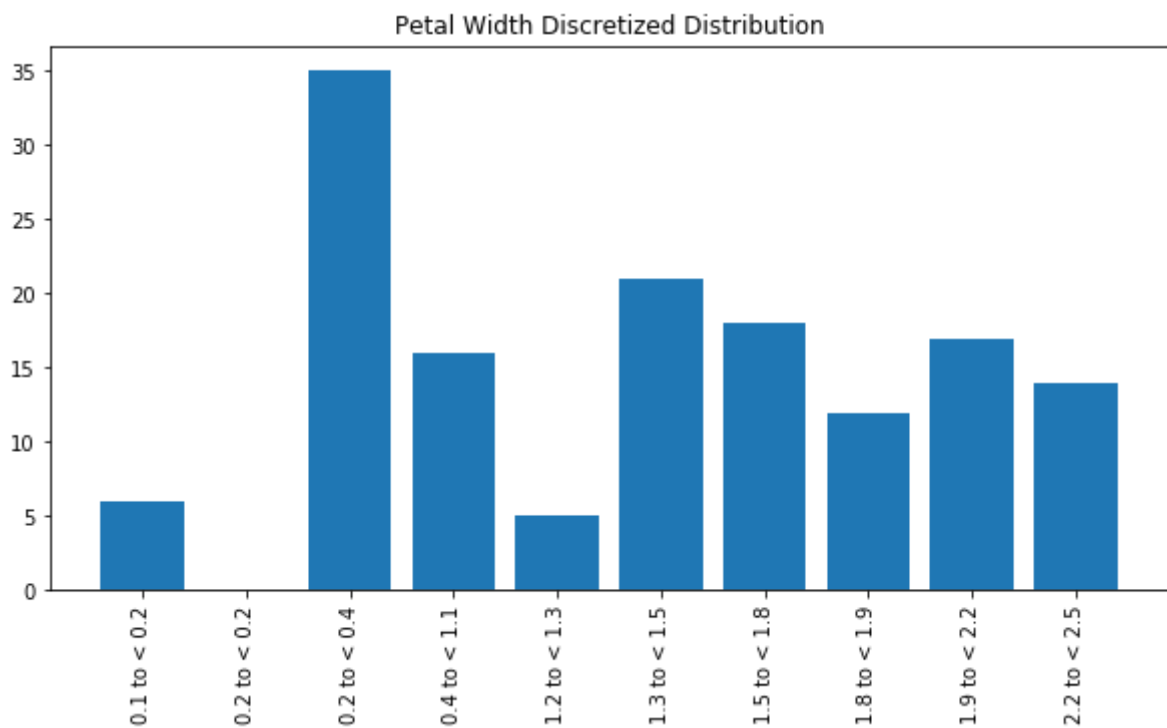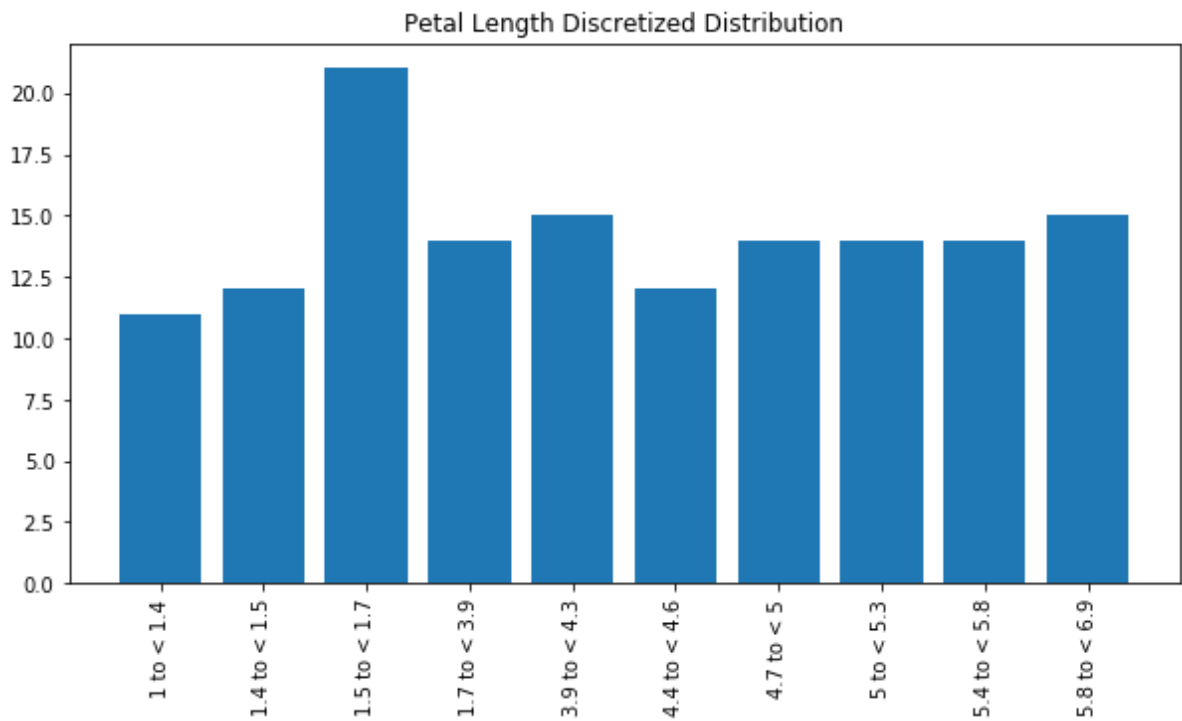
We will use >= the lower value, and < the upper value, so that no data point is used twice and so that identical values are not in multiple bins. At this point, all we need to do is show that our intervals are approximately even distributions, which we will do in histograms.

```python
In [11]: def get_vis_data(indexes,variable):
             histo_data = {}
             for i in range(len(indexes)):
                 name = str(indexes[i][0])+' to < '+str(indexes[i][1])
                 histo_data[name] = 0
                 for j in variable:
                     if j >= indexes[i][0] and j < indexes[i][1]:
                         histo_data[name]+=1
             return(histo_data)

         plength_hist_data = get_vis_data(plength_cats,petallength)
         pwidth_hist_data = get_vis_data(pwidth_cats,petalwidth)
```

```
In [12]:  def display_vis(histo_dict,string):
              x_labels,x = [],[]
              for k,v in histo_dict.items():
                  x_labels.append(k)
                  x.append(v)
              plt.figure(figsize=(10,5))
              plt.bar(x_labels,x)
              plt.xticks(x_labels,rotation='vertical')
              plt.title(string)
              plt.show()

          display_vis(plength_hist_data,'Petal Length Discretized Distribution')
          display_vis(pwidth_hist_data,'Petal Width Discretized Distribution')
```

## Petal Length Discretized Distribution



## Petal Width Discretized Distribution



Petal length discretizes easily into bins, while petal width has some apparent issues. For one, there is a category that was completely encompassed by a single measurement in the data. Since the script looks for >= a minimum and < a maximum, this means that it is an empty category. This could be dropped by adding a simple if value > 0 statement when building the visualizations.

Further, because of the quantity of values with a 0.2 measurement in Petal Width, there is a large category where they did get counted.

A manual distribution may be a better choice, but let's first consider the actual histogram.

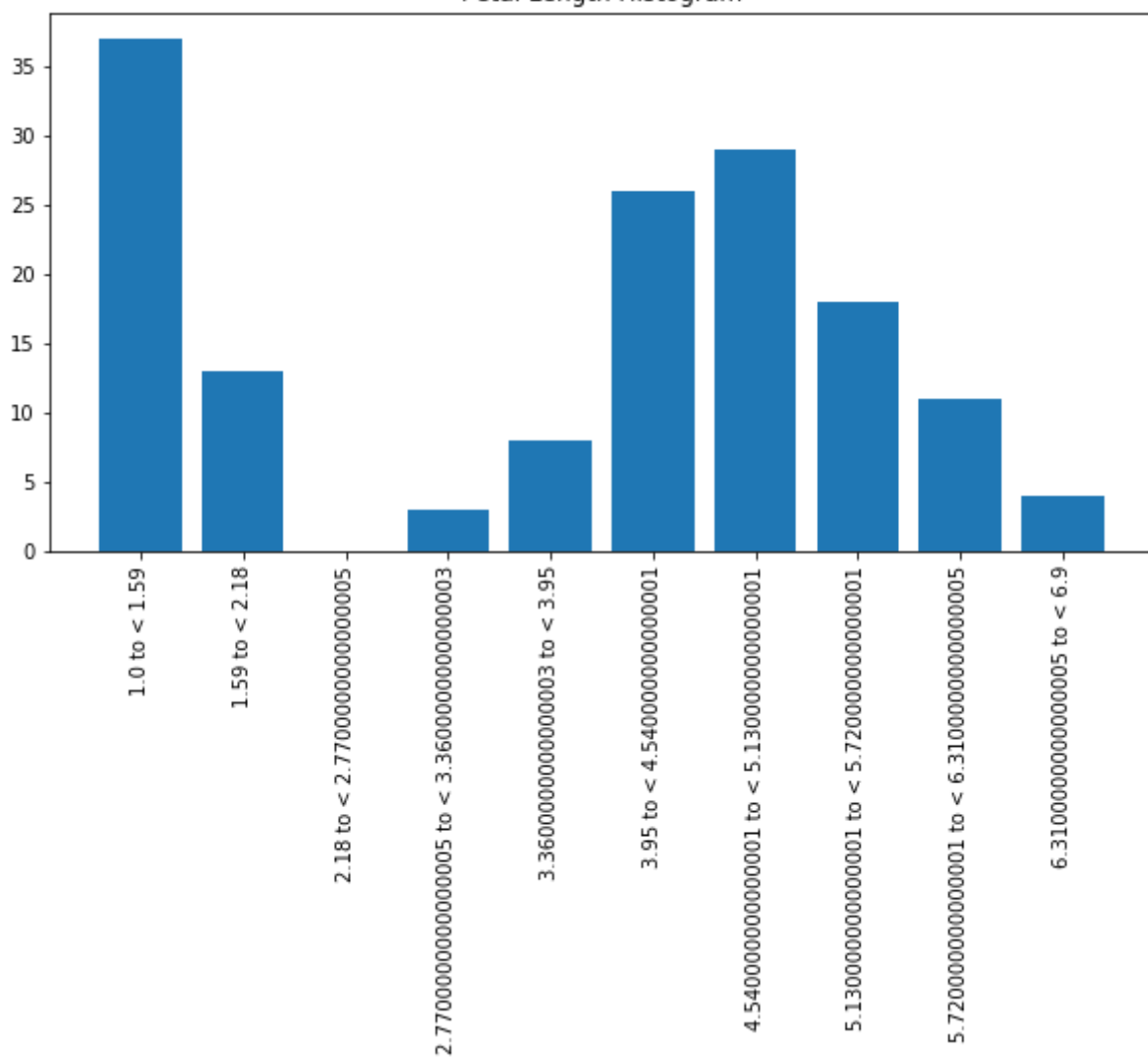In [ ]:

```
In [13]: def create_true_hist(variable,n,string):
             min_val = float(min(variable))
             max_val = float(max(variable))
             range_var = max_val-min_val
             step = range_var/n
             x,x_lab = [],[]
             for i in range(n):
                 min_step = min_val + (step*i)
                 max_step = min_val + (step *(i+1))
                 x_lab.append(str(min_step)+' to < '+str(max_step))
                 count = 0
                 for j in range(len(variable)):
                     if float(variable[j]) >= min_step and float(variable[j]) < m
         ax_step:
                         count+=1
                 x.append(count)
             plt.figure(figsize=(10,5))
             plt.bar(x_lab,x)
             plt.xticks(x_lab,rotation='vertical')
             plt.title(string)
             plt.show()

         n = 10
         create_true_hist(petallength,n,'Petal Length Histogram')
         create_true_hist(petalwidth,n,'Petal Width Histogram')
```
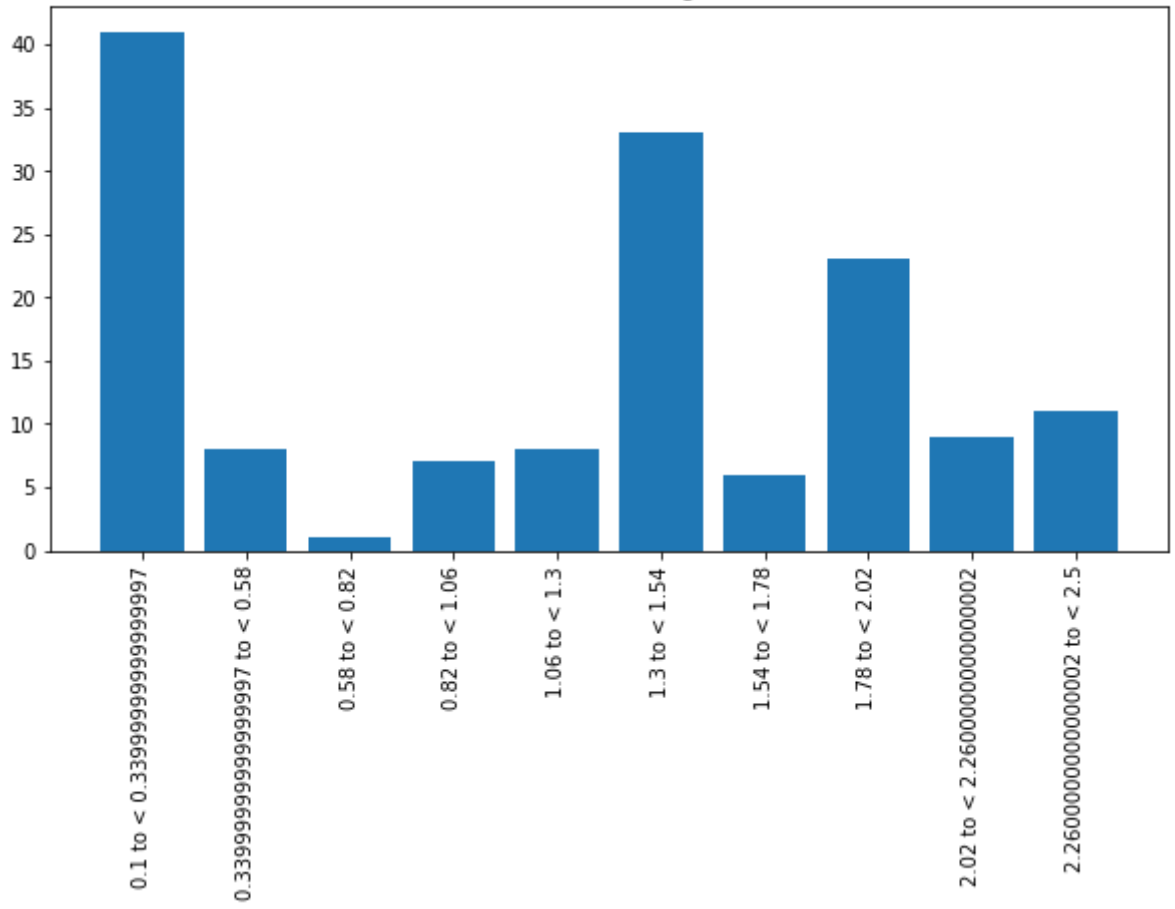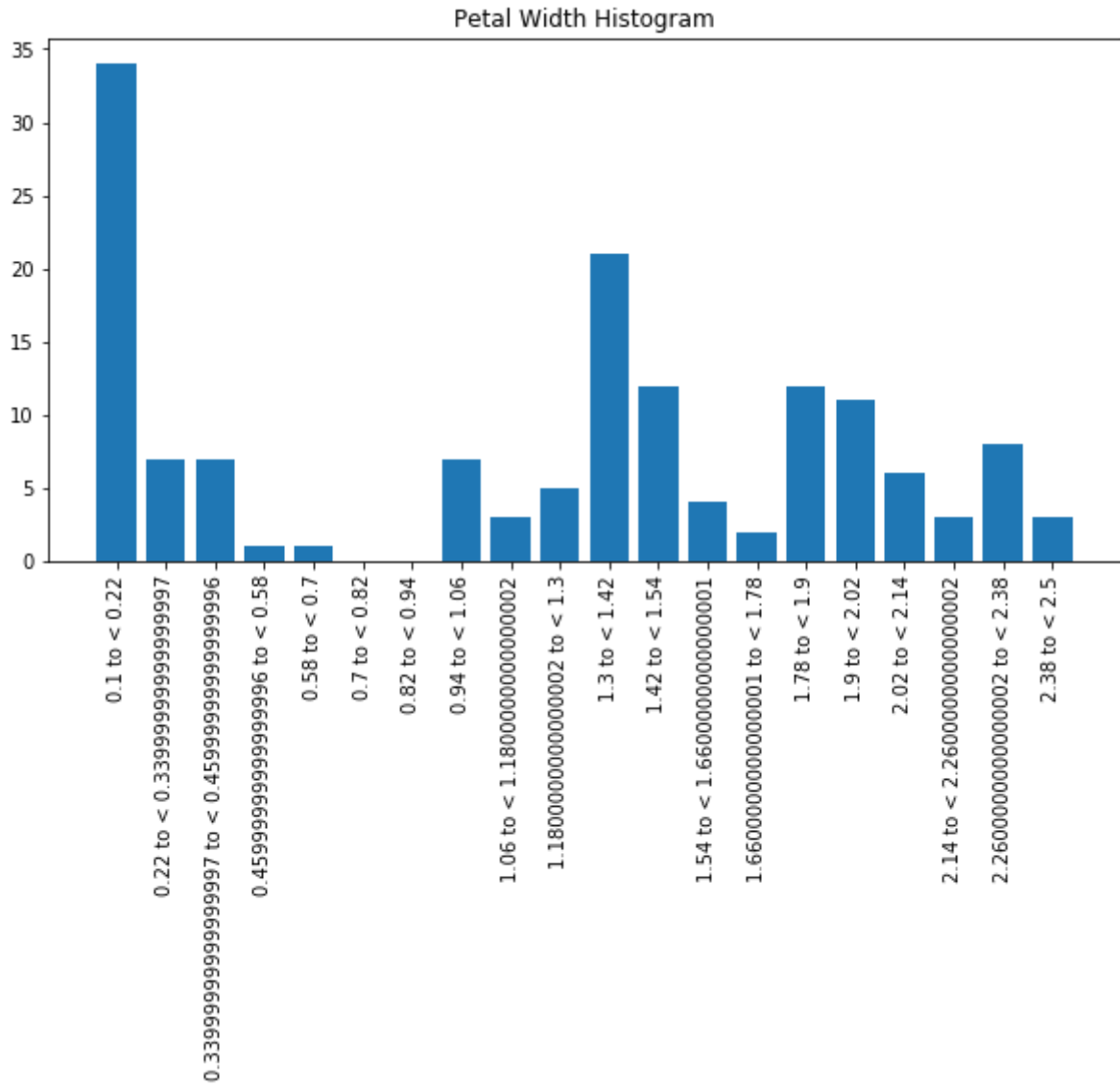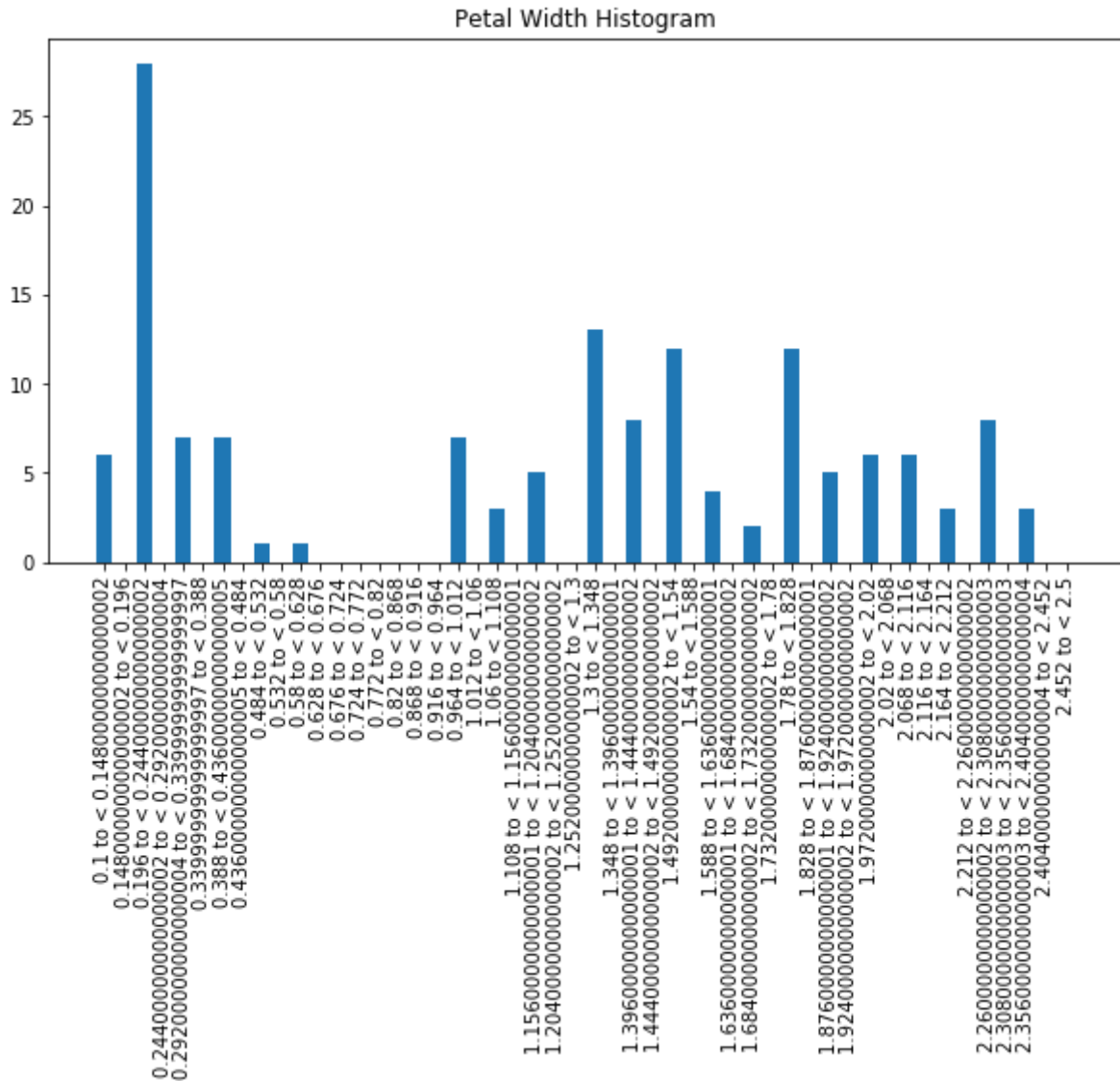
Petal Length Histogram

Petal Width Histogram

```
n=20
create_true_hist(petalwidth,n,'Petal Width Histogram')
```

Petal Width Histogram

```
n=50
create_true_hist(petalwidth,n,'Petal Width Histogram')
```



Petal Width Histogram

Here we have reached a point where we have more categories in the visualization than we have categories in the data. It is unlikely, then, that we will be able to separate the data for pedal width algorithmically into level distributions entirely without having attributes with the same measurement in multiple categories.

Just for fun, since I haven't done it before, let's look at a heat map comparing the quantities of data points using a heatmap with both variables.

```
In [16]:  def convert_to_numeric(data):
              for i in range(len(data)):
                  for j in range(len(data[i])):
                      try:
                          data_point = float(data[i][j])
                      except:
                          data_point = 'na'
                      if data_point != 'na':
                          data[i][j] = data_point
              return(data)

          def create_heatmap(dataset):
              # petal length will be x, at i[0]
              # petal width will be y, at i[1]
              headers = dataset[0]
              data = dataset[1:]
              x,y = [],[]
              for i in data:
                  x.append(i[0])
                  y.append(i[1])
              heatmap, xedges, yedges = np.histogram2d(x, y, bins=(30,60))
              extent = [xedges[0], xedges[-1], yedges[0], yedges[-1]]
              plt.clf()
              plt.title('Iris Petal Heatmap')
              plt.ylabel('Petal Width')
              plt.xlabel('Petal Lengh')
              plt.imshow(heatmap, extent=extent,cmap='Greys')
              plt.show()

          converted_iris_data = convert_to_numeric(iris_data)
          create_heatmap(converted_iris_data)
```
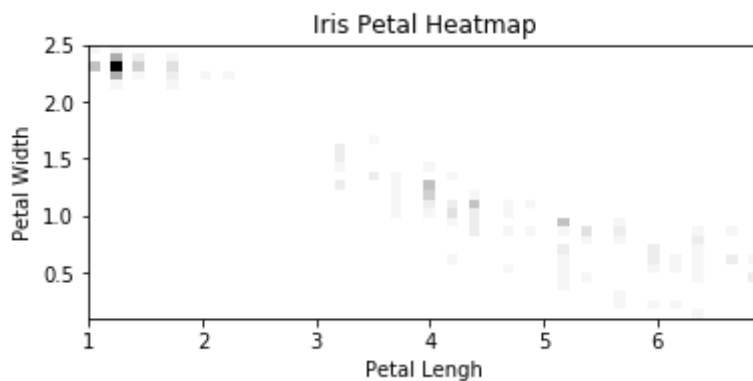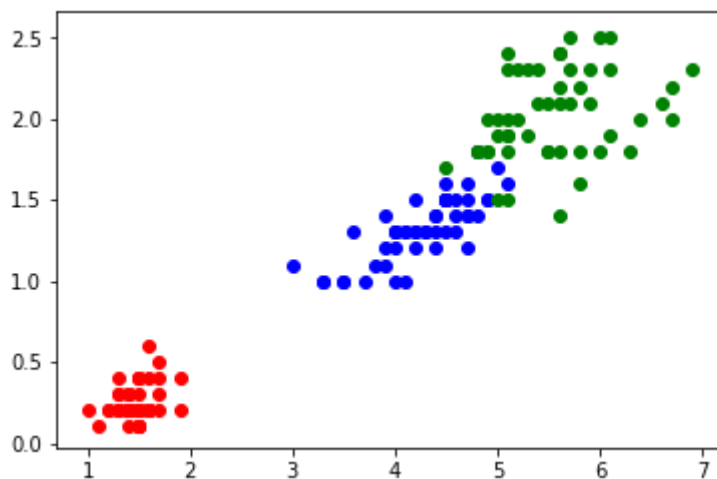
```
In [17]:  def create_plot(dataset):
              headers = dataset[0]
              data = dataset[1:]
              x,y,c = [],[],[]
              unique_species = []
              for i in data:
                  x.append(i[0])
                  y.append(i[1])
                  if i[2] == 'Iris-setosa':
                      c.append('r')
                  elif i[2] == 'Iris-versicolor':
                      c.append('b')
                  else: c.append('g')
              for i in range(len(x)):
                  plt.scatter(x[i], y[i], color=c[i])
              plt.show()

          create_plot(converted_iris_data)
```



While my heatmap doesn't actually have the Y axis correct (it is inverted relative to the Y label), it does show something I'd never noticed before. In scatter plots, a point appearring exactly on top of another is hidden. A heatmap gets color values from this type of occurence, so looking at iris this way shows just how many values must be at 0.2 width and just over 1 length.

In conclusion, I will stand by my initial reaction that there is likely not a good way to divide petal width to completely unbundle the 0.2 measurement.