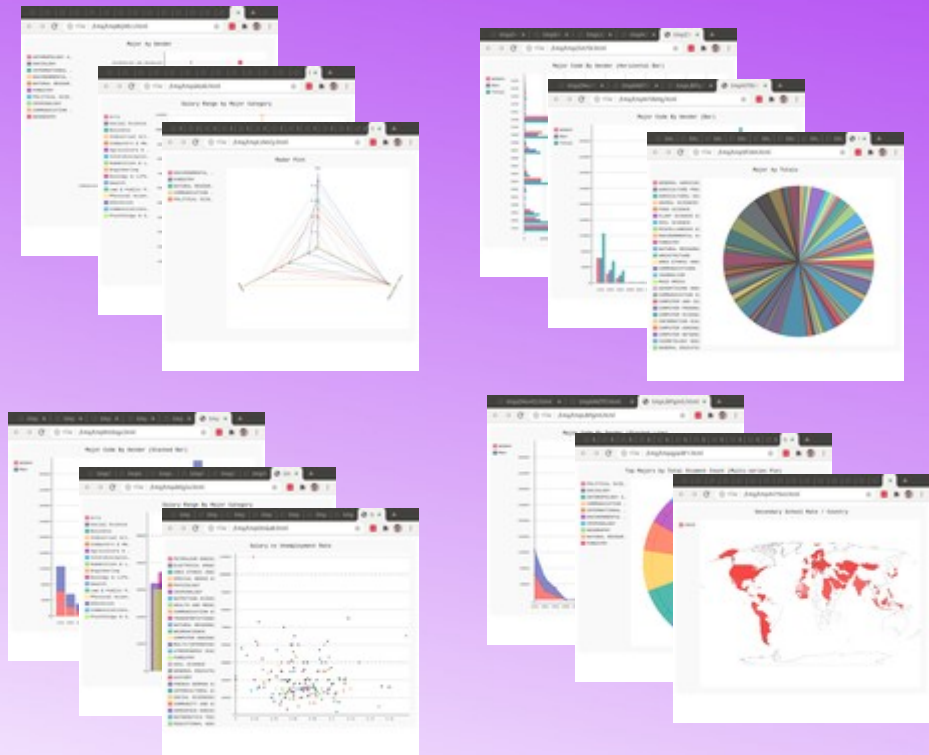


# Data Visualization with Pygal



# Agenda

- What is Pygal?
- Graphing Basics
- Chart/Graph Example Sampler
- Why?
  - Performing data analysis on debugging logs to attain system performance/behaviors has been an emphasis on last couple contracts
  - ‘Visualization’ of even modest data sets gives us a better understanding of the collective

# What is Pygal?

- Python module that creates **interactive** Scalable Vector Graphics (SVG) graphs/charts
- One of many data visualization modules (e.g. Matplotlib, Seaborn, Bokeh, Plotly...)
- In search of honing my data visualization chops, a course in Coursera introduced this module
- Simple, interactive graph/chart, readily integrated in web user interfaces and web pages

# Graphing Basics - Definitions



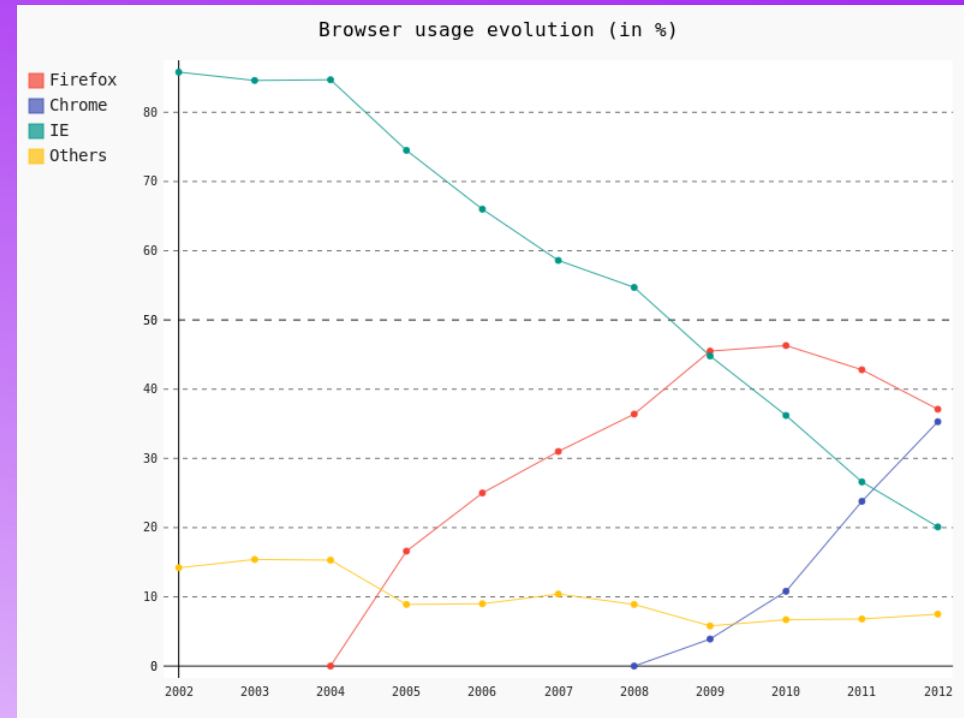
# Graphing Basics - Values

- 2-Dimensional Graphs require X + Y values to be plotted
- Subtle differences in how the values are specified
  - Scatter Point Graphs – Values come in the form  $[(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)]$
  - Most others – Values in the form of  $[y_1, y_2 \dots y_n]$ , x-value inferred by position in list/vector

# Simple Pygal Example

```
$ cat -n example.py
```

```
1  #!/usr/bin/python3
2  import pygal
3  chart = pygal.Line()
4  chart.title = 'Browser usage evolution (in %)'
5  chart.x_labels = map(str, range(2002, 2013))
6  chart.add('Firefox', [None, None, 0, 16.6, 25, 31, 36.4, 45.5, 46.3, 42.8, 37.1])
7  chart.add('Chrome', [None, None, None, None, None, None, 0, 3.9, 10.8, 23.8, 35.3])
8  chart.add('IE', [85.8, 84.6, 84.7, 74.5, 66, 58.6, 54.7, 44.8, 36.2, 26.6, 20.1])
9  chart.add('Others', [14.2, 15.4, 15.3, 8.9, 9, 10.4, 8.9, 5.8, 6.7, 6.8, 7.5])
10 chart.render_in_browser()
```

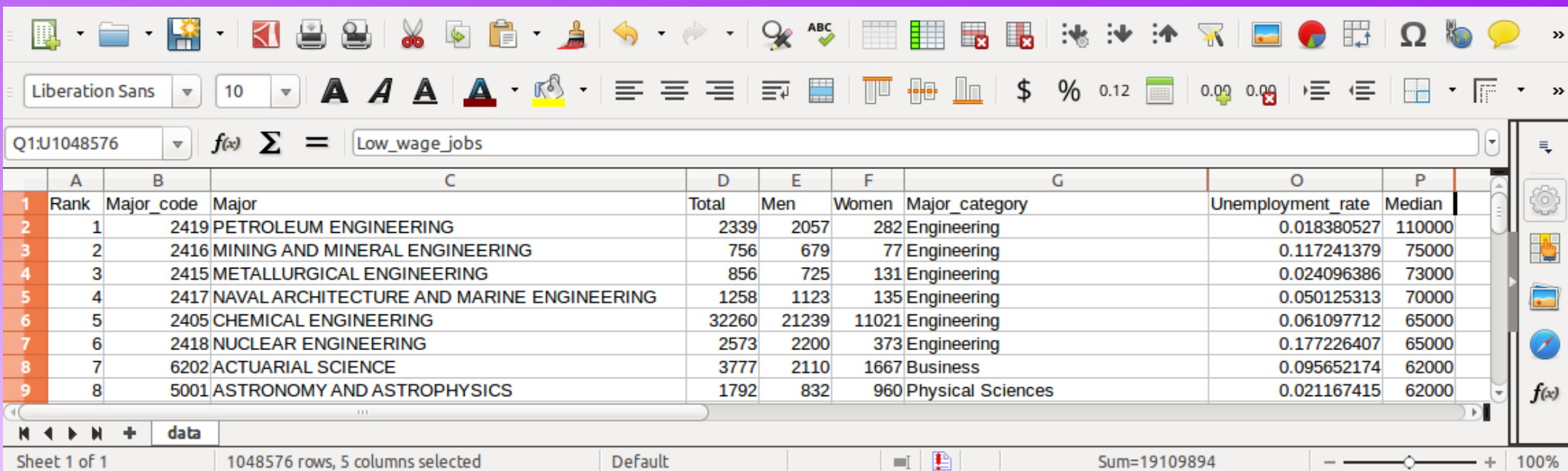




# Example Data Set

- FiveThirtyEight
  - The Economic Guide to Picking A College Major
    - <https://fivethirtyeight.com/features/the-economic-guide-to-picking-a-college-major/>
    - <https://raw.githubusercontent.com/fivethirtyeight/data/master/college-majors/recent-grads.csv>
  - Just an interesting dataset; not an endorsement of the paper
- Wanted a useful dataset that could be used to demonstrate a variety of means of plotting

# Data Overview



	A	B	C	D	E	F	G	O	P
	Rank	Major_code	Major	Total	Men	Women	Major_category	Unemployment_rate	Median
1	1	2419	PETROLEUM ENGINEERING	2339	2057	282	Engineering	0.018380527	110000
2	2	2416	MINING AND MINERAL ENGINEERING	756	679	77	Engineering	0.117241379	75000
3	3	2415	METALLURGICAL ENGINEERING	856	725	131	Engineering	0.024096386	73000
4	4	2417	NAVAL ARCHITECTURE AND MARINE ENGINEERING	1258	1123	135	Engineering	0.050125313	70000
5	5	2405	CHEMICAL ENGINEERING	32260	21239	11021	Engineering	0.061097712	65000
6	6	2418	NUCLEAR ENGINEERING	2573	2200	373	Engineering	0.177226407	65000
7	7	6202	ACTUARIAL SCIENCE	3777	2110	1667	Business	0.095652174	62000
8	8	5001	ASTRONOMY AND ASTROPHYSICS	1792	832	960	Physical Sciences	0.021167415	62000

- 174 Rows, 21 Columns of data organized by university major
- Focus out attention on 9 key columns in our examples



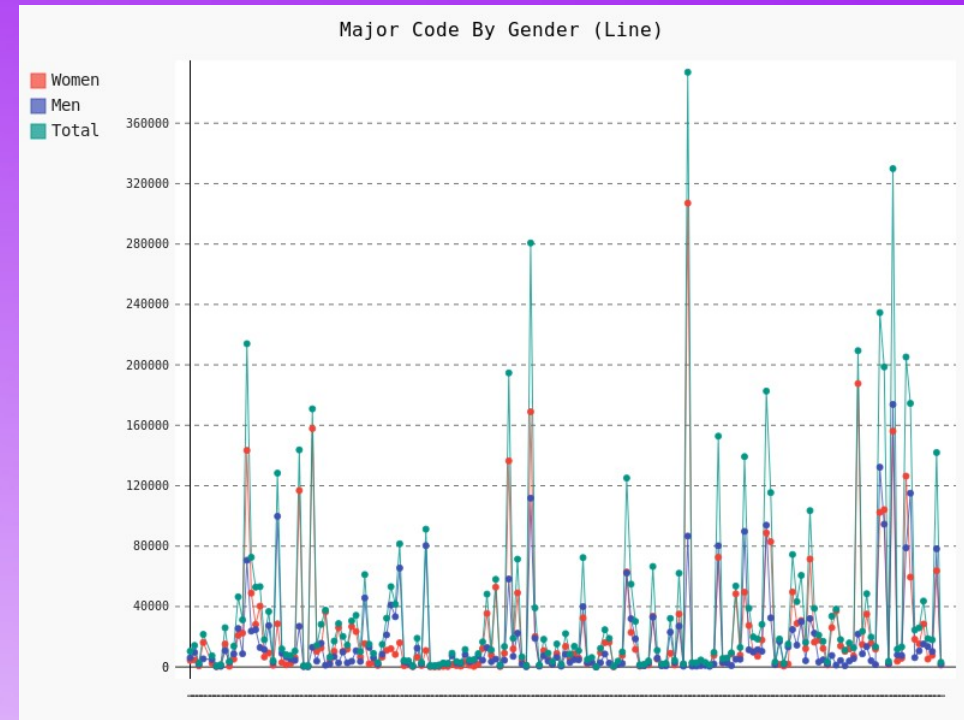
# CSV File Reader

- `def readCsvAsDict(fileName, keyField, separator=',', quote='"'):`
- `data = readCsvAsDict('data.csv',keyField='Major_code')`
  - Returns dictionary, keyed by 'Major\_code' column value, value is dictionary of all column field names
    - {"1301",
    - {
    - "Major":"ENVIRONMENTAL SCIENCE",
    - "Men":"10787",
    - "Unemployment\_rate":"0.078584681",
    - "Major\_code":"1301",
    - "Median":"35600",
    - "Rank":"93",
    - "Major\_category":"Biology & Life Science",
    - "Women":"15178"
    - }
    - ...
    - }
- In retrospect, I likely would have used Pandas csv reader

# Line

```
$ cat example.py

1  #!/usr/bin/python3
2  import pygal
3  import csv
4  data=readCsvAsDict('data.csv','Major_code')
5  Fields=['Women','Men','Total']
6  plotData=dict()
7  for key in Fields:
8      D=[(k,v[key]) for (k,v) in sorted(data.items())]
9      L=[int(el[1]) if el[1].isdigit() else None for el in D]
10     xLabel=[el[0] for el in D]
11     plotData[key]=L
12  chart=pygal.Line()
13  chart.title='Major Code By Gender (Line)'
14  for key in Fields:
15      chart.add(key,plotData[key])
16  chart.x_labels = xLabel
17  chart.render_in_browser()
```

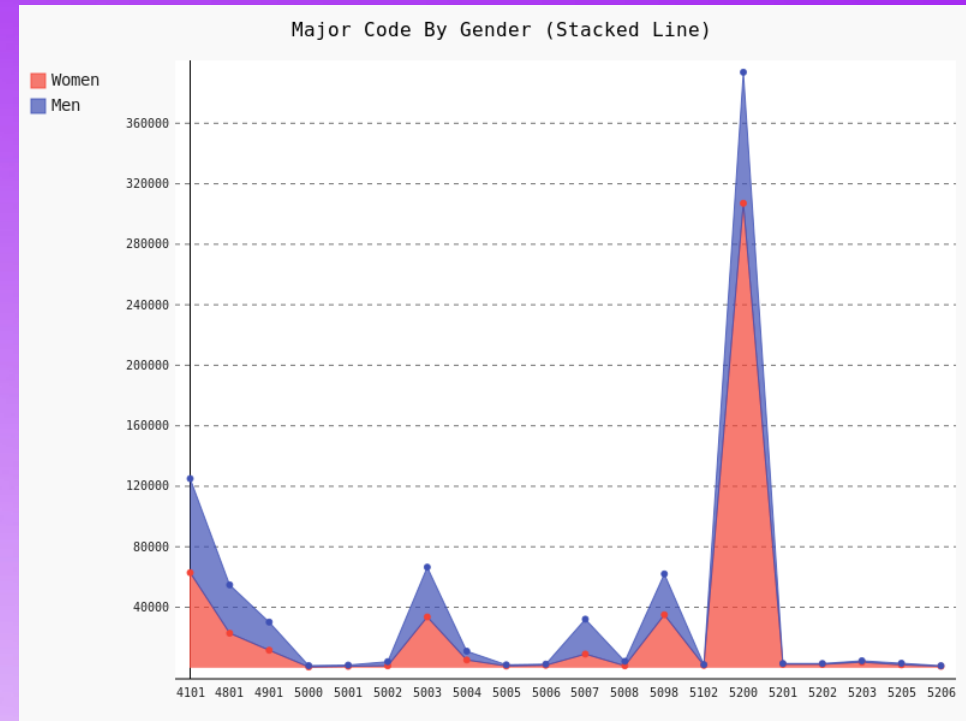


```
chart.add('Women', [77, 282, ...])
chart.add('Men', [679, 2057, ...])
chart.add('Total', [756, 2339, ...])
```

```
chart.x_labels=['2416', '2419', ...]
```

# Stacked Line

```
1  #!/usr/bin/python3
2  import pygal
3  import csv
4  data=readCsvAsDict('data.csv','Major_code')
5  Fields=['Women','Men']
6  plotData=dict()
7  for key in Fields:
8      D=[(k,v[key]) for (k,v) in sorted(data.items())[100:120]]
9      L=[int(el[1]) if el[1].isdigit() else None for el in D]
10     xLabel=[el[0] for el in D]
11     plotData[key]=L
12
13 chart=pygal.StackedLine(fill=True)
14 chart.title='Major Code By Gender (Stacked Line)'
15 for key in Fields:
16     chart.add(key,plotData[key])
17
18 chart.x_labels = xLabel
19 chart.render_in_browser()
```



`chart.add('Men', [62181, 31967, ...])`

# Bar

```
1  #!/usr/bin/python3
2  import pygal
3  import csv
4  data=readCsvAsDict('data.csv','Major_code')
5  Fields=['Women','Men','Total']
6  plotData=dict()
7  for key in Fields:
8      D=[(k,v[key]) for (k,v) in sorted(data.items())[100:120]]
9      L=([int(el[1]) if el[1].isdigit() else None for el in D])
10     xLabel=([el[0] for el in D])
11     plotData[key]=L
12  chart=pygal.Bar()
13  chart.title='Major Code By Gender (Bar)'
14  for key in Fields:
15      chart.add(key,plotData[key])
16  chart.x_labels = xLabel
17  chart.render_in_browser()
```

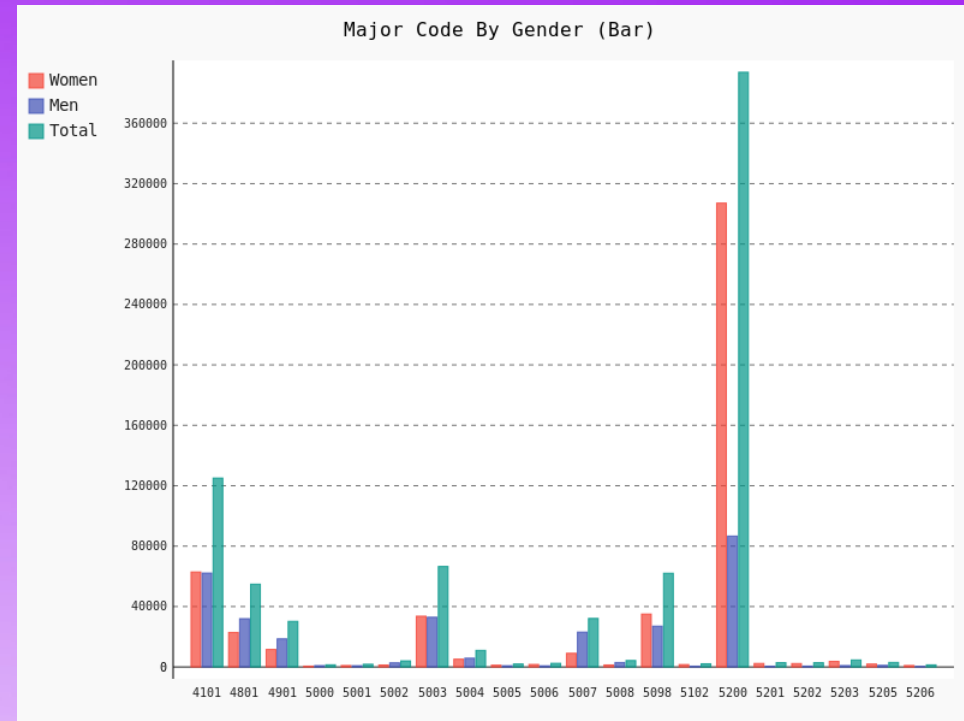
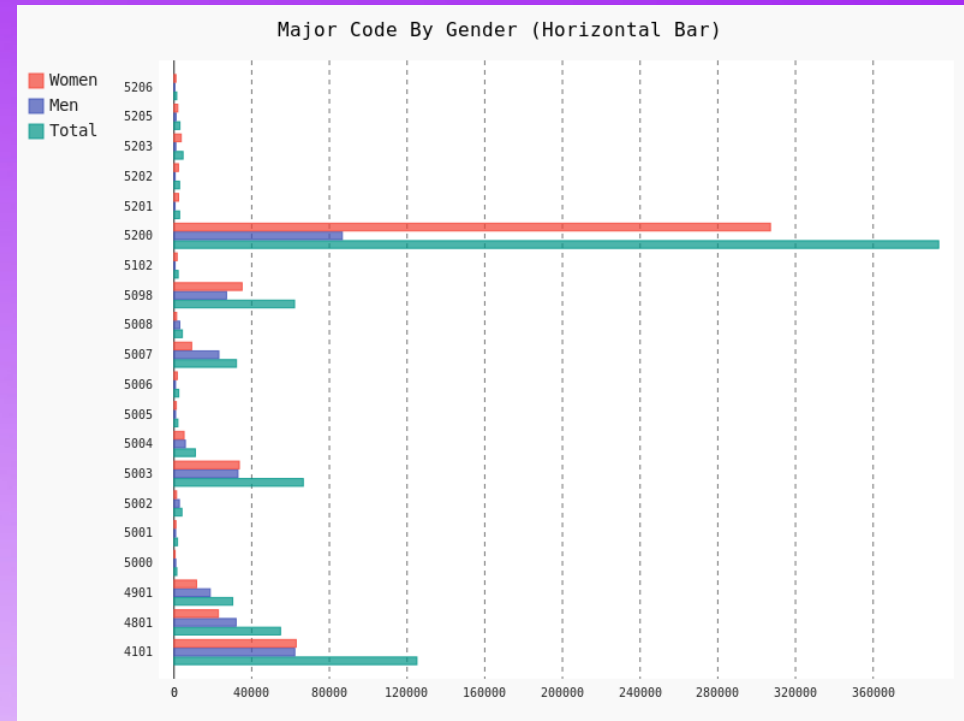


chart.add('Men', [62181, 31967, ...])

# Horizontal Bar

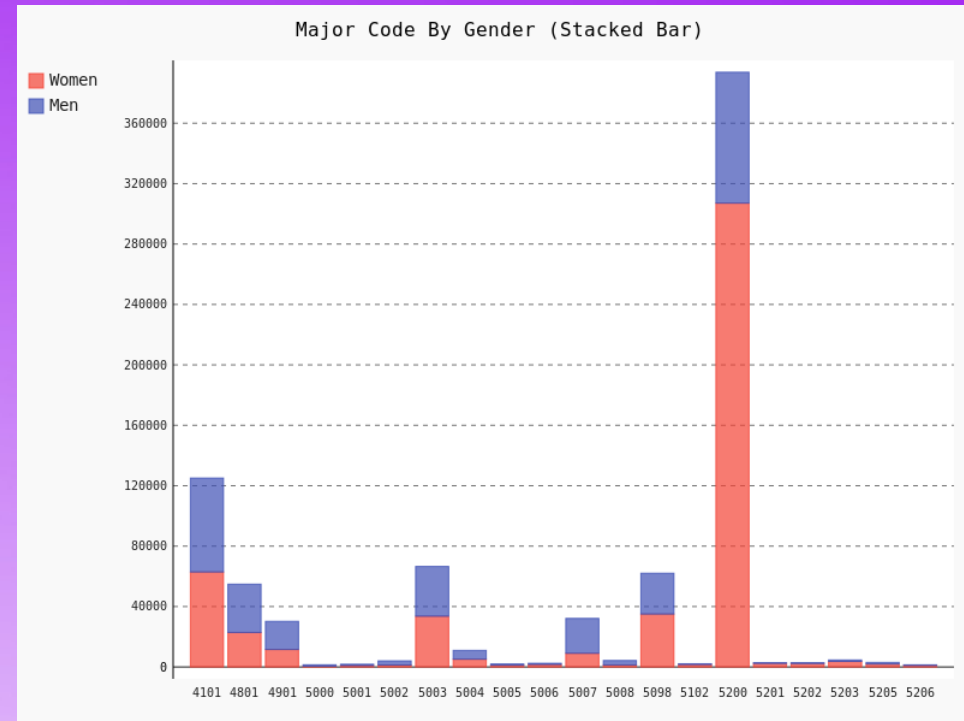
```
1  #!/usr/bin/python3
2  import pygal
3  import csv
4  data=readCsvAsDict('data.csv','Major_code')
5  Fields=['Women','Men','Total']
6  plotData=dict()
7  for key in Fields:
8      D=[(k,v[key]) for (k,v) in sorted(data.items())[100:120]]
9      L=([int(el[1]) if el[1].isdigit() else None for el in D])
10     xLabel=([el[0] for el in D])
11     plotData[key]=L
12  chart=pygal.HorizontalBar()
13  chart.title='Major Code By Gender (Horizontal Bar)'
14  for key in Fields:
15      chart.add(key,plotData[key])
16  chart.x_labels = xLabel
17  chart.render_in_browser()
```



`chart.add('Men', [62181, 31967, ...])`

# Stacked Bar

```
1  #!/usr/bin/python3
2  import pygal
3  import csv
4  data=readCsvAsDict('data.csv','Major_code')
5  Fields=['Women','Men']
6  plotData=dict()
7  for key in Fields:
8      D=[(k,v[key]) for (k,v) in sorted(data.items())[100:120]]
9      L=([int(el[1]) if el[1].isdigit() else None for el in D])
10     xLabel=[el[0] for el in D])
11     plotData[key]=L
12  chart=pygal.StackedBar()
13  chart.title='Major Code By Gender (Stacked Bar)'
14  for key in Fields:
15      chart.add(key,plotData[key])
16  chart.x_labels = xLabel
17  chart.render_in_browser()
```

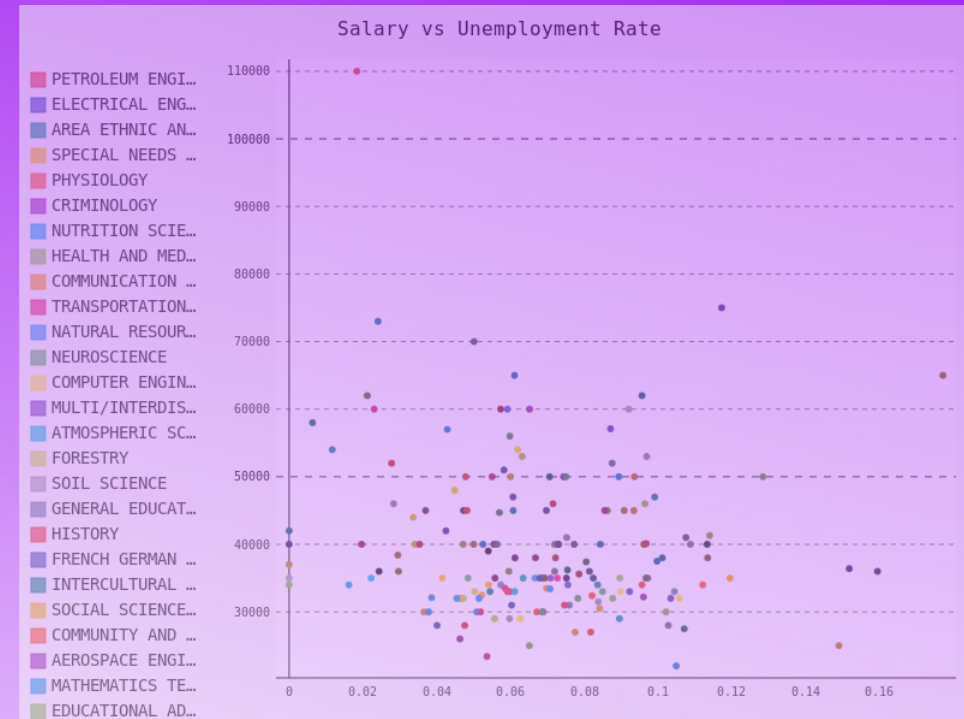


`chart.add('Men', [62181, 31967, ...])`



# XY

```
1  #!/usr/bin/python3
2  import pygal
3  import csv
4  data=readCsvAsDict('data.csv', 'Rank')
5  chart = pygal.XY()
6  chart.title='Salary vs Unemployment Rate'
7  for (k,v) in sorted(data.items()):
8      chart.add(v['Major'], [(float(v['Unemployment_rate']),int(v['Median']))])
9  chart.render_in_browser()
```



```
chart.add('FINE ARTS',
[(0.084186296, 30500)]
```

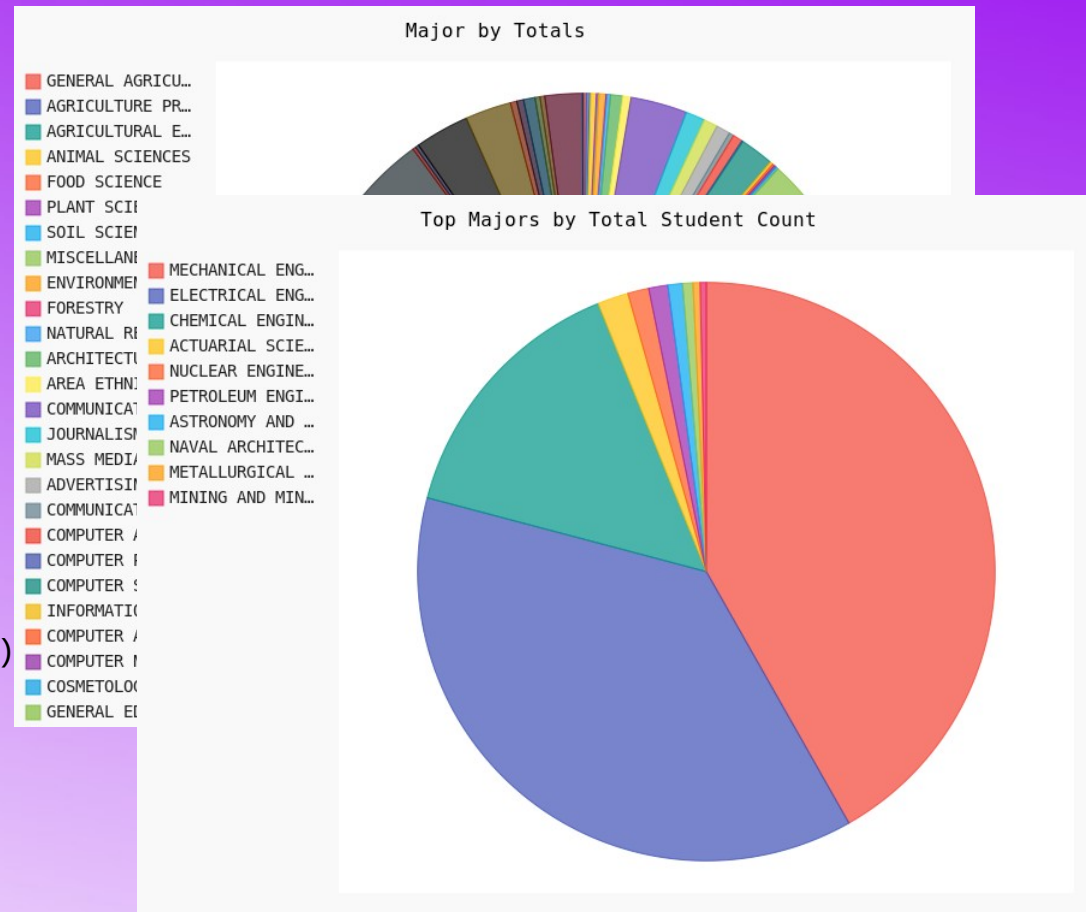
# Pie

```

1  #!/usr/bin/python3
2  import pygal
3  import csv
4  data=readCsvAsDict('data.csv','Major_code')
5  chart = pygal.Pie()
6  chart.title='Major by Totals'
7  for (k,v) in sorted(data.items()):
8      val = int(v['Total']) if v['Total'].isdigit() else None
9      chart.add(v['Major'],val)
10 chart.render_in_browser()

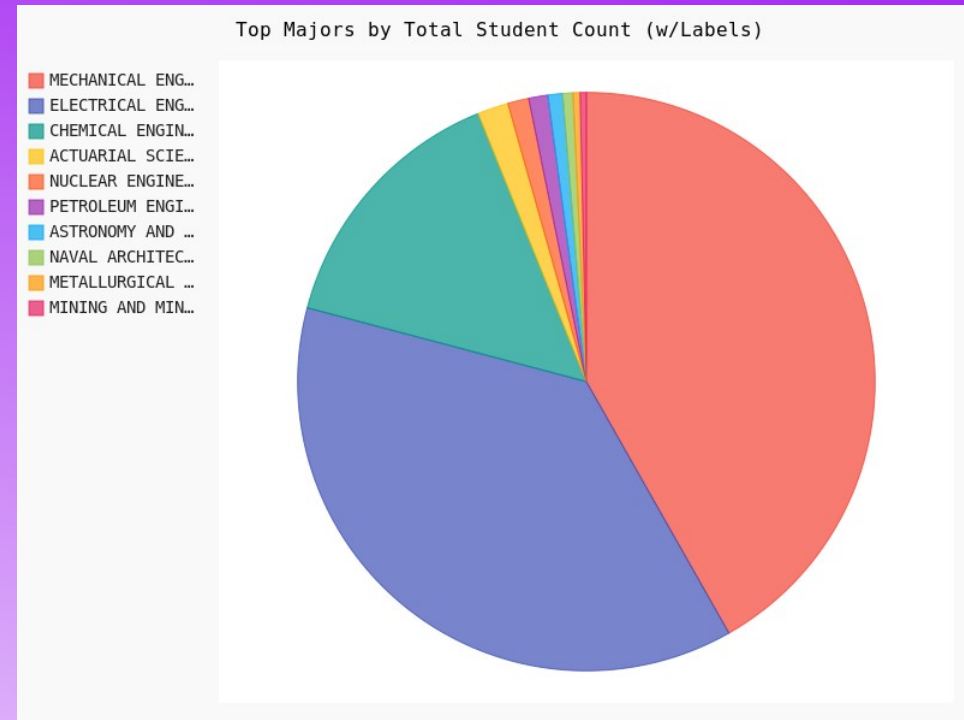
```

chart.add('ARCHITECTURE',46420)



# Pie w/Labels

```
1  #!/usr/bin/python3
2  import pygal
3  import csv
4  data=readCsvAsDict('data.csv','Major_code')
5  chart = pygal.Pie()
6  chart.title='Top Majors by Total Student Count (w/Labels)'
7  L=[(int(v['Total']),v['Major']) for (k,v) in data.items() if v['Total'].isdigit()]
8  L=L[0:10]
9  N=sum([v for (v,k) in L])
10 for (t,k) in sorted(L,reverse=True):
11     chart.add(k,[{'value': t, 'label': "%0.2f%%"%(float(100*t)/N)}])
12 chart.render_in_browser()
```



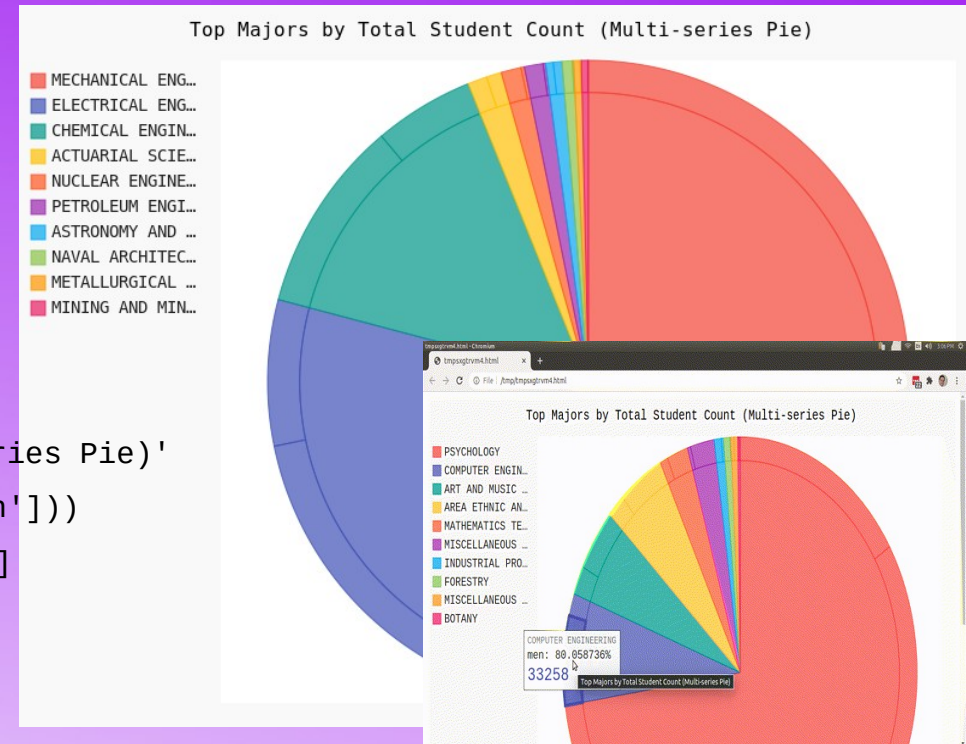
```
chart.add('ARCHITECTURE',
[{'value':46420,
'label':"0.06"}])
```

# Multi-Series Pie

```

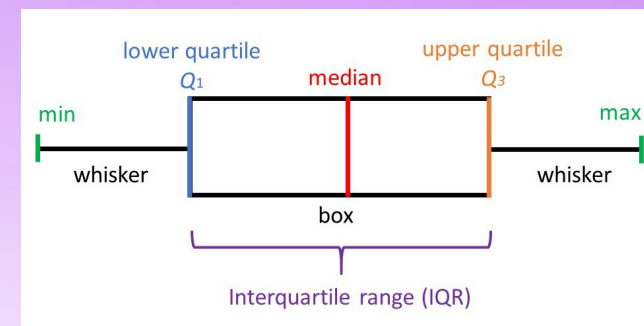
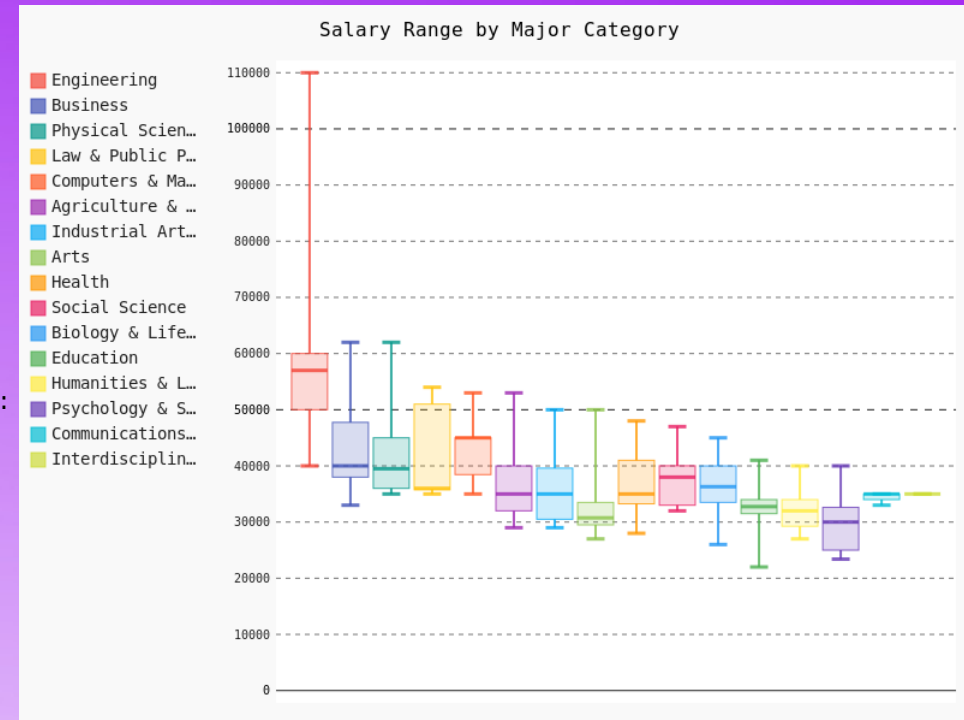
1  #!/usr/bin/python3
2  import pygal
3  import csv
4  data=readCsvAsDict('data.csv','Major_code')
5  chart = pygal.Pie()
6  chart.title='Top Majors by Total Student Count (Multi-series Pie)'
7  L=[(int(v['Total']),v['Major'],int(v['Men']),int(v['Women']))]
8  for (k,v) in data.items() if v['Total'].isdigit()
9  L=L[0:10]
10 N=sum([v for (v,t,m,w) in L])
11 for (t,k,m,w) in sorted(L,reverse=True):
12     chart.add(k, [{'value':m, 'label': 'men: %02f%%'%(float(100*m)/t)}, {'value':w, 'label': 'women:%02f%%'%(float(100*w)/t)}])
13
14 chart.add('ZOOLOGY' [{'value': 3050, 'label': 'men: 36.27%'},
15                      {'value': 5359, 'label': 'women:63.72%'}])
16 chart.render_in_browser()

```



# Box

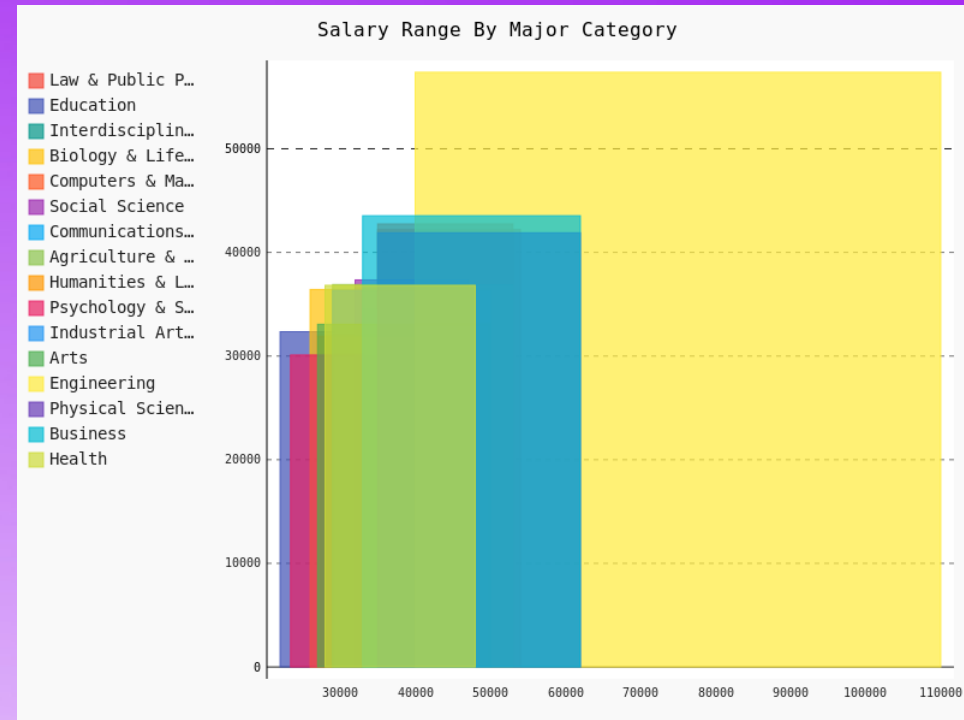
```
1  #!/usr/bin/python3
2  import pygal
3  import csv
4  data=readCsvAsDict('data.csv', 'Major_code')
5  plotData=dict()
6  for val in [v for (k,v) in data.items() if v['Median'].isdigit()]:
7      category=val['Major_category']
8      try:
9          plotData[category].append(int(val['Median']))
10     except(KeyError):
11         plotData[category]=list()
12         plotData[category].append(int(val['Median']))
13  chart = pygal.Box()
14  chart.title = 'Salary Range by Major Category'
15  for (k,v) in plotData.items():
16      chart.add(k,v)
17  chart.render_in_browser()
```





# Histogram

```
1  #!/usr/bin/python3
2  import pygal
3  import csv
4  data=readCsvAsDict('data.csv','Major_code')
5  plotData=dict()
6  for (key,val) in data.items():
7      try:
8          plotData[val['Major_category']].append(int(val['Median']))
9      except:
10         plotData[val['Major_category']]=list()
11         plotData[val['Major_category']].append(int(val['Median']))
12 categories=[val['Major_category'] for (k,val) in data.items()]
13 chart = pygal.Histogram()
14 chart.title='Salary Range By Major Category'
15 for k in set(categories):
16     x0=min(plotData[k])
17     x1=max(plotData[k])
18     y=sum(plotData[k])/float(len(plotData[k]))
19     chart.add(k,[(y,x0,x1)])
20 chart.render_in_browser()
```

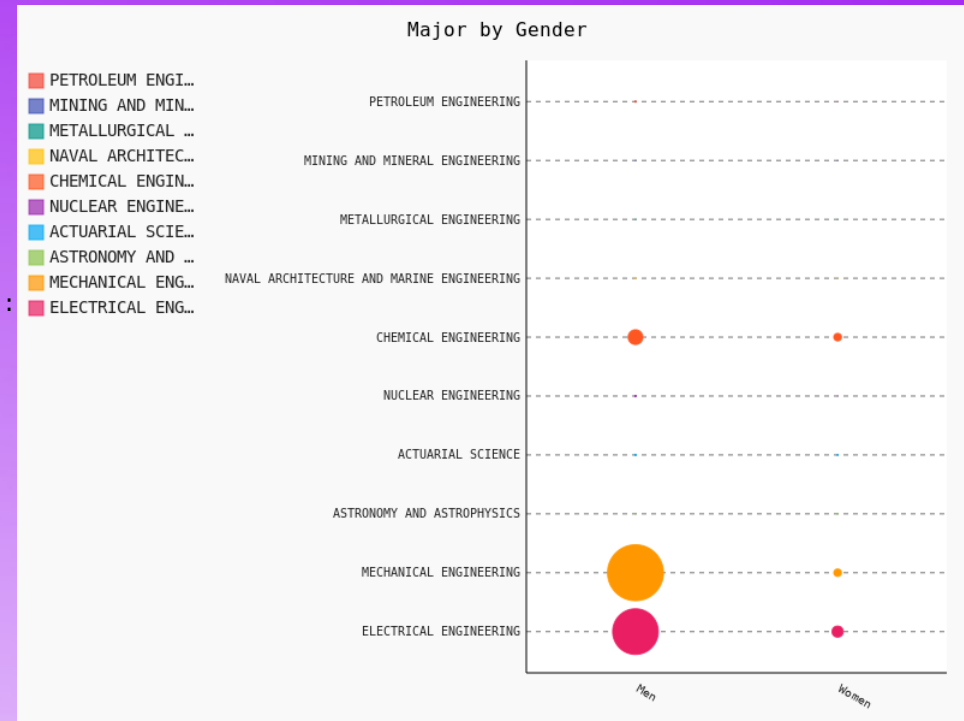


```
chart.add('Arts',  
[(33062.5,27000,50000),...])
```



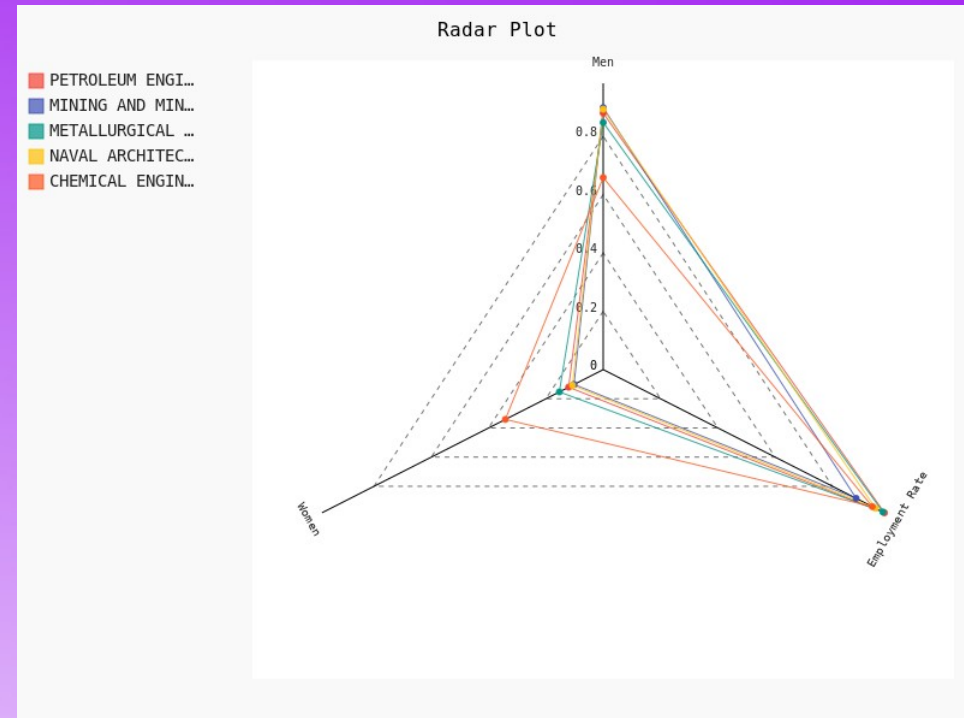
# Dot

```
1  #!/usr/bin/python3
2  import pygal
3  import csv
4  data=readCsvAsDict('data.csv','Major_code')
5  plotData=dict()
6  for val in [v for (k,v) in data.items() if v['Total'].isdigit()][0:10]:
7      category=val['Major']
8      try:
9          plotData[category].append(int(val['Men']))
10         plotData[category].append(int(val['Women']))
11     except(KeyError):
12         plotData[category]=list()
13         plotData[category].append(int(val['Men']))
14         plotData[category].append(int(val['Women']))
15  chart = pygal.Dot(x_label_rotation=30)
16  chart.title = 'Major by Gender'
17  chart.x_labels = ['Men', 'Women']
18  for (k,v) in plotData.items():
19      chart.add(k, v)
20  chart.render_in_browser()
```



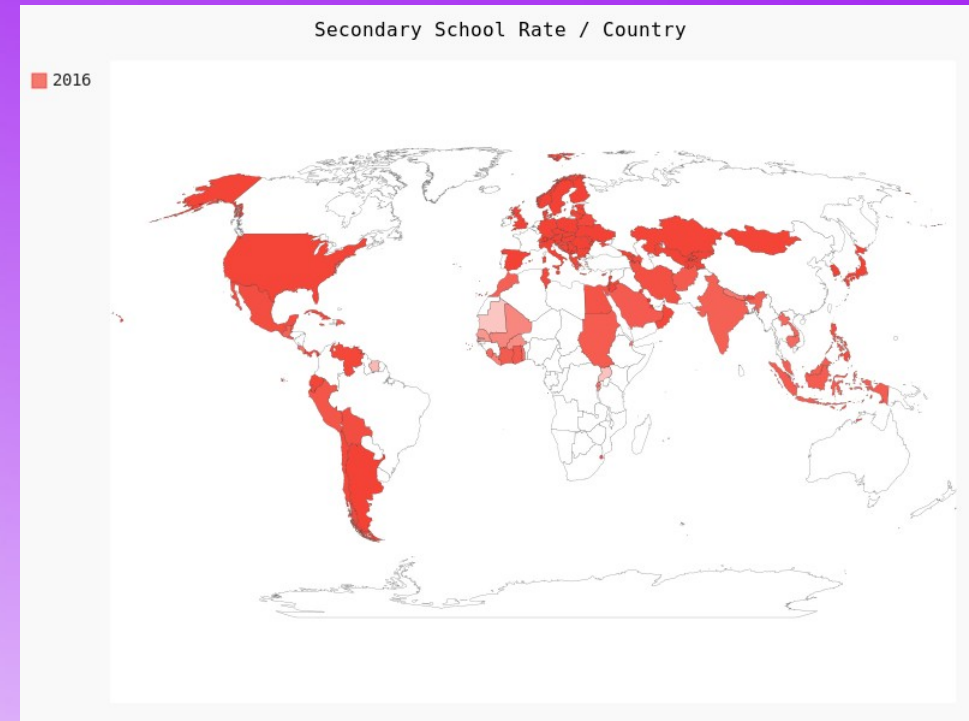
# Radar

```
1  #!/usr/bin/python3
2  import pygal
3  import csv
4  data=readCsvAsDict('data.csv','Major_code')
5  chart = pygal.Radar()
6  chart.title = 'Radar Plot'
7  chart.x_labels=['Men','Women','Employment Rate']
8  for val in [v for (k,v) in data.items() if v['Total'].isdigit()][0:5]:
9      L=[]
10     L.append(float(val['Men'])/float(val['Total']));
11     L.append(float(val['Women'])/float(val['Total']));
12     L.append(1.0-float(val['Unemployment_rate']));
13     chart.add(val['Major'],L)
14  chart.render_in_browser()
```



# World Map

```
1  #!/usr/bin/python3
2  import pygal
3  import csv
4  def convertCountryCodeToPygal(countryCode):
5      convertCountryCodeToPygal.data=readCsvAsDict('WDICountry.csv','Country Code')
6      return convertCountryCodeToPygal.data[countryCode]['2-alpha code'].lower()
7
8  data=readCsvAsDict('school.csv','Country Code')
9  chart = pygal.maps.world.World()
10 chart.title = 'Secondary School Rate / Country'
11 year=2016
12 plotData=dict()
13 for (k,v) in data.items():
14     try:
15         plotData[convertCountryCodeToPygal(k)]=float(v[str(year)])
16     except:
17         pass
18 chart.add(str(year),plotData)
19 chart.render_to_png('./example10.png')
20 chart.render_in_browser()
```



# References

- <http://www.pygal.org/>
  - Official Site
- <https://github.com/fivethirtyeight/data/tree/master/college-majors/>
- <https://datacatalog.worldbank.org/dataset/world-development-indicators/>

# Contact Info

- Slides:
  - <https://github.com/fsk-software/pub/>
- Blog: <http://dragonquest64.blogspot.com>
- Slack: [pymntos.slack.com](https://pymntos.slack.com) [lipeltgm](#)