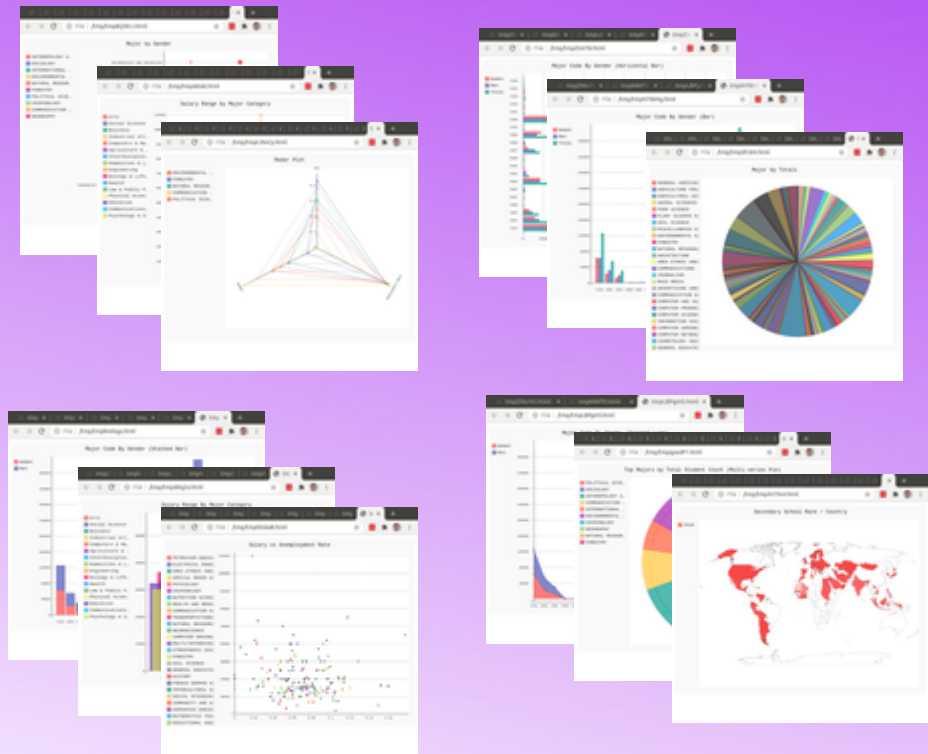# Data Visualization with Pygal

# 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, …)

- 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
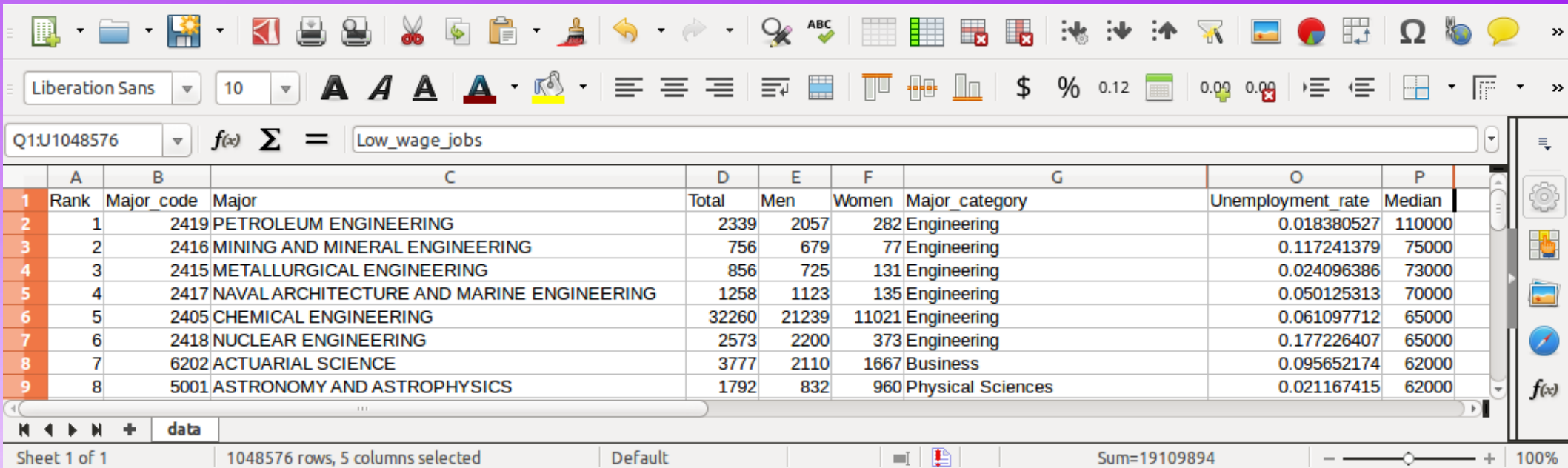
# Agenda

- What is Pygal?

- Chart/Graph Example Sampler

- 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

# Example Data Set

- https://raw.githubusercontent.com/fivethirtyeight/data/master/college-majors/recent-grads.csv

- FiveThirtyEight

    – The Economic Guide to Picking A College Major
        - https://fivethirtyeight.com/features/the-economic-guide-to-picking-a-college-major/

    – Just an interesting dataset; not and endorsement of the paper

    – Wanted a useful dataset that could be used to demonstrate a variety of means of plotting

# Data Overview



- 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"
    - }
    - ...
    - }

FSK INC. CONSULTING

# Simple Pygal Example

```
$ cat -n example.py

1    #!/usr/bin/python3

2    import pygal

3    import csv

4    chart = pygal.Line()

5    chart.title = 'Browser usage evolution (in %)'

6    chart.x_labels = map(str, range(2002, 2013))

7    chart.add('Firefox', [None, None,    0, 16.6,   25,   31, 36.4, 45.5, 46.3, 42.8, 37.1])

8    chart.add('Chrome',  [None, None, None, None, None, None,    0,  3.9, 10.8, 23.8, 35.3])

9    chart.add('IE',      [85.8, 84.6, 84.7, 74.5,   66, 58.6, 54.7, 44.8, 36.2, 26.6, 20.1])

10   chart.add('Others',  [14.2, 15.4, 15.3,  8.9,    9, 10.4,  8.9,  5.8,  6.7,  6.8,  7.5])

11   chart.render_in_browser()
```

# 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()
```

# Stacked Line

```python
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.StackedLine(fill=True)

13   chart.title='Major Code By Gender (Stacked Line)'

14   for key in Fields:

15     chart.add(key,plotData[key])

16   chart.x_labels = xLabel

17   chart.render_in_browser()
```

# Bar

```python
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()
```
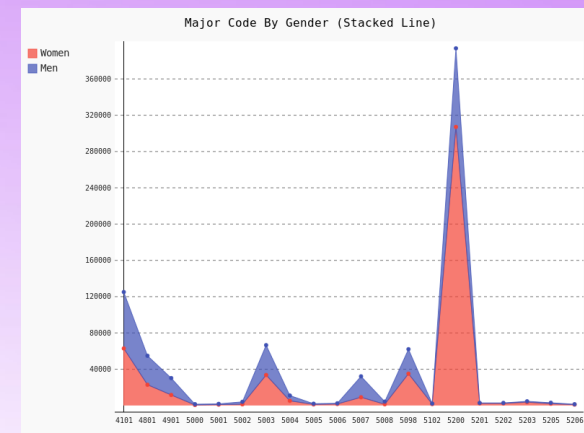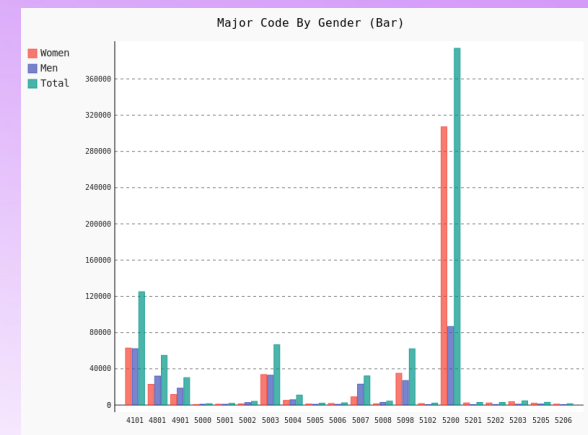
# Horizontal Bar

```python
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()
```
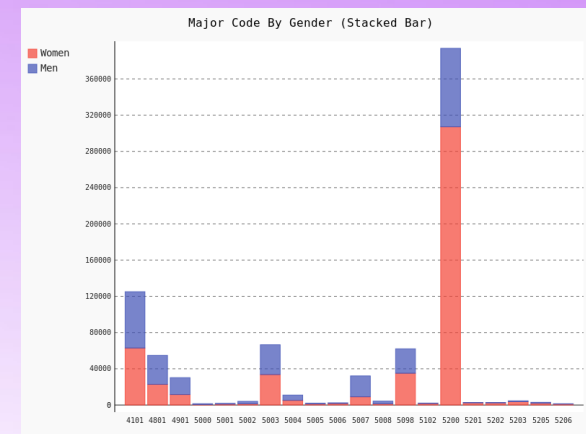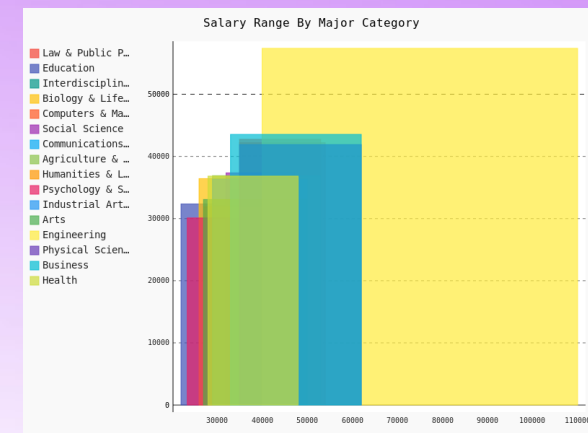
# Stacked Bar

```python
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()
```
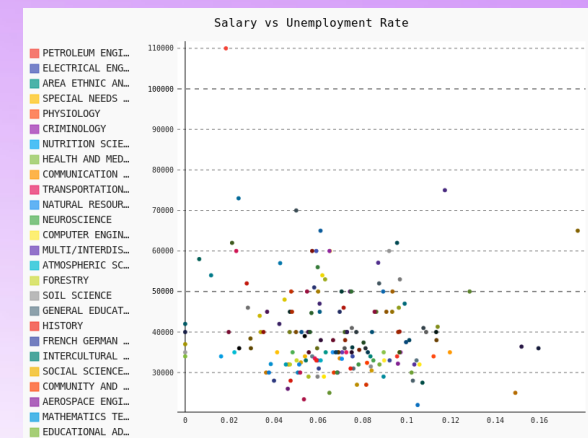


Major Code By Gender (Stacked Bar)

# Histogram

```python
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()
```





FSK INC.
CONSULTING

# XY

```python
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()
```
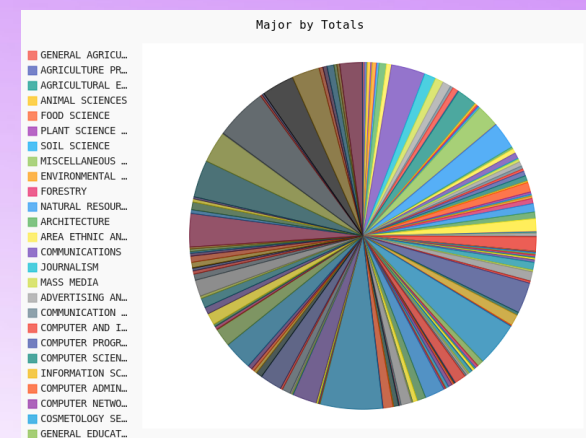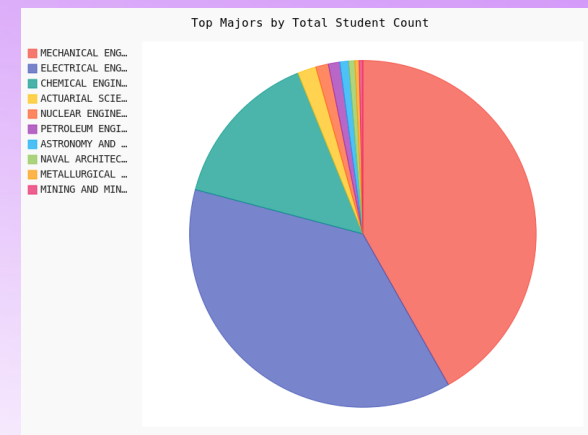
# Pie

```
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()
```
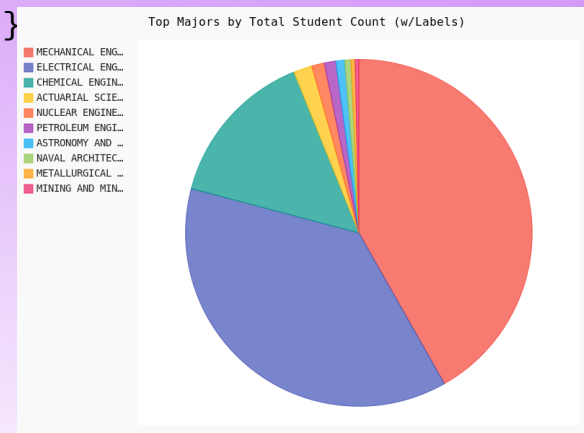


Major by Totals

# Pie

```python
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'

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,t)

12   chart.render_in_browser()
```
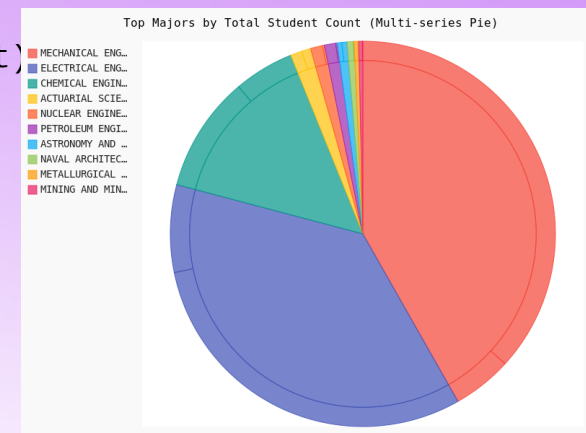


Top Majors by Total Student Count

MECHANICAL ENG...
ELECTRICAL ENG...
CHEMICAL ENGIN...
ACTUARIAL SCIE...
NUCLEAR ENGINE...
PETROLEUM ENGI...
ASTRONOMY AND ...
NAVAL ARCHITEC...
METALLURGICAL ...
MINING AND MIN...

# Pie w/Labels

```python
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()
```
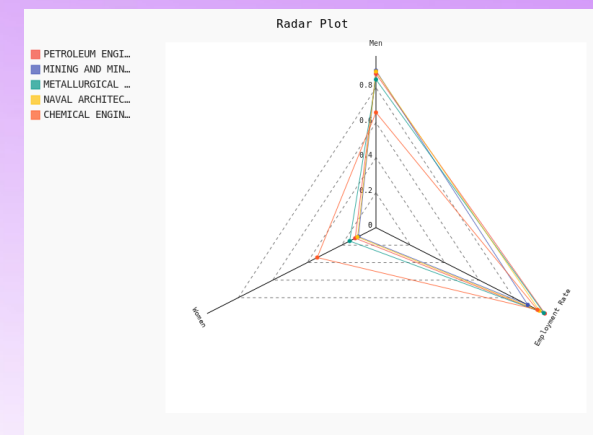
# Multi-Series Pie

```python
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'])) for (k,v) in data.items() if
     v['Total'].isdigit()]

8    L=L[0:10]

9    N=sum([v for (v,t,m,w) in L])

10   for (t,k,m,w) in sorted(L,reverse=True):

11       chart.add(k,[{'value':m,'label':'men: %02f%%'%(float(100*m)/t)
     %02f%%'%(float(100*w)/t)}])

12   chart.render_in_browser()
```
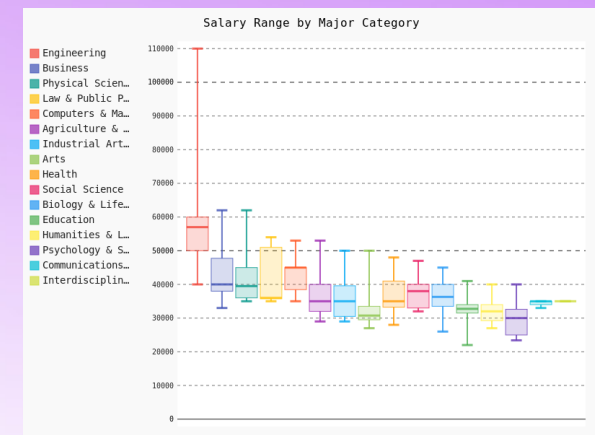
# Radar

```python
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()
```
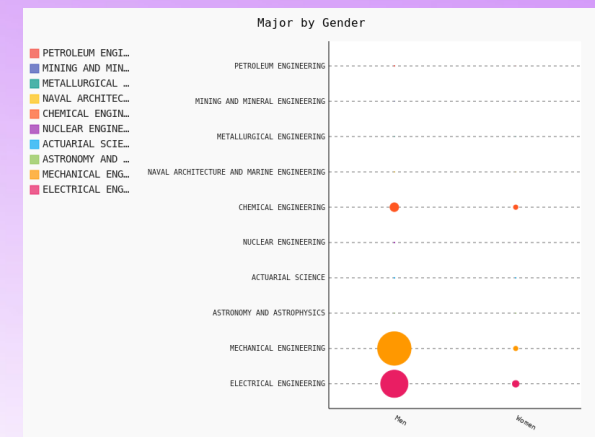


Radar Plot

# Box

```python
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()
```
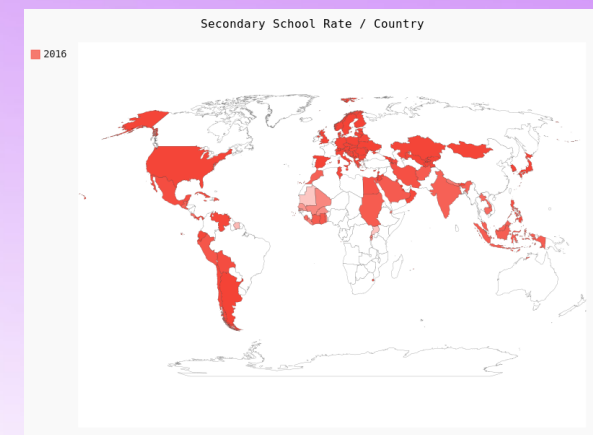
# Dot

```python
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()
```

# World Map

```python
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()
```


Secondary School Rate / Country

# 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 lipeltgm

**FSK** INC.
CONSULTING