

Christopher McClaskey
CSCI 4502 - Homework 1
SID: 810937948
christopher.mcclaskey@colorado.edu

Honor Code Pledge

On my honor, as a University of Colorado at Boulder student, I have neither given nor received unauthorized assistance on this work.

Question 1

Resource 1:

(a) <http://archive.ics.uci.edu/ml/datasets/Communities+and+Crime>

(b) Aggregation of socio-economic data, law-enforcement data, and crime data, by region. It has 1,994 objects, and 128 attributes (ie state, county, pctHSGrad, etc).

(c) Knowledge that can be mined would be patterns and associations between socio-economic status, region, and crime.

(d) This knowledge would be useful in mapping areas that need attention, whether that attention be better a better education system, economic stimulus, or law-enforcement improvement. This map could help in the decision to divvy out resources whether it be money or man-power.

Resource 2:

(a) <http://archive.ics.uci.edu/ml/datasets/Census+Income>

(b) Selection of tax and income data, and includes some socio-economic data. It has 48,842 objects, and 14 attributes (age, workclass, sex, etc).

(c) Knowledge that can be mined would be patterns and associations among socio-economic status and income.

(d) This knowledge would be useful in economic models. It could help judge the state of the economy, the employment situation, and see if there are any problems that could be fixed. Perhaps low income areas all have something in common that could be addressed.

Question 2 (extra credit)

(a) KDD'11. "Multiple Domain User Personalization" by Amr Ahmed, Yucheng Low, Mohamed Aly, Vanja Josifovski, Alexander J. Smola

<http://dl.acm.org/citation.cfm?id=2020433&CFID=176070325&CFTOKEN=86241249>

(b) The problem being addressed is the "Cold Start Problem". This problem is specific to websites that try to make recommendations to the user. Typically these websites require some user input before they can do anything. For example, there isn't a way for Netflix to know what kind of movies a user likes without asking the user or looking at their previously watched movies. The idea is to make recommendations from the start, without user input, but this is very difficult.

(c) The authors' proposed solution is to combine information from across multiple domains. A user's google searches or facebook likes can provide insight into a user's preferences elsewhere. The authors do this by constructing a Bayesian model which integrates multiple different recommendation systems.

(d) The solution is presented and evaluated with some statistical theory, implementation of the statistics with code, and testing. The test results were presented with metrics.

Question 3

Source Code and Output (Ruby):

```
1. # Christopher McClaskey
2. # CSCI 4502
3. # Homework 1
4. # Question 3 parts (a) and (b)
5.
6. require 'csv'
7.
8. wind = Array.new
9. headers = true
10. CSV.foreach("forestfires.csv") do |row|
11.     #skip the column name row
12.     if headers
13.         headers = false
14.         next
15.     end
16.     #add the wind value to my array
17.     wind.push(row[10].to_f)
18. end
19.
20. #Calculating Mean
21. mean = wind.inject { |sum, i| sum + i } / wind.size
22.
23. #Calculating Standard Deviation
24. std_dev = Math.sqrt( wind.inject { |sum, i| sum + (i-mean)**2 } /
    wind.size )
25.
26. #Calculating median and quartiles
27. wind = wind.sort
28. median_id = wind.size/2
29. lower = wind[0..median_id-1]
30. upper = wind[median_id+1..wind.size-1]
31. q1 = lower[lower.size/2]
32. median = wind[median_id]
33. q3 = upper[upper.size/2]
34. iqr = q3-q1
35.
36. puts "Min:          #{wind.min}"
37. puts "Max:          #{wind.max}"
38. puts "Mean:         #{mean}"
```

```

39. puts "Std Dev:      #{std_dev}"
40. puts "Q1:           #{q1}"
41. puts "Median:       #{median}"
42. puts "Q3:           #{q3}"
43. puts "IQR:          #{iqr}"
44.
45. # -OUTPUT-
46. # Min:             0.4
47. # Max:             9.4
48. # Mean:            4.017601547388782
49. # Std Dev:         1.789651406176351
50. # Q1:              2.7
51. # Median:           4.0
52. # Q3:              4.9
53. # IQR:             2.2

```

Scatter Plot of Attributes 10 (RH) and 11 (Wind):

