

Internet Measurement and Data Analysis (6)

Kenjiro Cho

2016-05-16

review of previous class

Class 5 Diversity and complexity (5/9)

- ▶ Long tail
- ▶ Web access and content distribution
- ▶ Power-law and complex systems
- ▶ exercise: power-law analysis

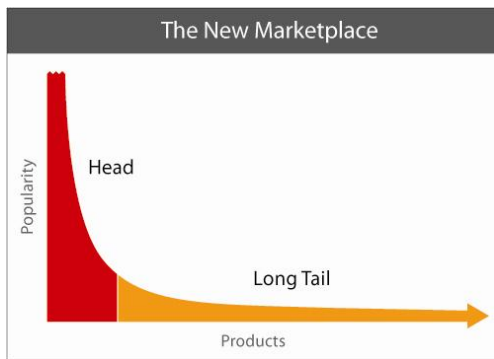
today's topics

Class 6 Correlation

- ▶ Online recommendation systems
- ▶ Distance
- ▶ Correlation coefficient
- ▶ exercise: correlation analysis

online recommender systems

- ▶ finding potential needs for long-tail users at EC sites
 - ▶ by recommending products which fit each user's taste
- ▶ widely used as the cost goes down by recomender package software



source: <http://longtail.com/>

recommender systems

- ▶ from user online behavior, infer useful information for users automatically
- ▶ EC sites: recommend products automatically from purchase or view records
- ▶ other applications: music, movies, search engine, etc

different approaches for database structure

- ▶ item based: compile data for each item
- ▶ user based: compile data for each user
- ▶ most systems combine both

prediction methods of recommender systems

- ▶ content based:
 - ▶ recommend items similar to the items the user used in the past
 - ▶ (manual) classifications of items
 - ▶ clustering items by machine learning methods
 - ▶ building rules from know-how
 - ▶ tend to recommend items in the same group, less surprising
- ▶ collaborative filtering: employed by amazon and others
 - ▶ e.g., "users who bought X also bought Y"
 - ▶ compute similarities among users from their online activities
 - ▶ recommend items bought by similar users
 - ▶ main feature: it does not use the information about items
 - ▶ could lead to surprising findings for user (serendipity)
- ▶ naive bayesian filter: often used for spam filtering
 - ▶ machine-learning technique to compute probabilities from a large number of item and user attributes

recent advances in targeted advertising

- ▶ targeted advertising
 - ▶ advertisements intended to reach specific consumer groups
 - ▶ so as to improve the effectiveness and cost-benefit
- ▶ online advertising networks
 - ▶ web services that connect advertisers to web publishers
 - ▶ e.g., a banner advertisement at a personal web site
- ▶ Real Time Bidding
 - ▶ platform for real-time auction of online advertisements
 - ▶ web publishers offer display space on user's visit
 - ▶ with user's attributes and activity history (tracked by cookies)
 - ▶ bid managers provides a platform for auction
 - ▶ advertisers place a bid for advertisement
 - ▶ decide the price based on the provided information
 - ▶ retargetting: for users who visited the advertising company in the past
 - ▶ RTB auction process completes in less than 100ms

collaborative filtering

- ▶ several well-known algorithms
- ▶ example: simple correlation analysis between users
 - ▶ compute correlation between users to find similar users
 - ▶ rate item as a sum of others' scores weighted by the similarity

example: purchase history

user	item						
	a	b	c	d	e	f	...
A	1		1		1		...
B			1	1			...
C	1	1					...
D	1		1		1		...
...							...

compute the scores of items that A does not have but A's similar users have

user	similarity σ	item						
		a	b	c	d	e	f	...
A	1	1		1		1		...
S	0.88		0.88		-		0.88	...
C	0.81		0.81		-		-	...
K	0.75		-		-		-	...
F	0.73		0.73		0.73		0.73	...
score			2.50		0.73		1.61	...

Example: Netflix Prize

- ▶ an open annual competition for collaborative filtering algorithms to predict user ratings for movies
- ▶ sponsored by Netflix, an online DVD-rental/download service company
- ▶ competition: data set
< *user_id, movie_id, date_of_grade, grade* >
 - ▶ training data set (100 million ratings)
 - ▶ qualifying data set (2.8 million ratings)
 - ▶ quiz data set (1.4 million)
 - ▶ test data set (1.4 million)
 - ▶ results are scored by root mean squared error
- ▶ competition started in 2006 and ended in 2009
 - ▶ criticized by privacy advocates

distances

various distances

- ▶ Euclidean distance
- ▶ standardized Euclidean distance
- ▶ Minkowski distance
- ▶ Mahalanobis distance

similarities

- ▶ binary vector similarities
- ▶ n-dimensional vector similarities

properties of distance

a metric of distance $d(x, y)$ between 2 points (x, y) in space
positivity

$$d(x, y) \geq 0$$

$$d(x, y) = 0 \Leftrightarrow x = y$$

symmetry

$$d(x, y) = d(y, x)$$

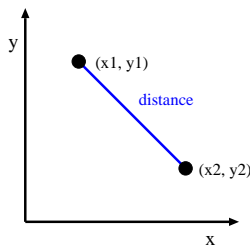
triangle inequality

$$d(x, z) \leq d(x, y) + d(y, z)$$

Euclidean distance

word “distance” usually means “Euclidean distance”
a distance of 2 points (x, y) in a n-dimensional space

$$d(x, y) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$$

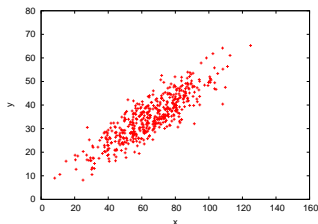


euclidean distance in 2-dimensional space

standardized Euclidean distance

- ▶ when variances are different among variables, distances are affected.
- ▶ standard Euclidean distance: normalized by dividing the Euclidean distance by the variance of each variable

$$d(x, y) = \sqrt{\sum_{k=1}^n \left(\frac{x_k}{s_k} - \frac{y_k}{s_k} \right)^2} = \sqrt{\sum_{k=1}^n \frac{(x_k - y_k)^2}{s_k^2}}$$

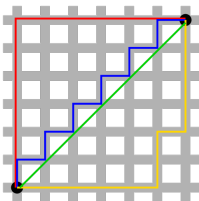


Minkowski distance

generalization of Euclidean distance: as parameter r grows, a short cut crossing different axes is preferred more

$$d(x, y) = \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{\frac{1}{r}}$$

- ▶ $r = 1$: Manhattan distance
 - ▶ Hamming distance: for 2 strings of equal length, the number of positions at which the corresponding symbols are different.
 - ▶ example: the hamming distance of 111111 and 101010 is 3
- ▶ $r = 2$: Euclidean distance



Manhattan distance vs. Euclidean distance

vector norm (1/2)

vector norm: the length of a vector

$$\|x\| \text{ where } x \text{ is a vector}$$

the l_n -norm of x is defined by Minkowski distance as

$$\|x\|_n = \sqrt[n]{\sum_i |x_i|^n}$$

l_0 -norm: the total number of non-zero elements in a vector

$$\|x\|_0 = \#(i|x_i \neq 0)$$

l_1 -norm: sum of absolute difference

$$\|x\|_1 = \sum_i |x_i|$$

l_2 -norm: Euclidean distance

$$\|x\|_2 = \sqrt{\sum_i |x_i|^2}$$

l_∞ -norm: the maximum entry's magnitude of a vector

$$\|x\|_\infty = \max(|x_i|)$$

vector norm (2/2)

For the example vector $x = (1, 2, 3)$

$$\|x\|_0 = 3 = 3.000$$

$$\|x\|_1 = 6 = 6.000$$

$$\|x\|_2 = \sqrt{14} = 3.742$$

$$\|x\|_3 = 6^{2/3} = 3.302$$

$$\|x\|_4 = 2^{1/4}\sqrt{7} = 3.146$$

$$\|x\|_\infty = 3 = 3.000$$



unit circles of l_p -norm with various values of p

Mahalanobis distance

a distance that takes correlations into account, when correlation exists between variables

$$\text{mahalanobis}(x, y) = (x - y)\Sigma^{-1}(x - y)^T$$

here, Σ^{-1} is the inverse matrix of its covariance matrix

similarities

similarity

- ▶ numerical measure of how alike 2 data objects are

properties of similarity

positivity

$$0 \leq s(x, y) \leq 1$$

$$s(x, y) = 1 \Leftrightarrow x = y$$

symmetry

$$s(x, y) = s(y, x)$$

in general, triangle inequality does not apply to similarities

similarity between binary vectors

Jaccard coefficient

- ▶ used for similarity between binary vectors in which the occurrences of 1 is much smaller than the occurrences of 0
- ▶ example: as a metric of similarity by occurrences of words in documents
- ▶ many words do not appear in both documents \Rightarrow not considered
- ▶ the following table shows the relationship of each item

		vector y	
		1	0
vector x	1	n_{11}	n_{10}
	0	n_{01}	n_{00}

Jaccard coefficient:

$$J = \frac{n_{11}}{n_{11} + n_{10} + n_{01}}$$

similarity between vectors

similarity between (non-binary) vectors

- ▶ example: similarity of documents where frequencies of words are also taken into consideration

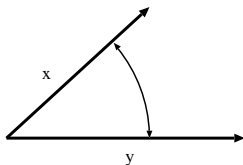
cosine similarity

- ▶ take the angle (cosine) of (x, y) of vectors
- ▶ normalized by the length of the vector \Rightarrow length is not considered

$$\cos(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

$x \cdot y = \sum_{k=1}^n x_k y_k$: product of vectors

$\|x\| = \sqrt{\sum_{k=1}^n x_k^2} = \sqrt{x \cdot x}$: length of the vector



example: cosine similarity

$$x = 3 \ 2 \ 0 \ 5 \ 0 \ 0 \ 0 \ 2 \ 0 \ 0$$

$$y = 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 2$$

$$x \cdot y = 3 * 1 + 2 * 1 = 5$$

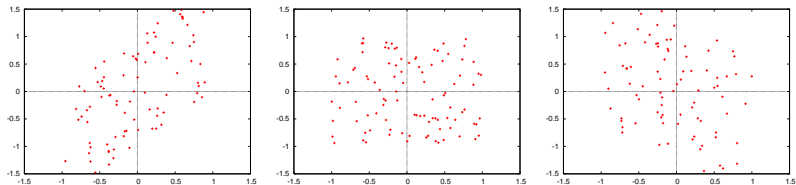
$$\|x\| = \sqrt{3 * 3 + 2 * 2 + 5 * 5 + 2 * 2} = \sqrt{42} = 6.481$$

$$\|y\| = \sqrt{1 * 1 + 1 * 1 + 2 * 2} = \sqrt{6} = 2.449$$

$$\cos(x, y) = \frac{5}{6.481 * 2.449} = 0.315$$

scatter plots and correlation

- ▶ explores relationships between 2 variables
 - ▶ X-axis: variable X
 - ▶ Y-axis: corresponding value of variable Y
- ▶ you can identify
 - ▶ whether variables X and Y related
 - ▶ no relation, positive correlation, negative correlation
- ▶ correlation coefficient: a measure of the strength and direction of correlation



examples: positive correlation 0.7 (left), no correlation 0.0 (middle), negative correlation -0.5 (right)

correlation

- ▶ covariance:

$$\sigma_{xy}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

- ▶ correlation coefficient:

$$\rho_{xy} = \frac{\sigma_{xy}^2}{\sigma_x \sigma_y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

- ▶ correlation coefficient: the covariance of 2 variables normalized by their product of their standard deviations, a value between -1 and $+1$ inclusive.
- ▶ sensitive to outliers. so, you should use a scatter plot to observe outliers.
- ▶ correlation and causality
 - ▶ correlation does not imply causal relationship
 - ▶ third factor C causes both A and B (e.g., break and test score)
 - ▶ coincidence

computing correlation coefficient (1)

sum of squares

$$\begin{aligned}\sum_{i=1}^n (x_i - \bar{x})^2 &= \sum_{i=1}^n (x_i^2 - 2x_i\bar{x} + \bar{x}^2) \\ &= \sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + n\bar{x}^2 \\ &= \sum_{i=1}^n x_i^2 - 2\bar{x} \cdot n\bar{x} + n\bar{x}^2 \\ &= \sum_{i=1}^n x_i^2 - n\bar{x}^2 = \sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n}\end{aligned}$$

sum of products

$$\begin{aligned}\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) &= \sum_{i=1}^n (x_i y_i - x_i \bar{y} - \bar{x} y_i + \bar{x} \bar{y}) \\ &= \sum_{i=1}^n x_i y_i - \bar{x} \sum_{i=1}^n y_i - \bar{y} \sum_{i=1}^n x_i + n\bar{x} \bar{y} \\ &= \sum_{i=1}^n x_i y_i - \bar{x} \cdot n\bar{y} - \bar{y} \cdot n\bar{x} + n\bar{x} \bar{y} \\ &= \sum_{i=1}^n x_i y_i - n\bar{x} \bar{y} = \sum_{i=1}^n x_i y_i - \frac{(\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i)}{n}\end{aligned}$$

computing correlation coefficient (2)

correlation coefficient

$$\begin{aligned}\rho_{xy} &= \frac{\sigma_{xy}^2}{\sigma_x \sigma_y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \\ &= \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sqrt{(\sum_{i=1}^n x_i^2 - n \bar{x}^2)(\sum_{i=1}^n y_i^2 - n \bar{y}^2)}} \\ &= \frac{\sum_{i=1}^n x_i y_i - \frac{(\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i)}{n}}{\sqrt{(\sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n})(\sum_{i=1}^n y_i^2 - \frac{(\sum_{i=1}^n y_i)^2}{n})}}\end{aligned}$$

other correlation coefficients

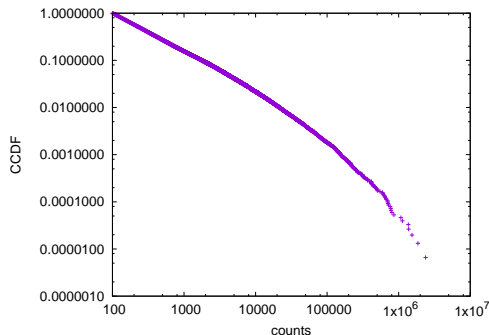
- ▶ Pearson's product-moment correlation coefficient
 - ▶ or simply "correlation coefficient" (what we have learned)
- ▶ rank correlation coefficient: relationships between different rankings on the same set of items
 - ▶ Spearman's rank correlation coefficient
 - ▶ Kendall's rank correlation coefficient
- ▶ others

previous exercise: CCDF plots

- ▶ plot the US surname distribution in CCDF
 - ▶ popular surnames (100 or more) in the US Census in 2000
 - ▶ csv format: comma-separated-variables

```
% head us-surnames.csv
name,rank,count,prop100k,cum_prop100k,pctwhite,pctblack,pctapi,pctaian,pct2prace,pcthispanic
SMITH,1,2376206,880.85,880.85,73.35,22.22,0.40,0.85,1.63,1.56
JOHNSON,2,1857160,688.44,1569.30,61.55,33.80,0.42,0.91,1.82,1.50
WILLIAMS,3,1534042,568.66,2137.96,48.52,46.72,0.37,0.78,2.01,1.60
BROWN,4,1380145,511.62,2649.58,60.71,34.54,0.41,0.83,1.86,1.64
JONES,5,1362755,505.17,3154.75,57.69,37.73,0.35,0.94,1.85,1.44
...

% ./make_ccdf.rb us-surnames.csv > ccdf.txt
```



script to convert the counts to CCDF

```
#!/usr/bin/env ruby

re = /^[A-Z]+,\d+,(\\d+),\\d+/ # for US surnames
#re = /^\\S+\\s+(\\d+)\\s+\\d+/ # for JAIST web server logs

n = 0
counts = Hash.new(0)
ARGF.each_line do |line|
  if re.match(line)
    counts[$1.to_i] += 1
    n += 1
  end
end

cum = 0
counts.sort.each do |key, value|
  comp = 1.0 - Float(cum) / n
  puts "#{key} #{value} #{comp}"
  cum += value
end

$stderr.puts "# #{n} entries matched"
```

cumulative surname counts

```
% cat ccdf.txt
```

```
100 1236 1.0
```

```
101 1108 0.9918507822853413
```

```
102 1084 0.9845454965022977
```

```
103 1149 0.977398447956432
```

```
104 1084 0.969822840226543
```

```
105 1103 0.9626757916806773
```

```
106 1061 0.9554034719887124
```

```
107 1028 0.9484080674618088
```

```
108 1031 0.9416302391360247
```

```
...
```

```
1380145 1 2.637287286300083e-05
```

```
1534042 1 1.9779654647278377e-05
```

```
1857160 1 1.3186436431444903e-05
```

```
2376206 1 6.593218215722452e-06
```

gnuplot script for plotting the counts in CCDF

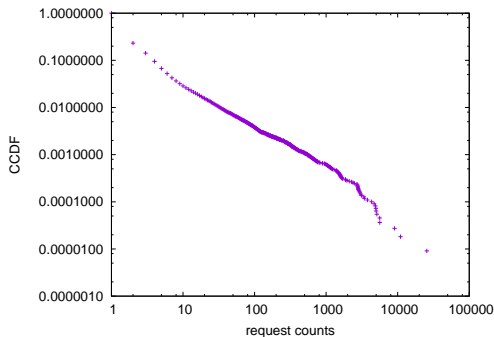
```
set logscale
set xlabel "counts"
set ylabel "CCDF"

plot "ccdf.txt" using 1:3 notitle with points
```

another exercise: CCDF plots

extract the access count of each unique content from the JAIST server access log, plot the access count distribution in CCDF

```
first edit the regular expression in the script (e.g., vim make_ccdf.rb)
% ./count_contents.rb sample_access_log > contents.txt
% ./make_ccdf.rb contents.txt > ccdf.txt
```



extracting the access count of each unique content

```
# output: URL req_count byte_count
# regular expression for apache combined log format
# host ident user time request status bytes referer agent
re = /^(S+) (\S+) (\S+) \[(.)*\] "(.*)" (\d+) (\d+|-) "(.*)" "(.*)" /
# regular expression for request: method url proto
req_re = /(\w+) (\S+) (\S+)/
contents = Hash.new([0, 0])
count = parsed = 0
ARGF.each_line do |line|
  count += 1
  if re.match(line)
    host, ident, user, time, request, status, bytes, referer, agent = $~.captures
    # ignore if the status is not success (2xx)
    next unless /2\d{2}/.match(status)
    if req_re.match(request)
      method, url, proto = $~.captures
      # ignore if the method is not GET
      next unless /GET/.match(method)
      parsed += 1
      # count contents by request and bytes
      contents[url] = [contents[url][0] + 1, contents[url][1] + bytes.to_i]
    else
      # match failed. print a warning msg
      $stderr.puts("request match failed at line #{count}: #{line.dump}")
    end
  else
    $stderr.puts("match failed at line #{count}: #{line.dump}") # match failed.
  end
end
contents.sort_by{|key, value| -value[0]}.each do |key, value|
  puts "#{key} #{value[0]} #{value[1]}"
end
$stderr.puts "# #{contents.size} unique contents in #{parsed} successful GET requests"
$stderr.puts "# parsed:#{parsed} ignored:#{count - parsed}"
```


access count of each unique content

```
% cat contents.txt
/project/linuxonandroid/Ubuntu/12.04/full/ubuntu1204-v4-full.zip 25535 17829045
/project/morefont/xiongmaozhongwen.apk 10949 13535294486
/project/morefont/zhongguoxin.apk 9047 9549531354
/project/honi/some_software/Windows/Office_Plus_2010_SP1_W32_xp911.com.rar 5616
/project/morefont/fangzhengyouyijian.apk 5609 2879391721
/pub/Linux/CentOS/5.9/extras/i386/repodata/repomd.xml 5121 12213484
/pub/Linux/CentOS/5.9/updates/i386/repodata/repomd.xml 5006 10969621
/pub/Linux/CentOS/5.9/os/i386/repodata/repomd.xml 4953 6832653
/project/npppluginmgr/xml/plugins.md5.txt 4881 1369547
/project/winpenpack/X-LenMus/releases/X-LenMus_5.3.1_rev5.zip 4689 990250462

...

/pub/Linux/openSUSE/distribution/12.3/repo/oss/suse/x86_64/gedit-3.6.2-2.1.2.x8
/pub/sourceforge/n/nz/nzbcatcher/source/?C=D;O=A 1 1075
/ubuntu/pool/universe/m/mmass/mmass_5.4.1.orig.tar.gz 1 3754849
```

cumulative access counts

```
% cat ccdf.txt
1 84414 1.0
2 9813 0.2315731022366253
3 5199 0.14224463601358184
4 3034 0.0949177537254331
5 1636 0.06729902688137779
6 1083 0.05240639764048316
7 663 0.04254776838138241
8 495 0.03651243024769468
9 367 0.03200640856417214
10 274 0.028665580366489807

...

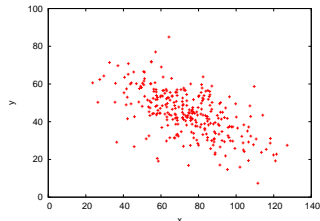
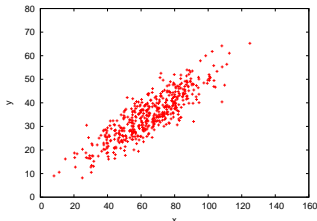
5616 1 3.6412296432475344e-05
9047 1 2.730922232441202e-05
10949 1 1.8206148216237672e-05
25535 1 9.103074108174347e-06
```

today's exercise: computing correlation coefficient

- ▶ compute correlation coefficient using the sample data sets
 - ▶ correlation-data-1.txt, correlation-data-2.txt

correlation coefficient

$$\rho_{xy} = \frac{\sigma_{xy}^2}{\sigma_x \sigma_y} = \frac{\sum_{i=1}^n x_i y_i - \frac{(\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i)}{n}}{\sqrt{(\sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n})(\sum_{i=1}^n y_i^2 - \frac{(\sum_{i=1}^n y_i)^2}{n})}}$$



data-1:r=0.87 (left), data-2:r=-0.60 (right)

script to compute correlation coefficient

```
#!/usr/bin/env ruby

# regular expression for matching 2 floating numbers
re = /([+]?[0-9]+\.[0-9]+)?\s+([+]?[0-9]+\.[0-9]+)?/

sum_x = 0.0      # sum of x
sum_y = 0.0      # sum of y
sum_xx = 0.0     # sum of x^2
sum_yy = 0.0     # sum of y^2
sum_xy = 0.0     # sum of xy
n = 0            # the number of data

ARGF.each_line do |line|
  if re.match(line)
    x = $1.to_f
    y = $2.to_f
    sum_x += x
    sum_y += y
    sum_xx += x**2
    sum_yy += y**2
    sum_xy += x * y
    n += 1
  end
end

r = (sum_xy - sum_x * sum_y / n) /
  Math.sqrt((sum_xx - sum_x**2 / n) * (sum_yy - sum_y**2 / n))

printf "n:%d r:%.3f\n", n, r
```

today's exercise 2: similarity

- ▶ compute similarity in data
 - ▶ data from “Programming Collective Intelligence” Section 2
 - ▶ movie rating scores of 7 people: scores.txt

```
% cat scores.txt
# A dictionary of movie critics and their ratings of a small set of movies
'Lisa Rose': 'Lady in the Water': 2.5, 'Snakes on a Plane': 3.5, 'Just My Luck': 3.0, 'Superman Returns':
'Gene Seymour': 'Lady in the Water': 3.0, 'Snakes on a Plane': 3.5, 'Just My Luck': 1.5, 'Superman Returns
'Michael Phillips': 'Lady in the Water': 2.5, 'Snakes on a Plane': 3.0, 'Superman Returns': 3.5, 'The Nigh
'Claudia Puig': 'Snakes on a Plane': 3.5, 'Just My Luck': 3.0, 'The Night Listener': 4.5, 'Superman Return
'Mick LaSalle': 'Lady in the Water': 3.0, 'Snakes on a Plane': 4.0, 'Just My Luck': 2.0, 'Superman Returns
'Jack Matthews': 'Lady in the Water': 3.0, 'Snakes on a Plane': 4.0, 'The Night Listener': 3.0, 'Superman
'Toby': 'Snakes on a Plane':4.5,'You, Me and Dupree':1.0,'Superman Returns':4.0
```

score data

- ▶ simplistic example: data is too small
- ▶ summarized in the following table

```
#name: 'Lady in the Water' 'Snakes on a Plane' 'Just My Luck' 'Superman Returns'
Lisa Rose:      2.5 3.5 3.0 3.5 3.0
Gene Seymour:   3.0 3.5 1.5 5.0 3.0
Michael Phillips: 2.5 3.0 - 3.5 4.0
Claudia Puig:   - 3.5 3.0 4.0 4.5
Mick LaSalle:   3.0 4.0 2.0 3.0 3.0
Jack Matthews:  3.0 4.0 - 5.0 3.0
Toby:           - 4.5 - 4.0 -
```

similarity computation

- ▶ create a similarity matrix using cosine similarity

```
% ruby similarity.rb scores.txt
Lisa Rose:      1.000 0.959 0.890 0.921 0.982 0.895 0.708
Gene Seymour:   0.959 1.000 0.950 0.874 0.962 0.979 0.783
Michael Phillips: 0.890 0.950 1.000 0.850 0.929 0.967 0.693
Claudia Puig:   0.921 0.874 0.850 1.000 0.875 0.816 0.695
Mick LaSalle:   0.982 0.962 0.929 0.875 1.000 0.931 0.727
Jack Matthews:  0.895 0.979 0.967 0.816 0.931 1.000 0.822
Toby:           0.708 0.783 0.693 0.695 0.727 0.822 1.000
```

similarity computation script (1/2)

```
# regular expression to read data
#       'name': 'title0': score0, 'title1': score1, ...
re = /'(.+?)':\s+(\S.*)/
name2uid = Hash.new      # keeps track of name to uid mapping
title2tid = Hash.new    # keeps track of title to tid mapping
scores = Hash.new       # scores[uid][tid]: score of title_id by user_id

# read data into scores[uid][tid]
ARGF.each_line do |line|
  if re.match(line)
    name = $1
    ratings = $2.split(",")

    if name2uid.has_key?(name)
      uid = name2uid[name]
    else
      uid = name2uid.length
      name2uid[name] = uid
      scores[uid] = {} # create empty hash for title and score pairs
    end
    ratings.each do |rating|
      if rating.match(/'(.+?)':\s*(\d.\d)/)
        title = $1
        score = $2.to_f
        if title2tid.has_key?(title)
          tid = title2tid[title]
        else
          tid = title2tid.length
          title2tid[title] = tid
        end
        scores[uid][tid] = score
      end
    end
  end
end
end
```


similarity computation script (2/2)

```
# compute cosine similarity between 2 users
def comp_similarity(h1, h2)
  sum_xx = 0.0 # sum of x^2
  sum_yy = 0.0 # sum of y^2
  sum_xy = 0.0 # sum of xy
  score = 0.0 # similarity score

  h1.each do |tid, score|
    sum_xx += score**2
    if h2.has_key?(tid)
      sum_xy += score * h2[tid]
    end
  end
  h2.each_value do |score|
    sum_yy += score**2
  end
  denom = Math.sqrt(sum_xx) * Math.sqrt(sum_yy)
  if denom != 0.0
    score = sum_xy / denom
  end
  return score
end

# create n x n matrix of similarities between users
n = name2uid.length
similarities = Array.new(n) { Array.new(n) }
for i in 0 .. n - 1
  printf "%-18s", name2uid.key(i) + ':'
  for j in 0 .. n - 1
    similarities[i][j] = comp_similarity(scores[i], scores[j])
    printf "%.3f ", similarities[i][j]
  end
  print "\n"
end
```

more realistic data set

- ▶ MovieLens:

- <http://grouplens.org/datasets/movielens/>

- ▶ dataset for collaborative filtering research by Univ. of Minnesota
 - ▶ movie ratings by users, data size:100K, 1M, 10M
 - ▶ u.data: rating data set
 - ▶ dataset includes other info (e.g., demographic info about the users, info about movies)

```
% head u.data
#user_id item_id rating timestamp
196      242      3      881250949
186      302      3      891717742
22       377      1      878887116
244      51       2      880606923
166      346      1      886397596
298      474      4      884182806
115      265      2      881171488
253      465      5      891628467
305      451      3      886324817
6        86       3      883603013
...
```

assignment 1: the finish time distribution of a marathon

- ▶ purpose: investigate the distribution of a real-world data set
- ▶ data: the finish time records from honolulu marathon 2015
 - ▶ <http://www.pseresults.com/events/741/results>
 - ▶ the number of finishers: 21,554
- ▶ items to submit
 1. mean, standard deviation and median of the total finishers, male finishers, and female finishers
 2. the distributions of finish time for each group (total, men, and women)
 - ▶ plot 3 histograms for 3 groups
 - ▶ use 10 minutes for the bin size
 - ▶ use the same scale for the axes to compare the 3 plots
 3. CDF plot of the finish time distributions of the 3 groups
 - ▶ plot 3 groups in a single graph
 4. discuss differences in finish time between male and female. what can you observe from the data?
 5. optional
 - ▶ other analysis of your choice (e.g., discussion on differences among age groups)
- ▶ submission format: a single PDF file including item 1-5
- ▶ submission method: upload the PDF file through SFC-SFS
- ▶ submission due: 2016-05-17

honolulu marathon data set

data format (compacted to fit in the slide)

Chip							Cat	Cat					Gndr	Gndr	
Place	Time	Number	Lname	Fname	Country	Category	Place	Total	5K	10K	40K	Place	Total	Pace	
1	2:11:43	3	Kiprotich	Filex	KEN	Melite	1	5	16:07	31:40	...	2:04:48	1	11346	5:02
2	2:12:46	1	Chebet	Wilson	KEN	Melite	2	5	16:07	31:41	...	2:05:57	2	11346	5:04
3	2:13:24	8	Limo	Daniel	KEN	Melite	3	5	16:06	31:41	...	2:06:13	3	11346	5:06
4	2:15:27	6	Kwambai	Robert	KEN	Melite	4	5	16:08	31:41	...	2:07:29	4	11346	5:10
5	2:18:36	4	Mungara	Kenneth	KEN	Melite	5	5	16:07	31:40	...	2:09:42	5	11346	5:18
6	2:27:58	11	Neuschwander	Florian	DEU	M30-34	1	1241	17:46	34:50	...	2:20:31	6	11346	5:39
7	2:28:34	F1	Chepkirui	Joyce	KEN	Welite	1	7	16:53	33:21	...	2:20:56	1	10207	5:40
8	2:28:42	28803	Takahashi	Koji	JPN	M25-29	1	974	16:54	33:22	...	2:20:52	7	11346	5:41
9	2:28:55	F5	Karimi	Lucy	KEN	Welite	2	7	16:54	33:22	...	2:20:58	2	10207	5:41
10	2:29:44	F6	Ochichi	Isabella	KEN	Welite	3	7	16:53	33:22	...	2:21:46	3	10207	5:43
...															

- ▶ Chip Time: finish time
- ▶ Number: bib number
- ▶ Category: MELite, WELite, M15-19, M20-24, ..., W15-29, W20-24, ...
 - ▶ note: 2 runners have "No Age" for Category, and num:18035 doesn't have cat/gender totals and its cat/gender placements are not reflected to the following entries
- ▶ Country: 3-letter country code: e.g., JPN, USA
- ▶ check the number of the total finishers when you extract the finishers

assignment 1: additional hints

- ▶ summary statistics: results can be in a table
- ▶ histograms:
 - ▶ X-axis: finish time (chip time) in 10min bin
 - ▶ Y-axis: the number of finishers for each bin
- ▶ CDF plot: (3 plots in a single figure)
 - ▶ X-axis: finish time
 - ▶ Y-axis: CDF [0:1]
- ▶ pages for the report: about 3-6 pages (source code not required)

sample code for extracting chip-time

```
# regular expression to read chiptime
re = /\d+\s+(\d{1,2}:\d{2}:\d{2})\s+/

ARGF.each_line do |line|
  if re.match(line)
    puts "#{$1}"
  end
end
```

summary

Class 6 Correlation

- ▶ Online recommendation systems
- ▶ Distance
- ▶ Correlation coefficient
- ▶ exercise: correlation analysis

next class

Class 7 Multivariate analysis (5/23)

- ▶ Data sensing and GeoLocation
- ▶ Linear regression
- ▶ Principal Component Analysis
- ▶ exercise: linear regression
- ▶ **assignment 2**