

インターネット計測とデータ解析 第7回

長 健二郎

2015年6月1日

前回のおさらい

第6回 相関 (5/25)

- ▶ オンラインお勧めシステム
- ▶ 距離と類似度
- ▶ 相関係数
- ▶ 演習: 相関

今日のテーマ

第7回 多変量解析

- ▶ データセンシング
- ▶ 地理的位置情報 (geo-location)
- ▶ 線形回帰
- ▶ 主成分分析
- ▶ 演習: 線形回帰

多変量データ解析

- ▶ 一変数解析 (univariate analysis)
 - ▶ 変数をひとつずつ独立して扱う
- ▶ 多変量解析 (multivariate analysis)
 - ▶ 複数の変数を同時に扱う
 - ▶ コンピュータの普及で発展
 - ▶ 隠れたトレンドを探る (データマイニング)

データセンシング

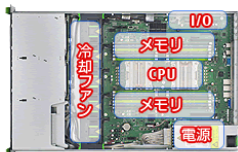
- ▶ データセンシング: 遠隔からデータを収集する
- ▶ インターネット経由でさまざまなセンサー情報が取得可能に
 - ▶ 気象情報、電力消費、その他さまざまな情報

省エネルギー技術

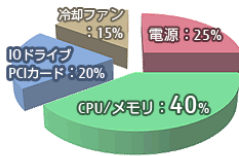
- ▶ 電力消費の削減: あらゆる技術分野で課題
 - ▶ センサー情報を用いた賢い制御で効率化
- ▶ 個別機器の効率化から全体の効率化へ
 - ▶ 例: PC サーバの効率化とデータセンタの効率化

PC サーバの効率化

- ▶ PC 内のセンサー情報を使った効率化制御
 - ▶ 温度、電圧、消費電力、ファン回転数
- ▶ PC サーバの電力消費内訳
 - ▶ CPU/メモリ: 40%
 - ▶ 高集積化・省電力化、クロック・電圧制御
 - ▶ 電源: 25%
 - ▶ ロス削減 (AC-DC, DC-DC)
 - ▶ IO: 20%
 - ▶ 省電力機能、省電力ディスク/SSD
 - ▶ 冷却ファン: 15%
 - ▶ 効率配置、エアフロー設計、制御最適化



PRIMERGY RX300 S7 の内部構造 (例)

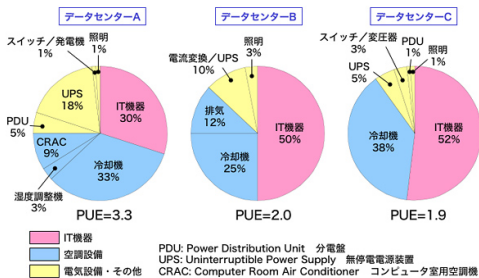


サーバ内の電力消費の割合の例

source: <http://jp.fujitsu.com/platform/server/primergy/solution/power-saving/option/>

データセンタの効率化

- ▶ データセンタ需要増加に伴う電力消費の急増
 - ▶ 空調電力や電力ロスが大きい
- ▶ IT 機器: 使用機器の効率化、高温仕様機器の利用
- ▶ 空調設備: 仕様見直し、エアフロー・熱負荷設計、空調器の効率化、外気冷却
- ▶ 給電設備: ロス削減、高圧給電、DC 給電、UPS 効率化、自然エネルギー利用
- ▶ 全体設計: 全体効率化、柔軟に消費電力変動に追従、入退室削減、アイドル機器の停止などの制御



source:

<http://librarytaisei.jp/feature/datacenter/002/>

位置情報サービス

- ▶ 場所に応じた情報の提供
- ▶ 地図サービス、ナビゲーション、時刻表
- ▶ 近隣のレストランや店舗検索 (広告への利用)
- ▶ その他、さまざまなサービスの可能性

例: 駅.Locky

- ▶ 名古屋大学 河口研が開発した時刻表サービス
 - ▶ WiFi 情報収集プロジェクトから派生した人気アプリ
- ▶ iPhone/Android 用 App
- ▶ 位置情報から最寄りの駅の時刻表を検索
 - ▶ GPS/WiFi による位置情報取得
 - ▶ 同時に、端末から見える WiFi 基地局情報を収集
- ▶ 次の出発までの時間をカウントダウン
 - ▶ 時刻表の閲覧も可能
- ▶ ユーザ提供型の時刻表データベース



GPS (Global Positioning System)

- ▶ 約 30 個の GPS 衛星
- ▶ 元来はアメリカ合州国の軍事用
 - ▶ 当初は意図的に誤差データを加え 100m 程度の精度にしていた
 - ▶ 2000 年に誤差混入が廃止され、10m 程度の精度になる
- ▶ さまざまな民生用途
 - ▶ カーナビ、携帯端末、デジカメ
- ▶ 測位: 3 個の GPS 衛星からの距離から位置を特定
 - ▶ GPS 信号には衛星の位置、時刻情報が含まれる
 - ▶ 距離は GPS 衛星からの時刻データのずれから計算
 - ▶ 受信機の時刻補正のため 4 個の衛星を捕捉する必要
 - ▶ より多くの衛星を捕捉すれば精度が向上
- ▶ 欠点
 - ▶ 衛星が見えないと使えない
 - ▶ 初期位置取得時間
- ▶ 高精度化: 加速度センサーや振動型ジャイロスコープと組合せ

基地局を利用した位置情報

- ▶ 端末は接続している基地局が分かる
 - ▶ 基地局側からも接続している端末が分かる
 - ▶ 接続していなくても電波を受信できる基地局が分かる
- ▶ 基地局がその緯度経度を発信するサービスも存在
- ▶ 屋内でも利用可能
 - ▶ 他にも、音波、可視光などによるアプローチも存在
- ▶ GPS との組合せによる精度向上

インターネットの特徴量

通信レベルの特徴量

- ▶ 回線容量、スループット
- ▶ 遅延
- ▶ ジッタ
- ▶ パケットロス

測定手法

- ▶ アクティブ計測: ping 等、計測パケットを注入
- ▶ パッシブ計測: 計測用パケットを使わない
 - ▶ 2点で観測して比較
 - ▶ TCP の挙動等から推測
 - ▶ トランスポート機能内部で情報収集

遅延

▶ 遅延成分

- ▶ 遅延 = 伝搬遅延 + キュー待ち遅延 + その他
- ▶ パケット毎に一定の遅延成分とパケット長に比例する成分
- ▶ 輻輳がなければ、遅延は伝搬遅延 + α

▶ 遅延計測

- ▶ RTT(round trip time) 計測: パケットの往復時間
- ▶ 一方向遅延計測: 両端の時刻同期が必要

- ▶ 遅延の平均
- ▶ 最大遅延: 例えば、一般に音声会話は 400ms 以下が必要
- ▶ ジッタ: 遅延値のばらつき
 - ▶ リアルタイム通信でのバッファサイズの決定
 - ▶ 下位層の影響: 無線での再送、イーサネットのコリジョン等

代表的な遅延値

- ▶ パケット伝送時間 (ワイヤースピード)
 - ▶ 1500 bytes at 10Mbps: 1.2 msec
 - ▶ 1500 bytes at 100Mbps: 120 usec
 - ▶ 1500 bytes at 1Gbps: 12 usec
- ▶ ファイバー中の伝搬速度: 約 200,000 km/s
 - ▶ 100km round-trip: 1 msec
 - ▶ 20,000km round-trip: 200 msec
- ▶ 衛星の RTT
 - ▶ LEO (Low-Earth Orbit): 200 msec
 - ▶ GEO (Geostationary Orbit): 600 msec

パケットロス

パケットロス率

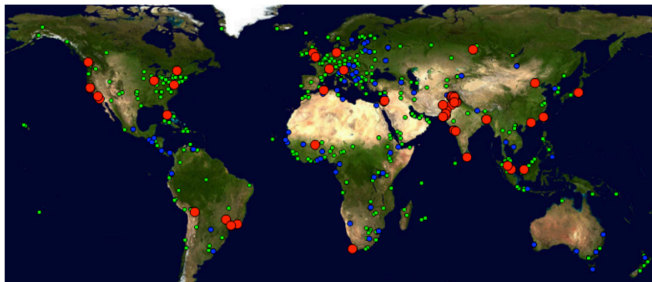
- ▶ パケットロスがランダムに発生すると見なせればロス率だけでいいが
- ▶ 一定間隔のプロブでは分からない傾向
 - ▶ バースト的なロス: バッファ溢れ等
 - ▶ パケット長による違い: 無線でのビット誤り等

pingER project

- ▶ the Internet End-to-end Performance Measurement (IEPM) project by SLAC
- ▶ using ping to measure rtt and packet loss around the world
 - ▶ <http://www-iepm.slac.stanford.edu/pinger/>
 - ▶ started in 1995
 - ▶ over 600 sites in over 125 countries

pingER project monitoring sites

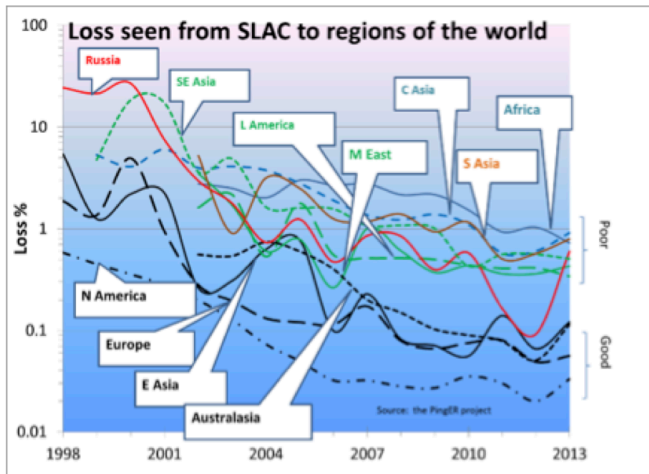
- ▶ monitoring (red), beacon (blue), remote (green) sites
 - ▶ beacon sites are monitored by all monitors



from pingER web site

pingER packet loss

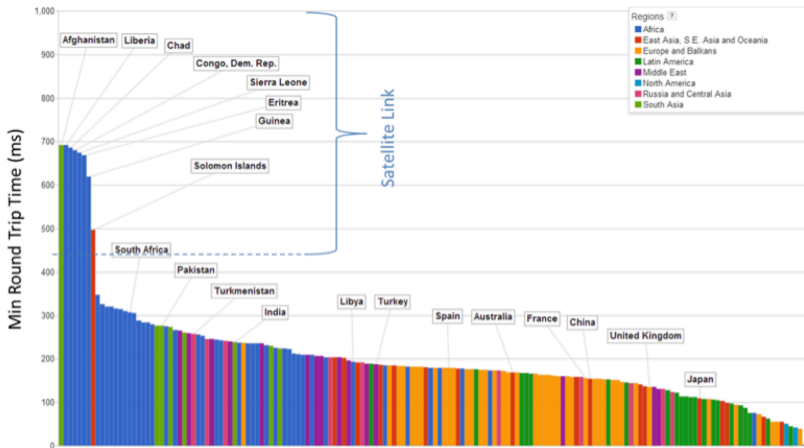
- ▶ packet loss observed from SLAC in the west coast
- ▶ exponential improvement in 15 years



from pingER web site

pinger minimum rtt

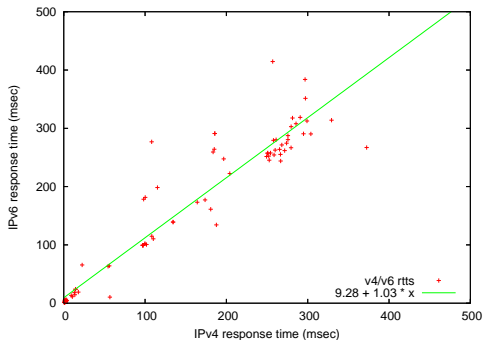
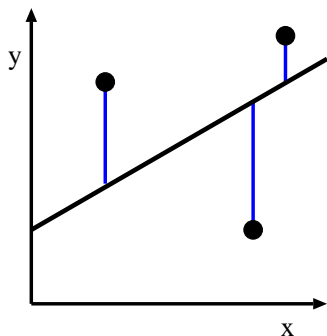
- ▶ minimum rtt observed from SLAC in the west coast



from pingER web site

線形回帰 (linear regression)

- ▶ データに一次関数を当てはめる
 - ▶ 最小二乗法 (least square method): 誤差の二乗和を最小にする



最小二乗法 (least square method)

誤差の二乗和を最小にする一次関数を求める

$$f(x) = b_0 + b_1x$$

切片と傾きの求め方

$$b_1 = \frac{\sum xy - n\bar{x}\bar{y}}{\sum x^2 - n(\bar{x})^2}$$

$$b_0 = \bar{y} - b_1\bar{x}$$

ここで

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\sum xy = \sum_{i=1}^n x_i y_i \quad \sum x^2 = \sum_{i=1}^n (x_i)^2$$

最小二乗法の導出

i 番目の変数の誤差 $e_i = y_i - (b_0 + b_1 x_i)$ 、 n 回の観測における誤差の平均は

$$\bar{e} = \frac{1}{n} \sum_i e_i = \frac{1}{n} \sum_i (y_i - (b_0 + b_1 x_i)) = \bar{y} - b_0 - b_1 \bar{x}$$

誤差平均が 0 になるようにすると $b_0 = \bar{y} - b_1 \bar{x}$

b_0 を b_1 で表現すると $e_i = y_i - \bar{y} + b_1 \bar{x} - b_1 x_i = (y_i - \bar{y}) - b_1 (x_i - \bar{x})$

誤差の二乗和 SSE は

$$SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n [(y_i - \bar{y})^2 - 2b_1(y_i - \bar{y})(x_i - \bar{x}) + b_1^2(x_i - \bar{x})^2]$$

分散に書き直す

$$\begin{aligned} \frac{SSE}{n} &= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 - 2b_1 \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) + b_1^2 \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= \sigma_y^2 - 2b_1 \sigma_{xy} + b_1^2 \sigma_x^2 \end{aligned}$$

SSE を最小にする b_1 は、この式を b_1 の 2 次式とみて b_1 について微分して 0 と置く

$$\frac{1}{n} \frac{d(SSE)}{db_1} = -2\sigma_{xy} + 2b_1 \sigma_x^2 = 0$$

$$\text{すなわち } b_1 = \frac{\sigma_{xy}^2}{\sigma_x^2} = \frac{\sum xy - n\bar{x}\bar{y}}{\sum x^2 - n(\bar{x})^2}$$

主成分分析 (principal component analysis; PCA)

主成分分析の目的

- ▶ 複数の変数間の関係を、少数の互いに独立な合成変数 (成分) で近似

共分散行列の固有値問題として解ける

主成分分析の応用

- ▶ 次元減少
 - ▶ 寄与率の大きい順に主成分を取る、寄与率の小さい成分は無視できる
- ▶ 主成分のラベル付け
 - ▶ 主成分の構成要素から、その意味を読みとる

注意点

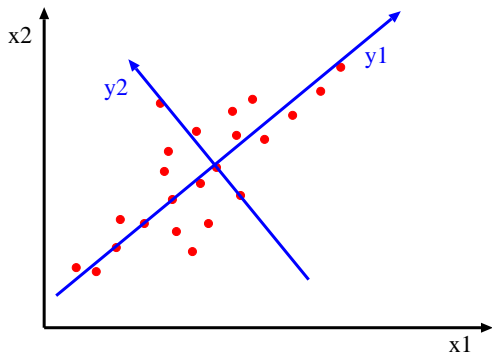
- ▶ あくまで、ばらつきの大きい成分を抜き出すだけ
 - ▶ とくに各軸の単位が違う場合は注意
- ▶ 機械的に複雑な関係を分析できる便利な手法であるが、それで複雑な関係が説明できる訳ではない

主成分分析の直観的な説明

座標変換の観点から 2 次元の図で説明すると

- ▶ データのばらつきが最も大きい方向に重心を通る直線 (第 1 主成分軸) を引く
- ▶ 第 1 主成分軸に直交し、次にばらつきが大きい方向に第 2 主成分軸を引く
- ▶ 同様に第 3 主成分軸以降を引く

例えば、「身長」と「体重」を「体の大きさ」と「太り具合」に変換。
「座高」や「胸囲」など変数が増えても同様



主成分分析 (おまけ)

主成分の単位ベクトルは、共分散行列の固有ベクトルとして求める
 X を d 次の変数、これを主成分 Y に変換する $d \times d$ の直交行列 P を求める

$$Y = P^T X$$

これを $cov(Y)$ は対角行列 (各変数が独立)、かつ P は直交行列 ($P^{-1} = P^T$) という制約のもとで解く
 Y の共分散行列は

$$\begin{aligned} cov(Y) &= E[YY^T] = E[(P^T X)(P^T X)^T] = E[(P^T X)(X^T P)] \\ &= P^T E[XX^T]P = P^T cov(X)P \end{aligned}$$

したがって

$$P cov(Y) = P P^T cov(X) P = cov(X) P$$

P を $d \times 1$ 行列でかくと、

$$P = [P_1, P_2, \dots, P_d]$$

また、 $cov(Y)$ は対角行列 (各変数が独立) なので

$$cov(Y) = \begin{bmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_d \end{bmatrix}$$

書き直すと

$$[\lambda_1 P_1, \lambda_2 P_2, \dots, \lambda_d P_d] = [cov(X)P_1, cov(X)P_2, \dots, cov(X)P_d]$$

$\lambda_i P_i = cov(X)P_i$ において、 P_i は X の共分散行列の固有ベクトルであることが分かる
したがって、固有ベクトルを見つければ求めていた変換行列 P が得られる

課題 1: ホノルルマラソン 2014 完走時間分布のプロット

- ▶ ねらい: 実データから分布を調べる
- ▶ データ: 2014 年のホノルルマラソンの記録
 - ▶ <http://www.pseresults.com/events/647/results>
 - ▶ 完走者 21,815 人
- ▶ 提出項目
 1. 全完走者、男性完走者、女性完走者それぞれの、完走時間の平均、標準偏差、中間値
 2. それぞれの完走時間のヒストグラム
 - ▶ 3つのヒストグラムを別々の図に書く
 - ▶ ビン幅は 10 分にする
 - ▶ 3つのプロットは比較できるように目盛を合わせる
 3. それぞれの CDF プロット
 - ▶ ひとつの図に 3つのプロットを書く
 4. オプション
 - ▶ 年代別や国別の CDF プロットなど自由
 5. 考察
 - ▶ データから読みとれることを記述
- ▶ 提出形式: レポートをひとつの PDF ファイルにして SFC-SFS から提出
- ▶ 提出〆切: 2015 年 5 月 27 日

ホノルルマラソンデータ

データフォーマット

Place	Num	Chip Time	Lname	Fname	Country	Division	Div Plc	Div Tot	Sex Plc	Sex Total	10Km	21Km	30Km	40Km	Pace
1	3	2:15:35	Chebet	Wilson	KEN	Melite	1	8	1	11507	0:31:50	1:08:45	1:37:47	2:09:49	5:11
2	6	2:16:04	Lonyangata	Paul	KEN	Melite	2	8	2	11507	0:31:50	1:08:45	1:37:47	2:10:05	5:12
3	7	2:16:27	Abraha	Geb	ETH	Melite	3	8	3	11507	0:31:49	1:08:44	1:37:46	2:10:24	5:13
4	5	2:16:37	Kolum	Benjamin	KEN	Melite	4	8	4	11507	0:31:50	1:08:44	1:37:47	2:10:32	5:13
5	4	2:17:54	Adhane	Yemane	ETH	Melite	5	8	5	11507	0:31:50	1:08:45	1:37:47	2:11:06	5:16
6	2	2:17:59	Chelimo	Nicholas	KEN	Melite	6	8	6	11507	0:31:51	1:08:46	1:37:48	2:11:32	5:16
7	25151	2:27:26	Harada	Taku	JPN	M30-34	1	1218	7	11507	0:32:26	1:11:15	1:43:20	2:20:27	5:38
8	8	2:28:23	Arile	Julius	KEN	Melite	7	8	8	11507	0:31:51	1:08:44	1:38:39	2:19:39	5:40
9	30300	2:29:52	Ito	Tatsuya	JPN	M30-34	2	1218	9	11507	0:34:36	1:13:04	1:44:37	2:22:16	5:43
10	F9	2:30:23	Chepkirui	Joyce	KEN	Welite	1	9	1	10308	0:34:37	1:13:07	1:44:50	2:22:43	5:45
...															

- ▶ Chip Time: 完走時間
- ▶ Category: Melite, Welite, M15-19, M20-24, ..., W15-29, W20-24, ...
 - ▶ "No Age" となっている人がいるので注意
- ▶ Country: 3-letter country code: e.g., JPN, USA
- ▶ 完走者を抽出したら、総数が合っているかチェックすること

課題 1 補足説明

- ▶ 要約統計量: 結果を表にすればよい
- ▶ ヒストグラム: X 軸は 10 分ピンの完走時間、Y 軸は各ピンの完走者数
- ▶ CDF プロット: X 軸は完走時間、Y 軸は CDF [0:1]
- ▶ ページ数: 3-6 ページ程度 (ソースコードは必要ない)

Chip Time を抜き出すサンプルコード

```
# regular expression to read chiptime
re = /^\\d+\\s+F?\\d+\\s+(\\d{1,2}:\\d{2}:\\d{2})\\s+/  
  
ARGF.each_line do |line|  
  if re.match(line)  
    puts "#{$1}"  
  end  
end
```

課題 1 問 1 平均、標準偏差、中間値の計算

- ▶ 要約統計量: 結果を表にすればよい
- ▶ 分単位での計算 (秒まで含めた値とは少し異なる)
- ▶ "No Age" は "Sex Total" の値を使って分類 (男性 1 名)

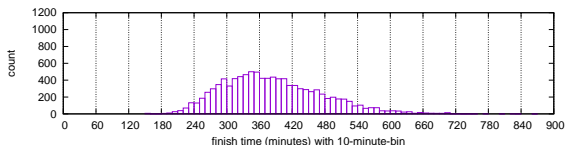
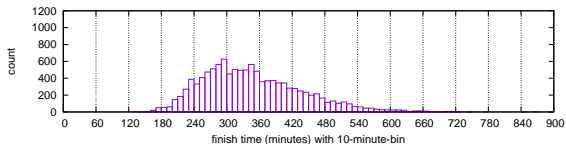
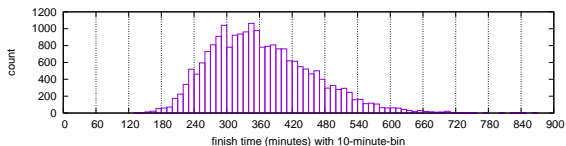
	n	mean	stddev	median
all	21,815	364.2	93.9	352
men	11,507	346.4	92.2	334
women	10,308	384.2	91.7	373

データ抽出用スクリプト例

```
# regular expression to read chiptime and category from honolulu2014.txt
re = /^\
```

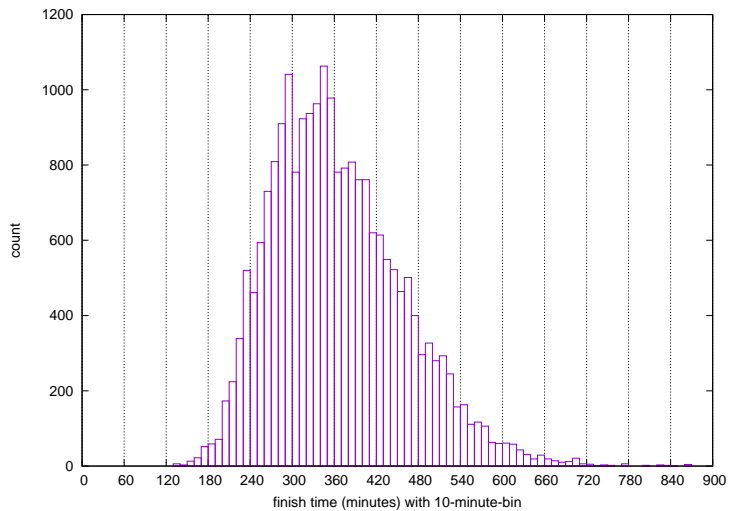
課題 1 問 2 完走時間のヒストグラム

- ▶ X 軸は 10 分ビンの完走時間、Y 軸は各ビンの完走者数
- ▶ 3 つのヒストグラムを別々の図に書く
- ▶ ビン幅は 10 分にする
- ▶ 3 つのプロットは比較できるように目盛を合わせること

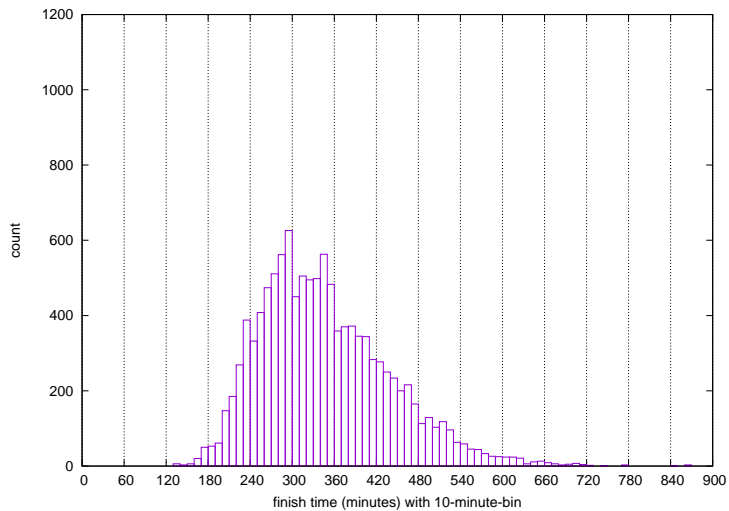


finish time histograms total(top) men(middle) women(bottom)

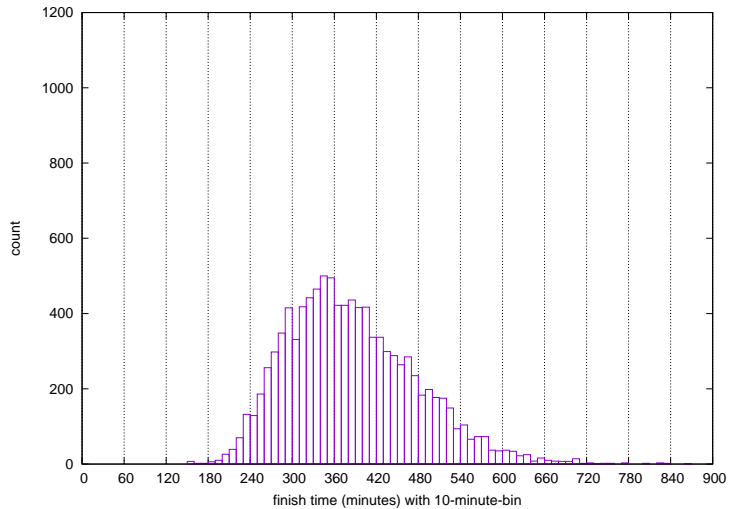
histograms for all



histograms for men

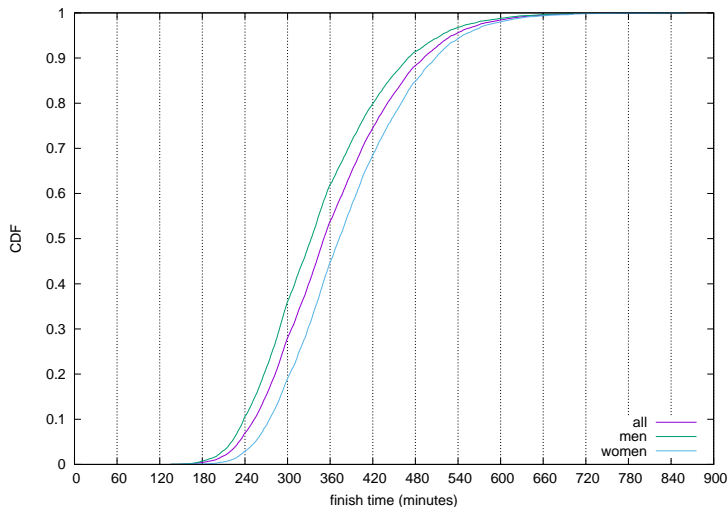


histograms for women



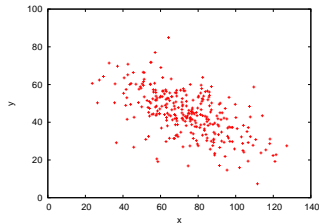
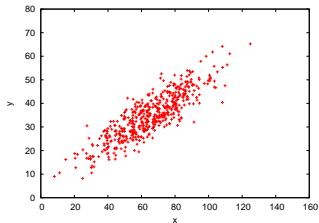
課題 1 問 3 CDF プロット

- ▶ ひとつの図に 3 つのプロットを書く
- ▶ X 軸は完走時間、Y 軸は CDF [0:1]



前回の演習: 相関係数の計算

- ▶ データの相関係数を計算する
 - ▶ correlation-data-1.txt, correlation-data-2.txt



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

演習: 相関係数の計算スクリプト

```
#!/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
cc = 0.0 # correlation coefficient
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

denom = (sum_xx - sum_x**2 / n) * (sum_yy - sum_y**2 / n)
if denom != 0.0
  cc = (sum_xy - sum_x * sum_y / n) / Math.sqrt(denom)
end

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

前回の演習 2: 類似度計算

- ▶ データの類似度を計算する
 - ▶ 「集合知プログラミング」2章の参考データ
 - ▶ 7人のユーザの映画評価スコア: 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
```

スコアデータ

- ▶ 簡単な例題: データが少なすぎる
- ▶ 一覧にすると以下のようなになる

```
#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 -
```


類似度計算の実行

- ▶ コサイン類似度を使ってユーザ間の類似度行列を作る

```
% 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
```

類似度計算スクリプト (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
```

類似度計算スクリプト (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
```

より本格的なデータセット

- ▶ MovieLens:

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

- ▶ ミネソタ大学が公開している協調フィルタリング用データセット
- ▶ ユーザの映画評価: 100K, 1M, 10M のスコアデータ
 - ▶ u.data: スコアデータ
 - ▶ 他にも各ユーザの属性情報や各映画の属性情報も含まれている

```
% 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
...
```

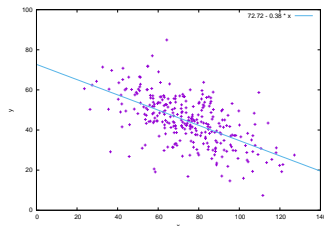
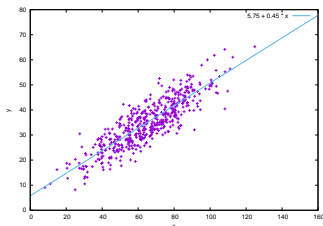
今回の演習: 線形回帰の計算

- ▶ 前回のデータを使い回帰直線を計算する
 - ▶ correlation-data-1.txt, correlation-data-2.txt

$$f(x) = b_0 + b_1x$$

$$b_1 = \frac{\sum xy - n\bar{x}\bar{y}}{\sum x^2 - n(\bar{x})^2}$$

$$b_0 = \bar{y} - b_1\bar{x}$$



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

演習: 回帰直線の計算スクリプト

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

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

sum_x = sum_y = sum_xx = sum_xy = 0.0
n = 0
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_xy += x * y
    n += 1
  end
end

mean_x = Float(sum_x) / n
mean_y = Float(sum_y) / n
b1 = (sum_xy - n * mean_x * mean_y) / (sum_xx - n * mean_x**2)
b0 = mean_y - b1 * mean_x

printf "b0:%.3f b1:%.3f\n", b0, b1
```

演習: 散布図に回帰直線を加える

```
set xrange [0:160]
set yrange [0:80]

set xlabel "x"
set ylabel "y"

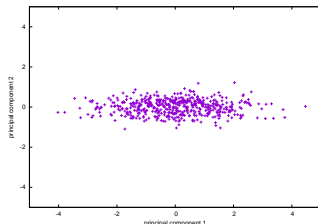
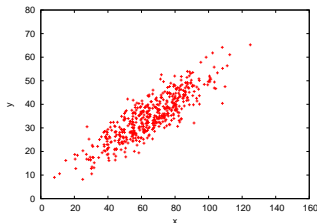
plot "correlation-data-1.txt" notitle with points, \
5.75 + 0.45 * x lt 3
```

今回の演習 2: 主成分分析

- ▶ 主成分分析: 線形回帰に使ったデータで実行

```
$ ruby pca.rb correlation-data-1.txt
PC1:
eigenval: 1.86477
proportion: 0.93239
cumulative proportion: 0.93239
eigenvector: [0.7071067811865475, 0.7071067811865475]

PC2:
eigenval: 0.13523
proportion: 0.06761
cumulative proportion: 1.00000
eigenvector: [0.7071067811865475, -0.7071067811865475]
```



data-1:r=0.87 (left), pca plot (right)

主成分分析: 3 変数の場合

```
$ cat pca-data.txt
7 4 3
4 1 8
6 3 5
8 6 1
8 5 7
7 2 9
5 3 3
9 5 8
7 4 5
8 2 2
$ ruby pca.rb -p pca-data.txt
-0.542660  0.664959  0.035324
2.803897  -0.066207  0.348792
0.615631   0.306325   0.165059
-2.158526  0.958839   0.386086
-0.931052  -1.044819   0.360013
1.142388  -1.273946   0.471245
0.803082   1.261879   0.472342
-1.246820  -1.655638  -0.023007
-0.286027  -0.024512   0.186799
-0.199912  0.873118   -1.460164

$ ruby pca.rb pca-data.txt
PC1:
eigenval: 1.76877
proportion: 0.58959
cumulative proportion: 0.58959
eigenvector: [-0.642004576349, -0.686361641360, 0.341669169247]

PC2:
eigenval: 0.92708
proportion: 0.30903
cumulative proportion: 0.89862
eigenvector: [-0.384672291688, -0.0971303301343, -0.917928606687]

PC3:
eigenval: 0.30415
proportion: 0.10138
cumulative proportion: 1.00000
eigenvector: [-0.663217424343, 0.720745028589, 0.20166618906]
```

演習: PCA スクリプト (1/4)

```
#!/usr/bin/env ruby
#
# usage: pca.rb [-p] datafile
# input datafile: row: variables, column: observations
# -p: convert input into principal components

# use matrix class for eigen vector computation
require 'matrix'
require 'optparse'

# normalize an array of array (m x n) into bb (m x n)
def normalize_matrix(aa)
  m = aa[0].length
  n = aa.length
  bb = Array.new(n) { Array.new(m) } # normalized array of array

  for i in (0 .. m - 1)
    sum = 0.0 # sum of data
    sqsum = 0.0 # sum of squares
    for j in (0 .. n - 1)
      v = aa[j][i]
      sum += v
      sqsum += v**2
    end
    mean = sum / n
    stddev = Math.sqrt(sqsum / n - mean**2)
    for j in (0 .. n - 1)
      bb[j][i] = (aa[j][i] - mean) / stddev # normalize
    end
  end
  bb # return the new array of array
end
```

演習: PCA スクリプト (2/4)

```
# make correlation matrix from data (array of array)
def make_corr_matrix(aa)
  m = aa[0].length
  n = aa.length
  corr_matrix = Array.new(m) { Array.new(m) }

  for i in (0 .. m - 1)
    for j in (0 .. m - 1)
      sum_x = 0.0
      sum_y = 0.0
      sum_xx = 0.0
      sum_yy = 0.0
      sum_xy = 0.0
      for k in (0 .. n - 1)
        x = aa[k][i]
        y = aa[k][j]
        sum_x += x
        sum_y += y
        sum_xx += x**2
        sum_yy += y**2
        sum_xy += x*y
      end
      cc = 0.0
      denom = (sum_xx - sum_x**2 / n) * (sum_yy - sum_y**2 / n)
      if denom != 0.0
        cc = (sum_xy - sum_x * sum_y / n) / Math.sqrt(denom)
      end
      corr_matrix[i][j] = cc
    end
  end
  corr_matrix
end
```

演習: PCA スクリプト (3/4)

```
do_projection = false
OptionParser.new {|opt|
  opt.on('-p') {|v| do_projection = true}
  opt.parse!(ARGV)
}

# read data into input (array of array)
input = Array.new
ARGF.each_line do |line|
  values = line.split
  if values.length > 0
    row = Array.new
    values.each do |v|
      row.push v.to_f
    end
    input.push row
  end
end

corr_aa = make_corr_matrix(input) # create correlation matrix
corrmatrix = Matrix.rows(corr_aa) # convert array of array into matrix class

# compute the eigenvalues and eigenvectors
# eigensystem returns v: eigenvectors, d: diagonal matrix of eigenvalues,
# v_inv: transposed matrix of v. corrmatrix = v * d * v_inv
v, d, v_inv = corrmatrix.eigensystem

# returned vectors are sorted in increasing order of eigenvals.
# so, re-order eigenvalues and eigenvectors in decreasing order
eigenvals = Array.new # array of eigenvalues
(d.column_size - 1).downto(0) do |i|
  eigenvals.push d[i,i]
end
eigenvectors = Matrix.columns(v.column_vectors.reverse)
```

演習: PCA スクリプト (4/4)

```
if do_projection != true
  # show summaries
  eig_sum = 0.0
  eigenvals.each do |val|
    eig_sum += val
  end
  cum = 0.0 # cumulative of eigenvalues
  eigenvals.each_with_index do |val, i|
    printf "PC%d:\n", i + 1
    printf "eigenval: %.5f\n", val
    printf "proportion: %.5f\n", val / eig_sum
    cum += val
    printf "cumulative proportion: %.5f\n", cum / eig_sum
    print "eigenvector: "
    print eigenvectors.column(i).to_a
    print "\n\n"
  end
else
  # project the input into new coordinate
  # first, normalize the input and then convert it to matrix
  normalized = Matrix.rows(normalize_matrix(input))
  # projected data = eigenvec.T x normalized.T
  projected = eigenvectors.transpose * normalized.transpose

  projected.column_vectors.each do |vec|
    vec.each do |v|
      printf "%.6f\t", v
    end
    print "\n"
  end
end
```

まとめ

第7回 多変量解析

- ▶ データセンシング
- ▶ 地理的位置情報 (geo-location)
- ▶ 線形回帰
- ▶ 主成分分析
- ▶ 演習: 線形回帰

次回予定

第 8 回 時系列データ (6/8)

- ▶ インターネットと時刻
- ▶ ネットワークタイムプロトコル
- ▶ 時系列解析
- ▶ 演習: 時系列解析
- ▶ 課題 2