

# インターネット計測とデータ解析 第6回

長 健二郎

2011年6月15日

# 前回のおさらい

## インターネットの構造を計る

- ▶ インターネットアーキテクチャ
- ▶ ネットワーク階層
- ▶ トポロジー
- ▶ グラフ理論
- ▶ 演習:トポロジ解析

# 今日のテーマ

## インターネットの特徴量を計る

- ▶ 遅延、パケットロス、ジッタ
- ▶ フロー計測
- ▶ 相関と多変量解析
- ▶ 主成分分析
- ▶ 演習:多変量と相関

# 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 +

## ▶ 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

# フローベースの計測

- ▶ SNMP によるインターフェイスカウンタ値による計測の限界
  - ▶ 総量は分かるが、それ以上の情報取得が困難
- ▶ フローベースの計測
  - ▶ フロー (5-tuple) 毎の統計情報
  - ▶ もともとは高速転送用のキャッシュ情報
  - ▶ プロトコル: NetFlow、sFlow、IPFIX、 ...
    - ▶ プロトコルバージョンや実装による違いも

# NetFlow の概要

- ▶ インターフェイス毎のキャッシュ情報を UDP でコレクタに送信
- ▶ パケットがインターフェイスに到着すると
  - ▶ 新規エントリを作成
  - ▶ または、既存のエントリをアップデート
    - ▶ バイトカウント、パケットカウント、エンドタイム、TCP フラグ (ORed)
  - ▶ エクスパイア条件 (4 種類) :
    - ▶ キャッシュがフル、TCP RST or FIN
    - ▶ 非アクティブフロー 15 秒、アクティブフロー 30 分
  - ▶ エクスパイアしたフローエントリはコレクタに送信される
- ▶ フロー情報
  - ▶ saddr, daddr, sport, dport, proto, ToS, input ifIndex byte count, packet count, start time, end time, output ifIndex TCP flags, next hop, src AS, dst AS

# フロー計測のサンプリング

情報量と負荷低減のために、 $N$  パケットに 1 回記録を取る機能

- ▶ 考慮すべき点

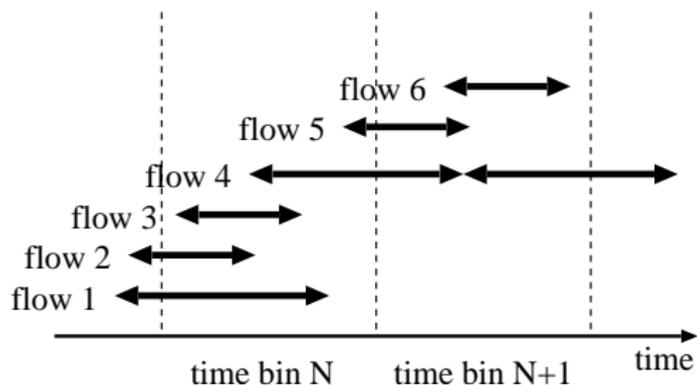
- ▶ ルータの負荷
- ▶ データ量
- ▶ コレクタの処理能力

- ▶ サンプリングの影響

- ▶ 測定結果は、測定値にサンプリング値の逆数を乗じて補正
- ▶ 使用量大きいフローはいいが、小さいフローは精度がでない
- ▶ 例：サンプリング値:1/100, 100 ユーザがそれぞれ 1KB パケットを 1 個送った
  - ▶ 測定結果: 100KB を送ったユーザが 1 人いると誤認
- ▶ 必要な精度に応じたサンプリング値の設定が必要
- ▶ 実際には、サンプリング値による精度の限界を理解して解析

## 時間粒度

- ▶ アクティブなフロー情報は 30 分に 1 度しかエクスポートされない
  - ▶ 単位時間 (ビンサイズ) は小さく出来ない
- ▶ 簡単のためエンドタイムでカウント
  - ▶ より正確にはスタートタイムも使い比例割り当て

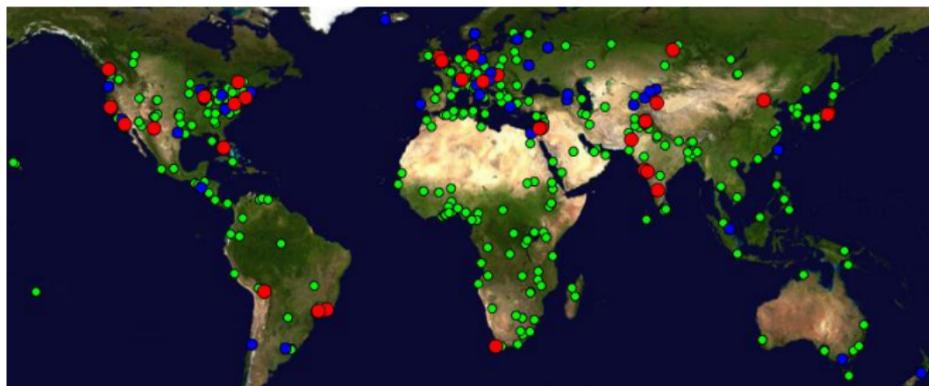


## pingER project

- ▶ 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 project monitoring sites in east asia

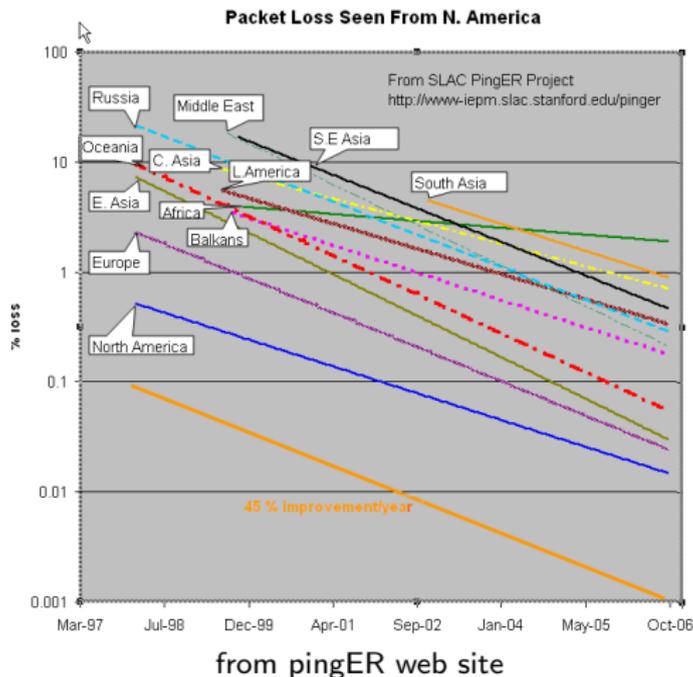
- ▶ monitoring (red) and remote (green) sites



from pingER web site

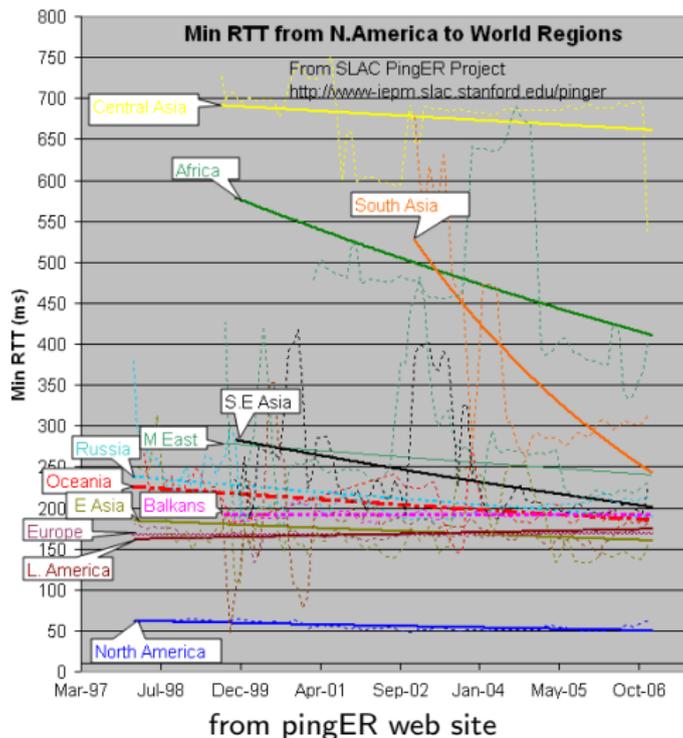
# pingER packet loss

- ▶ packet loss observed from N. America
- ▶ exponential improvement in 10 years



## pinger minimum rtt

- ▶ minimum rtt's observed from N. America
- ▶ gradual shift from satellite to fiber in S. Asia and Africa

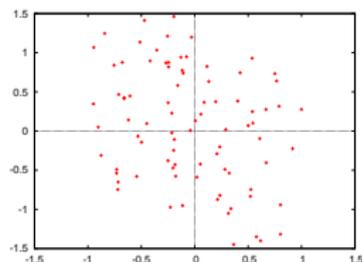
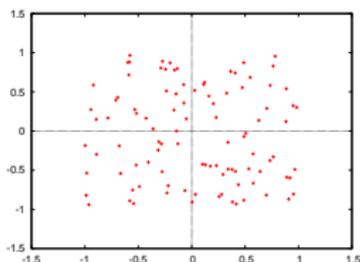
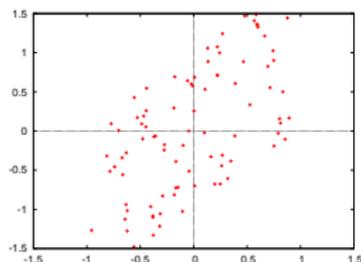


# 多変量データ解析

- ▶ 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)

# 散布図 (scatter plots)

- ▶ 2つの変数の関係を見るのに有効
  - ▶ X軸: 変数 X
  - ▶ Y軸: それに対応する変数 Y の値
- ▶ 散布図で分かる事
  - ▶ XとYに関連があるか
    - ▶ 無相関、正の相関、負の相関
  - ▶ 外れ値の存在があるか



例: (左) 正の相関 0.7 (中) 無相関 0.0 (右) 負の相関 -0.5

## 相関 (correlation)

- ▶ 共分散 (covariance):

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

- ▶ 相関係数 (correlation coefficient):

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

- ▶ 相関係数は共分散を正規化したもの。 -1 から 1 の値を取る。
- ▶ 相関係数は外れ値の影響を大きく受ける。 散布図と併用し、外れ値を確認する必要。
- ▶ 相関関係と因果関係
  - ▶ 相関関係が因果関係を示すとは限らない。
    - ▶ 未知の第 3 の共通の要因が存在する場合
    - ▶ 単なる偶然

# 相関と多変量解析

多変量解析: 互いに関係する複数の変数からなるデータを統計的に扱う手法

- ▶ 関係の視覚化
  - ▶ クラスタ分析: 変量間の距離 (類似度) を計算し、グループ (クラスタ) に分ける
- ▶ 次元減少
  - ▶ 主成分分析: 変数を減らす

# 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

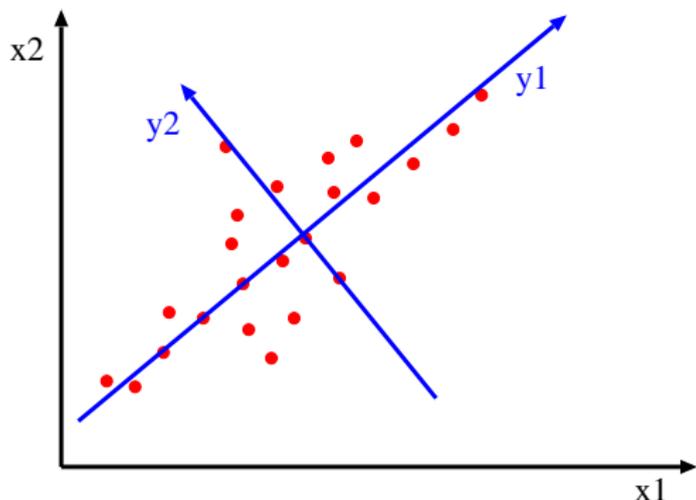
- ▶ あくまで、ばらつきの大きい成分を抜き出すだけ
  - ▶ とくに各軸の単位が違う場合は注意
- ▶ 機械的に複雑な関係を分析できる便利な手法であるが、それで複雑な関係が説明できる訳ではない

## 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まる  
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  $\text{cov}(Y)$  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

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

したがって

$$P \text{cov}(Y) = PP^T \text{cov}(X)P = \text{cov}(X)P$$

P を  $d \times 1$  行列でかくと、

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

また、 $\text{cov}(Y)$  は対角行列 (各変数が独立) なので

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

書き直すと

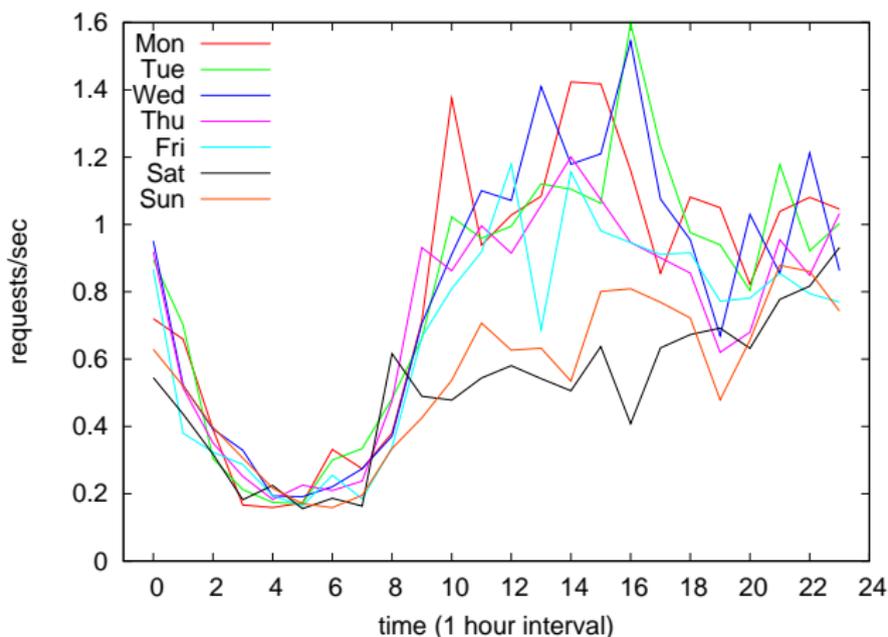
$$[\lambda_1 P_1, \lambda_2 P_2, \dots, \lambda_d P_d] = [\text{cov}(X)P_1, \text{cov}(X)P_2, \dots, \text{cov}(X)P_d]$$

$\lambda_i P_i = \text{cov}(X)P_i$  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

## 演習: 多変量と相関

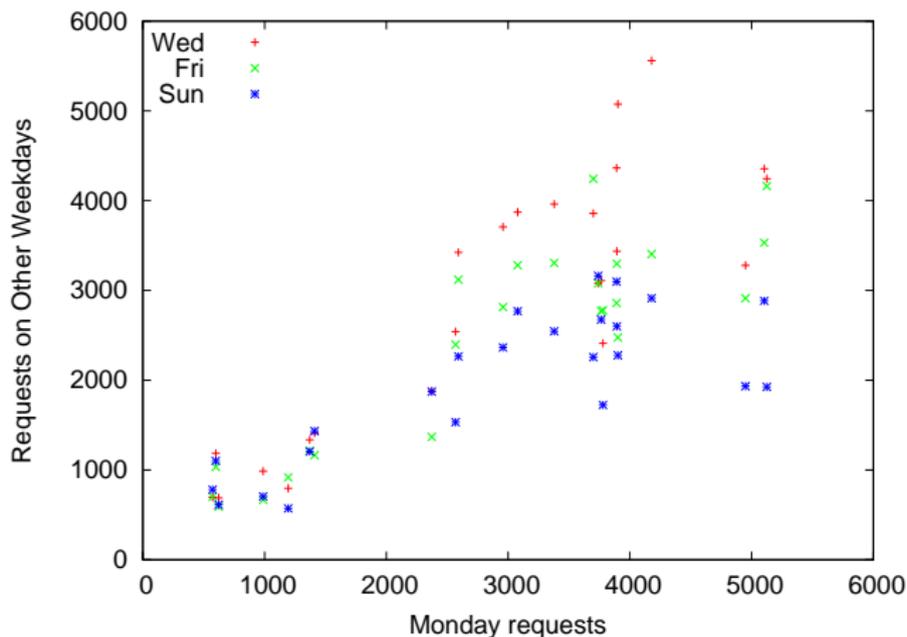
課題 1 の例: 各曜日の時間別リクエスト数プロット  
(2011-02-28/2011-03-06)

- ▶ 各曜日間には相関があり、平日と土日は少し傾向に違い
- ▶ 各曜日間の相関係数を求める



# 散布図による確認

月曜に対する、水、金、日曜の同時間帯別アクセス数散布図にする



# 共分散行列

共分散行列を計算する

	Mon	Tue	Wed	Thu	Fri	Sat	Sun
Mon	2052349	1732910	1766814	1481337	1438849	753142	912707
Tue	1732910	1853635	1791264	1387485	1345741	700134	956308
Wed	1766814	1791264	2047528	1475368	1408105	660976	990202
Thu	1481337	1387485	1475368	1325782	1182891	637794	768738
Fri	1438849	1345741	1408105	1182891	1264010	585177	737702
Sat	753142	700134	660976	637794	585177	573336	492188
Sun	912707	956308	990202	768738	737702	492188	642135

相関係数行列を計算する

	Mon	Tue	Wed	Thu	Fri	Sat	Sun
Mon	1.000	0.888	0.862	0.898	0.893	0.694	0.795
Tue	0.888	1.000	0.919	0.885	0.879	0.679	0.877
Wed	0.862	0.919	1.000	0.895	0.875	0.610	0.864
Thu	0.898	0.885	0.895	1.000	0.914	0.732	0.833
Fri	0.893	0.879	0.875	0.914	1.000	0.687	0.819
Sat	0.694	0.679	0.610	0.732	0.687	1.000	0.811
Sun	0.795	0.877	0.864	0.833	0.819	0.811	1.000

# まとめ

## インターネットの特徴量を計る

- ▶ 遅延、パケットロス、ジッタ
- ▶ フロー計測
- ▶ 相関と多変量解析
- ▶ 主成分分析
- ▶ 演習:多変量と相関

# 次回予定

## 第7回 インターネットの多様性と複雑さを計る (6/22)

- ▶ サンプルング
- ▶ 統計解析 (ヒストグラム、大数の法則)
- ▶ 演習: ヒストグラム、CDF
- ▶ 課題 2