

インターネット計測とデータ解析 第5回

長 健二郎

2010年10月27日

前回のおさらい

インターネットの速度を計る

- ▶ 速度計測
- ▶ 利用可能帯域の推測
- ▶ 平均 標準偏差
- ▶ 線形回帰
- ▶ 課題 1

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
- ▶ グラフによる可視化
- ▶ 相関と多変量解析

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

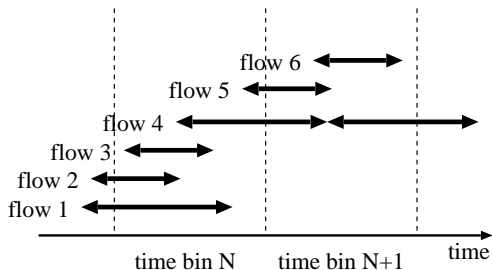
- ▶ ルータの負荷
- ▶ データ量
- ▶ コレクタの処理能力

- ▶ 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
- ▶ 使用量が多いフローはいいが、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

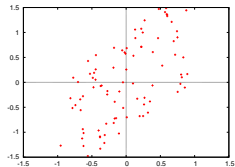
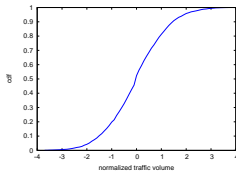
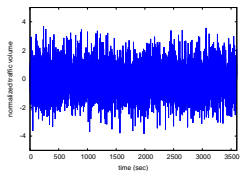
時間粒度

- ▶ 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



グラフ描画

直観的にデータの性質を把握するには、いくつかの統計的手法を用いてグラフを描画してみる



グラフ描画のガイドライン

読み手の立場にたって、分かり易いグラフを描画する

- ▶ XY 軸のラベルを明確に
- ▶ XY 軸の目盛りと単位を明確に
- ▶ 個々の直線曲線にもラベルを付ける
- ▶ 適切なフォントとサイズを使う
- ▶ 慣習に従う: 0 を起点にする、数学シンボルや略称の使用など
- ▶ ばらつきを示す (平均値だけでは不十分)
- ▶ グラフの範囲を適切か
- ▶ ひとつのグラフで多くを示さない
- ▶ 異なるデータを比較する場合は、適切な正規化を行う
- ▶ グラフ同士を比較する場合は、XY 軸のスケールを合わせる
- ▶ 技術系は円グラフや 3D 効果グラフは使わない
- ▶ 色を使う場合
 - ▶ 白黒印刷しても読めるように配慮
 - ▶ プロジェクタ投影も配慮 (例:黄色は避ける)

グラフ描画ツール

- ▶ gnuplot
 - ▶ コマンドラインツール、スクリプトで自動化し易い
 - ▶ <http://gnuplot.info/>
- ▶ grace
 - ▶ 使い易い GUI
 - ▶ 細かい仕上げ調整が可能
 - ▶ <http://plasma-gate.weizmann.ac.il/Grace/>

データ変数

- ▶ 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

生データのグラフ化

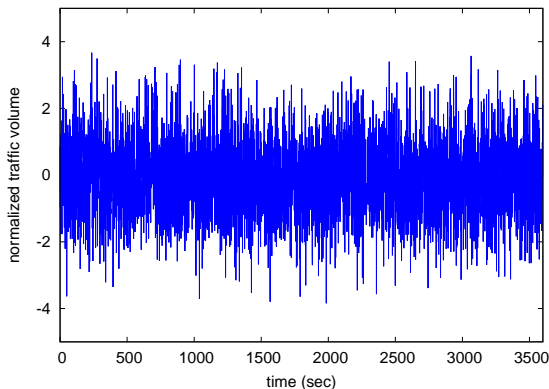
代表的なグラフ

- ▶ 時系列グラフ
- ▶ ヒストグラム
- ▶ 確率グラフ
- ▶ 散布図

時系列グラフ

変数の時間変化を見る

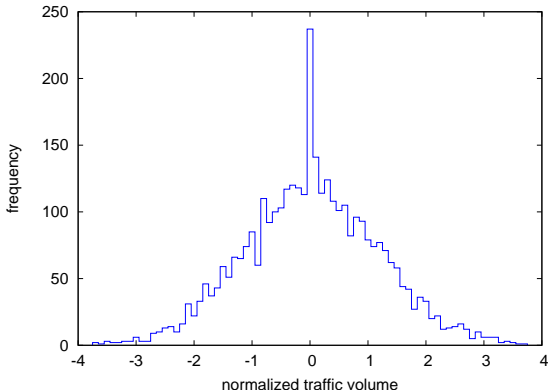
- ▶ X軸に時間、Y軸に変数値
- ▶ 時系列グラフから分かること
 - ▶ 位置の変化
 - ▶ ばらつきの変化
 - ▶ 外れ値の存在



ヒストグラム (1/2)

変数の分布の仕方を見る

- ▶ データを同じ幅のビンに分ける
- ▶ 各ビンのデータ数を数える
- ▶ X軸:ビンの値 Y軸:データ数



ヒストグラム (2/2)

ヒストグラムから分かる事

- ▶ 分布の中心 (位置)
- ▶ 分布の広がり
- ▶ 分布の偏り
- ▶ 外れ値の存在
- ▶ 複数のモードの存在 (山が複数あるか)

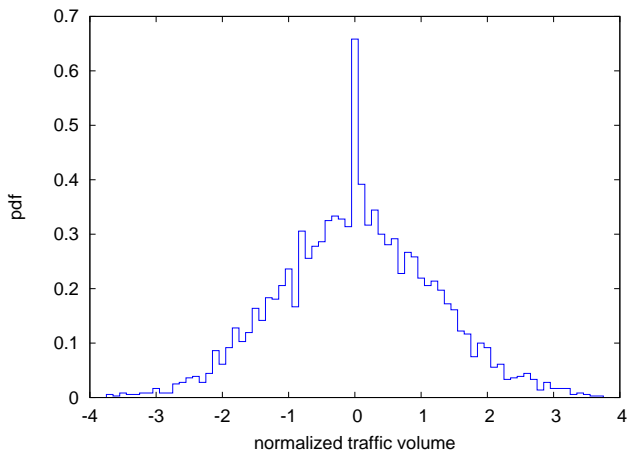
ヒストグラムの制約

- ▶ 適切なビン幅を選ぶ必要
 - ▶ 小さすぎると各ビンのサンプル数が足りなくなる
 - ▶ 大きすぎると分布の詳細が分からない
 - ▶ 偏りの大きい分布では適切なビン幅の選択は難しい
- ▶ 十分なサンプル数が必要

確率密度関数 (probability density function; pdf)

- ▶ 合計面積が1となるように出現数を正規化
 - ▶ 出現数を (総データ数 × ビン幅) で割る
- ▶ 確率密度関数: 確率変数 X が x という値をとる確率

$$f(x) = P[X = x]$$



累積分布関数 (cumulative distribution function; cdf)

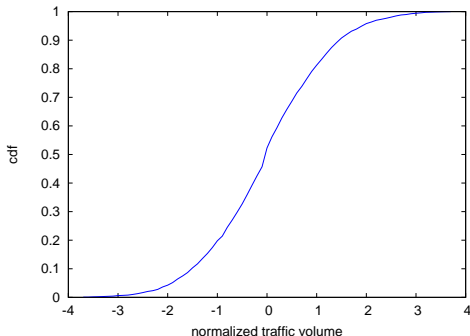
- ▶ 密度関数: x をいう値を観測する確率

$$f(x) = P[X = x]$$

- ▶ 累積分布関数: x 以下の値を観測する確率

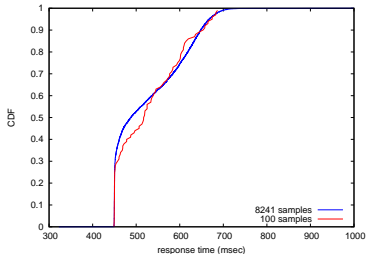
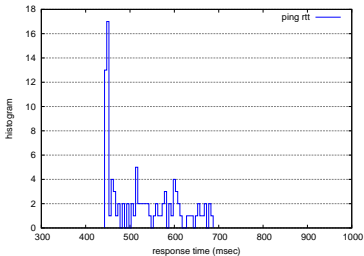
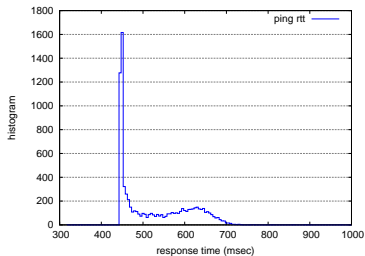
$$F(x) = P[X \leq x]$$

- ▶ 分布の偏りが大きい、サンプル数が少ない、外れ値が無視できない場合などは、ヒストグラムより有効



ヒストグラムとCDFの比較

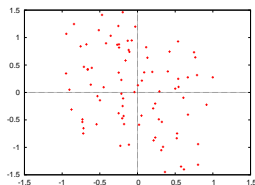
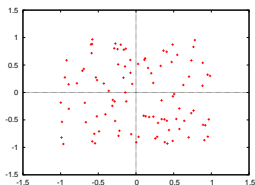
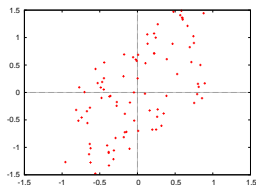
- ▶ CDFの場合、ビン幅やサンプル数不足を考慮しなくていい



(左) 元データ (右)100 サンプル (下)CDF

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

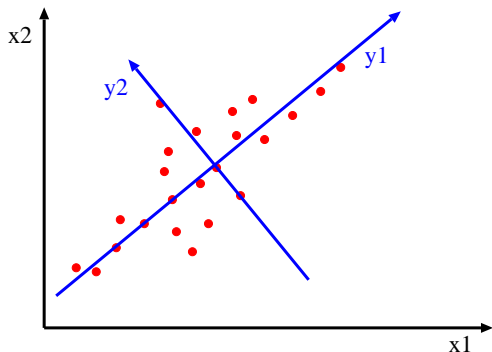
- ▶ あくまで、ばらつきの大きい成分を抜き出すだけ
 - ▶ とくに各軸の単位が違う場合は注意
- ▶ 機械的に複雑な関係を分析できる便利な手法であるが、それで複雑な関係が説明できる訳ではない

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

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

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

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
- ▶ グラフによる可視化
- ▶ 相関と多変量解析

次回予定

第6回 インターネットの多様性と複雑さを計る (11/10)

- ▶ ロングテールとさまざまな分布
- ▶ サンプリング
- ▶ 統計解析
 - ▶ 期待値と大数の法則、検定と信頼区間