

# SQC勉強会(3)

## ヒストグラムと代表値

1

QC勉強会資料

## 本の宣伝

- この勉強会の元本です
- サンプルデータもこの本からとっています
- 著者の先生もいい人です
- わかりやすくコンパクトな本でお勧めです  
(講演を聞いたのですが、人のよさそうなひとでした)
- **フリーソフトウェアRによる統計的品質管理入門**
- 編著：荒木孝治
- 出版社：(株)日科技連出版社



2

QC勉強会資料

## 本日の内容

- ヒストグラム
- 代表値と要約統計量
- 分布(ヒストグラムからモデルへ)

QC勉強会資料

## 品質の管理とは？

- 品質：「(モノの)利用における(効用の)適合度」  
(Montgomery, 2005)
- 4要素
  - 対象
    - 何(部門、工程)を対象とするのか
  - 特性
    - 対象のどのような性質を管理したいのか
  - **標準(水準)**
    - 許容できる特性の測定量はどの範囲か？
  - **要因**
    - 特性の計量値に影響を与える要素は何か？
- SQCでは特に標準の決め方と活動データと標準との関係の判断に関する手法を重視

QC勉強会資料

## QC7つ道具とは

2010年1月14日

- 定量データを基にしてSQCを行う際の基本的ツール
  - 視覚的に
    - 課題の整理
    - 問題要因の特定・標準の設定
    - 対策実施後の品質の判断
- ※新QC7つ道具:主に定性的なものに対して適用

QC勉強会資料

## QC7つ道具とは？

2010年1月14日

- グラフ・管理図:データの傾向を知る 勉強会で扱いません
- チェックシート:手順のチェック
- パレート図:重点対策する要因を決める
- ヒストグラム:データの特徴を知る
- 特性要因図:バラツキの要因の解析
- 散布図:要因間の関係を見る
- 層別:要因の存在を調べる

QC勉強会資料

## テーマ:「代表値」をどのように決めるのか

2010年1月14日

- 代表値:「グループ」(群)を代表する値
- グループの全体的な様子を1つの値に「代表」させる
- よく使うのは「平均値」
- しかし、平均値を用いるとグループを代表しない場合も・・・
- いくつかの代表値を学習し、適切なものを選択する
- **鉄則:事前にヒストグラムを描いて判断する**

QC勉強会資料

## ヒストグラム

2010年1月14日

- ヒストグラム
  - タテに値(区間)、ヨコに発生回数(頻度)をとった棒グラフ
- あるシステムの1日あたりのアクセス回数のデータ(架空)
- Access\_num.csv

QC勉強会資料

## ヒストグラム

2010年1月14日

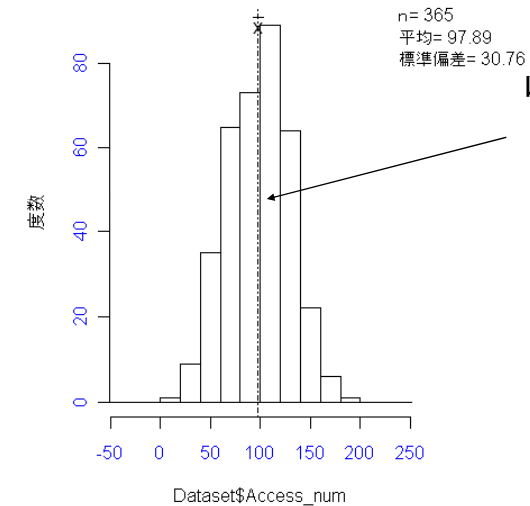
- データの読み込み
- 「データ」→「データのインポート」
- Macの人はクリップボードからコピー&ペーストしたほうがよい(らしいです)
- 「QCツール」→「QCヒストグラム」

QC勉強会資料

## ヒストグラム

2010年1月14日

Dataset\$Access\_numのヒストグラム



QC勉強会資料

## 本当に平均値は代表値？

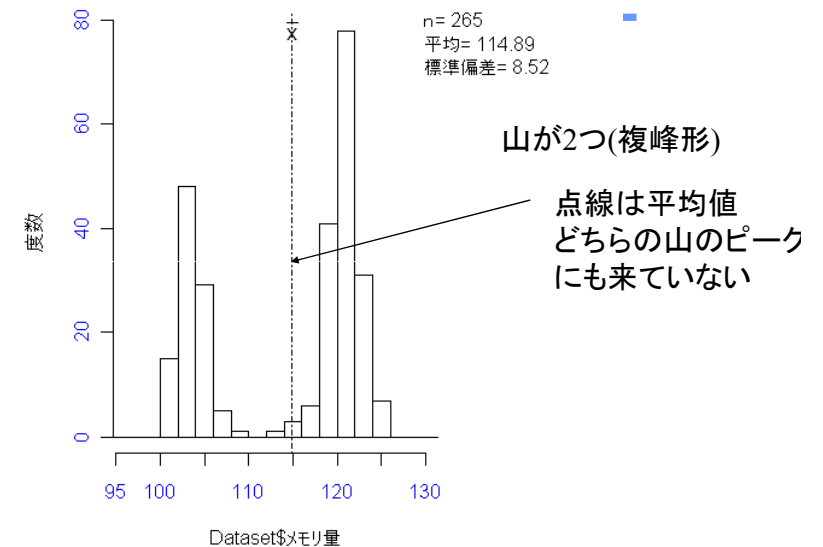
2010年1月14日

- 現実の場合は必ずしもそうでない場合も
- あるシステム(サーバ2台構成)のメモリ使用量の推移
- 全体のヒストグラムをしてみる

QC勉強会資料

Dataset\$メモリ量のヒストグラム

2010年1月14日



QC勉強会資料

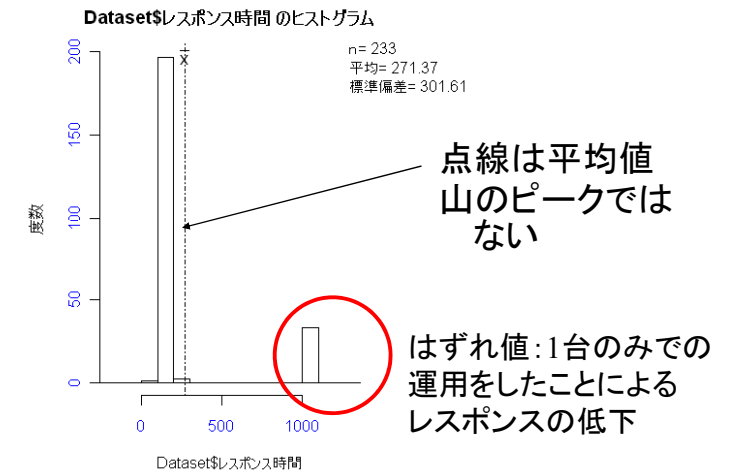
## 本当に平均値は代表値？

2010年1月14日

- あるシステムのレスポンス時間(ミリ秒)
- 2台構成のうち、途中で1台がハード故障
- その間は1台のみで運用
- response.csv

QC勉強会資料

2010年1月14日



QC勉強会資料

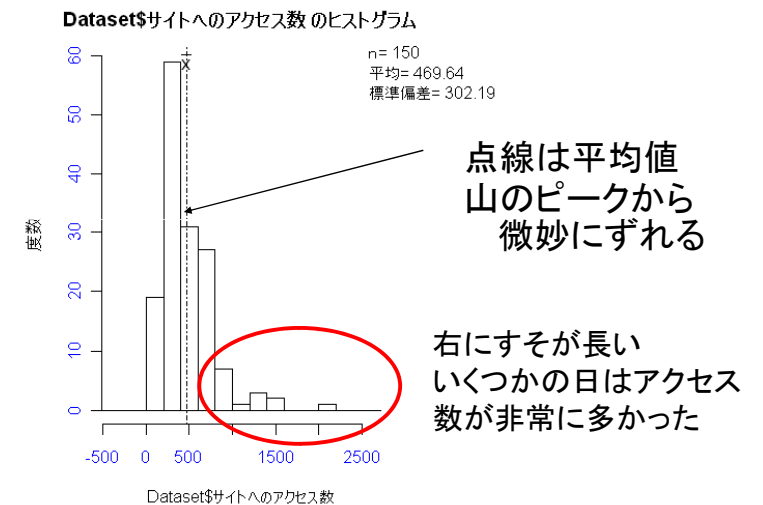
## 本当に平均値は代表値？

2010年1月14日

- あるECサイトの1日あたりの来訪者
- Access\_num2.csv

QC勉強会資料

2010年1月14日



QC勉強会資料

## 本当に平均値は代表値？

2010年1月14日

- ヒストグラムを見て、
  - 単峰: 山が1つ
  - 対称: 左右が対象

の場合は平均値と山のピークは一致  
平均値 = 代表値

- それ以外のときは別の手立てを考える

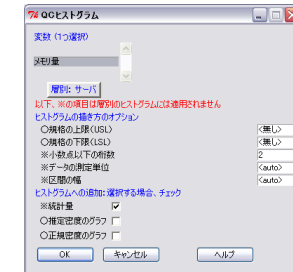
QC勉強会資料

## QC(データ解析)でよく用いられる手法

2010年1月14日

### その1: 層別

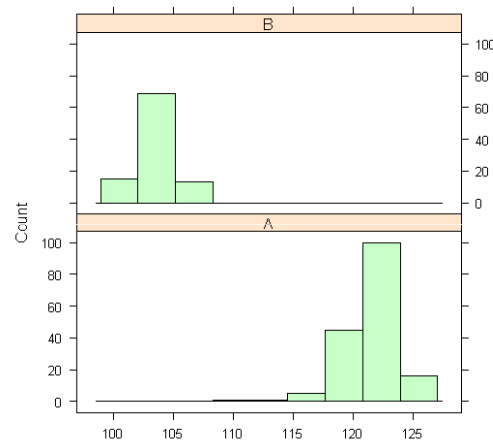
- memory.csvを再度見てみる
- サーバ情報がデータに存在
- サーバによってメモリ量に変化はないか？
- 層別の実施



QC勉強会資料

## 層別のヒストグラム

2010年1月14日



- 層別の結果
- 2つの山の原因はサーバの種類によるもの
- サーバごとに平均値を出し、代表値とみなすのが合理的

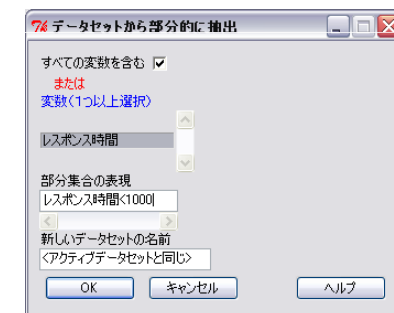
QC勉強会資料

## QC(データ解析)でよく用いられる手法

2010年1月14日

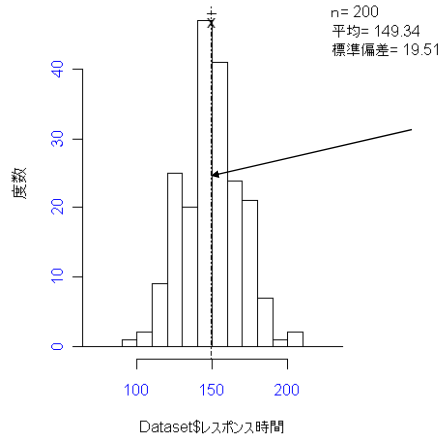
### その2: はずれ値を取り除く

- response.csvを見てみる
- 明らかに1000ミリ秒を超えたものは少し異常？
- システム設計上そうだった場合はこれをはずれ値とみなして取り除く



QC勉強会資料

Dataset\$レスポンス時間のヒストグラム



点線は平均値  
山のピークとほぼ  
一致

## その3: 他の代表値を使う

頑健性  
(ロバスト性の形に影響を受けない)

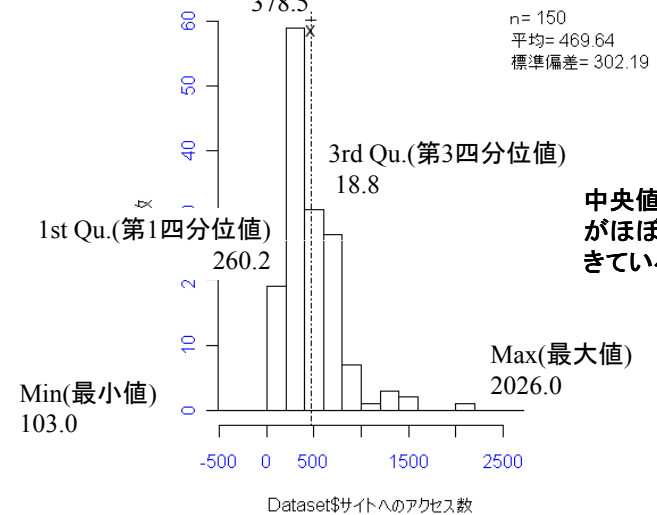
- 困ったらまず **中央値** を見てみる
- **平均値**
  - (すべての観測値の和) / (観測数)
- **中央値(メディアン)**
  - すべての観測値を昇順(降順)に並べたときの真ん中の値
- **最頻値(モード)**
  - すべての観測値でもっとも発生回数が高い値

統計的に扱いやすい

## 中央値の適用

- Access\_num2.csvをしてみる
- 「統計量」→「要約」→「アクティブデータセット」で色々な代表値が出る

Dataset Median(中央値) ヒストグラム



中央値  
がほぼ山のピークに  
きている

## 山のピーク=最頻値ではないの？

2010年1月14日

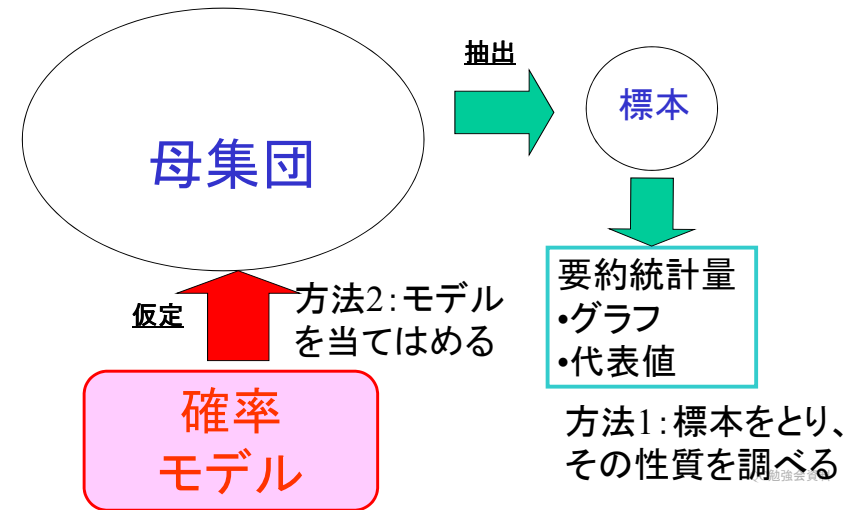
- 概略的にはそのとおりです。
- しかし、通常は
  - 観測データ(連続量)の場合、全く同じ値が2度でることほとんどない
  - データをある位で丸めると恣意性が入ってしまう
  - 後で出てくる検定(統計分析)が適用しづらいため、中央値までの適用にとどめます。

QC勉強会資料

## ヒストグラムからモデルへ

2010年1月14日

母集団の性質を理解するためにどうするのか？



## 確率モデルとは？

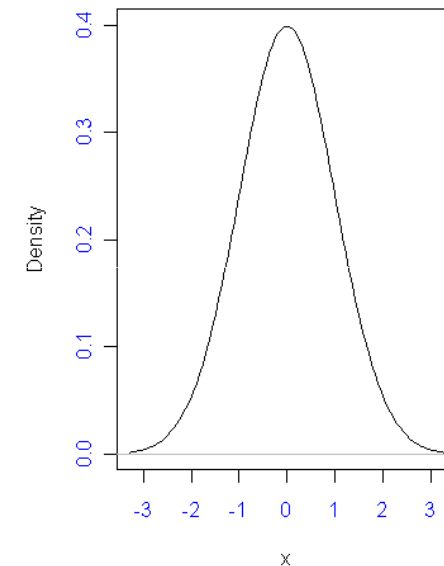
2010年1月14日

- 確率
  - ある事象(値)が発生する割合
- 確率密度関数
  - ヨコに事象(値)、タテに確率(割合)を採ったグラフ
- 「分布」→「正規分布」→「正規分布を描く」

QC勉強会資料

Normal Distribution:  $\mu = 0, \sigma = 1$

2010年1月14日



会資料

## ヒストグラムと確率密度関数

2010年1月14日

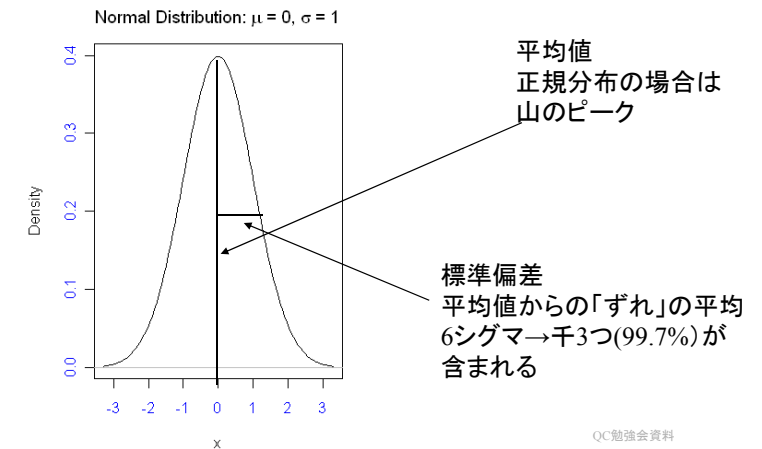
- ヒストグラム
  - ある事象(値)に対する**実際の発生回数(頻度)**
- ヒストグラムの発生回数を多くし、割合に直すと確率密度関数になる  
(頻度論的な考え方)

QC勉強会資料

## 正規分布:重要な分布

2010年1月14日

- 正規分布:  $N(\text{平均値}, \text{標準偏差})$



## 次回予告

2010年1月14日

- 正規分布を覚えて活用しよう
  - シックスシグマ(の概論)
  - 管理図

QC勉強会資料