

情報処理応用B 第8回

藤田 一寿

問題解決と統計処理

■ 分析と改善

- 問題が起こったときにどう解決すればよいのか.
- 現状を分析し改善する.
- 現状分析の手法
 - 統計処理
- 改善の手法
 - PDCAサイクル

PDCAサイクル

■ PDCAサイクル

- 改善のための4つの手順(PDCAサイクル)

- ① 計画(P: Plan)

- 目的を決め、達成に必要な計画を設定

- ② 実行(D: Do)

- 計画に基づき実施

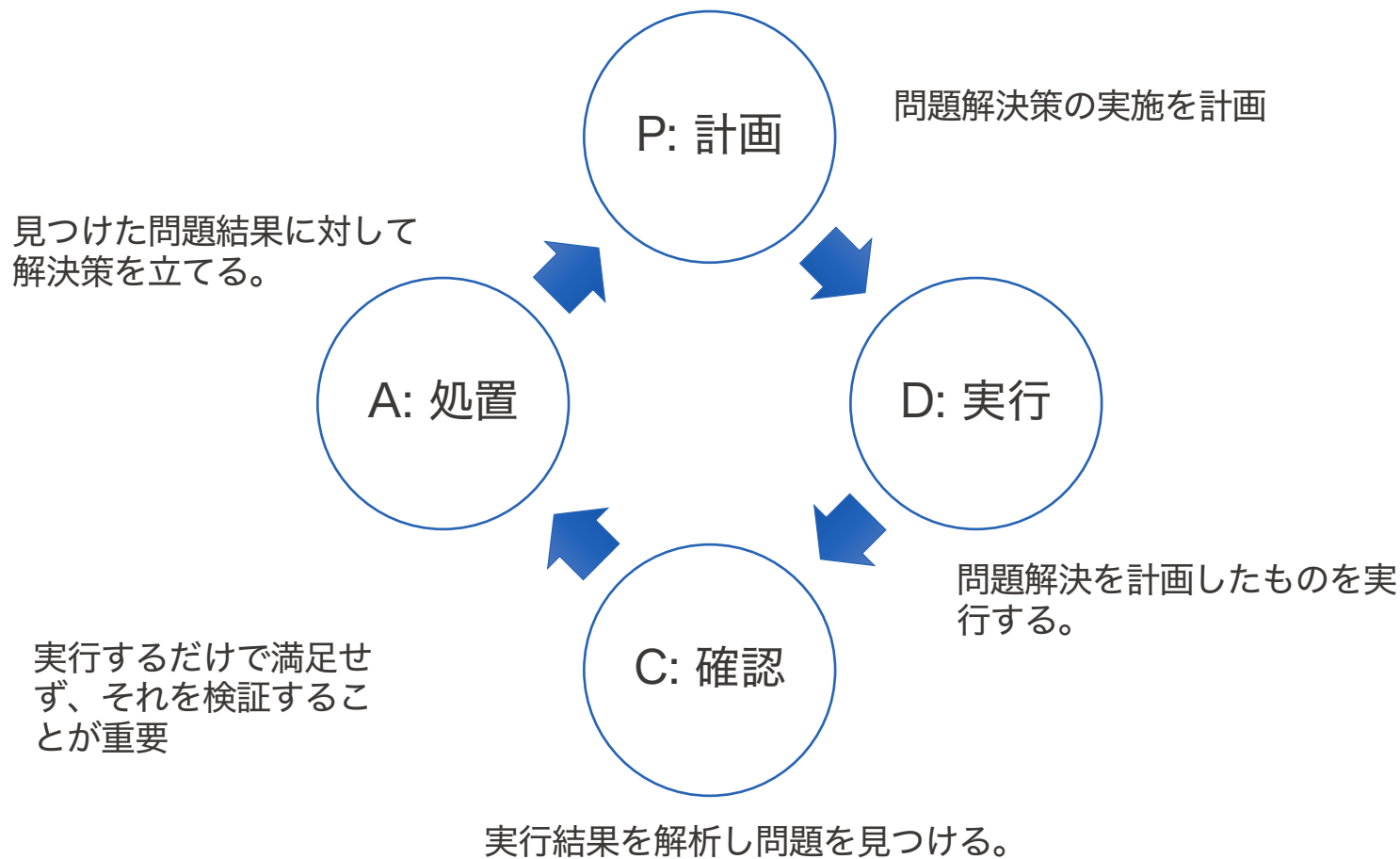
- ③ 確認(C: Check)

- 実施の結果を調べ評価

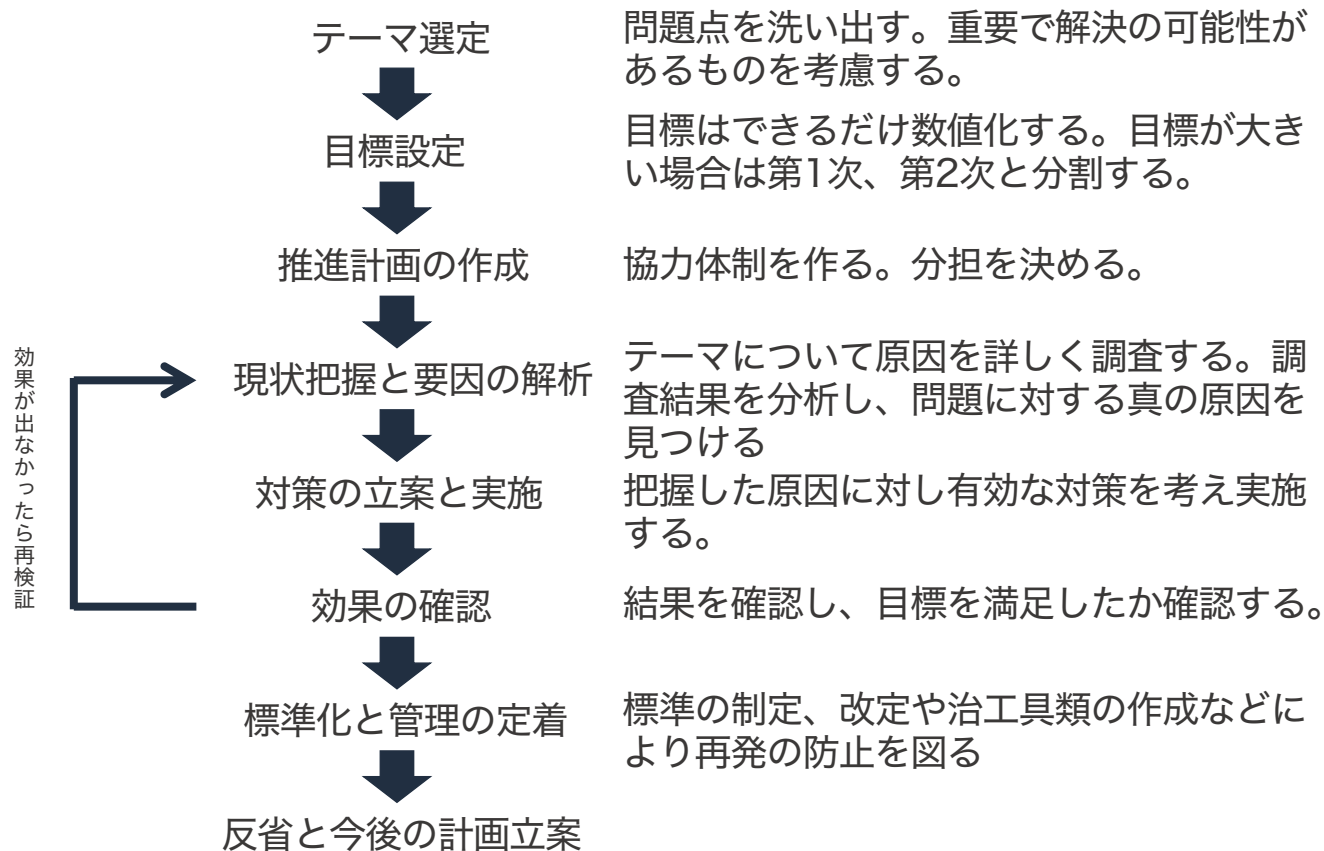
- ④ 処置(A: Action) 処置

- 必要により適切な処置

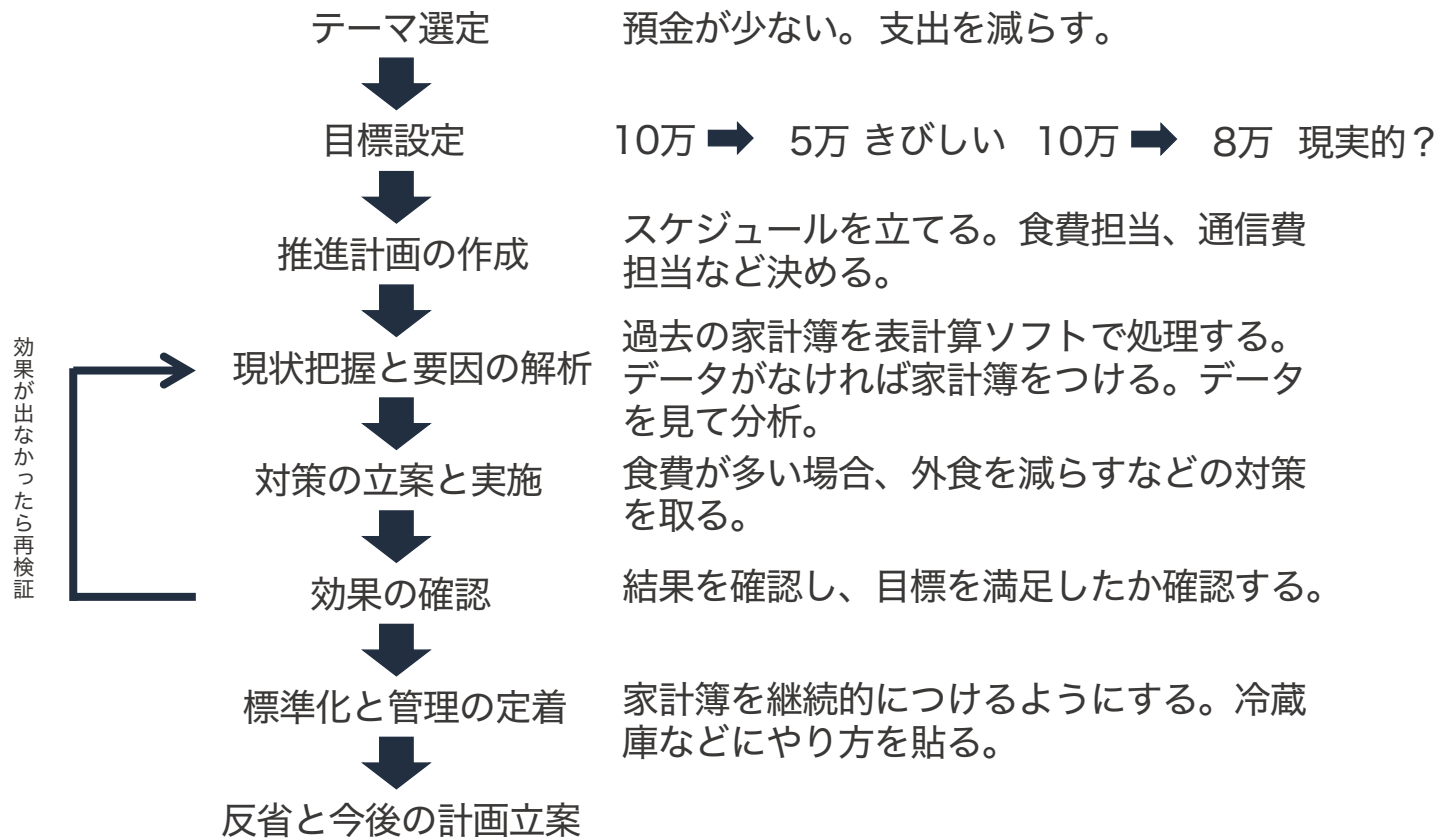
■ PDCAサイクル



■ 改善のアプローチ



■ 家計支出の例



統計の基礎

■ 科学的アプローチ

- 問題を解決するには、経験や感だけではなく、データや理論を用いて現状を認識する必要がある。
- 数値化する→統計
- 可視化してわかりやすく→グラフ

■ 問題解決で統計やグラフが必要なわけ

- 現状を数値化し評価したい
- すべての製品をチェックできないが製造工程全体を評価したい
- 現状を可視化してわかりやすく理解したい
- 図で表すことで問題点を発見しやすくしたい
- など

■ なぜ統計が必要か

- 無数にあるデータ一つ一つをチェックできない
- 個々のデータをみるだけでは全体の傾向がつかめない



統計的手法を用いる

■ 統計で用いる基礎的な数値（統計量）

- データ数
- 最大、最小、中央値
- 平均
- 分散、標準偏差(ばらつき具合)

■ 最大値、最小値、中央値（メディアン）

- 最大値
 - 最も大きい値
- 最小値
 - 最も小さい値
- 中央値（メディアン）
 - 順番的に中央の値
 - 平均より中央値の方が適切な場合もある。平均は外れ値に引っ張られるため。

データが奇数個の場合



データが偶数個の場合



■ 平均（算術平均）

- N 個のデータ x_1, x_2, \dots, x_N の平均は下記のように表される。

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_N}{N}$$

- データの重心の意味も持つ。

■ 平均と中央値

- 平均と中央値どちらが集団の特徴をより表しているのか？
- 右の例
 - 平均点：6.85点
 - 中央値：8点
 - 平均点は0点と1点の人に引っ張られていて、中央値にくらべ低めの値になっている。
- 統計量の特徴を知っておかないと、状況の把握を間違えることがあることに注意する。

あるテストの得点と人数

点数	人数
10	2
9	4
8	5
7	4
6	2
5	0
4	0
3	0
2	0
1	1
0	2

■ 分散、標準偏差

- 分布のばらつき具合を表す指標
 - 分散

$$V = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{x})^2$$

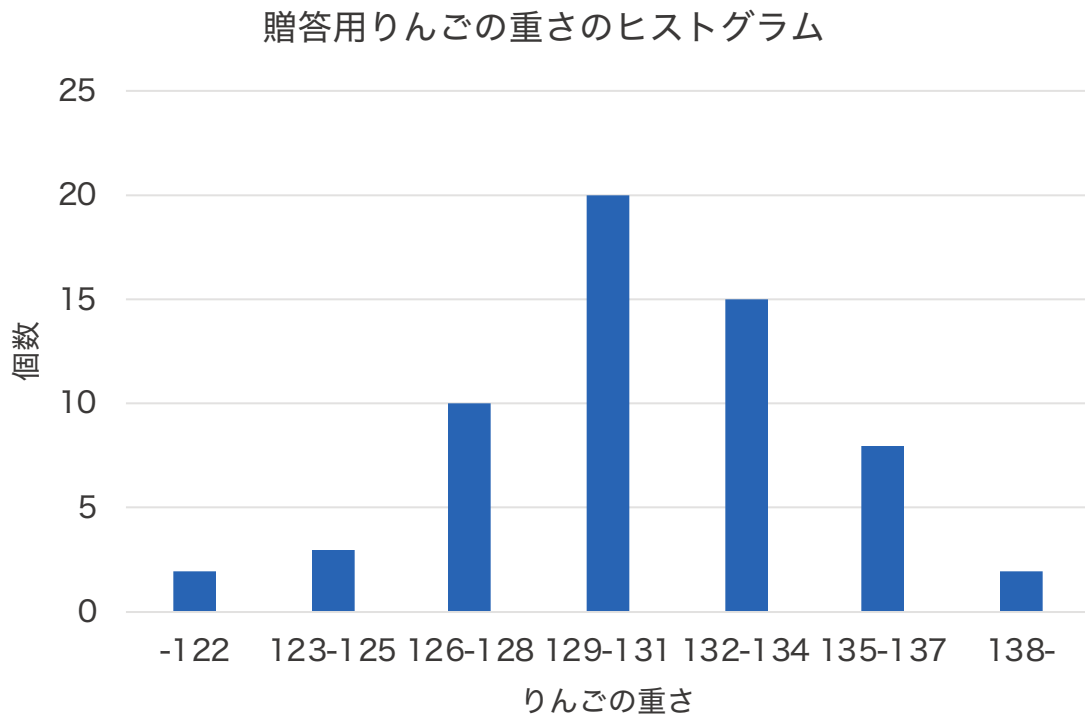
- 標準偏差

$$S = \sqrt{V}$$

(σ)

■ ヒストグラム(度数分布)

- データの分布の様子の把握に用いられる.



■ ヒストグラムの作り方

- 度数分布表を作る.
 - 観測値がとりうる値をいくつかの階級に分ける.
 - 観測値がそれぞれの階級でいくつあるか数える.
- 度数分布表にもとづき棒グラフを書く.
 - この棒グラフをヒストグラムという.

リンゴの重さのデータ表(g)					
133	130	127	121	137	130
132	130	129	130	130	137
135	133	121	129	132	130
140	133	132	129	129	132
126	132	132	127	129	129
130	124	135	137	127	132
126	129	130	135	137	132
130	130	127	133	135	124
126	127	130	132	133	126
124	127	140	130	132	129

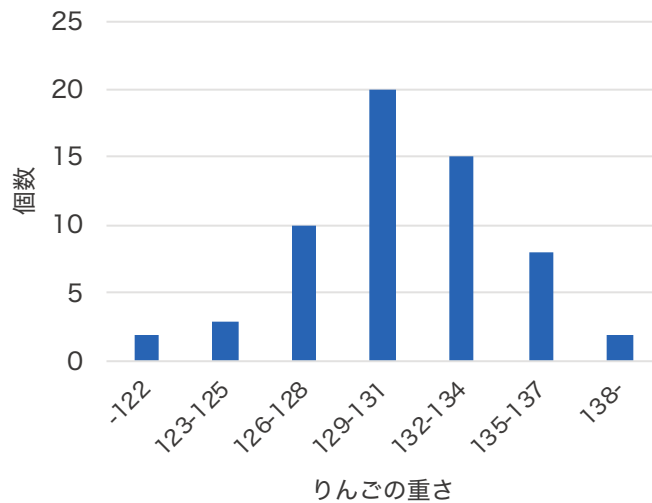


度数分布表

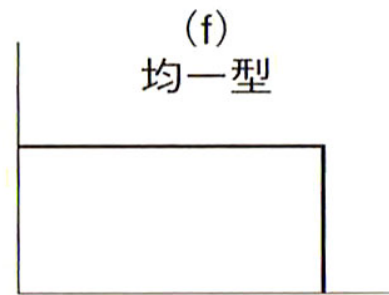
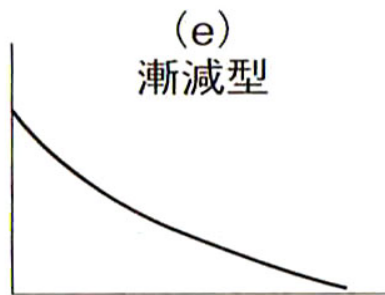
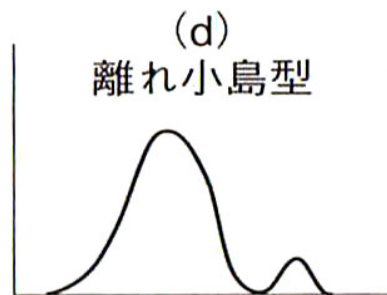
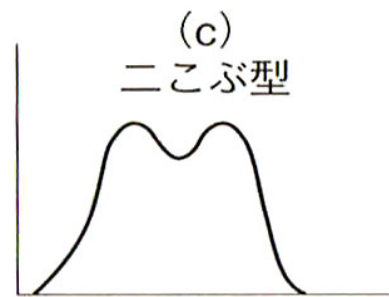
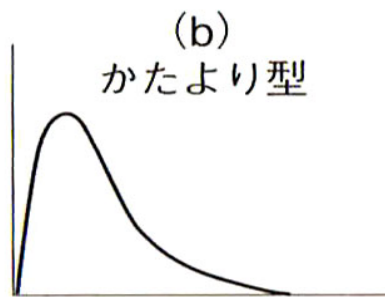
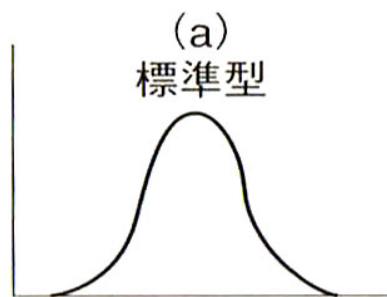
階級	度数
-122	2
123-125	3
126-128	10
129-131	20
132-134	15
135-137	8
138-	2



贈答用りんごの重さのヒストグラム

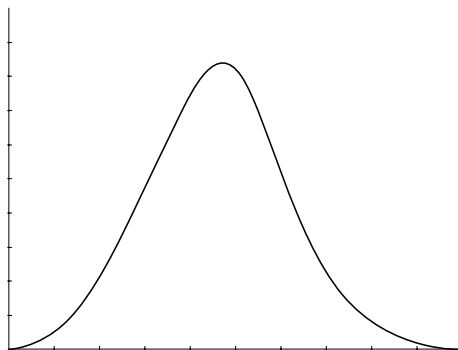


分布の形状



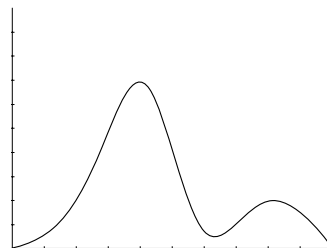
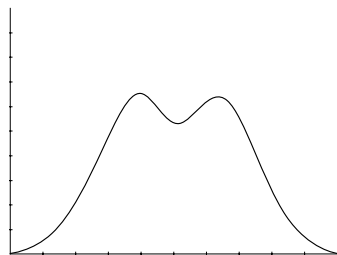
■ 正規分布

- 最も基礎的な分布の形
- 正規分布と呼ぶ
 - 標準型、一般型、ベル型、ガウス分布などと呼ばれることもある
- 何かを測定した場合、この分布になることが多い。



■ ふたこぶ型

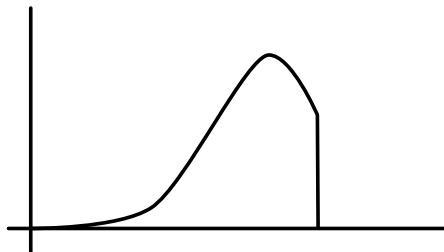
- 複数の要素を含む場合に生じる.
 - 成績の分布
 - 理解している人の集団と理解していない集団がある.
 - 製品の形状の分布
 - 一部が違う規格で作られている可能性がある.
 - このような分布が見られた場合、その原因を探る必要がある.
- データを適切にグループ分けすることで、峰が一つの単純な分布になることが多い. このグループ分けを層別と呼ぶ.



■ その他

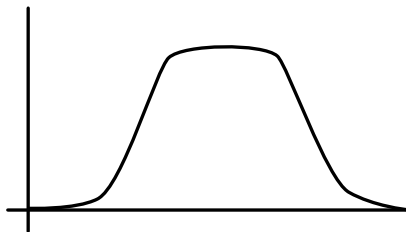
- 絶壁型

- ある値以下もしくは以上のものを選別して取り除いたときに現れる分布。



- 高原型

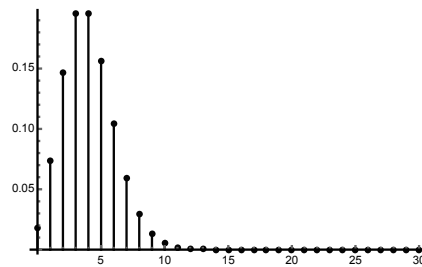
- ふたこぶ型の一種
- 平均値が少し異なるいくつかの分布が混在したときに現れる分布。
- 層別して原因を探る必要あり。



■ その他

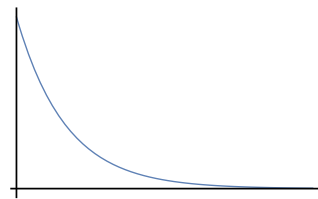
- ポアソン分布

- 品質管理の世界では偏り型などと呼ばれることもある。
- 交通事故件数, 大量生産の不良品件数、火災件数などはこの分布になる。



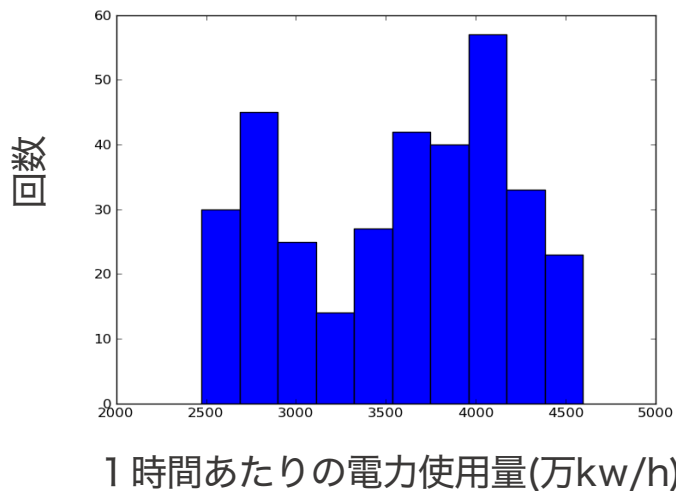
- 指数分布

- 品質管理の世界では漸減型などと呼ばれることがある。
- 待ち時間, 製品の故障、寿命などはこの分布になる。



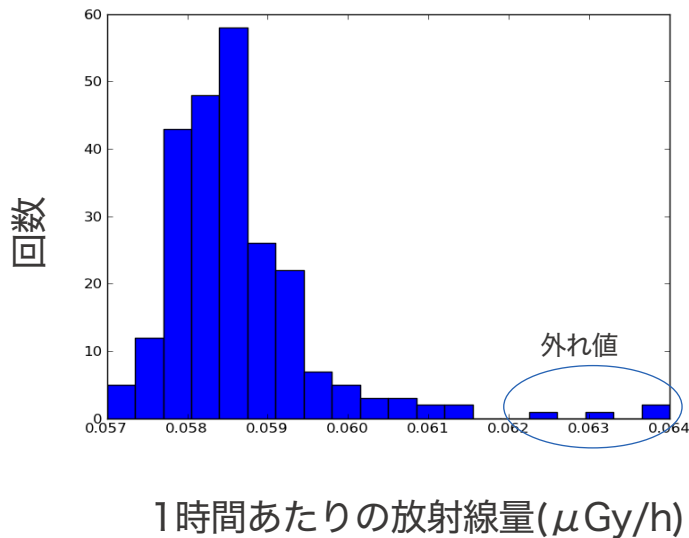
■ 例: 7月の1時間あたりの電力使用料

- ふたこぶ型の分布になっている
- 原因は夜と昼の電力使用量の性質が異なるためである。
- 昼と夜で層別が必要

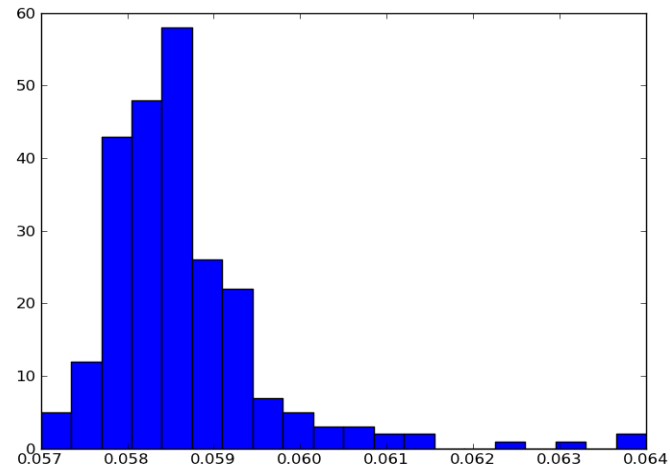
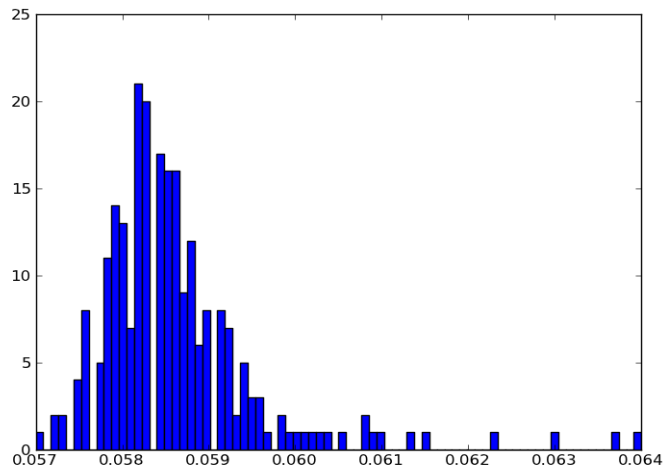


■ 例: 7月の新宿で観測された放射線量

- 基本的には正規分布ではあるが、外れ値が幾つか見られる。
- 外れ値がなぜ起こったか究明することが必要。



■ 作り方の悪いヒストグラム



グラフが歯抜けしているので良くない。
区間の設定(階級の幅)が不適切。

■ 標準偏差と分布の関係

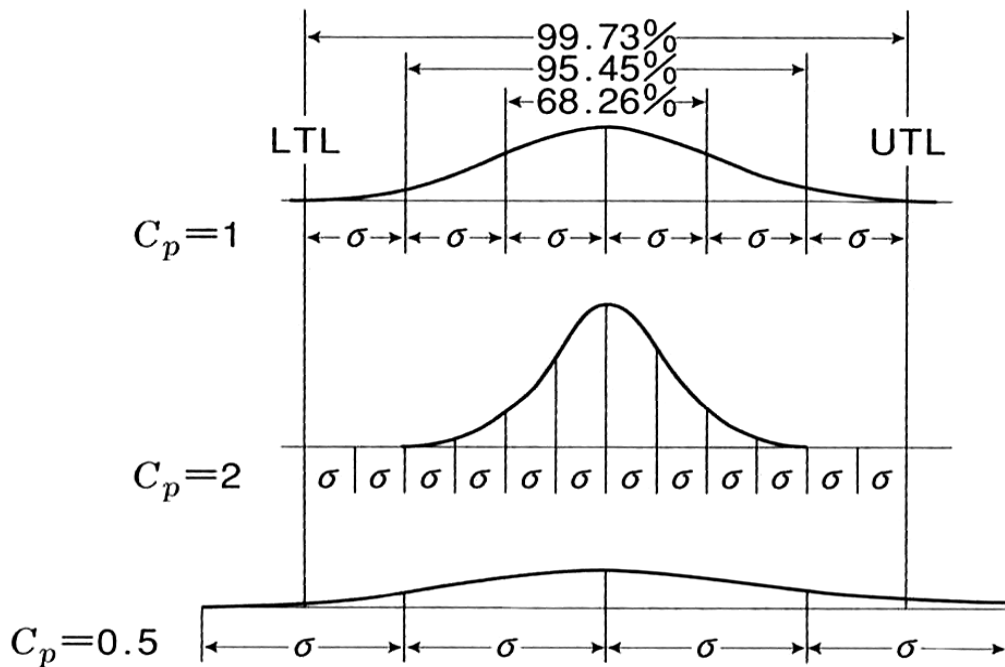
- 分散

$$V = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{x})^2$$

- 標準偏差

$$S = \sqrt{V}$$

(σ)



■ 工程能力指数

- 製品規格と分布の関係を表す指標
- 大きければ大きいほどよいが、大きすぎる場合過剰に対策をしている場合もある

- 上方許容限界(UTL)

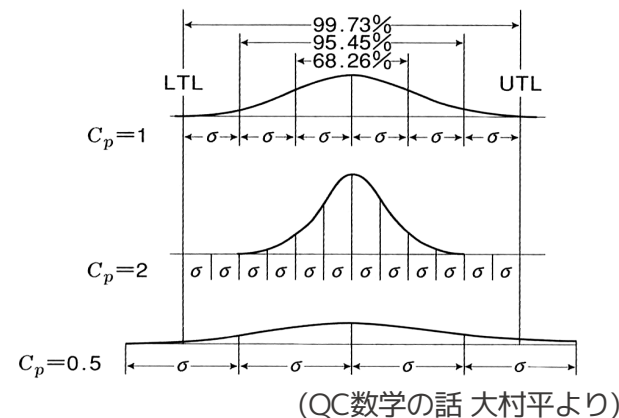
- 品質の上方限界

- 下方許容限界(LTL)

- 品質の下方限界

- 工程能力指数

$$C_p = \frac{UTL - LTL}{6\sigma}$$



工程能力指数は製品の分布が規格内にどれくらい収まっているかという指標

Cp=1なら生産された全製品のうち99.73%が規格内に収まっている。

■ 学力偏差値

- $(\text{得点} - \text{平均点}) \div \text{標準偏差} \times 10 + 50$
- で計算される.
- これは、テストの結果を平均点を50点にし、標準偏差を10点に正規化したという意味になる.

