

正規分布 ガウス分布 (ガウスのアイデア)

(天体の) x_i X ε_m
 観測値 = (真の値) + (誤差)



「誤差のバラツキ」が分かれば「真の値」を推定できる。

||
 ε 分布 確率密度 $f(\varepsilon) = f(x_i - X)$

1809 ガウス
 太陽の周りを楕円
 軌道で公転する天体
 の運行に関する理論

↓ ラプラスへ

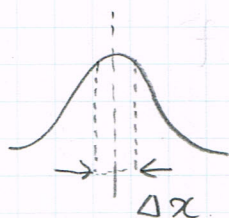
最小平乗法の前提と
 して応用

ガウスさんの経路から 3つの公理 (axiom)

A1. 絶対値の小さい誤差は
 絶対値の大きい誤差よりも多く発生

A2. 絶対値が同じ 正・負の誤差は
 同じくらいに発生

A3. 絶対値の^{とても}大きい誤差は非常に起こりにくい



$f(x_i - X) \cdot \Delta x$ は $x_i = X$ ($\varepsilon = 0$) で最大値

観測値全 n が Δx 内に入る同時確率は

$$f(x_1 - X) \Delta x \times f(x_2 - X) \cdot \Delta x \times \dots \times f(x_n - X) \cdot \Delta x$$

$$= (\Delta x)^n \prod_{i=1}^n f(x_i - X) = \left(\frac{d}{dx}\right)^n \cdot P(X) \text{ とする}$$

$$P(X) = f(x_1 - X) f(x_2 - X) \dots f(x_n - X) \dots \textcircled{A}$$

P も明らかに、 $x_i = X$ ($i=1 \sim n$) で最大値

$\geq P'(X = \bar{x}) = 0$ と仮定できる。(連続と仮定) ①

①の自然対数

$$\ln P(X) = \sum_{i=1}^n \ln(x_i - X)$$

$$X \text{ について微分} \quad \frac{P'(X)}{P(X)} = - \sum_{i=1}^n \frac{f'(x_i - X)}{f(x_i - X)} \quad \textcircled{B}$$

$$\left\{ \begin{array}{l} \ln(x) = \frac{1}{x} \text{ だと } \frac{d}{dx} (f(g(x))) \\ = f'(g(x)) \cdot g'(x) \text{ より} \end{array} \right.$$

正規分布 ガウス分布...ガウスのアイデア

⑧ $P'(x=\bar{x})=0$ より

$$\sum_{i=1}^n \frac{f'(x_i - \bar{x})}{f(x_i - \bar{x})} = 0 \quad \text{また} \quad \sum_{i=1}^n (x_i - \bar{x}) = 0$$

であることから、

$$\frac{f'(x)}{f(x)} = ax \quad (\text{注: 証明} \rightarrow \text{別紙} \boxed{3})$$

$$\Rightarrow f(x) = A \cdot e^{\frac{ax^2}{2}}$$

Aより、 $a < 0$ となるので $a = -\frac{1}{\sigma^2}$ と置くと、

$$f(x) = A \exp \left\{ -\frac{1}{2} \left(\frac{x}{\sigma} \right)^2 \right\}$$

全確率が1となるので

$$\begin{aligned} \int_{-\infty}^{\infty} f(x) dx = 1 &= \int_{-\infty}^{\infty} A \exp \left\{ -\frac{1}{2} \left(\frac{x}{\sigma} \right)^2 \right\} dx \\ &= A \sqrt{2\pi\sigma^2} \end{aligned}$$

(cf. ガウスの公式)
 $\int_{-\infty}^{\infty} e^{-ax^2} dx = \sqrt{\frac{\pi}{a}}$
 $(a > 0)$

$$\therefore A = \frac{1}{\sqrt{2\pi\sigma^2}}$$

$$\Rightarrow f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2} \left(\frac{x}{\sigma} \right)^2 \right\}$$

$$f(x-x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2} \left(\frac{x-x}{\sigma} \right)^2 \right\}$$

注) ただし、 $P'(x=\bar{x}) \neq 0$ という仮定は公理にはない。

$n(\geq 3)$ とし、 $\sum_{i=1}^n x_i = 0$ を満足するように x_i が変化するとき、

関数 $h(x)$ $\sum_{i=1}^n h(x_i) = 0$ が常に成立しているならば、

$h(x) = kx$ (k は任意の定数) である。

証明

$$x_n = -\sum_{i=1}^{n-1} x_i \quad \dots \textcircled{1}'$$

$$\frac{\partial}{\partial x_1} h(x_1) + \frac{\partial h(x_n)}{\partial x_n} \cdot \frac{\partial x_n}{\partial x_1} = 0$$

$$h'(x_1) + h'(x_n) \cdot \frac{\partial x_n}{\partial x_1} = 0$$

$$\textcircled{1}' \text{ から } \frac{\partial x_n}{\partial x_1} = -1 \text{ より } h'(x_1) = h'(x_n)$$

$$\text{同様に } h'(x_1) = h'(x_2) = \dots = h'(x_n)$$

$$h'(x_i) \text{ は } x \text{ によらないことから } h'(x) = k \text{ (定数)}$$

$$h(x) = kx + c$$

$$\textcircled{2} \text{ から } c = 0 \quad \text{よって } h(x) = kx \quad //$$

分散・共分散行列 と 機械学習 ML

1-1

cf. 多変数の正規分布

↑ (1変数の) 確率の分散の考え方を 多変数・多次元に拡張したもの

1変数 = 分散 $V(X) = \sigma_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ ($= \sigma_{xx}$)

2変数間 = 共分散

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

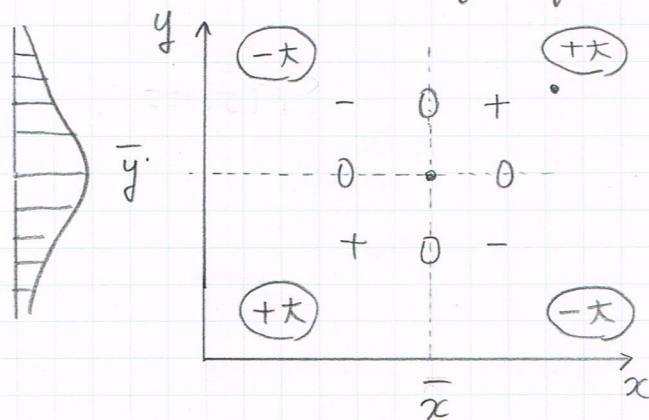
↑ 平均からの距離²の総和
データの回転質量(2次モーメント)

分散を最小
↓↑
最小二乗法

MLっぽい

これって何?

$(x_i - \bar{x})(y_i - \bar{y})$ は



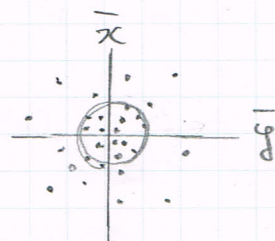
← このようになる。

↓
全てのデータについての総和

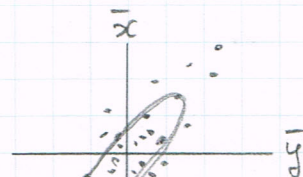
共分散

σ_{xy} の特徴

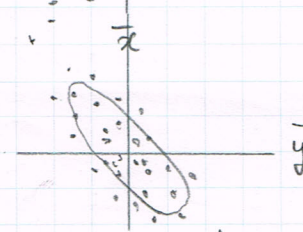
- ① 全体の (x, y) に関係性がなければ、
+点と-点が打ち消しあい、0になる。



- ② x増 → y増 の関係が多ければ、
+点が多くなり、「+」になる



- ③ x増 → y減 の関係が多ければ、
-点が多くなり、「-」になる。



分散・共分散行列 と 機械学習 ML

分散共分散行列の定義

$$\sigma_{xx} = (\sigma_x)^2, \sigma_{xy} = \sigma_{yx}$$

2変数

$$\Sigma = \Sigma_{xy} = \begin{pmatrix} \sigma_{xx} & \sigma_{xy} \\ \sigma_{xy} & \sigma_{yy} \end{pmatrix}$$

縮約記法的...

3変数

$$\Sigma = \Sigma_{xyz} = \begin{pmatrix} \sigma_{xx} & \sigma_{xy} & \sigma_{xz} \\ \sigma_{xy} & \sigma_{yy} & \sigma_{yz} \\ \sigma_{xz} & \sigma_{yz} & \sigma_{zz} \end{pmatrix}$$

「ことば」

対角成分は「分散」 それ以外は「共分散」の「行列」

例 x と y には正の相関があるが、
 z は独立しているデータセットの $\Sigma \rightarrow \begin{pmatrix} + & + & 0 \\ + & + & 0 \\ 0 & 0 & + \end{pmatrix}$

★ N 変数の中の任意の2変数の共分散($N C_2$ 個)と N 変数の(N 個)の分散が、対称な正方行列中に整理されている。

★ N 次元ベクトルに拡張した時の「分散」といえる。

→ 単に「分散(行列)」という人もいる
(共分散はビルイン)

↓ (バイズ勉強会の Σ)

$$f(\vec{x}) = \frac{1}{(2\pi)^{\frac{n}{2}} \sqrt{|\Sigma|}} \exp \left\{ -\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu}) \right\}$$

多変数 (n 次元ベクトル) \vec{x} の正規分布も Σ を使って

美しく表せる。 → cf. 1-1 多変数の正規分布

分散・共分散行列 と 機械学習 ML

それぞれの数値のスケールが異なるデータ間の比較

- 体重, 身長 ← 単位すらちがう
- 国語と算数の点数 ← 同じ国語の点数でも内容がちがえば...

標準化

$$\rightarrow \frac{x - \bar{x}}{s}$$

標準化する
事で異なる
スケールを
あわせる。

↓ データを標準化すると.

分散 $\Rightarrow 1$, 共分散 \Rightarrow 相関係数 になる。

例) 偏差値

分散 共分散行列 \Rightarrow 相関行列 になる。

|| Kaggleで良く出てくる Correlations テーブル

python: pandas-profiling.ProfileReport(df)

の Pearson の図

Pearson \rightarrow パラメトリック検定相関

Spearman \rightarrow データ順位相関

有意水準以上である

必要性がある。(正規分布仮定が可能)

ところで 相関係数って... 計算してみると.

$$r_{xy} = \frac{\sigma_{xy}}{\sigma_x \cdot \sigma_y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}} \xrightarrow{\text{ML}} \text{自然言語解析 画像解析}$$

$$= \frac{\vec{V}_x \cdot \vec{V}_y}{|\vec{V}_x| |\vec{V}_y|} = \cos \theta \quad \boxed{\text{コサイン類似度}}$$

なので、 $-1 \leq r_{xy} \leq 1$ であ。相関行列の成分の $[-1, 1]$ を色域に変換すれば。

多変数間のコサイン類似度を「見て分かる絵」にしてくれます。Kaggle