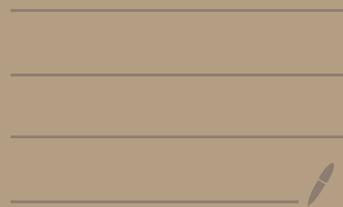


数学概論 (今野・2025/05/14)

最小2乗法と統計的回帰分析



講義の内容

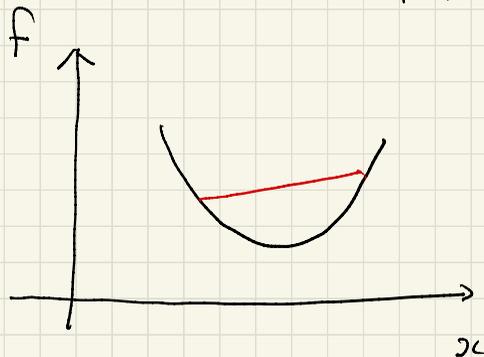
- 凸関数について
- 特異値分解と一般化逆行列
- 最小二乗法
- 最小二乗法の正則化: ridge 解と lasso 解

I. 凸関数の性質 (R上の関数)

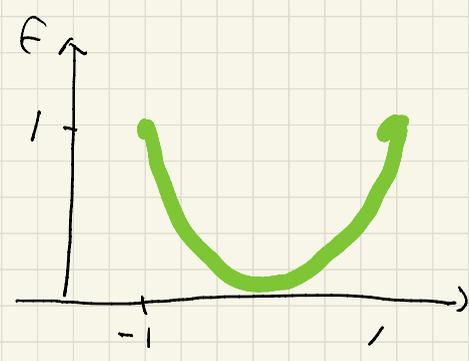
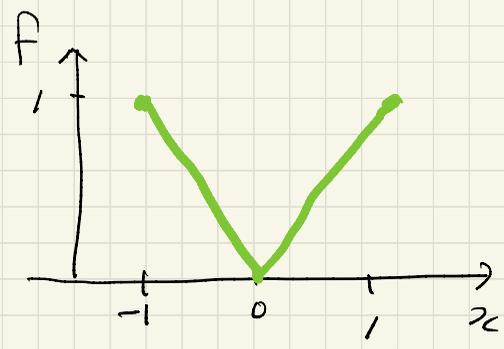
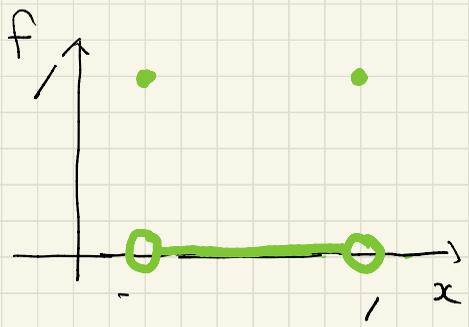
定義 C.1 $I \subset \mathbb{R}$ 区間とする.

関数 $f: I \rightarrow \mathbb{R}$ は凸 \iff $\forall x, y \in I$ と $0 \leq a \leq 1$ に対し

$$f(ax + (1-a)y) \leq af(x) + (1-a)f(y).$$



例 $I = [-1, 1]$ とする.



定理 C.3 $I \subset \mathbb{R}$ を区間とする. I 上の凸関数 f は
 I の内部で連続.

定理 C.5 区間 (a, b) 上で凸関数 f は次のように
 表現される.

$$f(x) = f(c) + \int_c^x g(t) dt \quad (a < x < b), \quad (a)$$

ただし, $c \in (a, b)$ で, g は単調増加の右連続.

注意 (a) の右辺の関数 g は凸.

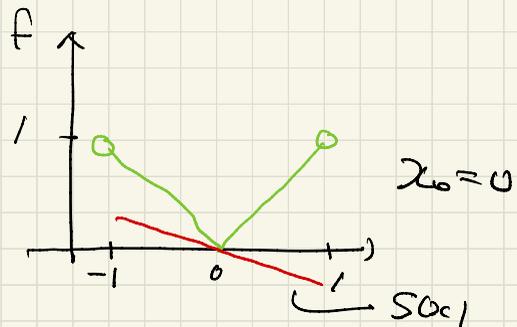
定義 C.8 f は区間 (a, b) 上の関数とする. 点 $x_0 \in (a, b)$ において, 適当に $m \in \mathbb{R}$ を定めると.

$$S(x) := m(x - x_0) + f(x_0) \leq f(x) \quad (\forall x \in (a, b))$$

が成り立つとき, $y = S(x)$ は点 x_0 において f を下から支える.

直線という

例



定理 C.9 $f: (a, b) \rightarrow \mathbb{R}$ とする. このとき

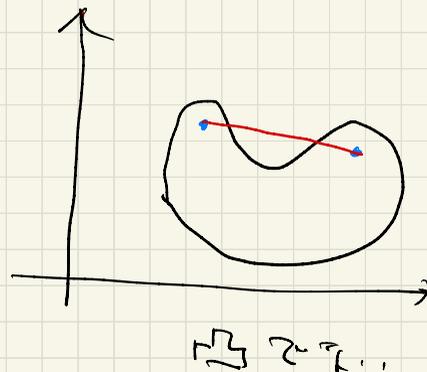
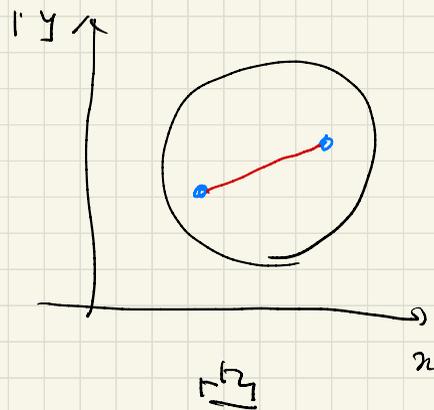
f は凸 \Leftrightarrow 下がまじる直線が x_0 で存在 ($\forall x_0 \in (a, b)$).

2. 凸関数と Young の不等式

定義 C.10 $C \subset \mathbb{R}^n$ の空でない部分集合とする.

C は凸集合 $\stackrel{\text{def}}{\iff} \forall x, y \in C \text{ と } 0 \leq a \leq 1 \text{ に対し}$

$$ax + (1-a)y \in C$$



定義 C.22 $C \subset \mathbb{R}^n$ は凸部分集合とする。

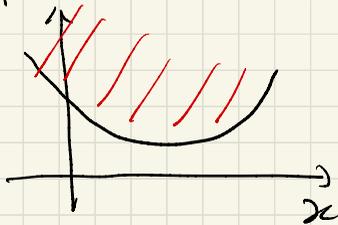
関数 $f: C \rightarrow \mathbb{R}$ は凸 $\iff \forall x, y \in \mathbb{R}^n$ と $0 \leq \alpha \leq 1$ に対し

$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$$

注意 C.23 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ は関数とし。

$$\text{epi}(f) := \left\{ \begin{bmatrix} x \\ y \end{bmatrix} : y \geq f(x) \quad (x \in \mathbb{R}^n) \right\} \subset \mathbb{R}^{n+1}$$

が f の IC^0 関数 \iff



すると

f は凸関数 $\Leftrightarrow \text{epi}(f)$ は凸集合.

□

定義 C.25 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ を関数とする. 各 $x \in \mathbb{R}^n$ に対して

$$\partial f(x) := \{ \mu \in \mathbb{R}^n; f(x) + \langle \mu, y - x \rangle_{2,n} \leq f(y) \} \\ (\forall y \in \mathbb{R}^n) |$$

とある. 尤も, $\langle \cdot, \cdot \rangle_{2,n}$ は \mathbb{R}^n 上の自然な内積とする.

この集合を点 x における関数 f の 劣微分 とする.

注意 • f は \mathbb{R}^n 上の凸関数 $\Leftrightarrow \partial f(x) \neq \emptyset$ ($\forall x \in \mathbb{R}^n$).

特に, 凸関数 f が点 $x_0 \in \mathbb{R}^n$ で微分可能ならば,

$$\partial f(x_0) = \{ \nabla f(x_0) \}.$$

$$\text{特に } \partial f(x) = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)^T, \quad x = (x_1, \dots, x_n)^T = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}.$$

例 $f(x) = |x|$ ($x \in \mathbb{R}$) とする。すると

$$\partial f(x) = \begin{cases} \{1\} & (x > 0) \\ [-1, 1] & (x = 0) \\ \{-1\} & (x < 0) \end{cases}$$

f は凸
 \Leftrightarrow

定理 C.28. $f: \mathbb{R}^n \rightarrow \mathbb{R}$ とする。各 $x \in \mathbb{R}^n$ に于いて, $\partial f(x) \neq \emptyset$

注意 (1) $f: \mathbb{R}^n \rightarrow \mathbb{R}$ は凸関数と $x_0 \in \mathbb{R}^n$ とする。

$$0_n \in \partial f(x_0) \iff \gamma = f(x_0) \text{ は } x = x_0 \text{ での最小値}$$

(2) $f: \mathbb{R}^n \rightarrow \mathbb{R}$ は凸と $x_0 \in \mathbb{R}^n$ とする。 $\epsilon > 0$ とする。

ある $\delta > 0$ が存在して

$$f(x_0) \leq f(x) \quad (\forall x \in B(x_0, \delta)) = \{y \in \mathbb{R}^n; |x_0 - y|_{2, n} \leq \delta\}$$

が成り立つ。

$$f(x_0) \leq f(x) \quad (\forall x \in \mathbb{R}^n)$$

仮定 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ は

$$\lim_{\|x\|_{2,n} \rightarrow \infty} \frac{f(x)}{\|x\|_{2,n}} = +\infty$$

をみたすとき、 f は超線型成長条件をみたすといふ。

例 $f(x) = |x|$ ($x \in \mathbb{R}$) は f である。

定義 C.31 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ は超線型成長条件をみたす凸関数

とする。このとき、 f の凸共役関数 f^v は

$$f^v(y) = \max_{x \in \mathbb{R}^n} \{ \langle x, y \rangle_{2,n} - f(x) \} \quad (y \in \mathbb{R}^n)$$

で定義される。

例 C.32 (2) $1 < p < \infty \leq 1$, $f(x) = |x|^p/p$ ($x \in \mathbb{R}$)

とすると、

$$f^v(y) = \frac{|y|^q}{q}$$

とすると、 $\tau \in \mathbb{R}$, $p^{-1} + q^{-1} = 1$ である。

補題 C.33 (Fenchel-Young 不等式) 定義 C.31 の設定を

仮定する。

$$\langle x, y \rangle_{2,n} \leq f(x) + f^v(y).$$

定理 C.34 定義 C.31 の設定を仮定する。

(1) $f^v: \mathbb{R}^n \rightarrow \mathbb{R}$ は凸関数

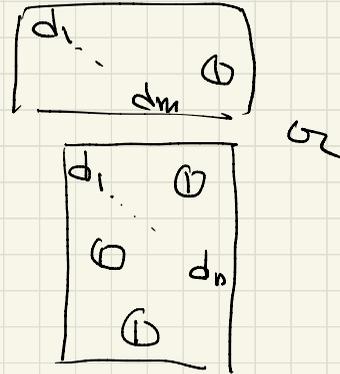
$$(2) \lim_{\|y\|_{2,n} \rightarrow \infty} \frac{|f^v(y)|}{\|y\|_{2,n}} = +\infty.$$

$$(3) (f^v)^v = f.$$

A: 特異値分解 (The singular Value Decomposition = SVD)

定理 A.7 $m, n \in \mathbb{N}$, $A \in \text{Mat}(m, n; \mathbb{R})$ は \mathbb{R} の

分解を持つ.



$$\underline{A} = \underline{U} \underline{D} \underline{V}^T,$$

$\underline{U} \in \text{Mat}(m; \mathbb{R})$, は 直交行列

$\underline{V}^T \in \text{Mat}(n; \mathbb{R})$ は 直交行列

$\underline{D} \in \text{Mat}(m, n; \mathbb{R})$ は 「対角行列」 である.

その対角成分は $d_1 \geq d_2 \geq \dots \geq d_{\min(m, n)} \geq 0$.

B. Moore-Penrose の一般化逆行列:

$\underline{A} \in \text{Mat}(m, n; \mathbb{R})$ とする。 A の Moore-Penrose の

一般化逆行列 $\underline{A}^{\dagger} \in \text{Mat}(n, m; \mathbb{R})$ を次の4条件を

みたすとする

$$A A^{\dagger} A = A, \quad (\text{B.3})$$

$$A^{\dagger} A A^{\dagger} = A^{\dagger}, \quad (\text{B.4})$$

$$(A A^{\dagger})^T = A A^{\dagger} \quad (\text{B.5})$$

$$(A^{\dagger} A)^T = A^{\dagger} A \quad (\text{B.6}).$$

注意 (1) 一意的に存在する.

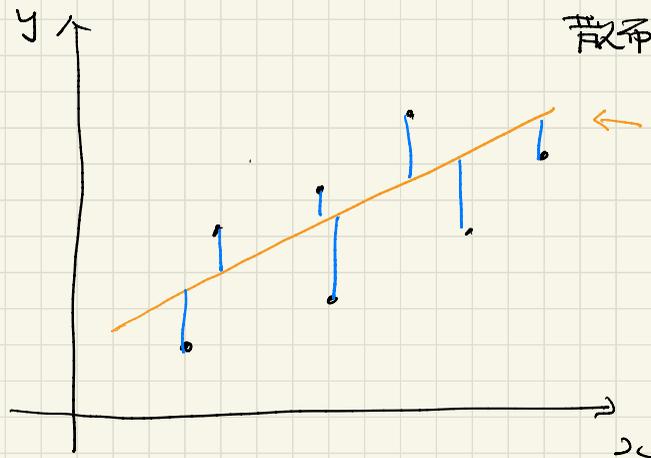
(2) $m \geq n$ かつ $A^T A$ は正則のとき.

$$A^+ = (A^T A)^{-1} A^T$$

と表す.

1. 最小2乗法の考え方

$n \geq 2$ とし、 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ を観測値とする



散布図

$$y = \beta_1 + \beta_2 x$$

と書く.

↑

x の情報を使って

y の動きを説明しよう!

↓

単回帰

1変量で1変量を説明

$$\hat{y}_j := \beta_1 + \beta_2 x_j \quad (j=1, 2, \dots, n)$$

$$e_j := y_j - \hat{y}_j : \text{残差}$$

最小二乗法

$$\min_{\beta_1, \beta_2} \sum_{j=1}^n \varepsilon_j^2 \quad \leftarrow \text{残差平方和の最小化}$$

可及的,

$$\begin{cases} n\beta_1 + \left(\sum_{j=1}^n x_j\right)\beta_2 = \sum_{j=1}^n y_j \\ \left(\sum_{j=1}^n x_j\right)\beta_1 + \left(\sum_{j=1}^n x_j^2\right)\beta_2 = \sum_{j=1}^n x_j y_j \end{cases}$$

$x_1 = \dots = x_n$ 以外の解が存在しない。

次元への拡張 (重回帰)

$(x_1, y_1), \dots, (x_n, y_n)$ を観測値として与える。

$x_j \in \mathbb{R}^d$ (線形ベクトル), $y_j \in \mathbb{R}$ ($j=1, \dots, n$)

記号

$$\underline{X}^{\text{data}} = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} \in \text{Mat}(n, d; \mathbb{R}), \quad \underline{y}^{\text{data}} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n, \quad \underline{\beta} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_d \end{pmatrix} \in \mathbb{R}^d$$

$n \geq d$ かつ $\underline{y}^{\text{data}} \notin \text{span}(x_1, \dots, x_n)$ となる。

つまり, $\underline{y}^{\text{data}} = \underline{X}^{\text{data}} \underline{\beta}$ となる $\underline{\beta} \in \mathbb{R}^d$ は存在しない!

残差ベクトル

$$e := \underline{X} \beta - y^{\text{data}} \quad \leftarrow \text{方程式 } y = \underline{X} \beta \text{ を定めたこと
況しては。}$$

と定める。 以後は「data」を扱う

最小2乗解 $\hat{\beta}^{\text{LS}}$ を

$$|\underline{X} \hat{\beta}^{\text{LS}} - y|_{2,n} \leq |\underline{X} \beta - y|_{2,n} \quad (\forall \beta \in \mathbb{R}^d)$$

で定める。(存在はするが - 竟とはかきこえぬ。)

最小2乗解を求めよう。

$$\beta = (\beta_1, \dots, \beta_d)^T$$

$$f(\beta) := \|\underline{X}\beta - \underline{y}\|_{2,n}^2, \quad \nabla := \left(\frac{\partial}{\partial \beta_1}, \dots, \frac{\partial}{\partial \beta_d} \right)^T$$

したがって、 $\hat{\beta}^{OLS}$ は

$$\nabla f(\beta) = 2\underline{X}^T(\underline{X}\beta - \underline{y}) = \mathbf{0}_d \iff \nabla f(\hat{\beta}^{OLS}) = \mathbf{0}_d$$

を得る。

$$\iff \underline{X}^T \underline{X} \hat{\beta}^{OLS} = \underline{X}^T \underline{y}$$

$\underline{X}^T \underline{X}$ が正則ならば

$$\hat{\beta}^{OLS} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} = \underline{X}^+ \underline{y}$$

と表す。

最小二乗法によるデータの当てはめ

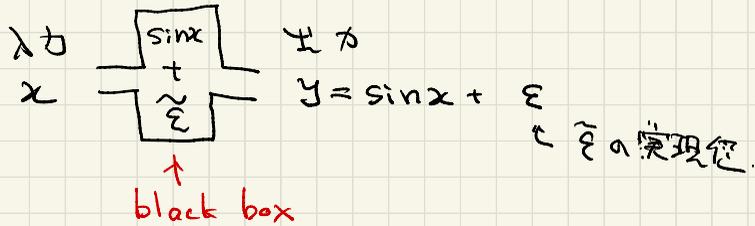
回帰モデル

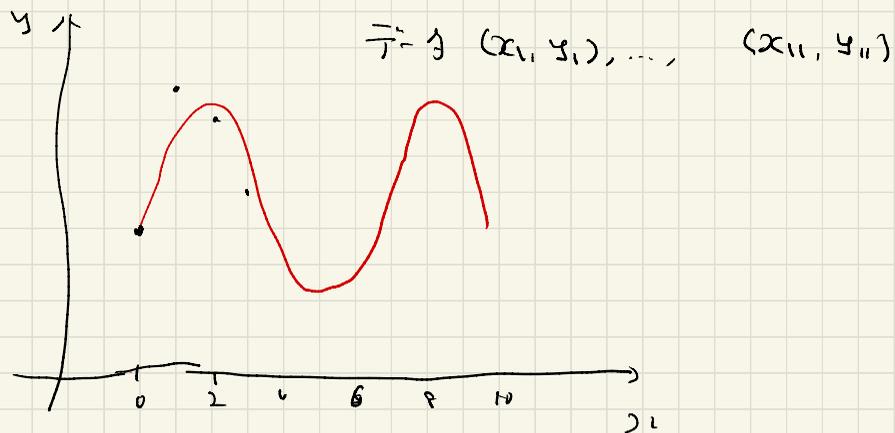
$$\tilde{y} = \sin x + \tilde{\varepsilon}, \quad \tilde{\varepsilon} \sim N(0, 0.04)$$

平均0, 分散(0.2)の
正規分布

人間には未知

を考慮. x は固定して仮, y と $\tilde{\varepsilon}$ は確率変数





の近似の式は

$$f(x) = \beta_1 + \beta_2 x + \dots + \beta_n x^{n-1}$$

$$y = X\beta + \text{誤差ベクトル} \quad \sim$$

$$y = (y_1, \dots, y_n)^T, \quad \beta = (\beta_1, \dots, \beta_n)^T$$

$$\underline{X}_d = \begin{pmatrix} 1 & x_1 & \dots & x_1^{d-1} \\ 1 & x_2 & \dots & x_2^{d-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & \dots & x_n^{d-1} \end{pmatrix} \in \text{Mat}(n, d; \mathbb{R}).$$

$$\underline{\beta}_d := (\beta_1, \dots, \beta_d) \in \mathbb{R}^d \quad (d=2, 3, \dots, n)$$

$$X_n = X, \quad \beta_n = \beta \text{ とおく.}$$

$$\hat{\beta}_d^{OLS} = (\underline{X}_d^T \underline{X}_d)^{-1} \underline{X}_d^T \underline{y} \text{ とし, } \hat{\beta}_d^{OLS} \in \mathbb{R}^d \text{ とする}$$

5 次の多項式を与えることが出来る.

↑ 答えを知っているから出来る.

$$\min_{\beta \in \mathbb{R}^d} \|X\beta - y\|_{2,11}^2 \quad \leftarrow \text{曲線のデータへの適合度}$$

を考へても可^い. 可^いわ^る, 残差平方和を最小化して可^い算^す.

変数選択 (AIC, BIC) $\Rightarrow \beta$ のいくつかの成分は「0」と見^らす.

$\hat{\beta}$: データ $(x_1, y_1), \dots, (x_n, y_n)$ から作^られた多項式の係数

$(y^{\text{future}}, x^{\text{future}})$: 将来のデータ

$$x^{\text{future}} \Rightarrow \mathcal{X}^{\text{future}} \Rightarrow |y^{\text{future}} - (\mathcal{X}^{\text{future}})^T \hat{\beta}|^2$$

$$x \mapsto \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^b \end{bmatrix}$$

を小さくする $\hat{\beta}$ がよいとの.

だが大抵は, brutal search!

◦ 正則化法

(曲線の予測の適合度) + (βの複雑さへの罰則)

βの複雑さ ↑ 罰則 ↑

最適化

βに課する罰則 $\beta = (\beta_1, \dots, \beta_d)$

(a) $|\beta|_{0,d} := (\beta \text{ の } 0 \text{ でない成分の個数})$

(b) $|\beta|_{1,d} := \sum_{j=1}^d |\beta_j|$

(c) $|\beta|_{2,d} := \sum_{j=1}^d \beta_j^2$

(a) は β に関する非凸問題だが, (b) と (c) は凸問題だ

$$\min_{\beta \in \mathbb{R}^n} \|X\beta - y\|_{2,1}^2 + r \|\beta\|_{2,1}^2 \rightarrow \hat{\beta}^{\text{ridge}}$$

$$\min_{\beta \in \mathbb{R}^n} \|X\beta - y\|_{2,1}^2 + 2r \|\beta\|_{1,1} \rightarrow \hat{\beta}^{\text{lasso}}$$

ただし, $r > 0$ は正則化変数 \leftarrow 人間が決めた。