

Part II

$$\theta^* \in \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n l(f_\theta(x_i), y_i)$$

where $f_i(\theta) := l(f_\theta(x_i), y_i)$

$$f(\theta) := \frac{1}{n} \sum_{i=1}^n f_i(\theta)$$

$$\begin{cases} \theta_0 \in \mathbb{R}^d \\ \theta_{t+1} = \theta_t - \eta \nabla f_t(\theta_t) \end{cases}$$

stochastic update

1. In general, is θ^* unique?

It's unconstrained optimization problem

$$\forall x \in \mathbb{R}^n \exists \nexists x^*, f(x^*) < f(x)$$

x^* is global optimizer

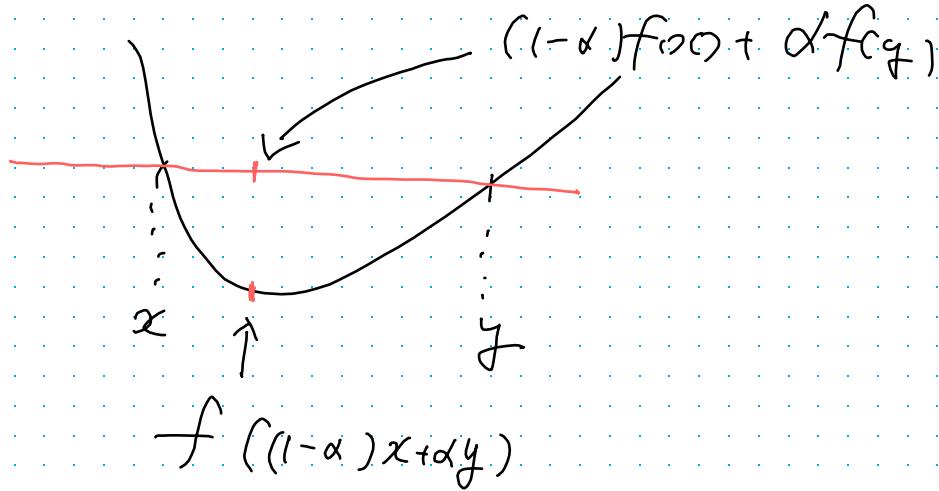
$$|x^* - x| < \varepsilon \quad \forall x \in \mathbb{R}^n \exists \nexists x^* f(x^*) < f(x)$$

x^* is local optimizer

Definition of Convex function

$f: \mathbb{R}^n \rightarrow (-\infty, \infty)$ $\forall x, y \in \mathbb{R}^n \quad \alpha \in [0, 1]$

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



if objective function is convex and
constraint function is also convex

convex problem

↓
local solution = global solution

Attendants

How to identify f is convex

\Rightarrow calc Hessian

$$\text{if } f(x) = x^2$$

$$f''(x) = \frac{\partial^2}{\partial x^2} f(x) = \frac{\partial^2}{\partial x^2} (2x) = 2 \geq 0$$

positive Semidefinite

Iterative Algorithm

$$x^{k+1} = x^k + \alpha^k d^k$$

step size , aka learning rate

Search direction

How to find step size

1. line search (逐次搜索)

find minimum $\alpha \geq 0$ where min loss value
when steps in α^k direction

$$f(x^k + \alpha^k d^k) = \min_{\alpha} f(x^k + \alpha d^k)$$

but, in practice it is not used due to computational cost

Attendents

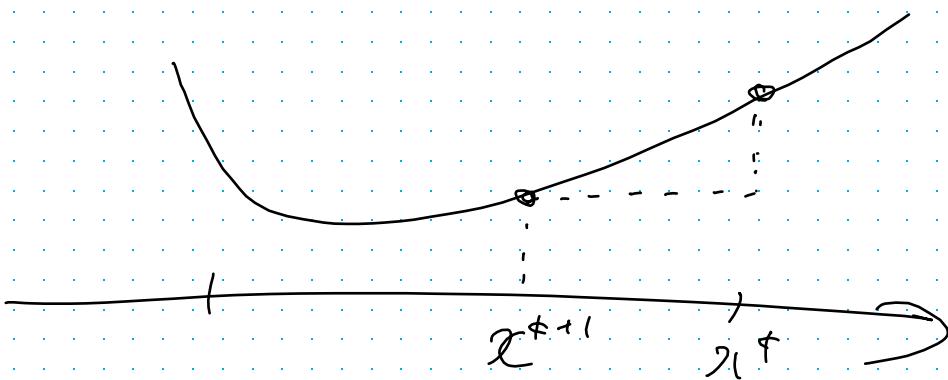
How to find good value of steps?

Armijo's condition

Wolf condition

How to find d^* direction

$$f(x^{k+1}) - f(x^k) = f(x^k + \alpha^k d^k) - f(x^k)$$



$$\lim_{t \rightarrow +\infty} \frac{f(x^k + t d^k) - f(x^k)}{t} = \nabla f(x^k) d^k$$

$d^k = -\nabla f(x^k)$ gradient
descent
direction

Attendants

Assumption

convex- $D^2 L(w)$ positive definite

Lipschitz condition

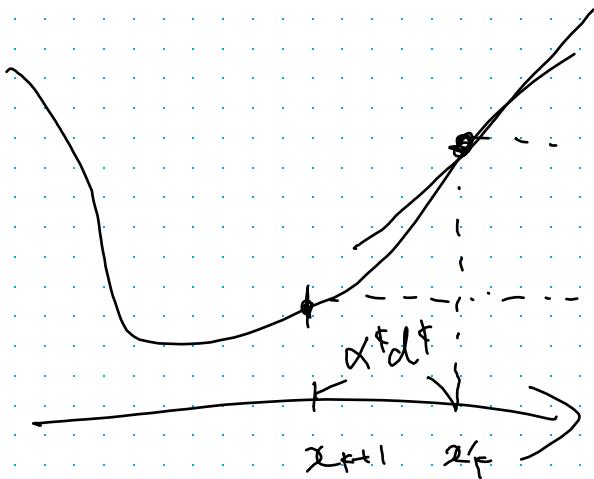
$$|L(w^k) - L(w^{k+1})| \leq Q \|w^k - w^{k+1}\|$$

$$\Rightarrow \|D L(w^k)\|_2 \leq Q$$

$$f(x_k + x) \approx f(x_k) + Df(x_k)^T x + \frac{1}{2} x^T H(x_k) x$$

f の \bar{x} は 2 次の T イテ - ~~近似~~ \bar{f}

$$\bar{f}(x) = f(\bar{x}) + f'(\bar{x})^T (x - \bar{x}) + \frac{1}{2} (x - \bar{x})^T f''(\bar{x}) (x - \bar{x})$$



$$f(x_{k+1}) - f(x_k)$$

$$= f(x_k + \alpha_k d_k) - f(x_k)$$

$$< 0$$

こうすると d_k を下方向

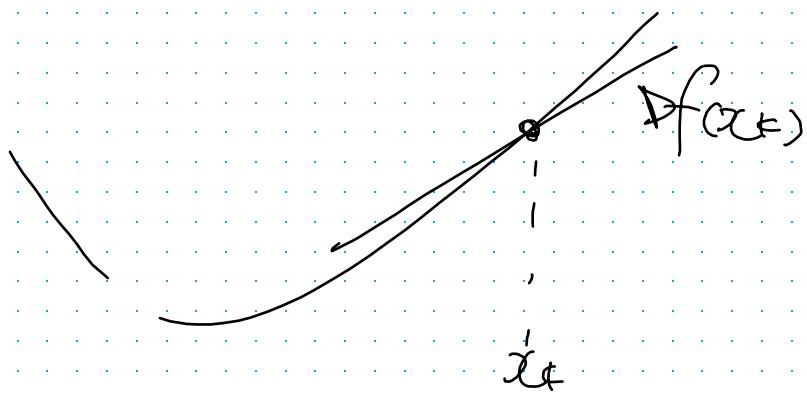
$$\lim_{\alpha_k \rightarrow 0} \frac{f(x_k + \alpha_k d_k) - f(x_k)}{\alpha_k d_k} = Df(x_k)$$

$$\Leftrightarrow \lim_{\alpha_k \rightarrow 0} \frac{f(x_k + \alpha_k d_k) - f(x_k)}{\alpha_k} = Df(x_k) d_k$$

方向微分係数

このときは d_k を下方向

$$\nabla f(x_k) d_k = \|\nabla f(x_k)\| \|d_k\| \cos \theta$$



$$d_k = \underset{dk}{\arg \min} \left\{ f(x_k + \alpha d_k) - f(x_k) \right\}$$

$$= \underset{\alpha k}{\arg \min} \frac{f(x_k + \alpha d_k) - f(x_k)}{\alpha k}$$

$$= \underset{d_k}{\arg \min} \underbrace{\left\{ \nabla f(x_k) d_k \right\}}$$

$$\min \left\{ \nabla f(x_k) d_k \right\} = \min \left\{ \|\nabla f(x_k)\| \|d_k\| \cos \theta \right\}$$

" " $\in \mathbb{R}^n$

$$\therefore d_k = -\nabla f(x_k)$$

Attendents

収束性

$$x^{k+1} = x^k - \alpha^k Df(x^k)$$

$\{x^k\}$ の点列が収束するならば $\|Df(x^k)\| \rightarrow 0$
 $(\alpha \rightarrow 0)$

PR (x^k + d^k)

$x_{outward}$ 条件)

$$\sum_{k=0}^{\infty} \left(\frac{Df(x^k)^T d^k}{\|d^k\|} \right) < \infty$$

$$\Leftrightarrow \sum_{k=0}^{\infty} \|Df(x^k)\| \cos \theta_k < \infty$$

$$\because \cos \theta_k = \left(\frac{-Df(x^k) d^k}{\|Df(x^k)\| \|d^k\|} \right)$$

$$d^k = -Df(x^k) \cdot \vec{v}$$

$$\cos \theta_k = 1$$

73)

Zoutendijk 定理

$$\frac{\infty}{2} \quad \frac{0}{0}$$

$$\|Df_{\text{act}}\| < \infty$$

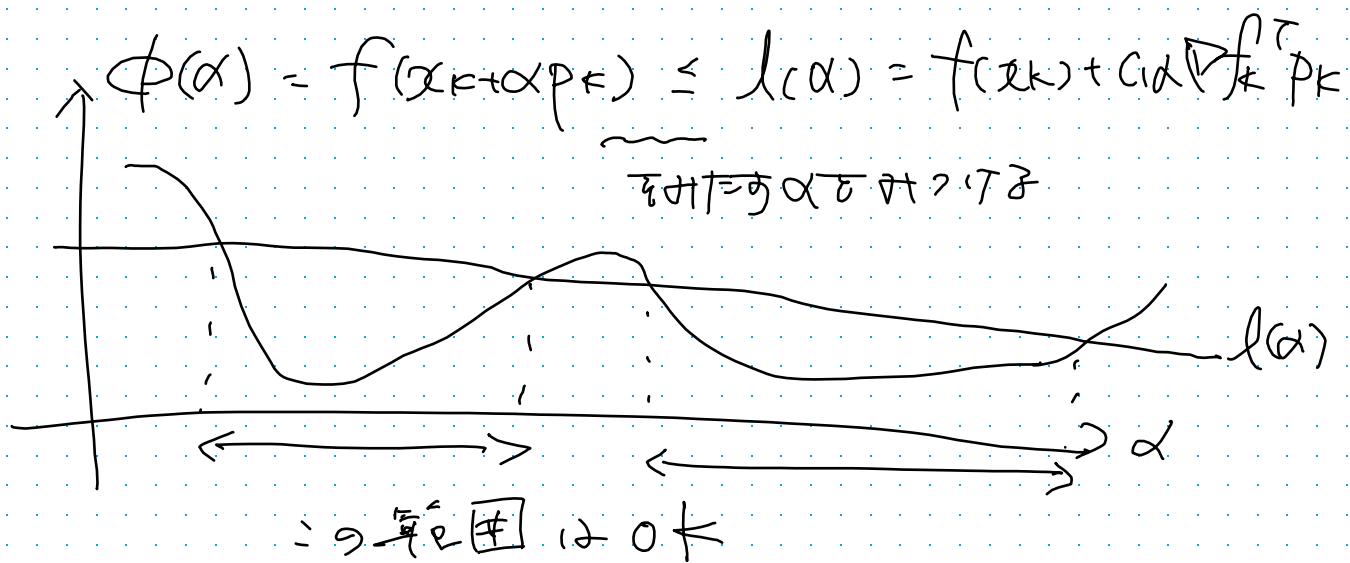
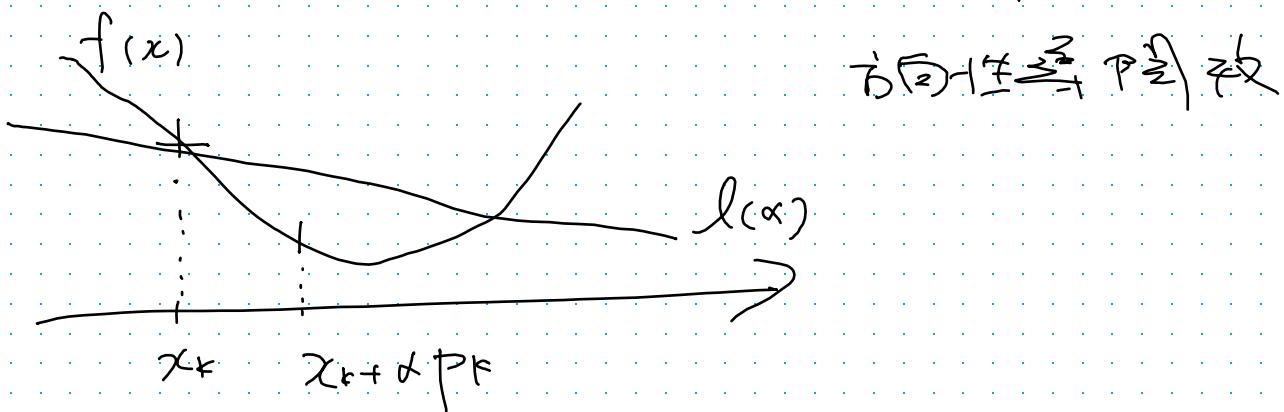
(無限級数収束定理)

$$\lim_{k \rightarrow \infty} \|Df_{\text{act}}^k\| = 0$$

無限級数収束定理

Armijo 条件 (不完全勾配法探索条件)

$$f(x_k + \alpha p_k) \leq f(x_k) + c_1 \alpha \underbrace{\nabla f_k^T p_k}_{\text{方向性一致}} = l(\alpha)$$



$\phi(\alpha) \leq l(\alpha)$ なら十分減少条件を満たす

Attendents

、~~小二すきいを~~ とての ~~条件~~ が ~~ない~~

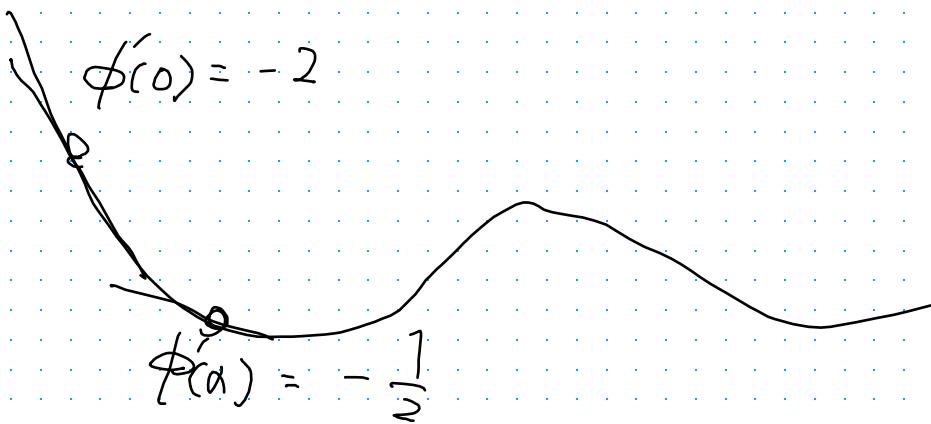
→ 曲率条件

$$\alpha \rightarrow 1^- \quad Df(x_k + \alpha_k p_k)^T p_k \geq C_2 \|Df\| p_k^T$$

"

$$\phi'(x_k)$$

$$\phi'(x_k) \neq \phi'(0) \text{ 大きい}$$



Attendents

Wolfe 条件

十分な減少条件と曲率条件

$$\phi(\alpha) < l(\alpha)$$

$$\epsilon + \tau = \alpha$$

$$\phi(\alpha) > l(\alpha)$$

$$\epsilon + \tau = \alpha$$

Zoutendijk 条件 (大域的収束原理)

P_k の降下方向, α_k が Wolfe 条件を満たす

$$\| Df(x) - Df(\bar{x}) \| \leq L \| x - \bar{x} \|$$

成り立つとすると

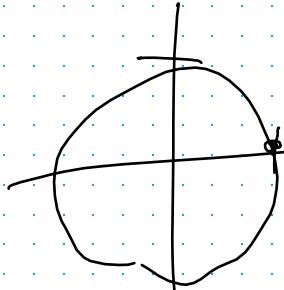
$$\sum_{k=0}^{\infty} \cos^2 \theta_k \| Df_k \|^2 < \infty \quad \text{が成立}$$

$$\text{where } \cos \theta_k = \frac{-Df_k^T P_k}{\| Df_k \| \| P_k \|}$$

探索方向 P_k と最急降下方向 $-Df_k$ 間の角 θ_k

$$\cos \theta = 0$$

$$\cos \theta \in \|\nabla f_k\|^2 \rightarrow 0$$



$$\sin \theta = y$$

$$\cos \theta = x$$

$$\lim_{k \rightarrow \infty} \|\nabla f_k\| = 0$$

$$\cos \theta \geq 0 \Leftrightarrow \theta + 90^\circ \text{ で } x \geq 0 \text{ の場合}$$

↑
大域的収束

が示すように

収束率

$$\frac{|x_{k+1} - \bar{x}|}{|x_k - \bar{x}|} \leq q, \quad \forall k > K \text{ かつ 成立}$$

\Rightarrow 一次収束

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - \bar{x}|}{|x_k - \bar{x}|} = 0 \Rightarrow \text{超一次収束}$$

$$\frac{|x_{k+1} - \bar{x}|}{|x_k - \bar{x}|^p} \leq M, \quad M > 0 \text{ かつ } p > 1 \Rightarrow p\text{-次収束}$$

Attendents

最急降下法の収束率

$$f(x) = \frac{1}{2} x^T Q x - b^T x$$

$$x^* = Q^{-1} b \quad \therefore f'(x) = 0$$

$$\Leftrightarrow Q x - b = 0$$

$$\Leftrightarrow x = Q^{-1} b$$

最急降下法

$$x_{k+1} = x_k - \alpha_k Df_k$$

$$f(x_{k+1}) = f(x_k - \alpha_k Df_k)$$

$$= \frac{1}{2} (x_k - \alpha_k Df_k)^T Q (x_k - \alpha_k Df_k) - b^T (x_k - \alpha_k Df_k)$$

$$\alpha_k \text{ は } \gamma \text{ の範囲 } , \quad \alpha_k = \frac{Df_k^T Df_k}{Df_k^T Q Df_k}$$

$$\Rightarrow x_{k+1} = x_k - \left(\frac{Df_k^T Df_k}{Df_k^T Q Df_k} \right) Df_k$$

$\gamma = 2^{-\frac{1}{2}} \lambda$, $\lambda = 1/\mu_n$

$$\|x\|_Q^2 = x^T Q x \leq \frac{1}{\lambda} \|x\|^2$$

Attendants

$$x^* = Q^{-1}b^*$$

$$\frac{1}{2} \|x - x^*\|_Q^2 = f(x) - f(x^*)$$

$$\leftarrow \frac{1}{2} (x - x^*)^T Q (x - x^*)$$

$$\frac{1}{2} x^T Q x - b^T x$$

$$f(x^*) = \frac{1}{2} Q^{-1} b^T Q Q^{-1} b - b^T Q^{-1} b$$

$$= -\frac{1}{2} Q^{-1} b^2 - Q^{-1} b^2$$

$$= -\frac{1}{2} Q^{-1} b^2$$

$$\rightarrow = \frac{1}{2} x^T Q x - b^T x + \frac{1}{2} Q^{-1} b^2$$

$$= f(x) - f(x^*)$$

$$\frac{1}{2} (x - x^*)^T Q$$

$$= \frac{1}{2} Q \left\{ x^2 - 2x x^* + x^{*2} \right\}$$

$$= \frac{1}{2} x^T Q x - x^T Q Q^{-1} b + \frac{1}{2} Q^{-1} b^T Q^{-1} b$$

Attendants

$$\nabla f_k = Q(x_k - x^*)$$

$$f(x) = \frac{1}{2} x^T Q x - b^T x$$

$$\nabla f(x) = Qx - b$$

$$= Q(x - Q^{-1}b)$$

$$\Rightarrow Q(x - x^*)$$

$$\|x_{k+1} - x^*\|_Q^2 = \|x_k - (-\nabla f_k - x^*)\|_Q^2$$

$$\frac{\nabla f_k^T \nabla f_k}{\nabla f_k^T Q \nabla f_k}$$

$$= \|(\nabla f_k + x^*) - (-\nabla f_k)\|_Q^2$$

$$= (\underbrace{\nabla f_k}_1)_Q$$

$$= \left\{ (\nabla f_k + x^*)^2 - 2(-\nabla f_k)(\nabla f_k + x^*) + (-\nabla f_k)^2 \right\}_Q$$

$$= \|x_k - x^*\|_Q^2 + \|(-\nabla f_k)\|_Q^2 - 2(-\nabla f_k)(\nabla f_k + x^*)$$

Attendents

$$\|x_k - x^*\|_Q^2 \left\{ 1 + \frac{(Df_k)^2}{Df_k(x - x^*)} - \frac{2(Df_k(x - x^*))}{Df_k(x - x^*)} \right\}$$

$$= \|x_k - x^*\|_Q^2 \left\{ 1 + (Df_k)^2 \times \frac{Df_k}{(x - x^*)} - 2(Df_k) \right\}$$

$$= \|\tilde{x}\|_Q^2 \left\{ 1 + (Df_k) \left[\frac{Df_k}{(x - x^*)} - \frac{2(x - x^*)}{(x - x^*)} \right] \right\}$$

$$= (\|\tilde{x}\|_Q^2) \left\{ 1 + (Df_k) \frac{Df_k - 2Q^{-1}Df_k}{Q^{-1}Df_k} \right\}$$

$$= \|x_k - x^*\|_Q^2 \left\{ 1 + \frac{Df_k^T Df_k}{Df_k^T Q Df_k} \frac{Df_k (1 - 2Q^{-1})}{Q^{-1} Df_k} \right\}$$

$$= \|\tilde{x}\|_Q^2 \left\{ 1 + \frac{Df_k^T Df_k (Q - 2)}{Df_k^T Q Df_k} \right\}$$

Attendents

-収束性と干渉の収束率

$$\|x_{k+1} - x^*\|_Q \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2 \|x_k - x^*\|_Q$$

誤差の減衰が漸下する

$\Rightarrow f_{k+1}, f_{k+2}$ 一次収束する

条件数 $C(Q) = \frac{\lambda_n}{\lambda_1}$ の大きさと

収束性が悪化する

$$\tau \in \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}, 1 \right) \quad \tau \rightarrow 0 \ (= 0)$$

$$f(x_{k+1}) - f(x^*) = \tau^2 |f(x_k) - f(x^*)|$$

$$\text{e.g. } \lambda_n - \lambda_1$$

$$0.5 - 0.1 = 0.4$$

$$\frac{2}{3} = 0.66.$$

$$0.5 + 0.1 = 0.6$$

$$\text{e.g. } \lambda_n - \lambda_1$$

$$0.5 - 0.3 = 0.2$$

$$\frac{1}{4} = 0.25$$

$$0.5 + 0.3 = 0.8$$

$\overline{0.25}$

Attendents

Newton Method.

$$p_k = - \underbrace{\nabla^2 f_k^{-1}}_{H} \nabla f_k$$

$H = \nabla^2 f_k$ の正定性と
PBM

→ trust region, modify H

2次収束でまとめて二階導関数を修正する

$$f(x_k + d) \approx f(x_k) + \nabla f(x_k) d + \frac{1}{2} \nabla^2 f(x_k) d^2$$

泰勒-展開 $a \rightarrow f(x)$ の尾関

$$f(x) = f(a) + f'(a)(x-a) + \frac{1}{2} f''(a)(x-a)^2 + \dots$$

マクロ-ソン尾関 a の $f(x)$ の尾関

x_k の尾関

$f(x) \approx x_k +$ 展開

$$f(x) = f(x_k) + \nabla f(x_k) (x - x_k) + \frac{1}{2} \nabla^2 f(x_k) (x - x_k)^2$$

$$x \leftarrow x_k + d \in \mathbb{R}^n \quad f(x_k + d) = f(x_k) + \nabla f(x_k) d + \dots$$

$$f(x_k + d) = f(x_k) + \nabla f(x_k) d + \frac{1}{2} \nabla^2 f(x_k) d^2 \dots$$

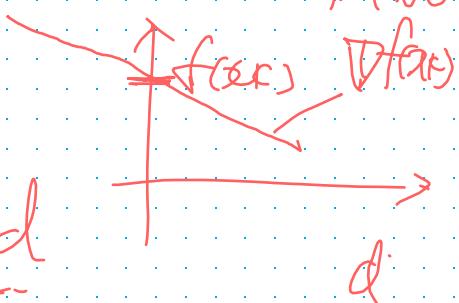
2次近似

$$\nabla f(x^*) = 0 \quad (\text{一次の最適性条件})$$

$$d = \underset{d}{\operatorname{argmin}} \quad f(x_k) + \nabla f(x_k) d + \frac{1}{2} \nabla^2 f(x_k) d^2 - J$$

$$\left(\frac{\partial J}{\partial d} = \nabla f(x_k) + \nabla^2 f(x_k) d = 0 \right) \quad \alpha \in \mathbb{R}$$

$$\Leftrightarrow d = -(\nabla^2 f(x_k))^{-1} \nabla f(x_k)$$



GD は 2 次近似

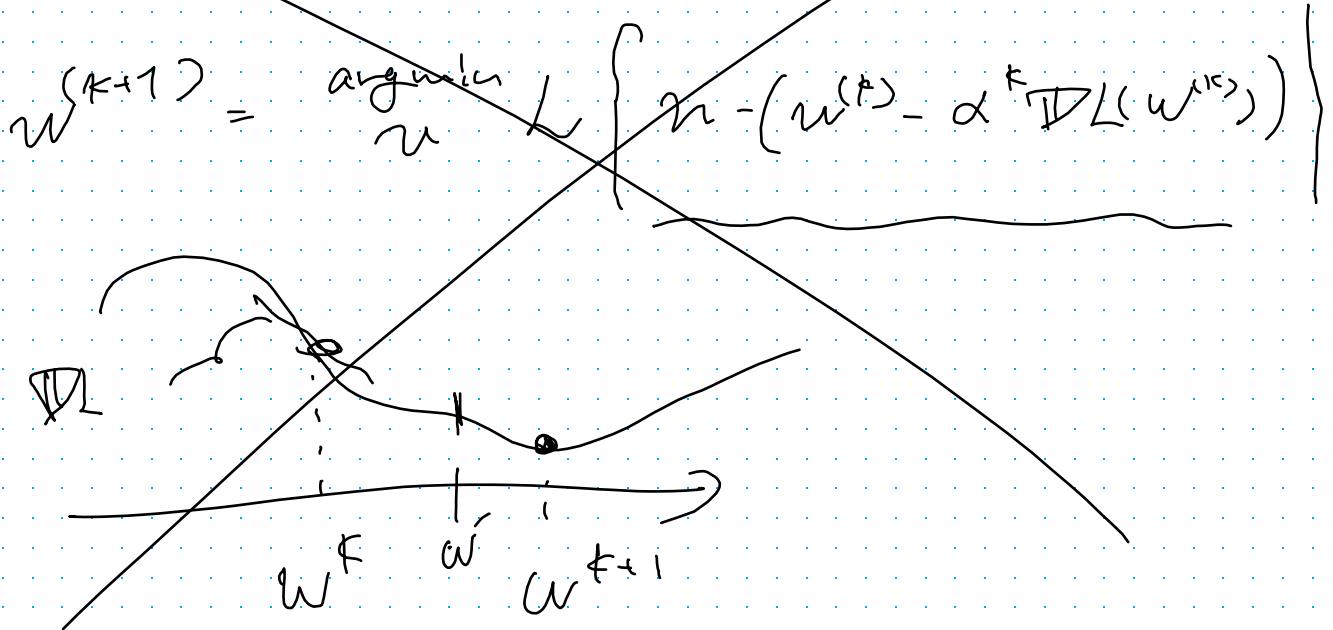
$$d = \underset{d}{\operatorname{argmin}} \quad f(x_k) + \nabla f(x_k) d$$

$$\nabla f(x_k) > 0 \quad \operatorname{Sign}(d) = -1,$$

Attendants

$$\bar{w}^* = \underset{w}{\operatorname{arg\,min}} \mathcal{L}(w)$$

~~梯度下降法~~



$$d = x_k - x^*$$

$$f(x) = f(\alpha) + f'(\alpha)(x - \alpha) + \dots$$

$$\alpha \leftarrow x_k - \text{ReLU}$$

$$f(x) = f(x_k) + \nabla f(x_k)(x - x_k)$$

$$x \leftarrow x_k + t d_k \in \mathbb{R}$$

$$f(x_k + t d_k) = f(x_k) + \nabla f(x_k) t d_k$$

Attendents

Iterative Method

$$\frac{f(x_k + \delta k) - f(x_k)}{\delta k} = Df(x_k) \delta k < 0$$

\Rightarrow もとよりの探索方向 $\delta k \leftarrow$

$f(x) \rightarrow$ 以下方向 と おこう

< Line - Search >

local & global convergence

q-linear convergence

$$\frac{\|x_k - x^*\|}{\|x_{k+1} - x^*\|} < c$$

q-quadratic convergence

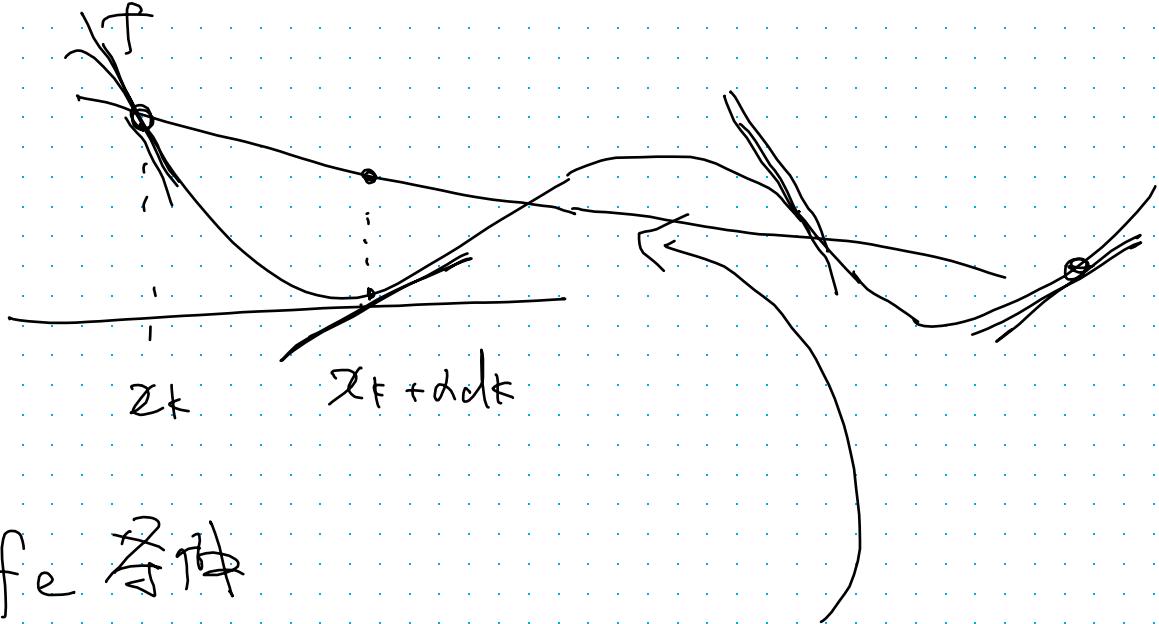
$$\frac{\|x_k - x^*\|}{\|x_{k+1} - x^*\|^3} < c$$

q-superlinear convergence

$$\lim_{k \rightarrow \infty} \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} = 0$$

α の探索

Armijo 条件



Wolfe 条件

$$f(x_k + \alpha d_k) \leq f(x_k) + \xi_1 \alpha \nabla f(x_k)^T d_k$$

$$= \text{true}$$

$$\xi_2 \nabla f(x_k)^T d_k \leq \nabla f(x_k + \alpha d_k)^T d_k$$

Zoutendijt 李博

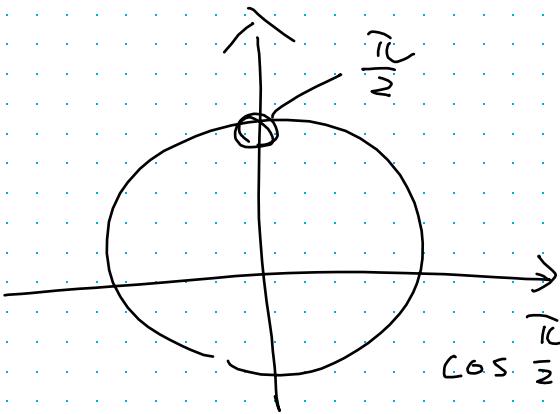
Assumption

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$$

Search direction

$$d_k \text{ s.t. } \nabla f(x_k)^T d_k < 0$$

"



$$\|\nabla f(x_k)\| \|d_k\| \cos \varphi < 0$$

$$\cos \frac{\pi}{2} = 0 \text{ 且 } \varphi \neq \frac{\pi}{2}$$

$$d_k = \arg \min_{d_k} \nabla f(x_k)^T d_k$$

$$= -\nabla f(x_k)$$

Use Wolfe condition

$$\sum_{k=0}^{\infty} \left(\frac{(\nabla f(x_k))^T d_k}{\|d_k\|} \right)^2 < \infty$$

$$\int_0^{\infty} (\|Df(x_k)\| \cos \varphi)^2 < \infty$$

for

Zoutendijk (ズーテンディック条件)

$$\lim_{k \rightarrow \infty} \frac{\|Df(x_k)\| \cos \varphi}{\|dx_k\|} = 0 \quad \leftarrow \begin{array}{l} \text{上記級数の} \\ \text{収束条件} \end{array}$$

$$\Leftrightarrow \lim_{k \rightarrow \infty} \|Df(x_k)\| \cos \varphi = 0$$

$\varphi \in [0, \pi]$ かつ $\cos \varphi > 0$ のとき

存在する ε が存在する

$$\lim_{k \rightarrow \infty} \|Df(x_k)\| \rightarrow 0$$

生成または更新

下記の図形を参考

Attendents

$$\cos \theta = \frac{-\nabla f(x_k)^\top d_k}{\|\nabla f(x_k)\| \|d_k\|}$$

$$*) \cos \theta = \frac{(-\nabla f(x_k))^\top}{\|(\nabla f(x_k))\|^2} = 1$$

最急降下法の収束率

$$\|x_{k+1} - x^*\|_A \leq \frac{\lambda_1 - \lambda_n}{\lambda_1 + \lambda_n} \|x_k - x^*\|_A$$

 $0 < \lambda_1 \leq \lambda_n$ は A の各行の A の

最小固有値と最大固有値

$$\|v\|_A = \sqrt{v^\top A v}$$

$$\frac{|\lambda_1 - \lambda_n|}{|\lambda_1 + \lambda_n|} = \left| \frac{\frac{\lambda_1}{\lambda_n} - 1}{\frac{\lambda_1}{\lambda_n} + 1} \right| = 1 \text{ となる} \rightarrow$$

左辺分子
分子分子
分子分母
分子分子

右辺分子
分子分子
分子分母
分子分子

Attendants

strong convex \rightarrow H is positive definite

$$\text{cig}(\nabla^2 f(w)) = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \ddots & 0 \\ 0 & & \lambda_n \end{pmatrix}$$
$$= Q \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} Q^T$$

where Q is orthogonal

$$z = Q^T w$$

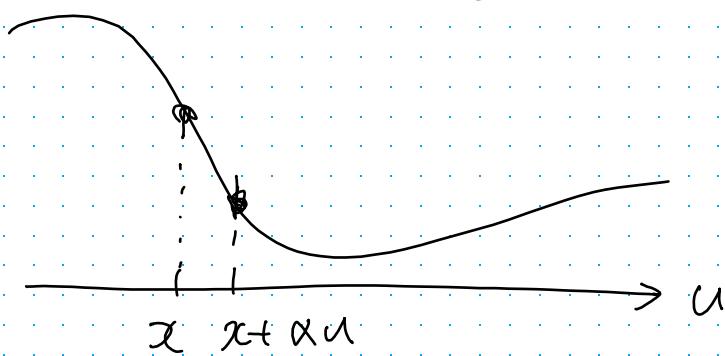
Attendents

方向微分 (directional derivative)

いわゆる f の傾きを表す

$$\frac{\partial}{\partial \alpha} f(x + \alpha u) \quad \alpha = 0 \text{ で } u^T D_x f(x)$$

chain rule $\frac{\partial f}{\partial \alpha} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial \alpha}$



$$= D_x f(x + \alpha u)$$

$$\times \frac{\partial}{\partial x} (x + \alpha u)$$

$$= u^T D_x f(x + \alpha u)$$

の定義は $\lim_{\alpha \rightarrow 0} \frac{f(x + \alpha u) - f(x)}{\alpha u} \quad \alpha \rightarrow 0 \text{ で}$

\downarrow
 $u^T D_x f(x)$

f を最小化する方向を考える

方向微分 $\Leftrightarrow f(x + du) \rightarrow d \parallel D_x f(x)$ で $\alpha = 0$ で

~~1つ目~~

\downarrow
 $\frac{\partial f}{\partial \alpha} = u^T D_x f(x)$

$$\min_{U, U^T U = 1} u^T D_x f(x) = \min_{U, \|U\|_2 = 1} \|U^T D_x f(x)\|_2 \cos \theta$$

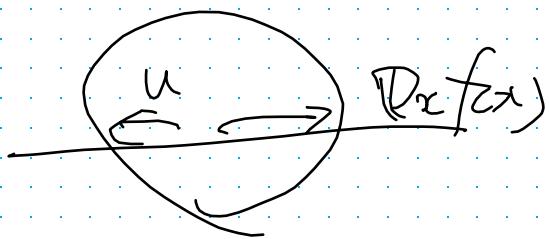
$$\|U\|_2 = \gamma^{\pm}$$

$$\min_u \|Df(x)\|_2 \cos \theta$$

↑

u は 依存しない

$$= \min_u \cos \theta \quad \text{条件 } u \in Df(x) \neq 0$$



反対側

$$\cos \theta = -1$$

$$\theta = 180^\circ$$

$$-1 - 1 = -2$$

したがって γ^-

$$x_{t+1} = x_t - \epsilon Df(x)$$

Jacobian

Matrix $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$

$$J \in \mathbb{R}^{n \times m} \quad J_{ij} = \frac{\partial}{\partial x_j} f_i(x)$$

$$\nabla^2 f(x) = \frac{\partial^2}{\partial x_i \partial x_j} : \text{curvature}$$

関数の導入の高次元の場合、2階微分の多次元存在する

\hookrightarrow 二重平均を Hessian Matrix と呼ぶ

$$H(f)(x)_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(x)$$

Jacobian of gradient = Hessian



実軸の点] \rightarrow 固有値の集合と固有ベクトルの
直交基底に分解する

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2} f''(x_0)(x - x_0)^2$$

新しい点 x (は $x_0 + g$ とする)

$$\begin{aligned} f(x_0 + \varepsilon g) &\approx f(x_0) - g^T \varepsilon g + \frac{1}{2} (\varepsilon g)^T H \\ &\approx f(x_0) - \varepsilon g^T g + \frac{1}{2} \varepsilon^2 g^T H g \end{aligned}$$

$g^T H g > 0$ の場合、最適な step 横

$$\epsilon^* = \frac{g^T g}{g^T H g}$$

最悪、場合は g が H の最大固有値 λ_{\max} に対応する
固有ベクトルと同一方向の場合

→ $\frac{1}{\lambda_{\max}}$ の step 横

Hessian

条件数が悪く場合、GD はあまり機能しない。

「ある方向 (= 慎激) で g が下がりづらい」と

「(= かゆい) g が上がる」

→ 勾配が長く負うとある方向を優先的に探索すべきと判断される。

二階導関数

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2} f''(x_0)(x - x_0)^2$$

$$f(x_0 + d) \approx f(x_0) + f'(x_0)d + \frac{1}{2} f''(x_0)d^2$$

$$\frac{\partial f}{\partial d} = f'(x_0) + f''(x_0)d = 0 \text{ のとき } f(x_0 + d) \text{ の } \frac{\partial f}{\partial d} \text{ が } 0.$$

$$f(x_0) + f'(x_0)d = 0$$

$$\begin{aligned} \Leftrightarrow d &= -f'(x_0)^{-1} f'(x_0) \\ &= -Df(x_0)^{-1} Df(x_0) \end{aligned}$$

Lipschitz Continuous

$$\forall x, \forall y, |f(x) - f(y)| \leq L \|x - y\|_2$$

Karush-Kuhn-Tucker 理論

constrained optimization (約束最適化)

一般的な解説

generalized Lagrangian
Lagrange function

$$\begin{aligned} L(x, \lambda, \alpha) &= f(x) + \sum_i \lambda_i g^{(i)}(x) + \sum_j \alpha_j h^{(j)}(x) \\ \text{s.t. } g^{(i)}(x) &= 0, h^{(j)}(x) \leq 0 \end{aligned}$$

Attendents

線形二乗法

最小

$$f(x) = \frac{1}{2} \|Ax - b\|_2^2$$

$$\nabla f(x) = A^T(Ax - b) = A^TAx - A^Tb$$

$$P(X) = \prod_{i=1}^n P(x_i | x_1, \dots, x_{i-1})$$

表面的 (i=1, 2, \dots, n) の分解にすることで $P(x)$ の教師学習で

なす \rightarrow supervised task は分解して解く。

$$P(y|x) = \frac{P(x,y)}{\sum_{y'} P(x,y')}$$

Linear Regression $\hat{y} = w^T x$

Mean Square Error

$$MSE_{test} = \frac{1}{m} \sum_i (\hat{y}^{(test)} - y^{(test)})^2$$

$$\nabla_w MSE_{train} = 0$$

Attendants

$$y = w^T x$$

$$\nabla_w \text{MSE}_{\text{train}} = 0$$

$$\Rightarrow \nabla_w \frac{1}{m} \| \hat{y} - y \|_2^2 = 0$$

$$\Leftrightarrow \nabla_w (w^T x - y)^T (w^T x - y) = 0$$

$$\Leftrightarrow x (w^T x - y) + (w^T x - y)^T x$$

$$\Leftrightarrow 2w^T x^2 - 2yx = 0$$

$$\Leftrightarrow w = yx (x^{-2})$$

$$= (X^T X)^{-1} X^T Y$$

Memo ① Zoom w/ Joannis

$$l(\hat{y}, y) = (\hat{y} - y)^2$$

$$f_i(\theta^*) = 0, \nabla f_i(\theta^*) = 0 \leftarrow \text{assumption 1}$$

$$f_i(\theta) = l(f_i(x_i), y_i) = (\hat{y} - y)^2$$

$$\nabla_{\theta} f_i(\theta) = 2 (\hat{y} - y)$$

finite sum
empirical risk

\leftarrow answer
G

$$y = x^3 - x$$

$$= J(x^2 - 1)$$

output
det

$f(\theta)$

Memo @Zoom w/ Journals

$$\frac{\partial}{\partial z} (y_i - z)^2 = 2(y_i - z)$$

$$z = f^*(x_i) = y_i$$

$$\rightarrow 0$$

interpolation $\hat{y} = f$

wake it

If grad $\neq 0 \rightarrow$

no constraints at problem

then we could

\Rightarrow Smaller value of Loss

Subject

Date

Attendants

Memo @ Zoom w/ Joann's

interpolation

possible with N

at least 2 layers

if we make it
wide enough

$P_\theta(y_i, f_\theta(x_i))$

Subject

Date

Attendants

Memo ② Zoom w/ Joann's
gradient LS always exists.



w/o Assum 1.

grad $\neq 0$



(oss given shell)

contrary (oss $\neq 0$).