$$X^T y = X^T X \beta$$

수선의 발.

$\beta \frac{?}{}$를 어떻게 찾는가?

---

## Least Square Estimate

$$(loss) = \| y - X\beta \|^2$$

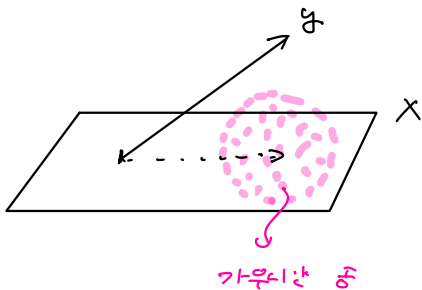$$= (y - X\beta)^T (y - X\beta)$$

$$= (y^T - \beta^T X^T)(y - X\beta)$$

$$= y^T y - y^T X\beta - \beta^T X^T y + \beta^T X^T X\beta$$

$$\nabla_\beta (loss) = -2 X^T y + 2 X^T X \beta = 0 \quad (극값)$$

hence $\quad X^T y = X^T X \beta \quad$ 수선의 발

$$\nabla_\beta \, \beta^T A \beta = \nabla_\beta \sum_{i,j} \beta_i A_{ij} \beta_j$$

$$= (A + A^T) \beta$$

---

## Maximum Likelyhood Estimate



가우시안 츰

let $Y = X\beta + \varepsilon \sim N(X\beta, \, \sigma^2 I)$

then $\quad p(Y) = \dfrac{1}{(2\pi)^{d/2} \cdot \sigma^d} \cdot \exp\left( \dfrac{-1}{2\sigma^2} \| y - X\beta \|^2 \right)$

likelyhood

squared L2 Norm

$p(\theta)$ : prior

$p(data|\theta)$ : likelyhood

$p(data)$ : evidence.

$p(\theta|data)$ : posterior $= \dfrac{p(data|\theta) \cdot p(\theta)}{p(data)}$

Bernoulli R.V     "동전 던지기"

$$p(x_i|\theta) = \begin{cases} \theta & x_i = 1 \\ 1-\theta & x_i = 0 \end{cases}$$

$$= \theta^{x_i}(1-\theta)^{1-x_i}$$

---

likelyhood 의 한계.

$$p(x|\theta) = \theta^{\sum x_i}(1-\theta)^{n-\sum x_i}$$

$$\log p(x|\theta) = \sum x_i \log\theta + (n-\sum x_i)\log(1-\theta)$$

$$\nabla_\theta \log p(x|\theta) = \frac{\sum x_i}{\theta} - \frac{(n-\sum x_i)}{(1-\theta)}$$

say     $\nabla_\theta \log p(x|\theta) = 0$     (MLE)

then ....     $\theta = \dfrac{\sum x_i}{n}$

---

$Beta(\alpha, \beta) = x^{\alpha-1} \cdot (1-x)^{\beta-1} \cdot \dfrac{\Gamma(\alpha+\beta)}{\Gamma(\alpha) \cdot \Gamma(\beta)}$

where     $0 \leq x \leq 1$

Note.

$\Gamma(n) = (n-1)!$

---

posterior 의 근사.

실험결과  Ⓗ Ⓣ
          7  3

$\nabla_\theta \log p(\theta|X) = 0$     (MAP)

Beta(1,1)

$p(\theta|X) \propto 1 \cdot \theta^7(1-\theta)^3 = \theta^7(1-\theta)^3$
              $p(\theta)$   $p(X|\theta)$
                            Beta(8,4)

Uniform
Distribution 이니까

$\Longrightarrow$     $\theta_{MAP} = \theta_{MLE}$
                      $= 0.7$

Beta(2,2)

$p(\theta|X) \propto \theta^1 \cdot (1-\theta)^1 \quad \theta^7(1-\theta)^3 = \theta^8(1-\theta)^4$
                $p(\theta)$        $p(X|\theta)$
                                    Beta(9,5)

$\Longrightarrow \theta_{MAP} = 0.666...$

- $h(x) = -\log p(x)$

- if $x, y$ independent

  $h(x, y) = h(x) + h(y)$

- $H(X) = E[-\log P(X)] \geq 0$

- $H(X, Y) = H(X) + H(Y|X)$

(추가)    $H \uparrow$        $H \Downarrow$

       동전 vs 주사위     동전 vs 찍기

---

## KL Divergence

$D_{KL}(P \| q)$

$$= \sum_x p(x) \log \frac{p(x)}{q(x)}$$

$$= \int p(x) \log \frac{p(x)}{q(x)} \, dx \geq 0$$

(X) symmetric

(X) triangular inequality

### (증명)

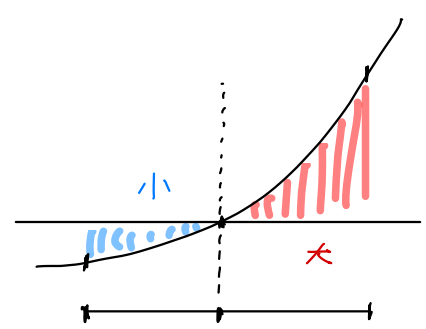$$\mathop{E}_{p(x)}\left[-\log \frac{q(x)}{p(x)}\right]$$

$$\geq -\log \mathop{E}_{p(x)}\left[\frac{q(x)}{p(x)}\right]$$   Const

$$= -\log \sum_x p(x) \cdot \frac{q(x)}{p(x)} \geq 0$$   Const $= 1$

---

## Jensen's Inequality

$f$ convex

then $E[f(x)] \geq f(E(x))$

$\Downarrow f(x) = $ Const



---

## Mutual Information

$I(X : Y) = D_{KL}(P(x, y) \| P(x) \cdot P(y))$

$$= \sum p(x, y) \log \frac{p(x, y)}{p(x) p(y)}$$

$$= H(X) + H(Y) - H(X, Y)$$

$$= H(X) - H(X|Y)$$

## Cross Entropy

- $H_p(q) = \mathop{E}_{p}[-\log q(x)]$

$$= -\sum_x p(x) \log q(x)$$

- $D_{KL}(P \| q) = \mathop{E}_{p}\left[\log \frac{p(x)}{q(x)}\right]$

$$= H_p(q) - H(P) \geq 0$$

           fixed

# overfitting

$$y = \begin{bmatrix} 1 & x^1 & \cdots & x^m \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_m \end{bmatrix}$$

단일 변수에 의존 (for $\begin{bmatrix} 1 & x^1 & \cdots & x^m \end{bmatrix}$)

불필요하게 큰 차원 (for the $\beta$ vector) $\Rightarrow$ overfitting

---

## Regularization

$$(loss) = \| Y - X\beta \|^2 + \boxed{\lambda \| \beta \|^2} \quad \text{Regularization}$$

$$\nabla_\beta (loss) = \nabla_\beta \left( y^T y - y^T X\beta - \beta^T X^T y + \beta^T X^T X\beta \right) + \nabla_\beta \, \lambda \beta^T \beta$$

$$= 2 \begin{bmatrix} \underbrace{-X^T(Y - X\beta)}_{\text{데이터 의존}} + \underbrace{\lambda \beta}_{\text{단순한 조절}} \end{bmatrix} = 0$$

---

## Kernel

$$K(x_1, x_2) = \underbrace{\phi(x_1)}_{\text{low dimension}}{}^T \underbrace{\phi(x_2)}_{\text{high dimension}}$$
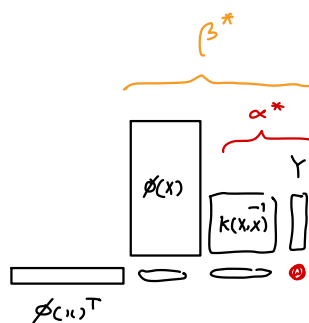
$$(loss) = \| Y - \phi(X)^T \beta \|^2 + \lambda \| \beta \|^2$$

$$= y^T y - y^T \phi(X)^T \beta - \beta^T \phi(X) Y + \beta^T \phi(X) \phi(X)^T \beta + \lambda \beta^T \beta$$

$$\nabla_\beta (loss) = 2 \begin{bmatrix} -\phi(X)(Y - \phi(X)^T \beta) + \lambda \beta \end{bmatrix} = 0$$

$$\beta^* = \left( \phi(X) \phi(X)^T + \lambda I \right)^{-1} \phi(X) Y$$

$$= \phi(X) \underbrace{\left( \phi(X)^T \phi(X) + \lambda I \right)^{-1} Y}_{= \alpha^*}$$

$$f^*(x) = \phi(x)^T \beta^*$$

$$= K(x, X) \alpha^*$$

**Sigmoid** 임의의 값을 확률 형태로 변환해줌.

↑

if probability $= \dfrac{e^z}{1+e^z}$    $[0, 1]$

then odds $= e^z$    $[0, \infty)$

log(odds) $= z$    $(-\infty, \infty)$

Note that $\sigma' = \sigma(1-\sigma)$



Sigmoid Function

$$\sigma(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1}$$

# Logistic Regression (MLE)

$$\sigma' = \sigma(1-\sigma)$$



아래로 당기면
신뢰도 상승

$$\nabla_w \sigma_i = \sigma_i(1-\sigma_i)\nabla_w(w^T x_i)$$

$$= \sigma_i(1-\sigma_i)\, x_i$$

$$\propto x$$

$$p(y_i | x_i, w) = \left(\sigma(w^T x_i)\right)^{y_i}\left(1 - \sigma(w^T x_i)\right)^{1-y_i}$$

$$\mathcal{L} = \prod_{i=1}^{\tilde{N}} p_i$$
$$\downarrow \text{data}$$

$$\log p_i = \left(y_i \log \sigma_i\right) + \left((1-y_i)\log(1-\sigma_i)\right)$$

$$\nabla_w(-\log p_i) = \left(y_i(\sigma_i - 1)x_i\right) + \left((1-y_i)\sigma_i x_i\right)$$

$$= \left(\begin{array}{l} y_i \sigma_i x_i \\ \quad - y_i x_i \end{array}\right) + \left(\begin{array}{l} \sigma_i x_i \\ \quad - y_i \sigma_i x_i \end{array}\right)$$

$$= (\sigma_i - y_i)\, x_i$$

예측   실제

$$\text{hence} \quad \nabla_w(-\log \mathcal{L}) = \sum_i \left(\sigma(w^T x_i) - y_i\right) x_i$$

예측   실제   최소정점 위치

---

$$\nabla_w \nabla_w(-\log p_i) = \left[\nabla_w(\sigma_i - y_i)\right] x_i^T$$

$$= \sigma_i'\, x_i\, x_i^T$$

$$\nabla_w \nabla_w(-\log \mathcal{L}) = \begin{bmatrix} | & \\ x_i & \cdots \\ | & \end{bmatrix} \begin{bmatrix} \sigma_1' & & \\ & \ddots & \\ & & \sigma_N' \end{bmatrix} \begin{bmatrix} - x_i - \\ & \ddots \end{bmatrix}$$

$X^T$    $S$    $X$

$\geq 0$   $\geq 0$

because $\sigma' = \sigma(1-\sigma) \geq 0$

for any $v$

$$v^T(X^T S X)v$$

$$= (Xv)^T S (Xv)$$

$$= \sum_i \left(\sigma_i'\right)\left((Xv)_i\right)^2 \geq 0$$

$\geq 0$    $\geq 0$

positive semi-definite $\longleftrightarrow$ $-\log \mathcal{L}$ convex

$\langle v, Hv \rangle \geq 0$    有 global min

# Multiclass classification

let $z_k = w_k^T x$



$$\nabla_{w_m}(w_k^T x) = \begin{cases} x & \text{if } m=k \\ \odot & \text{if } m \neq k. \end{cases}$$

$$= \delta_{mk} \cdot x$$

then $p_k = \dfrac{e^{z_k}}{\sum_j e^{z_j}}$

$$\nabla_{w_m}(p_k) = \frac{\sum_j e^{z_j}\left(e^{z_k} \cdot (\nabla z_k)\right) - e^{z_k}\left(\sum_j e^{z_j} \nabla z_j\right)}{\left(\sum_j e^{z_j}\right)^2}$$

where $\overbrace{\delta_{mk} \cdot x}$ and $\overbrace{e^{z_m} \cdot x}$

$$= p_k(\delta_{mk} - p_m)\, x$$

$$\nabla_{w_m}\left(-\log \begin{array}{l} \text{likelyhood} \\ \text{of } x \end{array}\right) = \sum_{k=1}^{c} y_k \nabla(-\log p_k)$$

$$= \sum_{k=1}^{c} y_k \cdot \frac{1}{p_k} \cdot p_k(p_m - \delta_{mk})\, x$$
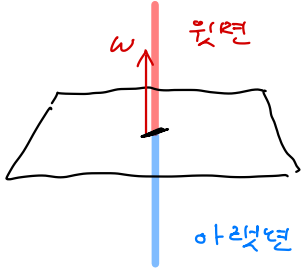
$$= (p_m - y_m)\, x$$

$$\nabla_{w_m}(-\log \mathcal{L}) = \sum_{i=1}^{N} (p_{im} - y_{im})\, x_i$$

SVM   경정면   $w^T x + b = 0$



조건 ①

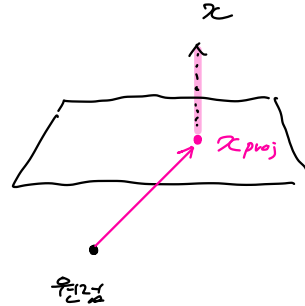$\forall i \quad y_i (w^T x_i + b) \geq 0$

$y = 1$
$w^T x + b \geq 0$

윗면

아랫면 $\quad w^T x + b \leq 0$

$y = -1$

조건 ②

$$\min_x |w^T x + b| = 1$$

$x = x_{proj} + d \dfrac{w}{\|w\|}$

$d = \dfrac{|w^T x + b|}{\|w\|}$

$= \dfrac{1}{\|w\|}$ (가장 가까운 점)

$x_{proj}$

원점

---

$\dfrac{1}{2} \|w\|^2$           $y_i (w^T x_i + b) \geq 1$

minimize   $\dfrac{1}{2} \|w\|^2 + \max_{\alpha_i \geq 0} \alpha_i (1 - y_i (w^T x_i + b)) = \begin{cases} 0 & y_i (w^T x_i + b) > 1 \\ \infty & else \end{cases}$

---

$\min_{w, b} \max_{\alpha_i \geq 0} \left[ \dfrac{1}{2} \|w\|^2 + \sum_i \alpha_i (1 - y_i (w^T x_i + b)) \right] \geq \max_{\alpha_i \geq 0} \min_{w, b} \left[ \,\, '' \,\, \right]$

$= \dfrac{1}{2} w^T w + \sum_i \alpha_i - w^T \left( \underset{w}{\underline{\sum_i \alpha_i y_i x_i}} \right) - b \underset{0}{\underline{\sum_i \alpha_i y_i}}$

maximize

$\sum_i \alpha_i - \dfrac{1}{2} w^T w$

with $\begin{cases} \nabla_w \text{⑭} = w - \sum_i \alpha_i y_i x_i = 0 \\ \dfrac{\partial}{\partial b} \text{⑭} = \sum_i \alpha_i y_i = 0 \end{cases}$

say $i \in I$ whenever $\alpha_i > 0$

$$\text{i.e } y_i(w^T x_i + b) = 1$$

then

$$w = \sum_i \alpha_i y_i x_i \qquad (i \in I)$$

$$b = y_i - w^T x_i \qquad (i \in I)$$

$$w^T w = \sum_i \alpha_i y_i (w^T x_i) \qquad (i \in I)$$

$$= \sum_i \alpha_i (1 - y_i b) \qquad \text{because } y_i(w^T x_i + b) = 1$$

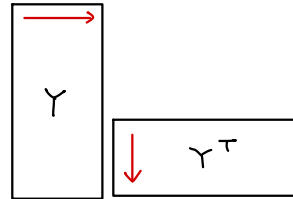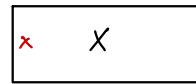$$= \sum_i \alpha_i - b \sum_i \alpha_i y_i$$

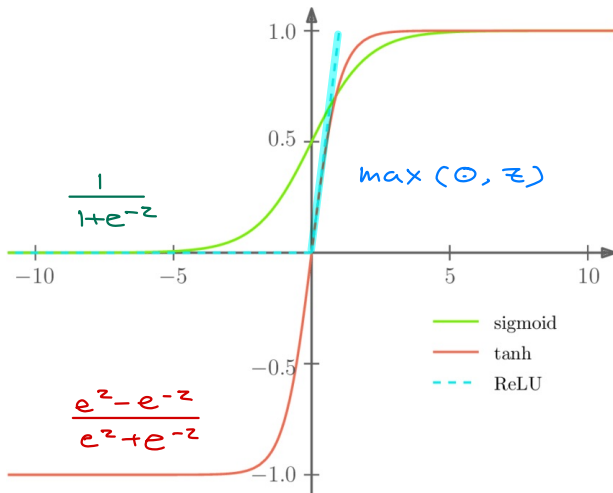$$\frac{|x^T w + b|}{\|w\|} = \frac{1}{\sqrt{\sum \alpha_i}}$$

say $W = XY$

Then $\dfrac{\partial L}{\partial X} = \dfrac{\partial L}{\partial W} \cdot Y^T$

$\dfrac{\partial L}{\partial Y} = X^T \dfrac{\partial L}{\partial W}$

---

Gradient Descent

minimize

$\theta_{k+1} = \theta_k - \alpha \nabla L(\theta_k)$

$\begin{cases} \text{Full Batch GD} & \tilde{\nabla} f(x) = \dfrac{1}{N} \sum_{i=1}^{N} \nabla f_i(x) \\ \\ \text{Stocastic GD} & \tilde{\nabla} f_{(x)} = \nabla f_i(x) \\ \\ \text{Mini Batch GD} \end{cases}$

---



$\dfrac{1}{1+e^{-z}}$

$\max(0, z)$

$\dfrac{e^z - e^{-z}}{e^z + e^{-z}}$

- sigmoid
- tanh
- ReLU

$\sigma' = \sigma(1-\sigma)$

$(\tanh)' = 1 - (\tanh)^2$

$(ReLU)' = \begin{cases} 1 & z > 0 \\ 0 \end{cases}$

Note. $\sigma$ 나 $\tanh$ 는

(CH$\frac{1}{2}$) linear

let $z = \sum_{i=1}^{Din} w_i x_i$     s.t
$$\begin{cases} w : \text{i.i.d}, & E[w_i] = 0 \quad Var[w_i] = \delta^2 \\ x : \text{i.i.d} & \text{independent of } w \\ & E[x_i] \neq 0 \quad (\text{in general}) \end{cases}$$

let $S = E[x^2]$   denote the 2nd moment of the input.

---

we know that $E[z] = 0$

then $Var(z) = \sum_{i=1}^{Din} Var(w_i x_i)$

$= \sum_{i=1}^{Din} E[w_i^2 x_i^2] - E[w_i x_i]^2 = Din\, \delta^2 S$

$\underbrace{E[w_i^2] E[x_i^2]}_{\delta^2 \quad S} \qquad \underbrace{E[w_i]^2 E[x_i]^2}_{0}$

in Xavier Init we want $E[z^2] = E[x^2]$

$$\cancel{Din\, \delta^2}\; S = S \qquad \delta = \frac{1}{\sqrt{Din}}$$

Note. independent

$\iint xy\, p(x,y)\, dx\, dy$

$= \iint xy\, p(x) p(y)\, dx\, dy$

---

now. let $h = ReLU(z)$

then $E[h^2] = \int_0^\infty z^2 p(z)\, dz$

$= \frac{1}{2} \int_{-\infty}^\infty z^2 p(z)\, dz = \frac{1}{2} E[z^2]$

$= \frac{1}{2} Din\, \delta^2 S$

in He Init we want $E[h^2] = E[x^2]$

$$\frac{1}{2} Din\, \delta^2 S = S \qquad \delta = \sqrt{\frac{2}{Din}}$$

$$\begin{cases} Xavier & \sigma_{ref}^2 = \frac{1}{Din} \\ He & \sigma_{ref}^2 = \frac{2}{Din} \end{cases}$$

then $w \sim N(0, \sigma_{ref}^2)$

or $w \sim Uni(-\sqrt{3}\sigma_{ref}, \sqrt{3}\sigma_{ref})$

as variance of $Uni(-b,b)$ becomes

$\int_{-b}^{b} x^2 \cdot \frac{1}{2b}\, dx$

$= \left[ \frac{x^3}{6b} \right]_{-b}^{b} = \frac{b^3}{3}$

leaning rate decay

① at a few fixed points ?

② cosine : $\frac{1}{2}\alpha_0\left(1 + \cos\left(\frac{t\pi}{T}\right)\right)$

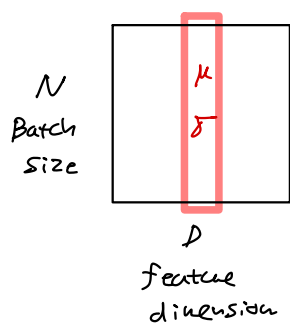③ linear : $\alpha_0\left(1 - \frac{t}{T}\right)$

④ linear sqrt : $\alpha_0/\sqrt{t}$

linear warmup : linearly increase
learning steps at $t \approx 0$

---

Batch Normalization

(training)



$$\tilde{x}_{ij} = \frac{x_{ij} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

$$\tilde{y}_{ij} = \gamma_j \tilde{x}_{ij} + \beta_j$$

$N$
Batch Size

$D$
feature dimension

$\mu^{run} = m \cdot \mu^{run} + (1-m)\mu^{batch}$

$(\sigma^2)^{run} = m(\sigma^2)^{run} + (1-m)(\sigma^2)^{batch}$

BN ⇒ Internal Convergence Shift : shift in the mean/var of hidden activation

⇒ Smooth loss landscape : Stable gradient.

---

GD minimizes $f(x)$

$x_{t+1} = x_t + \gamma v$

to minimize $f(x_{t+1}) \approx f(x_t) + \langle \nabla f(x_t), \gamma v \rangle$

set $v = -\frac{\nabla f(x_t)}{\|f(x_t)\|}$  GD idea

**Lemma 3.1**

- $f : \mathbb{R}^d \to \mathbb{R}$    continuously diffable function.

- let $f$   $\beta$-smooth :   $|\nabla f(x) - \nabla f(y)| \le \beta\|x - y\|$    $(\forall x, y)$

then    $f(y) \le f(x) + \underbrace{\langle f(x), y - x \rangle}_{\text{중간단}} + \underbrace{\frac{\beta}{2}\|y - x\|^2}_{\text{스무스}}$

---

(GD) :   $x_{t+1} = x_t - \eta \nabla f(x_t)$

$f(x_{t+1}) \le f(x_t) + \langle \nabla f(x_t), -\eta \nabla f(x_t)\rangle + \frac{\beta}{2}\|\eta \nabla f(x_t)\|^2$

$= (\cdots) = f(x_t) - \left(\eta - \frac{\beta}{2}\eta^2\right)\|\nabla f(x_t)\|^2$

$\qquad\qquad\qquad\qquad\quad \text{small step}$

---

(SGD) :   $x_{t+1} = x_t - \eta \tilde{\nabla} f(x_t)$

then    $f(x_{t+1}) \le f(x_t) - \eta \langle \nabla f(x_t), \tilde{\nabla} f(x_t)\rangle + \frac{\beta}{2}\cdot \eta^2 \|\tilde{\nabla} f(x_t)\|^2$

$E_t[f(x_{t+1})] \le f(x_t) - \eta\|\nabla f(x_t)\|^2 + \frac{\beta}{2}\eta^2 \underbrace{E_t[\|\tilde{\nabla} f(x_t)\|^2]}_{=: G}$

$E[\|\nabla f(x_t)\|^2] \le \frac{1}{\eta}\Big(E[f(x_t)] - E[f(x_{t+1})]\Big) + \frac{\beta}{2}\eta^2 G$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \overset{\text{(ower bound)}}{}$

$\sum_{t=0}^{T-1} E[\|\nabla f(x_t)\|^2] \le \frac{1}{\eta}(f_0 - f^*) + \frac{\beta}{2}\eta^2 G T$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{let } \eta = \frac{1}{\sqrt{T}}$

$\min_t \Big(E\|\nabla f(x_t)\|^2\Big) \le \frac{1}{\eta T}(f_0 - f^*) + \frac{\beta}{2}\eta^2 G$    then   $\min\Big(E[\|\nabla f(x_t)\|^2\Big) = O\left(\frac{1}{\sqrt{T}}\right)$

$$g(t) = f\left(x + t(y-x)\right)$$

$$g'(t) = (y-x)^T \nabla f(x + t(y-x))$$

$$g''(t) = (y-x)^T \nabla^2 f(x + t(y-x))(y-x) \quad \leq \quad \beta \|y-x\|^2$$

---

$$\int_0^1 (1-s) g''(s) ds = \left[(1-s) g'(s)\right]_0^1 + \int_0^1 g'(s) ds$$

$$= g(1) - g(0) - g'(0)$$

Hence $f(y) = f(x) + \langle \nabla f(x), y-x \rangle + \underline{\int_0^1 (1-s)(y-x)^T \nabla^2 f(x + s(y-x))(y-x) ds}$

$$\leq \int_0^1 (1-s) \beta \|y-x\|^2 ds$$

$$= \frac{\beta}{2} \|y-x\|^2$$

---

let $\phi(t) = \nabla f(x + t(y-x))$

$$\phi'(t) = \nabla^2 f(x + t(y-x))(y-x)$$

Then $\nabla f(y) - \nabla f(x) = \int_0^1 \phi'(t) dt$

$$= \int_0^1 \nabla^2 f(x + t(y-x)) dt \cdot (y-x)$$

$$(y-x)^T \int_0^1 \nabla^2 f(x + t(y-x)) dt (y-x)$$

$$= \langle \nabla f(y) - \nabla f(x), y-x \rangle \leq \|\nabla f(y) - \nabla f(x)\| \|y-x\|$$

$$\leq \beta \|y-x\|^2$$

## GD

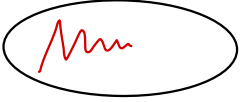$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$



상녹은함

## +Momentum

$$x_{t+1} = x_t + V_{t+1}$$

$$V_{t+1} = \rho V_t - \alpha \nabla f(x_t)$$

## Nesterov Momentum

$$x_{t+1} = x_t + U_{t+1}$$

$$V_{t+1} = \rho V_t - \alpha \nabla f(x_t + \rho V_t)$$

---

$\| \nabla f(x) \|$ . 가라을 정도.

" 가라르면 느리게

한받아면 아느게 "

## Ada Grad

$$x = x - \alpha \frac{\nabla f(x_t)}{\sqrt{\sum (\nabla f(x_t))^2}}$$

$$\to \infty$$

## RMSProp

$$x = x - \frac{\eta g_t}{\sqrt{E[g^2]_t}}$$

$$E[g^2]_t = \beta E[g^2]_{t-1} + (1-\beta) g_t^2$$

---

## Adam

$$x_{t+1} = x - \alpha \cdot \frac{\hat{m_t}}{\sqrt{\hat{V_t}}}$$

$$m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t \qquad m_0 = 0$$

$$V_t = \beta_2 V_{t-1} + (1-\beta_2) g_t \qquad V_0 = 0$$

$$\hat{m_t} = \frac{m_t}{1-\beta_1^t} \qquad \hat{V_t} = \frac{V_t}{1-\beta_2^t}$$

if $g$ const

then $\begin{cases} m_k = (1-\beta_1^k) g \\ V_k = (1-\beta_2^k) g^2 \end{cases}$