# Smooth optimiz°

$$\min_{x \in \mathbb{R}^d} f(x)$$
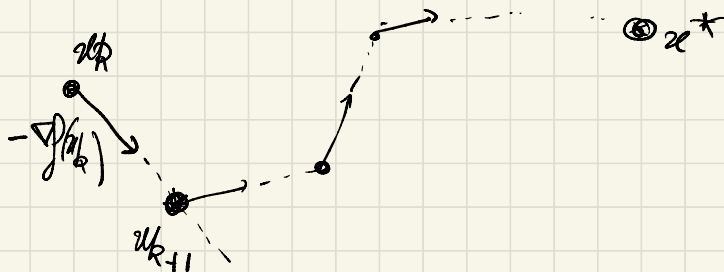
$f(x)$



$x^*$ $\to$ $x$

$\underline{\text{Gradient}}: \nabla f(x) \in \mathbb{R}^d$

Gradient Descent : (Batch)    $x_0 \leftarrow$ Init

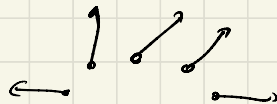$$x_{k+1} \triangleq x_k - \boxed{\tau} \cdot \nabla f(x_k)$$

$x_k$

$-\nabla f(x_k)$

$x_{k+1}$

$\circledcirc x^*$

If $x_k \to x^*$    $x^* = x^* - \tau \nabla f(x^*)$
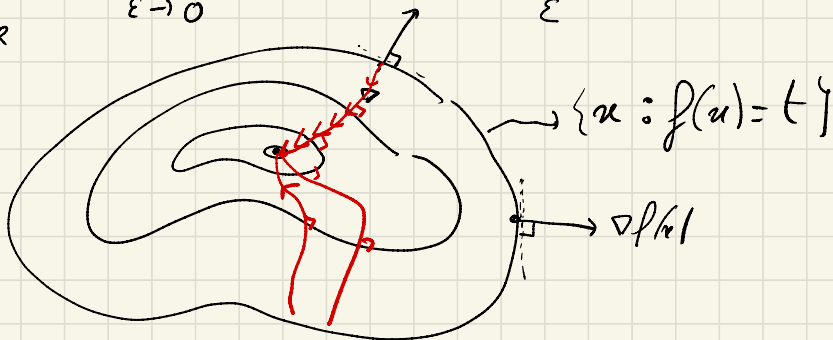
$\Rightarrow \nabla f(x^*) = 0$

① **Gradient** : $\nabla f(x) = \begin{pmatrix} \frac{\partial f}{\partial x_1}(x) \\ \vdots \\ \frac{\partial f}{\partial x_d}(x) \end{pmatrix} \in \mathbb{R}^d$

$f : x \in \mathbb{R}^d \to f(x) \in \mathbb{R}$

$\nabla f : x \in \mathbb{R}^d \to \nabla f(x) \in \mathbb{R}^d$

$$\frac{\partial f}{\partial x_k}(x) = \lim_{\varepsilon \to 0} \frac{\overbrace{f(x_1, \cdots x_{k-1}, x_k + \varepsilon, x_{k+1}, \cdots, x_d)}^{x + \varepsilon \delta_k} - f(x)}{\varepsilon}$$



$\{x : f(x) = t\}$

$\nabla f(x)$

"Smooth func°" : 1 time differentiable

Def : $f$ is diff. at $x$ if

$$f(x + \varepsilon v) = f(x) + \varepsilon \langle \nabla f(x), v \rangle + \underbrace{o(\varepsilon)}_{\xrightarrow[\varepsilon \to 0]{} 0}$$

$$\Updownarrow$$

$$\frac{f(x + \varepsilon v) - f(x)}{\varepsilon} \xrightarrow{\varepsilon \to 0} \langle \nabla f(x), v \rangle$$

$\triangle$ ! $Df(x)$ exists $\not\Rightarrow$ $f$ is diff.

$\Longleftarrow$

$$f(x,y) = \frac{xy\,(x+y)}{x^2+y^2} \qquad f(0)=0$$

Thm: if $Df(x)$ for $x$ in a ball around $x_0$
and $x \mapsto Df(x)$ is continuous
$\Longrightarrow$ $f$ is differentiable

Thm: if $x$ is a local minimizer
$\Longrightarrow$ $Df(x)=0$

Thm: if $f$ is convex; $x$ global min $\Longrightarrow$ $Df(x)=0$

1D : $f$ cvx $\iff$ $f''(x) \geq 0$ $\forall x$

$\partial^2 f(x) \succeq 0$

Composi: $\sum_{j=0}^{1} \underbrace{f(x)}_{cvx} + \underbrace{\mu \overset{\geq 0}{g(x)}}_{cvx}$ $\partial x$

$f \partial x,$ $\quad f(Ax+b)$ cvx

$A$ matrix

$\varphi : \mathbb{R} \to \mathbb{R}$ $\quad \varphi(f(x))$ $\partial x$

$\nearrow$ , $\partial x$

$\varphi(x) = \sum \varphi_i(x_i)$

$f : \mathbb{R}^d \to \mathbb{R}^p$ $\quad \underline{\varphi} : \mathbb{R}^p \to \mathbb{R}$

$f(x)$ cvx $\quad g(x,y) \overset{\downarrow}{=} y f\left(\frac{x}{y}\right)$ $\partial x$

perspective trans.

$KL(\underset{\uparrow}{x}|\underset{\uparrow}{y}) = \sum \left(\frac{x_i}{y_i}\right) \log\left(\frac{x_i}{y_i}\right) \times \underline{y_i}$

$\underline{ML}$: input : $(a_i^\circ, y_i)_{i=1}^M$   $a_i \xrightarrow{\psi} y_i$

$\underline{Supervised}$   $\underset{\mathbb{R}^d}{\uparrow}$   $\underset{\mathbb{R}}{\uparrow}$   $\psi(a) = y$.

$\{-1, +1\}$   $\psi(a_i) \approx y_i$

linear model   $\psi(a) = \langle a, \overset{\textcircled{\scriptsize 0}}{\underset{\uparrow}{x}} \rangle$

weight / need train

$\underline{Regression}$ : ERM

$$\min_{x \in \mathbb{R}} \sum_{i=1}^M \ell(\langle a_i^\circ, x \rangle, y_i^\circ) = f(x)$$

$\underset{\mathbb{R}}{\uparrow}$   $\underset{\mathbb{R}}{\uparrow}$

least square   $\ell(y, y') = (y - y')^2$

$\underline{Rmq}$ : if $\ell$ is convex, $f$ is convex

Design matrix :   $A = $



$\overset{\mathbb{R}^{M \times d}}{\underset{\cup}{}}$

$A \overset{\cup}{x} = \begin{pmatrix} \langle a_1, x \rangle \\ \langle a_2, x \rangle \\ \vdots \\ \langle a_m, x \rangle \end{pmatrix}$   $\overset{\mathbb{R}^d}{\underset{\cup}{}}$

$$\min_x \ f(x) = L(\overset{\large u}{\overbrace{Ax}})$$

$$L(u) = \sum_{i=1}^{m} \ell(u_i, y_i)$$

least square: $L(u) = \sum_{i=1}^{m} (u_i - y_i)^2 = \|u - y\|^2$
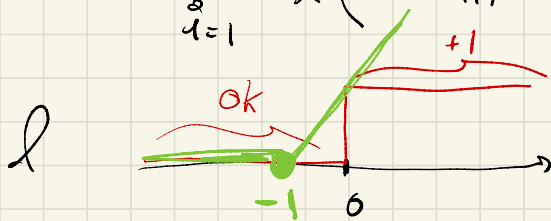
$$\min_x f(x) = \|Ax - y\|^2$$

Class$^{\circ}$: $y_i \in \{-1, +1\}$

predicta: $\text{sign}[\langle a, x \rangle] \in \{-1, 1\}$.

Ultimate goal: 0/1 loss.

$$\min_x \ \sum_{i=1}^{u} \text{Error}\left(\text{sign}(\langle a_i, x \rangle), y_o\right).$$
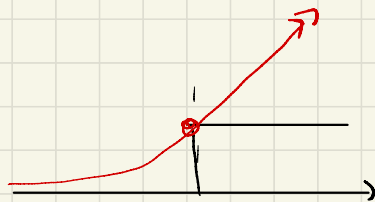
$$\sum_{i=1}^{u} \ell\left(-\langle a_i, x \rangle \ y_i\right) = \#error$$



$\ell(s) = \mathbb{1}_{\mathbb{R}^+}(s)$

① SVM: $\ell = $ HINGE LOSS $\quad \ell(s) \overset{\triangle}{=} (s+1)_+$

$f$ is non smooth

② Logistic loss / Cross entropy loss.



$$\ell(s) \overset{\Delta}{=} \frac{\log(1 + e^s)}{\log(2)}$$

$$\ell'(s) = \frac{e^s}{1 + e^s} \in [0,1]$$

$$f(x) = \sum_i \ell\left(-\, y_i \underbrace{\langle a_i, x \rangle}_{(Ax)_i}\right) = L(\underbrace{-\mathrm{diag}(y) A x}_{B})$$

$$\underbrace{-\mathrm{diag}(y) \cdot A \cdot x}$$

$$= L(Bx)$$

$$L(u) = \sum_{i=1}^n \ell(u_i)$$

Gradient comput° :   $\min_u f(x) = L(Bx)$.

$$\left[ x \in \mathbb{R}^d \underset{B \in \mathbb{R}^{n \times d}}{\longrightarrow} Bx \in \mathbb{R}^n \underset{\underset{\mathrm{diff.}}{\boxed{L}}}{\longrightarrow} L(Bx) \in \mathbb{R}\right.$$

"Proof": $\quad f(x + \varepsilon v) = (\cdots) = f(x) + \varepsilon \langle \underset{\underset{\nabla f(x)}{\parallel}}{??}, \boxed{v} \rangle + o(\varepsilon)$

$L(B(x+\varepsilon v)) = L(\underline{\underline{Bx}} + \varepsilon Bv) \quad\rangle \quad$ L is diff at $Bx$

$\overset{\circledast}{=} \underset{f(x)}{\underline{L(Bx)}} + \varepsilon \langle \nabla L(Bx), \boxed{Bv} \rangle + o(\varepsilon)$

Def: $B = (B_{ij})_{ij} \quad B^T = (B_{ji})_{ij} \qquad B \in \mathbb{R}^{u \times d}, \quad B^T \in \mathbb{R}^{d \times u}$

Prop: $\langle \underbrace{Bu}, v \rangle_{\mathbb{R}^u} = \langle u, B^T v \rangle_{\mathbb{R}^d}$

$f(x+\varepsilon v) \overset{\circledast}{=} f(x) + \varepsilon \langle \underbrace{B^T \cdot \nabla L(Bx)}_{= \nabla f(x)}, v \rangle + o(\varepsilon)$

Prop: $\nabla f(x) = \nabla (L \circ B)(x)$
$\qquad = B^T \times \nabla L(Bx)$

$$\boxed{\nabla(L \circ B) = B^T \circ \nabla L \circ B}$$

Examples: $\quad L(u) = \frac{1}{2} \| u - y \|^2 = \frac{1}{2} \sum (u_i - y_i)^2$

$$\nabla L(u) = u - y$$

$$\nabla \left( \frac{1}{2} \|.\|^2 \right)(u) = u$$

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

$$\nabla f(x) = A^T (Ax - y)$$

$$\min_u \|Ax - y\|^2 \implies A^T(Ax - y) = 0 = \nabla f(x) = 0$$

$$\implies \boxed{(A^TA)}u = A^T y \quad (\text{normal eq}^\circ)$$

⚠ $Ax = y$     Covariance
                $\in \mathbb{R}^{d \times d}$

$$\to \det \neq 0 \iff \text{Ker}(A) = \{0\}.$$

if $A^TA$ is invertible $("u > d")$, unique sol$^\circ$

$$\boxed{u = (A^TA)^{-1} A^T y} \quad \text{"Overdetermined"}$$

$\underbrace{\phantom{(A^TA)^{-1}A^T}}$
$A^+$  Moore-Penross
       Pseudo-Inverse

→ if $A^TA$ not invertible $(d \gg u)$  $\infty$ possibility
   → Ridge
   → Lasso

**Logistic:** $L(u) = \sum_i \ell(u_i)$

$$\ell(s) = \log(1 + e^s).$$
$$\ell'(s) = \frac{e^s}{1 + e^s}$$

$$\nabla L(u) = \begin{pmatrix} \ell'(u_1) = \frac{e^{u_1}}{1+e^{u_1}} \\ \ell'(u_2) \\ \vdots \\ \ell'(u_m) = \frac{e^{u_m}}{1+e^{u_m}} \end{pmatrix} = \text{sigmoid}(u)$$

$$\nabla f(x) = B^T \nabla L(Bx) = \underbrace{B^T}_{\text{Back Prop}} \times \text{sigmoid}(\underline{Bx})$$

**Grad descent:** $u_{k+1} = u_k - \boxed{\tau_k} \, \nabla f(u_k)$

$\hookrightarrow$ step size
$\hookrightarrow$ learning rate

$\tau_k$ large $\rightarrow$ fast.

$\tau_k$ small $\rightarrow$ avoid explosion

Baby case: $f(x) = \frac{1}{2} \| Ax - y \|^2$

$$C \triangleq A^T A \qquad C_{k\ell} = \begin{array}{l} \text{"how much} \\ \text{feature } k, \ell \\ \text{correlated} \end{array}$$
$$\underset{\mathbb{R}^{d \times d}}{}$$

$C$ is a symetric matrix: $C^T = C$

$$C^T = (A^T A)^T = A^T (A^T)^T = A^T A = C$$

$(AB)^T = B^T A^T$

Thm: Since $C$ is symetric, $(\underset{\in \mathbb{R}^d}{u_1}, \dots, u_d)$ eigenvectors

$(\underset{\geqslant 0}{\lambda_1}, \dots, \underset{\geqslant 0}{\lambda_d})$ eigenvalues

$$C u_i = \underset{\geqslant 0}{\lambda_i} u_i \qquad \qquad \lambda_1 \geqslant \lambda_2 \geqslant \dots \geqslant \lambda_d.$$



$$\underset{A^T A}{\underline{\underline{C}}} u_i = \lambda_i u_i \implies \langle \overset{\frown}{A^T A u_i}, u_i \rangle = \lambda_i \langle u_i, u_i \rangle$$

$$\langle A u_i, A u_i \rangle = \lambda_i \langle u_i, u_i \rangle$$

$$\| A u_i \|^2 = \lambda_i \| u_i \|^2$$

$$\lambda_i = \frac{\| A u_i \|^2}{\| u_i \|^2} \geqslant 0$$

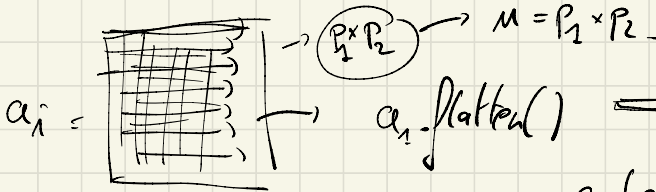Positive Semi-definite Matrix       SDP matrices.

Thm: For gradient descent if
$$0 < \tau_k < \frac{2}{\lambda_d}$$

then $u_k \longrightarrow$ sol$^c$ $u^* = (A^TA)^{-1}A^T y$.

$$A = \begin{pmatrix} \underline{a_1} \\ \underline{a_2} \\ \vdots \\ \underline{a_n} \end{pmatrix}$$

$$C_{k\ell} = \frac{1}{m}\sum_{i=1}^{m} a_i[k] \cdot a_i[\ell]$$

$$\approx \mathbb{E}_a\left(a[k]\, a[\ell]\right)$$

$\longrightarrow \boxed{P_1 \times P_2}\longrightarrow M = P_1 \times P_2$

$a_i = $ $\longmapsto a_1.\text{flatten}()$

$a_i = \left(a_i[1], a_i[2] \cdots a_i[d]\right)$

$a_i = (\text{weight, height, age, .....})$

$a = $ $\longrightarrow$ wave length

$C = $