

bit.ly / 3Pg h xme

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

ML (supervised)

notation: Data (a_i, y_i) $a_i \in \mathbb{R}^d$ $y_i \in \mathbb{R}$
 $\uparrow \quad \uparrow$
 feature label
 $y_i \in \{-1, +1\}$

Prediction model: $y_i \approx g_\pi(a_i)$
 $\pi = \text{params of model}$

Examples: linear $g_\pi(a) = \langle a, \pi \rangle$
 $\uparrow \quad \uparrow$
 $\mathbb{R}^d \quad \mathbb{R}^d$

def.

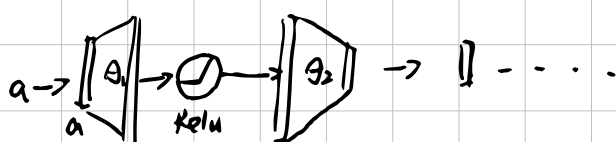
$$\langle a, \pi \rangle_{\mathbb{R}^d} \triangleq \sum_i a[i] \pi[i]$$

$d=1$.

with bias: $g_\pi(a) = \langle a, \pi^{(1)} \rangle + \pi^{(0)}$ $\pi = (\pi^{(1)}, \pi^{(0)})$
 $\uparrow \quad \uparrow$
 $\mathbb{R}^d \quad \mathbb{R}$

Neural networks.

MLP



$$\pi = (\theta_1, \theta_2, \theta_3)$$

Data: $(a_i, y_i)_{i=1}^n$ Model: g_π

ERM (Empirical Risk minimization)

$$\min_{\pi} \left[\frac{1}{n} \sum_{i=1}^n \ell(g_\pi(a_i), y_i) \triangleq f(\pi) \right]$$

$\uparrow \quad \uparrow$
 param input

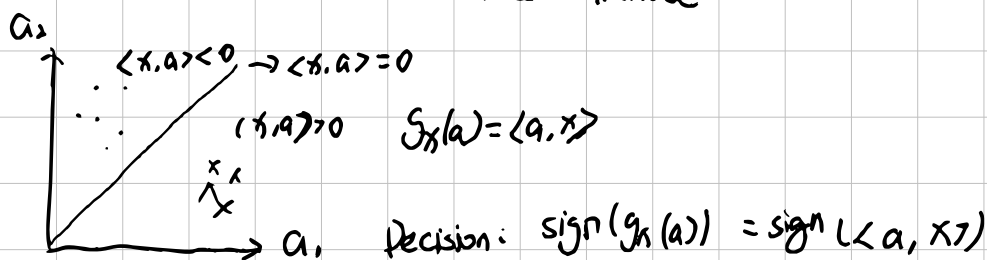
Loss: Regression: $y_i \in \mathbb{R}$

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

MSE (mean squared error)

(classification loss: $y \in \{-1, +1\}$ (binary class))

"Ideal": $\ell_0(y, y') = \begin{cases} 0 & y = y' \\ 1 & \text{otherwise} \end{cases}$



$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{ \langle a_i, x \rangle y_i \leq 0 \}} = f(\pi) \text{ not continuous}$$

Nice function: differentiable, convex --

$$\frac{1}{n} \sum_{i=1}^n \sigma(-y_i \langle a_i, x \rangle) = f(\pi)$$

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

→ SVM
 → logistic regression

hinge $\phi_{\text{sum}}(t) = \max(t + 1, 0)$

logistic $\phi_{\text{log}}(t) = \frac{\log(1 + e^t)}{\log(2)}$

$\log(1 + e^t) \approx t$

$\phi'(t) = \frac{e^t}{1 + e^t}$ sigmoid [2.1] Proof

Matrix notation: Design matrix

$$A = \underbrace{a_i}_{d} \left. \begin{array}{c} \square \\ \square \\ \square \end{array} \right\} n$$

Linear model: $(\langle a_i, x \rangle)_{i=1}^n = Ax$
 $A @ x \xrightarrow{\quad}$

MSE: $\min_x \frac{1}{n} \sum_i (\underbrace{\langle a_i, x \rangle}_{(Ax)_i} - y_i)^2 = f(x)$
 $f(x) = \frac{1}{n} \|Ax - y\|^2$

$\|z\|^2 = \langle z, z \rangle_{\mathbb{R}^n} = \sum_{i=1}^n z_i^2$ np.linalg.norm

Logistic: $\min_x \frac{1}{n} \sum_{i=1}^n \sigma_{\log}(-y_i \langle x, a_i \rangle)$

$(y_i \cdot \langle x, a_i \rangle)_{i=1}^n = \text{diag}(y) \cdot A \cdot x$
 $\boxed{A @ x * y}$ $J^*(A @ x)$

Global loss: $L(z) = \frac{1}{n} \sum_{i=1}^n \sigma_{\log}(z_i)$
 $L(z) = \frac{1}{n} \dots \dots \dots$
 $= L(\underbrace{-\text{diag}(y) \cdot A}_{B} \cdot x)$
 $= L(Bx)$

$B = \begin{bmatrix} \square \\ \square \\ -y_i a_i \\ \square \end{bmatrix}$

Summary: $\begin{cases} \text{MSE} & \min_x \frac{1}{n} \|Ax - y\|^2 \\ \text{log} & \min_x L(Bx) \end{cases}$
 $\hookrightarrow B=A \quad L(z) = \frac{1}{n} \|z - y\|^2$

Gradient of $f: \mathbb{R}^d \rightarrow \mathbb{R}$

Def: $\nabla f(x) = \begin{pmatrix} \frac{df}{dx_1}(x) \\ \vdots \\ \frac{df}{dx_d}(x) \end{pmatrix} \in \mathbb{R}^d$

partial derivatives:

$\frac{df}{dx_k}(x) = \lim_{\varepsilon \rightarrow 0} \frac{f(x + \varepsilon \delta_k) - f(x)}{\varepsilon}$
 $\delta_k = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \leftarrow k$

Def: f is differentiable at x iff

$\exists: f(x + \varepsilon u) = f(x) + \varepsilon \langle u, \nabla f(x) \rangle + o(\varepsilon)$

$\Leftrightarrow \frac{f(x + \varepsilon u) - f(x)}{\varepsilon} \xrightarrow{\varepsilon \rightarrow 0} \langle u, \nabla f(x) \rangle$

$\triangle!$ f is diff $\Rightarrow f$ has a gradient

$f(x_1, x_2) = \frac{2x_1 x_2 (x_1 + x_2)}{x_1^2 + x_2^2} \quad f(0) = 0?$

f is continuous

$\nabla f(0, 0) = 0$

Prop: f has a $\nabla f \Rightarrow f$ is diff. and ∇f is continuous

Why ∇ ? Theory: x is a ^{local} minimizer $\Rightarrow \nabla f(x) = 0$

Algorithm: $-\nabla f(x)$ is steepest descent direction

GD $x_0 \leftarrow \text{Init}$

$x_{k+1} = x_k - [\eta_k \nabla f(x_k)]$

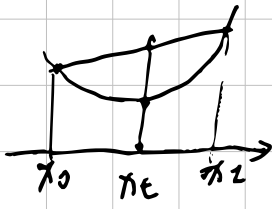
\hookrightarrow step size / learning rate

Thm: x^* is local minimizer of f

f is diff
at x^*

$$\nabla f(x^*) = 0$$

convexity:



$$(1-t)x_0 + tx_1 \quad t \in [0, 1]$$

Def: f is convex iff $\forall (x_0, x_1), \forall t \in [0, 1]$

$$f((1-t)x_0 + tx_1) \leq (1-t)f(x_0) + tf(x_1)$$

Thm: if f is conv and diff

$$\uparrow x \text{ is a minimizer} \Leftrightarrow \nabla f(x^*) = 0$$

$\downarrow f$ is a loc min $\Rightarrow f$ is a global min

$$f \text{ conv} \Leftrightarrow f'' \geq 0$$

f is above its tangent



Prop: f, g are convex $\left\{ \begin{array}{l} \lambda, \mu \geq 0 \end{array} \right\} \Rightarrow \lambda f + \mu g \text{ conv}$

A a matrix
 b a vector

$$f \text{ conv} \Rightarrow f(Ax + b) \text{ conv}$$

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

$$z \mapsto \|z\|^2 = \sum z_i^2$$

Logistic $f(x) = L(Bx)$

$$L(z) = \sum \sigma(z_i)$$



$$\sigma(s) = \log(1 + e^s)$$

$$\sigma'(s) = \frac{e^s}{1+e^s} \quad \sigma''(s) \geq 0$$

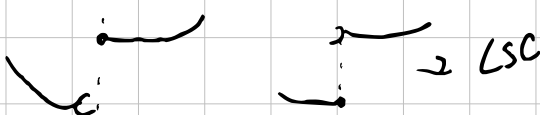
$$\textcircled{1} \quad \begin{array}{l} \varphi(x) = e^{-x} \geq 0 \\ f(x) = x \end{array} \quad \text{no minimizer}$$

Prop: f is continuous

f is $\sim f(x) \xrightarrow{\|x\| \rightarrow \infty} \text{two}$

Then there exists a minimizer

"Proof"



(lower semi-continuous (LSC))

$$x_k \rightarrow x \quad \liminf f(x_k) \geq f(x)$$

Gradient:

$$f(x + \varepsilon u) = f(x) + \varepsilon \langle \nabla f(x), u \rangle + o(\varepsilon)$$

MSG: $f(x) = \|Ax - y\|^2$

$$f(x + \varepsilon u) = \|A(x + \varepsilon u) - y\|^2$$

$$= \|(Ax - y) + \varepsilon Au\|^2 \quad \|a+b\|^2 = \|a\|^2 + 2\langle a, b \rangle + \|b\|^2$$

$$= \|Ax - y\|^2 + 2\langle Ax - y, \varepsilon Au \rangle + \varepsilon^2 \|Au\|^2$$

$$= f(x) + \varepsilon \cdot 2\langle Ax - y, Au \rangle_{\mathbb{R}^d} + o(\varepsilon)$$

$$\stackrel{?}{=} \dots \quad \langle \nabla f(x), u \rangle_{\mathbb{R}^d}$$

$$\langle Aw, z \rangle_{\mathbb{R}^n} = \langle w, A^T z \rangle_{\mathbb{R}^d}$$

$$\rightarrow = f(x) + \varepsilon \underbrace{\langle 2A^T(Ax - y), u \rangle}_{\nabla f(x)} + o(\varepsilon)$$

Conclu: $f(x) = \|Ax - y\|^2 \quad \nabla f(x) = 2A^T(Ax - y)$

corollary: $\min_{x \in \mathbb{R}^d} \|Ax - y\|^2$

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



may not be unique

$$A = \begin{bmatrix} \end{bmatrix} \quad \ker(A) = \{z: Az = 0\}$$

"under-determined" d big n small

$$x^* \text{ sol. } z \in \ker(A) \Rightarrow x^* + z \text{ sol.}$$

$$\|Ax - y\|^2 = \|A(x^* + z) - y\|^2$$

$$\underbrace{(A^T A)}_{\substack{\in \mathbb{C} \in \mathbb{R}^{d \times d} \text{ "correlation" }}} x^* = A^T y$$

Prop: $A^T A x = y$ has a unique solution

$$\Leftrightarrow \ker(A) = \{0\}$$

Proof: $\Leftrightarrow \ker(A^T A) = \{0\}$

"
 $\ker(A)$

$$\textcircled{1} \quad Az = 0 \Rightarrow A^T Az = 0$$

$$\textcircled{2} \quad A^T Az = 0 \quad \langle A^T Az, z \rangle = \langle 0, z \rangle = 0$$

\Downarrow

$$\langle Az, Az \rangle = 0 \Rightarrow \|Az\|^2 = 0 \Rightarrow Az = 0$$

$C = A^T A$ symmetric semi-definite matrix

$[C^T = C]$ diagonalizable in orthonormal basis

real eigenvalue

if $C = A^T A \Leftrightarrow$ eigenvalues are ≥ 0

Prop. if overdetermined, $\ker(A) = \{0\}$

$$n \gg d$$

$$x^* = \underbrace{(A^T A)^{-1}}_{\text{Moore-Penrose Pseudo Inverse}} \cdot A^T y \quad x^* = A^+ y$$

Ridge Penalty // weight decay:

$$\min_{x \in \mathbb{R}^d} \|Ax - y\|_{\mathbb{R}^n}^2 + \lambda \|x\|_{\mathbb{R}^d}^2$$

\downarrow
ridge param

Select λ : cross validation

$$f(x) = \|Ax - y\|^2 + \lambda \|x\|^2$$

$$\nabla f(x) = 2A^T(Ax - y) + \lambda 2x$$

$$\nabla f(x^*) = 0 \Leftrightarrow A^T Ax + \lambda x = A^T y$$

$$(A^T A + \lambda Id) x = A^T y$$

if u_i is eigenvalue of B

$$u_i + \lambda \quad \dots \quad B + \lambda Id$$

$$Bz = u_i z \Rightarrow (B + \lambda Id)z = (u_i + \lambda)z$$

$A^T A$ has eigenvalue ≥ 0

$$A^T A + \lambda Id \geq \lambda > 0$$

$$\Rightarrow (A^T A + \lambda Id) z = A^T y$$

have a unique solution

$$\underbrace{\|Ax - y\|^2}_{\text{parabola}} + \underbrace{\lambda \|x\|^2}_{\text{circle}} = \underbrace{\text{shifted parabola}}_{\text{unique min}}$$

Summary: $x_{\lambda}^* \stackrel{[F]}{=} (A^T A + \lambda Id_{\mathbb{R}^d})^{-1} A^T y$

Prop: kernel trick ||

$$x_{\lambda}^* \stackrel{[F]}{=} A^T (A A^T + \lambda Id_{\mathbb{R}^n})^{-1} y$$

[F] is better $d < n$

[E] is better $d > n$ $d = \text{trig}$

logistic: $f(x) = L(Bx)$ $\rightarrow \mathbb{R}^{n \times d}$
 $L(z) = \sum_{i=1}^n \sigma(z_i)$ \uparrow diff

$$f(x + \varepsilon u) = L(Bx + \varepsilon Bu)$$

Taylor exp. of L at point $z = Bx$

$$v = Bu$$

$$L(z + \varepsilon v) = L(z) + \varepsilon \langle \nabla L(z), v \rangle + o(\varepsilon)$$

$$f(x + \varepsilon u) = L(Bx) + \varepsilon \langle \nabla L(Bx), Bu \rangle + o(\varepsilon)$$

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

Prop: $f(x) = L(Bx) \quad \nabla f(x) = B^T \nabla L(Bx)$

$$\nabla (L \circ B) = B^T \circ \nabla L \circ B$$

example logistic: $L(z) = \sum_{i=1}^n \sigma(z_i)$

Jacobian : $F: \mathbb{R}^d \rightarrow \mathbb{R}^p$
 $x \mapsto f(x)$

Def: $\partial F(x) \in \mathbb{R}^{p \times d}$

$\partial F(x): \mathbb{R}^d \rightarrow \mathbb{R}^p$ Linear

$u \mapsto \partial F(x) \times u$

$F(x) = \begin{pmatrix} F_1(x) \\ F_2(x) \\ \vdots \\ F_p(x) \end{pmatrix} = (F_i(x))_{i=1}^p$

Ex. $f(x) = \begin{pmatrix} x_1^2 + x_2^2 \\ x_1 x_2 \end{pmatrix} \quad x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$

$\partial f(x) = \begin{pmatrix} 2x_1 & 2x_2 \\ x_2 & x_1 \end{pmatrix}$

$\partial F(x) \cong \left(\frac{\partial F_i}{\partial x_j}(x) \right)_{\substack{i=1 \dots p \\ j=1 \dots d}} \in \mathbb{R}^{p \times d}$

① $p=1 \quad F: \mathbb{R}^d \rightarrow \mathbb{R}$

$\partial F(x) \in \mathbb{R}^{1 \times d}$

$\nabla F(x) \in \mathbb{R}^{d \times 1}$

$\nabla F(x) = (\partial F(x))^\top$

Prop/Def: F differentiable iff

$F(x+u) = F(x) + \underbrace{\varepsilon}_{\mathbb{R}^{p \times d}} \underbrace{\partial F(x)}_{\mathbb{R}^d} \times u + o(\varepsilon)$

Chain rule $\left. \begin{array}{l} F: \mathbb{R}^d \rightarrow \mathbb{R}^p \\ G: \mathbb{R}^p \rightarrow \mathbb{R}^d \end{array} \right\} \text{diff}$

" $\partial (G \circ F) = \partial G \circ \partial F$ "

$\partial (G \circ F)(u) = \underbrace{\partial G(F(x))}_{q \times p} \times \underbrace{\partial F(x)}_{p \times d} u$

$f(x) = L(B(x))$

$\partial f(x) = \partial L(B(x)) \times \partial B(x)$

$\nabla f(x) = (\partial f(x))^\top$

$\nabla f(x) = \partial B(x)^\top \times \nabla L(B(x))$

$F(x) = Bx$

$F(x+\varepsilon u) = Bx + \underbrace{\varepsilon B u}_{\partial F(x) \times u} + o$

Conclu: $\nabla (L \circ B) = B^\top \circ \nabla L \circ B$