

Linear regression  
in data sciences.

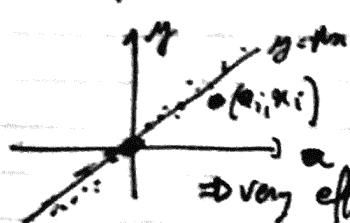
①

• Inverse problem in imaging: forward model  $y = Ax + w$  usually  $p \ll n$  (ill-posed).  $\begin{matrix} \text{inputting} \\ \text{deconvolve} \\ \text{Random tomography} \\ \text{IRM} \end{matrix}$

• Regression in ML: supervised learning  $(a_i, y_i)_{i=1}^n$   $a_i \equiv \text{features}$

lin. predic<sup>c</sup> model

$y = \langle a, x \rangle$  by fitting  $\forall i \quad y_i \approx \langle a_i, x \rangle$



$$\begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix} x = y.$$

$\Rightarrow$  very efficient in high dim!

• Why we don't want to solve exactly  $y = Ax$

$\hookrightarrow$  noise

$\hookrightarrow$  over-det  $y \notin \text{Im } A \rightarrow$  no sol<sup>c</sup>

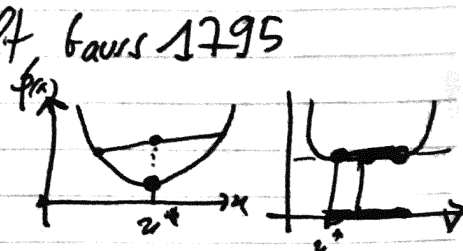
$\rightarrow$  least square fit Gauss 1795

• LS fit

$$\min_{x \in \mathbb{R}^p} f(x) = \frac{1}{2} \|Ax - y\|_2^2$$

$f$  is smooth (differentiable)

$f$  is convex



Thm:  $x \in \text{Argmin } f \Leftrightarrow \nabla f(x) = 0$

Refresher:  $f: \mathbb{R}^p \rightarrow \mathbb{R}$ ,  $\nabla f: \mathbb{R}^p \rightarrow \mathbb{R}^p$

$$\nabla f(x) = \left( \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_p} \right)^T \in \mathbb{R}^p$$



$$f(x + \epsilon) = f(x) + \langle \nabla f(x), \epsilon \rangle + o(\|\epsilon\|).$$

$$\text{Here: } f(x + \epsilon) = \frac{1}{2} \|Ax - y + A\epsilon\|_2^2 = \frac{1}{2} \|Ax - y\|_2^2 + \langle A\epsilon, Ax - y \rangle + \frac{1}{2} \|A\epsilon\|_2^2$$

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

$$x^* \in \text{Argmin } f \Leftrightarrow A^T (Ax^* - y) = 0 \Leftrightarrow (A^T A) x^* = A^T y \quad \text{normal equations}$$

Prop: if  $\ker A = \{0\}$  (over-determined)

$$\begin{bmatrix} A \\ I \end{bmatrix} = \begin{bmatrix} A \\ I \end{bmatrix}$$

$A^T A$  is positive and invertible

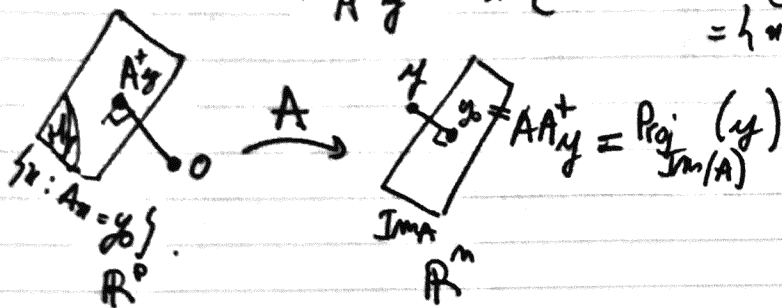
$$x^* = (A^T A)^{-1} A^T y \text{ is unique}$$

otherwise  $\text{Argmin } f$  is infinite (hyperplane  $\parallel$  to  $\ker(A)$ )

- Needs a select<sup>e</sup> method, also called Regularization. Use also least square  $\rightarrow$  Ridge regression.  $\rightarrow$  simplest method

$$A \min_{x \in \mathbb{R}^n} \frac{1}{2} \|x\|^2 : x \in \text{Argmin} \|Ax - y\|^2 \quad (\text{sol. unique})$$

$$= \{x : A^T A x = A^T y\} = \{x : Ax = y\}$$



Prop: ~~if~~ If  $\text{Im } A = \mathbb{R}^n$  (undetermined,  $\square \cdot \mathbb{I}$ )

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

Proof:  $\text{Im } A = \mathbb{R}^p \rightarrow y = y$  (surjective). Lagrange mult:  $\nabla \|u\|_2^2 \in (\text{constraint})^\perp$

$$\Rightarrow n \in \text{Ker}(A)^\perp = \text{Im}(A^T)$$

$$\Rightarrow n = A^T u \Rightarrow Ax = A A^T u = y$$

$$\Rightarrow u = (A A^T)^{-1} y$$

invertible

Summary:

- Over-det  $\leftrightarrow A$  surject  $\leftrightarrow A^T A$  inv.  $\Leftrightarrow A^+ = (A^T A)^{-1} A^T$
- Underdet  $\leftrightarrow A$  inj  $\leftrightarrow A A^T$  inv  $\Leftrightarrow A^+ = A^T (A A^T)^{-1}$
- $A$  invertible  $\leftrightarrow A^+ = A^{-1}$  (have the case)
- General case  $\rightarrow$  SVD (see below).

- Solving  $A^T A x = A^T y$  (over-det case): conjugate gradient. if  $p$  is large
- Choleski: if  $p$  small
- SVD: todo.

• Regular:  $\xrightarrow{\text{noise}} \text{ML (model error + small } m)$   $\min_{x \in \mathbb{R}^p} \frac{1}{2} \|y - Ax\|^2 + \frac{\lambda}{2} \|x\|^2$   $\Delta R(y)$

Prop<sup>6</sup>:  $x_1 \stackrel{\text{Q}}{=} (A^*A + \lambda \text{Id}_p)^{-1} A^* y \rightarrow$  for small  $p$  (underdetermined)  
 $\stackrel{\text{Q}}{=} A^* (AA^* + \lambda \text{Id}_n)^{-1} y \rightarrow$  ok for small  $n \ll p$  (underdetermined)

② allows to deal with  $\infty$  dimension feature space (kernel method)  
 $u = \underbrace{(AA^* + \lambda \text{Id}_n)^{-1}}_{\hat{K}} y \quad \hat{K} = (k(a_i, a_j))_{i,j}$

$\langle a, a \rangle$  becomes  $\langle a, A^* p \rangle \quad \sum_{i=1}^n u_i k(a_i, a)$

Thm:  $x_1 \xrightarrow{\text{noise}} A^* y$ , more generally:  $x_1 \xrightarrow{\text{noise}} \arg\min \{R(x) : Ax = y_0\}$

Proof:  $\|y - Ax\|^2 = \underbrace{\|y_0 - Ax\|^2}_{\in \text{Im } A} + \underbrace{\|y - y_0\|^2}_{\in (\text{Im } A)^\perp}$

• Stat vs Imaging: Rather normalize  $\min \frac{1}{n} \|y - Ax\|^2 + \lambda \|x\|^2$   
 $x_1 = \left( \underbrace{\frac{1}{n} A^*A}_{\hat{C}} + \lambda \text{Id}_p \right)^{-1} \underbrace{A^* y}_{\hat{u}}$   $\hat{C} = \frac{1}{n} \sum_{i=1}^n a_i a_i^T \in \mathbb{R}^{p \times p}$

under assumption  $(x_i, y_i) \stackrel{i.i.d.}{\sim} (x, y)$ :  $\begin{cases} \hat{C} \xrightarrow{P \rightarrow \infty} C = E(a a^T) \in \mathbb{R}^{p \times p} \\ \hat{u} \xrightarrow{P \rightarrow \infty} u = E(y a) \in \mathbb{R}^p \end{cases}$

heuristic:  $\begin{bmatrix} \|\hat{C} - C\| \\ \|\hat{u} - u\| \end{bmatrix} \sim \frac{1}{\sqrt{n}} \begin{bmatrix} \varepsilon \\ \varepsilon \end{bmatrix}$  VS Imaging IP  $\begin{cases} \|C - \hat{C}\| = 0 \\ \|u - \hat{u}\| = \|A^* w\| \stackrel{\Delta}{=} \varepsilon \end{cases} \quad y = Ax_0 + w$   
 $\underbrace{A^* A x_0}_{\text{SS}} \underbrace{A^* y}_{\text{SS}} \quad \|w\|$

ML is harder to study than IP bc  $\hat{C} \neq C$ .

• Focus on the IP setup  $y = Ax_0 + w \rightarrow$  see course note for ML

- Quest<sup>n</sup>: does  $x_1 \rightarrow x_0$  as  $\frac{1}{\|w\|} \rightarrow 0$ ? Rate? hyp on  $x_0$ ?
- No if  $x_0 \notin \text{Im}(A^*) = \text{Ker } A^\perp$ !! Because  $x_1 \in \text{Im } A^*$ !!  $\Rightarrow$  Hyp:  $x_0 \perp \text{Ker } A$  is  $x_0 \in \text{Im } A^*$
  - $x_1 - x_0 = (A^*A + \lambda \text{Id})^{-1} A^* (Ax_0 + w) - x_0$

$$= \left[ (A^*A + \lambda \text{Id})^{-1} A^*A - \text{Id}_p \right] x_0 + (A^*A + \lambda \text{Id})^{-1} (A^*w)$$

One can show (ive)  $\| \cdot \| \leq \frac{1}{\lambda} \| (A^*A)^{-1} \| \cdot \|x_0\| + \| (A^*A + \lambda \text{Id})^{-1} \| \cdot \|A^*w\|$

Plm:  $\sigma_{\min}(A)$  can be arbitrary small  
(and can be 0 in the  $\infty$  dimens<sup>n</sup>)

$\rightarrow$  Need some condit<sup>n</sup>, of course.