

Linear Regression reduces to the following optimisation problem

$$\min_{\omega} E(\omega, D) = \|X\omega - y\|^2$$

$\|A\omega - y\|^2$ is a convex fn ($\frac{d^2 f}{dx^2} \geq 0$)

$$\omega^* = (X^T X)^{-1} X^T y \text{ if } X^T X \text{ is invertible}$$

Otherwise we can solve for

$$\nabla_{\omega} E = 0 \Rightarrow X^T X \omega - X^T y = 0$$

Recall: A linear system of equations can be written as $Ax = b$

Now If $A = [a_1, \dots, a_d]$ where a_j is column vector

then $Ax = x_1 a_1 + \dots + x_d a_d$ when $x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$

Note that $Ax = b$ has a sol'n iff $b \in C(A)$: Column space of A

Using this we can deduce that $X^T X \omega = X^T y$ always has a sol'n

as $X^T y \in C(X^T X)$, hence can be solved to find ω^* .

each column of $X^T X \in C(X^T)$?

$y \in C(X)$

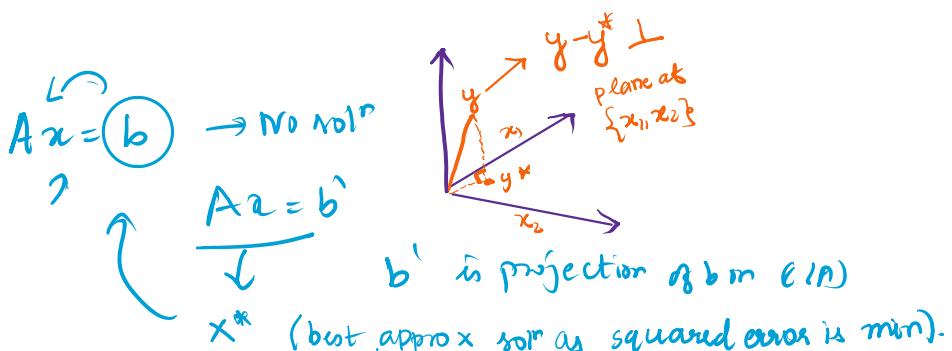
$$X^T y \in C(X^T X) ?$$

$$\begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} y = \begin{bmatrix} x_1^T y \\ \vdots \\ x_n^T y \end{bmatrix}$$

$$X^T (\underbrace{x_1, x_2, \dots, x_n}_y) = \begin{bmatrix} x_1^T x_1 & x_1^T x_2 & \dots & x_1^T x_n \\ x_2^T x_1 & x_2^T x_2 & \dots & x_2^T x_n \\ \vdots & \vdots & \ddots & \vdots \\ x_n^T x_1 & x_n^T x_2 & \dots & x_n^T x_n \end{bmatrix}$$

Case I: $X\omega^* = y$

Case II: Can't find a ω^* : $X^T X \omega - X^T y = 0 \Rightarrow X^T (X\omega - y) = 0$



Regression model with basis functions

Take $d=1$: $\hat{y}_i = w_0 + w_1 x_i + (w_2 x_i^2 + \dots + w_m x_i^m)$ (x_i 's are scalars)

We define basis fns:

$$\begin{aligned}\phi_0(x_i) &= 1 \\ \phi_1(x_i) &= x_i \quad \Rightarrow \quad \hat{y}_i = \sum_{j=0}^m w_j \phi_j(x_i) \\ &\vdots \\ \phi_m(x_i) &= x_i^m\end{aligned}$$

For d -dimensions, $\hat{y}_i = \sum_{j=0}^m \phi_j(x_i) w_j$, $m > d$

where, $\Phi = \begin{bmatrix} \phi_0(x_1) & \dots & \phi_m(x_1) \\ \vdots & \ddots & \vdots \\ \phi_0(x_n) & \dots & \phi_m(x_n) \end{bmatrix}_{n \times (m+1)}$, $w = \begin{pmatrix} w_0 \\ \vdots \\ w_m \end{pmatrix}_{(m+1) \times 1}$, $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$

$$\Rightarrow \arg \min \| \Phi w - y \|^2 \text{ with } w^* = (\Phi^\top \Phi)^{-1} \Phi^\top y$$

Probabilistic model of Linear Regression

In real life cases, some error/randomness is also observed.

$$\therefore y_i = \omega^\top x_i + \epsilon_i \text{ (noisy linear model)}$$

This leads to a probabilistic regression model.

In a noisy linear model, we have

- i) parameters = ω
- ii) noise: $\epsilon_i \sim N(0, \sigma^2)$, with $\text{cov}(\epsilon_i, \epsilon_j) = 0 \quad \forall i \neq j$ (iid's)
- iii) Data set: $D = \{(x_i, y_i)\}_{i=1}^n$

Our task is to estimate ω from this probabilistic model using following methods: (i) Maximum Likelihood Estimation

Maximum Likelihood Estimation (MLE)

Maximum Likelihood Estimation (MLE)

A set of independent and identically distributed observations $\{y_1, \dots, y_n\}$ are generated by a probabilistic model parametrised by θ

$$y_i \sim P(y_i | \theta) \quad \xrightarrow{\text{likelihood function}}$$

Log likelihood: $\log(P(y| \theta))$ (useful as log is monotonic fn
mathematically easier,
numerical advantage)

We've:

$$\hat{\theta}_{MLE} = \arg \max L(\theta) = \arg \max \sum_{j=1}^n \log P(y_j | \theta)$$

Ex: Toss a coin n times. Each one of the outcomes is a binary Random Variable with Bernoulli distribution with

$$P(y_i | \theta) = \theta^{y_i} (1-\theta)^{1-y_i}$$

$$L(\theta) = \log P(y | \theta) = \log \left(\prod_{j=1}^n P(y_j | \theta) \right)$$

$$= \sum_{j=1}^n \log P_j(y_j | \theta) = \sum_{j=1}^n (y_j \log \theta + (1-y_j) \log (1-\theta))$$

$$\Rightarrow \hat{\theta}_{\max} = \frac{\sum y_i}{n}$$

MLE for regression

$$y_j = \omega^\top x_j + \epsilon_j \sim N(0, \sigma^2)$$

$$y_j \sim N(\omega^\top x_j, \sigma^2)$$

$$P(y_j | x_j, \omega) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(y_j - \omega^\top x_j)^2}{2\sigma^2} \right\}$$

$$L(\omega) = (\text{const.}) - \sum_{j=1}^n \frac{(y_j - \omega^\top x_j)^2}{2\sigma^2}$$

$$\Rightarrow \arg \max_{\omega} L(\omega) = \arg \min_{\omega} \sum_{j=1}^n (y_j - \omega^\top x_j)^2 = \omega^*$$

$$\equiv \arg \min \|\mathbf{X}\omega - \mathbf{y}\|^2$$

(This is same as deterministic regression)

$$\equiv \arg \min \|Xw - y\|^2 \quad (\text{This is same as deterministic regression})$$

We know $w^* = (X^T X)^{-1} X^T y$, but this needs inverse which is computationally costly.

So we adopt an algorithmic way to optimise $\|Xw - y\|^2$, called

Gradient Descent.

- $w \leftarrow w_0$
 - Repeat until convergence: $\|\nabla E(w)\| < \epsilon$
- $$w \leftarrow w - \eta \nabla_w E \quad (\eta: \text{learning rate})$$