

Linear Regression, Regularization & the Bias-Variance Tradeoff.

Dataset $\mathcal{D} = \{(x_t, y_t)\}_{t=1}^n \subset \mathbb{R}^d \times \mathbb{R}$

real number (\mathbb{R})

Assume y_t is approximately an affine line of x_t , i.e.

$$y_t \sim \underline{\theta}^T \underline{x}_t + \theta_0 \quad \text{for some } (\underline{\theta}, \theta_0)$$

random not random

can use kernel trick on inputs x_t

More precisely, $E[y_t | x_t, \underline{\theta}, \theta_0] = \underline{\theta}^T \underline{x}_t + \theta_0$

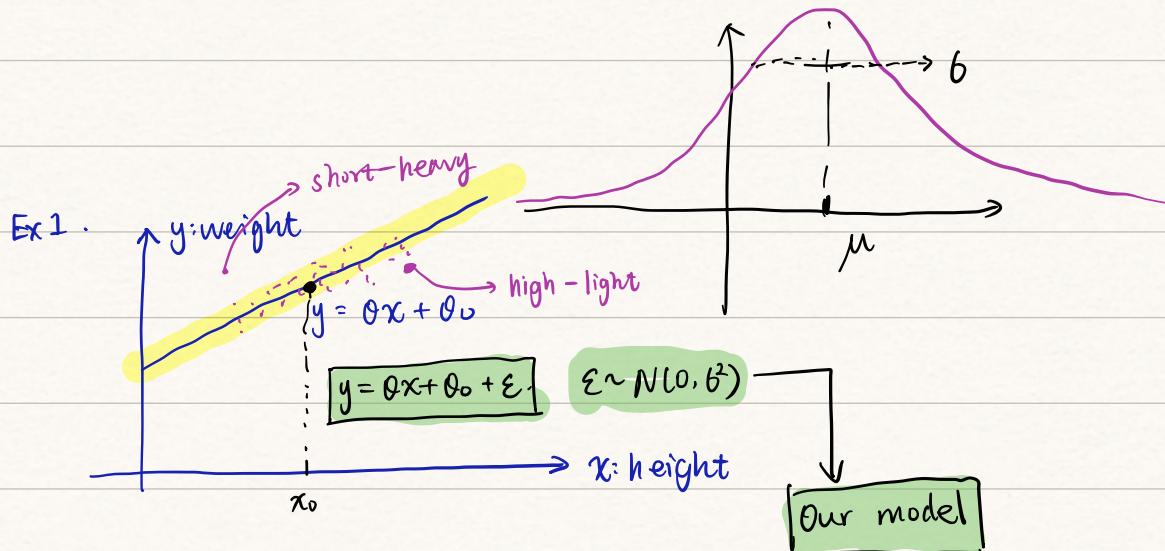
Associate a Prob. distribution of y given parameters $(\underline{\theta}, \theta_0)$ & input x

$$\rightarrow P(y | x, \underline{\theta}, \theta_0)$$

\Downarrow tractable

Gaussian: $P(y | x, \underline{\theta}, \theta_0) = N(y; \underline{\theta}^T \underline{x} + \theta_0, \sigma^2)$

$$N(y; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(y-\mu)^2}{2\sigma^2}\right\} \quad y \in \mathbb{R}$$



Clearly, $E[y | x, \underline{\theta}, \theta_0] = E[\underline{\theta}^T \underline{x} + \theta_0 + \epsilon | x, \underline{\theta}, \theta_0]$

$$= \underline{\theta}^T \underline{x} + \theta_0$$

Generative Model

Dataset $\mathcal{D} = \{(x_t, y_t)\}_{t=1}^n$

$$y_t = \underline{\theta}^T x_t + \theta_0 + \varepsilon_t \quad \varepsilon_t \sim N(0, \sigma^2)$$

$$\varepsilon_1, \dots, \varepsilon_n \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

Goal: Given $\mathcal{D} = \{(x_t, y_t)\}_{t=1}^n$, estimate the parameters $(\underline{\theta}, \theta_0) \in \mathbb{R}^d \times \mathbb{R}$.

Strategy: Maximum Likelihood Estimation

① Likelihood f^L : $L(\underline{\theta}, \theta_0, \sigma^2 | \mathcal{D}) = P(\mathcal{D} | \underline{\theta}, \theta_0, \sigma^2)$

$$\left. \begin{aligned} & (i.i.d.) = \prod_{t=1}^n P(y_t | x_t, \underline{\theta}, \theta_0, \sigma^2) \\ & = \prod_{t=1}^n N(y_t; \underline{\theta}^T x_t + \theta_0, \sigma^2) \end{aligned} \right\}$$

Note: How Likely We see the dataset under the current parameters $(\underline{\theta}, \theta_0, \sigma^2)$

Notation:

σ^2 : models error that are not captured by Linear Model.

② log-likelihood:

$$\ell(\underline{\theta}, \theta_0, \sigma^2 | \mathcal{D}) = \sum_{t=1}^n \log \{N(y_t; \underline{\theta}^T x_t + \theta_0, \sigma^2)\}$$

$$= \sum_{t=1}^n \log \left[\frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2\sigma^2} [y_t - (\underline{\theta}^T x_t + \theta_0)]^2 \right\} \right].$$

$$= \sum_{t=1}^n \left[\text{const} - \frac{1}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (y_t - (\underline{\theta}^T x_t + \theta_0))^2 \right]$$

$$= \text{const} - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{t=1}^n (y_t - \underline{\theta}^T x_t - \theta_0)^2$$

(Simplify)

Note: the estimation of $(\underline{\theta}, \theta_0)$ is decoupled from that of b^2 .

\Rightarrow we can estimate $(\underline{\theta}, \theta_0)$ first

↓

$$\text{Minimize}_{\underline{\theta}, \theta_0} \quad \boxed{\sum_{t=1}^n (y_t - \underline{\theta}^T x_t - \theta_0)^2} \quad \rightarrow \text{massage into matrix form!}$$

① Define the target vector $\underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$ y_t : t-th label/target.

② Define the design matrix $X = \begin{bmatrix} x_1^T & 1 \\ \vdots & \vdots \\ x_n^T & 1 \end{bmatrix}_{n \times d+1}$

$$\sum_{t=1}^n (y_t - \underline{\theta}^T x_t - \theta_0)^2 = \sum_{t=1}^n (y_t - [x_t^T \ 1] \begin{bmatrix} \underline{\theta} \\ \theta_0 \end{bmatrix})^2$$

$$= \left\| \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} x_1^T & 1 \\ \vdots & \vdots \\ x_n^T & 1 \end{bmatrix} \begin{bmatrix} \underline{\theta} \\ \theta_0 \end{bmatrix} \right\|^2.$$

$\underbrace{\underline{y}}_{\text{y}}$ $\underbrace{X}_{\text{x}}$ $\underbrace{\begin{bmatrix} \underline{\theta} \\ \theta_0 \end{bmatrix}}_{\text{parameters}}$.

$$= \| \underline{y} - X \begin{bmatrix} \underline{\theta} \\ \theta_0 \end{bmatrix} \|^2 = \| \underline{y} - X \widetilde{\theta} \|^2 \quad (\widetilde{\theta} = \begin{bmatrix} \underline{\theta} \\ \theta_0 \end{bmatrix}_{(d+1) \times 1})$$

(Understanding Of Decouple)

Firstly, we maximize log-likelihood:

$$\ell(\underline{\theta}, \theta_0, b^2 | \mathcal{D}) = \text{const} - \frac{n}{2} \log(b^2) - \frac{1}{2b^2} \sum_{t=1}^n (y_t - \underline{\theta}^T x_t - \theta_0)^2$$

$$= f_1(b^2) - f_2(b^2) \cdot f_3(\underline{\theta}, \theta_0)$$

$\underbrace{f_1(b^2)}_{\text{maximize}} \quad \underbrace{f_2(b^2)}_{\min} \quad f_3(\underline{\theta}, \theta_0)$

$\widetilde{\theta} = \begin{pmatrix} \underline{\theta} \\ \theta_0 \end{pmatrix}$

$$\frac{\partial}{\partial \underline{\theta}} l = 0 \Leftrightarrow \frac{\partial}{\partial \underline{\theta}} f_3 = 0 \quad \text{decouple!}$$

Now we want to: $\min_{\tilde{\theta}} \|y - X\tilde{\theta}\|^2$

$$\begin{aligned} &\Leftrightarrow \min_{\tilde{\theta}} (y - X\tilde{\theta})^T (y - X\tilde{\theta}) \\ &\Leftrightarrow \min_{\tilde{\theta}} \underbrace{\tilde{\theta}^T X^T X \tilde{\theta} - 2\tilde{\theta}^T X^T y}_{g(\tilde{\theta})}. \end{aligned}$$

$$\frac{\partial g}{\partial \tilde{\theta}} = 2X^T X \tilde{\theta} - 2X^T y = 0 \Rightarrow \tilde{\theta} = (X^T X)^{-1} X^T y$$

Actually is $\hat{\theta}$!

That is to say, $\left(\frac{\hat{\theta}}{\hat{\theta}_0}\right) = (X^T X)^{-1} X^T y$

X has full column rank

Question 1: Why is $\hat{\theta}$ the globally minimum?

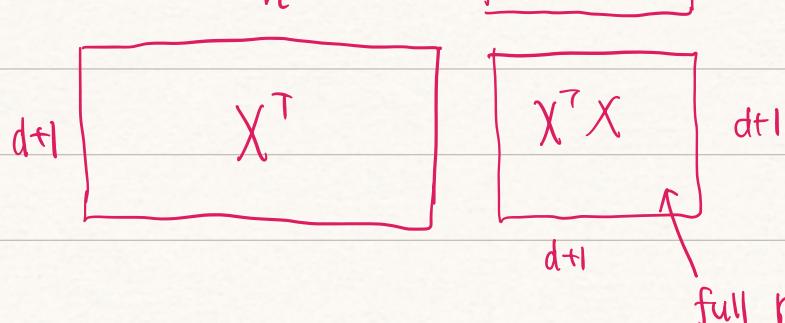
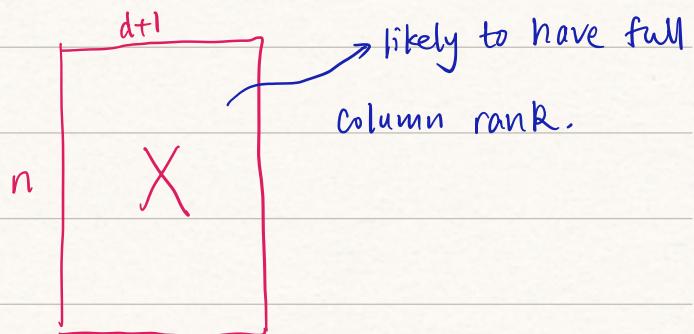
Quadratic Optimization

Question 2: What does X need to satisfy?

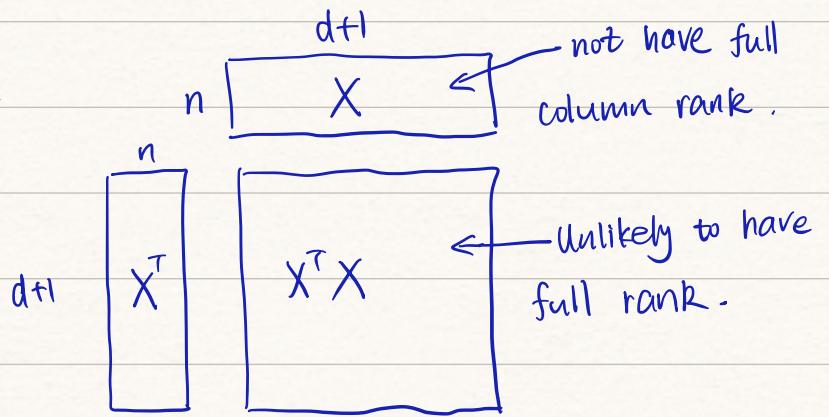
Q2. $X \rightarrow$ has full column rank. $\Rightarrow (X^T X)^{-1}$ exists

$$X \in \mathbb{R}^{n \times (d+1)} \Rightarrow$$

$(d+1) < n$



Problem: if X looks like



Rmk: i) $\begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_0 \end{bmatrix}$ is a linear function of target vector y .

ii) $\begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_0 \end{bmatrix}$ depends in a highly non-linear way on $[x_1, \dots, x_n]$

And it is coupled through the pseudo-inverse $(X^T X)^{-1} X^T$

iii) To be able to invert $X^T X$, we need X to have full column rank. (i.e. all columns in X are linearly independent).



Not likely to happen if $d+1 \gg n$



High-dimension curse.

$\left\{ \begin{array}{l} d \approx 19K \sim 20K \text{ protein coding genes} \\ n \approx 1,000 \text{ patients (subjects)} \end{array} \right.$

Now we have $\begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_0 \end{bmatrix} = (X^T X)^{-1} X^T y$

$$\ell(\underline{\theta}, \theta_0, \theta^2 | \underline{x}) = -\frac{n}{2} \log(\theta^2) - \frac{1}{2\theta^2} \sum_{t=1}^n (y_t - \underline{\theta}^T \underline{x}_t - \theta_0)^2$$

$$\left. \frac{d}{d\theta^2} \ell = 0 \Rightarrow \hat{\theta}^2 = \frac{1}{n} \sum_{t=1}^n (y_t - \underline{\hat{\theta}}^T \underline{x}_t - \hat{\theta}_0)^2 \right)$$

Meaning of this (Does it makes sense ?)

This is the average squared prediction error & captures

the "loss" non-ideal nature of the linear model!

Step 1: $\hat{\theta}, \hat{\theta}_0 \rightarrow$ Step 2: \hat{e}^2 .

Bias - Variance TradeOff.

Note that: $\hat{\theta}, \hat{\theta}_0, \hat{e}^2$ are RANDOM.

$$\left[\begin{array}{c} \hat{\theta} \\ \hat{\theta}_0 \end{array} \right] = \underbrace{(X^\top X)^{-1} X^\top}_{\text{not random}} \underbrace{y}_{\text{Random!}} \Rightarrow y_t = \underline{\theta}^\top \underline{x}_t + \theta_0 + \varepsilon_t$$

Suppose \mathcal{D} is indeed generated from some linear model.

$$\text{i.e. } y_t = \underline{\theta}^* \underline{x}_t + \theta_0^* + \varepsilon_t \quad t=1,2,\dots,n$$

for some $\underline{\theta}^* \in \mathbb{R}^d, \theta_0^* \in \mathbb{R}$

Here $\underline{\theta}^*$ & θ_0^* are the ground truth parameters

$$y = X \left[\begin{array}{c} \underline{\theta}^* \\ \theta_0^* \end{array} \right] + \varepsilon \quad \left\{ \begin{array}{l} E(\varepsilon) = 0 \\ \text{cov}(\varepsilon) = E(\varepsilon \varepsilon^\top) = \sigma^2 I_n \end{array} \right.$$

$$\left[\begin{array}{cc} E(\varepsilon^2) & E(\varepsilon \varepsilon^\top) \\ E(\varepsilon \varepsilon^\top) & E(\varepsilon^2) \end{array} \right] = \left[\begin{array}{cc} \sigma^2 & 0 \\ 0 & \sigma^2 \end{array} \right]$$

By the Previous result: $\left[\begin{array}{c} \hat{\theta} \\ \hat{\theta}_0 \end{array} \right] = (X^\top X)^{-1} X^\top y \rightsquigarrow y = X \left[\begin{array}{c} \underline{\theta}^* \\ \theta_0^* \end{array} \right] + \varepsilon$

$$\left[\begin{array}{c} \hat{\theta} \\ \hat{\theta}_0 \end{array} \right] = (X^\top X)^{-1} X^\top \left(X \left[\begin{array}{c} \underline{\theta}^* \\ \theta_0^* \end{array} \right] + \varepsilon \right)$$

$$\left[\begin{array}{c} \hat{\theta} \\ \hat{\theta}_0 \end{array} \right] = (X^\top X)^{-1} X^\top \left(X \left[\begin{array}{c} \underline{\theta}^* \\ \theta_0^* \end{array} \right] + \varepsilon \right)$$

$$= \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} + (X^T X)^{-1} X^T \underline{\varepsilon}$$

Q1: What's the BIAS?

$$\hookrightarrow \mathbb{E} \left[\left[\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} - \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} \right] \mid X \right] = \underline{0} \Rightarrow \text{BIAS is } \underline{0}.$$

for the least square sol²

↓ data.

Rmk: In Other Words, the estimate $\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} = (X^T X)^{-1} X^T \underline{\varepsilon}$ is UNBIASED!

Q2: What is the uncertainty of the Parameter Estimate?

$$\hookrightarrow \text{cov} \left(\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} \mid X \right) = \mathbb{E} \left[\left(\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} - \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} \right) \left(\begin{bmatrix} \hat{\theta} \\ \hat{\theta}_0 \end{bmatrix} - \begin{bmatrix} \theta^* \\ \theta_0^* \end{bmatrix} \right)^T \mid X \right]$$

Name of matrix:

(1) $X^T X \rightsquigarrow$ Gram Matrix

(2) $X \rightsquigarrow$ Design Matrix

Symmetry

$$= \mathbb{E} \left[(X^T X)^{-1} X^T \underline{\varepsilon} \underline{\varepsilon}^T X ((X^T X)^{-1})^T \mid X \right]$$

$$= (X^T X)^{-1} X^T \mathbb{E} [\underline{\varepsilon} \underline{\varepsilon}^T \mid X] X ((X^T X)^{-1})^T$$

$$= 6^2 (X^T X)^{-1}$$

[Prop] → if $A = A^T$, then $(A^{-1})^T = A^{-1}$

Rmk: Uncertainty depends on the inputs X

$$X^T X = \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix} \begin{bmatrix} x_1^T & 1 \\ \vdots & \\ x_n^T & 1 \end{bmatrix} = \sum_{t=1}^n [x_t] [x_t^T 1]$$

$$= n \left(\frac{1}{n} \sum_{t=1}^n [x_t] [x_t^T 1] \right)$$

Average of some matrices

①

(Large Number Thm)

Note: If $x_t \xrightarrow{i.i.d} p(x)$, then we have:

$$\frac{1}{n} \sum_{t=1}^n \begin{bmatrix} x_t \\ 1 \end{bmatrix} \begin{bmatrix} x_t^\top \\ 1 \end{bmatrix} \xrightarrow{P} C \rightsquigarrow C = \underset{x \sim P}{E} \left[\begin{bmatrix} x \\ 1 \end{bmatrix} \begin{bmatrix} x^\top \\ 1 \end{bmatrix} \right]$$

Roughly, $(X^\top X) \xrightarrow{P} n \cdot C \Rightarrow (X^\top X)^{-1} \xrightarrow{P} \frac{1}{n} \cdot \tilde{C}$

Therefore $\text{cov}[\varepsilon] \xrightarrow{P} \frac{b^2}{n} C^{-1}$

$$\tilde{C} = C^{-1}$$

(2)

Fact 1: $E[\|z - z^*\|^2] = \|E[z] - z^*\|^2 + E\|z - E[z]\|^2$

MSE (Mean Squared Error) $\cancel{\text{estimator (r.v.)}}$ \downarrow ground truth parameter $\underbrace{\text{Bias}^2}_{\text{Variance of } z \cdot \text{Tr}(\text{Cov}(z))}$

Pf: $E[\|z - z^*\|^2] = E[\|z - E[z] + E[z] - z^*\|^2]$ r.v. const
 $= E[\|z - E[z]\|^2 + \|E[z] - z^*\|^2 + 2 \langle z - E[z], E[z] - z^* \rangle]$
 $= E[\|z - E[z]\|^2 + \|E[z] - z^*\|^2] = 0$

(3)

Observation $\|a\|^2 = \text{Tr}(a \cdot a^\top) = a^\top a$.

Fact 2. Variance $= E[\|z - E[z]\|^2] = \text{Tr}(\text{Cov}(z))$

Pf: $E[\|z - E[z]\|^2]$
 $= E[(z - E[z])^\top (z - E[z])]$
 $= \text{Tr}(E[(z - E[z])^\top (z - E[z])])$
 $= \text{Tr}(E[(z - E(z))(z - E(z))^\top])$
 $= \text{Tr}(\text{Cov}(z))$

For this Problem, we already have:

$$C = \underset{x \sim p}{E} \left[\begin{bmatrix} x \\ 1 \end{bmatrix} \begin{bmatrix} x^\top \\ 1 \end{bmatrix} \right]$$

$$\begin{cases} E[\hat{\theta} - \hat{\theta}^*] = 0 \\ \text{Cov}[\hat{\theta}] = b^2 (X^\top X)^{-1} \end{cases}$$

$$\Rightarrow \text{MSE} = E[\|\hat{\theta} - \hat{\theta}^*\|^2] = \text{Tr}(b^2 (X^\top X)^{-1})$$

$$= b^2 \text{Tr}(X^\top X)^{-1}$$

Remind:

$$(X^\top X)^{-1} \xrightarrow{P} \frac{1}{n} \cdot C^{-1}$$

$$C \in \mathbb{R}^{(d+1) \times (d+1)}$$

$$\approx \frac{6^2}{n} \text{Tr}(C)$$

Rmk: i) If $n \rightarrow \infty$, MSE $\rightarrow 0$

ii) $\text{MSE} \propto 6^2$

iii) $\text{Tr}(C)$ is of order $d+1$ (is # of parameters in $(\underline{\theta}, \theta_0)$)

Penalized log-likelihood & Bayesian Inference.

Problem: $X^T X$ may not be invertible!



Idea: Put a prior probability distribution on $\left[\begin{array}{c} \underline{\theta} \\ \theta_0 \end{array} \right]$.

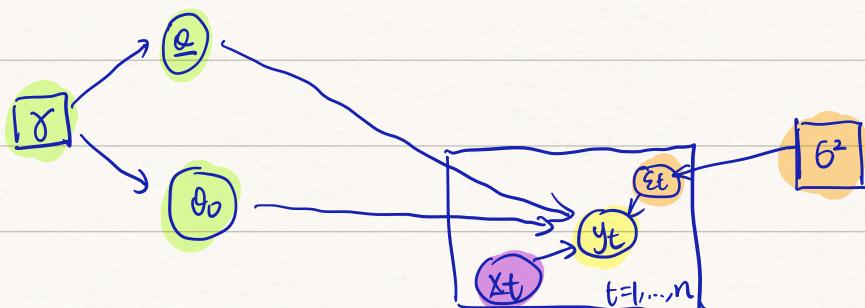
$$P(\underline{\theta}, \theta_0, \gamma) = N\left(\left[\begin{array}{c} \underline{\theta} \\ \theta_0 \end{array} \right]; \left[\begin{array}{c} \underline{\theta} \\ \theta_0 \end{array} \right], \gamma I\right)$$

$$\equiv N(\theta_0; 0, \gamma) \cdot \prod_{i=1}^d N(\theta_i; 0, \gamma)$$

Motivation:

Encouraging the parameters $(\underline{\theta}, \theta_0)$ shrunk to be $(\underline{\theta}, 0)$

Graphical Model (Bayesian Network)



- {
- Data
- Parameter
- Observation

$$y_t = \underline{\theta}^T x_t + \theta_0 + \varepsilon_t \quad t=1, 2, \dots, n$$

let $\underline{\omega}$ be all parameters!

{ Before, we use MLE \Rightarrow maximize $L(\underline{\theta}, \theta_0, \sigma^2 | \mathcal{D}) = P(\underline{y} | X, \underline{\omega})$

Now, we make $\underline{\omega} \rightarrow$ r.v. \Rightarrow MAP \Rightarrow max $P(\underline{\omega} | X, \underline{y})$

$$\Leftrightarrow \max \frac{P(\underline{\omega}, \underline{y} | X)}{P(\underline{y} | X)} \Leftrightarrow P(\underline{y} | X, \underline{\omega}) \cdot P(\underline{\omega})$$

$$\Leftrightarrow \max \underbrace{P(\underline{y} | X, \underline{\omega})}_{\text{posterior}} \cdot \underbrace{P(\underline{\omega})}_{\text{prior}}$$

Penalized log-likelihood (Regularization)

$$\hookrightarrow \ell'(\underline{\theta}, \theta_0, \sigma^2, \gamma) = \text{const} - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_t (y_t - \underline{\theta}^T \underline{x}_t - \theta_0)^2 - \frac{1}{2\gamma} (\theta_0^2 + \sum_{j=1}^d \theta_j^2) - \frac{d+1}{2} \log(\gamma)$$

Suppose σ^2 & γ are coupled as $\lambda = \frac{\sigma^2}{\gamma}$ (FIXED)

Prior of $(\underline{\theta}, \theta_0) \sim N([\underline{\theta}_0], [\theta_0^2] / \lambda I)$

$$= \text{const} - \frac{n+d+1}{2} \log(\sigma^2) + \frac{d+1}{2} \log\left(\frac{\sigma^2}{\lambda}\right) - \frac{1}{2\sigma^2} \left[\sum_t (y_t - \underline{\theta}^T \underline{x}_t - \theta_0)^2 + \lambda(\theta_0^2 + \sum_{j=1}^d \theta_j^2) \right]$$

$\downarrow \lambda$

Note that we have NO γ !

a function $f^* = f^*(\underline{\theta}, \theta_0, \sigma^2, \lambda)$

Also have the decoupled structure of $(\underline{\theta}, \theta_0)$

\hookrightarrow We have the solution: $\begin{bmatrix} \hat{\underline{\theta}} \\ \hat{\theta}_0 \end{bmatrix} = \underbrace{(X^T X + \lambda I)^{-1} X^T y}_{\text{guarantee } (X^T X + \lambda I) \text{ invertible}}$

Q1: But what's the BIAS? Suppose dataset is generated from the

noisy linear model: $y_t = \langle \underline{\theta}^*, \underline{x}_t \rangle + \theta_0^* + \varepsilon_t$

Answer is NO

$$\begin{aligned} \mathbb{E}\left[\begin{bmatrix} \hat{\underline{\theta}} \\ \hat{\theta}_0 \end{bmatrix} | X\right] &= \mathbb{E}\left[\left(X^T X + \lambda I\right)^{-1} X^T \left(X \cdot \begin{bmatrix} \underline{\theta}^* \\ \theta_0^* \end{bmatrix} + \varepsilon\right) | X\right] \\ &= (X^T X + \lambda I)^{-1} X^T X \begin{bmatrix} \underline{\theta}^* \\ \theta_0^* \end{bmatrix} \end{aligned}$$

$$= (X^T X + \lambda I)^{-1} (X^T X + \lambda I - \lambda I) \begin{bmatrix} \underline{\theta}^* \\ \theta_0^* \end{bmatrix}$$

$$= \begin{bmatrix} \underline{\theta}^* \\ \theta_0^* \end{bmatrix} - \lambda (X^T X + \lambda I)^{-1} \begin{bmatrix} \underline{\theta}^* \\ \theta_0^* \end{bmatrix}$$

$$\underbrace{\lambda \theta_0}_{\text{BIAS}} + \underbrace{\epsilon}_{\text{VARIANCE}}$$

Conclusion: If $\lambda \neq 0$, then BIAS $\neq 0$.

Q2: Recall $MSE = \text{bias}^2 + \text{variance}$

$$= \begin{cases} 0 + b^2 \text{tr}[(X^T X)^{-1}], & \lambda = 0 \\ \downarrow \quad \downarrow \\ \text{some bias} + \text{small variance}, & \lambda \neq 0 \end{cases}$$

Review

Linear regression \rightarrow Unregularized Case $\left\{ \begin{array}{l} \text{Model: } (X^T X)^{-1} X^T y \\ \text{Solution: } (X^T X)^{-1} X^T y \\ \text{Property of Solution: } \begin{cases} \text{Unbiased.} \\ \text{cov}(\hat{\theta}) = b^2 (X^T X)^{-1} \end{cases} \end{array} \right.$

Bias - Variance Trade-off

$$\hookrightarrow MSE = \text{bias}^2 + \text{variance}$$

\hookrightarrow In no regularization case:

$$\begin{cases} \text{bias} = 0 \\ \text{variance} = \text{tr}(b^2 (X^T X)^{-1}) \end{cases}$$

To deal with the
high dimension curse

\downarrow This may be quite big

Regularization Case.

\Rightarrow By introducing a prior distribution over (Ω, θ_0)

$\left\{ \begin{array}{l} \text{Model} \\ \text{Solution} \rightarrow (X^T X + \lambda I)^{-1} X^T y \\ \text{Property} \rightarrow \begin{cases} \text{Biased.} \\ \text{Smaller Variance} \end{cases} \end{array} \right. \Rightarrow \text{Bias - Variance Tradeoff}$

① The calculation of Penalized - LR.

$$\text{const} - \frac{n+d+1}{2} \log(\theta^2) + \frac{d+1}{2} \log\left(\frac{\theta^2}{\lambda}\right) - \frac{1}{2\theta^2} \left[\sum_t (y_t - \theta^T x_t - \theta_0)^2 + \lambda (\theta_0^2 + \sum_{j=1}^d \theta_j^2) \right]$$

$\downarrow \lambda$

$$\min \sum_t (y_t - \theta^T x_t - \theta_0)^2 + \lambda [\sum \theta_j^2 + \theta_0^2]$$

$$\Leftrightarrow \min \|y - X\theta\|_2^2 + \lambda \|\theta\|^2$$

$$\Leftrightarrow \min (y - X\theta)^T (y - X\theta) + \lambda \theta^T \theta$$

$$\Leftrightarrow \min y^T y - 2y^T X\theta + \theta^T X^T X\theta + \lambda \theta^T \theta$$

$$\Leftrightarrow \min y^T y - 2\theta^T X^T y + \theta^T (X^T X + \lambda I) \theta$$

$$\Rightarrow \frac{\partial J}{\partial \theta} = 0 \Rightarrow \theta = (X^T X + \lambda I)^{-1} X^T y$$

$$\textcircled{2} \quad \hat{\theta} = (X^T X + \lambda I)^{-1} X^T y$$

$$= (X^T X + \lambda I)^{-1} X^T (X\theta^* + \varepsilon)$$

$$= \underbrace{(X^T X + \lambda I)^{-1} X^T X\theta^*}_{\text{Expectation}} + \underbrace{(X^T X + \lambda I)^{-1} X^T \varepsilon}_{\text{Random}}$$

Expectation



$$\theta^* - \lambda (X^T X + \lambda I)^{-1} \theta^*$$



Random



{ } { }

Bias