# Bayesian Approches:

→ Freqnanties Approches: does not use prior prob.

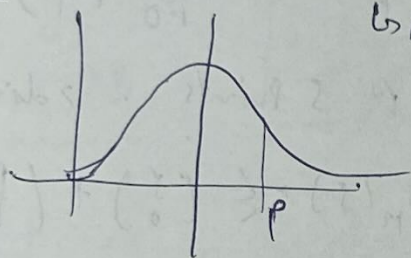but th Bayesian approches are use prior prob.

es $D = \{ H, T, 7H \dots \}$ find of bias of the coins.
$$P(H) = [0, 1]$$

+ **Bayesian**: distribution of values of bias
⟶ Means some perameter of that dist^n.



→ eth every inititeution prior will be updated.

$\Theta$ → bias, Parameters.

$$P(\Theta | D, \alpha) \quad \alpha \quad P(D | \Theta, \alpha) \cdot P(\Theta | \alpha)$$

Posterive         likelyhood        prior

# so Main thing is How get $\Theta$'s.
→ two Approches → ① MLE - use the Flernanties Approches
② MAP → Maxim A posteries

Q. Head with prob. $P$, tails w.p. $(1-P)$ n coins. tosses.

$y_i = 0$ tails
$y_i = 1$ Heads

$$D = \{ H T H H T \}$$
$$L(\Theta) = \prod_{j=1}^{n} P^{y_i} \cdot (1-P_i)^{(1-y_i)}$$

↳ $\Theta$ → parameter, likelyhood eqnestion
— this is not Prob we need to normalize them we get the prob.

$$\Theta_{MLE} = \underset{\Theta}{arg\,max} \ L(\Theta)$$

$$\therefore \underset{\theta}{\text{argmax}} \prod_{i=1}^{n} P(y_i \mid x_i, \theta)$$

$$= \underset{\theta}{\text{argmax}} \sum_{i=1}^{n} \log P(y_i \mid x_i, \theta)$$

$$= \underset{\theta}{\text{argmin}} \sum_{i=1}^{n} - \sum_{i=1}^{n} \log P(y_i \mid x_i, \theta)$$

$$\downarrow$$

(*) $\theta_{MAP} = \underset{\theta}{\text{argmax}} - E \ L(\theta) \cdot P(\theta)$

$$D = \{ (x_i, y_i) \}_{i=1 \ to \ n} \qquad x_i \in R^d$$
$$y_i \in R$$
$$y_i = w^T x_i + b \epsilon_i \qquad X = R^d$$
$$\underset{N(0, \ 1/\beta)}{\underbrace{\phantom{xxxx}}} \qquad Y = R$$
$$(\text{gaussian noise})$$

$$\text{Likehood} \quad \text{Likelihood} = \prod_{i=1}^{n} P(y_i \mid x_i, \overset{w}{\theta}, \beta)$$

→ here we assume that $y_i$ is also distr$^n$ on gaussian dist$^n$.

$$\text{mean is} \quad w^T x + \epsilon$$

$$= \prod_{j=1}^{n} \cancel{P(\epsilon_i)}$$

$$= \prod_{i=1}^{n} P(y_i \mid x_i, w, \beta) = \prod_{i=1}^{n} N(y_i \mid w^T x_i, \ 1/\beta)$$

$$\underset{\underset{\text{variance} \ r^2}{\uparrow}}{}$$

$$\boxed{= \prod_{i=1}^{n} \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \beta (y_i - w^T x_i)^2\right)}$$



$$\vdots$$

$$= \underset{(w, \beta)}{\text{argmin}} \sum_{i=1}^{n} - \log \sqrt{\frac{\beta}{2\pi}} + \sum_{i=1}^{n} \frac{1}{2} \beta (y_i - w^T x_i)^2$$

(if we know $\beta$ then we are to ERM problem)

→ but $\beta$ is unknown than find the $\beta$ & get the gradient that $\dfrac{\partial L(w, \beta)}{\partial w} = 0$ , $\dfrac{\partial L(w, \beta)}{\partial \beta} = 0$

No
auto
correlation

① $\varepsilon$ is independent $\omega^T x_i$

② $\varepsilon_1, \varepsilon \cdots \varepsilon_n$ are independent sample dist
③ Identical variant appose datapoints. $\Big\}$ homo skedusicity

$$\frac{\partial L(\omega, \beta)}{\partial \beta} = \sum_{i=1}^{n} \left( -\frac{1}{2\beta} + \frac{1}{2}(y_i - \omega^T x_i)^2 \right) = 0$$

$$\therefore \quad \frac{1}{\beta^4} = \frac{1}{n} \sum (y_i - \omega^T x_i)^2$$

$$\therefore \quad \sigma^{\frac{2}{4}}_{MLE} = \frac{1}{n} \sum (y_i - \omega^T x_i)^2$$

↳ is irreducible variance

$$\Rightarrow E_D[\omega_{MLE}] = E_D\left( (x^T x)^{-1} x^T y \right) \qquad y = x w + \varepsilon$$

↳ expectation of
`w`.

they is Randomness
because of Noise
$\varepsilon$ which random
Noise $N(0, 1/\beta)$

$\begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix} \; N(0, \beta$

$$= \left( (x^T x)^{-1} x^T E_D(y) \right)$$

$$= (x^T x)^{-1} x^T \left( \underline{E(x w) + E(\varepsilon)} \right)$$

$$= \underbrace{(x^T x)^{-1} x^T x}_{I} w$$

$$= W$$

$$\boxed{E_D(\omega_{MLE}) = W}$$

↳ not $\omega_{MLE}$ is unbias estimator of $w$.

$$* \; E_D \left( \frac{1}{\beta_{MLE}} \right) =$$

$$\overline{(xy)^T = y^T x^T}$$

$$\left( \frac{1}{\beta_{MLE}} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \omega^T x_i)^2 \right)$$

$$= E_D \left[ \frac{1}{n} (X W_{MLE} - y)^T (X W_{MLE} - y) \right]$$

$$= \frac{1}{n} E_D \left( W_{MLE}^T X^T X W_{MLE} - 2 W_{MLE}^T X^T y + y^T y \right)$$

$$= ((X^T X)^{-1} X^T y)^T = (y^T X (X^T X^{-1})^{-1})$$

$$= \frac{1}{n} E_D \left( y^T X (X^T X)^{-1} \underbrace{X^T X (X^T X)^{-1}}_{I} X^T y - 2 y^T X (X^T X)^{-1} X^T y + y^T y \right)$$

$$= \frac{1}{n} E_D \left( y^T X (X^T X)^{-1} X^T y - 2 y^T X (X^T X)^{-1} X^T y + y^T y \right)$$

$$\bar{=}$$

$$= \frac{1}{n} E_D \left( - y^T X (X^T X)^{-1} X^T y + y^T y \right)$$

$$= \frac{1}{n} E_D \left( y^T \underbrace{(I - X (X^T X)^{-1} X^T)}_{Z} y \right)$$

$$= \frac{1}{n} E_D \left[ (X w + \epsilon)^T Z (X w + \epsilon) \right]$$

this $\sigma$ is original 'w' not 'w$_{MLE}$'

$$= \frac{1}{n} E_D \left[ w^T X^T Z^{Xw} + \epsilon^T Z X w \!\!\!\!\diagup^0 + w^T X^T Z \epsilon \!\!\!\!\diagup^0 + \epsilon^T Z \epsilon \right]$$

$$= \frac{1}{n} \left[ w^T X^T Z X w + 0 + 0 + \epsilon^T Z \epsilon \right] \longrightarrow E(\epsilon^T \epsilon) \neq 0)$$

$$\longrightarrow \sum_{i=1}^{n} Z_{ij} \, E(\epsilon_i \epsilon_j)$$

$$\left( \because E(\epsilon_i \epsilon_j) = \underset{0}{E(\epsilon_i)} \; \underset{0}{E(\epsilon_j)} \right)$$

# scatter values

$\hookrightarrow$ this is for $i \neq j$ →but for $\boxed{i = j = 0}$ )

$$\text{tr}(ABC) = \text{tr}(CBA)$$

$$= \frac{1}{n}\left(w^T x^T Z x w + 0 + 0 + \frac{1}{\beta}\sum_{i=1}^{n}\text{diag}(Z)\right)$$

$$= \frac{1}{n}\left[w^T x^T x w - \underbrace{w^T x^T x (x^T x)^{-1}}_{I} x^T x w + \frac{1}{\beta}\sum_{i=1}^{n}\text{diag}(Z)\right]$$

$$= \frac{1}{n}\left(\frac{1}{\beta}\sum_{i=1}^{n}\text{diag}(Z)\right) = \boxed{\frac{1}{n}\left(\frac{1}{\beta}\text{Tr}(Z)\right)}$$

$$= \frac{1}{n}\left(\frac{1}{\beta}\text{Tr}(I - x(x^T x)^{-1}x^T)\right)$$

$$= \frac{1}{n}\left(\frac{1}{\beta}\underbrace{\text{Tr}(I)}_{3n} - \underbrace{\text{Tr}(x(x^T x)^{-1}x^T)}_{\substack{\downarrow\; x^T x \Rightarrow\; d\times d}}\right)$$

$$= \frac{1}{n}\left(\frac{1}{\beta}(n - D)\right)$$

$d$-dimension
$$\text{Tr}(x^T x) = d$$

$$= \left(\frac{n-D}{n}\right)\left(\frac{1}{\beta}\right)$$

$\downarrow$ dimension

$$\boxed{E_D\left(\frac{1}{\beta_{MLE}}\right) = \left(1 - \frac{D}{n}\right)\left(\frac{1}{\beta}\right)}$$

- Always under estimate the true variance
- If it is small rank so square matrix $E_D\left(\frac{1}{\beta_{MLE}}\right) = 0$
$\rightarrow$ if increase the $n$'s then bias will decrease.

$$\not{Var}(w_{MLE}) = \underbrace{E_D(w_{MLE}^2) - (E(w_{MLE}))^2}_{}$$ ( Var. of vector is variance cover matrix )

$$= E_D\left((w_{MLE} - E_D(w_{MLE}))(w_{MLE} - E_D(w_{ME})^T\right)$$

$$= E_D\left((w_{MLE} - w)(w_{ME} - w)^T\right)$$

$$= E_D\left(\underbrace{w_{MLE}^T w_{MLE} - 2 w_{MLE} w^T + w^T w}_{}\right)$$

$$= E_D\left(\overline{\left((X^TX)^{-1}X^Ty\right)^T}(X^TX)^{-1}X^Ty\right).$$

$$= E_D\left((X^TX)^{-1}X^Ty - w)(X^TX)^{-1}X^Ty - w)^T\right)$$

$$= E_D\left(\left((X^TX)^{-1}X^Ty\right)^T(X^TX)^{-1}X^TX - \right.$$

$$= E_D\left((X^TX)^{-1}X^T(\underset{w^Tx}{Xw^T}+E) - w) \cdot\right.$$

$$\left.\left((X^TX)^{-1}X^T(\underset{w^Tx}{Xw^T}+E) - w)^T\right)$$

$$= E_D\left((X^TX)^{-1}X^Tw^Tx + (X^TX)^{-1}X^TE - w) \cdot\right.$$

$$\left.\left((X^TX)^{-1}X^Tw^Tx + (X^TX)^{-1}X^TE_{\otimes} - w)^T\right)$$

$$= E_D\left([(X^TX)^{-1}X^Tw^Tx + (X^TX)^{-1}X^TE - w] \cdot\right.$$

$$\left.\left((X^TX)^{-1}X^Tw^Tx\right)^T + E^T(X^TX)^{-1}X - w^T\right)_{0}$$

$$= E_D\left(\left((X^TX)^{-1}X^Tw^Tx\right)^T(X^TX)^{-1}X^Tw^Tx) + (X^TX)^{-1}X^Tw^TxE^T(X^TX)^{-1}x\right.$$

$$- (X^TX)^{-1}X^Tw^Txw^T + 0 + 0 + 0 - w((X^TX)^{-1}X^Tw^Tx))^T$$

$$= \frac{H}{\beta}(X^TX)^{-1}$$

* Posterior $\propto$ prior $\cdot$ $\underset{\downarrow \text{maximize this}}{\underline{\text{likelyhood}}}$

* <u>Gauss markhov them</u>:

$W_{MLE}$ is BLUE $\downarrow$ $w$ (we want best $\underset{\text{(small)}}{\text{variance}}$

$\underset{\text{best linear unbians}}{\underbrace{\qquad}}$ $\rightarrow$ estimator

$\tilde{\omega}$ is another linear unbiased estimator

$$V(\tilde{\omega}) \ge \underset{\underbrace{\text{need to prove this}}}{Var(W_{MLE})}$$

$\rightarrow Var(\tilde{\omega}) = E\left((\tilde{\omega} - E(\tilde{\omega}))(\tilde{\omega} - E(\tilde{\omega}))^T\right)$

$= E\left((\tilde{\omega} - A\times w)(\tilde{\omega} - A\times w)^T\right)$ $\quad (\because E(\tilde{\omega}) = w)$

$= E\left((A\epsilon)(A\epsilon)^T\right)$

$= E(A\epsilon\epsilon^T A^T) = A E(\epsilon\epsilon^T) A^T$ $\quad$ because it is unbiased)

$\boxed{Var(\tilde{\omega}) = \dfrac{1}{\beta} A A^T}$

$\tilde{\omega} = A \cdot y$

$\underset{\text{some linear estimator}}{} = A(Xw + \epsilon)$

$\rightarrow \tilde{\omega} = A X w + A \epsilon$

$A X w = \tilde{\omega}$

$\boxed{A x = I}$ $\qquad \underset{= 0}{\overset{E(A\epsilon)}{\downarrow}}$

$\therefore$ Now: $Var(\tilde{\omega}) - Var(W_{MLE})$

$= \dfrac{1}{\beta} A A^T - \dfrac{1}{\beta}(x x^T)^{-1}$

$= \dfrac{1}{\beta}\left(A A^T - (x x^T)^{-1}\right)$

$\mathrel{\text{\large$\llcorner$}}$ some little settle up is here.

$= \dfrac{1}{\beta}\left(((x^T x)^{-1} x^T + B)((x^T x)^{-1} x^T + B)^T\right)$

$= \dfrac{1}{\beta}\left((x^T x)^{-1} x^T x (x^T x)^{-1} + \underset{T_1}{\underline{B x (x^T x)^{-1}}} + B B^T\right.$

$\left. + \underset{T_2}{\underline{\dfrac{(x^T x)^{-1} x^T B}{}}}\right)$

$Var(\tilde{\omega}) = \underset{Var(W_{MLE})}{\underbrace{\dfrac{1}{\beta}(x^T x)^{-1}}} + \underset{\text{this always a posi. sym. matrix}}{\underbrace{\dfrac{1}{\beta} B B^T}} + \dfrac{T_1}{\beta} + \dfrac{T_2}{\beta}$

$$( T_1 = BX (X^T X)^{-1} \quad \& \quad T_2 = (X^T X)^{-1} X^T B )$$

$$\text{So } A \Rightarrow (X^T X)^{-1} X^T + B$$

$$AX = (X^T X)^{-1} X^T X + BX$$

$$AX = I + BX$$

$$(AX - I) = BX$$

$$\quad\quad \overset{\hookrightarrow}{\cancel{AX=I}} = ( \because AX = I )$$

$$\boxed{BX = 0} \longrightarrow \text{So } T_1 = 0, \quad T_2 = 0$$

$$\therefore \text{Var}(\tilde{\omega}) - \text{Var}(\omega_{MLE}) = \frac{1}{\beta} B B^T$$

$$\therefore \quad \text{Var}(\tilde{\omega}) \gneq \text{Var}(\omega_{MLE})$$

$$\boxed{\text{Var}(\tilde{\omega}) - \text{Var}(\omega_{MLE}) \geq 0}$$

$\Rightarrow$ If we project $\omega_{MLE} \to \tilde{\omega}$ then point $t$ $\omega_{MLE}$ is tighter

**\* MAP :-**

$$y_i = \omega^T x_i + G_i \hookrightarrow N(0, {}^1/\beta)$$
$$\underset{\hookrightarrow \text{Prior}}{N(\mu_0, \Sigma_0)}$$

$$P(\omega \mid \underset{\text{Data}}{D}, {}^1/\beta, \text{otherParam}) \propto P(D \mid \omega, {}^1/\beta) \cdot P(\omega)$$

$$: P(\omega \mid (x_1 y_1) \dots (x_n, y_n), {}^1/\beta) \propto P(y_1, y_2 \dots y_n \mid \omega, {}^1/\beta \\ , x_1, x_2 \dots x_n) \\ \cdot P(\omega \mid \mu_0, \Sigma_0)$$

$$P(\omega \mid y_1 \dots y_n)(x_1 \dots x_n)) = \underbrace{\prod_{i=1}^{n} P(y_i \mid x_i, \omega, {}^1/\beta)}_{\text{likelyhood}} P(\omega \mid \mu_0, \Sigma_0)$$

$$= \prod_{i=1}^{n} N(y_i \mid \omega^T x_i, {}^1/\beta) \cdot P(\omega \mid \mu_0, \Sigma_0)$$

$$\cdot \prod_{i=1}^{n} \frac{\sqrt{\beta}}{\sqrt{2\pi}} \, exp\left(-\tfrac{1}{2}(y_i - w^T x_i)^2 \cdot \beta\right)$$

$$\cdot \frac{1}{2\pi^{d/2} |\Sigma_0|} \, exp\left(-\tfrac{1}{2}(w-\mu_0)^T \Sigma_0^{-1}(w-\mu_0)\right)$$

$$-\log P(w|D, \Sigma_0, \mu_0, \tfrac{1}{\beta}) = -n\log\sqrt{\beta} + n\log\sqrt{2\pi}$$

$$+ \frac{\beta}{2} \sum_{i=1}^{n} (y_i - w^T x_i)^2 + d\log\sqrt{2\pi} + \tfrac{1}{2}\log|\Sigma_0|$$

$$+ \tfrac{1}{2}(w-\mu_0)^T \Sigma^{-1}(w-\mu_0)$$

$$0 = 0 + 0 - \beta \sum_{i=1}^{n}(y_i - w^T x_i)x_i + 0 + 0$$

$$+ (w-\mu_0)\Sigma_0^{-1}$$

$$0 = (w - \mu_0)\Sigma_0^{-1} - \beta(y - w^T x_i)x_i$$

$$\frac{d}{dw}\left\{ \frac{\beta}{2}(Xw - y)^T(Xw - y) \right\} + \tfrac{1}{2}(w-\mu_0)^T \Sigma_0^{-1}(w-\mu_0)$$

$$\frac{\beta}{2}(2(X^T X + w) +$$

$$0 = \frac{\beta}{2}\left(2(X^T X)w - 2(X^T y)\right) + \tfrac{1}{2}\left(2\Sigma_0^{-1}w - 2\Sigma_0^{-1}\mu_0\right)$$

$$\beta(X^T X + \Sigma_0^{-1})w - \beta X^T y - \Sigma_0^{-1}\mu_0 = 0$$

$$\beta(X^T X + \Sigma_0^{-1})w = \beta X^T y + \Sigma_0^{-1}\mu_0$$

$$w_{ME} = \left(\beta(X^T X + \Sigma_0^{-1})\right)^{-1}(\beta X^T y + \Sigma_0^{-1}\mu_0)$$

\* 1 lecture is missing

$$\boxed{W_{MAP} = \underbrace{(\beta(X^TX) + \Sigma_0^{-1})^{-1}}_{\text{Always is invertible}} \cdot (\beta X^T y + \Sigma_0^{-1} \mu_0)}$$

$$X^T X = PSD \qquad x^T A x \geq \forall x \neq \vec{0}$$
$$\Sigma_0^{-1} = PD \qquad \vec{x}^T A x > 0 \; \forall x \neq \vec{0}$$

Cases: ① $n = 0$, $\boxed{W_{MAP} = \mu_0}$

② Prior is $\underline{\text{spherical}} \rightarrow \Sigma_0^{-1} = \frac{1}{c}[I]$

$$W_{MAP} = (\beta(X^TX) + \Sigma_0^{-1})^{-1}(\beta X^T y + \Sigma_0^{-1} \mu_0) \curvearrowleft \text{some scalled version of identity.}$$

$$\rightarrow (\beta(X^TX) + \frac{1}{c}I)^{-1} \beta X^T y + \frac{1}{c}\mu_0$$

③ → if spherical prior which is infinitely broad. $(c \to \infty)$

$$W_{MAP} = \beta(X^TX)^{-1} \beta X^T y$$

$$\hookrightarrow W_{MAP} = W_{MLE} = W_{ERM}$$

④ $\mu_0 = 0$ with spherical prior

$$(\beta(X^TX)^{-1} + \frac{1}{c}I)^{-1} \beta X^T y$$

$$= (X^TX)^{-1} + \underbrace{(\frac{1}{\beta c}I)}_{\searrow \lambda}^{-1} \beta X^T y$$

⑤ Laplacian prior

$$P(\omega / \sigma_0, \mu_0) = \frac{1}{2\sigma_0} \exp\left(-\frac{||w - \mu_0||_1}{\sigma_0}\right)$$

$$\Rightarrow -\log P(\omega/\sigma_0, \mu_0) = -\log \frac{1}{2\sigma_0} + \log \frac{||w - \cancel{\text{two}}\mu_0||_1}{\sigma_0}$$

$$\partial\left(\hat{J}(\omega)\right) = 0 + ||w_0 - \cancel{\text{two}} \mu_0|| \quad \cancel{-\log \text{too}} \hookrightarrow \frac{||w - \mu_0||_1}{\text{not}}$$
$$\text{not differ}$$

$\rightsquigarrow$ $w_{AP}$ ( $\mu_0 : 0$ , prior spherical )

$$w_{MAP} = \left( (X^T X)^* + \lambda I \right)^{-1} X^T y$$

$$E(w_{MAP}) = E\left( (X^T X)^* + \lambda I \right)^{-1} X^T y )$$

$(y = w^{*T} x + \epsilon)$

$$= E\left( (X^T X)^* + \lambda I \right)^{-1} X^T (w^* x + \epsilon) )$$

$$= E\left( (X^T X)^* + \lambda I \right)^{-1} w^* X^T x + X^T \epsilon )$$

$$= \cancel{E(X^T X)} \left( (X^T X)^* + \lambda I \right)^{-1} E(w^* X^T x) + 0$$

$$= \left( (X^T X)^* + \lambda I \right)^{-1} X^T x \, w^*$$

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

$$E(w_{MAP}) = w^* - \lambda (X^T x + \lambda I)^{-1} w^*$$

$\rightarrow$ this is bias estimator

$\not\ast$ $V(w_{MAP}) = \underline{\quad}$

( do by your self )

$$w_{MAP} = \underbrace{(x^T x + \lambda I)^{-1} x^T}_{} y \qquad \left( Z = (x^T x + \lambda I)^{-1} x^T \right)$$

$$w_{MAP} = Z y$$

$$= Z(x w^* + \epsilon)$$

$$\boxed{E_D(w_{MAP}) = Z x w^*}$$

$$Var(w_{MAP}) = E\left( (Zy - Zxw^*)(Zy - Zxw^*)^T \right)$$

$$\phantom{=} = \cancel{E[[Zy} - E$$

$$= E\left( (Zy - Zxw^*)(Zy - Zxw^*)^T \right)$$

$$\overset{\cancel{6}}{=} \cancel{E\left( (Zy - Zxw^*)(y^T Z^T - w^{*T} x^T Z^T) \right)}$$

$$= Z\, E\left( (y - xw^*)(y - xw^*)^T \right) Z^T$$

$$= Z\, E\left[ (\epsilon)(\epsilon)^T \right] Z^T$$

$$= Z \frac{1}{\beta} I Z^T$$

$$= (x^T x + \lambda I)^{-1} x^T \cdot \frac{1}{\beta} I \cdot \left( (x^T x + \lambda I)^{-1} x^T \right)^T$$

$$= \frac{1}{\beta} \left( (x^T x + \lambda I)^{-1} x^T x (x^T x + \lambda I)^{-1} \right)$$

$$= \frac{1}{\beta} V\, diag\left( \frac{\sigma_i^2}{(\lambda + \sigma_i^2)} \right) V^T$$

—

$$x^T x \qquad\qquad U \Sigma V^T \qquad x \in R^{n \times}$$
$$\underset{pxd \cdot nxd}{} \qquad\qquad U \quad \Sigma \quad V^T$$
$$\underset{nxn}{}$$

$$= \emptyset\ V \Sigma \underset{\frac{1}{I}}{U^T U^T} \Sigma V \emptyset$$

$$\boxed{x^T x = V \Sigma^2 V} = V\, diag(\sigma_i^2) V^T$$

$$(X^TX)^{-1} = (V\Sigma^2V^T)^{-1} \qquad (AB)^{-1} = B^{-1}A^{-1}$$

$$= V\Sigma^{-2}V^T$$

$$\boxed{(X^TX)^{-1} = V \, diag\left(\frac{1}{\sigma_i^2}\right) V^2}$$

$\hookrightarrow$ do this always
full rank matrix

$$(\cancel{X^TX + \lambda I})^{-1} = \cancel{(X^TX)^{-1} + (\lambda I)^{-1}}$$

$$= \cancel{V \, diag\left(\frac{1}{\sigma_i^2}\right)V^2} + \cancel{V \, diag\left(\frac{1}{\sigma_i^2}\right)} \, \textcircled{w}$$

$$(X^TX + \lambda I)^{-1} = \cancel{X}\left(V\Sigma^2V^T + \lambda VIV^T\right)^{-1}$$

$$= \left(V(\Sigma^2 + \lambda I)V^T\right)^{-1}$$

$$\boxed{(X^TX + \lambda I)^{-1} = V \, diag\left(\frac{1}{\sigma_i^2 + \lambda I}\right)V^T}$$

$$\ast \, (X^TX + \lambda I)^{-1}X^T = V \, diag\left(\frac{1}{\sigma_i^2 + \lambda I}\right)V^T \cdot \cancel{diag(\sigma_i^2)}$$
$$(V\Sigma V^T)$$

$$\downarrow \, get$$

$$= V \, diag\left(\frac{\sigma_i^2}{\sigma_i^2 + \lambda I}\right)V^T$$

$$= V \, diag\left(\frac{1}{\sigma_i^2 + \lambda I}\right) \underset{\underset{I}{\uparrow}}{V^T V} \Sigma V^T$$

$$= V \, diag\left(\frac{1}{\sigma_i^2 + \lambda I}\right) diag(\sigma_i^2)V^T$$

$$\boxed{Z = (X^TX + \lambda I)^{-1}X^T = V \, diag\left(\frac{\sigma_i^2}{\sigma_i^2 + \lambda I}\right)V^T}$$

$$Z Z^T = V \, diag \left( \frac{\sigma_i^{\circ}}{\sigma_i^2 + \lambda I} \right) V^T \, V \, diag \left( \frac{\sigma_i^{\circ}}{\sigma_i^2 + \lambda I} \right) V^T$$

$$\boxed{Z Z^T = V \, diag \left( \frac{\sigma_i^2}{(\sigma_i^2 + \lambda I)^2} \right) V^T}$$

※ Ensemble Methods: $\underset{\downarrow}{\underline{Bagging}}$ & $\underset{\downarrow}{\underline{Boosting}}$

Parallel          Sequential

↪ <u>Bagging</u>: - Randomforest
- Use Multiple classifier

⊛ <u>Bootstrapping</u>: sampling technique where sample are derived from whole population using replacement procedure.

* Aggregation: