

1. Probability

(1)

Let event $(H=h) = A$, $(D=d) = B$

Then, by Bayes' Rule

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

(i) it depends on the probability of $P(B|A)$ and $P(B)$ if $P(B|A) \geq P(B)$

$$\text{then } P(H=h|D=d) \leq P(H=h)$$

otherwise

$$P(H=h|D=d) > P(H=h)$$

(ii)

$$P(H=h) \geq P(D=d|H=h) P(H)$$

since $P(B) \leq 1$

(2)

$$(i) E_Y[E_X[X|Y]] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x p_{X|Y}(x|y) dx p(y) dy$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y p_{Y|X}(y|x) dy dx = \int_{-\infty}^{\infty} x p(x) dx = E[X]$$

(ii)

$$\text{Var}[X] = E[X^2] - E[X]^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x^2 p(x|y) dx p(y) dy - E[X]^2$$

$$= E_Y[E_X[X^2|Y]] - E[X]^2 = E_Y[E_X[X^2|Y] - E_X[X|Y]^2 + E_X[X|Y]^2] - E_Y[E_X[X|Y]]^2$$

$$= E_Y[E_X[X^2|Y] - E_X[X|Y]^2] + E_Y[E_X[X|Y]^2] - E_Y[E_X[X|Y]]^2$$

$$= E_Y[\text{Var}_X[X|Y]] + \text{Var}_Y[E[X|Y]]$$

2 Maximum likelihood.

(1)

$$\begin{aligned} \arg \max_{\theta} l(\theta) &= \arg \max_{\theta} \sum_{i=1}^n \log f(x_i; \theta) \\ &= \arg \max_{\theta} -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \alpha)^2 \end{aligned}$$

$$\theta = \frac{\partial l(\theta)}{\partial \alpha} = \sum_{i=1}^n (x_i - \alpha) = 0 \Rightarrow \alpha = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$$

$$\theta = \frac{\partial l(\theta)}{\partial \sigma^2} = \frac{n}{2\sigma^2} - \frac{1}{2\sigma^4} \left(\sum_{i=1}^n (x_i - \alpha)^2 \right) \Rightarrow \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \alpha)^2$$

if we rule out $\sigma^2 = 0$

(2)

$$l(\theta) = \sum_{i=1}^n \ln f(x_i; \mu, \Sigma) = -\frac{nd}{2} \ln 2\pi - \frac{n}{2} \ln |\Sigma| - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$$\Rightarrow \theta = \frac{\partial l(\theta)}{\partial \mu} = \sum_{i=1}^n (x_i - \mu) = 0 \Rightarrow \mu^* = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$$

3 Weighted Least Squares Regression

(1)

$$\text{Let } Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, X = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}, W = \begin{bmatrix} b \\ B^T \end{bmatrix}, C = \begin{bmatrix} c_1 & 0 \\ 0 & c_1 \\ \vdots & \vdots \\ 0 & c_n \end{bmatrix}$$

$$\Rightarrow \text{objective function become } \arg \min_w \|C(Y - Xw)\|_2^2$$

\Rightarrow by taking derivative w.r.t. w , then we have

$$X^T C^T C (Y - Xw) = 0$$

$$0 = C^T X^T C (Y - Xw) \Rightarrow w = (C^T X^T C X)^{-1} C^T X^T C Y \quad w = ((X^T C^T C X)^{-1}) X^T C^T C Y$$

when $c_i = 1 \quad \forall i \in 1 \sim n$

then solution become $w = (X^T X)^{-1} X^T Y$

Which is equivalent to the solution of $\arg \min_w \|Y - Xw\|_2^2$
 then if X is a R.V. who is i.i.d. on x_1, x_2, \dots, x_n
 based on the Maximum likelihood method, where $w^* \in \theta$
 $\arg \max_{\theta} p(\theta|X)$ if we don't have prior knowledge of $p(\theta)$.

$$\begin{aligned}
 &= \arg \max_{\theta} \log p(X|\theta) \\
 &= \log \left(\prod_{i=1}^n p(y_i | x_i, \theta) p(x_i | \theta) \right) \\
 &= \log \left(\prod_{i=1}^n p(w x_i + \varepsilon) p(x_i | \theta) \right) \quad , \text{ since } x_i \text{ is not depend on } \theta \\
 &\propto \log \left(\prod_{i=1}^n p(\varepsilon | x_i, \theta) \right) \\
 &\propto \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - w x_i)^2}{2\sigma^2}\right) \right) \propto -\sum_{i=1}^n \frac{1}{2\sigma^2} (y_i - w x_i)^2 \\
 &\propto -\sum_{i=1}^n (y_i - w x_i)^2
 \end{aligned}$$

$$\Rightarrow \arg \max_{\theta} \log p(X|\theta) = \arg \min_w \|y_i - w x_i\|_2^2$$

(2)

if each ε_i has different variance
 then

$$\arg \max_{\theta} \log p(X|\theta) = \arg \min_w \left\| \frac{C_i}{\sigma_i} (y_i - w x_i) \right\|_2^2$$

\Rightarrow the same form as equation (b) on problem 3 (1)

As a result, with the same variance, the answer will be general linear regression,
 but with different variance the answer will be local weighted linear regression.

$$\Rightarrow \text{Let } \hat{C} = \begin{bmatrix} \frac{C_1}{\sigma_1} & 0 \\ 0 & \frac{C_n}{\sigma_n} \end{bmatrix} \quad \text{by problem 3, (1)}$$

$$w = (\hat{C}^T X^T C X \hat{C})^{-1} \hat{C}^T X^T C Y$$

4. Bias-Variance Decomposition.

(1)

$$\begin{aligned} E_{Y|X,D}[(y - f_b(x))^2] &= E_{Y|X,D}[(y - f(x) + f(x) - f_b(x))^2] \\ &= E_{Y|X}[(y - f(x))^2] + E_D[(f(x) - f_b(x))^2] + 2E[(y - f(x))(f(x) - f_b(x))] \end{aligned}$$

$$\left(\text{where } E_{Y|X,D}[(y - f(x))(f(x) - f_b(x))] = E_{Y|X}[y]f(x) - f(x)^2 - E_{Y|X,D}[y f_b(x)] + f(x) E_D[f_b(x)] \right)$$

Since $y = f(x) + \varepsilon$, & ε is zero mean.

$$= f(x)^2 - f(x)^2 - f(x) E_D[f_b(x)] + f(x) E_D[f_b(x)] = 0$$

$$= E_{Y|X}[(y - f(x))^2] + E_D[(f(x) - f_b(x))^2] \quad \#$$

(2)

$$E_{Y|X}[(y - f(x))^2] = E_{\varepsilon}[(f(x) + \varepsilon - f(x))^2] = E[\varepsilon^2]$$

$$\text{also } \because E[\varepsilon] = 0$$

$$\therefore = E[\varepsilon^2] - E[\varepsilon]^2 = \text{var}[\varepsilon] \quad \#$$

(3)

$$\begin{aligned} E_D[(f(x) - f_b(x))^2] &= E_D[(f(x) - E_D[f_b(x)] + E_D[f_b(x)] - f_b(x))^2] \\ &= (f(x) - E_D[f_b(x)])^2 + E_D[(E_D[f_b(x)] - f_b(x))^2] + 2\left\{ f(x) E_D[f_b(x)] - f(x) E_D[f_b(x)] - E_D[f_b(x)]^2 + E_D[f_b(x)]^2 \right\} \\ &= (f(x) - E_D[f_b(x)])^2 + E_D[(E_D[f_b(x)] - f_b(x))^2] \end{aligned}$$

(4)

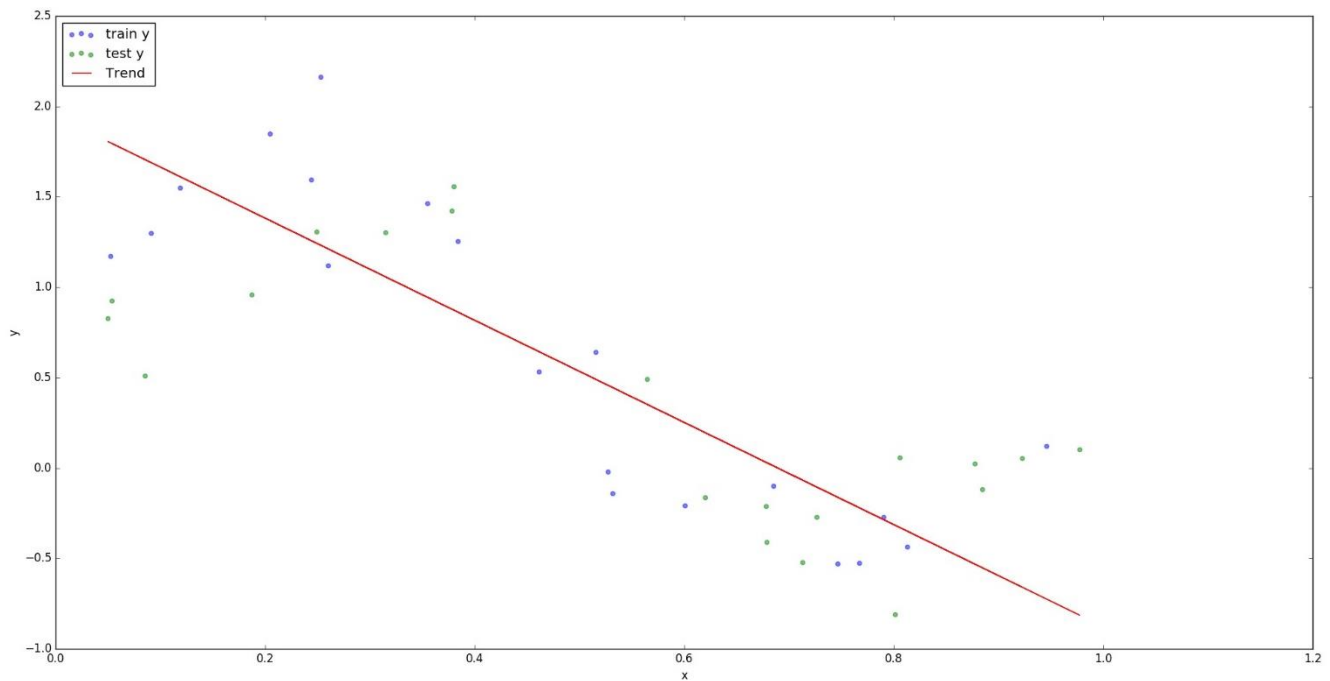
$$E_{X,Y,D}[(y - f_b(x))^2] \stackrel{\text{by (1)(2)}}{=} \text{var}[\varepsilon] + E_D[(f(x) - f_b(x))^2] =$$

$$\stackrel{\text{by (3)}}{=} \text{var}[\varepsilon] + \text{bias}(X) + \text{variance}(X) \quad \#$$

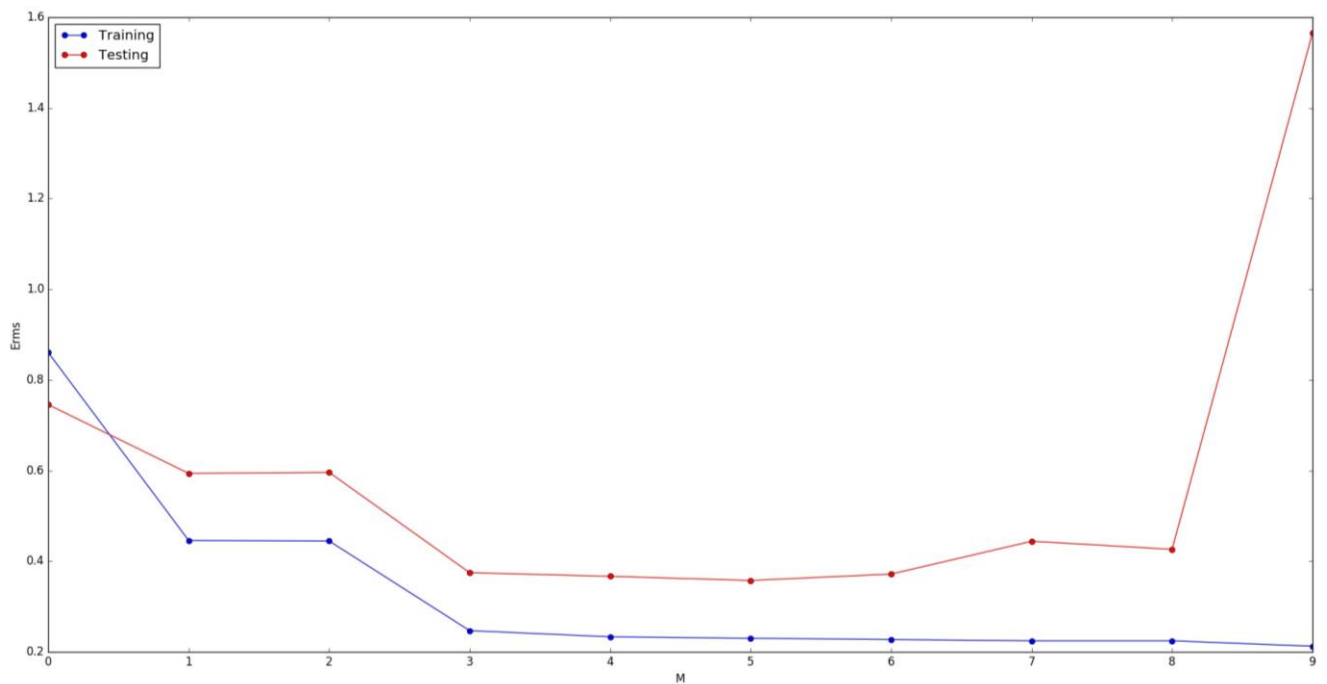
P5 Linear Regression

(1)

$$w = \begin{bmatrix} 1.9468 \\ -2.8241 \end{bmatrix}$$



(2)



(3)

In this question, I set $\lambda = 10^{-6}: 10^{13}$. The close form for regularized linear regression is

$$w = y\Phi^T(\Phi\Phi^T + \lambda I)^{-1}$$

Where Φ is a 9th degree polynomial of x .

