# 1 Questions from lecture notes

1. The key statistical assumption that we made to make the supervised inductive batch learning problem well-defined is:
TRAINING & DEPLOYMENT DATA ARE SAMPLES FROM SAME LIKELIHOOD.

[0.5 Mark]

2. The Bayes optimal classifier with 0-1 loss is the ___MODE___ of the underlying posterior likelihood (written as a function of input).

[0.5 Mark]

3. The component of the generalization error that is independent of the training set size is called: __model error__.

[0.5 Mark]

4. The formula for loss in logistic regression is given by: $l\,(\text{sign}(a), \text{sign}(b)) = \log\!\left(1 + e^{-ab}\right)$.

[0.5 Mark]

5. Suppose the underlying likelihood in a linear regression problem satisfies the equations:

$$Y = \sum_{i=1}^{n} \alpha_i X_i + N, \;\; \mathbb{E}[NX] = 0, \tag{1}$$

where $X, Y$ are the input,output. The necessary and sufficient conditions for parameters of a Bayes optimal linear regressor to be same as $\alpha$ are: __$\mathbb{E}[XX^T] \succ 0$__.

[1 Mark]

6. From bias-variance tradeoff discussion in linear models it is clear that $n \to \infty, m \to \infty$ is sufficient for good generalization. Here, $n, m$ are the number of parameters and training set size respectively. But should $n$ diverge like $\sqrt{m}$ or like $m$ or like $m^2$ ? __$\sqrt{m}$__. Fill in this blank with one of the three functions $\sqrt{m}, m, m^2$.

[0.5 Mark]

7. In applications of quantum information theory, the space (manifold) of positive definite (pd) matrices is often encountered and the standard loss function is the squared-Bures-Wasserstein metric: loss between $A$ and $B$ is given by $\left[\text{Tr } A + \text{Tr } B - 2\,\text{Tr}(A^{1/2}BA^{1/2})^{1/2}\right]$. While the standard loss on $\mathbb{R}$ is the squared-loss. Consider a regression problem with input space as that of pd matrices and real outputs, employing standard loss. The Bayes optimal, $f^*$, for this regression problem is defined by:

$$f^*(x) \equiv \underset{y \in \mathbb{R}}{\text{argmin}} \; \mathbb{E}\!\left[(y-Y)^2 / x\right]$$

Fill this blank with an expression involving the expression for the loss.

[1 Mark]

8. Consider a linear regression problem[1] where the underlying likelihood is such that $p(x, y) = p(x)p(y)$. Assume the mean and variance with $p(x)$ are 3,9 respectively. Assume the mean and variance with $p(y)$ are 4,36 respectively. Then, the simplified expression for the Bayes optimal linear regressor is:

$$f_{|L}^*(x) = \underline{2/3}\, x.$$

Fill this blank with a numeric constant.

[1.5 Marks]

9. Consider a linear regression problem[2] with training data (input,output pairs): $\mathcal{D} = \{(1, 2), (3, 4)\}$. Then, the simplified expression for the ERM linear regressor is:

$$\hat{f}_{|L}(x) = \underline{1\cdot4}\, x.$$

Fill this blank with a numeric constant.

[1.5 Marks]

## 2 Derivations done in lectures

10. Derive a simplified expression for the Bayes optimal restricted to the linear model over inputs with squared loss using only projection theorem as done in lectures.

$$\omega^* = \underset{\omega \in \mathbb{R}^n}{\arg\min}\; E\{(\omega^T x - y)^2\}$$
$$= \underset{\omega \in \mathbb{R}^n}{\arg\min}\; \|\omega^T x - y\|_H^2$$
$$= \underset{\omega \in \mathbb{R}^n}{\arg\min}\; \|\omega^T E\{xx^T\}^{1/2}\tilde{x} - y\|_H$$
$$\text{by } \tilde{x} = E\{xx^T\}^{-1/2} x$$

$$= E\{xx^T\}^{-1/2}\; \text{times}$$
$$\underset{\bar{\omega} \in \mathbb{R}^n}{\arg\min}\; \|\bar{\omega}^T \tilde{x} - y\|$$
$$= E\{xx^T\}^{-1/2} E\{\tilde{x}y\} \;\rceil$$
$$= E\{xx^T\}^{-1} E\{xy\} \quad \begin{array}{l}\text{Proj.}\\ \text{theorem.}\end{array}$$

[2 Marks]

## 3 Problems from other course page resources

11. Consider a problem where the underlying likelihood is defined by $p(x|y) \sim \mathcal{N}(\mu_y, \Sigma_y)$, $y = 1, 0$. Let $p(y) = \begin{cases} 0.4 & y = 1 \\ 0.6 & y = 0 \end{cases}$. Assume the loss is the 0-1 loss. Derive a simplified expression for the Bayes optimal.

---

[1] feature map $\phi(x) = x$.
[2] feature map $\phi(x) = x$.

$$p(\ell/n) \propto e^{-\frac{1}{2}(n-\mu_1)^T \Sigma_1^{-1}(n-\mu_1) + \ln 0.6}_{-\frac{1}{2}\ln|\Sigma_1|}$$

$$p(0/n) \propto e^{-\frac{1}{2}(n-\mu_0)^T \Sigma_0^{-1}(n-\mu_0) + \ln 0.6}_{-\frac{1}{2}\ln|\Sigma_0|}$$

$$f^*(n) = 1 \text{ if } sign(g(n)) \leq 0, \ 0 \text{ else}$$

$$g(n) = \frac{1}{2} x^T(\Sigma_1^{-1} - \Sigma_0^{-1}) - x^T(\Sigma_1^{-1}\mu_1 - \Sigma_0^{-1}\mu_0)$$

$$\ln 0.6/0.4 + \frac{1}{2}\mu_1^T \Sigma_1^{-1}\mu_1 - \frac{1}{2}\mu_0^T \Sigma_0^{-1}\mu_0 - \frac{1}{2}\ln\frac{|\Sigma_0|}{|\Sigma_1|}$$

[1 Mark]

If $\Sigma_1 = \Sigma_0$, then show that the model error with linear model over the feature map, $\phi(x) = \begin{bmatrix} x \\ 1 \end{bmatrix}$, is exactly zero.

$$g(n) \text{ becomes } -x^T(\Sigma_1^{-1}\mu_1 - \Sigma_0 \mu_0) + \ln 0.6/0.4 + \frac{1}{2}\mu_1^T \Sigma_1^{-1}\mu_1 - \frac{1}{2}\mu_0^T \Sigma_0^{-1}\mu_0 \ \underset{\text{linear}}{so}$$

$$-\frac{1}{2}\ln\frac{|\Sigma_0|}{|\Sigma_1|}$$

[0.5 Mark]

Prove that this condition is not necessary for the model error being zero[3]. In other words, provide a different condition on $\mu_1, \mu_0, \Sigma_1, \Sigma_0$, such that the model error is exactly zero.

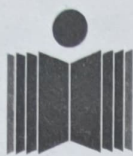model error is zero iff quadratic reduces to an affine or linear or constant in terms of its _sign_.

For alternate, $\Sigma_1 \neq \Sigma_0$, i.e. sign of quadratic be a constant

$\Longrightarrow$ discriminant $\leq 0$,

$x^T A x + 2 b^T n + c \gtrless 0 \ \forall n \Longleftrightarrow c - b^T A^{-1} b \gtrless 0$. So Alternate is

(A is invertible Aro)

$$\Sigma_1^{-1} - \Sigma_0^{-1} > 0, \ \mu_1^T \Sigma_1^{-1}\mu_1 - \mu_0^T \Sigma_0^{-1}\mu_0 + 2\ln\frac{0.6}{0.4} - \ln|\Sigma_0|/|\Sigma_1|$$

$$- (\Sigma_0^{-1}\mu_0 - \Sigma_1^{-1}\mu_0)^T (\Sigma_1^{-1} - \Sigma_0^{-1})^{-1}(\Sigma_0^{-1}\mu_0 - \Sigma_1^{-1}\mu_1) \gtrless 0$$

[4 Marks]

[3] is the solution provided for the practice problems partially wrong? Go home and think about the necessary and sufficient conditions for 0 model error. If you think you got them, meet me to discuss.

① First line on Pg 12 of summary notes

② First line, second para, Pg 16 of notes

③ First line, second para, Pg 20 of notes

④ Second last line, second last para, Pg 32 of notes

⑤ Last para Pg 23 – first para Pg 24 of notes

⑥ On
Second last para in Pg 23 of notes

⑥ ~~Necessity is not~~
$X_1, \ldots X_n$ are lin. independent $\iff$ $E(XX^T) \succ 0$

[Marks to either answer.

Based on Eqn. (10.2) and (11.1) in notes

$\dfrac{n}{m} \to 0$   $\dfrac{\sqrt{m}}{m} \to 0$   $\dfrac{m}{m} \to ?$   $\dfrac{m^3}{m} \to \infty$

⑦ Eqn (6.1) in notes . Substitute separate eqns.

⑧ Eqn (8.7) $\omega^* = E\{XX^T\}^{-1} E\{XY\}$  from notes

$E\{XY\} = E\{X\}E\{Y\}$   $(\because X \perp\!\!\!\perp Y)$

$= 3 \times 4 = 12$

$$E\{xx^T\} = E\{x^2\} = var(x) + (E\{x\})^2$$
$$= 9 + 3^2 = 18$$

$$\therefore \omega^* = \frac{12}{18} = \frac{2}{3} \qquad \therefore f_k^*(n) = \omega^{*T} x$$
$$= \frac{2}{3} x$$

⑨ From notes (8.8) $\hat{\omega} = \left(\frac{1}{m}\sum_i x_i x_i^T\right)^{-1}\left(\frac{1}{m}\sum_i x_i y_i\right)$

$$= \left(\frac{1}{2}(1\times1 + 3\times3)\right)^{-1}\left(\frac{1}{2}(1\times2 + 3\times4)\right)$$

$$= \frac{14}{10} = 1.4$$

$$\therefore f_k(n) = \hat{\omega}^T n = 1.4 n.$$

⑩ Theorem 3 on Pg 3 on hand-written derivations.

⑪ $p(y|n) \propto p(n|y)p(y)$

$p(1/n) \propto e^{-\frac{1}{2}(x-\mu_1)^T \Sigma_1^{-1}(x-\mu_1)} \quad 0.4 / |\Sigma_1|^{\frac{1}{2}}$

$= e^{-\frac{1}{2}(n-\mu_1)^T \Sigma_1^{-1}(n-\mu_1) + \ln 0.4 - \frac{1}{2}\ln|\Sigma_1|}$

$p(0/n) \propto e^{-\frac{1}{2}(n-\mu_0)^T \Sigma_0^{-1}(n-\mu_0) + \ln 0.6 - \frac{1}{2}\ln|\Sigma_0|}$
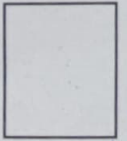
11) by $p(1/x) \geqslant p(0/n) \Longrightarrow -\frac{1}{2}(x-\mu_1)^T \Sigma_1^{-1}(n-\mu_1) + \ln 0.4 - \frac{1}{2}\ln|\Sigma_1|$

$$\geqslant -\frac{1}{2}(n-\mu_0)^T \Sigma_0^{-1}(n-\mu_0) + \ln 0.6 - \frac{1}{2}\ln|\Sigma_0|$$

$$\therefore f^*(x) = \begin{cases} 1 & \text{if } \frac{1}{2}x^T(\Sigma_1^{-1} - \Sigma_0^{-1})x - x^T(\Sigma_1^{-1}\mu_1 - \Sigma_0\mu_0) + \frac{1}{2}\mu_1^T\Sigma_1^{-1}\mu_1 \\ & \qquad -\frac{1}{2}\mu_0^T\Sigma_0^{-1}\mu_0 \\ & \qquad +\ln 0.6/0.4 \leq 0 \\ & \qquad -\frac{1}{2}\ln|\Sigma_0/\Sigma_1| \\ 0 & \text{else} \end{cases}$$

If $\Sigma_1 = \Sigma_0$, then

$$f^*(x) = \begin{cases} 1 & \text{if } \overbrace{-x^T(\Sigma_1^{-1}\mu_1 - \Sigma_0\mu_0)}^{\omega^*} + \overbrace{\frac{1}{2}\mu_1^T\Sigma_1^{-1}\mu_1 - \frac{1}{2}\mu_0^T\Sigma_0^{-1}\mu_0 + \ln\frac{0.6}{0.4}}^{b_0^*} \\ & \qquad\qquad\qquad\qquad -\frac{1}{2}\ln\frac{|\Sigma_0|}{|\Sigma_1|} \qquad\qquad \leq 0 \\ 0 & \text{else} \end{cases}$$

i.e. $f^*(x) = -\text{sign}(-\omega^{*T}x + b_0^*)$    $(-1 \to 0)$

$\therefore$ Model error is zero.

$\longrightarrow$  Can model error be zero $\begin{array}{l} \to \text{quadratic} \to \text{linear/affine} \\ \to \text{quadratic} \to \text{const.} \end{array}$

Can expression of quadratic behave like constant without $\Sigma_1 = \Sigma_0$?

Yes! if it's a non-negative quadratic $\left.\begin{array}{l} \\ \end{array}\right\}$ Discriminant $\leq 0$
     or non-positive quadratic

② marks if you write this.

$x^T A x + 2b^T x + c \geq 0 \quad \forall x$

$\Leftrightarrow \|A^{1/2}x + A^{-1/2}b\|^2 + c - b^T A^{-1} b \geq 0 \quad \forall x \quad$ (assuming $A \overset{>0}{\succ 0}$)

$\downarrow$

*If you also interested in knowing what happens if $A \succ 0$ then please meet me. or if $A$ is symmetric only.*

$\Leftrightarrow c - b^T A^{-1} b \geq 0$

$\therefore$ Alternate condition is

$\Sigma_1^{-1} - \Sigma_0^{-1} \succ 0, \quad \mu_1^T \Sigma_1^{-1} \mu_1 - \mu_0^T \Sigma_0^{-1} \mu_0 + 2 \ln \frac{0.6}{0.4} \approx \ln \frac{|\Sigma_0|}{|\Sigma_1|}$

$- \left( \Sigma_0^{-1} \mu_0 - \Sigma_1^{-1} \mu_1 \right)^T \left( \Sigma_1^{-1} - \Sigma_0^{-1} \right)^{-1} \left( \Sigma_0^{-1} \mu_0 - \Sigma_1^{-1} \mu_1 \right) \geq 0$

or ....

See Assignment Week 3 for comprehensive answer.