[PS 2]

#1.

(a) Training of data_B doesn't converge, whereas that of data_A does.

(b) Because data_B is linearly seperable, so scale of $\theta^T x$ gets larger to decrease loss until $|\theta|$ reaches $\infty$

(c) iii) $\to$ $\theta$ won't go $\infty$.

(e) { v) $\to$ this prevents data_B from being linearly seperable.

(d) No, because it maximises geometric margin, with $|\theta|$ fixed.

#2.

(a) Learned $\theta$ of logistic regression model meets:
$$\frac{\partial J(\theta)}{\partial \theta} = \sum (y^i - h_\theta(x^i)) x^i = 0$$

$\Rightarrow$ matrix form :

$$\begin{bmatrix} | & & | \\ \vdots & \cdots & \vdots \\ x_n^{(1)} & & x_n^{(m)} \end{bmatrix} (Y - h(X)) = 0$$
$\underbrace{\phantom{XXXXXXXXX}}_{\substack{X \\ (n+1 \times m)}}$  $(m \times 1)$

consider $1^{st}$ row of eq.

$$\sum (y^i - h_\theta(x^i)) = 0$$                              $\{ y^i \in (0,1) \}$

$\Rightarrow \sum P(y^i = 1 | x^i; \theta) = \sum h_\theta(x^i) = \sum y^i = \sum 1\{y^i = 1\}$

this is equivalent to question's condition $\left( (a,b) = (0,1) \right)$

(b) Perfect calibration doesn't necessarily imply perfect accuracy because calibration is about probability. Converse is necessarily true.

(c) $X(Y - h(X)) + 2\lambda\theta = 0$

Consider $1^{st}$ row of eq.

$\sum (y^i - h_\theta(x^i)) + 2\lambda\theta_o = 0 \Rightarrow \sum P(y^i = 1 | x^i; \theta) = \sum 1\{y^i = 1\}$

$\therefore$ prediction will be biased, proportional to $\theta_o$                      $+ 2\lambda\theta_o$

**#3.** Prove by contradiction.

assume $|\theta_{MAP}|_2 > |\theta_{ML}|_2$

then $P(\theta_{MAP}) < P(\theta_{ML})$    $\because \theta \sim N(0, \tau^2 I)$

$\Rightarrow P(\theta_{MAP}) \prod P(y^i | x^i; \theta)_{MAP} < P(\theta_{MA}) \prod P(y^i | x^i; \theta_{MAP})$

$(\because \theta_{ML} \text{ maximises } \prod P(y^i | x^i; \theta)) \to < P(\theta_{ML}) \prod P(y^i | x^i; \theta_{ML})$

$\Rightarrow$ contradicts to the fact that $\theta_{MAP}$ maximises $P(\theta) \prod P(y^i | x^i; \theta)$ "

**#4.**

(a) $u^T K u = \sum_{i,j} z_i z_j (K_1(x^i, x^j) + K_2(x^i, x^j))$

$\qquad = u^T K_1 u + u^T K_2 u \geq 0 \quad \leadsto$ Kernel "

(b) not Kernel

(c) $u^T K u = a\, u^T K_1 u \qquad \therefore$ Kernel

(d) not kernel

(e) Since $K_1(x, z), K_2(x, z)$ is Kernel, $\exists \phi_1 \in \mathbb{R}^{d_1}, \exists \phi_2 \in \mathbb{R}^{d_2}$

s.t $K_1(x, z) = \phi_1(x)^T \phi_1(z)$, $K_2(x, z) = \phi_2(x)^T \phi_2(x)$

Then $K(x, z)$

$\qquad = \sum_i \phi_1(x)_i \phi_1(z)_i \sum_j \phi_2(x)_j \phi_2(z)_j$

$\qquad = \sum_i \sum_j (\phi_1(x)_i \phi_2(x)_j)(\phi_1(z)_i \phi_2(z)_j) = \text{flat}(\phi_1(x)\phi_2(x)^T)^T \text{flat}(\phi_1(z)\phi_2(z)^T)$

$\qquad\qquad\qquad\qquad\qquad \underbrace{\qquad\qquad}_{\phi_3(x)} \qquad \underbrace{(\phi_1(z)\phi_2(z)^T)}_{\phi_3(z)}$

, $\phi_3 \in \mathbb{R}^{d_1 \cdot d_2} \Rightarrow K$ is Kernel "

(f) $u^T K u = \sum_{i,j} u_i u_j f(x^i) f(x^j) = (\sum u_i f(x^i))^2 \geq 0 \Rightarrow$ Kernel

(g) $u^T K u = \sum u_i u_j K_3(\phi(x^i), \phi(x^j)) \geq 0 \qquad \Rightarrow$ Kernel

(h) $u^T K u = \sum u_i u_j (\sum \alpha_k (K_1(x^i; x^j))^k) \geq 0 \quad (\because (a), (c), (e))$

**#5.** Since $\theta$ is too high dim, $\theta^T \phi(x)$ is represented as Kernel $K(\theta, \phi(x))$ always. (a) To predict, $h_{\theta^i}(x^{i+1}) = g(K(\theta^i, \phi(x^{i+1})))$ (b) And param update rule is modified to

$\qquad K(\theta^{i+1}, \phi(x^{i+1})) := K(\theta^i, \phi(x^{i+1})) + \alpha 1\{K(\theta^i, \phi(x^{i+1})) y^{i+1} < 0\} K(\theta^i, \phi(x^{i+1})) y^{i+1}$

(c). Only dot product in kernel form need: $O(n^2) \to O(n)$