

PS 2

Q1) (a) Code converges on data A but not on B

(b) B is perfectly linearly separable unlike A. The optimization is unbounded & $\theta \rightarrow \infty$. Regularization eg L2 can put a limit on θ . SVM might be more appropriate for linearly separable.

- (c) i. No - as $\theta \rightarrow \infty$ regardless of size of lr
ii. Depends - if $lr \rightarrow 0$ in finite time it will converge. If time too short, underfit
iii. No - as only changes values of θ as already linearly separable
iv. Yes - limits $\| \cdot \|_{L_2}$ of θ to be with a bound so training stop when bound achieved
v. Depends - if data is very large, ~~noise~~ noise cancels & problem persists. Here it may help in some instances due to $\sum \text{noise} \neq 0$ but depends if boundary is still linearly separable

(d) No - SVM not vulnerable

it maximizes the distance / margin b/w the two classes (distance of boundary from nearest pt)
it is implicitly regularizing θ as $\max_{\theta} \frac{1}{2} \|\theta\|^2$

Q2 @ 172 2

⑤ 0.978

⑥ class, non, prize, tone, night!

⑦ 0.1

Q. 9695 accuracy

Q3 @ Yes: $k_1 + k_2 = K = K^T$
$$z^T K z = z^T K_1 z + z^T K_2 z \geq 0$$

PSD

① Yes: $z^T a K z = a(z^T K z) \geq 0 \because a \in \mathbb{R}^+$

② No: by c, let $K_2 = 2K_1$ is kernel
$$z^T (K_1 - K_2) z = -z^T K_1 z \leq 0 \text{ for certain } z$$

 \therefore not PSD

(d) No: by b & c, $\nexists a \in \mathbb{R}^+$

③ Yes: $f(x) \in \mathbb{R} \therefore \langle f(x), f(z) \rangle = f(x)f(z)$
as scalar as ~~R~~ multiplication

(g) Yes: basic feature map
 $K_3 \text{ } p \times p$
 $\Phi: d \mapsto p \therefore K_3(\in \mathbb{R}^p, \mathbb{R}^p)$
 \therefore kernel

(h) Yes: By a, c, e, f
$$g_p(x) = ax^3 + cx^2 + d + x$$

Q2

$$\begin{aligned}
 k(x, z) &= k_1(x, z) k_2(x, z) \\
 &= \sum_i \phi_1^i(x) \phi_1^i(z) \sum_j \phi_2^j(x) \phi_2^j(z) \\
 &= \sum_i \sum_j \phi_1^i(x) \phi_1^i(z) \phi_2^j(x) \phi_2^j(z) \\
 &\quad \phi_1^i \phi_2^j \text{ is } \phi_3^{ij} \therefore \sum_{ij} \theta(x) \theta(z) \\
 &\quad \therefore \text{kernel}
 \end{aligned}$$

Q3 (i) $\theta^{(i)} = \sum_{j=1}^i \beta_j \phi(x^{(j)})$ $\theta^{(0)} = 0$ when $\beta_j = 0$
 β is y essentially at $t=0$

(ii) $h_{\theta}^{(i)}(x^{(i+1)}) = g(\theta^{(i)T} \phi(x^{(i+1)}))$ $g(z) = \text{sign}(z)$

$$\begin{aligned}
 &= g\left(\left[\sum_{j=1}^i \beta_j \phi(x^{(j)})\right]^T \phi(x^{(i+1)})\right) \\
 &= g\left(\sum_{j=1}^i \beta_j k(x^{(j)}, x^{(i+1)})\right)
 \end{aligned}$$

(iii) $\theta^{(i+1)} = \sum_{j=1}^{i+1} \beta_j \phi(x^{(j)}) = \underbrace{\sum_{j=1}^i \beta_j \phi(x^{(j)})}_{\theta^{(i)}} + \beta_{i+1} \phi(x^{(i+1)})$

$$\beta_{i+1} = \alpha (y^{(i+1)} - g(\theta^{(i)T} \phi(x^{(i+1)})))$$

$$\therefore \theta^{(i+1)} = \theta^{(i)} + \alpha (y^{(i+1)} - g(\theta^{(i)T} \phi(x^{(i+1)}))) \phi(x^{(i+1)})$$

Linear kernel (a.b) poor fit as data not linear or not approx linear

EBF can build circles to separate data

$$\begin{aligned}
 6. \quad \textcircled{a} \quad \theta_{\text{MAP}} &= \arg \max_{\theta} p(\theta | x, y) \\
 &= \arg \max_{\theta} \frac{p(y | x, \theta) p(\theta | x)}{p(y | x)} \\
 &= \arg \max_{\theta} p(y | x, \theta) \underbrace{p(\theta | x)}_{\text{constant } \theta} \\
 &= \arg \max_{\theta} p(y | x, \theta) p(\theta)
 \end{aligned}$$

$$\textcircled{b} \quad \text{Take } -\log \\
 \therefore \theta_{\text{MAP}} = \arg \min_{\theta} -\log p(y | x, \theta) - \log p(\theta)$$

$$\therefore \theta \sim N(\theta_0, \eta^2 I)$$

$$\therefore -\log p(\theta) = -\log(\text{const}) + \frac{1}{2} \theta^T \left(\frac{1}{\eta^2} I \right) \theta$$

$$\therefore \arg \min_{\theta} -\log p(y | x, \theta) + \underbrace{\frac{1}{2\eta^2} \|\theta\|_2^2}$$

$$\textcircled{c} \quad \therefore -\log p(y | x, \theta) = \frac{1}{2\sigma^2} \|\bar{y} - X\theta\|_2^2 \quad \text{ignore const}$$

$$\therefore J(\theta) = \frac{1}{2\sigma^2} \|\bar{y} - X\theta\|_2^2 + \frac{1}{2\eta^2} \|\theta\|_2^2$$

$$\nabla_{\theta} J(\theta) = 0 = \frac{1}{\sigma^2} X^T (X\theta - \bar{y}) + \frac{1}{\eta^2} \theta$$

$$\therefore \hat{\theta}_{\text{MAP}} = \left(X^T X + \frac{\sigma^2}{\eta^2} I \right)^{-1} X^T y$$

$$\textcircled{d} \quad d_L(\theta, \bar{0}, b) = \frac{1}{2b} e^{-\frac{\theta}{b}}$$

$$\therefore -\log f = -\text{const} + \frac{10\theta}{b}$$

using ϵ , ignore const

$$J(\theta) = \frac{1}{2\sigma^2} \|x\theta - y\|_2^2 + \frac{10\theta}{b}$$

$$\hat{\theta}_{\min} \arg \min J(\theta)$$

rescale to $J(\theta)$ in form $\|x\theta - y\|_2^2 + \frac{2\sigma^2}{b} 10\theta$

$$\gamma = \frac{2\sigma^2}{b}$$

$$\textcircled{Q7} \quad l(\theta) = \sum_{i=1}^m y^{(i)} \log h(x^{(i)}) + (1-y^{(i)}) \log (1-h(x^{(i)}))$$

$$\frac{\partial l(\theta)}{\partial \theta_j} = \sum_{i=1}^m (y^{(i)} - h(x^{(i)})) x_j^{(i)} = 0 \quad \text{--- } \textcircled{1}$$

$$\sum_{i=1}^m P(y=1|x^{(i)}; \theta) = \therefore \sum_{i=1}^m h(x^{(i)}) = \sum_i y = \sum_{i=1}^m \mathbb{I}\{y^{(i)}=1\} \quad \text{--- } \textcircled{2}$$

change to $(0,1) = (a,b)$ as $i=1$ to m

$$\therefore |\{i \in I_{a,b}\}| = m$$

$$\therefore \sum_{i \in I_{a,b}} \frac{P}{h(x^{(i)})} = \sum_{i \in I_{a,b}} \frac{\mathbb{I}\{y^{(i)}=1\}}{h(x^{(i)})}$$

\textcircled{B} No, no for converse

- model can misclassify based on threshold set so even though prob perfect, not 100% accuracy
- model can get correct accuracy with non-optimal probabilities, thus is not correct due to threshold. $P(y=1|x)=0.9$ when

\textcircled{C} may improve by reducing overfitting & overconfidence. too strong L_2 can underfit. truth is 0.8

$$\nabla J(\theta) = \textcircled{1} + \lambda \theta_2 \quad \therefore \textcircled{2} h + \lambda \theta \quad \therefore \text{not calibrated well}$$