

(a) dataset A : converges

dataset B : not converges

(b) dataset B is linear separable,

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y^{(i)} \theta^T x^{(i)}))$$
$$\nabla_{\theta} J(\theta) = -\frac{1}{m} \sum_{i=1}^m \frac{\exp(-y^{(i)} \theta^T x^{(i)})}{1 + \exp(-y^{(i)} \theta^T x^{(i)})}$$
$$= -\frac{1}{m} \sum_{i=1}^m \frac{y^{(i)} \theta^T x^{(i)}}{\exp(y^{(i)} \theta^T x^{(i)}) + 1}$$

$y^{(i)} \theta^T x^{(i)} > 0$  if the  $\theta$  is already decide the dataset B into 2 part.

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} J(\theta)$$

which makes  $J(\theta)$  decreases

$$\underbrace{\hspace{10em}}_{\text{logistic}} \rightarrow$$

it will converges slowly  
for  $y^{(i)} x^{(i)T} \theta > 0$

(c) (i) no, no help to solve the problem

(ii) yes, it can make  $\alpha \nabla_{\theta} J(\theta)$  decreases quickly

(iii) no, it just scalar  $\theta$

(iv) yes, it will make  $\|\theta\|_2^2$  smaller quickly

(v) yes, it will make data not linear separable

(d) no,  $[1-z]_+$  hinge function will be 0 as  $z > 1$  so it will be minimized to 0 if the dataset is linear separable

# use the  $\ell_2$  regularization

$$\frac{\partial}{\partial \theta_i} \|\theta\|_2^2 = \frac{\partial (\theta_1^2 + \theta_2^2 + \dots + \theta_n^2)}{\partial \theta_i} = 2\theta_i$$

$$\nabla_{\theta} \|\theta\|_2^2 = 2\theta$$

2.

(a)

$$\frac{\partial L(\theta)}{\partial \theta_j} = 0$$

$$\sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)} = 0$$

when  $j = 0$ 

$$\sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) = 0$$

$$\sum_{i=1}^m h_\theta(x^{(i)}) = \sum_{i=1}^m y^{(i)}$$

$$\sum_{i \in S_{ab}} P(y^{(i)} = 1 | x^{(i)}, \theta) = \sum_{i \in S_{ab}} I[y^{(i)} = 1]$$

(b) NO

for example, when  $a=0, b=0.5$ all prediction is 0, but  $P(y^{(i)} = 1 | x^{(i)}, \theta)$   
is not 0, the converse is also false

$$(c) \quad \frac{\partial L(\theta)}{\partial \theta_i} = \frac{\partial L(\theta)}{\partial \theta_j} + c\theta_i = 0$$

$$\frac{\partial L(\theta)}{\partial \theta_i} = -c\theta_i$$

so it will make it not  
perfectly calibrated

3. (a)

$$P(\theta | x, y) = \frac{P(\theta, x, y)}{P(x, y)} = \frac{P(y|x, \theta) P(x, \theta)}{P(x, y)}$$

$$= \frac{P(y|x, \theta) P(x) P(\theta)}{P(x, y)}$$

so  $\arg \max_{\theta} P(\theta | x, y) =$ 

$$\arg \max_{\theta} P(y|x, \theta) P(\theta)$$

$$(b) \quad \theta_{\text{MAP}} = \arg \max_{\theta} P(y|x, \theta) P(\theta)$$

$$= \arg \min_{\theta} -\log(P(y|x, \theta) P(\theta))$$

$$= \arg \min_{\theta} (-\log P(y|x, \theta) - \log(\frac{1}{2\pi})^{\frac{1}{2}} \eta^n)$$

$$\exp\left\{-\frac{1}{2\eta^2} \|\theta\|_2^2\right\}$$

$$= \arg \min_{\theta} \left(-\log P(y|x, \theta) + \frac{1}{2\eta^2} \|\theta\|_2^2\right)$$

$$\text{so for } \lambda = \frac{1}{2\eta^2}$$

it equivalent to  $L_2$  regularization

(c)

$$P(y^{(i)} | x^{(i)}, \theta) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(y^{(i)} - \theta^\top x^{(i)})^2}{2\sigma^2}\right)$$

$$\arg \min_{\theta} \left(-\log P(y|x, \theta) + \frac{1}{2\eta^2} \|\theta\|_2^2\right)$$

$$= \arg \min_{\theta} \left(\frac{1}{2\sigma^2} \|(y - X\theta)\|_2^2 + \frac{1}{2\eta^2} \|\theta\|_2^2\right)$$

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

$$\begin{aligned} D\theta J(\theta) &= \frac{-1}{\sigma^2} X^T (\vec{y} - X\theta) + \frac{1}{n} \theta \\ &= \frac{1}{\sigma^2} (X^T X \theta - X^T \vec{y}) + \frac{1}{n} \theta = 0 \\ \left( \frac{1}{\sigma^2} X^T X + \frac{1}{n} I \right) \theta &= \frac{1}{\sigma^2} X^T \vec{y} \\ \theta &= \left( \frac{1}{\sigma^2} X^T X + \frac{1}{n} I \right)^{-1} \frac{1}{\sigma^2} X^T \vec{y} \end{aligned}$$

(d)  $\log P(\theta)$

$$= -\log 2b - \frac{|z-\mu|}{b} \quad |_{\mu=0}$$

$$= -\log 2b - \frac{|z|}{b}$$

$$\theta = \underset{\text{MAP}}{\arg \min}_{\theta} \left( \frac{1}{2\sigma^2} \| \vec{y} - X\theta \|_2^2 + \frac{1}{b} \| \theta \| \right)$$

$$= \arg \min_{\theta} \left( \| \vec{y} - X\theta \|_2^2 + \frac{25^2}{b} \| \theta \|_1 \right)$$

$$SDJ(\theta) = \| X\theta - \vec{y} \|_2 + \gamma \| \theta \|_1.$$

$$\text{where } \gamma = \frac{25^2}{b}$$

q.

$$(a) K_{ij} = K(x^{(i)}, x^{(j)}) = k_1(x^{(i)}, x^{(j)}) + k_2(x^{(i)}, x^{(j)})$$

$$= k_{1,ij} + k_{2,ij}$$

$$\text{so } K = k_1 + k_2$$

$$z^T k z = z^T k_1 z + z^T k_2 z \geq 0$$

so it is a kernel

(b) it isn't a kernel

$$\text{for } k_2 = 2k_1$$

$$K = -k_1$$

$$z^T k z = -z^T k_1 z \leq 0$$

$$(c) k = \alpha k_1$$

$$z^T k z = \alpha z^T k_1 z \geq 0$$

so it's a kernel

(d) it isn't a kernel

$$(e) K_{ij} = k_{1,ij} k_{2,ij}$$

$$z^T k z = \sum_{i=1}^m \sum_{j=1}^m z_i k_{ij} z_j$$

$$= \sum_{i=1}^m \sum_{j=1}^m z_i k_{1,ij} k_{2,ij} z_j$$

$$= \sum_{i=1}^m \sum_{\alpha=1}^n \sum_{\beta=1}^m z_i \phi_{1,\alpha}(x_i) \phi_{1,\beta}(x_i) \phi_{2,\alpha}^T(x_i) \phi_{2,\beta}^T(x_i) z_j$$

$$= \sum_{\alpha} \sum_{\beta} \sum_i z_i \phi_{1,\alpha}(x_i) \phi_{1,\beta}(x_i) \phi_{2,\alpha}^T(x_i) \phi_{2,\beta}^T(x_i) z_j$$

$$= \sum_{\alpha} \sum_{\beta} \sum_i (z_i \phi_{1,\alpha}(x_i) \phi_{1,\beta}(x_i))^2 \sum_j (z_j \phi_{2,\alpha}(x_j) \phi_{2,\beta}^T(x_j))$$

$$= \sum_{\alpha} \sum_{\beta} \left( \sum_i (z_i \phi_{1,\alpha}(x_i) \phi_{1,\beta}(x_i))^2 \right) \geq 0$$

so it is a kernel

$$(f) K_{ij} = f(x^{(i)})^T f(x^{(j)})$$

$$\text{so } K = F(X) F^T(X)$$

$$\text{where } F: \mathbb{R}^n \rightarrow \mathbb{R}^n$$

$$\text{and } F\left[\begin{matrix} x_1 \\ \vdots \\ x_n \end{matrix}\right] = \left[\begin{matrix} f(x_1) \\ \vdots \\ f(x_n) \end{matrix}\right]$$

$$Z^T K Z = Z^T F(X) F^T(X) Z$$

$$= \|F(X)^T Z\|_2^2 > 0$$

so it's a kernel

$$(g) K_{ij} = k_3(\phi(x^{(i)}), \phi(x^{(j)}))$$

for any set  $\{y^{(1)}, \dots, y^{(m)}\}$

we can find  $\{x^{(1)}, \dots, x^{(m)}\}$   
where  $\phi(x^{(i)}) = y^{(i)}$

so it's a kernel

$$(h) k(x, z) = \sum_{b=0}^n C_b (k_b(x, z))^b$$

$$K_{ij} = \sum_{b=0}^n C_b (k_b(x^{(i)}, x^{(j)}))^b$$

according to (e),  $K(x, z)$  is valid  
(a), (c) and

5. (a)

$$\begin{aligned} (i) h^{(i+1)} &= \theta + \alpha(y^{(i+1)} - g(\theta + \phi(x^{(i)}))) \phi(x^{(i)}) \\ &= \sum_{j=1}^i \beta_j \phi(x^{(j)}) + \beta_{i+1} \phi(x^{(i+1)}) \\ &= \sum_{j=1}^i \beta_j \phi(x^{(j)}) \\ h^{(0)} &= \sum_{j=1}^0 \beta_j \phi(x^{(j)}) = 0 \end{aligned}$$

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

(iii)

$$\beta_{i+1} = \alpha(y^{(i+1)} - g\left(\sum_{j=1}^i \beta_j k(x^{(j)}, x^{(i+1)})\right))$$

(c) dot product kernel is bad, because  
this data set is not linearly separable

6.

$$(b) \log \frac{p(y=1) p(x|y=1)}{p(y=0) p(x^{(i)}|y=1)}$$

$$\begin{aligned} &= (\log p(y=1) + \sum_{j=1}^d \log p(x_j | y=1)) \\ &\quad - (\log p(y=0) + \sum_{j=1}^d \log p(x_j | y=0)) \\ &= \log \frac{p_y}{1-p_y} + \sum_{j=1}^d \log \frac{p_j|y=1}{p_j|y=0} \end{aligned}$$

using this method can easily predict