

PS 2

(a) algorithm converges on set A but doesn't on set B.

(b) math,

$$L(\theta) = \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta)$$

$$= \prod_{i=1}^m \left( \frac{1}{1 + e^{-y^{(i)} \times \theta^T x^{(i)}}} \right)$$

$$\log L(\theta) = \sum_{i=1}^m \log (1 + e^{-y^{(i)} \times \theta^T x^{(i)}}) = J(\theta)$$

$$\nabla_{\theta} (\log L(\theta)) = \sum_{i=1}^m \left( \frac{e^{-y^{(i)} \theta^T x^{(i)}}}{1 + e^{-y^{(i)} \theta^T x^{(i)}}} \times y^{(i)} \right) x^{(i)}$$

$$= -X^T (1 - \text{prob}_{\theta}) * Y$$

$$\nabla_{\theta} J(\theta) = \sum_{i=1}^m \left( \frac{1}{1 + e^{-y^{(i)} \theta^T x^{(i)}}} \times e^{-y^{(i)} \theta^T x^{(i)}} \times -y^{(i)} \right) x^{(i)}$$

$$= -X^T (1 - \text{prob}_{\theta}) Y$$

The formula in the code somehow works which I do not understand, but my derived equation also works well.

Note that B is linearly separable while A is not; also our loss function wants to maximize  $y^{(i)} \theta^T x^{(i)}$ , which is nothing but the functional margin which can be scaled to arbitrary values by scaling  $\theta$ , we cannot do so for A. This is why it fails to converge for dataset B.

note that if we make the learning rate very large then A fails to converge because it overshoots (the minimum) but B reaches the limit of diff.  $< 10^{-15}$  very fast. if we decrease this limit then we might get overflows but theoretically we can always keep changing theta for the dataset B acc. to our cond<sup>n</sup>.

Suppose we try to capture GM minimization,

$$J(\theta) = \sum_{i=1}^m \log \left( 1 + e^{-y^{(i)} \theta^T x^{(i)}} \right)$$

$$\nabla_{\theta} J(\theta) = \sum_{i=1}^m \left( \frac{1}{1 + e^{-y^{(i)} \theta^T x^{(i)}}} \times e^{-y^{(i)} \theta^T x^{(i)}} \times y^{(i)} \right) \frac{\left( m^{(i)} (\theta^T \theta) - 2 \theta (\theta^T x^{(i)}) \right)}{(\theta^T \theta)^2}$$

$$= \left\{ \frac{X^T (\theta^T \theta) - 2 \theta (x \theta)^T}{(\theta^T \theta)^2} \right\} (1-\text{prob}) Y$$

with this code however, the value with which  $\theta$  is initialized becomes crucial, I have observed this.

(C) (i) No, it will keep learning with arbitrarily large learning rates, although this will cause an overflow.

(C) (ii) Yes  $\times D$ , it will converge with a large  $\lambda$  ( $\sim 10^{10}$ )  
 $\because$  the loss will become 0 more quickly.

(C) (iii) Yes & No, if our convergence cond<sup>n</sup> is just the difference being  $< 10^{-15}$  then yes;  $\theta$  learning rate itself becomes  $< 10^{-15}$  but actually the fitted lines are not very accurate.

(iii) No, the features are already in a proper scale.

(iv) Yes, definitely, with regularization we add the required scaling constraint. (for me 0.001 reg. worked best).

(v) Yes, ∵ this will result in a linearly non-separable data.

(d) Nope, SVMs are robust against such datasets ∵ they minimize GMM with the constraint  $\sum \alpha_i = 1 \Rightarrow$  there is a scaling constraint on  $\theta$  & we would not be able to take it to arbitrarily high values indefinitely.

$$\begin{aligned}
 \text{(d) } \text{(a)} \quad & \sum_{i \in I_{a,b}} P(y^{(i)}=1 | m^{(i)}; \theta) = \sum_{i=1}^m P(y^{(i)}=1 | m^{(i)}; \theta) \\
 & \because (\alpha, b) = (0, 1) \\
 & \text{all eg.} \\
 & = \sum_{i=1}^m \left( \frac{1}{1 + \exp(-\theta^T m^{(i)})} \right) \\
 & = \sum_{i=1}^m \left( \frac{1}{1 + \exp(\theta_0^T m^{(i)} - \theta_0)} \right) \\
 & = \sum_{i=1}^m \left( \frac{e^{\theta_0}}{e^{\theta_0} + e^{-\theta_0^T m^{(i)}}} \right)
 \end{aligned}$$

$$\text{(d) } \text{(b)} \quad l(\theta) = \log L(\theta) = \sum_{i=1}^m y^{(i)} \log h(m^{(i)}) + (1-y^{(i)}) \log (1-h(m^{(i)}))$$

$$\frac{\partial l(\theta)}{\partial \theta_j} = - \sum_{i=1}^m (h(m^{(i)}) - y^{(i)}) m_j^{(i)} = 0$$

at  $j=0$ ,  $m_j^{(i)} = 1 \nabla j \in I_{a,b}$  ( $\leq$  full set)

$$\Rightarrow \sum_{i=1}^m h(m^{(i)}) = \sum_{i=1}^m y^{(i)} \neq$$

clearly,  $h(m^{(i)}) = P(y^{(i)}=1 | M^{(i)}; \theta)$

$$y^{(i)} = I\{y^{(i)}=1\}$$

& everything follows.

(b) No, suppose we have eg. in set  $(a, b) = (0.25, 0.75)$

$$\text{s.t. } h(m^{(1)}) = 0.66 \quad \& \quad y^{(1)} = 0$$

$$h(m^{(2)}) = 0.46 \quad \& \quad y^{(2)} = 1$$

$$h(m^{(3)}) = 0.66 \quad \& \quad y^{(3)} = 1$$

$$\Rightarrow \sum_{i \in I_{a,b}} h(m^{(i)}) = 0.66 \times 3$$

at all threshold of 0.5 our model doesn't get 100% accuracy but (a) holds!

converse is also not true, eg.-

$$h(m^{(1)}) = 0.51 \quad \& \quad y^{(1)} = 1$$

$$h(m^{(2)}) = 0.51 \quad \& \quad y^{(2)} = 1$$

$$h(m^{(3)}) = 0.51 \quad \& \quad y^{(3)} = 1$$

at 0.5 threshold, 100% accuracy  
but model isn't perfectly calibrated for  
range  $(0.5, 0.52)$ .

(c) with  $L_2$  regularization, model will not be well calibrated  $\because$ , our earlier proof fails to hold.

$$\text{i.e., it become } \sum_{i=1}^m h(m^{(i)}) + \lambda \theta_0 = \sum_{i=1}^m y^{(i)}$$

arg. term

PS2 contd.

$$\underline{Q3} \quad \theta_{MAP} = \underset{\theta}{\operatorname{argmax}} \ p(\theta | m, y)$$

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} \ p(y | m; \theta)$$

$$(a) \ p(\theta | m, y) \cdot p(y | m, \theta) = p(m, y | \theta) \frac{p(\theta)}{p(m, y)}$$

$$\in p(y | m, \theta) p(m | \theta) p(\theta)$$

$$(b) \ p(\theta | m, y) = \frac{p(m, y | \theta) p(\theta)}{p(m, y)}$$

$$\text{Also, } p(m, y | \theta) = p(y)$$

$$(a) \ p(\theta | m, y) = \frac{p(\theta, m, y)}{p(m, y)}$$

$$= \frac{p(y | \theta, m) p(\theta | m) p(m)}{p(m, y)}$$

$$= \frac{p(y | \theta, m) p(\theta | m)}{p(m, y)}$$

$$= \frac{p(y | \theta, m) p(\theta | m) p(m)}{p(y | m) p(m)}$$

$$p(\theta | y, m) = \frac{p(y | \theta, m) p(\theta | m)}{p(y | m)}$$

p(y | m)  
constant

$$\Rightarrow \theta_{MAP} = \underset{\theta}{\operatorname{argmax}} \ p(y | \theta, m) p(\theta)$$

$$(b) \theta \sim \mathcal{N}(0, \eta^2 I)$$

$$p(\theta) = \frac{1}{(2\pi)^{n/2} \sqrt{|\eta^2 I|}} \times \exp\left(-\frac{(\theta - 0)^T (\eta^2 I)^{-1} (\theta - 0)}{2}\right)$$

$$p(\theta) = \frac{1}{(2\pi)^{n/2} \times \eta^2} \times \exp\left(-\frac{\theta^T \theta}{2\eta^2}\right)$$

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} p(y|m, \theta) p(\theta)$$

$$\text{now, } p(y|m, \theta) p(\theta) = p(y|m, \theta) \times \frac{1}{(2\pi)^{n/2} \eta^2} e^{-\frac{\theta^T \theta}{2\eta^2}}$$

$$l(\theta) = -\log \{ p(y|m, \theta) p(\theta) \}$$

$$= -\log p(y|m, \theta) + \frac{n}{2} \log 2\pi + 2 \log \eta$$

$$+ \frac{1}{2\eta^2} \theta^T \theta$$

$$\theta^T \theta = \|\theta\|_2^2$$

$$\therefore \theta_{MAP} = \underset{\theta}{\operatorname{argmin}} \left( -\log p(y|m, \theta) + \frac{1}{2\eta^2} \|\theta\|_2^2 \right)$$

$$\lambda = \frac{1}{2\eta^2}$$

$$(c) \quad y = \theta^T m + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$\theta \sim \mathcal{N}(0, \eta^2 I)$$

$$\therefore \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$\Rightarrow (\vec{y} - \theta^T X) \sim \mathcal{N}(0, \sigma^2 I)$$

$$p(\vec{y} | \theta, X) = \frac{1}{(2\pi)^{n/2} \sigma^2} \exp \left( -\frac{\vec{y}^T - (\theta^T X - \theta)(\sigma^2 I)^{-1}(\theta^T X - \theta)}{2\sigma^2} \right)$$

$$= \frac{1}{(2\pi)^{n/2} \sigma^2} \exp \left( -\frac{(X^T \theta + \vec{y}^T - \frac{2}{\sigma^2} \vec{y}^T)(\theta^T X)}{2\sigma^2} \right)$$

$$p(\theta) = \frac{1}{(2\pi)^{n/2} \eta^2} \exp \left( -\frac{\theta^T \theta}{2\eta^2} \right)$$

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} p(\vec{y} | X, \theta) p(\theta)$$

as in prev. q,

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} \left( -\log p(\vec{y} | X, \theta) + \frac{n}{2} \log 2\pi + 2 \log \eta + \frac{1}{2} \theta^T \theta \right)$$

$$\begin{aligned} &= \underset{\theta}{\operatorname{argmax}} \left( \frac{n}{2} \log 2\pi + 2 \log \sigma + \frac{1}{2} \frac{(\theta^T X)^T}{\sigma^2} \right) \\ &= \underset{\theta}{\operatorname{argmin}} \left( \frac{n}{2} \log 2\pi + 2 \log \sigma + \frac{1}{2} \frac{(\vec{y} - \theta^T X)^T (\vec{y} - \theta^T X)}{\sigma^2} + 2 \log \eta + \frac{1}{2} \frac{\theta^T \theta}{\sigma^2} \right) \end{aligned}$$

$$\therefore \theta_{MAP} = \arg \min_{\theta} \left( \frac{(\bar{y} - \theta^T x)^T (\bar{y} - \theta^T x)}{2\sigma^2} + \frac{1}{2\eta^2} \theta^T \theta \right)$$

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

$$J(\theta) = \frac{1}{2\sigma^2} (y - x\theta)^T (y - x\theta) + \frac{1}{2\eta^2} \theta^T \theta$$

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \frac{1}{2\sigma^2} \nabla_{\theta} \{ (y^T - \theta^T x^T)(y - x\theta) \} + \frac{1}{2\eta^2} \theta \\ &= \frac{1}{2\sigma^2} \nabla_{\theta} (y^T y - y^T x \theta - \theta^T x^T y + \theta^T x^T x \theta) + \frac{1}{2\eta^2} \theta \\ &= \frac{1}{2\sigma^2} (x^T y - (x^T x) \theta + (x^T x) \theta) + \frac{1}{2\eta^2} \theta \\ &= 0 \end{aligned}$$

$$\Rightarrow (x^T x)^{-1} x^T y + \frac{\sigma^2}{\eta^2} \theta = 0$$

$$x^T x \theta - x^T y + \frac{\sigma^2}{\eta^2} \theta = 0$$

$$\left( x^T x + \frac{\sigma^2}{\eta^2} I \right) \theta = x^T y$$

$$\boxed{\theta = \left( x^T x + \frac{\sigma^2}{\eta^2} I \right)^{-1} x^T y}$$

(d)  $f_{\theta}(z|\mu, b) = \frac{1}{2b} \exp\left(-\frac{|z-\mu|}{b}\right)$

$$\theta \sim \mathcal{L}(0, bI)$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$y = \theta^T x + \epsilon$$

similar to prev,

$$P(\theta) = \frac{1}{2b} \exp\left(-\frac{|\theta|}{b}\right)$$

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmin}} \left( \frac{1}{2\sigma^2} \|(\mathbf{x}\theta - \bar{y})\|_2^2 + \log 2b + \frac{1}{b} \right)$$

$$\therefore J(\theta) = \|(\mathbf{x}\theta - \bar{y})\|_2^2 + \frac{2\sigma^2}{b} \|\theta\|_1$$

~~$\theta \sim \mathcal{N}(0, bI)$~~ 

$$\sigma^2/b \quad \tau = \frac{2\sigma^2}{b}$$

※※ L<sub>2</sub> regression  $\Rightarrow$  ridge regression  
 L<sub>1</sub>, "  $\Rightarrow$  Lasso "

also called weight decay OR shrinkage  $\because$  encourage parameters to stay close to their mean ( $= 0$ )  $\therefore$  resulting in shrinkage.

※※ L<sub>1</sub> reg. results in sparse parameters, i.e., most of the values are zero.

(d)  $\forall (a)$   $K(x, z) = K_1(m, z) + K_2(n, z)$

$$z^T K z = z^T (K_1 + K_2) z \\ = z^T K_1 z + z^T K_2 z \geq 0.$$

(b)  $z^T K z = z^T (K_1 - K_2) z$

let  $K_2 = \lambda K_1$ , ( $\lambda > 1$ )

$$= (1-\lambda) z^T K_1 z$$

$$\leq 0$$

= not PSD!

(c) is PSD :  $z^T K z = a (z^T K_1 z) \geq 0$

(d) is not PSD  $\because z^T K z = -a (z^T K_1 z) \leq 0$ .

(e)  $K = K_1 K_2$

$$K^T = K_2 K_1$$

$$z^T K z = \sum_i \sum_j z_i K(m^{(i)}, n^{(j)}) z_j$$

$$= \sum_i \sum_j z_i K_1(m^{(i)}, m^{(j)}) K_2(m^{(i)}, n^{(j)}) z_j$$

$$= \sum_i \sum_j z_i \left( \sum_k \phi_1(m^{(i)})_k \phi_1(m^{(j)})_k \right) \left( \sum_l \phi_2(m^{(i)})_l \phi_2(n^{(j)})_l \right)^T z_j$$

$$= \sum_i \sum_j \left( \sum_k \phi_1(m^{(i)})_k \phi_1(m^{(j)})_k \right) \left( \sum_l \phi_2(m^{(i)})_l \phi_2(n^{(j)})_l \right)^T z_j$$

$$= \sum_i \sum_j \sum_k \sum_l z_i \phi_1(m^{(i)})_k \phi_1(m^{(j)})_k \phi_2(m^{(i)})_l \phi_2(n^{(j)})_l z_j$$

$$z^T K z = \sum_j \sum_k \sum_l (z_j \phi_{1(m^{(j)})}^k \phi_{2(n^{(l)})}^l)^2$$

~~$$(z^T K z)^+ = z^T K^T z$$~~

~~$$\begin{aligned} \Rightarrow z^T K z &= z^T \left( \frac{1}{2} K + \frac{1}{2} K^T \right) z \\ &= \frac{1}{2} z^T (K_1 K_2 + K_2 K_1) z \end{aligned}$$~~

(f)  $K(m, z) = f(m) / f(z)$

$$\begin{aligned} z^T K z &= \sum_i \sum_j z_i k_{ij} z_j \\ &= \sum_i \sum_j z_i f(m^{(i)}) f(m^{(j)}) z_j \\ &= \sum_i (z_i f(m^{(i)}))^2 \quad \text{PSD} \end{aligned}$$

(g)  $K(m, z) = k_3 (\phi(m), \phi(z)) = \phi'(\phi(m))^T \phi'(\phi(z))$

let  $\phi'' = \phi'(\phi(\cdot))$  clearly  $K$  is PSD.

(h)  $K(m, z) = p(K_1(m, z)) = p(\phi(m)^T \phi(z))$

$$= \sum_{i=0}^n b_i (\phi(m)^T \phi(z))^i$$

(i)  $K(m, z) = p(K_1(m, z)) = \sum_i b_i (K_1(m, z))^i$

$\geq 0 \quad \hookrightarrow \text{PSD}$

$\sum_i K_i z \geq 0 \quad \text{PSD}$

$$(a) \theta^{(i+1)} := \theta^{(i)} + \alpha (y^{(i+1)} - h_{\theta^{(i)}}(m^{(i+1)})) m^{(i+1)}$$

$$(a)(i) \theta^{(i)} = \sum_{j=1}^i \beta_j \phi(x^{(j)}) *$$

$$\theta^{(0)} = \vec{0}$$

$$(b) (ii) h_{\theta^{(i)}}(m^{(i+1)}) = g \left( \sum_{j=1}^i \beta_j k(m^{(j)}, m^{(i+1)}) \right)$$

$$= g \left( \sum_{j=1}^i \beta_j \langle x^{(j)}, m^{(i+1)} \rangle \right)$$

$$(iii) \theta^{(i+1)} := \theta^{(i)} + \alpha (y^{(i+1)} - h_{\theta^{(i)}}(m^{(i+1)})) \phi(m^{(i+1)})$$

$$\therefore \beta_{i+1} = \alpha (y^{(i+1)} - h_{\theta^{(i)}}(m^{(i+1)}))$$

$$\& \theta^{(i+1)} = \sum_{j=1}^{i+1} \beta_j \phi(m^{(j)}).$$

(b) Code

(c) Code dot product kernel behaves poorly : it does not have the expressive power of higher order features & tries to classify points by making a linear decision boundary

Q6 Code