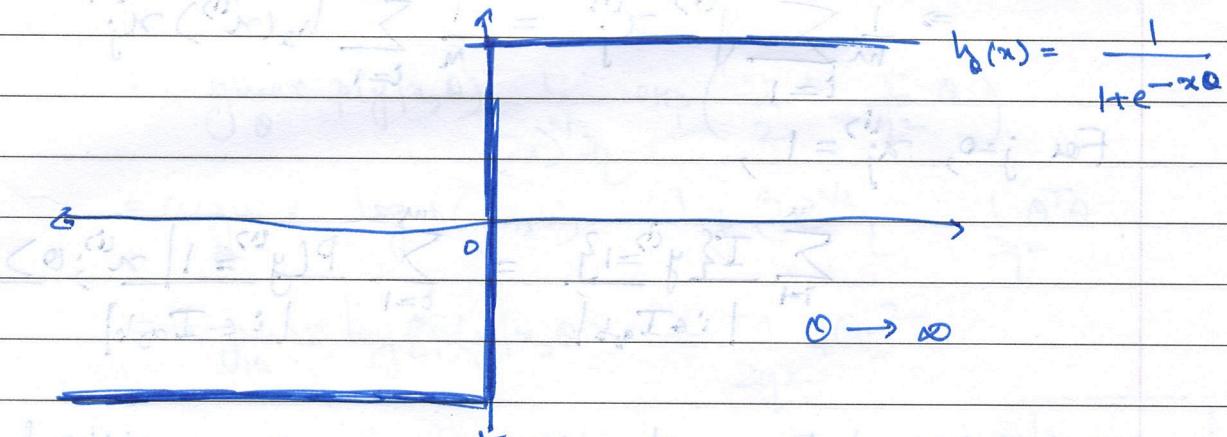


1. (a) Convergence is achieved on dataset A but not on dataset B.

(b) Dataset B is more separated than dataset A. This effectively means that in the loss equation ($J(\theta)$), some of the parameters are going towards infinity. This means that in $h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$, θ can be increased to an infinitely large value and yet the separation may be retained. In a setting where θ , x are 1-D, for example, θ can be set to an infinite value to get a curve like the one below:



Hence, the θ_j 's have the incentive to grow larger.

(c) (i) (ii) (iii) won't fix the issue as they don't modify the loss equation themselves.

(iv) (v) will work: (iv) penalizes the MLE params going towards infinity; (v) will make the dataset not readily separable.

(i) will ensure convergence eventually.

(d) SVMs do not fall prey to B-esque datasets as they tend to work with geometric margins, unlike classic logistic regression where functional margin is used.

FM has the risk of being manipulated by the inflation of parameters; GM uses normalization as a safeguard.

8.2. (a) Gradient calculation:

$$\frac{\partial L(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

When model is done training:

$$\begin{aligned} \frac{\partial L(\theta)}{\partial \theta_j} &= 0 \\ \Rightarrow \frac{1}{m} \sum_{i=1}^m [y^{(i)} - h_\theta(x^{(i)})] x_j^{(i)} &= 0 \end{aligned}$$

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

For $j=0$, $x_j^{(i)} = 1$,

$$\sum_{i=1}^m \frac{I\{y^{(i)}=1\}}{|I \cap I_{\text{test}}|} = \sum_{i=1}^m \frac{P(y^{(i)}=1|x^{(i)};\theta)}{|I \cap I_{\text{test}}|}$$

(b) No to both, if the range is modified to say $(0, 2, 1)$ then $0.2 < P(y^{(i)}=1|x^{(i)};\theta) < 1$

$$\Rightarrow \sum_i P(y|x;\theta) < \sum_i I\{y=1\}$$

Hence, perfect acc violated.

Similarly, it can be shown that perfectly accurate models needn't be perfectly calibrated.

(c)

$$\text{New } J = - \sum_i y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) + \frac{1}{2} \lambda \|\theta\|^2$$

$$\frac{\partial J}{\partial \theta_j} = \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)} + \lambda \theta_j \Rightarrow (\text{full training})$$

$$\Rightarrow \sum_i \frac{I\{y^{(i)}=1\}}{m} = \lambda \theta_j + \sum_i \frac{P(y^{(i)}|x^{(i)};\theta)}{m}$$

Next well-calibrated.

3. $\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} p(\theta | x, y)$
~~= $\underset{\theta}{\operatorname{argmax}} p(y | x, \theta)$~~

3. (a) $\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} p(\theta | x, y) = \underset{\theta}{\operatorname{argmax}} \frac{p(\theta, x, y)}{p(x, y)}$

$$= \underset{\theta}{\operatorname{argmax}} \frac{p(y | x, \theta) p(\theta | x) p(x)}{p(y, x)}$$

$$= \underset{\theta}{\operatorname{argmax}} p(y | x, \theta) p(\theta | x) = \underset{\theta}{\operatorname{argmax}} p(y | x, \theta) p(\theta)$$

(b) $\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} p(y | x, \theta) p(\theta), \quad \theta \sim N(0, \eta^2 I)$

$$= \underset{\theta}{\operatorname{argmax}} p(y | x, \theta) \frac{1}{(2\pi)^{n/2} \eta} \exp\left(-\frac{1}{2} \frac{\theta^T \theta}{\eta^2}\right)$$

$$= \underset{\theta}{\operatorname{argmax}} \log p(y | x, \theta) - \log (2\pi)^{n/2} \eta - \frac{1}{2\eta^2} \theta^T \theta$$

$$= \underset{\theta}{\operatorname{argmax}} \log p(y | x, \theta) - \frac{1}{2\eta^2} \|\theta\|_2^2$$

$$= \underset{\theta}{\operatorname{argmin}} -\log p(y | x, \theta) + \frac{1}{2\eta^2} \|\theta\|_2^2 \Rightarrow \boxed{\lambda = \frac{1}{2\eta^2}}$$

(c) $y = \theta^T x + \epsilon, \quad \epsilon = y - \theta^T x, \quad \epsilon \sim N(0, \sigma^2)$
 $\Rightarrow P(\epsilon) = P(y | x, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y - \theta^T x)^2\right)$

$$\Rightarrow \hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmin}} -\log p(y | x, \theta) + \frac{1}{2\eta^2} \|\theta\|_2^2$$

$$= \underset{\theta}{\operatorname{argmin}} -\log \frac{1}{\sqrt{(2\pi\sigma^2)^n}} \exp\left(-\frac{1}{2} \frac{\|\vec{y} - \vec{x}\theta\|_2^2}{\sigma^2}\right) + \frac{1}{2\eta^2} \|\theta\|_2^2$$

$$= \underset{\theta}{\operatorname{argmin}} + \log \sqrt{2\pi\sigma^2} + \frac{1}{2\sigma^2} \|\vec{y} - \vec{x}\theta\|_2^2 + \frac{1}{2\eta^2} \|\theta\|_2^2$$

$$= \underset{\theta}{\operatorname{argmin}} \left[\frac{1}{2\sigma^2} \|\vec{y} - \vec{x}\theta\|_2^2 + \frac{1}{2\eta^2} \|\theta\|_2^2 \right] \boxed{J(\theta)}$$

argmax _{θ}

$$(1) \quad \Theta_{MAP} = p(y|x, \theta) p(\theta)$$

$$= \text{argmax } \log p(y|x, \theta) + \log p(\theta)$$

$$= \text{argmax}_{\theta} \frac{1}{(2\pi)^m} \exp(-\frac{\|y - x\theta\|_2^2}{2\sigma^2})$$

$$= \text{argmin}_{\theta} \frac{1}{2\sigma^2} \|y - x\theta\|_2^2 + \frac{\|\theta\|_1}{b}$$

$$= \text{argmin}_{\theta} \|y - x\theta\|_2^2 + \underbrace{\frac{2\sigma^2}{b} \|\theta\|_1}_{r}$$

4.

(a) Yes $K(x, z) = k_1(x, z) + k_2(x, z)$

$$\begin{aligned} \Rightarrow y^T K y &= y^T (k_1 + k_2) y \\ &= (y^T k_1 + y^T k_2) y = y^T k_1 y + y^T k_2 y \geq 0 \quad \cancel{\geq 0} \end{aligned}$$

(b) No, if $k_2 = 5k_1$, $y^T K y = -y^T 4k_1 y \leq 0$ ~~No~~

(c) Yes, $y^T a K y = a y^T K y \geq 0$

(d) No, similar to above and (b)

(e) $y^T K y = y^T k_1 k_2 y$

$$= \sum_{i=1}^m \sum_{j=1}^m y_i^T k_1(x^{(i)}, x^{(j)}) k_2(x^{(i)}, x^{(j)}) y_j$$

$$= \sum_{i=1}^m \sum_{j=1}^m y_i^T \phi_1(x^{(i)})^T \phi_1(x^{(j)}) \phi_2(x^{(i)})^T \phi_2(x^{(j)}) y_j$$

$$= \sum_{i=1}^m \sum_{j=1}^m y_i^T \sum_k \phi_{1k}(x^{(i)})^T \phi_{1k}(x^{(j)}) \sum_k \phi_{2k}(x^{(i)})^T \phi_{2k}(x^{(j)}) y_j$$

$$= \sum_k \sum_{i=1}^m \sum_{j=1}^m y_i^T \phi_{1k}(x^{(i)})^T \phi_{1k}(x^{(j)}) \phi_{2k}(x^{(i)})^T \phi_{2k}(x^{(j)}) y_j$$

$$= \sum_{k=1}^n \sum_{i=1}^m \sum_{j=1}^m (y_i^T \phi_{1k}(x^{(i)}) \phi_{2k}(x^{(j)}))^2 \geq 0 \quad (\text{Yes})$$

(f) $y^T K y = y^T f(x) f(z) y = f(x)^T f(z) \|y\|^2 \geq 0 \quad (\text{Yes})$

(g) $y^T K y = y^T k_3(\phi(x), \phi(z)) y$
~~= y^T \phi(\phi(x))^T \phi(\phi(z)) y~~
 $= y^T \phi(x)^T \phi(z) y \geq 0 \quad (\text{Yes})$

(h) From (e) and (g), it can be inferred. (Yes)

5. (a) (i) $\phi^{(i)} = \sum_{j=1}^i \beta_j \phi(x^{(j)})$, $\phi^{(0)} = \sum_{j=1}^0 \beta_j \phi(x^{(j)}) = \vec{0}$

(ii) $h_{\phi^{(i)}}(x^{(i+1)}) = g(\phi^{(i)} \phi(x^{(i+1)}))$

$$= g\left(\left(\sum_{j=1}^i \beta_j \phi(x^{(j)})\right)^T \phi(x^{(i+1)})\right)$$

$$= g\left(\sum_{j=1}^i \beta_j \phi(x^{(j)})^T \phi(x^{(i+1)})\right)$$

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

$$= \text{sign}\left(\sum_{j=1}^i \beta_j K(x^{(j)}, x^{(i+1)})\right)$$

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

$$\begin{aligned} &:= \sum_{j=1}^i \beta_j \phi(x^{(j)}) + \underbrace{\alpha (y^{(i+1)} - \text{sign}\left(\sum_{j=1}^i \beta_j K(\phi^{(i)}, \phi^{(i+1)})\right) \phi^{(i+1)})}_{\beta_{i+1}} \\ &:= \sum_{j=1}^i \beta_j \phi(x^{(j)}) + \alpha y^{(i+1)} \phi^{(i+1)} - \alpha \text{sign}\left(\sum_{j=1}^i \beta_j K(\phi^{(i)}, \phi^{(i+1)})\right) \phi^{(i)} \end{aligned}$$

(b) Coded

(c) The Dot kernel performs worse. RBF kernel defines a larger parameter space hence improving performance. Dot kernel essentially does a plain regression which is considerably fast for the dataset.

6. (a) (b) (c) (d) Coded