

1. (b)

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} - \sum_{n=1}^N \{ y_n \log h_{\theta}(x_n) + (1-y_n) \log (1 - h_{\theta}(x_n)) \}$$

$$J(\theta) = \sum_{n=1}^N \{ y_n \log h_{\theta}(x_n) + (1-y_n) \log [1 - h_{\theta}(x_n)] \}$$

$$= \left(\sum_{n=1}^N y_n x_n \right)^T \theta - \sum_{n=1}^N \log (1 + e^{\theta^T x_n})$$

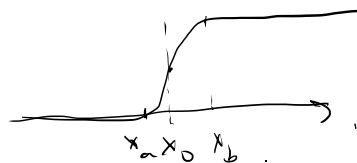
$$\nabla_{\theta} J(\theta) = \left(\sum_{n=1}^N y_n x_n \right) - \sum_{n=1}^N \frac{x_n \cdot e^{\theta^T x_n}}{1 + e^{\theta^T x_n}}$$

$$= \sum_{n=1}^N x_n \cdot \left(y_n - \frac{e^{\theta^T x_n}}{1 + e^{\theta^T x_n}} \right)$$

$$= \sum_{n=1}^N x_n [y_n - h_{\theta}(x)] = 0.$$

If the two classes are linearly separable, the optimal θ should switch from $y_n = 0$ to $y_n = 1$ instantaneously to satisfy the gradient being 0.

Assume x_0 is the switching point.



$$h_{\theta}(x_a) = \frac{1}{1 + e^{-\theta^T(x_a - x_0)}}$$

$$h_{\theta}(x_b) = \frac{1}{1 + e^{\theta^T(x_b - x_0)}} \quad \lim_{\substack{x_b \rightarrow x_0 \\ x_a \rightarrow x_0}} \frac{h_{\theta}(x_b) - h_{\theta}(x_a)}{x_b - x_a} \rightarrow \infty.$$

As the steepness of the curve goes to infinity, $\|\theta\| \rightarrow \infty$.

Hence $\|w\| \rightarrow \infty$ and $\|w_0\| \rightarrow \infty$.

i) As shown in i).

$$\nabla_{\theta} J(\theta) = \sum_{n=1}^N x_n [y_n - h_{\theta}(x_n)] = 0,$$

Since the steepness of $h_{\theta}(x_n)$ only tends to ∞ , but never reaches it, the gradient descent would never converge unless it reaches ∞ , which is impossible.

ii) If we set a maximum steepness on the logistic curve, the gradient descent will stop. once it reaches there,

Another way to counter this problem is to add regularization term in the loss function.

iv) We haven't met non-convergence issue with other linear classifiers, because if the data are linearly separable, then we can invert the matrix and perform good's linear regression.

$$2. J(\theta) = -\frac{1}{N} \sum_{n=1}^N \left\{ y_n \log h_{\theta}(x_n) + (1-y_n) \log(1-h_{\theta}(x_n)) \right\}$$

where $h_{\theta}(x) = \frac{1}{1 + \exp(-\theta^T x)}$. Show $J(\theta)$ convex by

its Hessian, & show it's PSD.

$$J(\theta) = -\frac{1}{N} \sum_{n=1}^N \left\{ \underbrace{y_n \log \left(\frac{h_{\theta}(x_n)}{1-h_{\theta}(x_n)} \right)}_{(1)} + \underbrace{\log(1-h_{\theta}(x_n))}_{(2)} \right\}.$$

$$(1): y_n \log \left(\frac{h_{\theta}(x_n)}{1-h_{\theta}(x_n)} \right)$$

$$= y_n \log \left(\frac{\frac{1}{1+e^{-\theta^T x_n}}}{\frac{e^{-\theta^T x_n}}{1+e^{-\theta^T x_n}}} \right) = y_n \log(e^{\theta^T x_n}) = \underbrace{y_n \theta^T x_n}_{\text{linear}},$$

thus the first term is convex.

$$(2) \quad \nabla_{\theta} [-\log(1-h_{\theta}(x))] = -\nabla_{\theta} \left[\log \left(1 - \frac{1}{1+e^{-\theta^T x}} \right) \right]$$

$$= -\nabla_{\theta} \left[\log \frac{e^{-\theta^T x}}{1+e^{-\theta^T x}} \right] = h_{\theta}(x) x,$$

$$\nabla_{\theta}^2 [-\log(1-h_{\theta}(x))] = h_{\theta}(x) [1-h_{\theta}(x)] x x^T$$

For any $v \in \mathbb{R}^d$,

$$v^T \nabla_{\theta}^2 [-\log(1-h_{\theta}(x))] v = v^T [h_{\theta}(x) (1-h_{\theta}(x)) x x^T] v = h_{\theta}(x) (1-h_{\theta}(x)) \|v^T x\|^2 \geq 0$$

thus the Hessian is PSD., so $-\log(1-h_{\theta}(x))$ is convex in θ .