

Ex. 4.5

$$l(\beta) = \sum_i \{ y_i \beta^T x_i - \log(1 + \exp(\beta^T x_i)) \}$$

Now  $y_i = 1$  if  $x_i > x_0$   
 $y_i = 0$  if  $x_i < x_0$  (since classes separated by  $x_0$ )

$$\Rightarrow \sum_{x_i < x_0} -\log(1 + \exp(\beta^T x_i)) + \sum_{x_i > x_0} \beta^T x_i - \log(1 + \exp(\beta^T x_i))$$

Now  $\beta^T = (\beta_0, \beta_1)$

and notice  $x_i - x_0$  is positive when  $x_i > x_0$  (\*)  
and  $x_i - x_0$  is negative when  $x_i < x_0$

$$l(\beta) = \textcircled{1} \sum_{x_i < x_0} -\log(1 + \exp(\beta_0 + \beta_1 x_0 + \beta_1 (x_i - x_0)))$$

$$+ \textcircled{2} \sum_{x_i > x_0} \beta_0 + \beta_1 x_0 + \beta_1 (x_i - x_0) - \log(1 + \exp(\beta_0 + \beta_1 x_0 + \beta_1 (x_i - x_0)))$$

Now  $\textcircled{1}$  is maximised at 0 as  $\max(-\log) \rightarrow 0$ .  
and  $\textcircled{2}$  may approach  $\infty$ .

Setting  $\beta_0 = -\beta_1 x_0$  we obtain:

$$l(\beta) = \sum_{x_i < x_0} -\log(1 + \exp(\beta_1 (x_i - x_0))) + \sum_{x_i > x_0} \beta_1 (x_i - x_0) - \log(1 + \exp(\beta_1 (x_i - x_0)))$$

and using (\*) above this is maximised as  $\beta_1 \rightarrow \infty$  as:

$$-\log(1 + 1/\infty) \rightarrow 0$$

and

$$\infty - \log(1 + e^\infty) \rightarrow \infty$$

Thus MLE's are  $\beta_0 = -\beta_1 x_0$  and  $\beta_1 = \infty$

(a) Now we consider the case where  $x_i \in \mathbb{R}^p$ .

If the classes are separable then there exists some HyperPlane,  $\alpha$ , such that:

$$\begin{aligned} \alpha x_i &> 0 & \text{if } y_i = 1 \\ \alpha x_i &< 0 & \text{if } y_i = 0 \end{aligned}$$

Now minimising  $\alpha x_i$  when  $y_i = 1$  and  $-\alpha x_i$  when  $y_i = 0$   
We find there exists some minimum value  $\epsilon > 0$  such that:

$$\begin{aligned} \alpha x_i &\geq \epsilon & \text{if } y_i = 1 \\ \alpha x_i &\leq -\epsilon & \text{if } y_i = 0 \end{aligned} \quad (*)$$

The Problem is reduced to analyzing what happens in the direction of  $\alpha$ .

Now  $h$  is logistic regression's link function that continuously maps the interval of possible probabilities  $(0, 1)$  with the real line (e.g. 4.18) and  $h^{-1}$  exists so

$$\lim_{x \rightarrow \infty} h^{-1}(x) = 1 \quad \text{and} \quad \lim_{x \rightarrow -\infty} h^{-1}(x) = 0$$

Then the likelihood for coefficient vector  $\beta$  is the Product, over all observations, of the chance that a Bernoulli Variable with Parameter  $\beta^T x_i = y_i$ :

$$L(\beta) = \prod_{i: y_i = 1} h^{-1}(\beta^T x_i) \prod_{i: y_i = 0} (1 - h^{-1}(\beta^T x_i))$$

Consider a positive real number  $\lambda$ , using (\*) and the fact that  $h^{-1}$  is increasing, then

$$L(\alpha\lambda) = \prod_{i: y_i=1} h'(\lambda \alpha x_i) \prod_{i: y_i=0} (1 - h'(\lambda \alpha x_i))$$

$$\geq \prod_{i: y_i=1} h'(\lambda \epsilon) \prod_{i: y_i=0} (1 - h'(-\lambda \epsilon))$$

Now, Since  $\epsilon > 0$ , as  $\lambda \rightarrow \infty$

$$\epsilon \lambda \rightarrow \infty$$

$$\text{and } -\epsilon \lambda \rightarrow -\infty$$

giving:

$$\lim_{\epsilon \rightarrow \infty} h'(\lambda \epsilon) = 1 \quad \text{and} \quad \lim_{\epsilon \rightarrow -\infty} 1 - h'(-\lambda \epsilon) = 1$$

So,

$$\lim_{\epsilon \rightarrow \infty} L(\alpha\lambda) \geq \prod_{i: y_i=1} 1 \prod_{i: y_i=0} 1 = 1$$

And Since this is a Product of Probabilities, this is the global Maximum.

Finally we notice that it is Possible For components of  $\alpha$  to be 0 and when multiplied by large  $\lambda$  would remain 0 (not undefined).

To deal with these cases we notice that we can adjust the non-zero components of  $\alpha$  so that all components are greater than some  $\epsilon^*$  where  $\epsilon > \epsilon^* > 0$  and therefore  $\lambda \epsilon^* \rightarrow \infty$  as  $\lambda \rightarrow \infty$  making the above analysis still valid. Now we can find another hyperplane (of the infinite that exist) such that all components are non-zero and still separate the two classes.

Therefore we conclude, when a separating hyperplane exists, the likelihood can be maximised in a manner that causes all of the coefficients of  $\beta$  to diverge.



There may exist separating hyperplanes in which some components of  $\beta$  are finite, but in these cases:

1. Any such component must be zero.
2. At least one component will diverge.
3. There can be infinitely many directions  $\alpha$  for which the likelihood of  $\beta = \lambda\alpha$  is maximised as  $\lambda \rightarrow \infty$ .

(b) In this case we know there exists hyperplanes separating each individual class from all of the others. In this case we may find these  $K-1$  hyperplanes sequentially by labelling the current class as  $y_i = 1$  and all other classes as  $y_i = 0$  and following the exact steps as described in Part (a). There are  $K-1$  rather than  $K$  hyperplanes as once we reach the  $K^{\text{th}}$  class we can simply use the previous  $K-1$  hyperplanes to classify the observations and any remaining observations that do not belong to these  $K-1$  classes must belong to this final  $K^{\text{th}}$  class.