

XCS299i Problem Set #2

JUAN RICARDO PEDRAZA ESCOBAR

Machine Learning

Stanford Center for Professional Development

August, 2021

1.a

Since  $g(z) = g(z)(1 - g(z))$  and  $h(x) = g(\theta^T x)$ , it follows  $\partial h(x)/\partial \theta_k = h_\theta(x^{(i)})(1 - h_\theta(x^{(i)}))$ .

Letting  $h_\theta(x^{(i)}) = g(\theta^T x^{(i)}) = 1/(1 + \exp(-\theta^T x^{(i)}))$ , we have

$$\frac{\partial \log h_\theta(x^{(i)})}{\partial \theta_k} = \frac{1}{h_\theta(x^{(i)})} \frac{\partial h_\theta(x^{(i)})}{\partial \theta_k}$$

Replacing  $\partial h(x)/\partial \theta_k$ ,

$$\frac{\partial \log h_\theta(x^{(i)})}{\partial \theta_k} = \frac{1}{h_\theta(x^{(i)})} h_\theta(x^{(i)}) (1 - h_\theta(x^{(i)})) x_k = (1 - h_\theta(x^{(i)})) x_k$$

$$\frac{\partial \log (1 - h_\theta(x^{(i)}))}{\partial \theta_k} = \frac{1}{1 - h_\theta(x^{(i)})} \frac{\partial (1 - h_\theta(x^{(i)}))}{\partial \theta_k} = -h_\theta(x^{(i)}) x_k$$

Substituting into our equation for  $J(\theta)$ , we have

$$\frac{\partial J(\theta)}{\partial \theta_k} = -\frac{1}{n} \sum_{i=1}^n y^{(i)} (1 - h_\theta(x^{(i)})) x_k - (1 - y^{(i)}) h_\theta(x^{(i)}) x_k$$

$$\frac{\partial J(\theta)}{\partial \theta_k} = -\frac{1}{n} \sum_{i=1}^n (y^{(i)} - h_\theta(x^{(i)})) x_k$$

Consequently, the  $(k, l)$  entry of the Hessian is given by,

$$H_{kl} = \frac{\partial^2 J(\theta)}{\partial \theta_k \partial \theta_l}$$

$$\begin{aligned} \frac{\partial^2 J(\theta)}{\partial \theta_k \partial \theta_l} &= -\frac{1}{n} \sum_{i=1}^n \frac{\partial J(\theta)}{\partial \theta_l} (y^{(i)} - h_\theta(x^{(i)})) x_k \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\partial J(\theta)}{\partial \theta_l} (y^{(i)} - h_\theta(x^{(i)})) x_k \end{aligned}$$

By deriving with respect to  $\theta_l$ ,

$$\frac{\partial^2 J(\theta)}{\partial \theta_k \partial \theta_l} = \frac{1}{n} \sum_{i=1}^n h_\theta(x^{(i)}) (1 - h_\theta(x^{(i)})) x_k^{(i)} x_l^{(i)}$$

Using the fact that  $X_{ij} = x_i x_j$  if and only if  $X = x x^T$ , we have

$$H = \frac{1}{n} \sum_{i=1}^n h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) x^{(i)} x^{(i)T}$$

To prove that  $H$  is positive semi-definite, show  $z^T H z \geq 0$  for all  $z \in \mathbb{R}^d$ .

$$z^T H z = \frac{1}{n} z^T \left( \sum_{i=1}^n h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) x^{(i)} x^{(i)T} \right) z$$

$$z^T H z = \frac{1}{n} \sum_{i=1}^n h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) (z^T x^{(i)}) (x^{(i)T} z)$$

$$z^T H z = \frac{1}{n} \sum_{i=1}^n h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) (z^T x^{(i)})^2$$

- Note that the function  $h_{\theta}(x^{(i)})$  follows  $0 \leq h_{\theta}(x^{(i)}) \leq 1$  and the square term  $(z^T x^{(i)})^2$  is always positive so ,

$$\frac{1}{n} \sum_{i=1}^n h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) (z^T x^{(i)})^2 \geq 0$$

$$z^T H z \geq 0$$

1.c

For shorthand, we let  $H = \{\phi, \Sigma, \mu_0, \mu_1\}$  denote the parameters for the problem. Since the given formulae are conditioned on  $y$ , use Bayes rule to get:

$$\begin{aligned} p(y = 1|x; H) &= \frac{p(x|y = 1; H)p(y = 1; H)}{p(x; H)} \\ &= \frac{p(x|y = 1; H)p(y = 1; H)}{p(x|y = 1; H)p(y = 1; H) + p(x|y = 0; H)p(y = 0; H)} \end{aligned}$$

We replace the equations that we use to model the distribution of  $(x, y)$ ,

$$\begin{aligned} p(y = 1|x; H) &= \frac{\exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right) \phi}{\exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right) \phi + \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right) (1 - \phi)} \\ &= \frac{1}{1 + \exp\left(\log\left(\frac{1-\phi}{\phi}\right) - \frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0) + \frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right)} \end{aligned}$$

We take the quadratic term located in the previous expression denominator and apply properties of distribution,

$$\begin{aligned} &(x - \mu_0)^T \Sigma^{-1} (x - \mu_0) + (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \\ &= x^T \Sigma^{-1} x - \mu_0^T \Sigma^{-1} x - x^T \Sigma^{-1} \mu_0 + \mu_0^T \Sigma^{-1} \mu_0 - x^T \Sigma^{-1} x + \mu_1^T \Sigma^{-1} x + x^T \Sigma^{-1} \mu_1 - \mu_1^T \Sigma^{-1} \mu_1 \end{aligned}$$

Simplifying and reorganizing the term,

$$= 2(\mu_1 - \mu_0)^T \Sigma^{-1} x + \mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1$$

Therefore,

$$p(y = 1|x; H) = \frac{1}{1 + \exp\left(\log\left(\frac{1-\phi}{\phi}\right) + (\mu_1 - \mu_0)^T \Sigma^{-1} x + \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 - \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0\right)}$$

Where,

- $\theta = -(\mu_0 - \mu_1) \Sigma^{-1}$
- $\theta_0 = \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 - \log\left(\frac{1-\phi}{\phi}\right)$

Finally gives,

$$p(y = 1|x; \phi, \Sigma, \mu_0, \mu_1) = p(y = 1|x; H) = \frac{1}{1 + \exp(-(\theta^T x + \theta_0))}$$

1.d

First, derive the expression for the log-likelihood of the training data:

$$\begin{aligned}\ell(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^n p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi) \\ &= \sum_{i=1}^n \log p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) + \sum_{i=1}^n \log p(y^{(i)}; \phi)\end{aligned}$$

Using the equations for this model,

$$\begin{aligned}&= \sum_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x^{(i)} - (1 - y^{(i)})\mu_0 - y^{(i)}\mu_1)^2}{2\sigma^2}\right) + \sum_{i=1}^n \log \phi^{y^{(i)}} (1 - \phi)^{1-y^{(i)}} \\ &= -m \log(\sqrt{2\pi}\sigma) + \sum_{i=1}^n \left[ y^{(i)} \log \phi + (1 - y^{(i)}) (\log(1 - \phi)) - \frac{(x^{(i)} - (1 - y^{(i)})\mu_0 - y^{(i)}\mu_1)^2}{2\sigma^2} \right]\end{aligned}$$

Now, the likelihood is maximized by setting the derivative (or gradient) with respect to each of the parameters to zero.

- For  $\phi$  :

$$\begin{aligned}\frac{\partial \ell}{\partial \phi} &= \sum_{i=1}^n y^{(i)} \frac{1}{\phi} + (1 - y^{(i)}) \frac{1}{1 - \phi} \\ \frac{\partial \ell}{\partial \phi} &= \sum_{i=1}^n y^{(i)} \frac{1}{\phi} + (1 - y^{(i)}) \frac{1}{1 - \phi} \\ &= \frac{\sum_{i=1}^n 1(y^{(i)} = 1)}{\phi} + \frac{m - \sum_{i=1}^n 1(y^{(i)} = 1)}{1 - \phi}\end{aligned}$$

Setting this equal to zero and solving for  $\phi$  gives the maximum likelihood estimate.

$$\begin{aligned}0 &= \frac{\sum_{i=1}^n 1(y^{(i)} = 1)}{\phi} + \frac{m - \sum_{i=1}^n 1(y^{(i)} = 1)}{1 - \phi} \\ 0 &= \sum_{i=1}^n 1(y^{(i)} = 1) (1 - \phi) + \left(m - \sum_{i=1}^n 1(y^{(i)} = 1)\right) \phi\end{aligned}$$

$$\phi = \frac{\sum_{i=1}^n 1(y^{(i)} = 1)}{n}$$

$$\phi = \frac{1}{n} \sum_{i=1}^n 1(y^{(i)} = 1)$$

- For  $\mu_0$ :

Hint: Remember that  $\Sigma$  (and thus  $\Sigma^{-1}$ ) is symmetric.

$$\begin{aligned} \nabla_{\mu_0} \ell &= \frac{\partial \ell}{\partial \mu_0} \sum_{i=1}^n -\frac{(x^{(i)} - (1 - y^{(i)})\mu_0 - y^{(i)}\mu_1)^2}{2\sigma^2} \\ &= \sum_{i=1}^n -\left(\frac{(x^{(i)} - (1 - y^{(i)})\mu_0 - y^{(i)}\mu_1)}{\sigma^2}\right) \frac{\partial \ell}{\partial \mu_0} (x^{(i)} - (1 - y^{(i)})\mu_0 - y^{(i)}\mu_1) \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n (x^{(i)} - \mu_0) 1(y^{(i)} = 0) \end{aligned}$$

Setting this gradient to zero gives the maximum likelihood estimate for  $\mu_0$ .

$$\begin{aligned} 0 &= \sum_{i=1}^n (x^{(i)} - \mu_0) 1(y^{(i)} = 0) \\ 0 &= \sum_{i=1}^n 1(y^{(i)} = 0)x^{(i)} - \sum_{i=1}^n 1(y^{(i)} = 0)\mu_0 \\ \mu_0 &= \frac{\sum_{i=1}^n 1(y^{(i)} = 0)x^{(i)}}{\sum_{i=1}^n 1(y^{(i)} = 0)} \end{aligned}$$

- For  $\mu_1$ :

Hint: Remember that  $\Sigma$  (and thus  $\Sigma^{-1}$ ) is symmetric.

$$\begin{aligned} \nabla_{\mu_1} \ell &= \frac{\partial \ell}{\partial \mu_1} \sum_{i=1}^n -\frac{(x^{(i)} - (1 - y^{(i)})\mu_0 - y^{(i)}\mu_1)^2}{2\sigma^2} \\ &= \sum_{i=1}^n -\left(\frac{(x^{(i)} - (1 - y^{(i)})\mu_0 - y^{(i)}\mu_1)}{\sigma^2}\right) \frac{\partial \ell}{\partial \mu_1} (x^{(i)} - (1 - y^{(i)})\mu_0 - y^{(i)}\mu_1) \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n (x^{(i)} - \mu_1) 1(y^{(i)} = 1) \end{aligned}$$

Setting this gradient to zero gives the maximum likelihood estimate for  $\mu_1$ .

$$\begin{aligned}
 0 &= \sum_{i=1}^n (x^{(i)} - \mu_0) \mathbb{1}(y^{(i)} = 1) \\
 0 &= \sum_{i=1}^n \mathbb{1}(y^{(i)} = 1)x^{(i)} - \sum_{i=1}^n \mathbb{1}(y^{(i)} = 0)\mu_1 \\
 \mu_1 &= \frac{\sum_{i=1}^n \mathbb{1}(y^{(i)} = 1)x^{(i)}}{\sum_{i=1}^n \mathbb{1}(y^{(i)} = 1)}
 \end{aligned}$$

- For  $\Sigma$ , note that  $\Sigma = \sigma^2$ ,

$$\begin{aligned}
 \frac{\partial \ell}{\partial \sigma^2} &= \frac{\partial \ell}{\partial \sigma^2} - m \log(\sqrt{2\pi}\sigma) + \sum_{i=1}^n \left[ y^{(i)} \log \phi + (1 - y^{(i)}) (\log(1 - \phi)) - \frac{(x^{(i)} - (1 - y^{(i)})\mu_0 - y^{(i)}\mu_1)^2}{2\sigma^2} \right] \\
 &= \frac{\partial \ell}{\partial \sigma^2} \left[ -m \log(\sqrt{2\pi}\sigma) - \frac{(x^{(i)} - (1 - y^{(i)})\mu_0 - y^{(i)}\mu_1)^2}{2\sigma^2} \right] \\
 &= -\frac{m}{2\sigma^2} + \sum_{i=1}^n \mathbb{1}(y^{(i)} = 0) \frac{(x^{(i)} - \mu_0)^2}{\sigma^4} + \sum_{i=1}^n \mathbb{1}(y^{(i)} = 1) \frac{(x^{(i)} - \mu_1)^2}{\sigma^4}
 \end{aligned}$$

Setting this gradient to zero gives the maximum likelihood estimate for  $\sigma^2$

$$\sigma^2 = \frac{1}{n} \left( \sum_{i=1}^n \mathbb{1}(y^{(i)} = 0) (x^{(i)} - \mu_0)^2 + \sum_{i=1}^n \mathbb{1}(y^{(i)} = 1) (x^{(i)} - \mu_1)^2 \right)$$

$$\Sigma = \sigma^2 = \frac{1}{n} \left( \sum_{i=1}^n (x^{(i)} - y_{(i)}) (x^{(i)} - y_{(i)})^T \right)$$

1.f

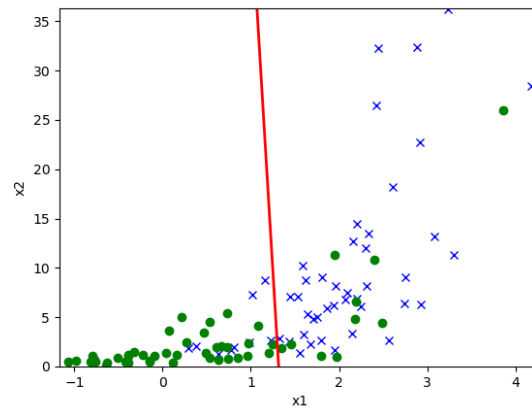


Figure 1 . GDA decision boundary on data set 1

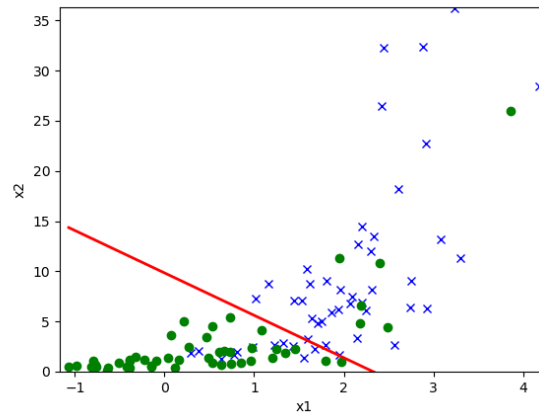


Figure 2. Logistic regression decision boundary dataset 1.

- Comment

The decision boundary from logistic regression seems more reasonable on Dataset1. The logistic decision boundary puts weight on both  $x_1$  and  $x_2$ . However, with the GDA decision boundary the accuracy is lower on Dataset 1.



1.g

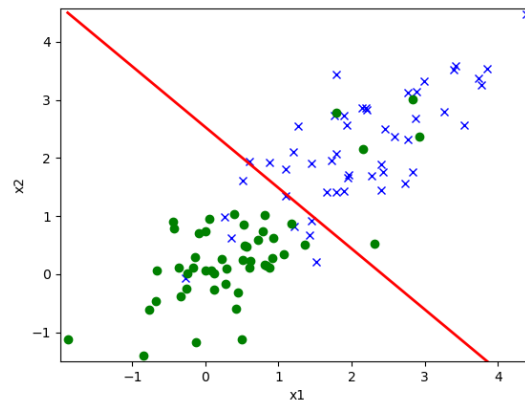


Figure 3. GDA Decision Boundary on Dataset 2.

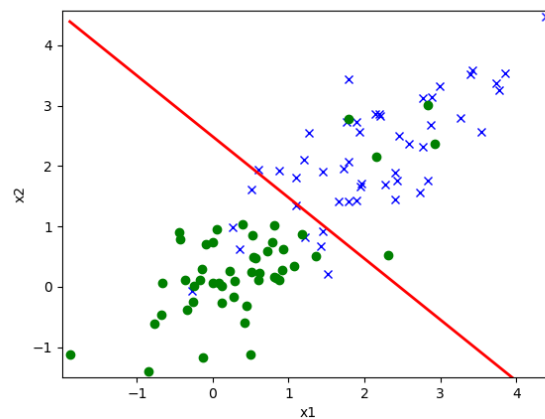


Figure 4. Logistic regression Decision Boundary on Dataset 2.

- Comment

Decision boundaries from logistic regression and GDA are nearly identical for Dataset 2. However, in the first dataset (Dataset 1) GDA seem to perform worse than logistic regression. This probably relates to the fact that the training data( $x_i$ ) are not Gaussian and make the GDA model performs worse. Also, as we already know logistic regression is more powerful when the data set is not drawn from a multivariate Gaussian.

1.h

- Comment.

By setting,

$$x_2^{(i)} := \log x_2^{(i)}$$

all  $x_2^{(i)}$  in the Dataset 2 become Gaussian, therefore the GDA model performs significantly better.

2.a

$$p(y; \lambda) = \frac{e^{-\lambda} e^{y \log \lambda}}{y!}$$

$$p(y; \lambda) = \frac{1}{y!} \exp(y \log \lambda - \lambda)$$

Using the standard form for the exponential family to compare,

$$\eta = \log \lambda$$

$$b(y) = \frac{1}{y!}$$

$$T(y) = y$$

$$\alpha(\eta) = e^{-\eta}$$

2.b

For GML model the canonical response function will be,

$$g(\eta) = E[Y; \eta]$$

$$g(\eta) = \lambda$$

$$g(\eta) = e^{\eta}$$

2.c

The log-likelihood of an example  $(\mathbf{x}^{(i)}, y^{(i)})$  is defined as  $\ell(\boldsymbol{\theta}) = \log p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta})$ . To derive the stochastic gradient ascent rule, use the results in part (a) and the standard GLM assumption that  $\boldsymbol{\eta} = \boldsymbol{\theta}^T \mathbf{x}$ .

$$\begin{aligned}
 \frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_j} &= \frac{\partial \log(p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}))}{\partial \theta_j} \\
 &= \frac{\partial \log\left(\frac{1}{y!} \exp(\boldsymbol{\eta}^T \mathbf{y}^{(i)} - e^{\boldsymbol{\eta}})\right)}{\partial \theta_j} \\
 &= \frac{\partial \log(\exp(\boldsymbol{\eta}^T \mathbf{y}^{(i)} - e^{\boldsymbol{\eta}}))}{\partial \theta_j} + \frac{\partial \log\left(\frac{1}{y!}\right)}{\partial \theta_j} \\
 &= \frac{\partial \left( \left( \sum_k \theta_k x_k^{(i)} \right) y^{(i)} - e^{\sum_k \theta_k x_k^{(i)}} \right)}{\partial \theta_j} \\
 &= \left( y^{(i)} - e^{\sum_k \theta_k x_k^{(i)}} \right) x_j^{(i)}
 \end{aligned}$$

Thus, the stochastic gradient ascent update rule should be:

$$\theta_j := \theta_j + \alpha \frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_j}$$

which reduces here to:

$$\theta_j := \theta_j + \alpha \left( y^{(i)} - e^{\sum_k \theta_k x_k^{(i)}} \right) x_j^{(i)}$$