

$$\textcircled{1} \quad a) \frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - g(\theta^T x^{(i)})) x_j^{(i)}$$

$$\frac{\partial^2 J(\theta)}{\partial \theta_j \partial \theta_k} = -\frac{1}{m} \sum_{i=1}^m g(\theta^T x^{(i)}) (1 - g(\theta^T x^{(i)})) x_j^{(i)} x_k^{(i)}$$

For any  $z$ ,

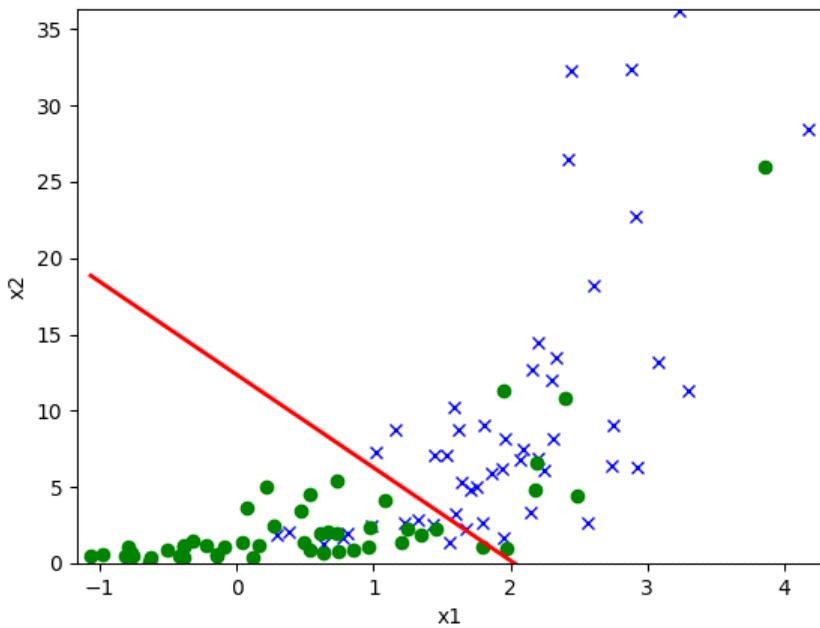
$$z^T H z = \frac{1}{m} \sum_{i=1}^m \sum_{j,k=1}^m g(\theta^T x^{(i)}) (1 - g(\theta^T x^{(i)})) x_j^{(i)} x_k^{(i)} z_j z_k$$

$$= \frac{1}{m} \sum_{i=1}^m \sum_{j,k=1}^m g(\theta^T x^{(i)}) (1 - g(\theta^T x^{(i)})) (x^{(i)})^T z^2 \geq 0$$

# PS1 - 1b

Friday, July 8, 2022 7:48 PM

b)



$$\begin{aligned}
 \text{Q) } p(y=1 | x, \emptyset, \mu_0, \mu_1, \Sigma) &= \frac{p(x|y=1) p(y=1)}{p(x|y=1) p(y=1) + p(x|y=0) p(y=0)} \\
 &= \frac{1}{1 + \frac{p(x|y=0) p(y=0)}{p(x|y=1) p(y=1)}} \\
 &= \frac{1}{1 + \exp\left(-\frac{1}{2}[(x-\mu_0)^T \Sigma^{-1} (x-\mu_0) - (x-\mu_0)^T \Sigma^{-1} (x-\mu_1)]\right)} \quad \frac{1-\phi}{\phi} \\
 &= \frac{1}{1 + \exp\left(-\frac{1}{2}[(\mu_1 - \mu_0)^T \Sigma^{-1} x] + \frac{1}{2}(\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1)\right) - \log \frac{1-\phi}{\phi}}
 \end{aligned}$$

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

$$\begin{aligned}
 \text{d) } \ell(\phi, \mu_0, \mu_1, \varepsilon) &= \log \prod_{i=1}^n p(x^{(i)} | y^{(i)}, \mu_0, \mu_1, \varepsilon) p(y^{(i)}; \phi) \\
 &= \sum_{i=1}^n \log p(x^{(i)} | y^{(i)}, \mu_0, \mu_1, \varepsilon) + \sum_{i=1}^n \log p(y^{(i)}; \phi) \\
 &= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} \exp^{-\left\{\frac{(x^{(i)} - (1-y^{(i)})\mu_0 - y^{(i)}\mu_1)^2}{2\sigma^2}\right\}} + \sum_{i=1}^n \log (\phi^{y^{(i)}} (1-\phi)^{1-y^{(i)}}) \\
 &= -m \log (\sqrt{2\pi}\sigma) + \sum_{i=1}^n \left[ -\left\{\frac{(x^{(i)} - (1-y^{(i)})\mu_0 - y^{(i)}\mu_1)^2}{2\sigma^2}\right\} + y^{(i)} \log(\phi) + (1-y^{(i)}) \log(1-\phi) \right] \\
 \frac{\partial \ell}{\partial \phi} &= \sum_{i=1}^n y^{(i)} \frac{1}{\phi} + (1-y^{(i)}) \frac{1}{1-\phi} = \underbrace{\sum_{i=1}^n \mathbb{I}(y^{(i)}=1)}_{\phi} + \underbrace{m - \sum_{i=1}^n \mathbb{I}(y^{(i)}=1)}_{1-\phi} = 0 \\
 \sum_{i=1}^n \mathbb{I}(y^{(i)}=1) &+ \frac{m - \sum_{i=1}^n \mathbb{I}(y^{(i)}=1)}{1-\phi} = 0 \\
 \sum_{i=1}^n \mathbb{I}(y^{(i)}=1)(1-\phi) + (m - \sum_{i=1}^n \mathbb{I}(y^{(i)}=1))\phi &= 0 \\
 \phi &= \frac{\sum_{i=1}^n \mathbb{I}(y^{(i)}=1)}{n} \\
 \frac{\partial \ell}{\partial \mu_0} &= \frac{\partial}{\partial \mu_0} \sum_{i=1}^n -\left\{\frac{x^{(i)} - (1-y^{(i)})\mu_0 - y^{(i)}\mu_1)^2}{2\sigma^2}\right\} \\
 &= \sum_{i=1}^n -\frac{x^{(i)} - (1-y^{(i)})\mu_0 - y^{(i)}\mu_1}{\sigma^2} \frac{\partial}{\partial \mu_0} (x^{(i)} - (1-y^{(i)})\mu_0 - y^{(i)}\mu_1) \\
 &= \frac{1}{\sigma^2} \sum_{i=1}^n \mathbb{I}(y^{(i)}=0) (x^{(i)} - \mu_0) \Rightarrow \phi \Rightarrow \\
 \mu_0 &= \frac{\sum_{i=1}^n \mathbb{I}(y^{(i)}=0) x^{(i)}}{\sum_{i=1}^n \mathbb{I}(y^{(i)}=0)} \\
 \text{By symmetry, } \mu_1 &= \frac{\sum_{i=1}^n \mathbb{I}(y^{(i)}=1) x^{(i)}}{\sum_{i=1}^n \mathbb{I}(y^{(i)}=1)}
 \end{aligned}$$

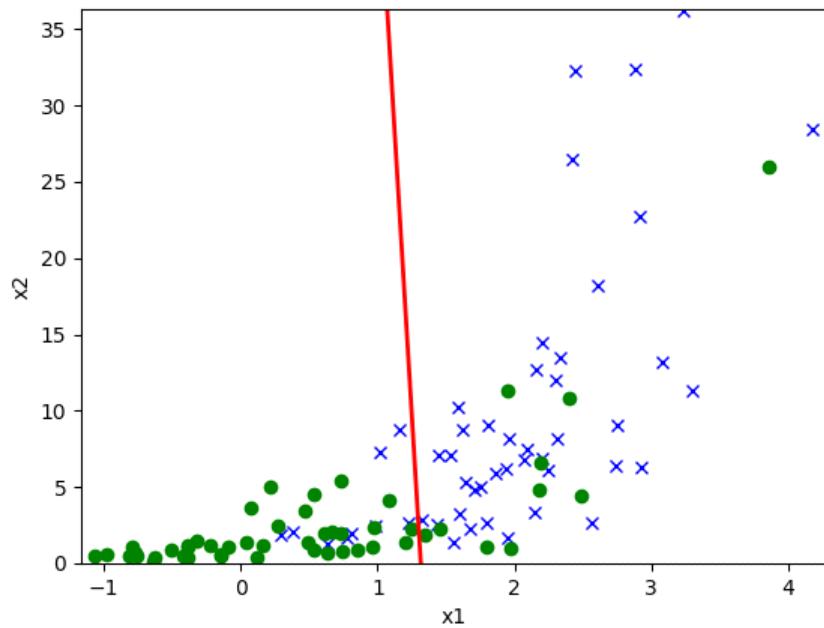
$$\begin{aligned}
 \varepsilon &= \sigma^2 \\
 \frac{\partial \ell}{\partial \sigma^2} &= \frac{\partial}{\partial \sigma^2} \left[ -m \log \sqrt{2\pi\sigma^2} - \sum_{i=1}^n \left\{ \frac{(x^{(i)} - (1-y^{(i)})\mu_0 - y^{(i)}\mu_1)^2}{2\sigma^2} \right\} \right]
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial L}{\partial \sigma} &= \sum_{i=1}^n \left[ \frac{\partial \ln L}{\partial \sigma} \right] \\
 &= -\frac{m}{2\sigma^2} + \sum_{i=1}^n \left[ I(y^{(i)} = 0) \left( \frac{(x^{(i)} - \mu_0)^2}{2\sigma^4} \right) + I(y^{(i)} = 1) \left( \frac{(x^{(i)} - \mu_1)^2}{2\sigma^4} \right) \right] = 0 \Rightarrow \\
 \frac{m}{\sigma^2} &= \frac{1}{2(\sigma^2)^2} \sum_{i=1}^n \left[ I(y^{(i)} = 0) (x^{(i)} - \mu_0)^2 + I(y^{(i)} = 1) (x^{(i)} - \mu_1)^2 \right] \\
 \sigma^2 &= \frac{1}{n} \sum_{i=1}^n \left[ I(y^{(i)} = 0) (x^{(i)} - \mu_0)^2 + I(y^{(i)} = 1) (x^{(i)} - \mu_1)^2 \right] \\
 &= \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu_{y^{(i)}})^2
 \end{aligned}$$

# PS1 - 1e

Friday, July 8, 2022 8:11 PM

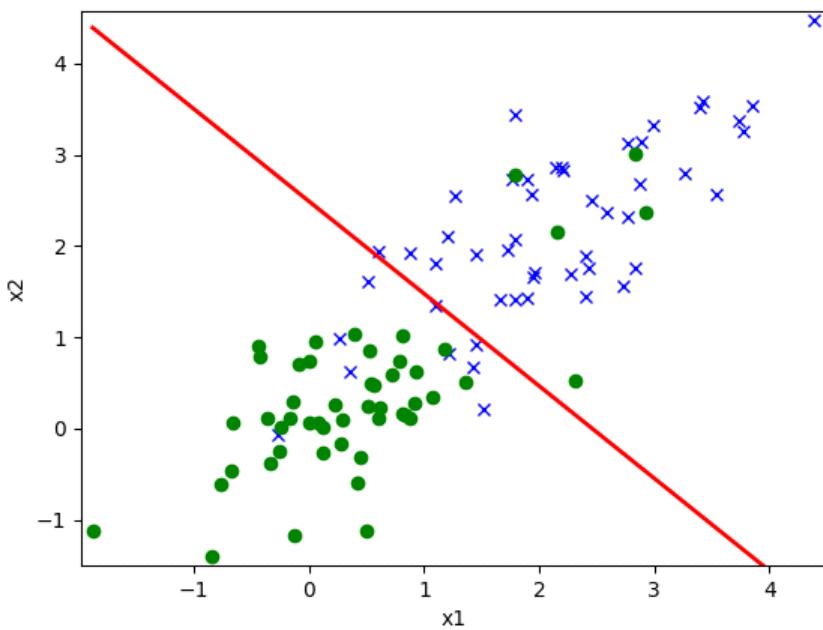
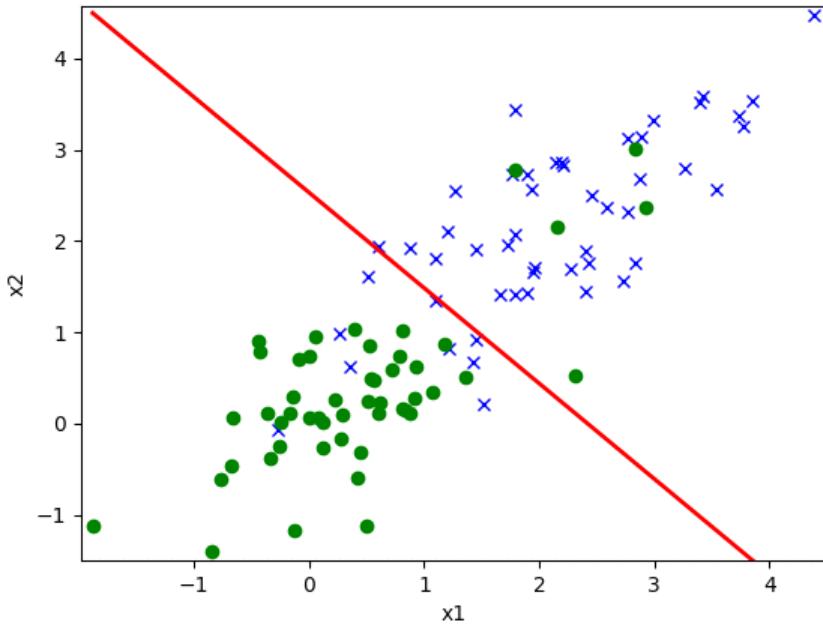
e)



# PS1 - 1fgh

Friday, July 8, 2022 8:11 PM

- f) For dataset 1, the decision boundaries for logistic regression and GDA are similar, however for logistic regression the line is more diagonal whereas for GDA it is closer to vertical. GDA has slightly more green dots that are misclassified whereas logistic regression has a few more blue data points that are misclassified. Overall, logistic regression appears to perform slightly better
- g) Logistic Regression



GDA

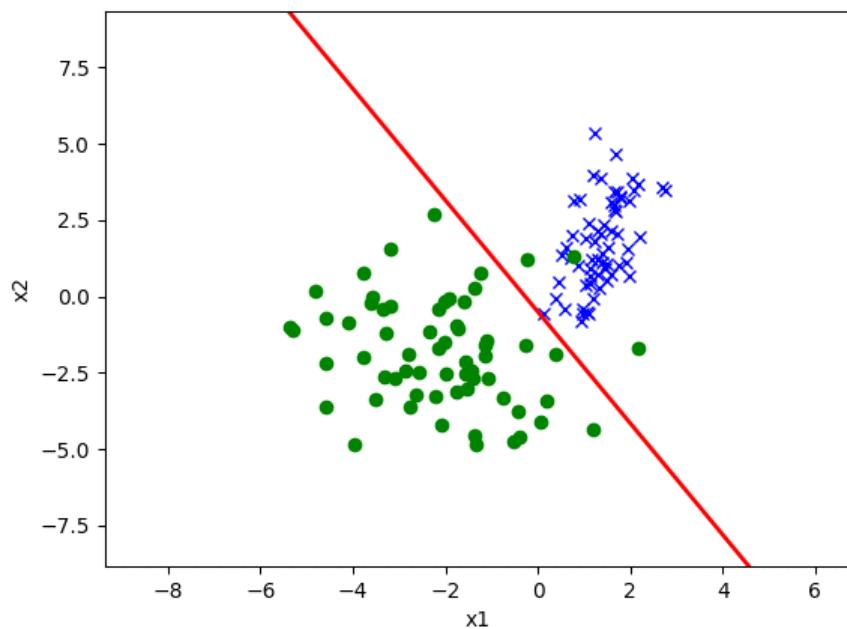
For dataset 2, the two decision boundaries are nearly identical. For dataset 1, GDA appears to perform slightly worse than logistic regression as there are a few more green dots that are misclassified. This may be because the distribution of the data is not Gaussian, which is one of the assumptions made by GDA.

- h) Any transformation which causes the data to become closer to Gaussian distributed would work.  
An example may be taking the logarithm of the data.

# PS1 - 2a

Friday, July 8, 2022 8:11 PM

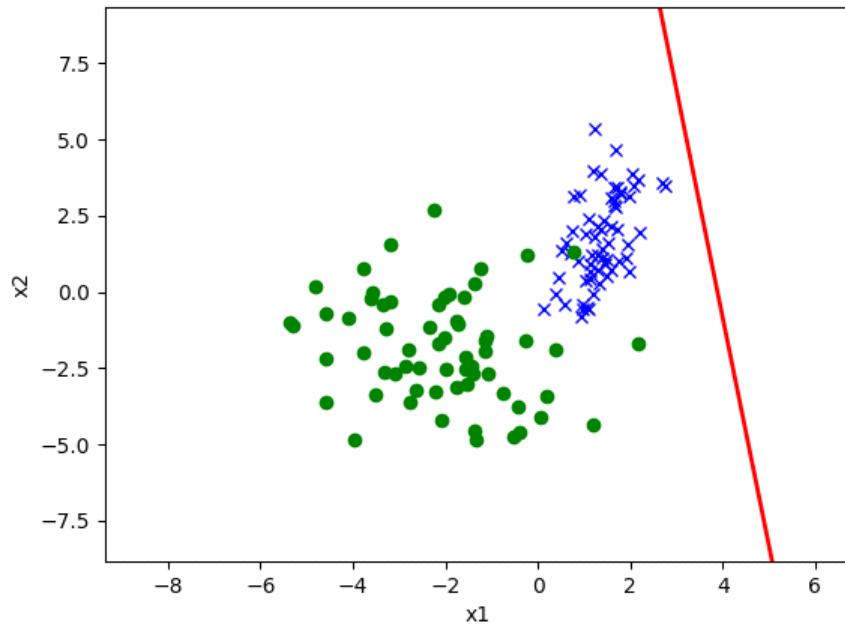
2. a)



# PS1 - 2b

Friday, July 8, 2022 8:11 PM

b)



$$\begin{aligned}
 c) p(f^{(1)}=1 | y^{(i)}=1, x^{(i)}) &= \frac{p(y^{(i)}=1 | f^{(i)}=1, x^{(i)}) p(f^{(i)}=1 | x^{(i)})}{p(y^{(i)}=1 | x^{(i)})} \\
 &= \frac{p(y^{(i)}=1 | f^{(i)}=1, x^{(i)}) p(f^{(i)}=1 | x^{(i)})}{p(y^{(i)}=1 | f^{(i)}=0, x^{(i)}) p(f^{(i)}=0 | x^{(i)}) + p(y^{(i)}=1 | f^{(i)}=1, x^{(i)}) p(f^{(i)}=1 | x^{(i)})} \\
 &= \frac{p(y^{(i)}=1 | f^{(i)}=1, x^{(i)}) p(f^{(i)}=1 | x^{(i)})}{p(y^{(i)}=1 | f^{(i)}=1, x^{(i)}) p(f^{(i)}=1 | x^{(i)})} \\
 &= \frac{0 \cdot p(f^{(i)}=0 | x^{(i)}) + p(y^{(i)}=1 | f^{(i)}=1, x^{(i)}) p(f^{(i)}=1 | x^{(i)})}{p(y^{(i)}=1 | f^{(i)}=1, x^{(i)}) p(f^{(i)}=1 | x^{(i)})} \\
 &= \frac{p(y^{(i)}=1 | f^{(i)}=1, x^{(i)})}{p(y^{(i)}=1 | f^{(i)}=1, x^{(i)})} \\
 &= 1
 \end{aligned}$$

$$\begin{aligned} d) \quad p(y^{(i)}=1 | x^{(i)}) &= p(y^{(i)}=1 \wedge t^{(i)}=1 | x^{(i)}) \text{ from c)} \\ &= p(t^{(i)}=1 | x^{(i)}) p(y^{(i)}=1 | t^{(i)}=1, x^{(i)}) \\ &= p(t^{(i)}=1 | x^{(i)}) \quad \text{Q.E.D.} \\ \frac{1}{\alpha} p(y^{(i)}=1 | x^{(i)}) &\geq p(t^{(i)}=1 | x^{(i)}) \end{aligned}$$

e) Let  $g(x) = p(f(i)=1 | x^{(i)})$

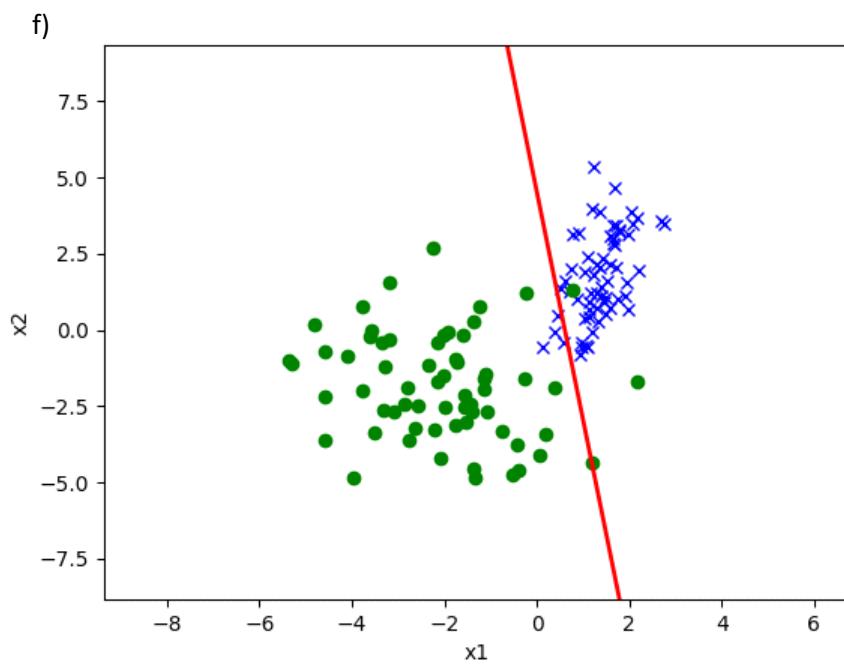
From d,  $h(x^{(i)}) = \alpha g(x^{(i)})$

$$\begin{aligned} \mathbb{E}[h(x^{(i)}) | y^{(i)}=1] &= \mathbb{E}(h(x^{(i)}))_{y^{(i)}=1} \\ &= \frac{\mathbb{E}[h(x^{(i)}) \mathbf{1}\{y^{(i)}=1\}]}{p(y^{(i)}=1)} \\ &= \frac{\mathbb{E}[\mathbb{E}[\mathbf{1}\{y^{(i)}=1\} | x^{(i)}] \cdot h(x^{(i)})]}{\mathbb{E}(p(y^{(i)}=1 | x^{(i)}))} \\ &= \frac{\mathbb{E}[h(x^{(i)}) \cdot h(x^{(i)})]}{\mathbb{E}(h(x^{(i)}))} \\ &= \alpha \frac{\mathbb{E}[g(x^{(i)}) \cdot g(x^{(i)})]}{\mathbb{E}[g(x^{(i)})]} \\ &= \alpha \end{aligned}$$

 $\geq \alpha$

# PS1 - 2f

Saturday, July 9, 2022 1:56 AM



$$\textcircled{3} a) p(y; \lambda) = \frac{e^{-\lambda} \lambda^y}{y!} = \frac{1}{y!} \exp(y \log \lambda - \lambda)$$

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

$$n = \log \lambda$$

$$T(y) = t$$

$$a(n) = e^n$$

$$b) g(n) = E[y; n] = \lambda = e^n$$

$$c) l(\theta) = \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}, \theta)$$

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

$$\begin{aligned} \frac{\partial l(\theta)}{\partial \theta_j} &= \sum_{i=1}^m y^{(i)} x_j^{(i)} - x_j^{(i)} g(\theta^T x^{(i)}) \\ &= \sum_{i=1}^m (y^{(i)} - g(\theta^T x^{(i)})) x_j^{(i)} \end{aligned}$$

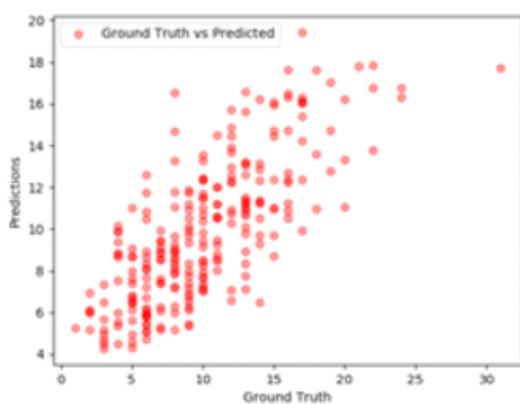
The update rule is  $\theta_j := \theta_j + \alpha (y^{(i)} - g(\theta^T x^{(i)})) x_j^{(i)}$

$$\text{where } g(n) = e^n$$

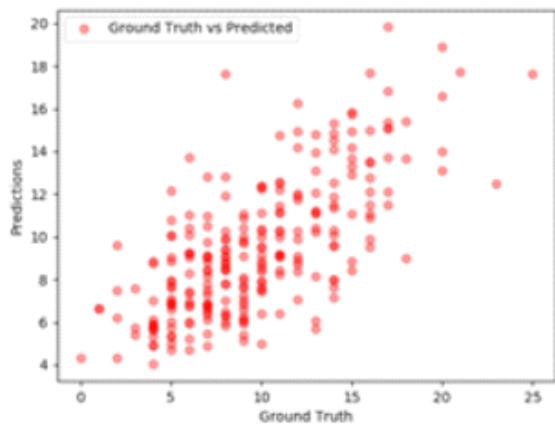
# PS1 - 3d

Saturday, July 9, 2022 2:11 AM

d)



Test set



Validation set

$$(4) \text{ a) } \frac{\partial}{\partial n} p(y; n) = p(y; n) y - \frac{\partial}{\partial n} a(n)$$

$$\int \frac{\partial}{\partial n} p(y; n) dy = \int p(y; n) y - \frac{\partial}{\partial n} a(n) dy$$

$$\frac{\partial}{\partial n} \int p(y; n) dy = \int y p(y; n) dy - \int p(y; n) \frac{\partial}{\partial n} a(n) dy$$

$$\frac{\partial}{\partial n} 1 = \mathbb{E}[y; n] - \frac{\partial}{\partial n} a(n) \int p(y; n) dy$$

$$\mathbb{E}[y; n] = \frac{\partial}{\partial n} a(n)$$

$$\begin{aligned}
 b) \frac{\partial^2}{\partial \eta^2} a(n) &= \frac{\partial}{\partial n} \mathbb{E}[y|x;\theta] = \frac{\partial}{\partial n} \int y p(y;n) dy \\
 &= \int y p(y;n) y - \frac{\partial}{\partial n} a(n) dy \\
 &= \int y^2 p(y;n) dy - \frac{\partial}{\partial n} a(n) \int y p(y;n) dy \\
 &= \mathbb{E}[y^2;n] - \mathbb{E}[y;n]^2 \\
 &= \text{Var}(y;n)
 \end{aligned}$$

$$\begin{aligned}
 c) \quad J(\theta) &= -\log [p(y; \theta)] \\
 &= -\log [b(y) \exp(ny - a(\theta))] \\
 &= a(\theta) - ny + C \\
 &= a(\theta^T x) - \theta^T x + C
 \end{aligned}$$

$$\nabla \phi(\theta) = \frac{\partial}{\partial \theta} \alpha(\theta) \nabla \theta^n - y^*$$

$$= \frac{\partial}{\partial n} a(n)x - yx$$

$$\nabla_{\theta'}^2(\theta) = \nabla_\theta(\nabla_\theta'(\theta))$$

$$= \nabla_{\theta} \frac{\partial}{\partial \theta} a(\theta) x - y x$$

$$= x \nabla \phi \frac{\partial}{\partial n} a(n)$$

$$= x \frac{\partial}{\partial n} \frac{\partial}{\partial n} a(n) \nabla^2$$

$$= \hat{\partial}^2_{\alpha\beta} a^{(\eta)} x x^T$$

$$-\frac{\partial \mathcal{H}}{\partial p_i} \quad \text{and} \quad T$$

$$\approx \text{Var}(\gamma; n) \times \times$$

For any  $\theta$ ,  $\text{Var}(\gamma; \theta)$  is positive, and the gradient  $\nabla_{\theta} \ell_{\theta}(\gamma) = \text{Var}(\gamma; \theta) \times x^T$  of the GLM NLL loss is PSD.

and convex

$$\textcircled{5} \text{ a) } J(\theta) = \frac{1}{2} \sum_{i=1}^N (\theta^T \hat{x}^{(i)} - y^{(i)})^2$$

$$\nabla_{\theta} J(\theta) = \sum_{i=1}^N (\theta^T \hat{x}^{(i)} - y^{(i)}) \hat{x}^{(i)}$$

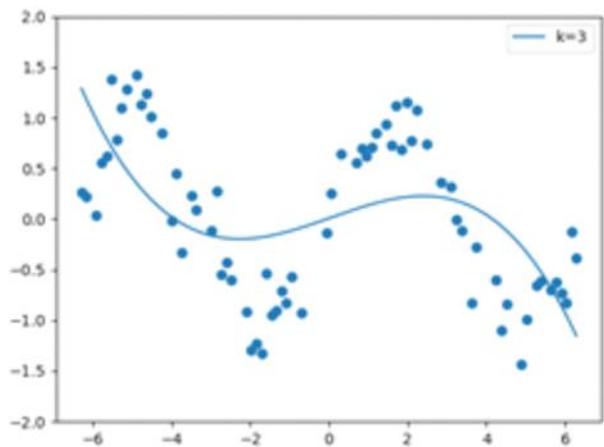
$$\alpha \lambda \sum_{i=1}^N (\theta^T \hat{x}^{(i)} - y^{(i)}) \hat{x}^{(i)}$$

$$\text{Update rule: } \theta := \theta - \lambda \sum_{i=1}^N (\theta^T \hat{x}^{(i)} - y^{(i)}) (\hat{x}^{(i)})$$

# PS1 - 5b

Saturday, July 9, 2022 3:48 PM

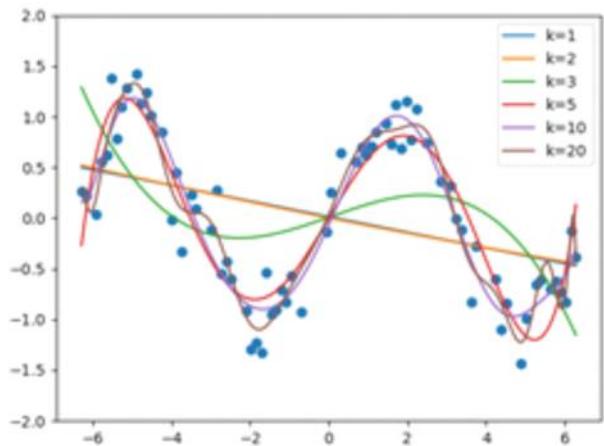
b)



# PS1 - 5c

Saturday, July 9, 2022 3:49 PM

c)

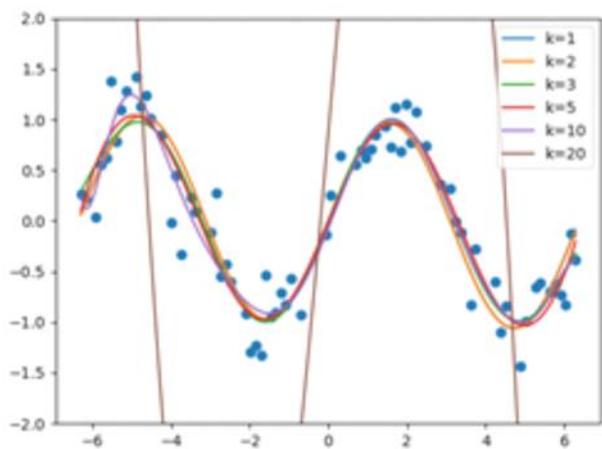


The higher degree of the polynomial, the closer the fit to the data. However, higher degree polynomials are less smooth and less simple by nature of having more parameters

# PS1 - 5d

Saturday, July 9, 2022 3:51 PM

d)

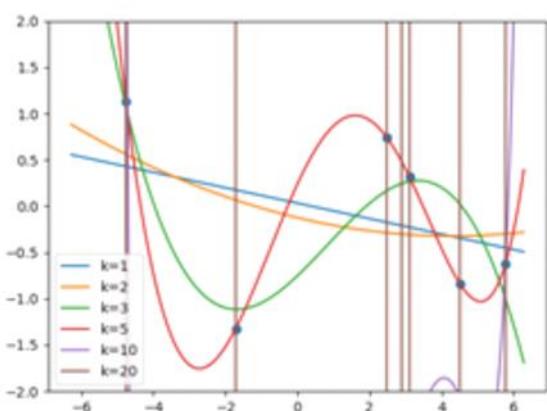


With the  $\sin(x)$  feature, the polynomials seem to fit the data closer, with the exception of the  $k=20$  polynomial

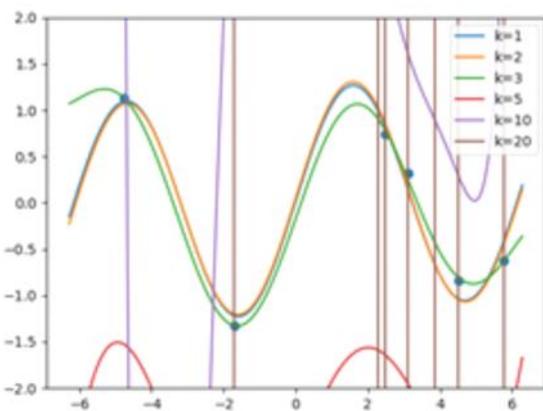
# PS1 - 5e

Saturday, July 9, 2022 3:53 PM

e)



Polynomial



With  $\sin(x)$  feature

With limited data, the polynomial approximations can easily pass through all the data points but vary greatly in their shape. The expanded feature version seems to fit better. However, in both cases as the degree becomes too large there is overfitting.