

Answer 1:

Part a

We will start by taking the first derivative of $J(\theta)$

$$\begin{aligned}\nabla_{\theta} J(\theta) &= -\frac{1}{n} \sum_{i=1}^n \left[y^{(i)} \frac{g'(\theta^T x^{(i)})}{g(\theta^T x^{(i)})} \cdot x^{(i)} + (1 - y^{(i)}) \cdot \frac{-g'(\theta^T x^{(i)})}{1 - g(\theta^T x^{(i)})} \cdot (1 - g(\theta^T x^{(i)})) \cdot x^{(i)} \right] \\ &= -\frac{1}{n} \sum_{i=1}^n \left[y^{(i)} \frac{g(\theta^T x^{(i)})}{g(\theta^T x^{(i)})} \cdot (1 - g(\theta^T x^{(i)})) x^{(i)} + (1 - y^{(i)}) \cdot \frac{1 - g(\theta^T x^{(i)})}{1 - g(\theta^T x^{(i)})} \cdot g(\theta^T x^{(i)}) x^{(i)} \right] \\ &= -\frac{1}{n} \sum_{i=1}^n \left[y^{(i)} - y^{(i)} g(\theta^T x^{(i)}) - g(\theta^T x^{(i)}) + y^{(i)} g(\theta^T x^{(i)}) \right] \cdot x^{(i)} \\ &= -\frac{1}{n} \sum_{i=1}^n \left[y^{(i)} - g(\theta^T x^{(i)}) \right] \cdot x^{(i)}\end{aligned}$$

Then we take the Hessian of the loss function

$$\begin{aligned}H = \nabla_{\theta}^2 J(\theta) &= \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} g(\theta^T x^{(i)}) \cdot \nabla_{\theta} (\theta^T x^{(i)}) \cdot x^{(i)} \\ &= \frac{1}{n} \sum_{i=1}^n g(\theta^T x^{(i)}) (1 - g(\theta^T x^{(i)})) x^{(i)} [x^{(i)}]^T\end{aligned}$$

Given that $h_{\theta}(x) = g(\theta^T x)$, for any vector $z \in R^d$, define $x \in R^n$

$$\begin{aligned}z^T (Hz) &= z^T \cdot \begin{pmatrix} z^{(1)} \cdot H \\ z^{(2)} \\ \vdots \\ z^{(d)} \cdot H \end{pmatrix} \\ &= \sum_{j=1}^d [z^{(j)}]^T z^{(j)} \cdot \frac{1}{n} \sum_{i=1}^n \left[x^{(i)} [x^{(i)}]^T g(\theta^T x^{(i)}) (1 - g(\theta^T x^{(i)})) \right] \\ &= \sum_j [z^{(j)}]^T z^{(j)} x^T x \cdot \frac{1}{n} \sum_{i=1}^n \left[h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) \right]\end{aligned}$$

Since $x^T x$ is a scalar, define $z \in R^d$, $\sum_j [z^{(j)}]^T z^{(j)} x^T x$ can be written as

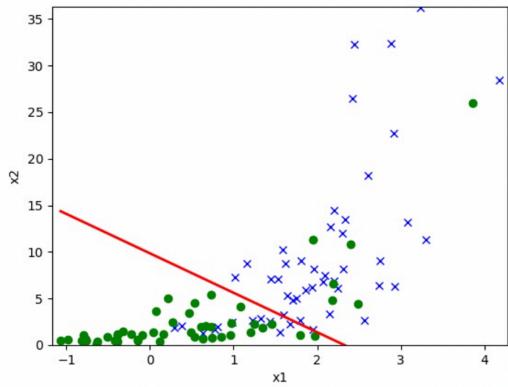
$$(x^T x) \cdot (z^T z) = z^T (x^T x) z = (x^T z)^T (x^T z) = (x^T z)^2$$

So, the quadratic form becomes

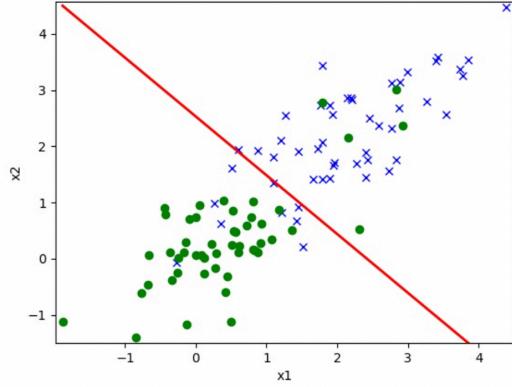
$$z^T (Hz) = (x^T z)^2 \cdot \frac{1}{n} \sum_{i=1}^n \left[h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) \right]$$

Observe that all the terms are positive because $h(\theta)$ is a possibility function and $(x^T z)^2 \geq 0$. We conclude that matrix H is positive semi-definite.

part b :



Dataset 1



Dataset 2

Logistic Regression Results given two datasets.

Part C.

Simplify $P(Y=1 | X=x; \phi, w_0, w_1, \Sigma)$

$$\begin{aligned} &= \frac{P(X=x | Y=1) P(Y=1)}{P(X=x | Y=1) P(Y=1) + P(X=x | Y=0) P(Y=0)} \\ &= \frac{1}{1 + \frac{P(X=x | Y=0) P(Y=0)}{P(X=x | Y=1) P(Y=1)}} \end{aligned}$$

Take a closer look at the function in the denominator:

$$\begin{aligned} \frac{P(X=x | Y=0) P(Y=0)}{P(X=x | Y=1) P(Y=1)} &= \frac{\exp(-\frac{1}{2}(x-w_0)^T \Sigma^{-1} (x-w_0))}{\exp(-\frac{1}{2}(x-w_1)^T \Sigma^{-1} (x-w_1))} \cdot \frac{1-\phi}{\phi} \\ &= \frac{1-\phi}{\phi} \exp\left[\frac{1}{2}(x-w_1)^T \Sigma^{-1} (x-w_1) - \frac{1}{2}(x-w_0)^T \Sigma^{-1} (x-w_0)\right] \end{aligned}$$

as both the nominator and denominator share the common constant factor $[(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}]^{-1}$

As $x - \mu_0 = x - \mu_1 + \mu_1 - \mu_0$

The exponent term $\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1) - \frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0)$ can be written as

$$\begin{aligned} & \frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (\cancel{x - \mu_1} + \mu_1 - \mu_0) + \frac{1}{2}(\mu_1 - \mu_0)^T \Sigma^{-1} (x - \mu_1 + \mu_1 - \mu_0) - \frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (\cancel{x - \mu_1}) \\ &= \frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (\mu_1 - \mu_0) + \frac{1}{2}(\mu_1 - \mu_0)^T \Sigma^{-1} (x - \mu_0) \\ &= (\mu_1 - \mu_0)^T \Sigma^{-1} \left(x - \frac{\mu_0 + \mu_1}{2} \right) \end{aligned}$$

Hence, the above fraction equals

$$\exp \left[(\mu_1 - \mu_0)^T \Sigma^{-1} x + \log \left(\frac{1-\phi}{\phi} \right) - \frac{(\mu_1 - \mu_0)^T \Sigma^{-1} (\mu_0 + \mu_1)}{2} \right]$$

$$\text{if we let } \theta^T = (\mu_1 - \mu_0)^T \Sigma^{-1} \Rightarrow \theta = \left[(\mu_1 - \mu_0)^T \Sigma^{-1} \right]^T = \left[\Sigma^{-1} \right]^T (\mu_1 - \mu_0)$$

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

If we set the threshold at 0.5, the decision boundary will have the equation $\theta^T x + \theta_0 = 0$

part d.

$$\begin{aligned} \text{log-likelihood function: } \ell(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^n p(x^{(i)} | y^{(i)}) p(y^{(i)}) \\ &= \sum_{i=1}^n [\log p(x^{(i)} | y^{(i)}) + \log p(y^{(i)})] \\ &= \sum_{i=1}^n \left[-\log(2\pi)^{d/2} - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x^{(i)} - \mu y^{(i)})^T \Sigma^{-1} (x^{(i)} - \mu y^{(i)}) \right. \\ &\quad \left. + y^{(i)} \log \phi + (1-y^{(i)}) \log (1-\phi) \right] \end{aligned}$$

$\hat{\phi}$ is achieved by $\nabla_\phi \ell(\phi, \mu_0, \mu_1, \Sigma) = 0$,

$$\nabla_\phi \ell(\phi, \mu_0, \mu_1, \Sigma) = \sum_{i=1}^n \left[\frac{y^{(i)}}{\phi} - \frac{1-y^{(i)}}{1-\phi} \right] = 0.$$

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

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

$$\begin{aligned}
\nabla_{w_0} \ell(\theta) &= \nabla_{w_0} \sum_{i=1}^n \left[-\frac{1}{2} \{y^{(i)} = 0\} (x^{(i)} - w_0)^T \Sigma^{-1} (x^{(i)} - w_0) \right] \\
&= \sum_{i=1}^n \left[-\Sigma^{-1} (x^{(i)} - w_0) \times (-1) \right] \quad (\nabla x^T A x = 2Ax) \\
&= \Sigma^{-1} \sum_{i=1}^n \{y^{(i)} = 0\} (x^{(i)} - w_0)
\end{aligned}$$

Let $\nabla_{w_0} \ell(\theta) = 0$ we obtain

$$\begin{aligned}
\Sigma^{-1} \cdot \sum_{i=1}^n x^{(i)} &= \Sigma^{-1} \cdot \sum_{i=1}^n \hat{w}_0 \\
\hat{w}_0 &= \frac{\sum_{i=1}^n \{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^n \{y^{(i)} = 0\}}
\end{aligned}$$

Similarly, taking the gradients w.r.t. w_1 , we obtain

$$\nabla_{w_1} \ell(\theta) = \Sigma^{-1} \sum_{i=1}^n \{y^{(i)} = 1\} (x^{(i)} - w_1)$$

Let $\nabla_{w_1} \ell(\theta) = 0$, there is

$$\hat{w}_1 = \frac{\sum_{i=1}^n \{y^{(i)} = 1\} x^{(i)}}{\sum_{i=1}^n \{y^{(i)} = 0\}}$$

Claim that for any vector $z \in \mathbb{R}^d$, $\nabla_{\Sigma} z^T \Sigma^{-1} z = -\Sigma^{-1} z z^T \Sigma^{-1}$

$$\begin{aligned}
\frac{d}{d\Sigma} (z^T \Sigma^{-1} z) &= \frac{d}{d\Sigma} \text{Tr}[z^T \Sigma^{-1} z] = \text{Tr}[z z^T \frac{d\Sigma^{-1}}{d\Sigma}] = \text{Tr}[-z z^T \Sigma^{-1} (\Sigma)' \Sigma^{-1}] \\
&= \text{Tr}[-\Sigma^{-1} z z^T \Sigma^{-1} (\Sigma)']
\end{aligned}$$

Given $\Sigma \Sigma^{-1} = I$, we have $\frac{d}{d\Sigma} (\Sigma^{-1}) = -\Sigma^{-1} (\Sigma)' \Sigma^{-1}$

$$\text{Hence } \nabla_{\Sigma} (z^T \Sigma^{-1} z) = -\Sigma^{-1} z z^T \Sigma^{-1}$$

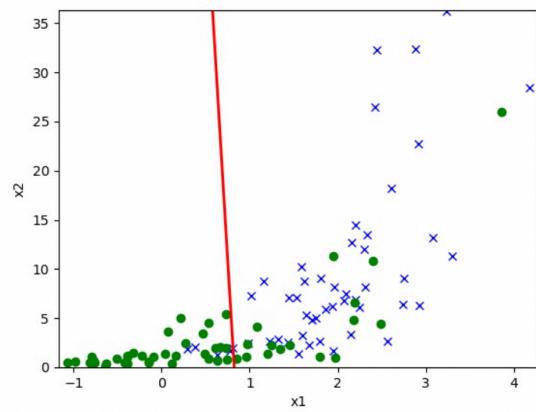
Let $z = (x^{(i)} - w_0 y^{(i)})$ as $x^{(i)} - w_0 y^{(i)}$ is a vector

$$\begin{aligned}
\nabla_{\Sigma} \ell(\theta) &= \nabla_{\Sigma} \sum_{i=1}^n \left(-\frac{1}{2} \log |\Sigma| - \frac{1}{2} (x^{(i)} - w_0 y^{(i)})^T \Sigma^{-1} (x^{(i)} - w_0 y^{(i)}) \right) \\
&= \sum_{i=1}^n \left(-\frac{1}{2} \Sigma^{-1} + \frac{1}{2} (x^{(i)} - w_0 y^{(i)}) (x^{(i)} - w_0 y^{(i)})^T \Sigma^{-1} \right) \quad [\nabla_{\Sigma} \log |\Sigma| = (\Sigma)^{-1}] \\
&= 0.
\end{aligned}$$

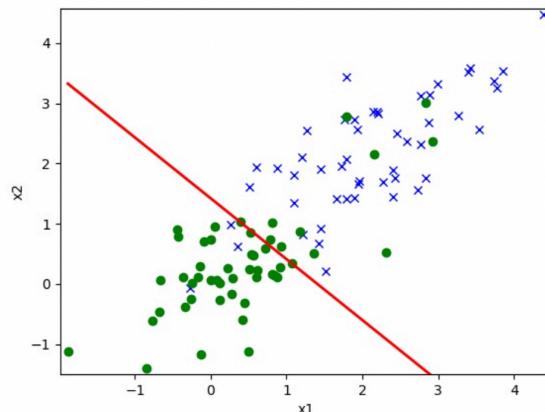
So, we have $n - \sum_{i=1}^n (x^{(i)} - w_0 y^{(i)}) (x^{(i)} - w_0 y^{(i)})^T \Sigma^{-1} = 0$

$$\Rightarrow \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - w_0 y^{(i)}) (x^{(i)} - w_0 y^{(i)})^T$$

part f:



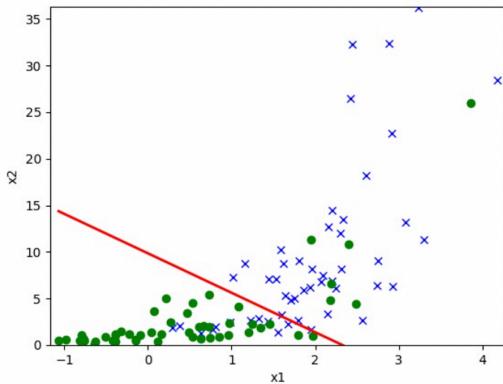
Dataset 1



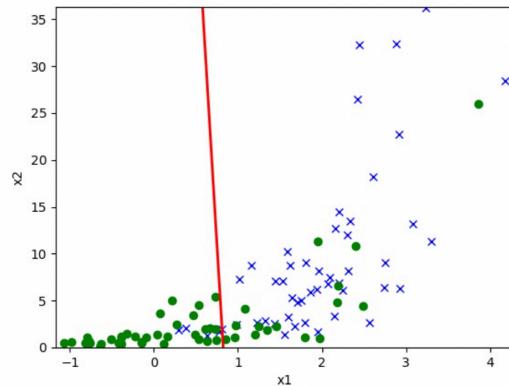
Dataset 2

GDA Results on two given datasets

part f:



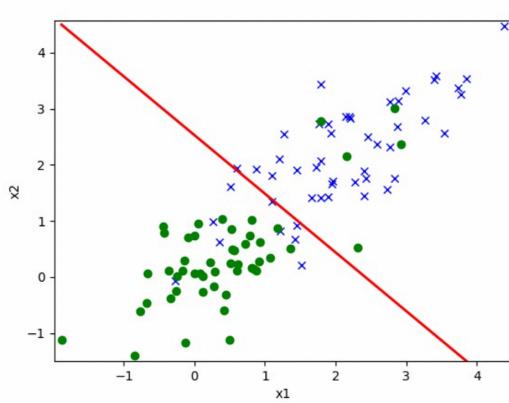
Logistic Regression



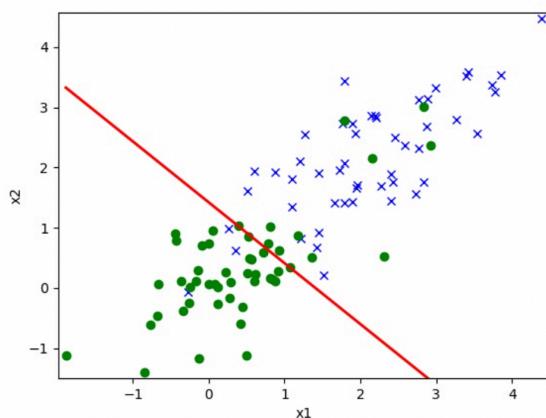
GDA .

for dataset 1, Logistic Regression obtains a better decision boundary than GDA , so the classification model using logistic regression yields more accurate results:

part g



logistic regression



GDA .

GDA performance in validation set 2 is similar to the logistic regression. But in dataset 1, GDA's performance is worse than Logistic Regression.

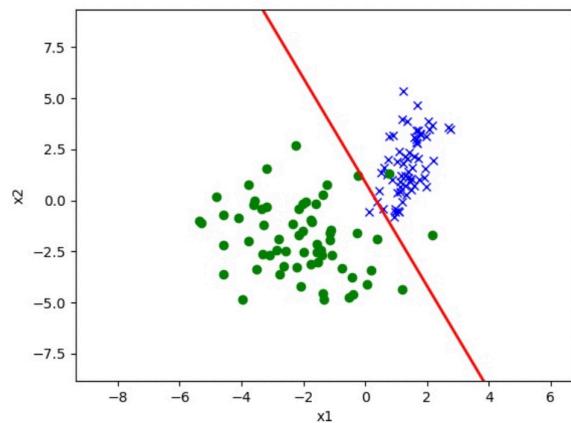
GDA makes more specific assumptions about the dataset than LR.
[e.g. Gaussian Dist] when Dataset is gaussian dist, GDA better.

part h.

We can apply the power transform by taking the square root and the log. to $x^{(i)}$, so that the transformed data is gaussian distributed.

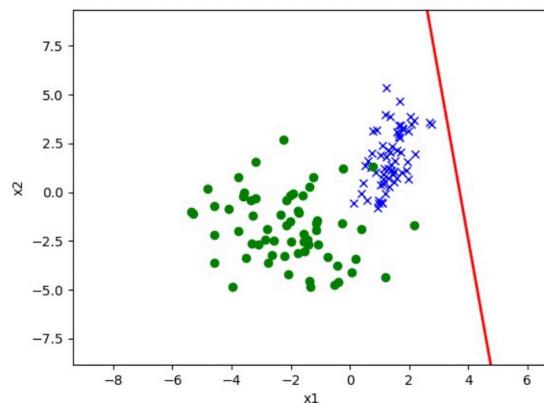
Question 2:

Part a:



fully observed result on the test data

part b:



naive method on partial labels

Part C

since $\#x^{(i)}$, $P(y^{(i)}=1 | t^{(i)}=0, x) = 0$

$$P(t^{(i)}=1 | y^{(i)}=1, x^{(i)}) = \frac{P(y^{(i)}=1 | t^{(i)}=1, x) P(t^{(i)}=1 | x)}{P(y^{(i)}=1 | t^{(i)}=1, x) P(t^{(i)}=1 | x) + P(y^{(i)}=1 | t^{(i)}=0, x) P(t^{(i)}=0 | x)}$$

$$= 1$$

(2d)

$$P(y^{(i)}=1 | x^{(i)}) = P(y^{(i)}=1, x^{(i)}, t^{(i)}=1) / P(y^{(i)}=1, x^{(i)}, t^{(i)}=0)$$

$$= P(y^{(i)}=1 | x^{(i)}, t^{(i)}=1) [P(t^{(i)}=1 | x^{(i)})]$$

$$+ P(y^{(i)}=1 | t^{(i)}=0, x^{(i)}) P(t^{(i)}=0 | x^{(i)})$$

$$= \alpha P(t^{(i)}=1 | x^{(i)}) + 0.$$

$$\therefore P(t^{(i)}=1 | x^{(i)}) = \frac{1}{\alpha} P(y^{(i)}=1 | x^{(i)})$$

Part C

We will show that $h(x^{(i)}) = \alpha$ when $y^{(i)}=1$, and $h(x^{(i)}) = 0$ when $y^{(i)}=0$.

$h(x^{(i)}) =$

$P(y^{(i)}=1 | x^{(i)})$ from ad, we obtain that

$$P(y^{(i)}=1 | x^{(i)}) = \alpha \cdot P(t^{(i)}=1 | x^{(i)}) + 0 \cdot P(t^{(i)}=0 | x^{(i)})$$

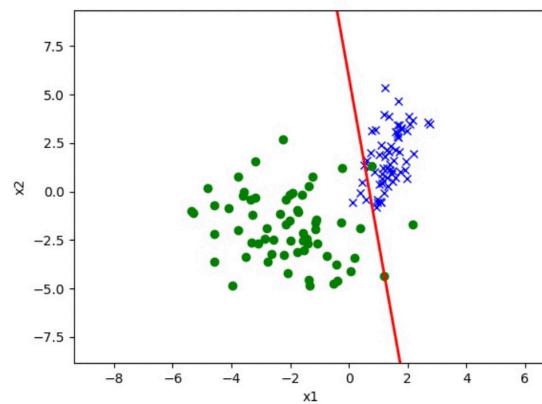
Given the assumption $P(t^{(i)}=1 | x^{(i)}) = 1$ and $P(t^{(i)}=0 | x^{(i)}) = 0$.

$$h(x^{(i)}) = \alpha \cdot 1 + 0$$

Hence, given $y^{(i)}=1$, $h(x^{(i)})$ always take the value of α , \Rightarrow expected value = α .

Therefore, $E(h(x^{(i)}) | y^{(i)}=1) = \alpha$

part f



Rescaled result on the test set.

Question 3

Part a:

We can write the Poisson Distribution in the form

$$P(y; \eta) = b(y) \exp(\eta T(y) - a(\eta)) \quad \text{s.t.}$$

$$P(y; \eta) = \frac{1}{y!} \exp(\log(\lambda) \cdot y - e^n)$$

$$\text{where } \eta = \log(\lambda), T(y) = y, a(\eta) = \lambda = e^{\log \lambda} = e^n, b(y) = \frac{1}{y!}$$

Part b.

Because $E[y; \eta] = \lambda = e^n$, so the canonical response function is

$$u = g(\eta) = e^n$$

Part c

$\ell(\theta) = \log P(y^{(i)} | x^{(i)}, \theta)$, picking the j th $x^{(i)}$, $y^{(i)}$ and given that $\theta^T x = y$

$$\frac{\partial}{\partial \theta_j} \ell(\theta) = \frac{\partial}{\partial \theta_j} \log \left[\frac{1}{y!} \exp(\eta \cdot y - e^n) \right] = \frac{\partial}{\partial \theta_j} \log \left[\frac{1}{y!} \exp(y \theta^T x^{(i)} - e^{\theta^T x^{(i)}}) \right]$$

$$= \frac{\partial}{\partial \theta_j} \left[\log \frac{1}{y!} + y \cdot \theta^T x^{(i)} - e^{\theta^T x^{(i)}} \right]$$

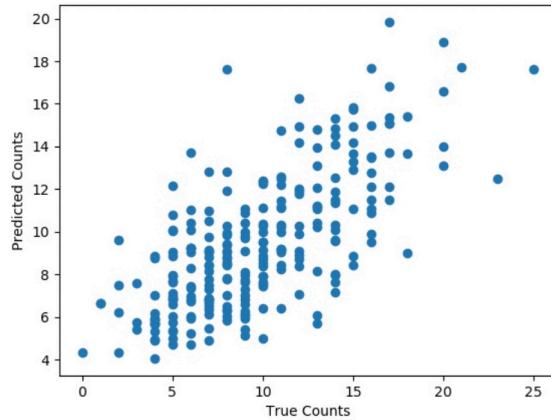
$$= y x^{(i)j} - e^{\theta^T x^{(i)}} \cdot x^{(i)j}$$

$$= (y - e^{\theta^T x^{(i)}}) x^{(i)j}$$

Hence, the SGD update rule for the GLM is that

$$\theta_{j+1} = \theta_j + \alpha (y - \theta^T x^{(i)}) x^{(i)j} \quad \text{for any } j$$

Part d.



Learned Poisson result on the validation set

4. Note that (1) $\int_{-\infty}^{\infty} p(y; \eta) dy = 1$ and (2) $\int_{-\infty}^{\infty} \frac{\partial^k}{\partial \eta^k} \cdot p(y; \eta) dy = 0$. for $k=1, 2, \dots$

Hence we have

$$\int_{-\infty}^{\infty} p'(y; \eta) dy = 0 \quad \text{and} \quad \int p''(y; \eta) dy = 0$$

(3) Let $\ell = \log p(y; \eta)$, then

$$\ell' = \frac{p'}{p} \quad \text{and} \quad \ell'' = -\left(\frac{p'}{p}\right)^2$$

$$\begin{aligned} (4) \quad E(\ell') &= \int \ell' p(y; \eta) dy = \int \frac{p'(y; \eta)}{p(y; \eta)} \cdot p(y; \eta) dy \quad (\text{By (3)}) \\ &= \int p'(y; \eta) dy \quad (\text{By (2)}) \\ &= 0 \end{aligned}$$

$$(5) \quad E(\ell'') = \int p''(y; \eta) dy - E\left[\left(\frac{p'(y; \eta)}{p(y; \eta)}\right)^2\right]$$

$$= -E[(\ell')^2]$$

$$= -E[(\ell')^2] - [E(\ell')]^2$$

$$= -\text{Var}(\ell')$$

$$(4b) \quad \ell''(\eta) = -a''(\eta)$$

$$E(\ell'') = -\text{Var}(\ell')$$

(4a) Because $\eta = \eta^T$, $\eta \in \mathbb{R}$

$$\ell(\eta) = \log p(y; \eta) = \log b(y) + \eta^T a(\eta) - a(\eta)$$

$$\text{so, } \text{Var}(T(\eta) - a'(\eta))$$

$$\ell'(\eta) = T(\eta) - a'(\eta)$$

$$= -E(-a''(\eta)) \text{ const.}$$

$$E(\ell') = E(T(\eta) - a'(\eta)) = E(T(\eta)) - a'(\eta)$$

$$= 0$$

$$\therefore E(T(\eta)) = a'(\eta)$$

4C

The NLL of the distribution is written as =

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^n \log p(y^{(i)}; \theta) = -\frac{1}{n} \sum_{i=1}^n \left(\log(b(y^{(i)})) + y^{(i)} \theta^T x^{(i)} - a(\theta^T x^{(i)}) \right)$$

Take the gradient w.r.t. θ :

$$\nabla_{\theta} J(\theta) = -\frac{1}{n} \sum_{i=1}^n (y^{(i)} x^{(i)} - a'(\theta^T x^{(i)}) \cdot x^{(i)})$$

Then take the hessian : $H_{J(\theta)}$

$$H_{J(\theta)} = -\frac{1}{n} \sum_{i=1}^n (-a''(\theta^T x^{(i)}) x^{(i)} [x^{(i)}]^T)$$

#

we already know that $a''(\theta^T x^{(i)}) = a''(\eta) = \text{Var}(y^{(i)} | x^{(i)})$

$$\text{Hence, } H_{J(\theta)} = \frac{1}{n} \sum_{i=1}^n \text{Var}(y^{(i)} | x^{(i)}) x^{(i)} [x^{(i)}]^T$$

and $E(\ell''(\eta)) = \text{Var}(\ell')$
 $= \text{Var}(\gamma)$ by part b.
 $= \text{Var}(\gamma | X; \theta)$

consider any vector $a \in \mathbb{R}^d$,

$$\begin{aligned} a^T H_{J(\theta)} a &= a^T \frac{1}{n} \sum_{i=1}^n \text{Var}(y^{(i)} | x^{(i)}) x^{(i)} [x^{(i)}]^T \cdot a \\ &= \frac{1}{n} \sum_{i=1}^n \text{Var}(y^{(i)} | x^{(i)}) [a^T x^{(i)}]^2 \end{aligned}$$

since $\text{Var}(y^{(i)} | x^{(i)}) \geq 0$ for all i and $[a^T x^{(i)}]^2$ is always positive.

$a^T H_{J(\theta)} a \geq 0$ i.e. The hessian of the loss is positive semi-definite.

which concludes that NLL loss of GLM is convex.

Question 5

- (5a) The objective function is written as
- $$J(\theta) = \sum_{i=1}^n (\theta^T \hat{x}^{(i)} - y^{(i)})^2$$

And the update gradient descent equation is

$$\underline{\theta_{t+1}} = \underline{\theta_t}$$

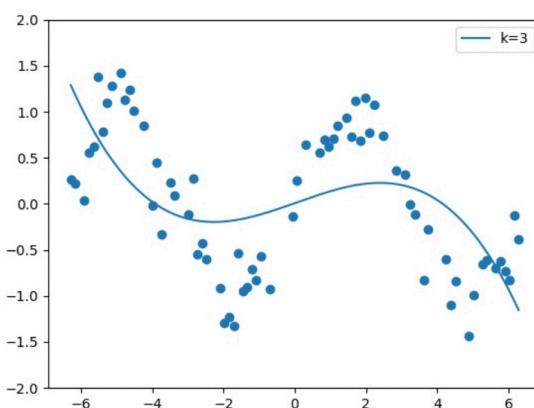
$$\theta^{t+1} = \theta^t - \alpha' \nabla_{\theta} J(\theta)$$

$$= \theta^t - \alpha \sum_{i=1}^n 2(\theta^T \hat{x}^{(i)} - y^{(i)}) \hat{x}^{(i)}$$

for some α , there is the update rule =

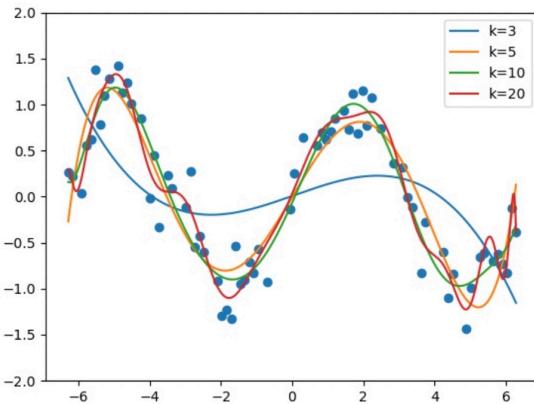
$$\theta^{t+1} = \theta^t - \alpha \sum_{i=1}^n (\theta^T \hat{x}^{(i)} - y^{(i)}) \hat{x}^{(i)}$$

part b :



Degree-3 Polynomial Regression

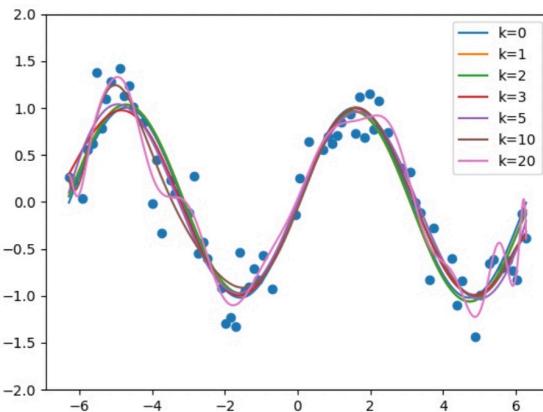
Part C:



Degree- k Polynomial Regression

The model fits the dataset better as k increases. However at $k=20$ onwards the model starts to overfit the data.

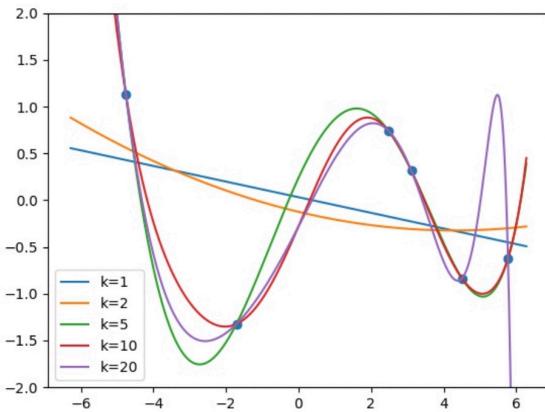
Part d:



Degree- k Polynomial Regression with Sine-feature-map.

The model fits the training data perfectly at lower number of k . As k increases, the model tends to overfit the dataset.

Part e:



Overfitting with small data

Since there are only 6 sets of data in this case. The model fits the data into sine feature map when $k \geq 5$. However, due to the limited amount of training data, when the degree increases, the model does not conform to the sine shape but some strange and asymmetric shape.
It shows the result of overfitting.