

Stanford CS229: Problem Set #1, math parts

1) Linear Classifiers (Logistic Regression, GDA)

(a) The cost function for logistic regression was

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))$$

recalling $h_\theta(x) = \frac{1}{1+\exp(-\theta^T x)}$ & that
 $(x^{(i)}, y^{(i)})$ denote the training examples.

- Find the Hessian of J , H .
- Show H is positive semidefinite.

↳ This implies that J is convex \Rightarrow local extrema are global extrema!
 Note that H is already symmetric.

Start w/ gradient. (Hint: it's the same as update rule for lin reg.)

$$\frac{\partial J}{\partial \theta_k} = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(h_\theta(x^{(i)})) \right] + \left[(1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right]$$

Break into pieces.

$$y^{(i)} \cdot \frac{\partial J}{\partial \theta_k} \left[\log(h_\theta(x^{(i)})) \right] = y^{(i)} \left(\frac{1}{h_\theta(x^{(i)})} \right) \cdot \frac{\partial h_\theta(x^{(i)})}{\partial \theta_k}$$

$$\begin{aligned} \frac{\partial h_\theta(x^{(i)})}{\partial \theta_k} &= \frac{\partial}{\partial \theta_k} \left[\frac{1}{1+\exp(-\theta^T x)} \right] \\ &= -\frac{\frac{2}{\partial \theta_k} \left[1+e^{-\theta^T x^{(i)}} \right]}{(1+e^{-\theta^T x^{(i)}})^2} = \frac{-e^{-\theta^T x^{(i)}} \cdot \frac{\partial}{\partial \theta_k} (\theta^T x^{(i)})}{(1+e^{-\theta^T x^{(i)}})^2} \end{aligned}$$

$$= \frac{e^{-\theta^T x^{(i)}} \cdot x_k^{(i)}}{(1+e^{-\theta^T x^{(i)}})^2}$$

$$\dots = y^{(i)} \left(\frac{1 + e^{-\theta^T x^{(i)}}}{1} \right) \cdot \left[\frac{e^{-\theta^T x^{(i)}}}{(1 + e^{-\theta^T x^{(i)}})^2} \chi_k^{(i)} \right]$$

$$= y^{(i)} \cdot \frac{e^{-\theta^T x^{(i)}} \chi_k^{(i)}}{(1 + \exp(-\theta^T x^{(i)}))}$$

$$(1-y^{(i)}) \frac{\partial J}{\partial \theta_k} \left[\log(1 - h_{\theta}(x^{(i)})) \right]$$

$$= (-) \left(\frac{1}{1 - h_{\theta}(x^{(i)})} \right) \cdot \frac{\partial}{\partial \theta_k} \left[-h_{\theta}(x^{(i)}) \right]$$

$$= (-) \frac{1}{1 - h_{\theta}(x^{(i)})} \cdot \frac{\partial}{\partial \theta_k} \left[-h_{\theta}(x^{(i)}) \right]$$

$$= -(-) \left(\frac{1 + e^{-\theta^T x^{(i)}}}{e^{-\theta^T x^{(i)}}} \right) \frac{e^{-\theta^T x^{(i)}} \cdot \chi_k^{(i)}}{(1 + e^{-\theta^T x^{(i)}})^2} = -(1-y^{(i)}) h_{\theta}(x^{(i)}) \cdot \chi_k^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \cdot \frac{e^{-\theta^T x^{(i)}} \chi_k^{(i)}}{(1 + \exp(-\theta^T x^{(i)}))} - h_{\theta}(x^{(i)}) \chi_k^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m y^{(i)} (1 - h_{\theta}(x^{(i)}) \chi_k^{(i)}) - (1-y^{(i)}) h_{\theta}(x^{(i)}) \chi_k^{(i)}$$

$$\chi_k^{(i)} [y^{(i)} - y^{(i)} h_{\theta}(x^{(i)}) - h_{\theta}(x^{(i)}) + y^{(i)} h_{\theta}(x^{(i)})]$$

$$\frac{\partial J}{\partial \theta_k} = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} - h_{\theta}(x^{(i)})] \chi_k^{(i)}$$

$$\frac{\partial^2 J}{\partial \theta_k^2} = -\frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \theta_k} [y^{(i)} - h_{\theta}(x^{(i)})] \chi_k^{(i)}$$

2.1) derivatives

cancel $\rightarrow \frac{\partial}{\partial \theta_k} [h_{\theta}(x^{(i)})]$

$$= \frac{1}{m} \sum_{i=1}^m \frac{e^{-\theta^T x^{(i)}} \cdot (\chi_k^{(i)})^2}{(1 + e^{-\theta^T x^{(i)}})^2} = \frac{1}{m} \sum_{i=1}^m h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) (\chi_k^{(i)})^2$$

Similarly,

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

So the Hessian!

$$H_{kj} = \begin{cases} \frac{1}{m} \sum_{i=1}^m h_\theta(x^{(i)}) (1 - h_\theta(x^{(i)})) (x_k^{(i)} x_j^{(i)})^2 & \text{when } k=j \\ \frac{1}{m} \sum_{i=1}^m h_\theta(x^{(i)}) (1 - h_\theta(x^{(i)})) (x_k^{(i)} x_j^{(i)}) & \text{when } k \neq j \end{cases}$$

Or even more succinctly, recalling

$$g'(z) = g(z) (1 - g(z)), \quad \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} [x_1 \dots x_n]$$

$$H_{kj} = \frac{1}{m} \sum_{i=1}^m g'(\theta^\top x^{(i)}) x_k^{(i)} x_j^{(i)}$$

$$H = \frac{1}{m} \sum_{i=0}^m g'(\theta^\top x^{(i)}) x^{(i)} x^{(i)\top} \quad (? ? ?)$$

(Express H as sum
of outer products?)

Now, show $H \in \text{PSD}$

$\hookrightarrow H$ is already symmetric

$$\rightarrow \text{WTS: } z^\top H z \geq 0 \quad \forall z \in \mathbb{R}^n$$

$$z^\top H z = \frac{1}{m} z^\top \left(\sum_{i=0}^m g'(\theta^\top x^{(i)}) x^{(i)} (x^{(i)})^\top \right) z$$

(matrix mult distributes over addition)

$$= \frac{1}{m} \left(\sum_{i=0}^m g'(\theta^\top x^{(i)}) z^\top x^{(i)} (x^{(i)})^\top z \right) \quad g'(z) = g(z)(1 - g(z))$$

$$= \frac{1}{m} \left(\sum_{i=0}^m g'(\theta^\top x^{(i)}) ((x^{(i)})^\top z)^\top ((x^{(i)})^\top z) \right)$$

\hookrightarrow since in
 $(\theta^\top x^{(i)})^\top g'(z)$ can
never be negative.

$$= \frac{1}{m} \left(\sum_{i=0}^m g'(\theta^\top x^{(i)}) ((x^{(i)})^\top z)^2 \right)$$

\hookrightarrow squared for n
always pos.

The typical way to write the Hessian

$$X = \begin{bmatrix} -x^{(1)} - \\ | \\ | \\ -x^{(m)} - \end{bmatrix} \quad \boxed{H = X^T D X}$$

where $D = \text{diag}(d)$

$$\& d = \begin{bmatrix} g'(\theta^T x^{(1)}) \\ | \\ | \\ g'(\theta^T x^{(m)}) \end{bmatrix} = \begin{bmatrix} h_{\theta}(x^{(1)}) (1 - h_{\theta}(x^{(1)})) \\ | \\ | \\ h_{\theta}(x^{(m)}) (1 - h_{\theta}(x^{(m)})) \end{bmatrix}$$

Let $d^{(i)} = g'(\theta^T x^{(i)})$ (a scalar!)

To do: Show this coincides with earlier derivation.

- Let's get an expression for H_{ij} .

$$X^T D = \begin{bmatrix} | & | \\ x^{(1)} & \dots & x^{(m)} \\ | & | \end{bmatrix} \begin{bmatrix} d^{(1)} \\ \ddots \\ d^{(m)} \end{bmatrix}$$

$$= \begin{bmatrix} | & | \\ d^{(1)} x^{(1)} & \dots & d^{(m)} x^{(m)} \\ | & | \end{bmatrix}$$

$$(X^T D)_{ij} = x_i^{(j)} \cdot d^{(j)}$$

$$X^T D X = \begin{bmatrix} | & | \\ d^{(1)} x^{(1)} & \dots & d^{(m)} x^{(m)} \\ | & | \end{bmatrix} \begin{bmatrix} -x^{(1)} - \\ | \\ | \\ -x^{(m)} - \end{bmatrix}$$

$$= \begin{bmatrix} 1 & & & & & \\ \vdots & d^{(1)} x_k^{(1)} & \cdots & d^{(m)} x_k^{(m)} & & \\ & \vdots & & \vdots & & \\ & & & & x_j^{(1)} & \\ & & & & \vdots & \\ & & & & & x_j^{(m)} \end{bmatrix} \begin{bmatrix} x_j^{(1)} \\ \vdots \\ x_j^{(m)} \end{bmatrix}$$

$$H_{kj} = \sum_{c=1}^m d^{(c)} x_k^{(c)} x_j^{(c)} = \sum_{c=1}^m g'(\theta^T x^{(c)}) x_k^{(c)} x_j^{(c)}$$

which is exactly the expression we derived earlier.

Now, WTS: $z^T H z \geq 0$ for all $z \in \mathbb{R}^m$

what does an element of that matrix look like?

$H \in \mathbb{R}^{m \times m}$?

$$\begin{bmatrix} z_1 & \dots & z_m \end{bmatrix} - h_i - \begin{bmatrix} z_1 \\ \vdots \\ z_m \end{bmatrix}$$

$$\left(\sum_{i=1}^m z_i h_i \right) = \left[\sum_{i=1}^m z_i H_{1i} \dots \sum_{i=1}^m z_i H_{mi} \right] \begin{bmatrix} z_1 \\ \vdots \\ z_m \end{bmatrix}$$

$$\sum_{i=1}^m z_i H_{1i} z_1 + \dots + \sum_{i=1}^m z_i H_{mi} z_m$$

$$z^T H z = \sum_{j=1}^m \sum_{i=1}^m z_i H_{ij} z_j$$

$$= \sum_{j=1}^m \sum_{i=1}^m z_i \left(\sum_{k=1}^m g'(\theta^T x^{(k)}) x_i^{(k)} x_j^{(k)} \right) z_j$$

$$= \sum_{k=1}^m g'(\theta^T x^{(k)}) \sum_{j=1}^m \sum_{i=1}^m z_i x_i^{(k)} x_j^{(k)} z_j$$

$$= ((x^{(k)})^T z)^2$$

$$= \sum_{k=1}^m g(z) (1 - g(z)) ((x^{(k)})^T z)^2$$

$\underbrace{g(z) \in (0, 1)}_{\text{squared always } \geq 0}$

That is
remarkably
similar to
what I got
earlier. Nice.

Need to verify: $\sum_i \sum_j z_i x_i x_j z_j = (x^T z)^2$

$$(x^T z) \cdot (x^T z) = (\sum_i x_i z_i) (\sum_j x_j z_j)$$

$$= \sum_c \sum_j x_i z_i x_j z_j$$

for all possible
 i & j pairings
in the mult.

(b) implementing logistic regression using Newton's method

Vector valued version:

$$\theta^{(t+1)} := \theta^{(t)} + H^{-1} \nabla_{\theta} l(\theta)$$

$$X\theta = \begin{bmatrix} \rightarrow x^{(1)} \rightarrow \\ \vdots \\ \rightarrow x^{(m)} \rightarrow \end{bmatrix} \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}$$

Plotting decision boundary

$$\frac{1}{1 + \exp(-\theta^T x)} = 0.5$$

$$\sigma = -\theta^T x - \theta_0$$

$$\theta_0 = -\theta^T x$$

$$1 = 0.5 (1 + \exp(-\theta^T x))$$

$$\theta_0 = -\theta^T x_1 - \theta^T x_2$$

$$1 = \frac{1}{2} + \frac{1}{2} \exp(-\theta^T x)$$

$$\frac{1}{2} = \frac{1}{2} \exp(-\theta^T x)$$

$$1 = \exp(-\theta^T x)$$

$$\theta^T x = 0$$

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 = 0$$

$$\theta_2 x_2 = -\theta_0 - \theta_1 x_1$$

$$x_2 = \frac{-\theta_0 - \theta_1 x_1}{\theta_2}$$

Figure 1



Set 1: GDA

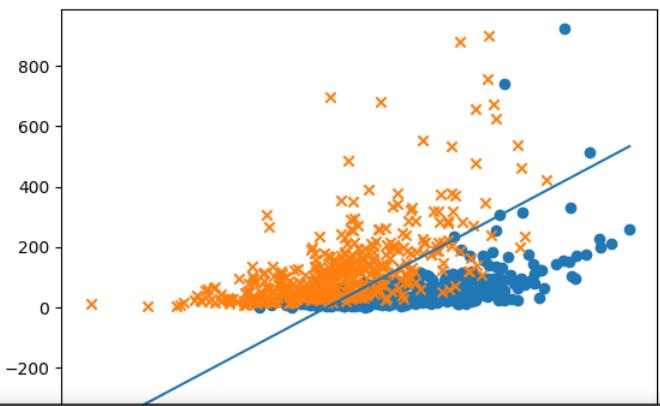


Figure 2



Set 1: LR

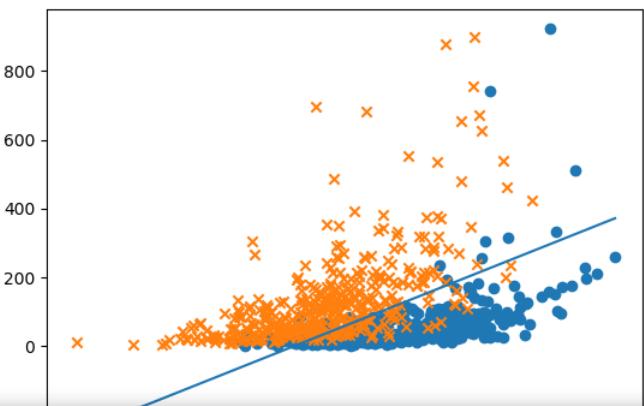


Figure 3



Set 2: GDA

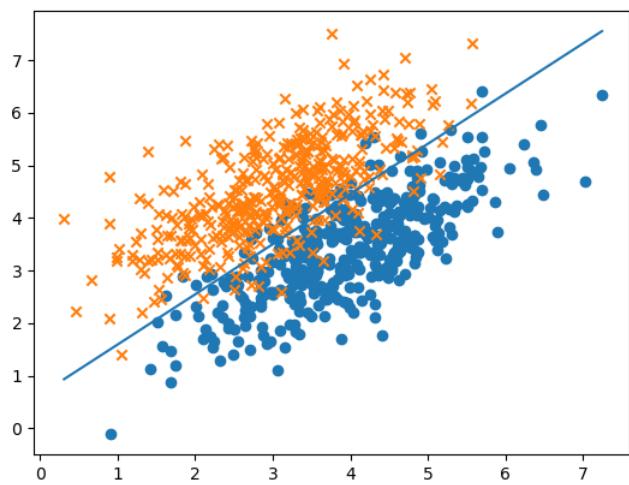
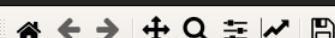
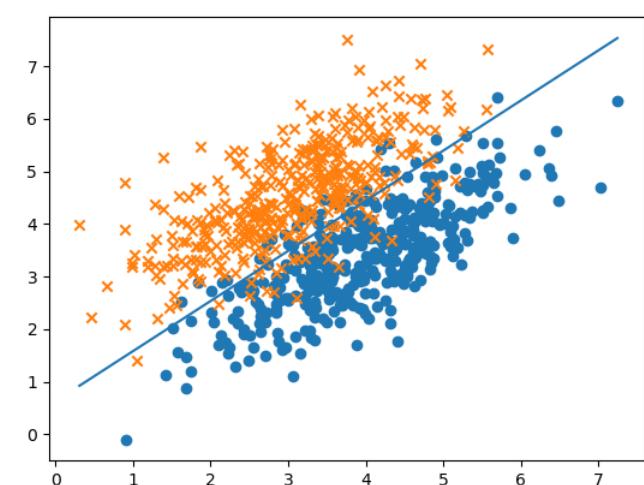


Figure 4



Set 2: LR



```
In [120]: %cpaste -q
print("set 2")
print("LR\n",confusion_mat(df_2_train["y"],LR_2_preds))
print("GDA\n",confusion_mat(df_2_train["y"],GDA_2_preds))
-
set 2
LR
[[367  33]
 [ 35 365]]
GDA
[[368  32]
 [ 37 363]]
```

Confusion matrices

```
-- 
set 1
LR
[[358  42]
 [ 52 348]]
GDA
[[351  49]
 [ 58 342]]
```

GDA/LR performance is about the same on set 2,
w/ GDA performing slightly better.

GDA does worse on set 1.

→ The data does not appear to be Gaussian.

(c) Recall GDA

$$x|y=0 \sim N(\mu_0, \Sigma)$$

$$x|y=1 \sim N(\mu_1, \Sigma)$$

$$y \sim \text{Bernoulli}(\phi)$$

Show GDA has a linear decision boundary, namely;

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

where $\theta \in \mathbb{R}^n$ & $\theta_0 \in \mathbb{R}$ are functions of the given parameters.

I think I'm
really off the mark
here...

$$P(y=1|x) = \frac{P(x|y=1) P(y=1)}{P(x)}$$

$$(A \text{ little incorrect}) = \frac{P(x|y=1) \cancel{P(y=1)} = \phi}{P(y=1) \cdot P(x|y=1) + P(y=0) P(x|y=0)}$$

Strictly speaking,
density fn...

$$= \frac{\frac{1}{(2\pi)^{n/2}} |\Sigma|^{1/2} \exp\left(-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1)\right) \cdot \phi}{\frac{1}{(2\pi)^{n/2}} |\Sigma|^{1/2} \exp\left(-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1)\right) + (1-\phi) \frac{\exp\left(-\frac{1}{2}(x-\mu_0)^T \Sigma^{-1} (x-\mu_0)\right)}{M_0}}$$

$$= \frac{M_1 \phi}{\phi M_1 - (1-\phi) M_0} = \frac{M_1 \phi}{\phi M_1 - M_0 + \phi M_0} = \frac{M_1 \phi}{\phi(M_1 + M_0) - M_0}$$

$$= \frac{1}{\frac{\phi(M_1 + M_0)}{M_1 \phi} - \frac{M_0}{M_1 \phi}} = \frac{1}{1 + \frac{M_0}{M_1} - \frac{M_0}{M_1} \cdot \frac{1}{\phi}} = \frac{1}{1 + \left(1 - \frac{1}{\phi}\right) \frac{M_0}{M_1}}$$

$$\left(1 - \frac{1}{\phi}\right) \frac{M_0}{M_1} = \exp\left(\ln\left(1 - \frac{1}{\phi}\right)\right) \exp\left(-\frac{1}{2}(x-\mu_0)^T \Sigma^{-1} (x-\mu_0) - \frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1)\right)$$

Can it be combined into one term?

doesn't seem correct / need to identify θ_0, θ_1

$$P(x|y=1) \quad P(y=1) = 0$$

(Attempt 2)

$$= \frac{P(x|y=1)}{P(y=1) \cdot P(x|y=1) + P(y=0) \cdot P(x|y=0)}$$

$$= \frac{\phi P(x|y=1)}{\phi P(x|y=1) + (1-\phi) P(x|y=0)}$$

$$= \frac{\phi M_1}{\phi M_1 + (1-\phi) M_0} = \frac{1}{\frac{\phi M_1}{\phi M_1} + \frac{1-\phi}{\phi} \frac{M_0}{M_1}}$$

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

↑
 $\exp \left(\ln \left(\frac{1-\phi}{\phi} \right) \right)$

(distribute into ...)

$$\dots \exp \left[\frac{-1}{2} \left((x - \mu_0)^T \Sigma^{-1} x - (x - \mu_0)^T \Sigma^{-1} \mu_0 \right. \right.$$

$$\left. \left. - \left((x - \mu_1)^T \Sigma^{-1} x - (x - \mu_1)^T \Sigma^{-1} \mu_1 \right) \right] \right]$$

$$\dots (x^T \Sigma^{-1} - \mu_0^T \Sigma^{-1}) x - \boxed{(x^T \Sigma^{-1} - \mu_0^T \Sigma^{-1}) \mu_0}$$

$$\quad \quad \quad - (x^T \Sigma^{-1} - \mu_1^T \Sigma^{-1}) x \quad \boxed{+ (x^T \Sigma^{-1} - \mu_1^T \Sigma^{-1}) \mu_1} \dots$$

↓

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

$$\Rightarrow (\mu_1 - \mu_0)^T \Sigma^{-1} x$$

$$\boxed{x^T \Sigma^{-1} \mu_1 - x^T \Sigma^{-1} \mu_0 + \mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1}$$

\therefore Inverse of symm. matrix is symm $\Rightarrow \Sigma^{-1}$ symm.
 $\therefore x^T A y = y^T A x$ for any symm. matrix A .

("Symmetric Bilinear Form")

$$\begin{aligned} &\Rightarrow \mu_1^\top \Sigma^{-1} x - \mu_0^\top \Sigma^{-1} x \\ &= (\mu_1 - \mu_0)^\top \Sigma^{-1} x \\ &= \dots \exp\left(-\frac{1}{2} (2(\mu_1 - \mu_0)^\top \Sigma^{-1} x + \mu_0^\top \Sigma^{-1} \mu_0 - \mu_1^\top \Sigma^{-1} \mu_1)\right) \\ &\sim \exp\left(\frac{1-\phi}{\phi}\right) \end{aligned}$$

Another fact you didn't know:

$$(\mu_0 + \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1)$$

$$\mu_0^\top \Sigma^{-1} (\mu_0 - \mu_1) + \mu_1^\top \Sigma^{-1} (\mu_0 - \mu_1)$$

$$\mu_0^\top \Sigma^{-1} \mu_0 - \mu_0^\top \Sigma^{-1} \mu_1 + \underbrace{\mu_1^\top \Sigma^{-1} \mu_0 - \mu_1^\top \Sigma^{-1} \mu_1}_{\downarrow \text{symmetric}}$$

cancel

$$\Rightarrow \mu_0^\top \Sigma^{-1} \mu_0 - \mu_1^\top \Sigma^{-1} \mu_1$$

$$\Rightarrow \dots \exp\left[-\frac{1}{2}(2(\mu_1 - \mu_0)^\top \Sigma^{-1} x + (\mu_0 + \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1)) + \ln\left(\frac{1-\phi}{\phi}\right)\right]$$

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

$$\Rightarrow \Theta = (\mu_1 - \mu_0)^\top \Sigma^{-1}$$

$$\Theta_0 = \frac{1}{2}(\mu_0 + \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1) - \ln\left(\frac{1-\phi}{\phi}\right)$$

(d) Derive the MLE estimates for GDA.

$$\ell(\phi, \mu_0, \mu_1, \Sigma) = \ln \prod_{i=1}^m P(x^{(i)} | y^{(i)}, \dots) P(y^{(i)} | \phi)$$

Taking partials w.r.t each parameter assume all is const
Scalar (so no matrix derivatives?)

Recall here, $y^{(i)} \in \{0, 1\}$. or, $y^{(i)} = 1 \Leftrightarrow y^{(i)} = 1$.
(a little redundant.)

$$= \ln \prod_{i=1}^m \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_{y^{(i)}})^T \Sigma^{-1} (x - \mu_{y^{(i)}})\right)$$

$$= \sum_{i=1}^m \ln \left[\dots \right]$$

(simplify the log)
likelihood

$$= \sum_{i=1}^m \ln \left(\frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \phi^{y^{(i)}} (1-\phi)^{1-y^{(i)}} \right) + \boxed{-\frac{1}{2}(x - \mu_{y^{(i)}})^T \Sigma^{-1} (x - \mu_{y^{(i)}})}$$

$$\boxed{\ln \left(\frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \right)} + \ln(\phi^{y^{(i)}} (1-\phi)^{1-y^{(i)}})$$

$$\boxed{+} \quad \boxed{\ln(\phi) y^{(i)} + \ln(1-\phi) (1-y^{(i)})} \quad \boxed{+} \quad \boxed{-}$$

$$\frac{\partial \ell(\dots)}{\partial \phi} = \sum_{i=1}^m \frac{y^{(i)}}{\phi} - \frac{1-y^{(i)}}{1-\phi} \quad \text{Estimator of } \phi:$$

Set 0, solve for ϕ

$$0 = \sum_{i=1}^m \frac{y^{(i)}}{\phi} - \frac{1-y^{(i)}}{1-\phi} = \sum_{i=1}^m \frac{(1-\phi)y^{(i)} - \phi(1-y^{(i)})}{\phi(1-\phi)}$$

$$= \sum_{i=1}^m y^{(i)} - \phi y^{(i)} - \phi + \phi y^{(i)} = \sum_{i=1}^m y^{(i)} + \phi(-y^{(i)} - 1 + y^{(i)}) \\ = \sum_{i=1}^m (y^{(i)} - \phi)$$

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

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

Estimators of μ_0/μ_1 :

$$\frac{\partial l}{\partial \mu_0} = \sum_{i=1}^m \frac{\partial}{\partial \mu_0} \left[\frac{1}{2} (x^{(i)} - \mu_2 c^{(i)})^T \Sigma^{-1} (x^{(i)} - \mu_2 c^{(i)}) \right]$$

Note if $y^{(i)} = 1$, the derivative of the inner term is 0.

But, if $y^{(i)} = 0$, we can get a nonzero derivative.
Also rewrite for univariate case:

$$\Sigma = \sigma^2 I, \quad \Sigma^{-1} = \frac{1}{\sigma^2} I$$

$$= \sum_{i=1}^m \frac{\partial}{\partial \mu_0} \left[\frac{-1}{2\sigma^2} (x^{(i)} - \mu_0)^2 \right]$$

$$0 = \sum_{\substack{i \in \{1 \leq i \leq m \\ |y^{(i)} = 0\}}} \frac{-1}{2\sigma^2} \cdot 2(x^{(i)} - \mu_0) (-1) - \sum_{\substack{i \in \{1 \leq i \leq m \\ |y^{(i)} = 0\}}} \frac{1}{\sigma^2} (x^{(i)} - \mu_0)$$

$$= \sum_{\substack{i \in \{1 \leq i \leq m \\ |y^{(i)} = 0\}}} (x^{(i)} - \mu_0) = \sum_{i=1}^m x^{(i)} \mathbb{1}_{\{y^{(i)} = 0\}} - \mu_0$$

$$0 = \sum_{i=1}^m x^{(i)} \mathbb{1}_{\{y^{(i)} = 0\}} - \mu_0 \sum_{i=1}^m \mathbb{1}_{\{y^{(i)} = 0\}}$$

$$\mu_0 = \frac{\sum_{i=1}^m x^{(i)} \mathbb{1}_{\{y^{(i)} = 0\}}}{\sum_{i=0}^m \mathbb{1}_{\{y^{(i)} = 0\}}}$$

And the same derivations
should hold for μ_1 :
just change the label.

Estimator of $\Sigma = \sigma^2$: $\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_2 c^{(i)})^2$

$$\frac{\partial l}{\partial \sigma^2} = \sum_{i=0}^m \frac{\partial}{\partial \sigma^2} \left[\ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) + \frac{-1}{2\sigma^2} (x^{(i)} - \mu_2 c^{(i)})^2 \right]$$

$$\left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^{-1} \cdot \left(\frac{1}{2} \right) (2\pi\sigma^2)^{-\frac{3}{2}} \cdot 2\sigma + \left(\frac{1}{2} \right) (-1)(\sigma^2)^{-2} (x^{(i)} - \mu_2 c^{(i)})^2$$

$$- \left(\frac{1}{2\sigma^2} \right)^{-1} \cdot 2\sigma \cdot \left(\frac{1}{2} \right) + \left(\frac{1}{2} \right) (\sigma^2)^{-2} (x^{(i)} - \mu_2 c^{(i)})^2$$

$$- \frac{1}{2} (\sigma^2)^{-2} (x^{(i)} - \mu_2 c^{(i)})^2 = \frac{1}{2} (\sigma^2)^{-2} \sigma^2$$

$$\sum_{i=1}^m \frac{1}{2} (\sigma^2)^{-2} \left[(x^{(i)} - \mu_{y(i)})^2 - \sigma^2 \right]$$

$$\sum_{i=1}^m \frac{(x^{(i)} - \mu_{y(i)})^2 - \sigma^2}{2(\sigma^2)^2}$$

$$\frac{1}{2(\sigma^2)^2} \left[\sum_{i=1}^m (x^{(i)} - \mu_{y(i)})^2 - m\sigma^2 \right]$$

Set equal to 0:

$$\frac{m\sigma^2}{2(\sigma^2)^2} = \frac{1}{2(\sigma^2)^2} \sum_{i=1}^m (x^{(i)} - \mu_{y(i)})^2$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y(i)})^2$$

which is equivalent to the expression given
in the univariate case, i.e.

$$(x^{(i)} - \mu_{y(i)})^2 = (x^{(i)} - \mu_{y(i)})^T (x^{(i)} - \mu_{y(i)})$$

when $x \in \mathbb{R}$.

(2) Incomplete, Positive Only Labels.

Scenario: Binary labeled data, but only a subset of the positive examples are labeled. The rest are unlabeled. They could be pos. or neg.

Goal: Construct a binary classifier that matches the true label as closely as possible.

Notation

$$y^{(i)} = \begin{cases} 1 & \text{if example } i \text{ is labeled} \\ 0 & \text{otherwise} \end{cases}$$

$$\ell^{(i)} = \begin{cases} 1 & \text{if the true label is 1} \\ 0 & \text{if the true label is 0} \end{cases}$$

(a) Suppose the labeled examples were selected uniformly randomly from the set of positive examples. Then,

$$P(y^{(i)}=1 | \ell^{(i)}=1, x^{(i)}) = P(y^{(i)}=1 | \ell^{(i)}=1)$$

In other words, it didn't matter what the features of the examples were when we selected the ones to be the labeled examples.

↳ e.g., violated if the only labeled ones were

$x > 100$, due to measurement sensitivity?

Show: $P(\ell^{(i)}=1 | x^{(i)}) = \frac{1}{\alpha} \cdot P(y^{(i)}=1 | x^{(i)})$

The probability of being labeled differs by a constant factor from the prob. of having a true label of 1.

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

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

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

constant w.r.t $x^{(i)}$
 but might vary
 i ?

Constant, roughly equal to the proportion

$$= \frac{\text{labeled examples}}{\text{all positive examples}}$$

(b) Estimate α using some trained model h validation set V

$$\text{Let } V_+ := \{x^{(i)} \in V \mid y^{(i)} = 1\}$$

\hookrightarrow the labeled examples in V .

Assume:

$$\circ h(x^{(i)}) \approx P(y^{(i)}=1 | x^{(i)})$$

$$\circ x^{(i)} \in V_+ \Rightarrow P(t^{(i)}=1 | x^{(i)}) \approx 1 \quad \text{if in labeled set,}$$

the classifier predicts prob of being labeled

Show: $h(x^{(i)}) \approx \alpha \quad \forall x^{(i)} \in V_+$

that $x^{(i)}$'s true label is 1.

Assume $x^{(i)} \in V_+$. Then, by the last problem,

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

$\hookrightarrow \approx h(x^{(i)})$ $\hookrightarrow \approx 1$ by assumption
 α by last problem

$$\Rightarrow h(x^{(i)}) \approx \alpha$$

Plotting decision boundary for α -corrected prediction

$$\frac{1}{\alpha} \cdot \frac{1}{1 + \exp(-\theta^T x)} = 0.5$$

$$\frac{1}{1 + \exp(-\theta^T x)} = \alpha \cdot 0.5$$

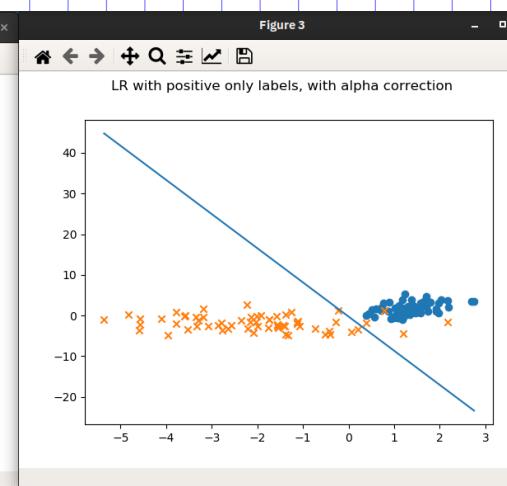
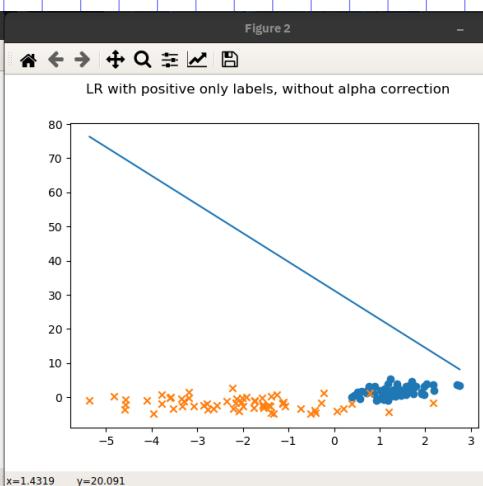
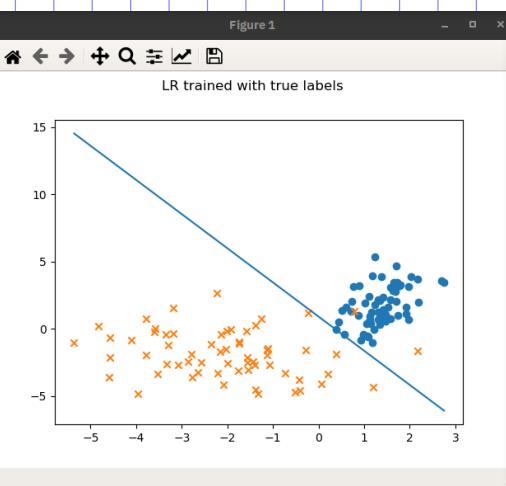
$$1 + \exp(-\theta^T x) = \frac{1}{\alpha \cdot 0.5}$$

$$\exp(-\theta^T x) = \frac{1}{\alpha \cdot 0.5} - 1$$

$$-\theta^T x = \ln\left(\frac{1}{\alpha \cdot 0.5} - 1\right)$$

$$-\theta_0 + -\theta_1 x_1 - \theta_2 x_2 = \ln\left(\frac{1}{\alpha \cdot 0.5} - 1\right)$$

$$= \frac{\ln\left(\frac{1}{\alpha \cdot 0.5} - 1\right)}{-\theta_2} + \theta_1 x_1 + \theta_0$$



LR trained with access to true labels

```
[[14  0]
 [50 60]]
```

LR trained with access to partial positive labels, no rescale

```
[[ 0 14]
 [ 0 110]]
```

LR trained with access to partial positive labels

```
[[14  0]
 [52 58]]
```

(3) Poisson Regression

(a) Show that the Poisson dist. is in the exponential family.

$$\begin{aligned} p(y; \lambda) &= \frac{e^{-\lambda} \lambda^y}{y!} = \frac{1}{y!} \cdot \exp(\ln(\lambda^y)) e^{-\lambda} \\ &= \frac{1}{y!} \exp(y \ln(\lambda) - \lambda) \\ &= \frac{1}{y!} \exp(y \ln(\lambda) - \exp(\ln(\lambda))) \end{aligned}$$

The Poisson distribution is in the Exponential Family w/...

$$\begin{array}{ll} T(y) = y & a(n) = e^n \\ b(y) = \frac{1}{y!} & n = \ln(\lambda) \Rightarrow \boxed{\lambda = e^n} \end{array}$$

(b) What is the canonical response fn when we do GLM regression w/ a Poisson R.V?

& in fact that's the prediction fn, since Poisson reg. would assume:

$$\begin{aligned} y|x &\sim \text{Poisson}(\lambda), \text{ but } \lambda = e^{\eta} = e^{\theta^T x} \\ \Rightarrow h_{\theta}(x) = \mathbb{E}[y|x; \theta] &= \lambda = e^{\eta} = e^{\theta^T x} \Rightarrow h_{\theta}(x) = e^{\theta^T x} \end{aligned}$$

(c) Derive the stochastic gradient ascent rule by maximizing the log likelihood.

Gradient ascent:

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

$$l(\theta) = \log \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta) = \log \prod_{i=1}^m \frac{e^{-\lambda(\theta)}}{y^{(i)!}} \frac{y^{(i)}}{\lambda(\theta)}$$

$$= \sum_{i=1}^m \log(-\cdot) = \sum_{i=1}^m -\lambda(\theta) + y^{(i)} \log(\lambda(\theta)) + \log(\frac{1}{y^{(i)!}})$$

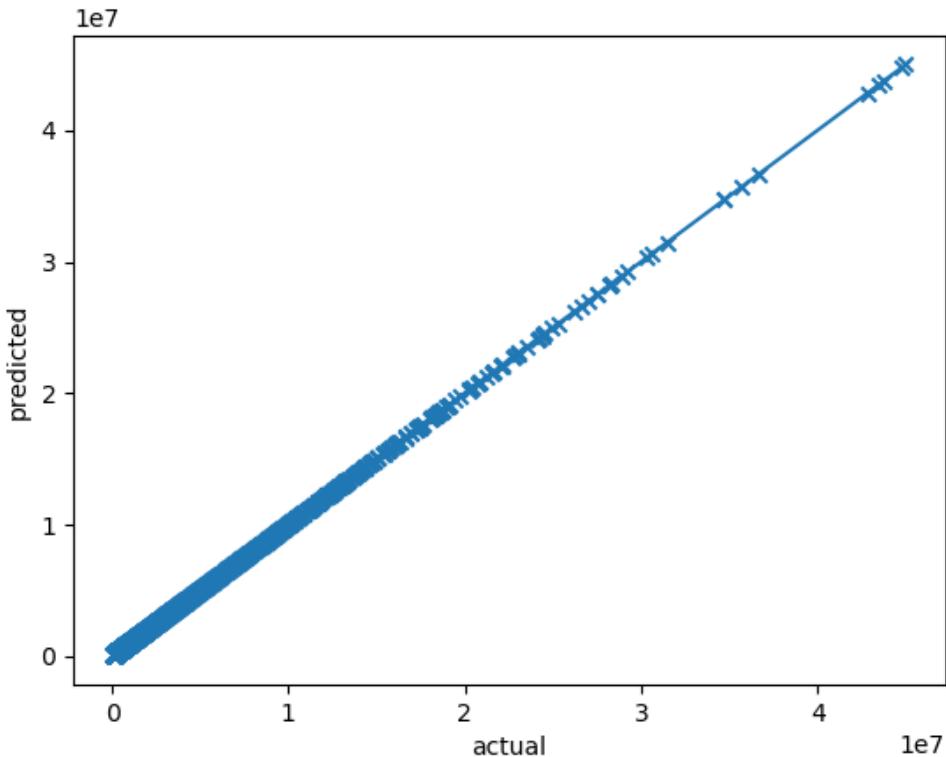
$$= \sum_{i=1}^m -e^{\theta^T x^{(i)}} + y^{(i)} \log(e^{\theta^T x^{(i)}}) + \log(\frac{1}{y^{(i)!}})$$

$$\begin{aligned}
 \frac{\partial L}{\partial \theta_j} &= \sum_{i=1}^m \frac{\partial}{\partial \theta_j} \left[-e^{\theta^T x^{(i)}} + y^{(i)} e^{\theta^T x^{(i)}} \right] \\
 &= \sum_{i=1}^m -e^{\theta^T x^{(i)}} \cdot x_j^{(i)} + y^{(i)} x_j^{(i)} \\
 &= \sum_{i=1}^m x_j^{(i)} (y^{(i)} - \underbrace{e^{\theta^T x^{(i)}}}_{\log(x^{(i)})})
 \end{aligned}$$

Figure 1



Predicted vs. Actual



x=4.19589e+07 y=4.6118e+07

[7.36659576 3.93332251 3.43327325 1.999993815 4.39995804]

Final θ
vector

Key: need to choose a pretty extreme learning rate $\alpha \approx 10^{-8}$ since otherwise gradient vectors are too large: since magnitude of y too large.

Using a fixed number of SGD iterations worked okay in this case. (~15 iters)

(4) Convexity of GLMs

Goal: Show NLL (Negative Log-Likelihood) of a GLM is convex in the natural parameter η .

⇒ local maxima are global maxima.

Show this by showing the Hessian is PSD.

Setup: $Y|X \sim \text{ExpFamily}(\eta)$

where $\eta \in \mathbb{R}$ & the sufficient statistic is $T(Y) = y$.

$$\Rightarrow f_{Y|X}(y|\eta) = b(y) \exp(\eta y - a(\eta))$$

Recall that $\eta(X) = \theta^T X$. But we'll get there in a second.

$$(a) \text{Show } E[Y|X;\eta] = \frac{\partial a}{\partial \eta}.$$

Let A be the support set of $Y|X$.

$$= \int_A y \cdot f_{Y|X}(y|\eta) dy = \int_A y \cdot b(y) \exp(\eta y - a(\eta)) dy$$

$$= \int_A b(y) \cdot y \exp(\eta y - a(\eta)) dy$$

$$\begin{aligned} \frac{\partial}{\partial \eta} [\exp(\eta y - a(\eta))] &= \exp(\eta y - a(\eta)) \cdot [y - \frac{\partial a}{\partial \eta}] \\ &= \exp(\eta y - a(\eta)) y - \frac{\partial a}{\partial \eta} \exp(\eta y - a(\eta)) \end{aligned}$$

$$= \int_A b(y) \left[\frac{\partial}{\partial \eta} [\exp(\eta y - a(\eta))] + \frac{\partial a}{\partial \eta} \exp(\eta y - a(\eta)) \right] dy$$

$$\stackrel{1}{=} \frac{\partial}{\partial \eta} \int_A b(y) \exp(\eta y - a(\eta)) dy + \frac{\partial a}{\partial \eta} \int_A b(y) \exp(\eta y - a(\eta)) dy$$

$$= \frac{\partial}{\partial \eta} \left[\int_A f_{Y|X}(y|\eta) dy \right] + \frac{\partial a}{\partial \eta} \int_A f_{Y|X}(y|\eta) dy$$

$$= \frac{\partial}{\partial \eta} [1] + \frac{\partial a}{\partial \eta} \cdot 1 = \frac{\partial a}{\partial \eta}$$

$$(b) \text{ Show } \text{Var}(Y|X; \theta) = \frac{\partial^2 a}{\partial n^2}.$$

$$\begin{aligned} \mathbb{E}[(Y|X)^2 | \theta] &= \int_A y^2 b(y) \exp(ny - a(n)) dy \\ &= \frac{\partial^2}{\partial n^2} [\exp(ny - a(n))] \\ &= \frac{\partial}{\partial n} \left[\exp(ny - a(n)) y - \frac{\partial a}{\partial n} \exp(ny - a(n)) \right] \\ &= y \frac{\partial}{\partial n} [\exp(ny - a(n))] - \left[\frac{\partial^2 a}{\partial n^2} \exp(ny - a(n)) + \frac{\partial a}{\partial n} \frac{\partial}{\partial n} [\exp(ny - a(n))] \right] \end{aligned}$$

Short hand: let $M = \exp(ny - a(n))$

$$\Rightarrow \frac{\partial M}{\partial n} = My - \frac{\partial a}{\partial n} M$$

$$\begin{aligned} &= y \left(My - \frac{\partial a}{\partial n} M \right) - \frac{\partial^2 a}{\partial n^2} M + \frac{\partial a}{\partial n} \left(My - \frac{\partial a}{\partial n} M \right) \\ &= My^2 - \cancel{\frac{\partial a}{\partial n} M y} - \frac{\partial^2 a}{\partial n^2} M + \cancel{\frac{\partial a}{\partial n} My} - \left(\frac{\partial a}{\partial n} \right)^2 M \\ \Rightarrow & My^2 = \frac{\partial^2 M}{\partial n^2} + \frac{\partial^2 a}{\partial n^2} M + \left(\frac{\partial a}{\partial n} \right)^2 M \\ \Rightarrow & \dots = \int_A b(y) \left[\frac{\partial^2 M}{\partial n^2} + \frac{\partial^2 a}{\partial n^2} M + \left(\frac{\partial a}{\partial n} \right)^2 M \right] dy \\ &= \frac{\partial^2}{\partial n^2} \left[\int_A b(y) M dy \right] + \frac{\partial^2 a}{\partial n^2} \int_A b(y) M dy + \left(\frac{\partial a}{\partial n} \right)^2 \int_A b(y) M dy \\ \Rightarrow & \mathbb{E}[(Y|X)^2 | n] = \frac{\partial^2 a}{\partial n^2} + \left(\frac{\partial a}{\partial n} \right)^2 \end{aligned}$$

$$\begin{aligned} \text{Var}(Y|X; n) &= \mathbb{E}[(Y|X)^2 | n] - (\mathbb{E}[Y|X; n])^2 \\ &= \frac{\partial^2 a}{\partial n^2} + \left(\frac{\partial a}{\partial n} \right)^2 - \left(\frac{\partial a}{\partial n} \right)^2 = \frac{\partial^2 a}{\partial n^2}. \end{aligned}$$

(c) $\ell(\theta)$: NLL. Then show $\ell(\theta)$'s Hessian is always PSD, which shows the convexity of GLMs.

$$\begin{aligned}\ell(\theta) &= -\log \prod_{i=1}^m f_{Y|X}(y^{(i)} | x^{(i)}, \theta) \\ &= -\log \prod_{i=1}^m b(y^{(i)}) \exp(n y^{(i)} - a(n)) \\ &= -\sum_{i=1}^m \log(b(y^{(i)})) + [n y^{(i)} - a(n)]\end{aligned}$$

& recall: $\eta = \theta^T x^{(i)}$

$$= -\sum_{i=1}^m \log(b(x^{(i)})) + \theta^T x^{(i)} \cdot y^{(i)} - a(\theta^T x^{(i)})$$

Compute gradient + hessian:

$$\begin{aligned}\frac{\partial \ell}{\partial \theta_j} &= -\sum_{i=1}^m \frac{\partial}{\partial \theta} [\theta^T x^{(i)}] y^{(i)} - \left[\frac{\partial a}{\partial \eta} \Big|_{\eta=\theta^T x^{(i)}} \right] \frac{\partial}{\partial \theta} [\theta^T x^{(i)}] \\ &= -\sum_{i=1}^m x_j^{(i)} y^{(i)} - h_\theta(x^{(i)}) x_j^{(i)} \\ &= -\sum_{i=1}^m x_j^{(i)} (y^{(i)} - h_\theta(x^{(i)}))\end{aligned}$$

where $h_\theta(x^{(i)}) = \frac{\partial a}{\partial \eta} \Big|_{\eta=\theta^T x^{(i)}} = \mathbb{E}[Y | X=x^{(i)}]$

$$\begin{aligned}\frac{\partial^2 \ell}{\partial \theta_j \partial \theta_k} &= \frac{\partial}{\partial \theta_k} \left[-\sum_{i=1}^m x_j^{(i)} y^{(i)} - x_j^{(i)} \cdot \frac{\partial a}{\partial \eta} \Big|_{\eta=\theta^T x^{(i)}} \right] \\ &= \sum_{i=1}^m x_j^{(i)} x_k^{(i)} \frac{\partial^2 a}{\partial \eta^2} \Big|_{\eta=\theta^T x^{(i)}} \\ &= \sum_{i=1}^m x_j^{(i)} x_k^{(i)} \text{Var}(Y | X=x^{(i)})\end{aligned}$$

$$H = X^T D X \quad \text{where}$$

$$X = \begin{bmatrix} -x^{(1)} \\ \vdots \\ -x^{(m)} \end{bmatrix}$$

$m \times n$

$$D = \text{diag}([\text{Var}(Y|X=x^{(1)}) \dots \text{Var}(Y|X=x^{(m)})])$$

(so D is a diagonal matrix w/ all positive entries on diag.)

To show PSD, wts $\forall z \in \mathbb{R}^n \quad z^T H z \geq 0$.

- Since D is a diagonal matrix w/ only positive entries on its diagonal, D is PSD.

Easy to prove. Suppose all $d_i \geq 0$.

$$\begin{bmatrix} x_1 & \dots & x_m \end{bmatrix} \begin{bmatrix} d_1 \\ \vdots \\ d_m \end{bmatrix} = [x_1 d_1 \dots x_m d_m]$$

$$[x_1 d_1 \dots x_m d_m] \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} = \sum x_i^2 d_i \geq 0.$$

$$\Rightarrow z^T (X^T D X) z = (Xz)^T D (Xz)$$

Let $y = Xz$. $y \in \mathbb{R}^m$.

$$y^T D y \geq 0 \text{ since } D \text{ PSD.}$$

$$\Rightarrow (Xz)^T D (Xz) \geq 0 \Rightarrow z^T H z \geq 0.$$

H is PSD.

(5) Locally Weighted Linear Regression

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m w^{(i)} (\theta^\top x^{(i)} - y^{(i)})^2$$

(i) Show that $J(\theta)$ can be rewritten as

$$J(\theta) = (X\theta - y)^\top W (X\theta - y)$$

for an appropriate matrix W .

$$X = \begin{bmatrix} -x^{(1)} & \cdots \\ \vdots & \ddots \\ -x^{(m)} & \cdots \end{bmatrix}$$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$X\theta - y = \begin{bmatrix} \theta^\top x^{(1)} - y^{(1)} \\ \vdots \\ \theta^\top x^{(m)} - y^{(m)} \end{bmatrix}_{m \times 1}$$

Guess:

$$W = \text{diag}[w^{(1)} \dots w^{(m)}]$$

i.e. $w_{ij} = \begin{cases} w^{(i)} & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$

$$(X\theta - y)^\top W = [w^{(1)} \theta^\top x^{(1)} \dots w^{(m)} \theta^\top x^{(m)}]$$

$$\begin{aligned} \Rightarrow (X\theta - y)^\top W (X\theta - y) &= [w^{(1)} \theta^\top x^{(1)} \dots w^{(m)} \theta^\top x^{(m)}] [\theta^\top x^{(1)} \\ &\quad \vdots \\ &\quad \theta^\top x^{(m)}] \\ &= w^{(1)} (\theta^\top x^{(1)})^2 + \dots + w^{(m)} (\theta^\top x^{(m)})^2 \\ &= \sum_{i=1}^m w^{(i)} (\theta^\top x^{(i)})^2 \end{aligned}$$

(ii) Generalize the normal equations to the weighted setting.

(Take the derivative $\nabla_\theta J(\theta)$ & set it equal to 0.)

$$\begin{aligned} \frac{\partial J}{\partial \theta_j} &= \frac{1}{2} \sum_{i=1}^m w^{(i)} \frac{\partial}{\partial \theta_j} \{ (\theta^\top x^{(i)} - y^{(i)})^2 \} = \frac{1}{2} \sum_{i=1}^m w^{(i)} \cdot 2(\theta^\top x^{(i)} - y^{(i)}) x_j^{(i)} \\ &= \sum_{i=1}^m w^{(i)} (\theta^\top x^{(i)} - y^{(i)}) x_j^{(i)} \end{aligned}$$

Somehow express in matrix form.

$$\langle X\theta - y \rangle_i = \theta^T x^{(i)} - y^{(i)}$$

$$WX = \begin{bmatrix} w^{(1)} \\ \vdots \\ w^{(m)} \end{bmatrix} \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(m)} \end{bmatrix} = \begin{bmatrix} -w^{(1)}x^{(1)} - \\ \vdots \\ -w^{(m)}x^{(m)} - \end{bmatrix}$$

$$\begin{bmatrix} \overbrace{\quad}^{m \times m} \overbrace{\quad}^{m \times n} \\ \overbrace{\quad}^{n \times n} \end{bmatrix} \quad mx1$$

No dimension mismatch $\rightarrow x^T W (X\theta - y)$ (this is right)

$$x^T w = \begin{bmatrix} w^{(1)} & \dots & w^{(m)} \\ | & \dots & | \\ 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} \theta^T x^{(1)} - y^{(1)} \\ \vdots \\ \theta^T x^{(m)} - y^{(m)} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^m w^{(i)} x_i^{(i)} (\theta^T x^{(i)} - y^{(i)}) \\ \vdots \end{bmatrix}$$

$$\Rightarrow x^T w x \theta - x^T w y = 0$$

$$x^T w x \theta = x^T w y$$

Or, you could use
matrix calc rules
I guess.

$$\boxed{\theta = (X^T W X)^{-1} X^T W y}$$

(iii) Dataset $\{(x^{(i)}, y^{(i)})\}_{i \in \{1, \dots, m\}}$

but the probability model of each y has a different variance $(\sigma^{(i)})^2$.

$$\text{i.e., } y^{(i)} | x^{(i)} \sim N(\theta^T x^{(i)}, (\sigma^{(i)})^2)$$

Show: The ML estimate of θ in this scenario is a weighted linear regression problem.

What is the weight in $w^{(i)}$?

$$\ell(\theta) = \ln \prod_{i=1}^m P(y^{(i)} | x^{(i)}, \theta)$$

$$= \ln \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^{(i)}}} \exp \left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^{(i)})^2} \right)$$

$$= \sum_{i=1}^m \ln\left(\frac{1}{\sqrt{\sigma^{(i)}}}\right) - \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^{(i)})^2}$$

$$\frac{\partial l}{\partial \theta_j} = \sum_{i=1}^m -\frac{1}{2(\sigma^{(i)})^2} \cdot 2(y^{(i)} - \theta^T x^{(i)}) \cdot (-x_j^{(i)})$$

$$= \sum_{i=1}^m \frac{1}{(\sigma^{(i)})^2} (y^{(i)} - \theta^T x^{(i)}) x_j^{(i)}$$

$$\Rightarrow w^{(i)} = (\sigma^{(i)})^{-2}$$

Figure 1

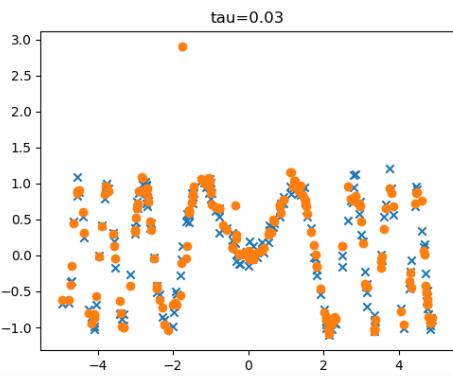


Figure 2

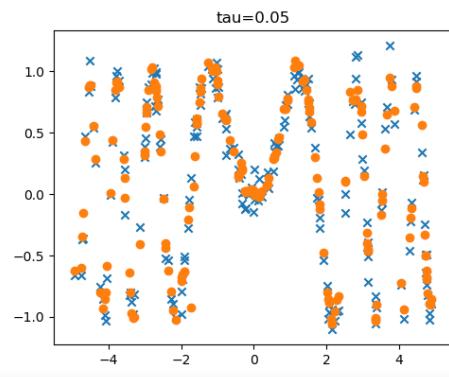


Figure 3

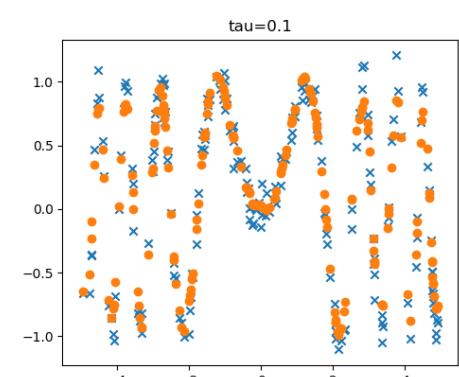


Figure 4

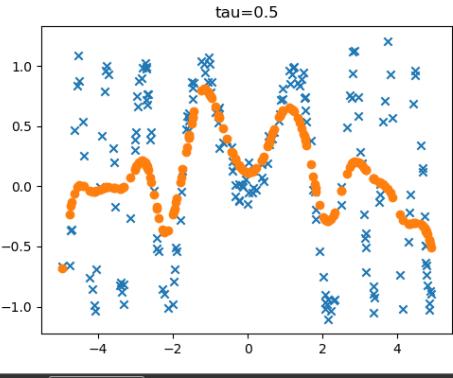


Figure 5

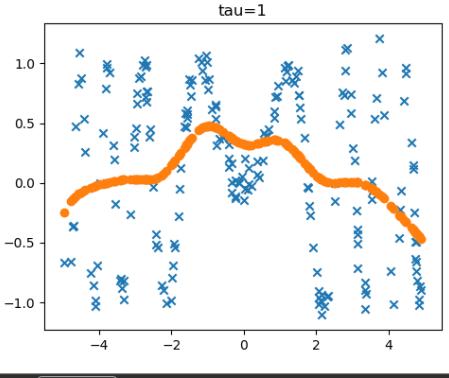
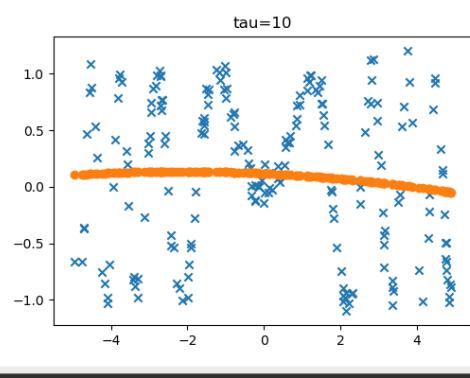


Figure 6



```
--  
validation set mse for tau=0.03: 0.059038566452005385  
validation set mse for tau=0.05: 0.01894440603986488  
validation set mse for tau=0.1: 0.014461728527284943  
validation set mse for tau=0.5: 0.2578133818426481  
validation set mse for tau=1: 0.3794515567158331  
validation set mse for tau=10: 0.42989004449194135
```