

Name :- Prateek Mahajan

~~EST~~

EE511 : Homework 2

Problem 1: -

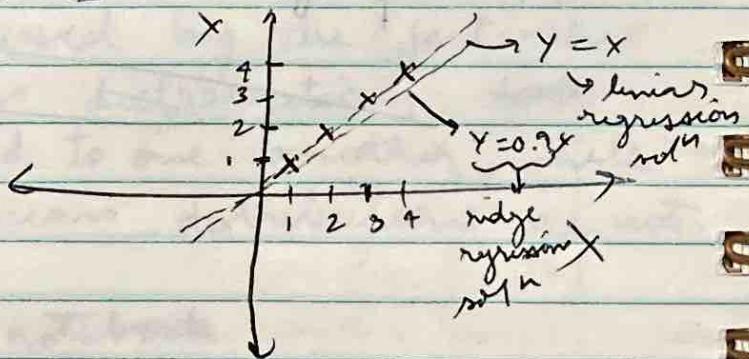
general

- 1a) I have read and understood the instructions for HW2 at the top of HW2 and I formally declare that all the work I turn in for everything in this course will not contain or involve any cheating at all.

Problem 2: -

- 2a) Consider a simple dataset as below: -

X	Y
1	1
2	2
3	3
4	4



For a dataset like the above, linear regression would give a solⁿ of $y = x$ with zero error.

For ridge regression:- (setting $\lambda = 0.4 / n = 0.1$) and assuming 4 datapoints

let's assume that we get $y = 0.9x$

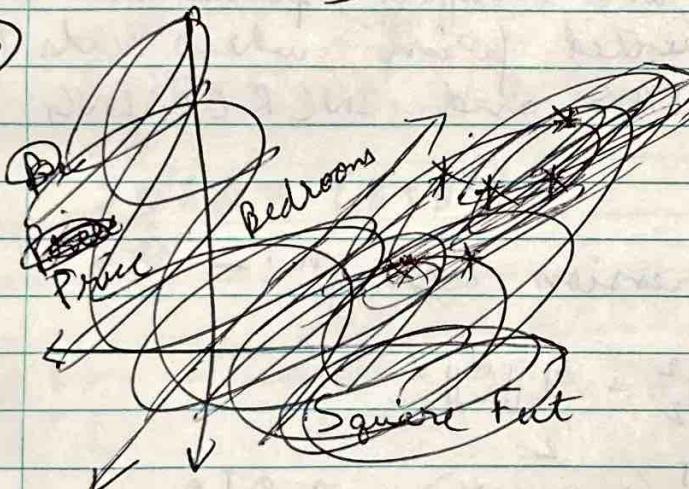
Using the ridge regression solⁿ, Error($y = x$) = 4, Error($y = 0.9x$) = 3.9
so clearly, it would be preferred despite $y = x$

fitting the data better

In cases like the above, where there is only 1 feature or there is NO correlation amongst features, the linear regression is better than ridge. The reason for this is that ridge regression adds a bias that makes it good at dealing with collinearity amongst various features, but makes it perform worse when the focus is on interpreting the individual effects of features & when there isn't collinearity amongst them.

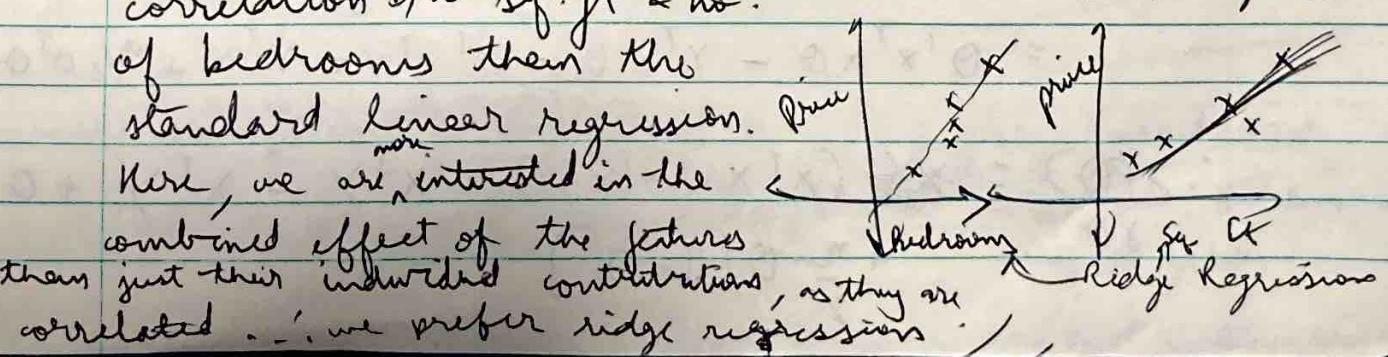
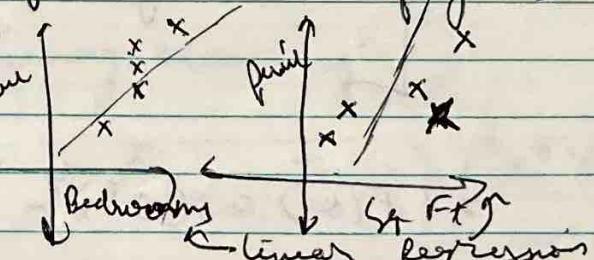
Ridge

b)



	Price	Bedrooms	Sq. ft.
	400	3	2104
	330	3	1600
	369	3	2400
	232	2	1416
	540	4	3000
	:	:	,
	:	:	,
	:	:	,

In the given example, there are ~~not~~ 2 correlated features (no. of bedrooms & sq. ft.). The addition of a regularization term ~~alone~~ adds a bias to the model that helps it better represent the correlation b/w sq. ft & no. of bedrooms than the standard linear regression. Here, we are interested in the combined effect of the features than just their individual contributions, as they are correlated. ∴ we prefer ridge regression.



~~c) Increasing η modifies the optimization set by increasing regularization and therefore, increases bias & reduces variance.~~

c) Increasing η tries to push the parameters of the optimisation problem to zero thereby simplifying the model and INCREASING BIAS and DECREASING VARIANCE by pushing the model to a more general form.

Decreasing η tries to give more importance to square error and therefore, pushes the model to a more complicated form which leads to DECREASING BIAS and INCREASING VARIANCE //

d) Ridge regression cost f^n :-

$$\begin{aligned} F(\theta) &= \|x\theta - y\|_2^2 + \frac{\eta}{2} \|\theta\|_2^2 \\ &= (x\theta - y)^T (x\theta - y) + \frac{\eta}{2} \theta^T \theta \end{aligned}$$

[in matrix form, if $x \rightarrow n \times m$ design matrix
 $y \rightarrow n \times 1$ output matrix
 $\theta \rightarrow m \times 1$ parameter feature matrix]

$$\begin{aligned} \Rightarrow F(\theta) &= (\theta^T x^T - y^T)(x\theta - y) + \frac{\eta}{2} \theta^T \theta \\ &= \theta^T x^T x \theta - y^T x \theta - \theta^T x^T y + y^T y + \frac{\eta}{2} \theta^T \theta \end{aligned}$$

$$\begin{aligned} \therefore \frac{\partial F(\theta)}{\partial \theta} &= (x^T x + x^T x) \theta - x^T y - x^T y + \theta \\ &+ \frac{\eta}{2} \cdot \theta [I + I] \end{aligned}$$

$$\Rightarrow \frac{\partial f(\theta)}{\partial \theta} = 2X^T X \theta - 2X^T Y + 2\theta \cancel{+ 2\theta}$$

$$\text{Setting } \frac{\partial F(\theta)}{\partial \theta} = 0 \Rightarrow 2X^T X \theta + 2\theta = 2X^T Y \\ \Rightarrow \left(X^T X + \frac{nI}{2}\right) \theta = X^T Y$$

$$\Rightarrow \hat{\theta}_{\text{optimal}} = (X^T X + \frac{nI}{2})^{-1} \cdot X^T Y$$

closed form solⁿ of
ridge regression

Formulas used:-

$$\rightarrow \frac{\partial b^T \theta^T \theta}{\partial \theta} = \theta (b c^T + c b^T)$$

$$\rightarrow \frac{\partial \theta^T a}{\partial \theta} = \frac{\partial a^T \theta}{\partial \theta} = a$$

$$\rightarrow \frac{\partial a^T \theta b}{\partial \theta} = a b^T$$

$$\rightarrow \frac{\partial \theta^T B \theta}{\partial \theta} = (B + B^T) \theta$$

e) For a vanilla linear regression,

$$1) F(\theta) = \|X\theta - y\|_2^2 = (X\theta - y)^T (X\theta - y) \\ = (\theta^T X^T - y^T) \cdot (X\theta - y) \\ = \theta^T X^T X \theta - y^T X \theta - \theta^T X^T y + y^T y$$

, where $X \rightarrow n \times m$ design matrix, $\theta \rightarrow m \times 1$ matrix, $y \rightarrow n \times 1$ matrix

$$\therefore \frac{\partial F}{\partial \theta} = 2X^T X \theta - 2X^T Y \quad (\text{using above eqn})$$

$$\Rightarrow \frac{\partial F}{\partial \theta} = 0 \Rightarrow 2X^T X \theta = 2X^T Y \Rightarrow \hat{\theta}_{\text{optimal}} = (X^T X)^{-1} X^T Y$$

$X^T X \rightarrow m \times m$ matrix

If $m > n$ and features of X are highly correlated,
 $X^T X$ may NOT be invertible. (since the determinant
may be zero)

$X^T X$ is
(if not invertible)

In that case, the $\hat{\theta}$ optimal / closed form soln of the vanilla linear regression may not be computable.

- 2) comparing the vanilla linear regression closed form soln to ridge regression's :-
the matrix to be inverted in ridge regression is $(X^T X + \lambda I)$, not $X^T X$.
 \therefore as long as we make sure $\lambda > 0$, $(X^T X + \lambda I)$ is always invertible \Rightarrow the closed form soln for ridge regression always exists (which is the benefit of ridge regression)

f)

$$\begin{aligned} (i) F(\theta, \theta_0) &= \|X\theta + \theta_0 1 - y\|_2^2 + \frac{\eta}{2} \|\theta\|_2^2 \\ &= (X\theta + \theta_0 1 - y)^T (X\theta + \theta_0 1 - y) + \frac{\eta}{2} \theta^T \theta \\ &= (\theta^T X^T + 1^T \theta_0^T - y^T) \cdot (X\theta + \theta_0 1 - y) + \frac{\eta}{2} \theta^T \theta \\ &= \theta^T X^T X \theta + \theta^T X^T \theta_0 1 - \theta^T X^T y + 1^T \theta_0^T X \theta + 1^T \theta_0^T 1 \\ &\quad - 1^T \theta_0^T y - y^T X \theta - y^T \theta_0 1 + y^T y + \frac{\eta}{2} \theta^T \theta \end{aligned}$$

$$\therefore \frac{\partial F(\theta, \theta_0)}{\partial \theta_0} = 0 + 1^T X \theta_0 - 0 + \cancel{1^T X^T X \theta} + \cancel{1^T X^T \theta_0 1} - \cancel{1^T y} + \cancel{0} - \cancel{y^T \theta_0} + 0 + 0 + 1^T X \theta$$

$$\begin{aligned} \Rightarrow \frac{\partial F(\theta, \theta_0)}{\partial \theta_0} &\approx 2 X \theta_0 1^T + 2 \theta_0 1 1^T - 2 y 1^T \\ &= 2 1^T X \theta + 2 1^T \theta_0 1 - 2 1^T y \end{aligned}$$

Setting $\frac{\partial F(\theta, \theta_0)}{\partial \theta} = 0$

$$\Rightarrow \cancel{1^T y} = \cancel{1^T \theta_0} + 1^T \theta_0 \cdot 1 - \cancel{1^T \theta_0} - 0 \quad \text{--- (1)}$$

$$\Rightarrow \cancel{y} = \cancel{x \theta_0} + \theta_0 \cdot 1 - \cancel{\theta_0} \quad \text{--- (2)}$$

$$\Rightarrow \cancel{y} = \cancel{x \theta_0} + \theta_0 \cdot 1 \quad \text{--- (3)}$$

Now consider $\frac{\partial F(\theta, \theta_0)}{\partial \theta}$

$$\therefore \frac{\partial F(\theta, \theta_0)}{\partial \theta} = \frac{\partial}{\partial \theta} (2x^T x \theta + x^T \theta_0 \cdot 1 - x^T y + x^T \theta_0 \cdot 1 + 0 - 0) \\ = -x^T y - 0 + 0 + \eta \theta$$

$$\Rightarrow \frac{\partial F(\theta, \theta_0)}{\partial \theta} = 2x^T x \theta + 2x^T \theta_0 \cdot 1 - 2x^T y + \eta \theta$$

Setting $\frac{\partial F(\theta, \theta_0)}{\partial \theta} = 0$

$$\Rightarrow 2x^T y = 2x^T x \theta + 2x^T \theta_0 \cdot 1 + \eta \theta$$

On substituting θ : —

~~$$2x^T y = 2x^T x \theta + 2x^T \theta_0 \cdot 1 - 2x^T y + \eta \theta$$~~

~~$$2x^T y = 2x^T x \theta + 2x^T x \theta - 2x^T y + \eta \theta$$~~

~~$$\Rightarrow 4x^T y = 4x^T x \theta + \eta \theta$$~~

~~$$\Rightarrow 2x^T y \cdot 1^T = 2x^T x \theta \cdot 1^T + 2x^T (y \cdot 1^T - x \theta \cdot 1^T) + \eta \theta \cdot 1^T$$~~

$$\Rightarrow (2x^T x + \eta I) \cdot \theta = 2x^T y - 2x^T \theta_0 \cdot 1$$

$$\Rightarrow \theta = \frac{(x^T x + \eta I)^{-1}}{2} (x^T) (y - \theta_0 \cdot 1) \quad \text{--- (2)}$$

Taking (1): —

~~$$y \cdot 1^T = x \theta_0 \cdot 1^T + \theta_0 \cdot 1 \cdot 1^T$$~~

However θ_0 is a scalar bias.

~~$$\therefore \theta_0 \cdot 1 \cdot 1^T = \{\theta_0\}_{n \times n}$$~~

This eqⁿ is only satisfied if can be interpreted as θ_0 optimal

Taking ① : —

$$1^T y = 1^T x \theta + 1^T \theta_0 1$$

$$\Rightarrow 1^T \theta_0 1 = 1^T y - 1^T x \theta$$

$$\Rightarrow \theta_0 1^T 1 = 1^T y - 1^T x \theta \quad \text{--- ③}$$

$$\Rightarrow n \theta_0 = \sum_{i=1}^n y_i - \sum_{i=1}^n \sum_{j=1}^m x_{i,j} \theta_{j,1}$$

$$\Rightarrow \theta_0 = \frac{1}{n} \left(\sum_{i=1}^n y_i - \sum_{j=1}^m x_{i,j} \theta_{j,1} \right)$$

= mean of prediction error (say μ_e)

Taking ② : —

$$\theta = (x^T x + \frac{n}{2} I)^{-1} x^T (y - \theta_0 1)$$

$$\Rightarrow \theta = (x^T x + \frac{n}{2} I)^{-1} x^T (y - \mu_e 1) //$$

If we manipulate differently (i.e. ~~without~~ substituting θ before inverse) :

And ~~$\theta_0 = \mu_y - \frac{1^T x \theta}{n}$~~ (from ③)

$$(x^T x + \frac{n}{2} I) \theta = x^T y + x^T \frac{1^T x \theta 1}{n} - x^T \mu_y 1$$

$$[\mu_y = \frac{1}{n} \sum_{i=1}^n y_i]$$

Note that $(1^T x \theta)$ is a 1×1 result & ~~so~~ can be treated as a scalar

$$\Rightarrow (x^T x + \frac{n}{2} I) \theta = x^T y + \frac{x^T 1 (1^T x \theta)}{n} - x^T \mu_y 1$$

$$\Rightarrow (x^T x + \frac{n}{2} I - \frac{x^T 1 1^T x}{n}) \theta = x^T (y - \mu_y 1)$$

$$\Rightarrow \theta = (x^T x + \frac{n}{2} I - \frac{x^T 1 1^T x}{n})^{-1} x^T (y - \mu_y 1)$$

NOTE : - 2f(ii) is after Problem 3 b.)

~~no sorry~~ sorry about the confusion

~~model, and so regularizing it would~~

Problem 3 :-

a) Consider a gaussian noise linear least squares regression model, where $\vec{y} = X\theta + \vec{\epsilon}$, where X is a $n \times m$ design matrix & $\vec{\epsilon}$ is a length- n vector of gaussians, $\epsilon_i \in N(0, \sigma^2)$.

$$\text{MLE parameter estimate} = \hat{\theta} = (X^T X)^{-1} X^T \vec{y}$$

Let θ^* be the "true" parameter estimate for the gaussian LLS model. ($\theta^* \in \mathbb{R}^{m \times 1}$)

~~Let $Y|n$ be n instances of x from a dataset~~

~~Let $Y|n$ be n instances of x from a dataset~~

$\therefore E(Y|n) = \text{true prediction of model}$
(as $E(Y|x)$ is the best fit soln)
= $x^T \theta^*$, where n instances of x from a dataset
 $(x \in \mathbb{R}^{m \times n})$. — (1)

~~Let $h_\theta(x)$~~

$$\begin{aligned} \text{Consider } h_\theta(x) &= x^T \hat{\theta} \\ &= x^T ((X^T X)^{-1} X^T \vec{y}) \\ &= x^T (X^T X)^{-1} x^T (X \theta^* + \vec{\epsilon}) \quad (\text{since } \vec{y} = X \theta^* + \vec{\epsilon}) \\ &= x^T \theta^* + x^T (X^T X)^{-1} x^T \vec{\epsilon} \end{aligned}$$

$$\therefore E_D(h_\theta(x)) = E_D(x^T \theta^* + x^T (X^T X)^{-1} x^T \vec{\epsilon})$$

$$= E_D(x^T \theta^*) + E_D(x^T (X^T X)^{-1} x^T \vec{\epsilon})$$

n, x are fixed matrices which implies this term is ~~$x^T (X^T X)^{-1} x^T E_D(\vec{\epsilon}) = 0$~~

$$\therefore E_D(h_\theta(x)) = E_D(x^T \theta^*)$$

$$= x^T \theta^* \quad (\text{as } n, \theta^* \text{ are not variables}) \quad (2)$$

$$\therefore (1) = (2) \Rightarrow E_D(h_\theta(x)) = E(Y|n) \Rightarrow \text{this model is unbiased}$$

b) Consider the LNS = $E_D[(h_0(x) - E_D(h_0(x)))^2]$

$$\therefore \text{LNS} = E_D[h_0(x) - n^T \theta^*]^2$$

~~Since $E_D(h_0(x)) = n^T \theta^*$ (as per 3a)~~

$$h_0(x) = n^T \theta^*$$

$$= n^T (x^T x)^{-1} x^T y$$

$$= n^T (x^T x)^{-1} x^T (x^T \theta^* + e)$$

$$= n^T (x^T x)^{-1} x^T \theta^* + n^T (x^T x)^{-1} x^T e$$

$$= n^T \theta^* + n^T (x^T x)^{-1} x^T e$$

$$\therefore \text{LNS} = E_D(n^T \theta^* + n^T (x^T x)^{-1} x^T e - n^T \theta^*)^2$$

~~$= E_D((n^T (x^T x)^{-1} x^T e)^T (n^T (x^T x)^{-1} x^T e))$~~

~~$= E_D(n^T (x^T x)^{-1} x^T e e^T x (x^T x)^{-1} n^T)$~~

~~$= E_D(n^T (x^T x)^{-1} x^T e e^T x (x^T x)^{-1})$~~

this is

~~$= E_D(n^T (x^T x)^{-1} x^T E E^T x (x^T x)^{-1} n^T)$~~

$$\Rightarrow \text{LNS} = E_D[(n^T (x^T x)^{-1} x^T e)^2]$$

$$= E_D[(n^T (x^T x)^{-1} x^T e)^T (n^T (x^T x)^{-1} x^T e)]$$

$$= E_D[(n^T (x^T x)^{-1} x^T e)^T (n^T (x^T x)^{-1} x^T e)]$$

$$= E_D[n^T (x^T x)^{-1} x^T E E^T x (x^T x)^{-1} n^T]$$

n, x are not random variables

$$\Rightarrow \text{LNS} = E_D[n^T (x^T x)^{-1} x^T E E^T x (x^T x)^{-1} n^T]$$

[Note:- $(x^T x)$ is symmetric as it is a covariance matrix $\Rightarrow (x^T x)^T = x^T x$

$$\Rightarrow (x^T x)^{-1} = ((x^T x)^T)^{-1} = (x^T x)^{-1}$$

Similarly, $E(E E^T) = \sigma^2 I$ (as e is a GRV

& $E^T E^T$ is a covariance matrix

$$\begin{aligned} \Rightarrow LNS &= x^T (x^T x)^{-1} x^T (\sigma^2 I) x (x^T x)^{-1} x \\ &= \sigma^2 x^T (x^T x)^{-1} x^T x (x^T x)^{-1} x \\ &= \sigma^2 x^T (x^T x)^{-1} x = RNS \end{aligned}$$

$\therefore LNS = RNS \Rightarrow 3b)$ is correct

Problem 2: - (continued)

f)
(ii) $\hat{\theta} = (x^T x + \frac{\gamma}{2} I - \underbrace{x^T 1 1^T x}_{\text{bias ridge}})^{-1} x^T (y - \mu_y I))$

$$\hat{\theta}_{\text{normal ridge}} = (x^T x + \frac{\gamma}{2} I)^{-1} x^T y$$

On comparing these clearly, it looks like
when we do not penalize the bias / offset /
intercept parameter by regularising it,
the resultant answer seems to
center the "x" & "y" terms ~~more~~, and ~~may~~
by subtracting their means from them.

$\mu_y = \text{mean of } y$, $\frac{x^T 1 1^T x}{n}$ is kind of like

the mean of $x^T x$ as we are essentially ~~not~~
adding a multiplicant of $\frac{1 1^T}{n}$ between $x^T x$,

which is somewhat like that the average of
 $x^T x$]

\therefore not penalizing the offset parameter "normalizes"
- is "the effect of x & y " ~~not~~ by ~~the~~ subtracting
their means

~~Further, bias is an inherent property~~

Further bias is an inherent property of the model. The
parameter is not responsible for changing the shape of the

model and so, regularising it would not prevent it from overfitting (or have any function whatsoever).

Problem 4:-

LSTAT
DIS

- a) According to the scatter plot, LSTAT & RM are ~~the~~ the most correlated to MEDV.
- b) According to the correlation matrix RM & ~~RATIO~~ are the most correlated to MEDV. The answer ~~would~~ is ~~different from~~ q3, as we are trying to correlate house prices from the same dataset. The 3rd most correlated feature is not clear in the graphs.

c) linear Regression :-

$$\text{closed form soln} : - \hat{\theta} = (X^T X)^{-1} X^T y$$

Ridge Regression :-

$$\text{closed form soln} : - \hat{\theta} = (X^T X + \frac{\eta I}{2})^{-1} X^T y$$

Please refer to the end of this pdf or the ipynb file for coefficient values (note that I used $\eta = 20.0$ for this problem)

- d) Please refer to the .ipynb file attached at the end of this pdf for the final RMSE outputs I got for the linear & ridge regression.

Training RMSE :-

linear < Ridge (linear is better)

This is because the ridge regression performs regularisation to reduce the tendency of the model to overfit to the training data and ~~seems~~ generalise it more (i.e. increase bias & reduce variance)

since the linear regression model is trained to fit the training data better, it performs better on that data

Testing RMSE :-

linear \rightarrow Ridge (ridge is better)

The linear model does not regularise & so, has higher variance ~~as~~ than the ridge model, which regularises the cost fn to reduce the tendency of the model to overfit the training data. For this reason, the ridge model is more generalised & performs better on another dataset (i.e. the testing one)

e) Please refer to the .ipynb file for my RMSE outputs (NOTE $\eta = 20$)

On comparing the RMSE outputs of top 3 features ~~as~~ that of all 13, we can see that the RMSE is higher when we use just 3 features (albeit not significantly).

This indicates that while a large chunk of correlation ^{with me} originates in these features, there are some other features in the model that have significance & ~~should~~ be ignored.

Problem 5:-

$$a) F(\theta) = \frac{1}{n} \sum_{i=1}^n -\log [P(y=y_i | x=x_i, \theta)] + \frac{\gamma}{2} \| \theta \|^2_2$$

~~First let us consider -~~
 where $n \in \mathbb{R}^{1 \times m}$, $\theta \in \mathbb{R}^{m \times c}$

First, let us consider

$$P(y=y_i | x_i, \theta) = \frac{\exp(n \cdot \theta^{(k)})}{\sum_{j=1}^c \exp(n \cdot \theta^{(j)})}$$

(where $\theta^{(k)}$ is the K^{th} column of θ)

$$\Rightarrow \log(P(y=y_i | x_i, \theta)) = \underbrace{n \theta^{(k)}}_{\text{take this as } z^{(k)}} - \log \left(\sum_{j=1}^c \exp(z^{(j)}) \right)$$

which implies
that this is $z^{(j)}$

$$\begin{aligned} \therefore \frac{\partial}{\partial z^{(k)}} (\log(P(y=y_i | x_i, \theta))) &= \frac{\partial (n \theta^{(k)})}{\partial z^{(k)}} - \frac{\partial}{\partial z^{(k)}} \left(\sum_{j=1}^c \exp(n \cdot \theta^{(j)}) \right) \\ &= \frac{\partial (z^{(k)})}{\partial z^{(k)}} - \frac{\partial}{\partial z^{(k)}} \left(\sum_{j=1}^c \exp(z^{(j)}) \right) \end{aligned}$$

(NOTE: — k is any value b/w 1 & c)

$$\frac{\partial z^{(k)}}{\partial z^{(l)}} = \begin{cases} 1, & \text{if } k=l \\ 0, & \text{otherwise} \end{cases}$$

$$\therefore \frac{\partial}{\partial z^{(k)}} (\log(P(y=y_i | x_i, \theta))) = 1 \{k=1\} - \frac{1}{\sum_{j=1}^c \exp(z^{(j)})} \sum_{j=1}^c e^{z^{(j)}}$$

$$\text{But } \frac{\partial}{\partial z^{(k)}} \sum_{j=1}^c e^{z^{(j)}} \approx \frac{\partial}{\partial z^{(k)}} e^{z^{(k)}} = e^{z^{(k)}}$$

(as all other terms will give zero)

(CS cont.)

$$\therefore \frac{\partial}{\partial z^{(i)}} (\log(P(y=y_i|x_i, \theta))) = 1_{\{k=1\}} - \frac{e^{z^{(i)}}}{\sum_{j=1}^C \exp(z^{(i)})}$$

~~$$\Rightarrow \frac{\partial}{\partial z^{(i)}} (\log(P(y=y_i|x_i, \theta))) = 1_{\{k=1\}} - P(y=1|x_i, \theta)$$~~

~~$$\therefore \frac{\partial}{\partial \theta} (\log(P(y=y_i|x_i, \theta))) = \cancel{\frac{\partial}{\partial z^{(i)}}} \cdot \cancel{\frac{\partial}{\partial z^{(i)}}} \cdot \cancel{\frac{\partial}{\partial z^{(i)}}} (\log(P(y=y_i|x_i, \theta)))$$~~

~~$$\text{Consider } \frac{\partial}{\partial \theta} = \frac{\partial}{\partial \theta} (n \cdot \theta^{(i)}) = \sum_{j=1}^m n_j \theta_j^{(i)}$$~~

~~$$= \begin{bmatrix} 0 \\ 0 \\ \vdots \\ n_1 & n_2 & \dots & n_m \end{bmatrix}$$~~

$$\therefore \frac{\partial}{\partial \theta} (\log(P(y=y_i|x_i, \theta))) = \frac{\partial}{\partial \theta} \cdot \frac{\partial}{\partial z^{(i)}} (\log(P(y=y_i|x_i, \theta)))$$

~~$$\text{Consider } \frac{\partial}{\partial \theta} = \frac{\partial}{\partial \theta} (n \cdot \theta^{(i)}) = x_i \text{ (for a single column)}$$~~

~~$$\therefore \frac{\partial}{\partial \theta} (\log(P(y=y_i|x_i, \theta))) = x_i (1_{\{k=1\}} - P(y=1|x_i, \theta))$$~~

For $F(\theta)$ we ~~can't~~ sum this log probability for all n classes.

$$\begin{aligned} \Rightarrow \frac{\partial F(\theta)}{\partial \theta} &= \frac{\partial}{\partial \theta} \left[\frac{1}{n} \sum_{i=1}^n -\log(P(y=y_i|x_i, \theta)) + \frac{n}{2} \|\theta\|_2^2 \right] \\ &= -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^C n_i (1_{\{k=j\}} - P(y_i=j|x_i, \theta)) \\ &\quad + n \theta^\top \end{aligned}$$

$$\begin{aligned} \Rightarrow \frac{\partial F(\theta)}{\partial \theta} &= -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^C n_i (1_{\{k=j\}} - P(y_i=j|x_i, \theta)) \\ &\quad + n \theta^\top \end{aligned}$$

Gradient Descent Rule for θ :-

$$\theta \leftarrow \theta - \alpha \cdot \nabla F(\theta)$$

$$\Rightarrow \theta \leftarrow \theta + \frac{\alpha}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} n_i (1_{\{y_i=j\}} - p_n(y_i=j | x_i, \theta)) - \alpha \eta \quad \text{--- (2)}$$

For a matrix formulation :-

Invert back to

$$\frac{\partial F}{\partial \theta} = \frac{\partial}{\partial \theta} \log(p(y=y_i | x_i, \theta)) = \frac{\partial z^{(i)}}{\partial \theta} \cdot \frac{\partial}{\partial z^{(i)}} \log(p(y=y_i | x_i, \theta))$$

(consider $\frac{\partial z^{(i)}}{\partial \theta}$)

Branching it out further (to $\frac{\partial z}{\partial \theta}$ instead of $\frac{\partial z^{(i)}}{\partial \theta}$)

$$\frac{\partial}{\partial \theta} \log(p(y=y_i | x_i, \theta)) = \frac{\partial z}{\partial \theta} \cdot \frac{\partial}{\partial z} \log(p(y=y_i | x_i, \theta))$$

$$\frac{\partial z}{\partial \theta} \approx X^T \quad (\text{where } z = n\theta)$$

using θ (extrapolating (1)) :-

$$\frac{\partial}{\partial z} \log(p(y=y_i | x_i, \theta)) = Y - P,$$

where P is a ~~nxC~~ matrix, containing the values

$$P_{i,j} = \exp(x_{i,j} \cdot \theta_j)$$

$$P_{i,j} = \frac{\exp(x_{i,j} \cdot \theta_j)}{\sum_{l=1}^C \exp(x_{i,l} \cdot \theta_l)} \quad \text{and assuming } Y \text{ is a } 1 \text{ hot encoding matrix of size } n \times C$$

The matrix form

$$\frac{\partial}{\partial \theta} \log(p(y=y_i | x_i, \theta)) = X^T (Y - P)$$

$$\therefore \frac{\partial F}{\partial \theta} \text{ (in matrix form)} = -\frac{1}{N} X^T (Y - P) + \alpha \eta \theta$$

The least mean squares gradient descent rule is :-

$$\theta \leftarrow \theta + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) (x^{(i)})$$

On comparing ~~(1)~~ to (2) to the above:-

clearly, both are of the form

$\theta \leftarrow \theta + \alpha \times \text{error}$, where error is the deviation of the prediction from the original output in the data.

- b) Please refer to the .ipynb file at the end of this pdf for this qn.
- c) Based on the convergence curve, clearly as we increase the learning rate, the model gets trained faster / requires fewer training epochs.

However we also see that the final cost J^* is ~~lower~~ lower & training / testing precision is higher to a certain point as we reduce the learning rate.

The above are the tradeoffs to consider. However clearly, 0.01 is the sweet spot of this model as it is the point where cost J^* is lowest, precision is highest and training time is acceptable / not too long.

- d) Please refer to the .ipynb file attached at the end of the pdf for this qn.
- e) No, different batch sizes show different convergence rates with the same learning

rate. From the fig. convergence sped reduces as batch size increases.

Please find the curves with the tuned learning rates at the end of the .ipynb notebook.

Since the new curves are tuned better they give a better precision o/p than the earlier ones.

A higher learning rate worked better for smaller batch sizes because increasing λ reduces the impact of regularisation, which can become overpowering when you reduce your batch size.

A high learning rate of 0.3 yielded the overall fastest convergence in terms of wall clock time

