

# COL774 – Machine Learning

## Assignment 1

*Nilaksh Agarwal*  
2015PH10813

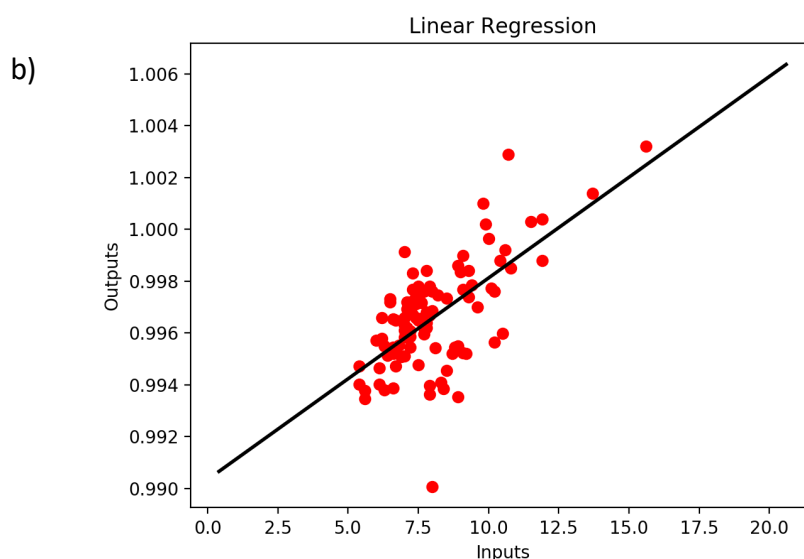
### Part 1:

In this part, we implement linear regression in 1 dimension. So, our curve fitted turns out to be a line. Since a line can be perfectly represented by two parameters (slope, intercept) we minimize our cost function to find values of this.

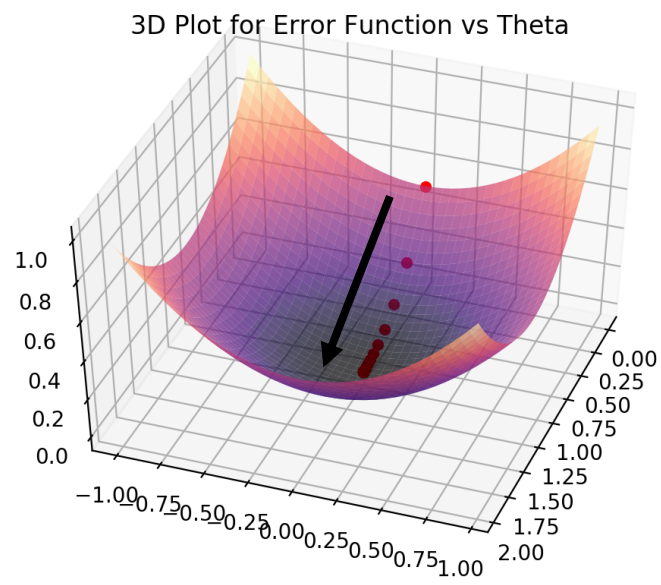
To speed up this process, I use numpy's matrix multiplication to directly get our predicted values and gradients for all parameters (instead of summing over all examples individually)

The stopping criteria is two-fold. First we have a max\_iter which provides a hard limit on the number of iterations we update our theta parameters over. Secondly we have a condition over the change in our cost function ( $J[i+1] - J[i]$ ) which has to be less than  $10^{-10}$  times our learning rate. This is done so that even a very small learning rate will still iterate and not instantly converge just due to minimal increments in J.

- a) Learning Rate: 0.3  
Stopping Criteria:  $(J[i+1] - J[i]) / LR < 10^{-10}$   
Final Parameters ( $\theta$ ): [0.99034, 0.00078]

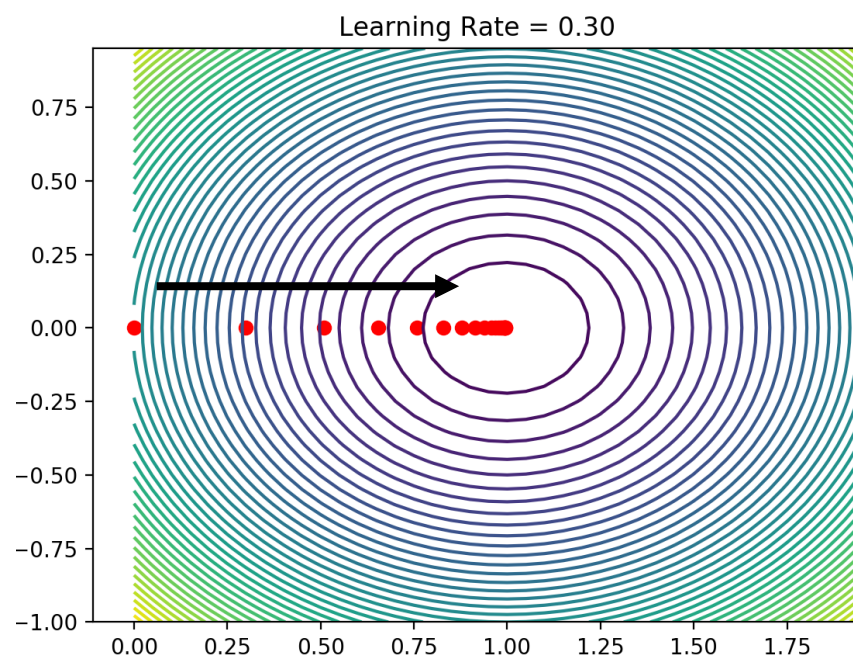


c)



Arrow Shows direction of moving cost function

d)



Arrow Shows direction of moving cost function

- e) On increasing the learning rate from 0.1 to 2.5, we see that initially, for [0.1, 0.5, 0.9] the cost function moves faster (in fewer iterations) towards the minima.

After that for [1.3, 1.7], the cost function overshoots the minima, then comes back but overshoots it again, and keeps spiraling towards the minima.

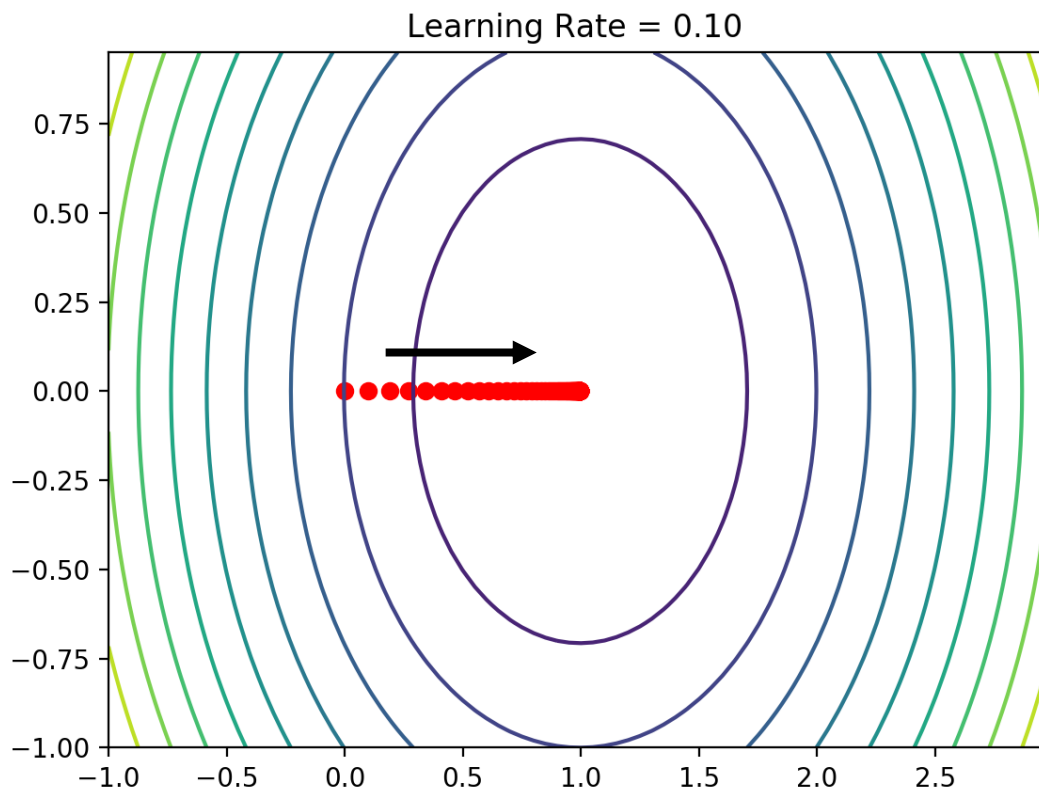
For [2.1, 2.5], the cost function overshoots the minima and keeps increasing the overshoot, progressively, the cost function keeps increasing at each iteration. In this way it never reaches the minima.

Plot always starts from [0,0]

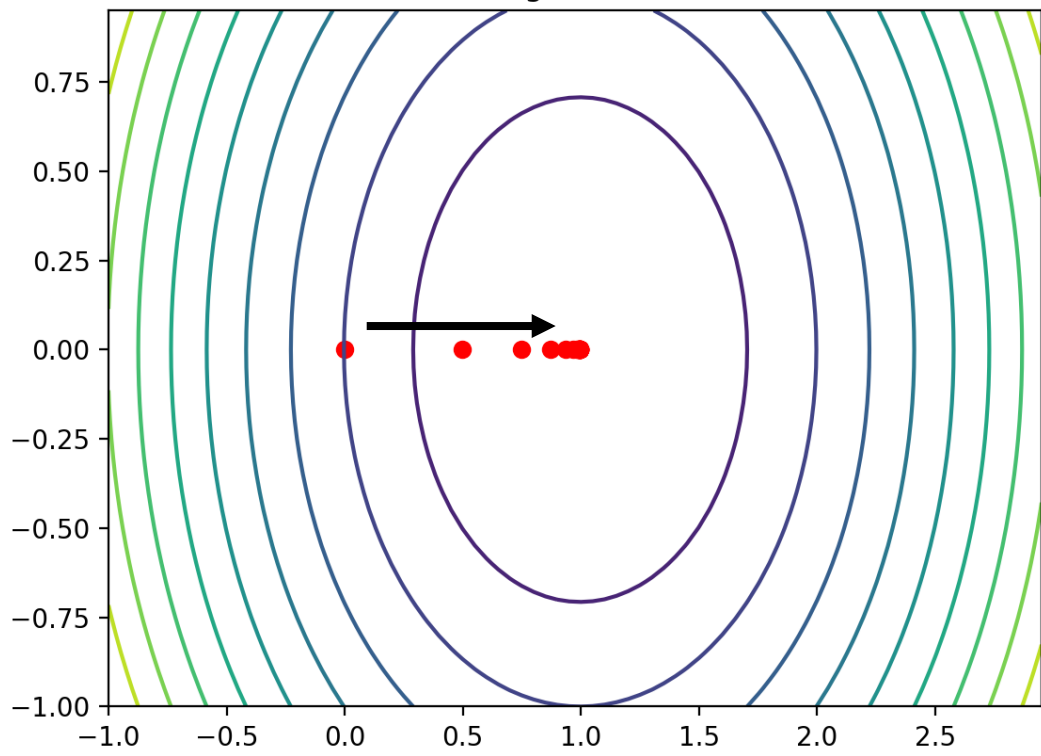
[0.1, 0.5, 0.9] – linear reducing cost function

[1.3, 1.7] – reducing spiral cost function

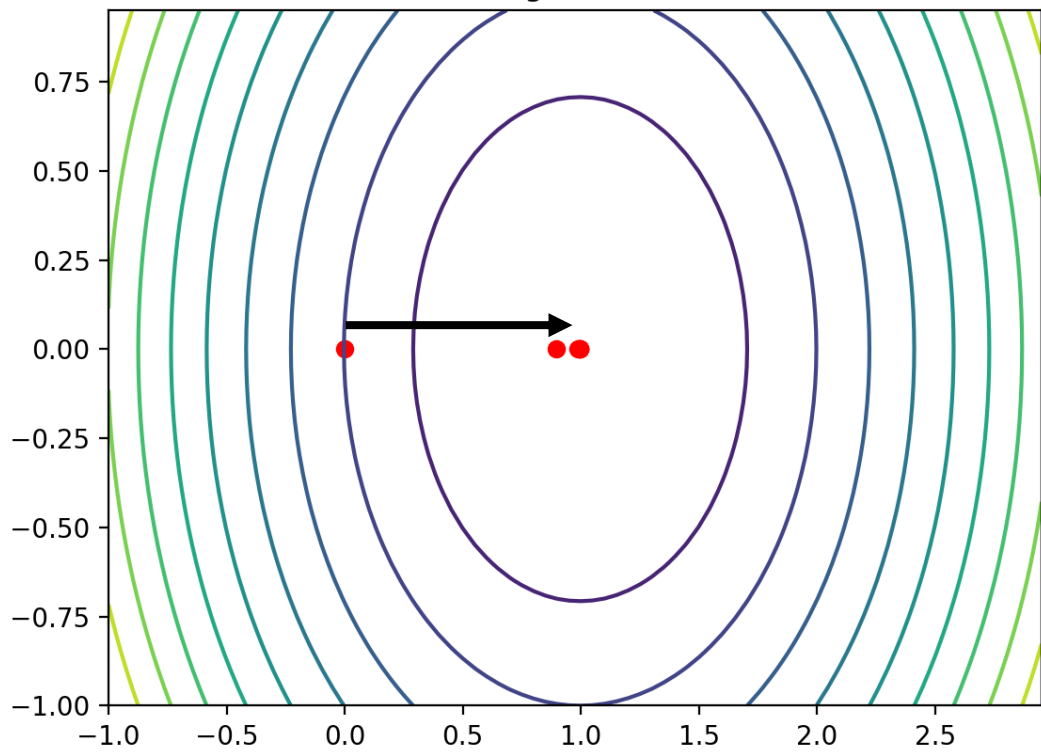
[2.1, 2.5] – increasing spiral cost function



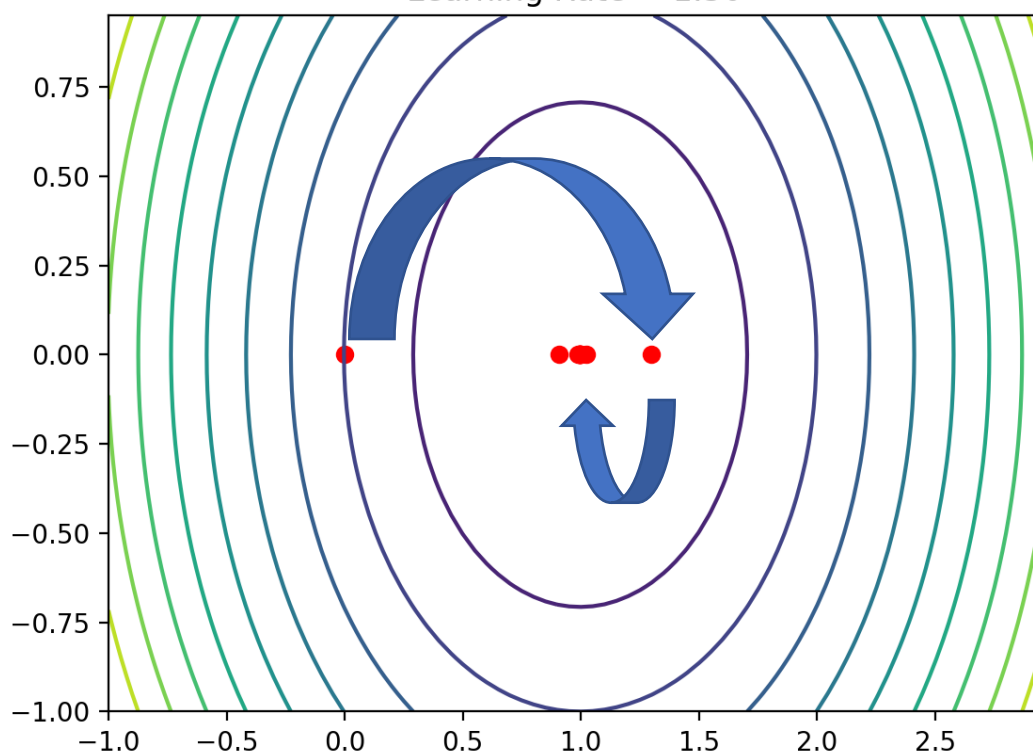
Learning Rate = 0.50



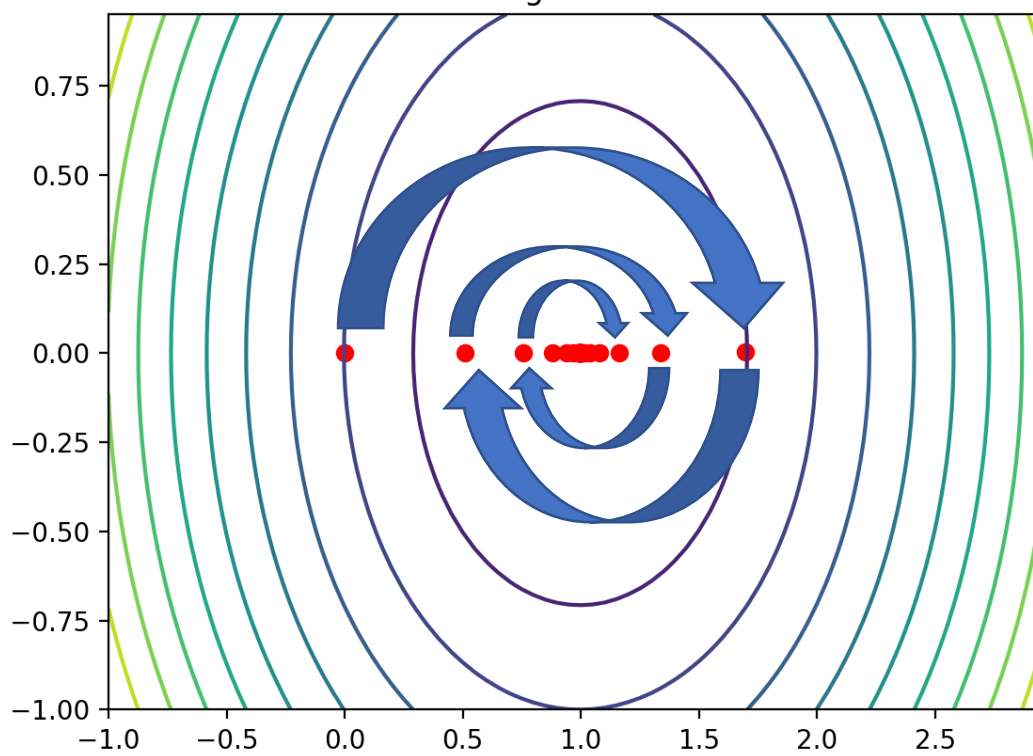
Learning Rate = 0.90



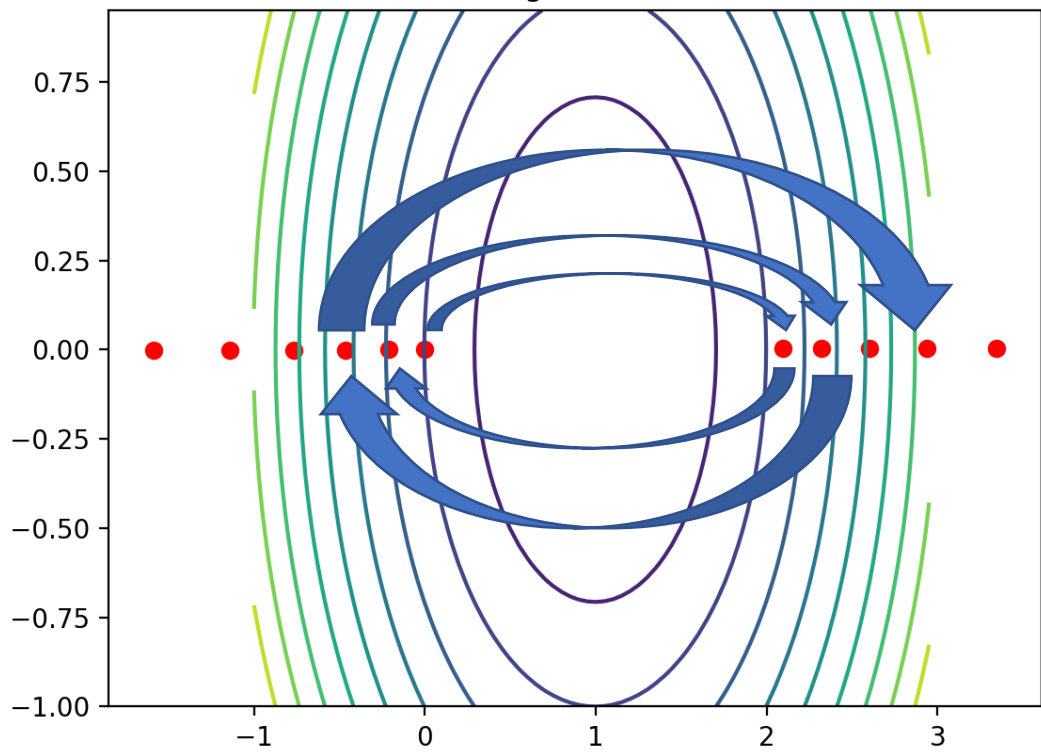
Learning Rate = 1.30



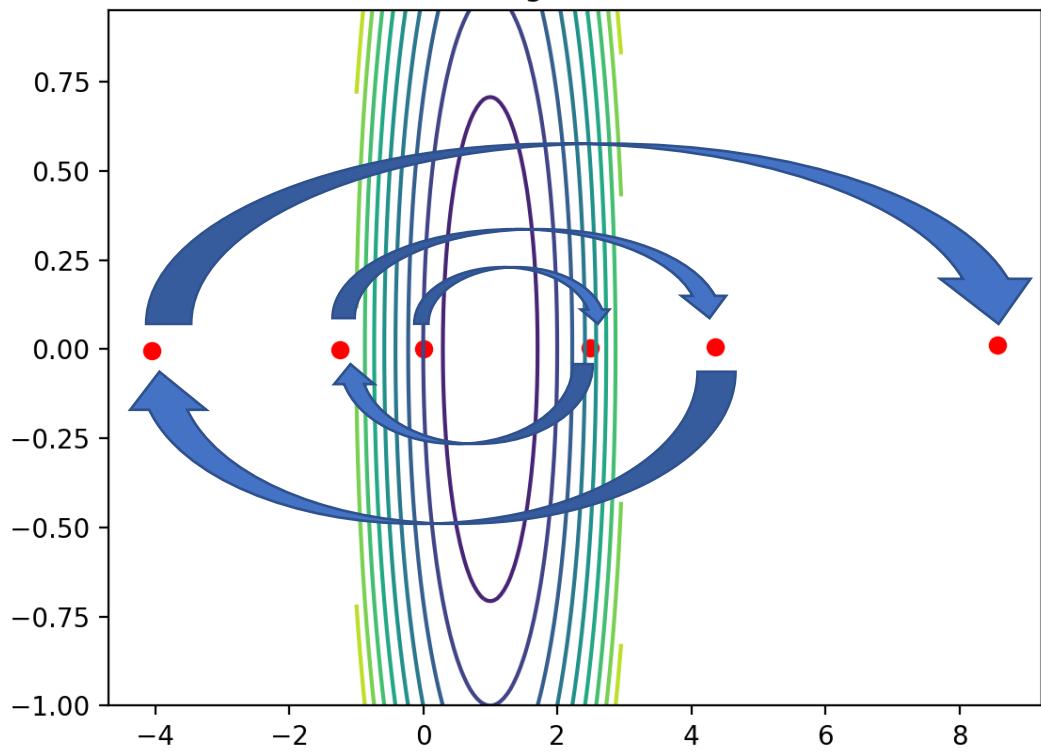
Learning Rate = 1.70



Learning Rate = 2.10



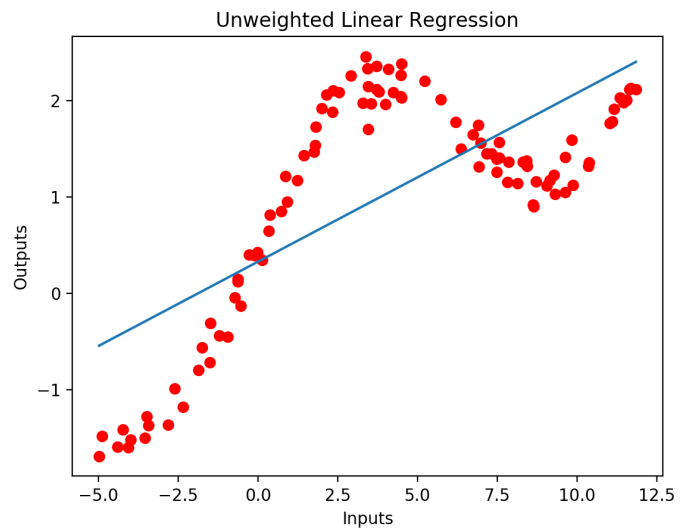
Learning Rate = 2.50



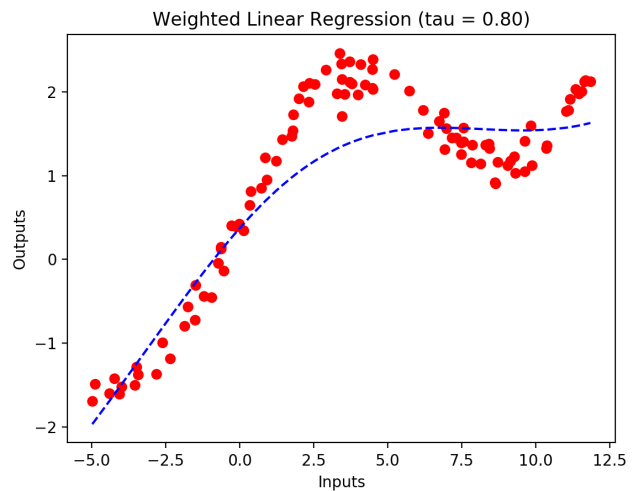
## Part 2:

In this part, we implement weighted linear regression. We calculate the weights matrix for each individual  $X[i]$  point with respect to all  $m$  points in the space. Using this calculated  $W$  matrix, we put it into the derived normal equation formula to get the value of the parameters for each point individually. So, we have multiple sets of  $\theta$ , one for each point in our  $X$ -space.

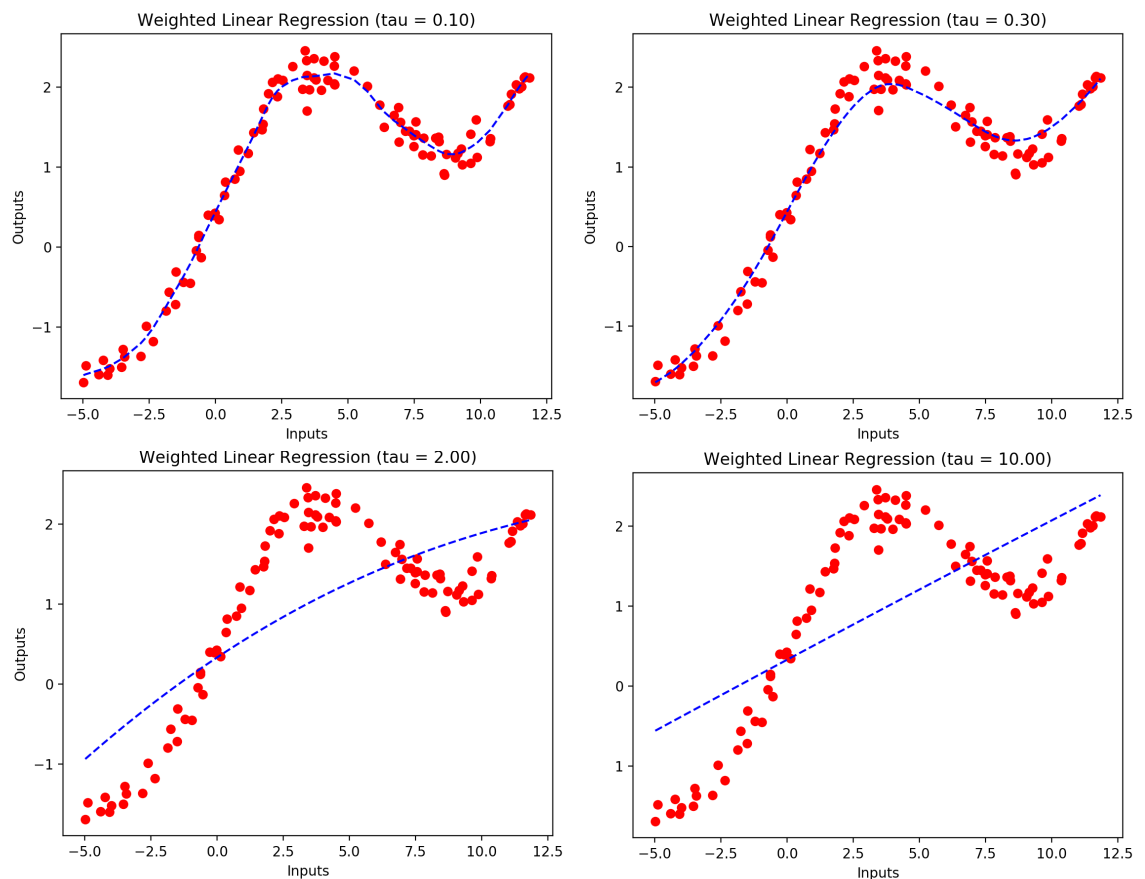
$$\begin{aligned} \text{a) } J(\theta) &= \frac{1}{2m} (X\theta - Y)^T (X\theta - Y) \\ \theta &= (X^T X)^{-1} X^T Y \end{aligned}$$



$$\begin{aligned} \text{b) } J(\theta) &= \frac{1}{2m} (X\theta - Y)^T W (X\theta - Y) \\ \theta &= (X^T (W + W^T) X)^{-1} X^T (W + W^T) Y \end{aligned}$$



- c) When  $\tau$  is too small, we overfit to the data, i.e., the curve follows the data too exactly, and it might not be able to generalize to new points.  
When  $\tau$  is too large, we underfit to the data, i.e., the curve is not a good general representation, and the curve is unable to capture some features of the data.



### Part 3:

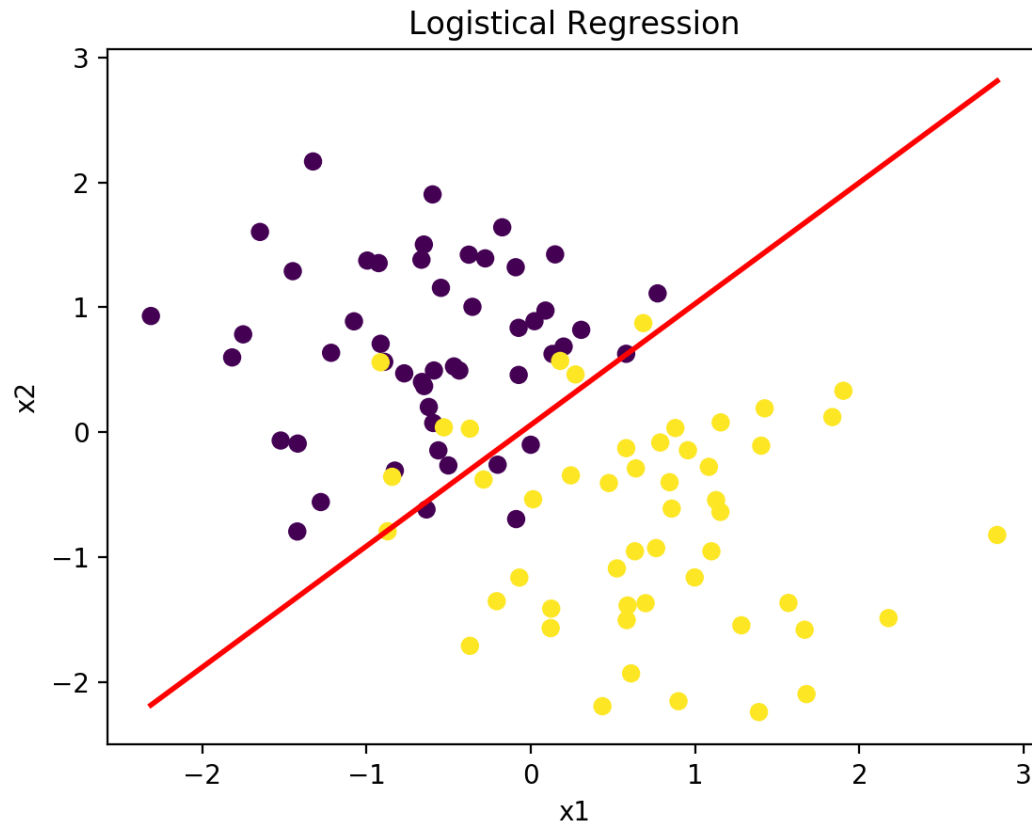
Here we implement the newtons method to calculate the values of the parameters. We define two functions to calculate the gradient and hessian of  $J$  w.r.t. the parameters. For both gradient and hessian, we use the right side derivative since for our chosen initial  $\theta = 0$ , the functions on  $0+\epsilon$  and  $0-\epsilon$  were equal and cancelled out. Similarly for the hessian. This method is much faster (no. of iterations wise) than regular gradient descent but can be costly if there are a large number of parameters.

a)  $\theta := \theta - H^{-1} \nabla_{\theta} \ell(\theta).$

Final Parameters ( $\theta$ ): [0.15547, 2.59393]



b) Data plot with decision boundary



#### Part 4:

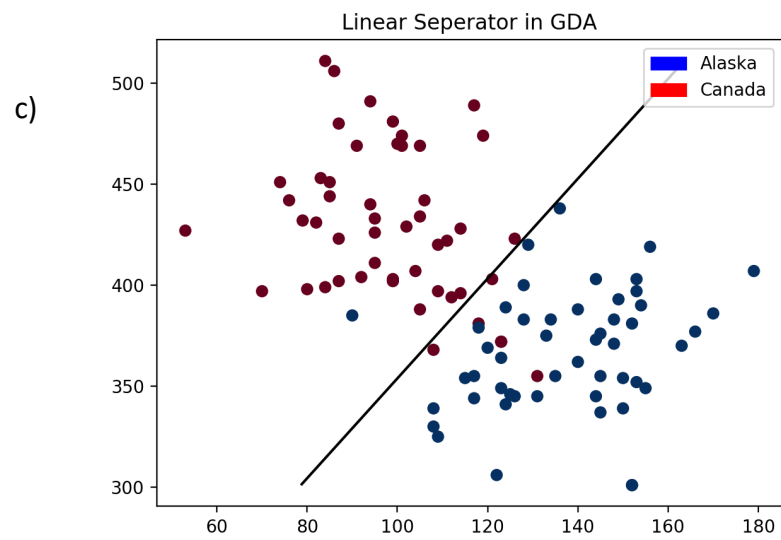
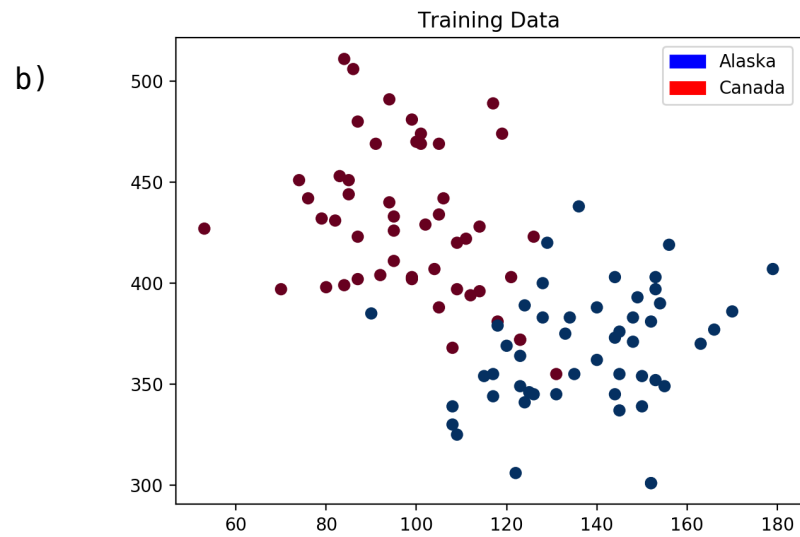
Here we implement Gaussian Discriminant Analysis using both a quadratic and linear separator. First we calculate the values of means and co-variance(s) from the data. Now, using these values we obtain the parameters of the line(linear) and parabola(quadratic) that becomes the decision boundary for the same.

a) Parameter Values:

$$\mu_0 = [98.38, 429.66]$$

$$\mu_1 = [137.46, 366.62]$$

$$\Sigma = \begin{bmatrix} 287.482 & -26.748 \\ -26.748 & 1123.25 \end{bmatrix}$$



d) Parameter Values:

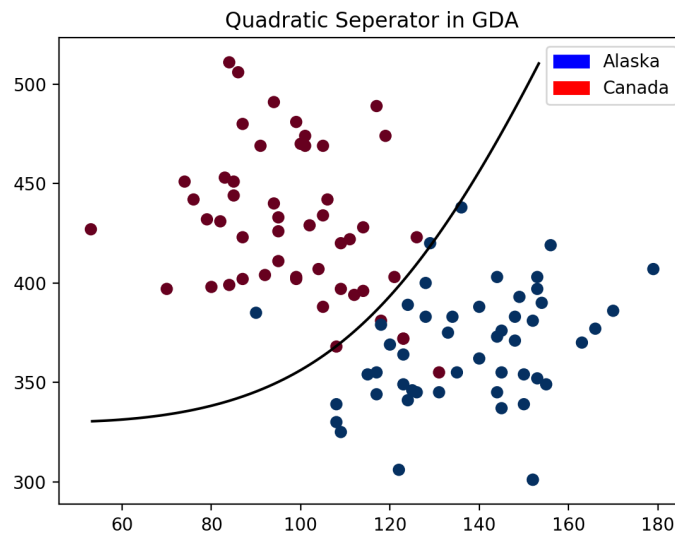
$$\mu_0 = [98.38, 429.66]$$

$$\mu_1 = [137.46, 366.62]$$

$$\Sigma_0 = \begin{bmatrix} 255.3956 & -184.3308 \\ -184.3308 & 1371.1044 \end{bmatrix}$$

$$\Sigma_1 = \begin{bmatrix} 319.5684 & 130.8348 \\ 130.8348 & 875.3956 \end{bmatrix}$$

e)



f) On analysing the two separators, we see that the Quadratic Separator is able to give a slightly better distinction between the two classes.

This is because the linear separator assumes that the observations vary consistently across classes. This might not be the case, and a more general separator (quadratic) which models the covariance of each class independently is able to do a slightly better job.

So we see, that for a generalisation perspective, the quadratic separator is able to make better predictions about the underlying distribution of the individual classes.

---

Files included with this:

1\_part.py  
2\_part.py  
3\_part.py  
4\_part.py  
run.sh