# ECE 595: Homework 4
## Rahul Deshmukh, Class ID: 58
## (Spring 2019)
## In collaboration with: Murali Krishnan & Anmol Singh Guram

**Exercise 1**

**a** Convergence of Logistic Regression:

**(i)** Given: The two datasets are linearly separable.

To Prove: The magnitude of slope and intercept parameters $\boldsymbol{w}$ and $w_0$ would tend to $\infty$

Proof: If the data is linearly separable then we would want the Loss function to be zero. The Loss function for Logistic regression is given by:

$$J(\boldsymbol{\theta}) = -\sum_j (y_j log(h_{\boldsymbol{\theta}}(\boldsymbol{x}_j)) + (1 - y_j)log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}_j)))$$

$$= -\sum_{\boldsymbol{x}_j \in C_1} log(h_{\boldsymbol{\theta}}(\boldsymbol{x}_j)) - \sum_{\boldsymbol{x}_j \in C_0} log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}_j))$$

$$\text{where } h_{\boldsymbol{\theta}}(\boldsymbol{x}_j) = \frac{1}{1 + e^{\boldsymbol{\theta}^T \boldsymbol{x}_j}}$$

We can observe that:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}_j) < 1 \text{ for finite value of } \boldsymbol{\theta}$$
$$\Rightarrow log(h_{\boldsymbol{\theta}}(\boldsymbol{x}_j)) < 0 \ \& \ log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}_j)) < 0$$
$$\Rightarrow J(\boldsymbol{\theta}) > 0$$

Also, we can say the following:

$$\boldsymbol{x}_j \in C_1 \lim_{\boldsymbol{\theta} \to \infty} h_{\boldsymbol{\theta}}(\boldsymbol{x}_j) = 1$$
$$\boldsymbol{x}_j \in C_0 \lim_{\boldsymbol{\theta} \to \infty} h_{\boldsymbol{\theta}}(\boldsymbol{x}_j) = 0$$
$$\Rightarrow \lim_{\boldsymbol{\theta} \to \infty} J(\boldsymbol{\theta}) = 0$$

Therefore, we will attain the minimum zero loss only when $\boldsymbol{\theta} \to \infty$. This means that for any solution $\boldsymbol{\theta}_k$ we can obtain a better solution $\alpha\boldsymbol{\theta}_k$ where $\alpha > 0$ with a smaller loss value. This indicates that the logistic regression will suffer from nonconvergence in case of linearly separable data.

To Prove: Gradient descent iterates would not converge in a finite number of steps if we have the stopping criteria as $||\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)}|| = 0$

Proof:

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \alpha_k \left( \sum_{n=1}^{N} (h_{\boldsymbol{\theta}^{(k)}}(\boldsymbol{x}_n) - y_n) \boldsymbol{x}_n \right)$$

$$\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)} = -\alpha_k \left( \sum_{n=1}^{N} (h_{\boldsymbol{\theta}^{(k)}}(\boldsymbol{x}_n) - y_n) \boldsymbol{x}_n \right)$$

$$\text{if } ||\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)}|| = 0 \Rightarrow h_{\boldsymbol{\theta}^{(k)}}(\boldsymbol{x}_n) = y_n$$

$$\Rightarrow \text{only when } \boldsymbol{\theta} \to \infty$$

$$\Rightarrow k \to \infty$$

$$\text{ie never converges}$$

**(ii)** When we have $||\boldsymbol{w}|| < c_1$ and $|w_0| < c_2$ for some $c_1, c_2 > 0$, then we can say that $\boldsymbol{\theta}$ is no longer unbounded and cannot reach to $\infty$. Then these limits will be used as a termination criteria and thus we will have a converged solution.

Some other ways to counter the non-convergence issue can be to add a regularization term $(\lambda ||\theta||^2)$ to the loss function and solve the unconstrained problem. **Or** To Rather minimize the function $||\theta||^2$ subject to the constraint $J(\boldsymbol{\theta}) < \epsilon$.

**(iii)** We don't face the issue of non-convergence of other linear classifiers when the data is linearly separable this is because for the other methods we can attain a zero loss for linearly separable data within finite number of iterations (which we have proved in class for several methods)

**b** Proof of convergence of online mode of perceptron algorithm:

Given: $\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} + \alpha_k y_j \boldsymbol{x_j}$ with $\alpha_k = 1$ for $j \in \mathcal{M}_k$

**(i)** To prove: The move from $\boldsymbol{w^{(k)}}$ to $\boldsymbol{w^{(k+1)}}$ is in the right direction i.e. decreases the misclassification error.

We need to show that $y_j(\boldsymbol{w^{(k+1)}})^T \boldsymbol{x_j} > y_j(\boldsymbol{w^{(k)}})^T \boldsymbol{x_j}$

Proof:

$$
\begin{aligned}
y_j(\boldsymbol{w}^{(k+1)})^T \boldsymbol{x}_j &= y_j(\boldsymbol{w}^{(k)} + \alpha_k y_j \boldsymbol{x}_j)^T \boldsymbol{x}_j \\
&= y_j(\boldsymbol{w}^{(k)})^T \boldsymbol{x}_j + y_j(\alpha_k y_j \boldsymbol{x}_j)^T \boldsymbol{x}_j \\
&= y_j(\boldsymbol{w}^{(k)})^T \boldsymbol{x}_j + (\alpha_k y_j^2 ||\boldsymbol{x}_j||^2) \\
&\geq y_j(\boldsymbol{w}^{(k)})^T \boldsymbol{x}_j \quad \text{............as } \alpha_k y_j^2 ||\boldsymbol{x}_j||^2 \geq 0 \\
\Rightarrow y_j(\boldsymbol{w}^{(k+1)})^T \boldsymbol{x}_j &\geq y_j(\boldsymbol{w}^{(k)})^T \boldsymbol{x}_j
\end{aligned}
$$
$$\text{hence proved}$$

**(ii)** Proof of convergence:

Given: $\rho = \min_j y_j(\boldsymbol{w^*})^T \boldsymbol{x}_j$

**(1)** To prove: $(\boldsymbol{w^{(k)}})^T \boldsymbol{w}^* \geq k\rho$ (Please note: this is a tighter bound than the homework prompt)

Proof:

$$
\begin{aligned}
(\boldsymbol{w^{(k)}})^T \boldsymbol{w}^* &= (\boldsymbol{w^{(k-1)}} + y_j \boldsymbol{x}_j)^T \boldsymbol{w}^* \\
&= (\boldsymbol{w^{(k-1)}})^T \boldsymbol{w}^* + (y_j \boldsymbol{x}_j)^T \boldsymbol{w}^* \\
&= (\boldsymbol{w^{(k-1)}})^T \boldsymbol{w}^* + y_j(\boldsymbol{w}^*)^T \boldsymbol{x}_j \\
&\text{by definition: } y_j(\boldsymbol{w}^*)^T \boldsymbol{x}_j \geq \rho \\
&\geq (\boldsymbol{w^{(k-1)}})^T \boldsymbol{w}^* + \rho \\
&\text{By induction we can say:} \\
(\boldsymbol{w^{(k)}})^T \boldsymbol{w}^* &\geq (\boldsymbol{w^{(0)}})^T \boldsymbol{w}^* + k\rho \\
&\text{where } \boldsymbol{w^{(0)}} = \boldsymbol{0} \\
\Rightarrow (\boldsymbol{w^{(k)}})^T \boldsymbol{w}^* &\geq k\rho
\end{aligned}
$$
$$\text{hence proved}$$

**(2)** Given: $R = max_j ||\boldsymbol{x}_j||_2$, To prove: $||\boldsymbol{w^{(k)}}||_2^2 \leq k^2 R^2$

Proof:

$$\|\boldsymbol{w}^{(k)}\| = \|\boldsymbol{w}^{(k-1)} + y_j \boldsymbol{x}_j\|$$

using Triangle inequality we get:

$$\leq \|\boldsymbol{w}^{(k-1)}\| + \|y_j \boldsymbol{x}_j\|$$

$$\leq \|\boldsymbol{w}^{(k-1)}\| + \|\boldsymbol{x}_j\| \quad \text{........ as } y_j = \pm 1$$

By induction we have:

$$\leq \|\boldsymbol{w}^{(0)}\| + k\|\boldsymbol{x}_j\|$$

where $\boldsymbol{w}^{(0)} = \boldsymbol{0}$

$$\leq k\|\boldsymbol{x}_j\|$$

$$\leq kR$$

$$\Rightarrow \|\boldsymbol{w}^{(k)}\|^2 \leq k^2 R^2$$

hence proved

**(3)** To Prove: $\dfrac{(\boldsymbol{w}^{(k)})^T \boldsymbol{w}^*}{\|\boldsymbol{w}^{(k)}\|_2} \geq \dfrac{\rho}{R}$ and $k \leq \dfrac{R^2 \|\boldsymbol{w}^*\|}{\rho^2}$ (Please note: this is a tighter bound than the homework prompt)

$$(\boldsymbol{w}^{(k)})^T \boldsymbol{w}^* \geq k\rho \ \text{(proved in part (1))}$$

$$\boldsymbol{w}^{(k)} \leq kR \ \text{(proved in part (2))}$$

$$\Rightarrow \dfrac{(\boldsymbol{w}^{(k)})^T \boldsymbol{w}^*}{\|\boldsymbol{w}^{(k)}\|_2} \geq \dfrac{k\rho}{kR} = \dfrac{\rho}{R}$$

hence proved

Now let us try to prove that the sequence $\boldsymbol{w}^{(k)}$ converges to $\boldsymbol{w}^*$. Therefore we need to prove the following:

$$\exists\, \epsilon > 0 \mid \|w^{(k)} - w^*\| < \epsilon$$

4

Proof:

$$||\boldsymbol{w}^{(k)} - \boldsymbol{w}^*||^2 = ||\boldsymbol{w}^{(k-1)} + \alpha_k y_j \boldsymbol{x}_j - \boldsymbol{w}^*||^2$$
$$= ||\boldsymbol{w}^{(k-1)} - \boldsymbol{w}^* + \alpha_k y_j \boldsymbol{x}_j||^2$$
$$= ||\boldsymbol{w}^{(k-1)} - \boldsymbol{w}^*||^2 + ||\alpha_k y_j \boldsymbol{x}_j||^2 + 2\alpha_k (\boldsymbol{w}^{(k-1)} - \boldsymbol{w}^*)^T y_j \boldsymbol{x}_j$$

we know that at k-1 iteration the misclassification error is greater that 0 ie $(\boldsymbol{w}^{(k-1)})^T y_j \boldsymbol{x}_j \leq 0$

$$\leq ||\boldsymbol{w}^{(k-1)} - \boldsymbol{w}^*||^2 + \alpha_k^2 ||y_j \boldsymbol{x}_j||^2 - 2\alpha_k (\boldsymbol{w}^*)^T y_j \boldsymbol{x}_j$$
$$\leq ||\boldsymbol{w}^{(k-1)} - \boldsymbol{w}^*||^2 + \alpha_k (\alpha_k ||y_j \boldsymbol{x}_j||^2 - 2(\boldsymbol{w}^*)^T y_j \boldsymbol{x}_j)$$

if $\alpha_k ||y_j \boldsymbol{x}_j||^2 - 2(\boldsymbol{w}^*)^T y_j \boldsymbol{x}_j \leq 0$

ie $\alpha_k \leq \dfrac{2(\boldsymbol{w}^*)^T y_j \boldsymbol{x}_j}{||y_j \boldsymbol{x}_j||^2}$

choosing $\alpha_k = \dfrac{(\boldsymbol{w}^*)^T y_j \boldsymbol{x}_j}{||y_j \boldsymbol{x}_j||^2}$ we get:

$$||\boldsymbol{w}^{(k)} - \boldsymbol{w}^*||^2 \leq ||\boldsymbol{w}^{(k-1)} - \boldsymbol{w}^*||^2 - \dfrac{((\boldsymbol{w}^*)^T y_j \boldsymbol{x}_j)^2}{||y_j \boldsymbol{x}_j||^2}$$

By induction we get:

$$||\boldsymbol{w}^{(k)} - \boldsymbol{w}^*||^2 \leq ||\boldsymbol{w}^{(0)} - \boldsymbol{w}^*||^2 - k\dfrac{((\boldsymbol{w}^*)^T y_j \boldsymbol{x}_j)^2}{||y_j \boldsymbol{x}_j||^2}$$

where $\boldsymbol{w}^{(0)} = \boldsymbol{0}$

$$||\boldsymbol{w}^{(k)} - \boldsymbol{w}^*||^2 \leq ||\boldsymbol{w}^*||^2 - k\dfrac{((\boldsymbol{w}^*)^T y_j \boldsymbol{x}_j)^2}{||y_j \boldsymbol{x}_j||^2} = \epsilon^2$$

Therefore, The sequence the sequence $\boldsymbol{w}^{(k)}$ converges to $\boldsymbol{w}^*$ if $\epsilon^2 > 0$:

$$\Rightarrow ||\boldsymbol{w}^*||^2 - k\dfrac{((\boldsymbol{w}^*)^T y_j \boldsymbol{x}_j)^2}{||y_j \boldsymbol{x}_j||^2} > 0$$
$$||\boldsymbol{w}^*||^2 > k\dfrac{((\boldsymbol{w}^*)^T y_j \boldsymbol{x}_j)^2}{||y_j \boldsymbol{x}_j||^2}$$
$$k < \dfrac{||\boldsymbol{w}^*||^2 ||y_j \boldsymbol{x}_j||^2}{((\boldsymbol{w}^*)^T y_j \boldsymbol{x}_j)^2}$$
$$k < \dfrac{||\boldsymbol{w}^*||^2 ||\boldsymbol{x}_j||^2}{((\boldsymbol{w}^*)^T y_j \boldsymbol{x}_j)^2}$$

using earlier definitions $||\boldsymbol{x}_j|| \leq R$ and $(\boldsymbol{w}^*)^T y_j \boldsymbol{x}_j \geq \rho$

$$k \leq \dfrac{||\boldsymbol{w}^*||^2 R^2}{\rho^2}$$

hence proved

**c** Soft-Margin SVM:

$$\operatorname*{argmin}_{\boldsymbol{\theta},\boldsymbol{\xi}} \frac{1}{2}||\boldsymbol{w}||_2^2 + \frac{C}{2}||\boldsymbol{\xi}||_2^2$$
$$\text{subject to } y_j g_{\boldsymbol{\theta}}(\boldsymbol{x}_j) \geq 1 - \xi_j$$
$$\xi_j \geq 0, j = 1, ...., N$$

We need to find the convex dual problem for it

**(i)** The constraint $\xi_j \geq 0$ can be removed without affecting the solution to the optimization problem, why? lets write ou the lagrangian of the function:

$$L(\theta, \lambda, \xi, \mu) = \frac{1}{2}||\boldsymbol{w}||_2^2 + \frac{C}{2}||\boldsymbol{\xi}||_2^2 + \sum_j (\lambda_j(1 - \xi_j - y_j(\boldsymbol{w}^T\boldsymbol{x}_j + w_0)) + \mu_j(-\xi_j))$$

Form complementary slackness condition we have:

$$\mu_j(\xi_j) = 0$$
$$\Rightarrow \mu_j = 0, \xi_j \geq 0$$
$$\text{or } \mu_j \geq 0, \xi_j = 0$$

Therefore the constraint $\xi_j \geq 0$ is always satisfied and using the above fact we can simplify the lagrangian to:

$$L(\theta, \lambda, \xi) = \frac{1}{2}||\boldsymbol{w}||_2^2 + \frac{C}{2}||\boldsymbol{\xi}||_2^2 + \sum_j (\lambda_j(1 - \xi_j - y_j(\boldsymbol{w}^T\boldsymbol{x}_j + w_0)))$$

**(ii)** Finding the optimal solutions:

$$L(\theta, \lambda, \xi) = \frac{1}{2}||\boldsymbol{w}||_2^2 + \frac{C}{2}||\boldsymbol{\xi}||_2^2 + \sum_j (\lambda_j(1 - \xi_j - y_j(\boldsymbol{w}^T\boldsymbol{x}_j + w_0)))$$
$$\nabla_{\boldsymbol{w}} L = 0 \Rightarrow \boldsymbol{w}^* - \sum_j \lambda_j y_j \boldsymbol{x}_j = 0$$
$$\Rightarrow \boldsymbol{w}^* = \sum_j \lambda_j y_j \boldsymbol{x}_j$$
$$\nabla_{w_0} L = 0 \Rightarrow \sum_j \lambda_j y_j = 0$$
$$\nabla_{\boldsymbol{\xi}} L = 0 \Rightarrow C\boldsymbol{\xi} - \boldsymbol{\lambda} = 0$$
$$\Rightarrow C\boldsymbol{\xi} = \boldsymbol{\lambda} \Rightarrow \boldsymbol{\xi}^* = \frac{\boldsymbol{\lambda}}{C}$$

hence proved

**(iii)** Finding the convex dual of the problem:

$$\underset{\boldsymbol{w}, w_0}{\text{minimize}}\, \mathcal{L}(\boldsymbol{w}, w_0, \lambda)$$

$$= \underset{\boldsymbol{w}, w_0}{\text{minimize}} \left\{ \frac{1}{2}||\boldsymbol{w}||_2^2 + \frac{C}{2}||\boldsymbol{\xi}||_2^2 + \underset{\lambda \geq 0}{\max} \left\{ \sum_j (\lambda_j (1 - \xi_j - y_j(\boldsymbol{w}^T \boldsymbol{x}_j + w_0))) \right\} \right\}$$

$$= \underset{\boldsymbol{w}, w_0}{\text{minimize}} \{ \underset{\lambda}{\text{maximize}}\, \mathcal{L}(\boldsymbol{w}, w_0, \lambda) \}$$

$$= \underset{\lambda}{\text{maximize}} \{ \underset{\boldsymbol{w}, w_0}{\text{minimize}}\, \mathcal{L}(\boldsymbol{w}, w_0, \lambda) \}$$

$$= \underset{\lambda}{\text{maximize}} \{ \mathcal{L}(\boldsymbol{w}^*, w_0^*, \lambda) \}$$

$$= \underset{\lambda}{\text{maximize}} \{ \frac{1}{2}||\boldsymbol{w}^*||_2^2 + \frac{C}{2}||\boldsymbol{\xi}^*||_2^2 + \sum_j (\lambda_j (1 - \xi_j^* - y_j(\boldsymbol{w}^{*T} \boldsymbol{x}_j + w_0^*))) \}$$

$$= \underset{\lambda}{\text{maximize}} \left\{ + \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \boldsymbol{x}_i^T \boldsymbol{x}_j + \frac{C}{2} \sum_k \frac{\lambda_k}{C}^2 + \sum_j \lambda_j \right.$$

$$\left. - \sum_k \frac{\lambda_k^2}{C} - \sum_i \sum_j \lambda_i \lambda_j y_i y_j \boldsymbol{x}_i^T \boldsymbol{x}_j - w_0^* \sum_j (\lambda_j y_j) \right\}$$

$$= \underset{\lambda}{\text{maximize}} \left\{ - \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \boldsymbol{x}_i^T \boldsymbol{x}_j - \frac{1}{2} \sum_k \frac{\lambda_k^2}{C} + \sum_j \lambda_j \right\}$$

$$\text{subject to } \sum_j \lambda_j y_j = 0$$

hence proved

**d** Hard- Margin SVM:

$$\operatorname*{argmin}_{\boldsymbol{\theta}} \frac{1}{2}||\boldsymbol{w}||_2^2$$

$$\text{st } y_j g_{\boldsymbol{\theta}}(\boldsymbol{x_j}) \geq 0, j = 1, ......, N$$

**(i)** To prove: if the constraint is changed to $y_j g_{\boldsymbol{\theta}}(\boldsymbol{x}_j) \geq \gamma$ such that $\gamma \geq 0$ then the problem remains the same.

Proof:

$$\operatorname*{argmin}_{\boldsymbol{\theta}} \frac{1}{2}||\boldsymbol{w}||_2^2$$

$$y_j g_{\boldsymbol{\theta}}(\boldsymbol{x}_j) \geq \gamma$$
$$\Rightarrow y_j(\boldsymbol{w}^T \boldsymbol{x}_j + w_0)/\gamma \geq 1$$
$$\Rightarrow y_j(\frac{\boldsymbol{w}^T}{\gamma}\boldsymbol{x}_j + \frac{w_0}{\gamma}) \geq 1$$

We can scale the objective function by a constant $\frac{1}{\gamma^2}$ without changing the problem

$$\Rightarrow \operatorname*{argmin}_{\boldsymbol{\theta}} \frac{1}{2}\frac{||\boldsymbol{w}||_2^2}{\gamma^2}$$

$$y_j(\frac{\boldsymbol{w}^T}{\gamma}\boldsymbol{x}_j + \frac{w_0}{\gamma}) \geq 1$$

$$\Rightarrow \operatorname*{argmin}_{\hat{\boldsymbol{\theta}}=\frac{\boldsymbol{\theta}}{\gamma}} \frac{1}{2}||\hat{\boldsymbol{w}}||_2^2$$

$$y_j(\hat{\boldsymbol{w}}^T \boldsymbol{x}_j + \hat{w}_0) \geq 1$$

This is the same problem as earlier. Hence proved.

**(ii)** To prove : SVM for two points $\boldsymbol{x}_1 \in C_+$ and $\boldsymbol{x}_2 \in C_-$ is solvable.

Proof:

The Dual Problem is:

$$\max_{\lambda \geq 0}\left\{\sum_{j=1}^{2}\lambda_j - \frac{1}{2}\sum_{i=1}^{2}\sum_{j=1}^{2}(\lambda_i\lambda_j y_i y_j \boldsymbol{x}_i^T \boldsymbol{x}_j)\right\}$$

$$\text{st} \sum_{j=1}^{2}\lambda_j y_j = 0$$

$$\Rightarrow \max_{\lambda \geq 0}\left\{\sum_{j=1}^{2}\lambda_j - \frac{1}{2}\sum_{i=1}^{2}\sum_{j=1}^{2}(\lambda_i\lambda_j y_i y_j \boldsymbol{x}_i^T \boldsymbol{x}_j)\right\}$$

$$\text{st} \sum_{j=1}^{2}\lambda_j y_j = 0$$

$$\Rightarrow \lambda_1 - \lambda_2 = 0$$

$$\Rightarrow \lambda_1 = \lambda_2 = \lambda$$

$$\Rightarrow \max_{\lambda \geq 0}\left\{2\lambda - \frac{\lambda^2}{2}(\boldsymbol{x}_1^T \boldsymbol{x}_1 + \boldsymbol{x}_2^T \boldsymbol{x}_2 - 2\boldsymbol{x}_1^T \boldsymbol{x}_2)\right\}$$

taking derivative and equating to zero we get

$$2 - \lambda(||\boldsymbol{x}_1 - \boldsymbol{x}_2||^2) = 0$$

$$\lambda = \frac{2}{||\boldsymbol{x}_1 - \boldsymbol{x}_2||^2} \geq 0$$

$$\boldsymbol{w}^* = \sum_{j=1}^{2}\lambda_j y_j \boldsymbol{x}_j$$

$$\Rightarrow \boldsymbol{w}^* = \lambda(\boldsymbol{x}_1 - \boldsymbol{x}_2)$$

$$= \frac{2(\boldsymbol{x}_1 - \boldsymbol{x}_2)}{||\boldsymbol{x}_1 - \boldsymbol{x}_2||^2}$$

$$w_0 = -\frac{\boldsymbol{w}^{*T}(\boldsymbol{x}_1 + \boldsymbol{x}_2)}{2}$$

$$\Rightarrow w_0 = -\frac{1}{2}\lambda(\boldsymbol{x}_1 - \boldsymbol{x}_2)^T(\boldsymbol{x}_1 + \boldsymbol{x}_2)$$

$$w_0 = -\frac{(||\boldsymbol{x}_1||^2 - ||\boldsymbol{x}_2||^2)}{||\boldsymbol{x}_1 - \boldsymbol{x}_2||^2}$$
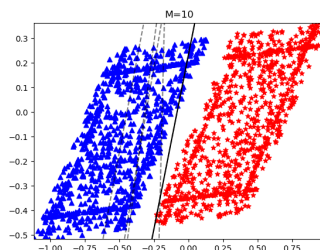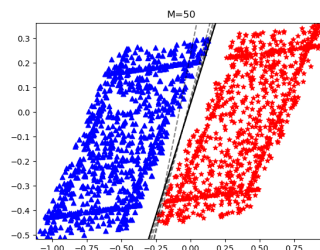
Hence proved.

# Exercise 2

**(a)** Please find the code at page 15

**(b)** Please find the code at page 15

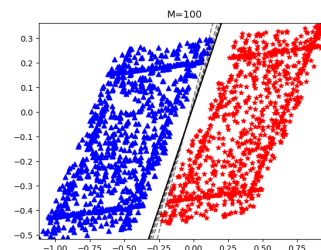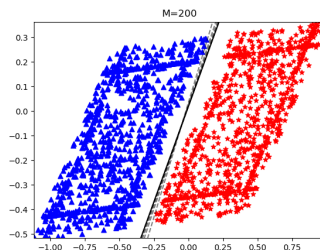The Learning rate was chosen as $10^{-1}$



(a) M=10
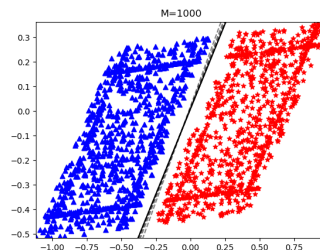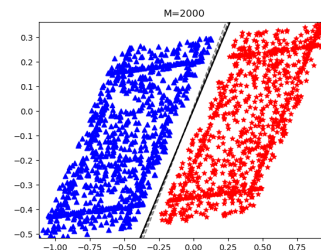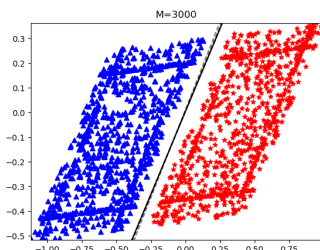
(b) M=50
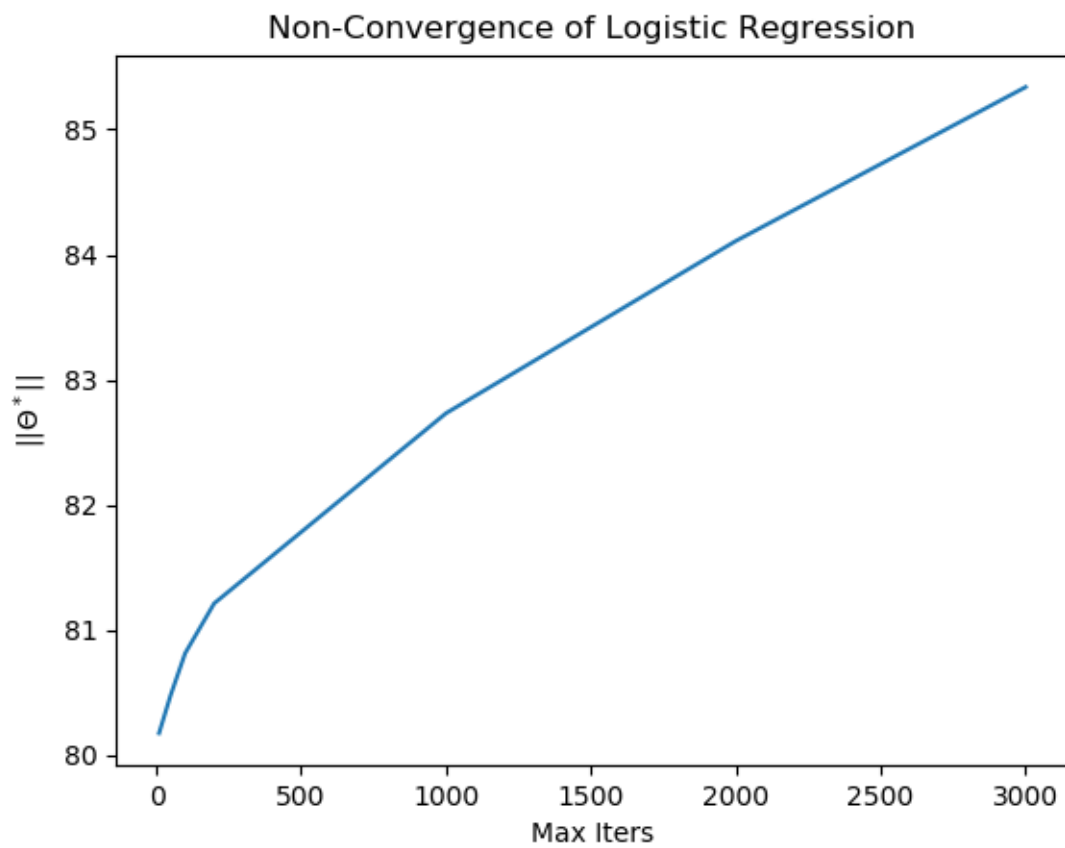
(c) M=100

(d) M=200

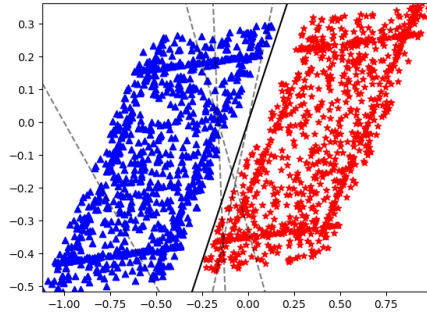(e) M=1000

(f) M=2000

(g) M=3000

As can be observed, the decision boundary improves as the maximum number of iterations increases however, I observed that after attaining a good separating hyperplane there is no change in the decision boundary as we further increase the number of iterations.
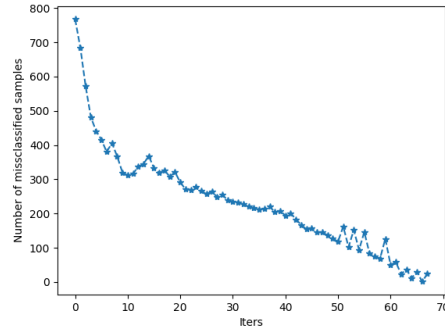
Non-Convergence of Logistic Regression

We can observe from the above plot that the magnitude of theta $||\boldsymbol{\theta}||$ monotonically increases as the number of iterations increases this further confirms our proof in exercise 1(a).

**b** Perceptron: Please find the code at page 15
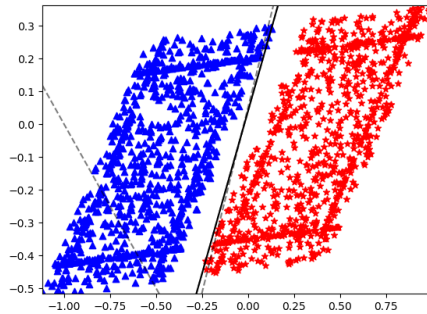
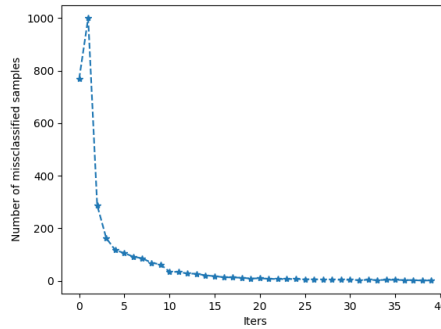**(i)** Online Mode



(h) decision boundary plot

(i) Convergence
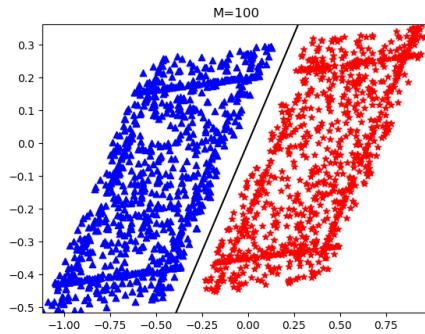
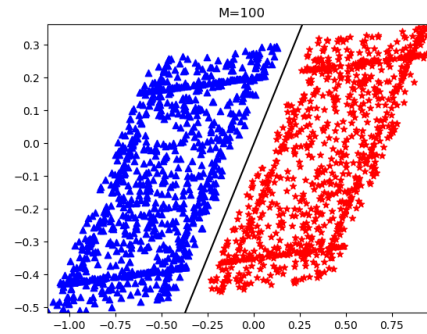**(i)** Batch Mode



(j) decision boundary plot

(k) Convergence

From the above plot for convergence of both the modes we observe that convergence is faster for Batch Mode. However, in terms of cost of computation online mode is cheaper than batch mode.
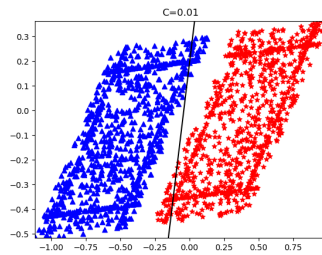
**c** SVM, Please find the code at page 15



(l) Hard



(m) soft

We can observe that the decision boundaries for the two cases are nearly. More insights can be gained from calculating the margin for the two cases. We can observe from the above plots that



(n) C=0.01



(o) C=0.1



(p) C=1



(q) C=10

for the soft margin as C increases the decision boundary starts resembling the hard SVM decision boundary. Also, as the data is linearly separable with a good margin, therefore using a hard margin SVM is better suited.

**d** Classification Error: For all of the classifiers we saw(from earlier plots) that we were able to get correct classification.



(r) Computational Time
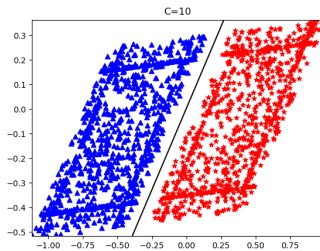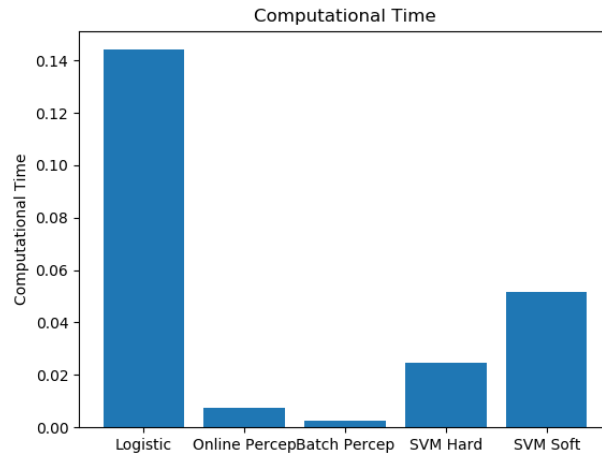
Computational Time: We can observe that Logistic regression takes the maximum time to get a good solution as it requires more number of iterations. Also, The perceptron methods are cheaper than SVM because in perceptron we are not doing additional work of maximizing the margin.

Effect of Learning rate: It was observed that a smaller learning rate required more iterations to converge and a faster learning rate can give erroneous results.



(s) $\gamma_{unsigned}$

From the above plot we can observe that Perceptron Methods are the least robust and can be a poor method to classify a test dataset. Of all the classifiers SVM hard seems to be the most robust. This is because the dataset is linearly seperable and has no outliers, had there been outliers then SVM soft would have been the best option.

14

## Main Code for Exercise 2

Listing 1: Source Code

```python
#!/usr/bin/env python3
# -*- coding: utf-8 -*-
"""
Created on Mon Mar 18 16:47:50 2019
@author: rahul
"""
# import libraries
import numpy as np , csv
from functions import *
import matplotlib.pyplot as plt
#%% Read Data
#(a)
read_path = '../data/'
label_filename = 'hw04_labels.csv'
sample_filename = 'hw04_sample_vectors.csv'

samples = []
with open(read_path+sample_filename) as csvfile:
    readcsv = csv.reader(csvfile,delimiter=',')
    for row in readcsv:
        vals = list(map(float ,row))
        samples.append(vals)
samples_copy = np.array(samples)
samples = np.ones((np.shape(samples_copy)[0],np.shape(samples_copy)[1]+1))
samples[:,:-1]=samples_copy
samples = samples.T

labels = []
with open(read_path+label_filename) as csvfile:
    readcsv = csv.reader(csvfile,delimiter=',')
    for row in readcsv:
        vals = float(row[0])
        labels.append(vals)
labels = np.array(labels)

scaled_labels = 2*labels-1 #percep labels in range -1 to +1
# Declare constants

rate = 0.1
#%%
#(a) Logistic regression based classification
log_star_store =[]
Ms= [10,50,100,200,1000,2000,3000]

for M in Ms:
    [log_theta_star ,log_theta_store] = logistic(samples,labels,rate ,M)
    log_star_store.append(np.linalg.norm(log_theta_star))
    freq = int(0.2*M)
#      plotdata(samples,labels,log_theta_store,freq ,M)
#      plotdata_single(samples,labels,log_theta_star,M)
    log_gamma = min(scaled_labels*(log_theta_star.T@samples)/np.linalg.norm(
        log_theta_star))

#plt.figure()
#plt.plot(Ms,log_star_store)
```

15

```python
55  #plt.xlabel('Max Iters')
56  #plt.ylabel('||$\Theta^*$||')
57  #plt.title('Non-Convergence of Logistic Regression')
58  #plt.show()
59
60  #%%(b) Perceptron Online mode
61  M = 100 # max number of iters
62  freq = int(0.2*M)
63  scaled_labels = 2*labels-1 #percep labels in range -1 to +1
64  ##(i)
65  [ol_percp_theta_star,ol_percp_theta_store]= perceptron(samples,scaled_labels,rate,M,
        convplot=0)
66  #plotdata(samples,labels,ol_percp_theta_store,freq,[M])
67  ol_gamma = min((scaled_labels*(ol_percp_theta_star.T@samples))/np.linalg.norm(
        ol_percp_theta_star))
68
69
70  ##(ii)
71  [bt_percp_theta_star,bt_percp_theta_store]= perceptron(samples,scaled_labels,rate,M,
        convplot=0,online=False)
72  #plotdata(samples,labels,bt_percp_theta_store,freq)
73  bt_gamma = min(scaled_labels*(bt_percp_theta_star.T@samples)/np.linalg.norm(
        bt_percp_theta_star))
74
75  #%%(c) SVM
76  #(i) Hard Margin
77  svm_hard_theta_star = SVM_hard(samples,scaled_labels)
78  #plotdata_single(samples,labels,svm_hard_theta_star,M)
79  hard_gamma = min((scaled_labels*(np.array(svm_hard_theta_star).T@samples))/np.linalg
        .norm(svm_hard_theta_star))
80
81  ##(ii) Soft Margin
82  C = 1
83  #for C in [10**(-2),10**(-1),1,10]:
84  svm_soft_theta_star = SVM_soft_L1(samples,scaled_labels,C)
85  #plotdata_single(samples,labels,svm_soft_theta_star,C)
86  soft_gamma = min((scaled_labels*(np.array(svm_soft_theta_star).T@samples))/np.linalg
        .norm(svm_hard_theta_star))
87
88  #plt.figure()
89  #plt.bar(['Logistic','Online Percep','Batch Percep','SVM Hard','SVM Soft'],\
90  #        [log_gamma,ol_gamma,bt_gamma,hard_gamma,soft_gamma])
91  #plt.ylabel('$\gamma_{unsigned}$')
92  #plt.title('Robustness')
93  #plt.show()
```

## Functions for Exercise 2

Listing 2: Functions

```python
1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Created on Mon Mar 18 17:07:13 2019
5
6  @author: rahul
7  """
8  # import libraries
```

```python
 9  import numpy as np
10  import matplotlib.pyplot as plt
11  import cvxpy as cvx
12  import time
13  #%% SVM
14  # SVM Soft Margin
15  def SVM_soft_L1(samples, labels, C):
16      """
17      Input:
18      sample: np array in the format [[x1],[x2],....,[xn]] xi as col vectors
19      labels: np array [y1,y2,....,yn]
20      """
21      t_in = time.time()
22      dim, Nsamples = samples.shape
23      theta = cvx.Variable(dim) # declaring dimension of variable
24      lam = 1/C
25      Jump_term = cvx.max_elemwise(0,1-cvx.mul_elemwise(labels,(theta.T@samples).T))
26      Jump_term = cvx.sum_entries(Jump_term)
27      obj_expr = Jump_term +(lam/2)*(cvx.sum_squares(theta[:-1]))
28      obj = cvx.Minimize(obj_expr)
29      prob = cvx.Problem(obj)
30      prob.solve(solver = cvx.ECOS)
31      theta_star = theta.value
32      theta_star = theta_star.tolist()
33      theta_star= [x[0] for x in theta_star]
34      print('SVM soft time '+str(time.time() - t_in)+'\n')
35      return(theta_star)
36
37  # SVM Hard Margin
38  def SVM_hard(samples, labels):
39      """
40      Input:
41      sample: np array in the format [[x1],[x2],....,[xn]] xi as col vectors
42      labels: np array [y1,y2,....,yn]
43      """
44      t_in = time.time()
45      dim, Nsamples = samples.shape
46      theta = cvx.Variable(dim) # declaring dimension of variable
47      obj = cvx.Minimize(cvx.sum_squares(theta[:-1]))
48      temp = cvx.mul_elemwise(labels,(theta.T@samples).T)
49      const = [ temp >= 1 ]
50      prob = cvx.Problem(obj, const)
51      prob.solve(solver = cvx.ECOS)
52      theta_star = theta.value
53      #format change
54      theta_star = theta_star.tolist()
55      theta_star= [x[0] for x in theta_star]
56      print('SVM hard time '+str(time.time() - t_in)+'\n')
57      return(theta_star)
58
59
60  #%% Perceptron Method
61  def  percept_batch_tangent(theta, samples, true_labels):
62      """
63      Input:
64          theta: decision boundary
65          sample: np array in the format [[x1],[x2],....,[xn]] xi as col vectors
66          labels: np array [y1,y2,....,yn]
67      Out: J: Jacobian
```

```python
68          """
69          # predict samples labels using theta
70          gx = theta.T@samples
71          ygx = np.multiply(true_labels,gx)
72          miss_idx = np.where(ygx<0)[0]
73          N_miss = miss_idx.size
74          all_miss_labels = true_labels[miss_idx]
75          all_miss_x = samples[:,miss_idx]
76          J = np.sum(all_miss_labels*all_miss_x,axis=1)
77          return(J,N_miss)
78
79    def  percept_online_tangent(theta,samples,true_labels):
80          """
81          Input:
82              theta: decision boundary
83              sample: np array in the format [[x1],[x2],....,[xn]] xi as col vectors
84              labels: np array [y1,y2,....,yn]
85          Out: J: Jacobian
86          """
87          # predict samples labels using theta
88          gx = theta.T@samples
89          ygx = np.multiply(true_labels,gx)
90          miss_idx = np.where(ygx<0)[0]
91          N_miss = miss_idx.size
92          if N_miss>0:
93              picked_idx = miss_idx[np.random.permutation(N_miss)]
94              picked_idx = picked_idx[0]
95              J = true_labels[picked_idx]*samples[:,picked_idx]
96          else:
97              J = np.zeros_like(samples[:,0]) #just to pass some value
98          return(J,N_miss)
99
100   def perceptron(samples,labels,rate,Max_iter,convplot,online = True):
101         """
102         Input:
103             samples: np array in the format [[x1],[x2],....,[xn]] xi as col vectors
104             labels: np array [y1,y2,....,yn]
105             rate: learning rate or the step length
106             Max_iter: for the gradient descent
107             online: if online mode then True(default) else False
108         """
109         t_in = time.time()
110         xdim,Nsamples= samples.shape
111         theta_k = np.ones(xdim) # initial guess
112         theta_store = []
113         N_store = []
114         theta_store.append(np.copy(theta_k))
115         if online:
116             for k in range(Max_iter):
117                 grad_k,N_miss = percept_online_tangent(theta_k,samples,labels)
118                 if N_miss == 0:
119                     break
120                 theta_k += rate*grad_k
121                 theta_store.append(np.copy(theta_k))
122                 N_store.append(N_miss)
123             print('online percep time '+str(time.time() - t_in)+'\n')
124         else:
125             for k in range(Max_iter):
126                 grad_k,N_miss = percept_batch_tangent(theta_k,samples,labels)
```

```python
127                    if N_miss == 0:
128                        break
129                    theta_k += rate*grad_k
130                    theta_store.append(np.copy(theta_k))
131                    N_store.append(N_miss)
132                print('batch percep time '+str(time.time() - t_in)+'\n')
133        if convplot==1:
134            plt.plot(np.arange(len(N_store)),N_store,'*--')
135            plt.xlabel('Iters')
136            plt.ylabel('Number of missclassified samples')
137            plt.show()
138
139        return(theta_k,theta_store)
140
141  #%% function for plotting data
142  def plotdata_single(samples,labels,theta_star,*argv):
143        """
144        Input:
145            sample: np array in the format [[x1],[x2],....,[xn]] xi as col vectors
146            labels: np array [y1,y2,....,yn]
147            theta_star: decision boundary params
148            N: frequency of plots for decision boundary
149        """
150
151        plt.figure()
152        #plot the training set
153        for i in range(np.shape(samples)[1]):
154            if labels[i] == 1:
155                plt.scatter(samples[0,i],samples[1,i],c='red',marker='*')
156            else:
157                plt.scatter(samples[0,i],samples[1,i],c='blue',marker='^')
158
159        # plot the decision boundary
160        x_min = min(samples[0,:])
161        x_max = max(samples[0,:])
162        y_min = min(samples[1,:])
163        y_max = max(samples[1,:])
164        x_line = np.linspace(x_min,x_max,10)
165        y_line = -(theta_star[0]*x_line +theta_star[2])/theta_star[1]
166        plt.plot(x_line,y_line,'k-')
167        plt.xlim(x_min,x_max)
168        plt.ylim(y_min,y_max)
169        if argv!=None:
170            plt.title('C='+str(argv[0]))
171        plt.show()
172        return()
173
174  def plotdata(samples,labels,theta_store,N,*argv):
175        """
176        Input:
177            sample: np array in the format [[x1],[x2],....,[xn]] xi as col vectors
178            labels: np array [y1,y2,....,yn]
179            theta_store: decision boundary params
180            N: frequency of plots for decision boundary
181        """
182
183        plt.figure()
184        #plot the training set
185        for i in range(np.shape(samples)[1]):
```

```python
186             if  labels[i] == 1:
187                 plt.scatter(samples[0,i],samples[1,i],c='red',marker='*')
188             else:
189                 plt.scatter(samples[0,i],samples[1,i],c='blue',marker='^')
190
191         # plot the decision boundary
192         x_min = min(samples[0,:])
193         x_max = max(samples[0,:])
194         y_min = min(samples[1,:])
195         y_max = max(samples[1,:])
196         x_line = np.linspace(x_min,x_max,10)
197         count=0
198         while count<len(theta_store):
199             theta_k = theta_store[count]
200             y_line = -(theta_k[0]*x_line +theta_k[2])/theta_k[1]
201             plt.plot(x_line,y_line,'k--',alpha=0.5)
202             count+=N
203         theta_k = theta_store[-1]
204         y_line = -(theta_k[0]*x_line +theta_k[2])/theta_k[1]
205         plt.plot(x_line,y_line,'k-')
206         plt.xlim(x_min,x_max)
207         plt.ylim(y_min,y_max)
208         if argv!=None:
209             plt.title('M='+str(argv[0]))
210         plt.show()
211         return()
212
213 #%% Logistic regression function
214 # tangent for logistic function
215 def logistic_tangent(theta,samples,labels):
216     """
217     Input:
218         theta: decision boundary
219         sample: np array in the format [[x1],[x2],....,[xn]] xi as col vectors
220         labels: np array [y1,y2,....,yn]
221     Out: J: Jacobian
222     """
223     thetaTx  = theta.T@samples
224     h_theta_x = 1/(1+np.exp(-thetaTx))
225     temp = (h_theta_x - labels)*samples
226     J  = np.sum(temp,axis=1)
227     return(J)
228 # main logistic function
229 def logistic(samples,labels,rate,Max_iter):
230     """
231     Input:
232         samples: np array in the format [[x1],[x2],....,[xn]] xi as col vectors
233         labels: np array [y1,y2,....,yn]
234         rate: learninig rate or the step length
235         Max_iter: for the gradient descent
236     """
237     t_in = time.time()
238     xdim,Nsamples= samples.shape
239     theta_k = np.zeros(xdim) # initial guess
240     theta_store = []
241     theta_store.append(np.copy(theta_k))
242     for k in range(Max_iter):
243         theta_k -= rate*logistic_tangent(theta_k,samples,labels)
244         theta_store.append(np.copy(theta_k))
```

```
245         print('SVM logistic time '+str(time.time() - t_in)+'\n')
246         return(theta_k, theta_store)
```