

Student Name: Musale Krushna Pavan
 Roll Number: 20111268
 Date: November 27, 2020

Given logistic regression model with where y_n in $\{0, 1\}$ with no regularizer. Then the loss function is

$$\mathcal{L}(\omega) = - \sum_{n=1}^N (y_n \omega^T \mathbf{x}_n) - \log(1 + \exp(\omega^T \mathbf{x}_n)) \quad (1)$$

second-order optimization based update

$$\omega^{(t+1)} = \omega^{(t)} - (\mathbf{H}^{(t)})^{-1} \mathbf{g}^{(t)} \quad (2)$$

Where \mathbf{H} is hessian and \mathbf{g} denotes the gradient

Need to show that the second-order optimization is equal to

$$\begin{aligned} \omega^{(t+1)} &= \arg \min_{\omega} \sum_{n=1}^N \gamma_n^{(t)} (\hat{y}_n^{(t)} - \omega^T \mathbf{x}_n)^2 \\ \Rightarrow \sum_{n=1}^N 2\gamma_n^{(t)} (\hat{y}_n^{(t)} - \omega^T \mathbf{x}_n) \mathbf{x}_n &= 0 \\ \sum_{n=1}^N \gamma_n^{(t)} (\hat{y}_n^{(t)} \mathbf{x}_n) &= \sum_{n=1}^N \gamma_n^{(t)} \mathbf{x}_n \mathbf{x}_n^T \omega \\ \omega^{(t+1)} &= \left(\sum_{n=1}^N \gamma_n^{(t)} \mathbf{x}_n \mathbf{x}_n^T \right)^{-1} \sum_{n=1}^N \gamma_n^{(t)} \hat{y}_n^{(t)} \mathbf{x}_n \end{aligned}$$

Now considering the the equation (1)

$$\begin{aligned} \mathbf{g} &= \frac{\partial \mathcal{L}(\omega)}{\partial \omega} = - \sum_{n=1}^N (y_n \mathbf{x}_n) - \frac{\exp(\omega^T \mathbf{x}_n)}{1 + \exp(\omega^T \mathbf{x}_n)} \mathbf{x}_n \\ &= - \sum_{n=1}^N (y_n \mathbf{x}_n) - \sigma(\omega^T \mathbf{x}_n) \mathbf{x}_n \quad \text{where } \sigma(x) = \frac{1}{1 + \exp(-x)} \end{aligned}$$

$$\begin{aligned} \mathbf{H} &= \frac{\partial^2 \mathcal{L}(\omega)}{\partial^2 \omega} = \frac{\partial \mathbf{g}(\omega)}{\partial \omega} = \frac{\partial}{\partial \omega} \left[- \sum_{n=1}^N (y_n \mathbf{x}_n) - \sigma(\omega^T \mathbf{x}_n) \mathbf{x}_n \right] \\ &= \sum_{n=1}^N \sigma(\omega^T \mathbf{x}_n) \sigma(1 - \omega^T \mathbf{x}_n) \mathbf{x}_n \mathbf{x}_n^T \end{aligned}$$

Now substituting \mathbf{g} and \mathbf{H} in equation (2)

$$\begin{aligned}
\boldsymbol{\omega}^{(t+1)} &= \boldsymbol{\omega}^{(t)} + \left(\sum_{n=1}^N \sigma(\boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \sigma(1 - \boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \mathbf{x}_n \mathbf{x}_n^T \right)^{-1} \sum_{n=1}^N (y_n \mathbf{x}_n - \sigma(\boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \mathbf{x}_n) \\
&= \left(\sum_{n=1}^N \sigma(\boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \sigma(1 - \boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \mathbf{x}_n \mathbf{x}_n^T \right)^{-1} \\
&\quad \left[\left(\sum_{n=1}^N \sigma(\boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \sigma(1 - \boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \mathbf{x}_n \mathbf{x}_n^T \right) \boldsymbol{\omega}^{(t)} + \sum_{n=1}^N (y_n \mathbf{x}_n - \sigma(\boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \mathbf{x}_n) \right] \\
&= \left(\sum_{n=1}^N \sigma(\boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \sigma(1 - \boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \mathbf{x}_n \mathbf{x}_n^T \right)^{-1} \\
&\quad \left[\sum_{n=1}^N [\sigma(\boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \sigma(1 - \boldsymbol{\omega}^{(t)T} \mathbf{x}_n)] [\mathbf{x}_n^T \boldsymbol{\omega}^{(t)} + \frac{y_n - \sigma(\boldsymbol{\omega}^{(t)T} \mathbf{x}_n)}{\sigma(\boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \sigma(1 - \boldsymbol{\omega}^{(t)T} \mathbf{x}_n)}] \mathbf{x}_n \right]
\end{aligned}$$

Let

$$\begin{aligned}
\gamma_n^{(t)} &= \sigma(\boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \sigma(1 - \boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \\
\hat{y}_n^{(t)} &= \mathbf{x}_n^T \boldsymbol{\omega}^{(t)} + \frac{y_n - \sigma(\boldsymbol{\omega}^{(t)T} \mathbf{x}_n)}{\sigma(\boldsymbol{\omega}^{(t)T} \mathbf{x}_n) \sigma(1 - \boldsymbol{\omega}^{(t)T} \mathbf{x}_n)}
\end{aligned}$$

Then the

$$\boldsymbol{\omega}^{(t+1)} = \left(\sum_{n=1}^N \gamma_n^{(t)} \mathbf{x}_n \mathbf{x}_n^T \right)^{-1} \sum_{n=1}^N \gamma_n^{(t)} \hat{y}_n^{(t)} \mathbf{x}_n$$

from the equation of γ_n we can make observe that the points near the boundary have more weight than the points far away from the boundary.

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: November 27, 2020

My solution to problem 2

Suppose that there is a transformation function ϕ that convert x to its huber space.
 And let there is kernal function $k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$

We know that the weight vector in the perceptron is in the form $\omega = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$

Now kernalized algorithms for SGD

1. Initialize $w = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$ and $\alpha_i = 0 \quad \forall i$ and set $\eta_t = 1$ – Initialization
2. Pick a point (\mathbf{x}_i, y_i) randomly
- 3.

$$\hat{y} = \text{sign}\left(\left(\sum_{n=1}^N \alpha_n y_n \phi(\mathbf{x}_n)\right)^T \phi(\mathbf{x}_i)\right)$$

$$\hat{y} = \text{sign}\left(\sum_{n=1}^N \alpha_n y_n \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_i)\right)$$

$$\hat{y} = \text{sign}\left(\sum_{n=1}^N \alpha_n y_n k(\mathbf{x}_n, \mathbf{x}_i)\right) \quad \text{– mistake finding}$$

4. If $\hat{y} \neq y_i$ then $\alpha_i + = 1$ – update if mistake
5. If not converged go to step 2.

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: November 27, 2020

My solution to problem 3

As sometimes its costs more to wrongly classify as negative than positive we can incorporate this information in the SVM equation as follows:

$$\mathcal{L}(\omega, b, \xi) = \min_{\omega, b, \xi} \frac{\|\omega\|^2}{2} + \sum_{n=1}^N C_{y_n} \xi_n$$

$$\text{such that } y_n(\omega^T \mathbf{x}_n + b) \geq 1 - \xi_n \text{ and } \xi_n \geq 0, \forall n$$

Using legrangian method to solve the constraint optimization problem. whose dual can be formed as following

$$\mathcal{L}(\omega, b, \xi, \alpha, \beta) = \max_{\alpha \geq 0, \beta \geq 0} \min_{\omega, b, \xi} \frac{\|\omega\|^2}{2} + \sum_{n=1}^N C_{y_n} \xi_n + \sum_{n=1}^N \alpha_n (1 - \xi_n - y_n(\omega^T \mathbf{x}_n + b)) - \sum_{n=1}^N \beta_n \xi_n \quad (3)$$

$$= \max_{\alpha \geq 0, \beta \geq 0} \min_{\omega, b, \xi} \frac{\omega^T \omega}{2} + \sum_{n=1}^N C_{y_n} \xi_n + \sum_{n=1}^N \alpha_n (1 - \xi_n - y_n(\omega^T \mathbf{x}_n + b)) - \sum_{n=1}^N \beta_n \xi_n \quad (4)$$

Now for solving for ω, b, ξ we partially derivate $\mathcal{L}(\omega, b, \xi, \alpha, \beta)$ with each of the ω, b, ξ and equating to 0

First derivating \mathcal{L} with ω and equating to 0

$$\frac{\partial \mathcal{L}}{\partial \omega} = 0 \quad (5)$$

$$\omega - \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n = 0 \quad (6)$$

$$w = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \quad (7)$$

Now derivating \mathcal{L} with b and equating to 0

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \quad (8)$$

$$\sum_{n=1}^N \alpha_n y_n = 0 \quad (9)$$

Now derivating \mathcal{L} with ξ_n and equating to 0, which is similar $\forall n$

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = 0 \quad (10)$$

$$C_{y_n} - \alpha_n - \beta_n = 0 \quad (11)$$

$$C_{y_n} \geq \alpha_n \quad \text{using } \beta_n \geq 0 \quad (12)$$

By substituting $w = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$ in the Lagrangian, we get dual as

$$\begin{aligned}
\mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\beta}) &= \max_{\boldsymbol{\alpha} \geq 0, \boldsymbol{\beta} \geq 0} \frac{\boldsymbol{\omega}^T \boldsymbol{\omega}}{2} + \sum_{n=1}^N \mathcal{C}_{y_n} \xi_n + \sum_{n=1}^N \alpha_n (1 - \xi_n - y_n (\boldsymbol{\omega}^T \mathbf{x}_n + b)) - \sum_{n=1}^N \beta_n \xi_n \\
&= \max_{\boldsymbol{\alpha} \geq 0, \boldsymbol{\beta} \geq 0} \frac{1}{2} \left(\sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \right)^T \left(\sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \right) + \sum_{n=1}^N \mathcal{L}_{y_n} \xi_n + \sum_{n=1}^N \alpha_n - \sum_{n=1}^N \alpha_n \xi_n \\
&\quad - \sum_{n=1}^N \alpha_n y_n \left(\sum_{m=1}^N \alpha_m y_m \mathbf{x}_m \right)^T \mathbf{x}_n + \sum_{n=1}^N \alpha_n y_n b - \sum_{n=1}^N \beta_n \xi_n
\end{aligned}$$

using $\mathcal{C}_{y_n} - \alpha_n - \beta_n = 0$ and $\sum_{n=1}^N \beta_n \geq 0$ we get

$$\mathcal{L}(\boldsymbol{\alpha}) = \max_{\alpha_n \leq \mathcal{C}_{y_n}} \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n,m=1}^N \alpha_n \alpha_m y_n y_m \mathbf{x}_m^T \mathbf{x}_n \quad (13)$$

When we compare the above equation with the standard SVM equation there is only difference in the condition of $\alpha_n \leq \mathcal{C}_{y_n}$. This specifies that the 2 supporting vectors for the need not be same distance form the decision boundary can used to give more weight on negatively classify on appropriate setting of \mathcal{C}_{y_n} .

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: November 27, 2020

Now making the above batch k-means to online.

Consider the following initialization:

let the data points be $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_n$

There be K clusters where $K \in 1, 2, \dots, K$

Means of the cluster k is randomly initialized be μ_k .

Number of nodes in cluster k be $\eta_k = 0 \forall k \in \{1, 2, 3, \dots, k\}$

Consider a point random point \mathbf{x}_n

(1) Assigning \mathbf{x}_n "greedily" to the "best" cluster by solving the following function.

$$\mathbf{z}_n = \arg \min_{k \in 1, 2, \dots, K} \|\mathbf{x}_n - \mu_k\|^2$$

\mathbf{z}_n is the one hot encoding whose values is 0 or 1. 1 is at the index where the point is assigned.
 The loss is given by

$$\begin{aligned} \ell(\mu, \mathbf{x}_n, \mathbf{z}_n) &= \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2 \\ &= \|\mathbf{x}_n - \mu_{k_n}\|^2 \quad \text{when assigned to } k^{\text{th}} \text{ cluster.} \end{aligned}$$

Now the gradient \mathbf{g} of the loss function is given by:

$$\begin{aligned} \mathbf{g} &= 2(\mu_{k_n} - \mathbf{x}_n) \\ \mathbf{g}^{(t)} &= 2(\mu_{k_n}^{(t)} - \mathbf{x}_n) \end{aligned}$$

Update rule according to stochastic gradient decent is :

$$\begin{aligned} \mu_{k_n}^{(t+1)} &= \mu_{k_n}^{(t)} - n_t \mathbf{g}^{(t)} \quad n_t \text{ is learning rate.} \\ &= \mu_{k_n}^{(t)} - n_t 2(\mu_{k_n}^{(t)} - \mathbf{x}_n) \\ &= (1 - 2n_t) \mu_{k_n}^{(t)} + (2n_t) \mathbf{x}_n \quad \text{--- (1)} \\ \eta_k &= \eta_k + 1 \end{aligned}$$

Suggesting a good choice of the step size.

In general when we assign a new point x_n to a cluster the mean μ_{k_n} and if the cluster already has η_k points in it then the updated mean will be

$$\mu_{k_n}' = \frac{\mu_{k_n} \eta_k + \mathbf{x}_n}{\eta_k + 1} \quad \text{--- (2)}$$

By comparing eqn (1) and (2) we get the learning rate

$$n_t = \frac{1}{2(\eta_k + 1)}$$

The above learning rate also gives the intuition that when there less points in the cluster it has larger steps to update the mean and similarly when there are more points in the cluster it has very small steps to update the mean, which should be the general case of updating.

Student Name: Musale Krushna Pavan

Roll Number: 20111268

Date: November 27, 2020

let the data points be $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_n$. There be K clusters where $K \in 1, 2, \dots, K$. let ϕ be a mapping function that maps to hilbert space. and kernel $k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$. Then the cluster mean in hilbert space is given by $\phi(\mu_k) = \frac{1}{|C_k|} \sum_{n:z_n=k} \phi(\mathbf{x}_n)$

Algorithm:

(1) Let randomly assign each point to one of the cluster

$$\mathbf{z}_n = \text{random}(k) \quad \forall n, k \in 1, 2, \dots, K$$

(2) Update the each point to its nearest cluster in the hilbert space

Consider:

$$\begin{aligned} \|\phi(\mathbf{x}_n) - \phi(\mu_k)\|^2 &= (\phi(\mathbf{x}_n) - \phi(\mu_k))^T (\phi(\mathbf{x}_n) - \phi(\mu_k)) \\ &= (\phi(\mathbf{x}_n)^T - \phi(\mu_k)^T) (\phi(\mathbf{x}_n) - \phi(\mu_k)) \\ &= \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_n) - 2\phi(\mathbf{x}_n)^T \phi(\mu_k) - \phi(\mu_k)^T \phi(\mu_k) \\ &= \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_n) - 2\phi(\mathbf{x}_n)^T \left(\frac{1}{|C_k|} \sum_{m:z_m=k} \phi(\mathbf{x}_m) \right) \\ &\quad - \left(\frac{1}{|C_k|} \sum_{m:z_m=k} \phi(\mathbf{x}_m) \right)^T \left(\frac{1}{|C_k|} \sum_{m:z_m=k} \phi(\mathbf{x}_m) \right) \\ &= \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_n) - 2 \frac{1}{|C_k|} \sum_{m:z_m=k} \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) - \frac{1}{|C_k|^2} \sum_{\substack{i,j:z_i=k \\ z_j=k}} \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \\ &= k(\mathbf{x}_n, \mathbf{x}_n) - \frac{2}{|C_k|} \sum_{m:z_m=k} k(\mathbf{x}_n, \mathbf{x}_m) + \frac{1}{|C_k|^2} \sum_{\substack{i,j:z_i=k \\ z_j=k}} k(\mathbf{x}_i, \mathbf{x}_j) \end{aligned}$$

Now updating each point to its nearest cluster as

$$\mathbf{z}_n = \arg \min_{k \in 1, 2, \dots, K} k(\mathbf{x}_n, \mathbf{x}_n) - \frac{2}{|C_k|} \sum_{m:z_m=k} k(\mathbf{x}_n, \mathbf{x}_m) + \frac{1}{|C_k|^2} \sum_{\substack{i,j:z_i=k \\ z_j=k}} k(\mathbf{x}_i, \mathbf{x}_j)$$

We do not compute as mean as its is infinite dimensional(kernel)

(3) If not converge i.e the cluster assignments does not change then go to step 2.

Instead of storing the cluster mean as in the case of normal k-means, in kernel-K means we store which point belongs to which cluster.

Finally, assuming each input to be D-dimensional in the original feature space, and N to be the number of inputs. The cost of input to clusters mean distance calculation

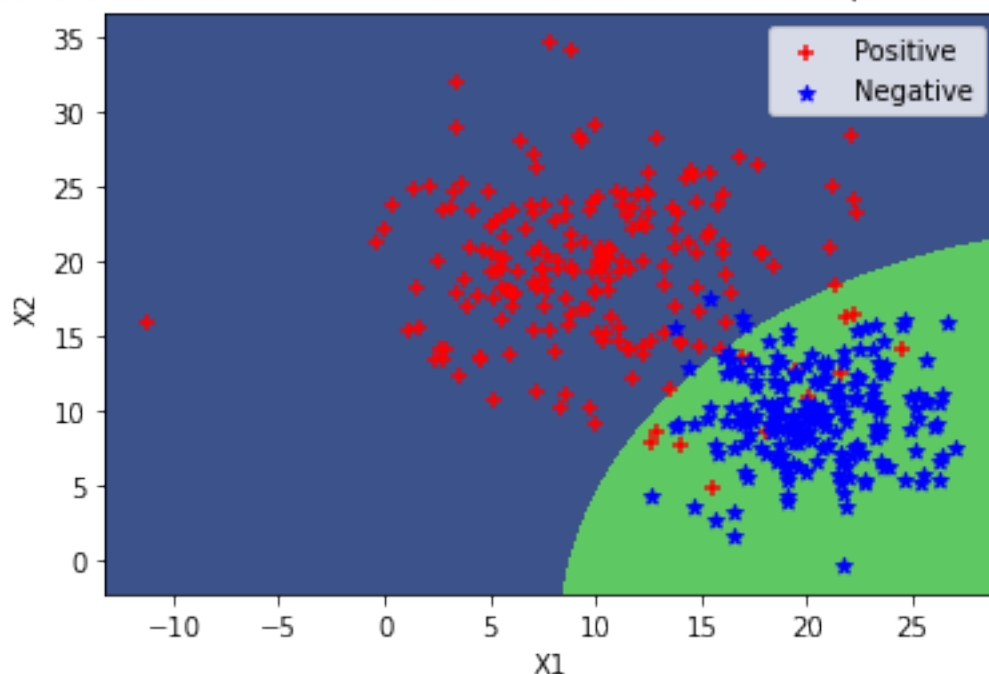
(1) standard K-means just $\|\mathbf{x}_n - \mu_k\|^2 \quad \forall k$ is $O(KD)$

(2) Kernel k-means it is $O(D + DN + D(N - K + 1)(N - K)) = O(DN^2 + DNK + DK^2) = O(DN^2)$ when $N \gg K$ from the update equation we got.

1 Part1 - binclass.txt

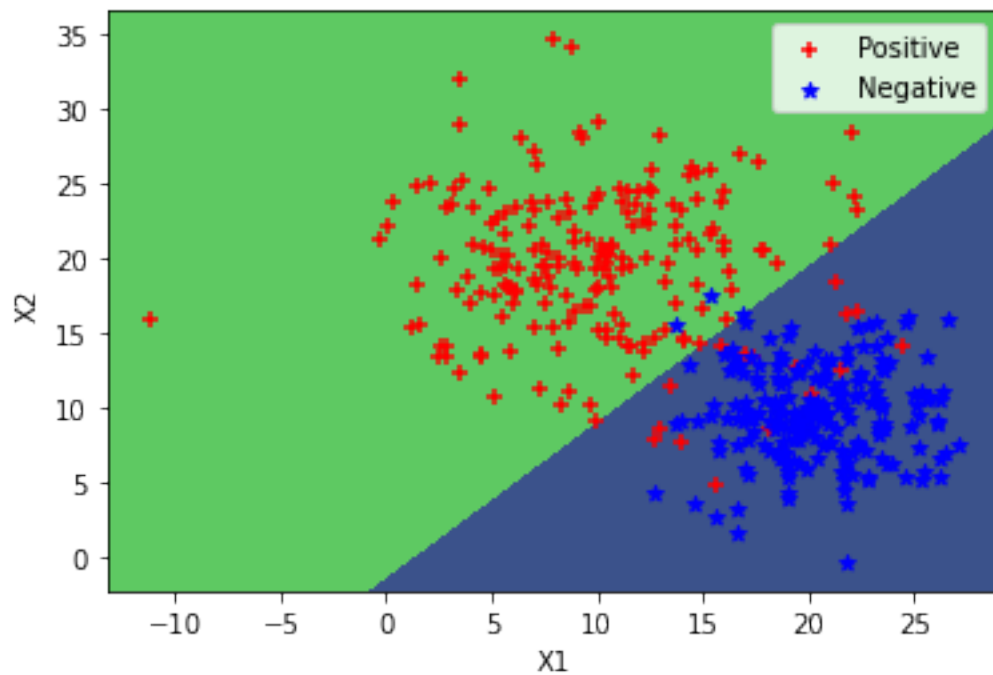
The dataset that is used is in binclass.txt with gaussian class-conditional distributions with different covariance matrices for positive and negative class examples the decision boundary that we get is as follows

Generative Classification with Gaussian Class-Conditionals (different covariances)

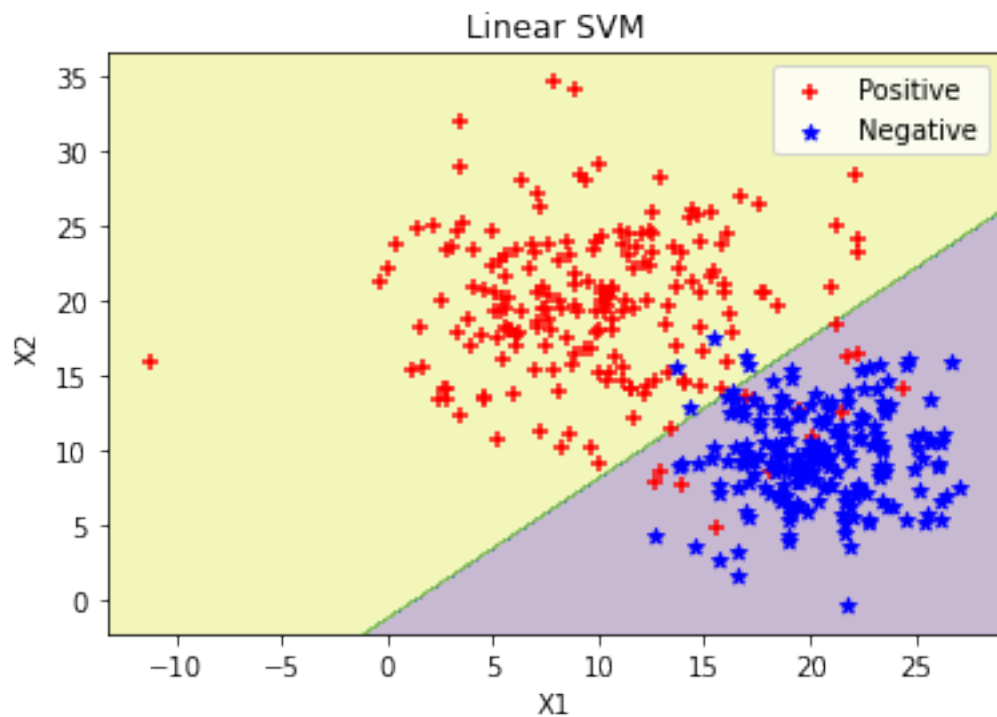


with gaussian class-conditional distributions with same covariance matrices for positive and negative class examples the decision boundary that we get is as follows

Generative Classification with Gaussian Class-Conditionals ((same covariances))



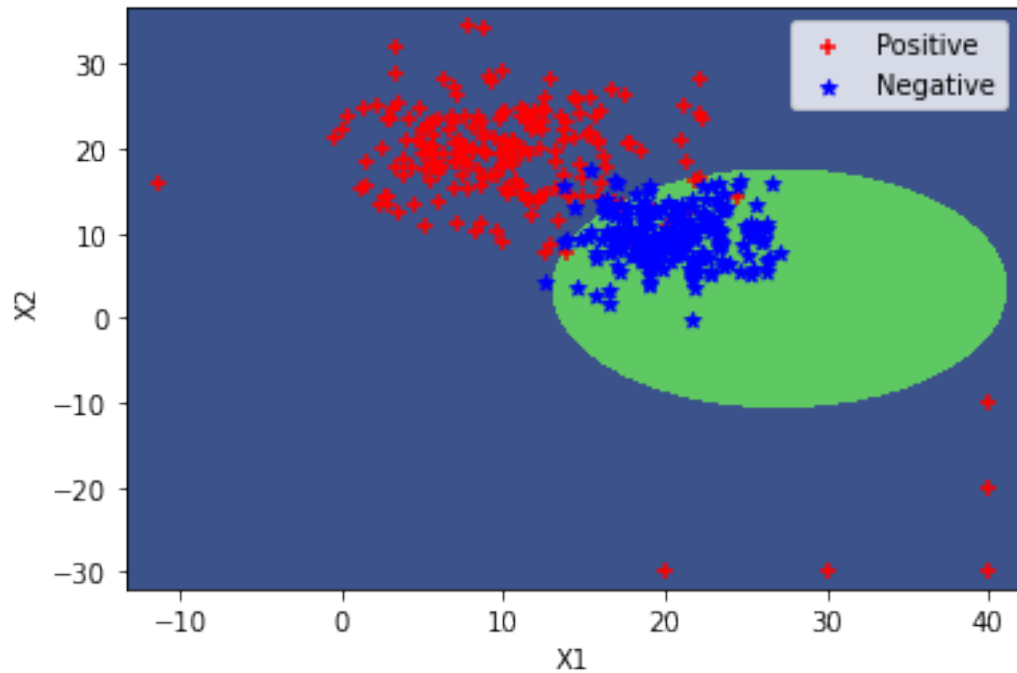
with SVM (from scikit-learn)



2 Part2

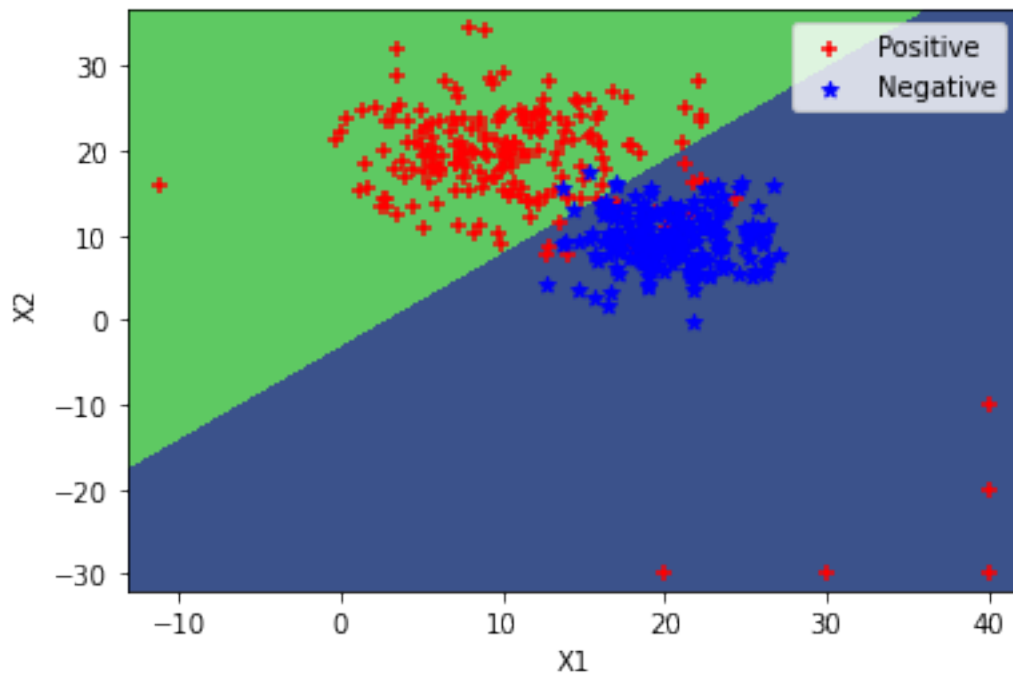
The dataset that is used is in `binclassv2.txt` with gaussian class-conditional distributions with different covariance matrices for positive and negative class examples the decision boundary that we get is as follows

Generative Classification with Gaussian Class-Conditionals (different covariances)

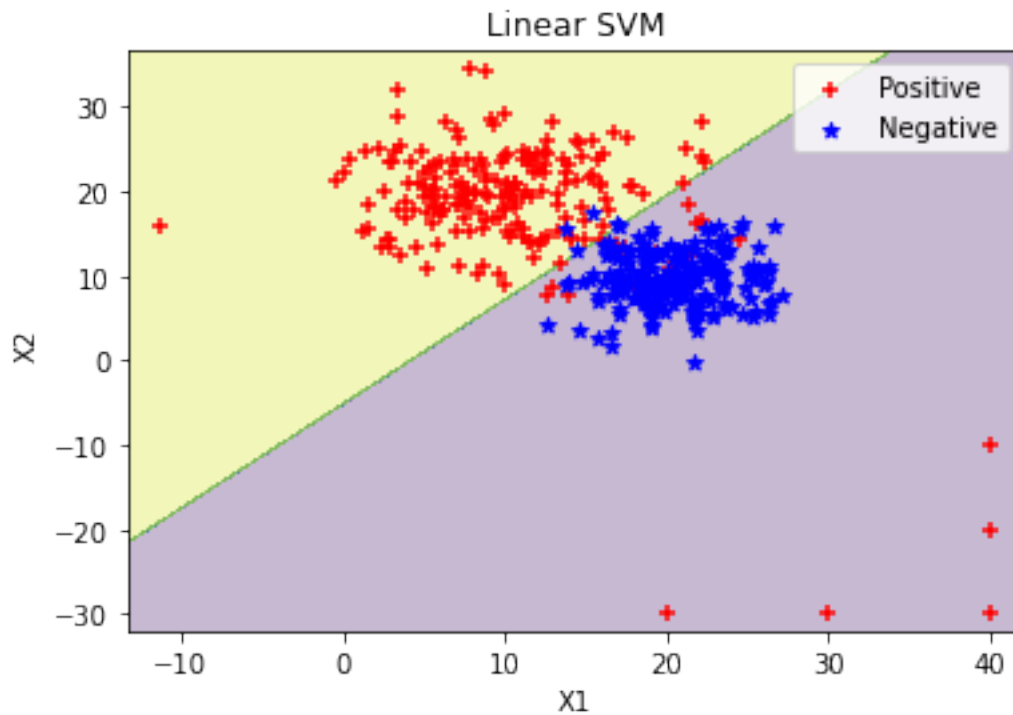


with gaussian class-conditional distributions with same covariance matrices for positive and negative class examples the decision boundary that we get is as follows

Generative Classification with Gaussian Class-Conditionals ((same covariances))



with SVM (from scikit-learn)



Observation:

First dataset: All the 3 decision boundaries are equally good

Second dataset: We see that there are some outliers and our generative classifier with different covariance was trying to overfit the data, we also have to keep in mind that we are using MLE estimation which is prone to overfit. But the other 2 models are good.

In general SVM with kernel is good as it can learn non linear boundaries with maximum margin.