

▼ udaylunawat@gmail.com

```
1 import matplotlib.pyplot as plt
2 from sklearn.linear_model import SGDClassifier
3 from sklearn.linear_model import LogisticRegression
4 import pandas as pd
5 import numpy as np
6 from sklearn.preprocessing import StandardScaler, Normalizer
7 import matplotlib.pyplot as plt
8 from sklearn.svm import SVC
9 import warnings
10 warnings.filterwarnings("ignore")
```

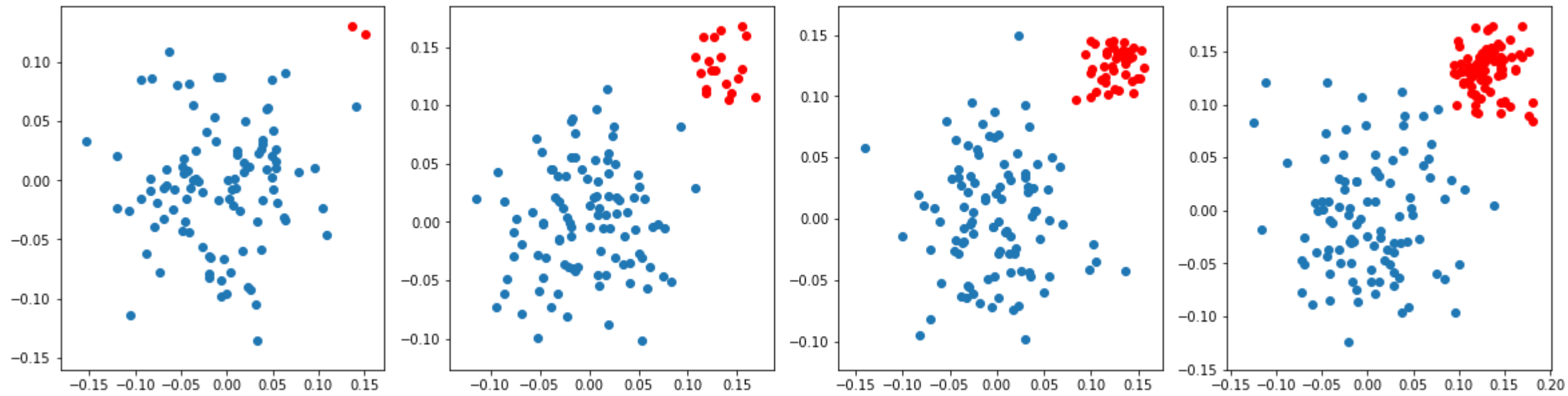
```
1 def draw_line(coef, intercept, mi, ma):
2     # for the separating hyper plane ax+by+c=0, the weights are [a, b] and the intercept is c
3     # to draw the hyper plane we are creating two points
4     # 1. ((b*min-c)/a, min) i.e ax+by+c=0 ==> ax = (-by-c) ==> x = (-by-c)/a here in place of y we are keeping the minimum value of y
5     # 2. ((b*max-c)/a, max) i.e ax+by+c=0 ==> ax = (-by-c) ==> x = (-by-c)/a here in place of y we are keeping the maximum value of y
6     points = np.array([((-coef[1]*mi - intercept)/coef[0]), mi], [((-coef[1]*ma - intercept)/coef[0]), ma])
7     plt.plot(points[:,0], points[:,1])
```

1

▼ What if Data is imabalanced

1. As a part of this task you will observe how linear models work in case of data imbalanced
2. observe how hyper plane is changes according to change in your learning rate.
3. below we have created 4 random datasets which are linearly separable and having class imbalance
4. in the first dataset the ration between positive and negative is 100 : 2, in the 2nd data its 100:20, in the 3rd data its 100:40

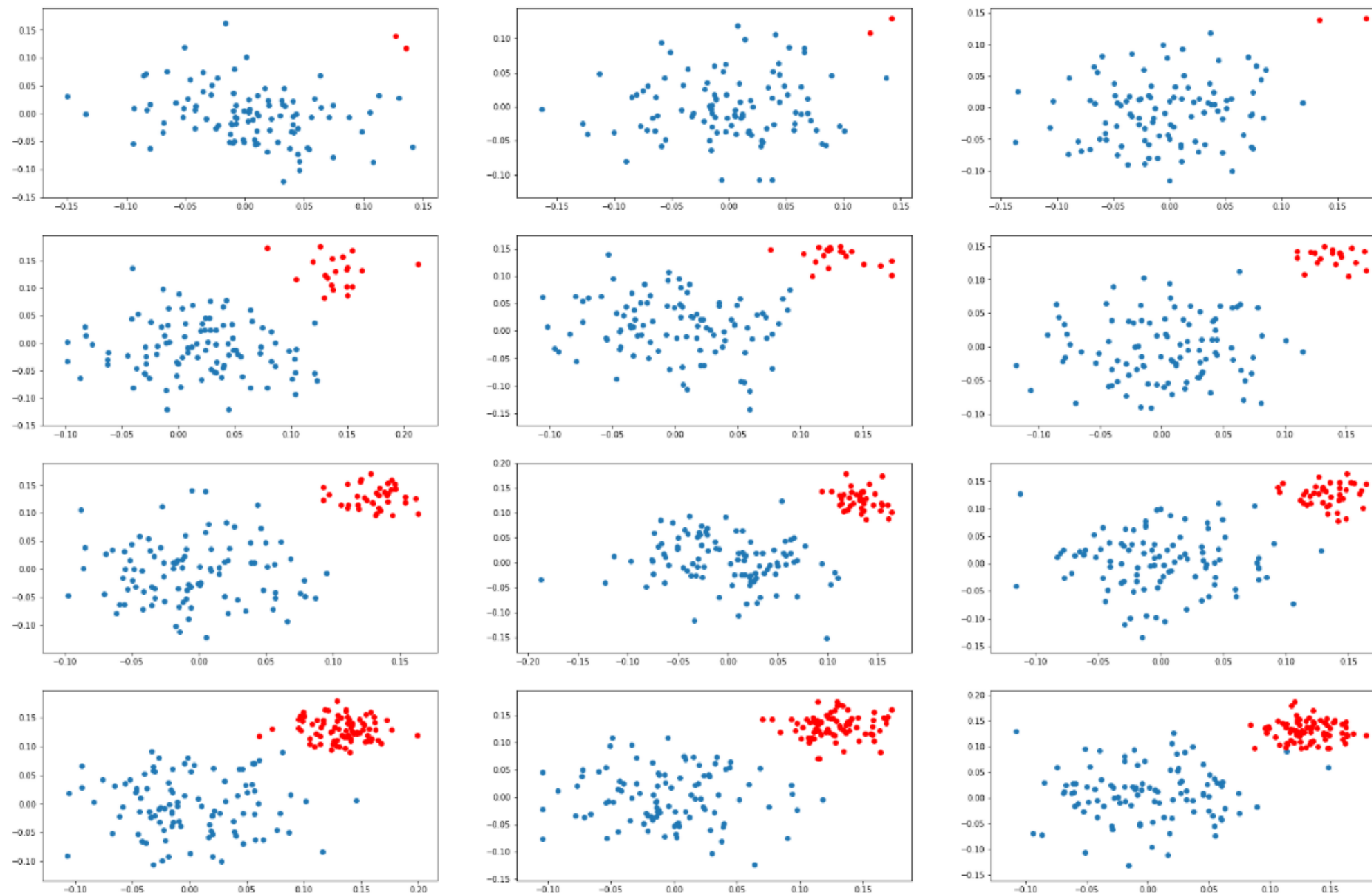
```
1 # here we are creating 2d imbalanced data points
2 ratios = [(100,2), (100, 20), (100, 40), (100, 80)]
3 plt.figure(figsize=(20,5))
4 for j,i in enumerate(ratios):
5     plt.subplot(1, 4, j+1)
6     X_p=np.random.normal(0,0.05,size=(i[0],2))
7     X_n=np.random.normal(0.13,0.02,size=(i[1],2))
8     y_p=np.array([1]*i[0]).reshape(-1,1)
9     y_n=np.array([0]*i[1]).reshape(-1,1)
10    X=np.vstack((X_p,X_n))
11    y=np.vstack((y_p,y_n))
12    plt.scatter(X_p[:,0],X_p[:,1])
13    plt.scatter(X_n[:,0],X_n[:,1],color='red')
14 plt.show()
```



your task is to apply SVM ([sklearn.svm.SVC](#)) and LR ([sklearn.linear_model.LogisticRegression](#)) with different regularization strength [0.001, 1, 100]

▼ Task 1: Applying SVM

1. you need to create a grid of plots like this



in each of the cell[i][j] you will be drawing the hyper plane that you get after applying [SVM](#) on ith dataset and jth learning

i.e

```
Plane(SVM().fit(D1, C=0.001)) Plane(SVM().fit(D1, C=1)) Plane(SVM().fit(D1, C=100))
Plane(SVM().fit(D2, C=0.001)) Plane(SVM().fit(D2, C=1)) Plane(SVM().fit(D2, C=100))
Plane(SVM().fit(D3, C=0.001)) Plane(SVM().fit(D3, C=1)) Plane(SVM().fit(D3, C=100))
Plane(SVM().fit(D4, C=0.001)) Plane(SVM().fit(D4, C=1)) Plane(SVM().fit(D4, C=100))
```

if you can do, you can represent the support vectors in different colors, which will help us understand the position of hyper plane

Write in your own words, the observations from the above plots, and what do you think about the position of the hy

check the optimization problem here <https://scikit-learn.org/stable/modules/svm.html#mathematical-formulation>

if you can describe your understanding by writing it on a paper and attach the picture, or record a video upload it in assignment.

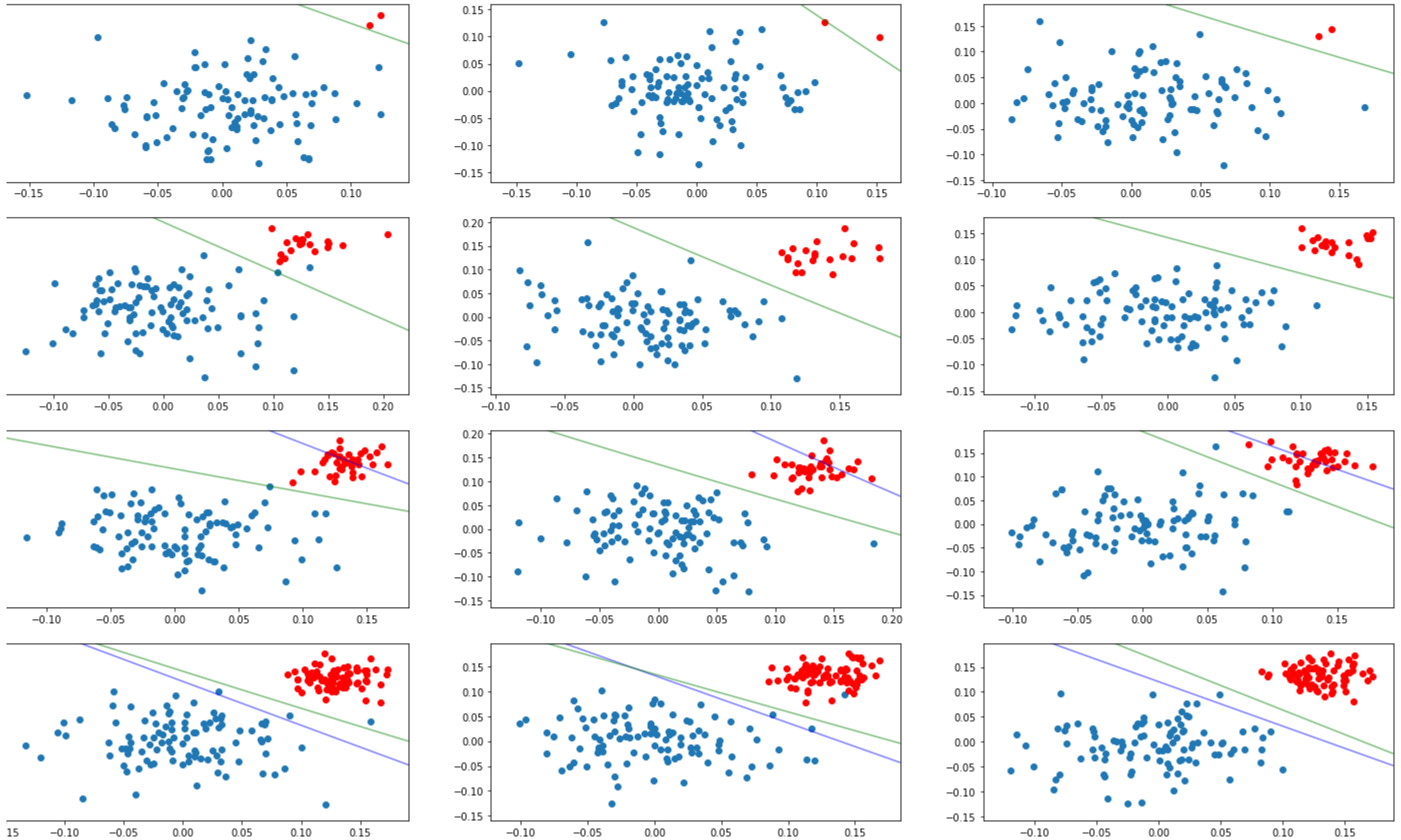
```
1 #Source - https://jakevdp.github.io/PythonDataScienceHandbook/05.07-support-vector-machines.html
2 def plot_svc_decision_function(model, color, ax=None, plot_support=False):
3     """Plot the decision function for a 2D SVC"""
4     if ax is None:
5         ax = plt.gca()
6     xlim = ax.get_xlim()
7     ylim = ax.get_ylim()
8
9     # create grid to evaluate model
10    x = np.linspace(xlim[0], xlim[1], 30)
11    y = np.linspace(ylim[0], ylim[1], 30)
12    Y, X = np.meshgrid(y, x)
13    xy = np.vstack([X.ravel(), Y.ravel()]).T
14    P = model.decision_function(xy).reshape(X.shape)
15
16    # plot decision boundary and margins
17    ax.contour(X, Y, P, colors=color,
18              levels=[0], alpha=0.5,
19              linestyles=['-'])
20
21    ax.set_xlim(xlim)
22    ax.set_ylim(ylim)
```

```
1 #SVC
2 ratios = [(100,2), (100,2), (100,2), (100, 20), (100, 20), (100, 20), (100, 40), (100, 40), (100, 40), (100, 80), (100, 80), (100, 80)]
3 plt.figure(figsize=(25, 15))
4 plt.suptitle('SVM', fontsize =25)
5
6 for j,i in enumerate(ratios):
7     plt.subplot(4, 3, j+1)
8     X_p = np.random.normal(0,0.05,size=(i[0],2))
9     X_n = np.random.normal(0.13,0.02,size=(i[1],2))
10    y_p = np.array([1]*i[0]).reshape(-1,1)
11    y_n = np.array([0]*i[1]).reshape(-1,1)
12    X = np.vstack((X_p,X_n))
13    y = np.vstack((y_p,y_n))
14    plt.scatter(X_p[:,0],X_p[:,1])
15    plt.scatter(X_n[:,0],X_n[:,1],color='red')
16
17
18
19    clf1 = SVC(C=0.001)
20    clf1.fit(X, y)
21    l1 = plot_svc_decision_function(clf1, color = 'red');
22
```

```
23     clf2 = SVC(C=1)
24     clf2.fit(X, y)
25     l2 = plot_svc_decision_function(clf2, color = 'blue');
26
27     clf3 = SVC(C=100)
28     clf3.fit(X, y)
29     l3 = plot_svc_decision_function(clf3, color = 'green');
30
31
32 plt.show()
```



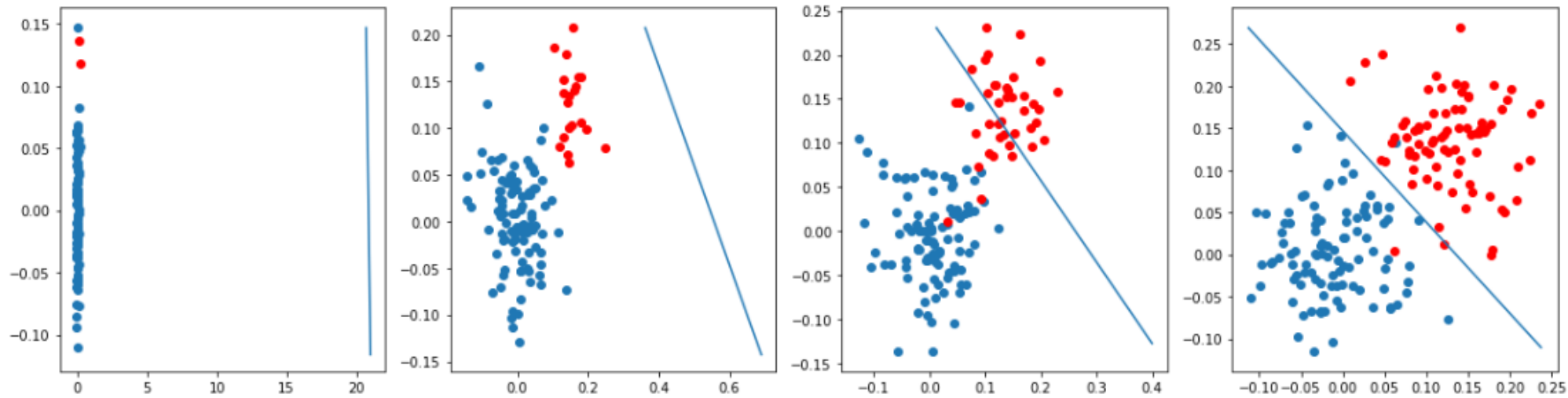
SVM



▼ Task 2: Applying LR

you will do the same thing what you have done in task 1.1, except instead of SVM you apply [logistic regression](#)

these are results we got when we are experimenting with one of the model

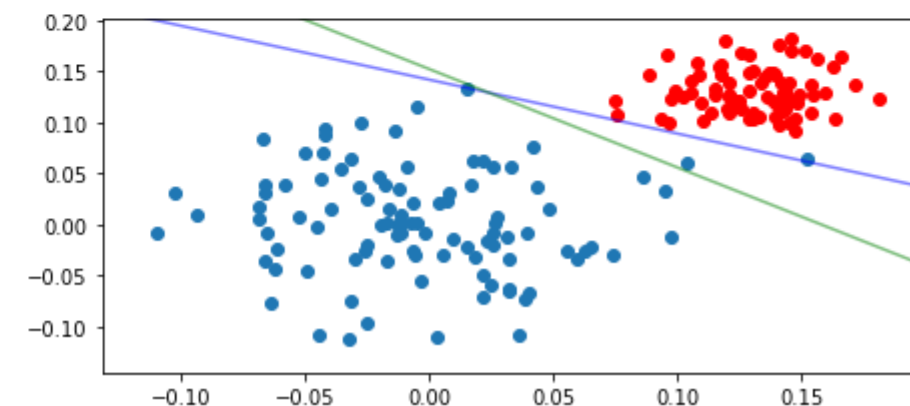
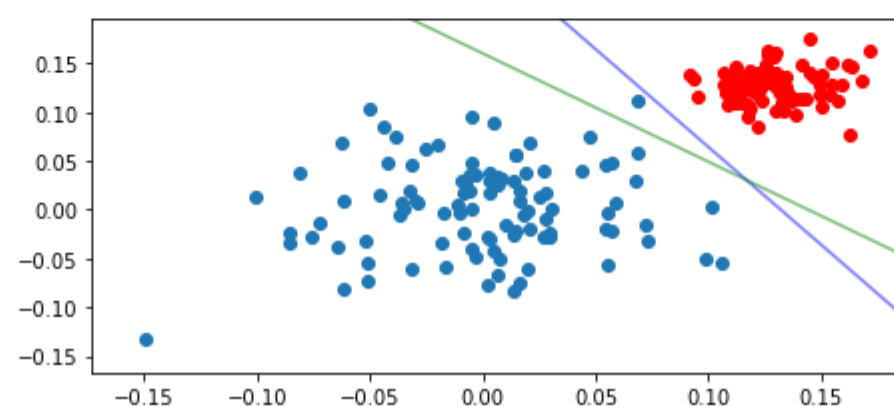
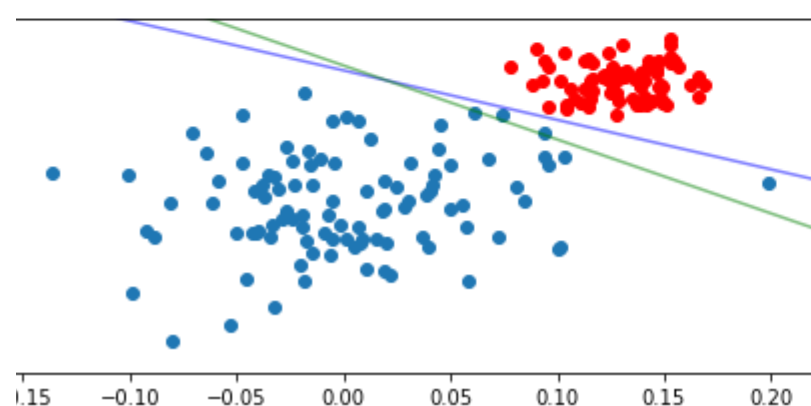
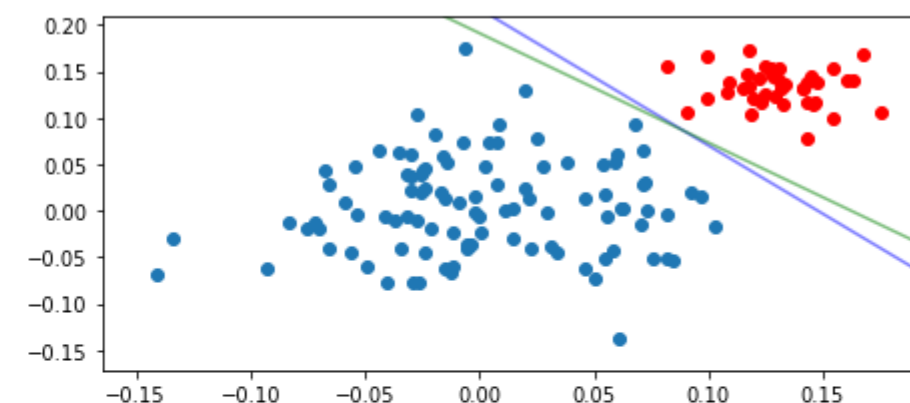
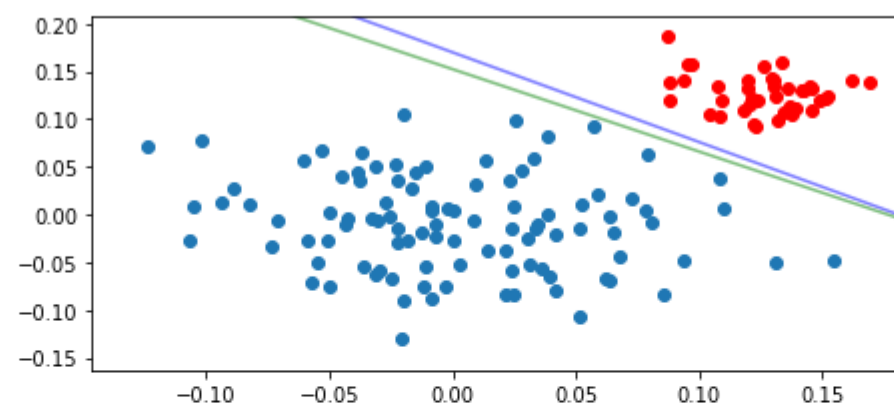
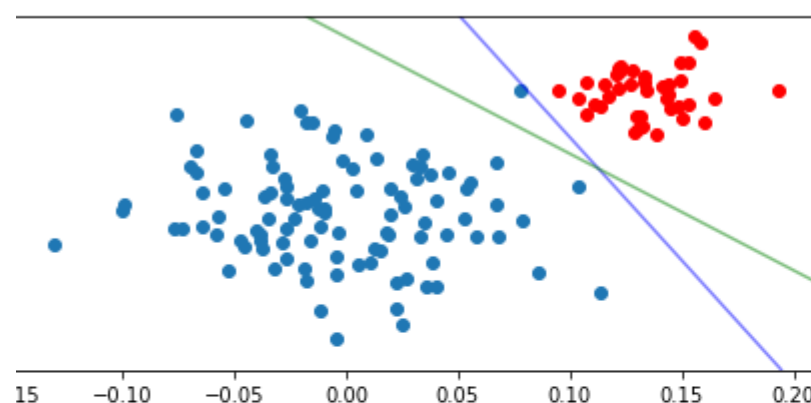
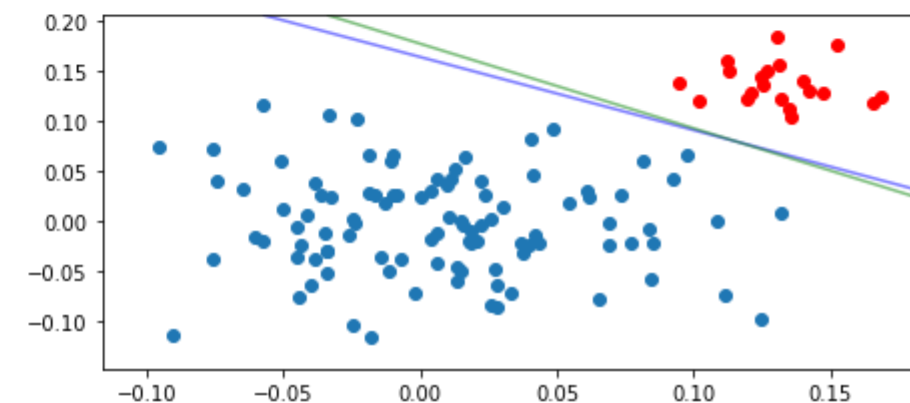
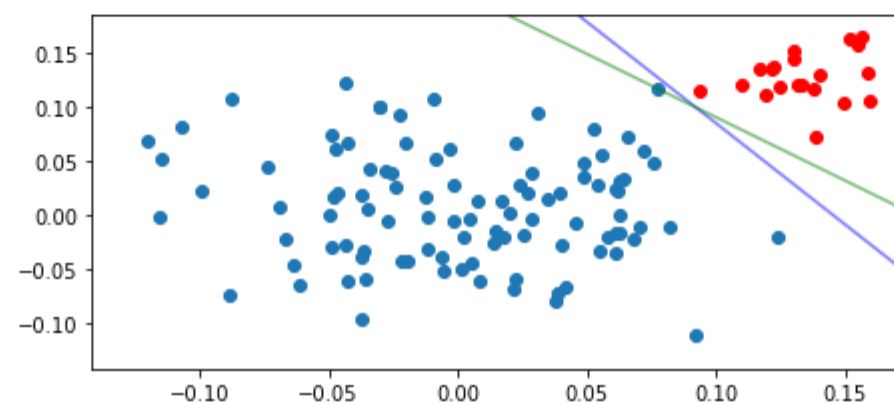
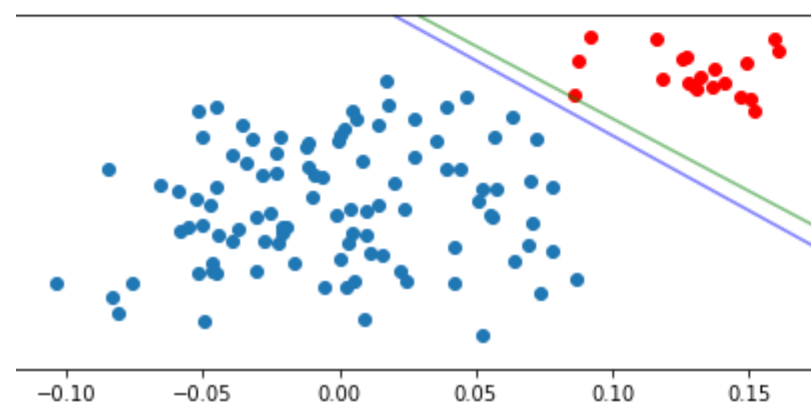
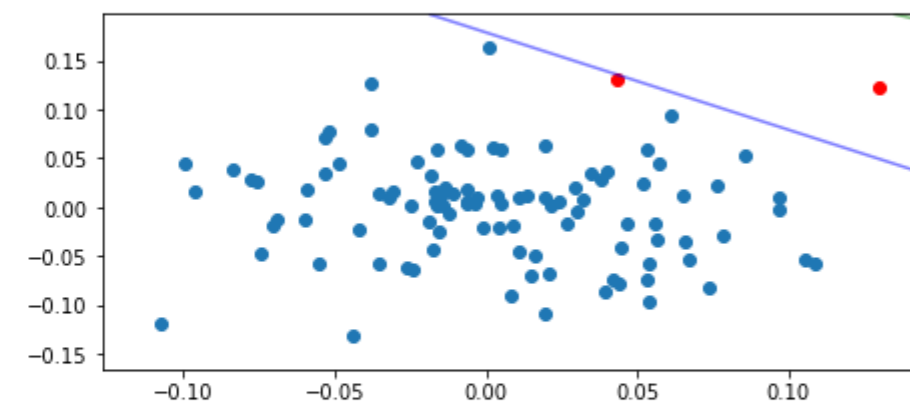
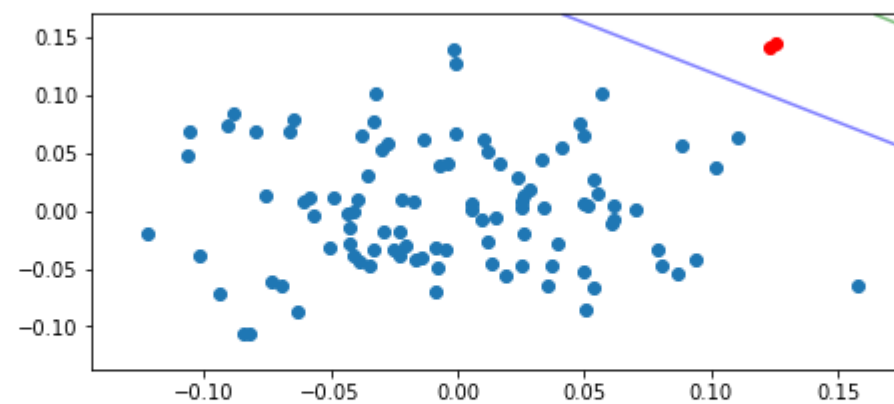
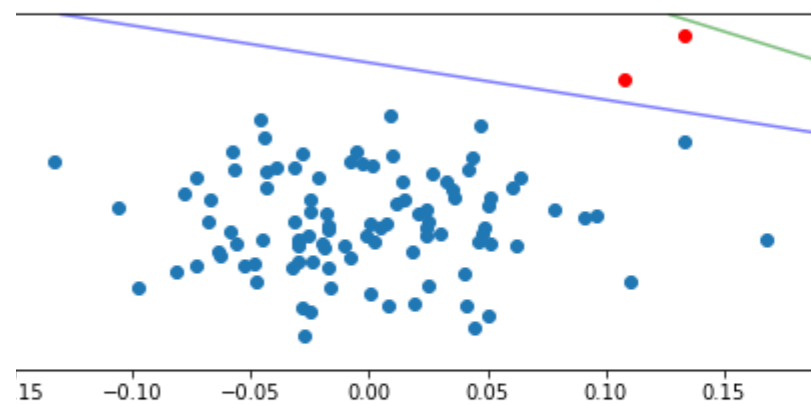


```
1 #Logistic Regression
2
3 ratios = [(100,2), (100,2), (100,2), (100, 20), (100, 20), (100, 20), (100, 40), (100, 40), (100, 40), (100, 80), (100, 80), (100, 80)]
4 plt.figure(figsize=(25,15))
5 plt.suptitle('Logistic Regression', fontsize =25)
6 for j,i in enumerate(ratios):
7     plt.subplot(4, 3, j+1)
8     X_p=np.random.normal(0,0.05,size=(i[0],2))
9     X_n=np.random.normal(0.13,0.02,size=(i[1],2))
10    y_p=np.array([1]*i[0]).reshape(-1,1)
11    y_n=np.array([0]*i[1]).reshape(-1,1)
12    X=np.vstack((X_p,X_n))
13    y=np.vstack((y_p,y_n))
14    plt.scatter(X_p[:,0],X_p[:,1])
15    plt.scatter(X_n[:,0],X_n[:,1],color='red')
16
17
18    clf1 = LogisticRegression(C=0.001)
19    clf1.fit(X, y)
20    plot_svc_decision_function(clf1, color = 'red');
21
22    clf2 = LogisticRegression(C=1)
23    clf2.fit(X, y)
24    plot_svc_decision_function(clf2, color = 'blue');
25
26    clf3 = LogisticRegression(C=100)
```

```
27     clf3.fit(X, y)
28     plot_svc_decision_function(clf3, color = 'green');
29
30 plt.show()
```



Logistic Regression



Observation

- Green Line ($C = 100$) seems to be the best fitting line.
- Red Line ($C = 0.001$) seems nowhere to be seen.
- Blue Line ($C = 1$) seems to be working just OK.
- In case of very Imbalanced data, none of the models seems to work.
- With increase in C , the model also improves.