# 8E and 8F: Finding the Probability P(Y==1|X)

## 8E: Implementing Decision Function of SVM RBF Kernel

After we train a kernel SVM model, we will be getting support vectors and their corresponding coefficients
$\alpha_i$
Check the documentation for better understanding of these attributes:

https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html



As a part of this assignment you will be implementing the `decision_function()` of kernel SVM, here decision_function() means based on the value return by `decision_function()` model will classify the data point either as positive or negative

Ex 1: In logistic regression After traning the models with the optimal weights$w$ we get, we will find the value $\frac{1}{1+\exp(-(wx+b))}$, if this value comes out to be < 0.5 we will mark it as negative class, else its positive class

Ex 2: In Linear SVM After traning the models with the optimal weights$w$ we get, we will find the value of $sign(wx+b)$, if this value comes out to be -ve we will mark it as negative class, else its positive class.

Similarly in Kernel SVM After traning the models with the coefficients$\alpha_i$ we get, we will find the value of $sign(\sum_{i=1}^{n}(y_i\,\alpha_i K(x_i,x_q)) + intercept)$, here $K(x_i,x_q)$ is the RBF kernel. If this value comes out to be -ve we will mark $x_q$ as negative class, else its positive class.

RBF kernel is defined as: $K(x_i,x_q) = exp(-\gamma||x_i - x_q||^2)$

For better understanding check this link: https://scikit-learn.org/stable/modules/svm.html#svm-mathematical-formulation </font>

## Task E

1. Split the data into $X_{train}$(60), $X_{cv}$(20), $X_{test}$(20)
2. Train $SVC(gamma$ on the $(X_{train}, y_{train})$
   $= 0.001, C$
   $= 100.)$
3. Get the decision boundry values $f_{cv}$ on the $X_{cv}$ data i.e. `fcv` `= decision_function(`
   $X_{cv}$ `)` you need to implement this decision_function()

In [1]:

```python
import numpy as np
import pandas as pd
from sklearn.datasets import make_classification
import numpy as np
from sklearn.svm import SVC
```

In [2]:

```python
X, y = make_classification(n_samples=5000, n_features=5, n_redundant=2,
                           n_classes=2, weights=[0.7], class_sep=0.7, random_state=1
5)
```

## Pseudo code

clf = SVC(gamma=0.001, C=100.)
clf.fit(Xtrain, ytrain)

def decision_function(Xcv, ...): #use appropriate parameters
    for a data point $x_q$ in Xcv:
        #write code to implement
(                     , here the values
$\sum_{i=1}^{\text{all the support vectors}}$
$(y_i \alpha_i K(x_i, x_q))$
$+ intercept)$

$y_i$,
$\alpha_i$, and
$intercept$ can be obtained from the trained model
return # the decision_function output for all the data points in the Xcv

fcv = decision_function(Xcv, ...) # based on your requirement you can pass any other parameters

**Note**: Make sure the values you get as fcv, should be equal to outputs of clf.decision_function(Xcv)

In [3]:

```python
# you can write your code here
from sklearn.model_selection import train_test_split
X_train,X_test,Y_train,Y_test = train_test_split(X,y,stratify = y,test_size = 0.2)
X_train,X_cv,Y_train,Y_cv = train_test_split(X_train,Y_train,stratify = Y_train,tes
t_size = 0.25)


clf = SVC(gamma = 0.001,C = 100)
clf.fit(X_train,Y_train)
f_cv_master = clf.decision_function(X_cv)
intercept = clf.intercept_
```

In [4]:

```
SV = clf.support_vectors_
alpha_y_i = clf.dual_coef_
alpha_y_i = alpha_y_i.reshape(-1,1)
alpha_y_i.shape
```

Out[4]:

```
(577, 1)
```

In [5]:

```
from scipy.spatial import distance
from tqdm import tqdm
import math

def decision_function(X_cv,SV,alpha_y_i,gamma,intercept):
    f_cv = []
    gamma = 0.001

    for x_q in tqdm(X_cv):
        sum1 = 0
        for k,x_i in enumerate(SV):
            sum1 += alpha_y_i[k]*(math.exp(-gamma*(distance.euclidean(x_q,SV[k])**2))
)

        sum1+=intercept
        f_cv.append(round(sum1[0],8))

    return f_cv

gamma = 0.001
f_cv_slave = decision_function(X_cv,SV,alpha_y_i,gamma,intercept)
```

```
100%|
████████████████████████████████████████████████████████████████| 100
0/1000 [00:10<00:00, 91.73it/s]
```

In [6]:

```
print(f_cv_master[:10])
print(f_cv_slave[:10])
```

```
[ 0.62865451 -1.32893918 -3.08056443 -4.11320069 -1.43270507 -3.71791514
 -1.59637489 -2.7666321  -5.32448941 -4.31882056]
[0.62865451, -1.32893918, -3.08056443, -4.11320069, -1.43270507, -3.71791514, -1.59
637489, -2.7666321, -5.32448941, -4.31882056]
```

# 8F: Implementing Platt Scaling to find P(Y==1|X)

Check this PDF

Let the output of a learning method be $f(x)$. To get calibrated probabilities, pass the output through a sigmoid:

$$P(y = 1|f) = \frac{1}{1 + exp(Af + B)} \tag{1}$$

where the parameters $A$ and $B$ are fitted using maximum likelihood estimation from a fitting training set $(f_i, y_i)$. Gradient descent is used to find $A$ and $B$ such that they are the solution to:

$$argmin_{A,B}\{-\sum_i y_i log(p_i) + (1 - y_i)log(1 - p_i)\}, \quad (2)$$

where

$$p_i = \frac{1}{1 + exp(Af_i + B)} \quad (3)$$

Two questions arise: where does the sigmoid train set come from? and how to avoid overfitting to this training set?

If we use the same data set that was used to train the model we want to calibrate, we introduce unwanted bias. For example, if the model learns to discriminate the train set perfectly and orders all the negative examples before the positive examples, then the sigmoid transformation will output just a 0,1 function. So we need to use an independent calibration set in order to get good posterior probabilities. This, however, is not a draw back, since the same set can be used for model and parameter selection.

To avoid overfitting to the sigmoid train set, an out-of-sample model is used. If there are $N_+$ positive examples and $N_-$ negative examples in the train set, for each training example Platt Calibration uses target values $y_+$ and $y_-$ (instead of 1 and 0, respectively), where

$$y_+ = \frac{N_+ + 1}{N_+ + 2}; \; y_- = \frac{1}{N_- + 2} \quad (4)$$

For a more detailed treatment, and a justification of these particular target values see (Platt, 1999).

# TASK F

1. Apply SGD algorithm with ($f_{cv}$, $y_{cv}$) and find the weight $W$ intercept $b$ `Note: here our data is of one dimensional so we will have a one dimensional weight vector i.e W.shape (1,)`

Note1: Don't forget to change the values of $y_{cv}$ as mentioned in the above image. you will calculate y+, y- based on data points in train data

Note2: the Sklearn's SGD algorithm doesn't support the real valued outputs, you need to use the code that was done in the `'Logistic Regression with SGD and L2'` Assignment after modifying loss function, and use same parameters that used in that assignment.

```python
def log_loss(w, b, X, Y):
    N = len(X)
    sum_log = 0
    for i in range(N):
        sum_log += Y[i]*np.log10(sig(w, X[i], b)) + (1-Y[i])*np.log10(1-sig(w, X[i], b))
    return -1*sum_log/N
```

if Y[i] is 1, it will be replaced with y+ value else it will replaced with y- value

1. For a given data point from $X_{test}$, $P(Y = 1|X)$ where $f_{test} =$

$$= \frac{1}{1 + exp(-(W*f_{test}+b))}$$

`decision_function(` $X_{test}$ `)`, W and b will be learned as metioned in the above step

**Note: in the above algorithm, the steps 2, 4 might need hyper parameter tuning, To reduce the complexity of the assignment we are excluding the hyerparameter tuning part, but intrested students can try that**

If any one wants to try other calibration algorithm istonic regression also please check these tutorials

1. http://fa.bianp.net/blog/tag/scikit-learn.html#fn:1
2. https://drive.google.com/open?id=1MzmA7QaP58RDzocBORBmRiWfl7Co_VJ7
3. https://drive.google.com/open?id=133odBinMOIVb_rh_GQxxsyMRyW-Zts7a
4. https://stat.fandom.com/wiki/Isotonic_regression#Pool_Adjacent_Violators_Algorithm

In [16]:

```python
import math

def sigmoid(z):
    sig = 1/(1+math.exp(-z))

    return sig
```

In [17]:

```python
def gradient_dw(x,y,w,b,alpha,N):
    '''In this function, we will compute the gardient w.r.to w '''

    a = w[0]*x
    temp = x *(y-sigmoid(a + b))
    dw = temp-(w* (alpha/N))

    return dw
```

In [18]:

```python
def gradient_db(x,y,w,b):
    '''In this function, we will compute gradient w.r.to b '''
    a = w*x

    db = y-sigmoid(a+b)
    return db
```

In [19]:

```python
from tqdm import tqdm
def train(X_train,y_train,epochs,alpha,eta0):

    w = np.array([0])
    b = 0

    for i in tqdm(range(epochs)):
        for j in range(len(X_train)):
            dw = gradient_dw(X_train[j],y_train[j],w,b,alpha,len(X_train))
            db = gradient_db(X_train[j],y_train[j],w,b)

            w = w + eta0*dw
            b = b + eta0*db

    return w,b
```

In [20]:

```python
Y_cv1 = []

for i in range(len(Y_cv)):
    Y_cv1.append(Y_cv[i])
```

```
y_p = Y_cv1.count(1)
y_n = Y_cv1.count(0)

pos = float(y_p + 1)/(y_p+2)
neg = 1/float(y_n+2)

for i in range(len(Y_cv1)):
    if(Y_cv1[i]==1):
        Y_cv1[i] = pos
    else:
        Y_cv1[i] = neg
print(pos,neg)
Y_cv2 = np.array(Y_cv1)
```

0.9967213114754099 0.001430615164520744

```
alpha=0.0001
eta0=0.0001
N=len(Y_cv2)
epochs=50

w,b = train(Y_cv2,f_cv_slave,epochs,alpha,eta0)
```

```
100%|
██████████████████████████████████████████████████████████████████████████████|
50/50 [00:00<00:00, 114.45it/s]
```

```
print(w,b)
```

[1.50877492] -7.065664560900679

```
import math
f_test = decision_function(X_test,SV,alpha_y_i,gamma,intercept)
P_X_1 = [1/(1+math.exp(-(w*x + b))) for x in f_test]
```

```
100%|
███████████████████████████████████████████████████████████████████████████| 100
0/1000 [00:10<00:00, 92.95it/s]
```

```
print(P_X_1[:10])
```

[2.655744200366622e-05, 0.012780531936096448, 2.3849817855945682e-05, 3.34702752351
2707e-05, 0.387939916985917, 6.770887444003207e-06, 9.030394258307848e-05, 9.659566
356423562e-06, 0.03677936537330718, 0.000301762277604018303]

OBSERVATIONS The above are the probability values of every test query point after caliberation.