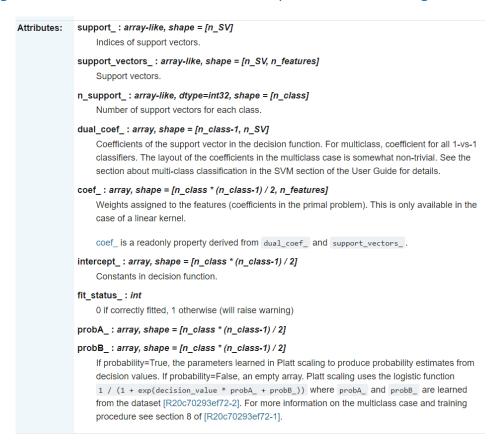# 8E and 8F: Finding the Probability P(Y==1|X)

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

After we train a kernel SVM model, we will be getting support vectors and their corresponsing coefficients $\alpha_i$

Check the documentation for better understanding of these attributes:
https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html (https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html)

| Attributes: | |
| --- | --- |
| | **support_ : array-like, shape = [n_SV]** |
| | Indices of support vectors. |
| | **support_vectors_ : array-like, shape = [n_SV, n_features]** |
| | Support vectors. |
| | **n_support_ : array-like, dtype=int32, shape = [n_class]** |
| | Number of support vectors for each class. |
| | **dual_coef_ : array, shape = [n_class-1, n_SV]** |
| | Coefficients of the support vector in the decision function. For multiclass, coefficient for all 1-vs-1 classifiers. The layout of the coefficients in the multiclass case is somewhat non-trivial. See the section about multi-class classification in the SVM section of the User Guide for details. |
| | **coef_ : array, shape = [n_class * (n_class-1) / 2, n_features]** |
| | Weights assigned to the features (coefficients in the primal problem). This is only available in the case of a linear kernel. |
| | `coef_` is a readonly property derived from `dual_coef_` and `support_vectors_` . |
| | **intercept_ : array, shape = [n_class * (n_class-1) / 2]** |
| | Constants in decision function. |
| | **fit_status_ : int** |
| | 0 if correctly fitted, 1 otherwise (will raise warning) |
| | **probA_ : array, shape = [n_class * (n_class-1) / 2]** |
| | **probB_ : array, shape = [n_class * (n_class-1) / 2]** |
| | If probability=True, the parameters learned in Platt scaling to produce probability estimates from decision values. If probability=False, an empty array. Platt scaling uses the logistic function `1 / (1 + exp(decision_value * probA_ + probB_))` where `probA_` and `probB_` are learned from the dataset [R20c70293ef72-2]. For more information on the multiclass case and training procedure see section 8 of [R20c70293ef72-1]. |

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 (https://scikit-learn.org/stable/modules/svm.html#svm-mathematical-formulation)

## Task E

1. Split the data into $X_{train}$ (60), $X_{cv}$ (20), $X_{test}$ (20)

2. Train $SVC(gamma = 0.001, C = 100.)$ on the $(X_{train}, y_{train})$

3. Get the decision boundry values $f_{cv}$ on the $X_{cv}$ data i.e. $f_{cv}$ = `decision_function(` $X_{cv}$ `)` you need to implement this decision_function()

In [105]:
```python
import numpy as np
import pandas as pd
from sklearn.datasets import make_classification
import numpy as np
from sklearn.svm import SVC
import math
from sklearn.model_selection import train_test_split
from tqdm import tqdm
import random
```

In [88]:
```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=15)

train_X,test_X,train_y,test_y  = train_test_split(X,y,test_size=0.20,random_state=42)
train_X,train_CV,train_y,test_CV  = train_test_split(train_X,train_y,test_size=0.20,random_state=42)
```

In [89]:
```python
s = pd.DataFrame(test_CV)
```

In [90]:
```python
s[0].value_counts()
```

Out[90]:
```
0    569
1    231
Name: 0, dtype: int64
```

## 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 $(\sum_{i=1}^{\text{all the support vectors}} (y_i \alpha_i K(x_i, x_q)) + intercept)$, here the values $y_i$, $\alpha_i$, 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 [91]:

```python
model = SVC(gamma=0.001, C=100.0, probability=True)
model.fit(train_X, train_y)

model_dec_value = model.decision_function(train_CV)

prams= model.get_params()
b =  model.intercept_
sv = model.support_vectors_
a = model.dual_coef_

custom_dec_val = []

def decision_function(x_cv, c, gamma):
    #print(x_cv)
    for cv in x_cv:
    #for i, cv in tqdm(enumerate(x_cv)):
        value = 0
        for j in range(len(sv)):
            l2_norm = np.linalg.norm(cv - sv[j])
            kernel = np.exp(-prams['gamma'] * (l2_norm**2))
            value += a[0][j] * kernel
        value = value + b
        custom_dec_val.append(value[0])

fcv = decision_function(train_CV, 100.0, 0.001)
```

```
In [92]: print(prams)
         print(b)
         print(sv)
         print(a)
```

```
{'C': 100.0, 'break_ties': False, 'cache_size': 200, 'class_weight': None, 'coef0': 0.0, 'decision_function_shape': 'ovr', 'degree': 3, 'gamma':
0.001, 'kernel': 'rbf', 'max_iter': -1, 'probability': True, 'random_state': None, 'shrinking': True, 'tol': 0.001, 'verbose': False}
[-1.80151873]
[[-0.45249961 -0.55624639 -0.18877168 -0.26346616  0.11359858]
 [ 0.71918122 -0.90416915 -0.18571414 -0.27647336 -0.18785388]
 [-0.97685553  0.1003584  -0.09967914 -0.12004848  0.39077859]
 ...
 [ 1.16313433 -0.29206536 -0.05199842 -0.07929317 -0.08525518]
 [ 0.21190377 -0.01954856 -0.01010565 -0.01360925  0.014668  ]
 [ 0.48453654  0.00724127  0.03051939  0.03859359 -0.08777591]]
[[-100.        -100.        -100.        -100.        -100.
  -100.        -100.        -100.        -100.        -100.
  -100.        -100.        -100.        -100.        -100.
  -100.        -100.        -100.        -100.        -100.
  -100.        -100.        -100.        -100.        -100.
  -100.        -100.        -100.        -100.        -100.
  -100.        -100.        -100.        -100.        -100.
  -100.        -100.        -100.        -100.        -100.
  -100.        -100.        -100.        -100.        -100.
```

```
In [93]: compare_decision_fun = pd.DataFrame(data= [model.decision_function(train_CV),custom_dec_val]).T
         #sample results
         compare_decision_fun.rename({0:'Cf_inbuilt',1:'manual_cf'},axis=1)[:5]
```

Out[93]:

|   | Cf_inbuilt | manual_cf |
|---|-----------|-----------|
| 0 | 0.877196  | 0.877196  |
| 1 | -1.120046 | -1.120046 |
| 2 | -3.563691 | -3.563691 |
| 3 | -0.631567 | -0.631567 |
| 4 | -1.192167 | -1.192167 |

Type *Markdown* and LaTeX: $\alpha^2$

**As we have binary target variables, Decision function concludes as when the data points gives postive points which in the postive side of hyperplane , simlarly negative**

**means Hyperplane in negative direction**

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

Check this PDF (https://drive.google.com/open?id=133odBinMOIVb_rh_GQxxsyMRyW-Zts7a)

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

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

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:

$$\underset{A,B}{argmin}\{-\sum_i y_i log(p_i) + (1 - y_i)log(1 - p_i)\}, \qquad (2)$$

where

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

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 ex-
ample, if the model learns to discriminate the train set per-
fectly and orders all the negative examples before the posi-
tive examples, then the sigmoid transformation will output
just a 0,1 function. So we need to use an independent cali-
bration 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 train-
ing example Platt Calibration uses target values $y_+$ and $y_-$
(instead of 1 and 0, respectively), where

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

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

## TASK F

4. 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

5. For a given data point from $X_{test}$, $P(Y = 1|X) = \frac{1}{1+exp(-(W*f_{test}+b))}$ where $f_{test}$ = `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 (http://fa.bianp.net/blog/tag/scikit-learn.html#fn:1)
2. https://drive.google.com/open?id=1MzmA7QaP58RDzocB0RBmRiWfl7Co_VJ7 (https://drive.google.com/open?id=1MzmA7QaP58RDzocB0RBmRiWfl7Co_VJ7)
3. https://drive.google.com/open?id=133odBinMOIVb_rh_GQxxsyMRyW-Zts7a (https://drive.google.com/open?id=133odBinMOIVb_rh_GQxxsyMRyW-Zts7a)
4. https://stat.fandom.com/wiki/Isotonic_regression#Pool_Adjacent_Violators_Algorithm (https://stat.fandom.com/wiki/Isotonic_regression#Pool_Adjacent_Violators_Algorithm)

```python
In [94]: #Calculate Y+ and y-
         N_pos =0
         N_neg = 0
         for i in test_CV:
             if i ==1:
                 N_pos+=1
             else:
                 N_neg+=1
```

```python
In [95]: Y_pos = (N_pos+1)/(N_pos+2)
         Y_neg  =1/(N_neg+2)
```

```python
In [96]: # Replace Postive test CV with Y_pos and 0 with Y_neg
         Updated_testCV = []
         for i in test_CV:
             if i ==1:
                 Updated_testCV.append(Y_pos)
             else:
                 Updated_testCV.append(Y_neg)
```

In [97]:

```python
def initialize_weights():
    w = np.zeros(1)
    b = 0
    return w,b

def sigmoid(z):
    return 1 / (1 + np.exp(-z))

def logloss(y_true, y_pred):
    loss = 0
    for index in range(len(y_true)):
        a = (y_true[index] * math.log(y_pred[index], 10)) + \
        (1 - y_true[index]) * math.log(1 - y_pred[index], 10)
        b = (-1/len(y_true))
        loss = loss + a * b
    return loss

def gradient_dw(x, y, w, b):
    error = y - sigmoid(np.dot(x, w.T) + b)
    dw = x * error
    return dw

def gradient_db(x, y, w, b):
    db = y - sigmoid(np.dot(x, w.T) + b)
    return db

def probability(x, w, b):
    predicted = []
    probability = []
    for i in range(len(x)):
        z = np.dot(w, x[i]) + b
        sig = sigmoid(z)
        if sig >= 0.5:
            predicted.append(1)
        else:
            predicted.append(0)
        probability.append(sig)
    return np.array(predicted), np.array(probability)

def find_accuracy(actual, predicted):
    return round((np.sum((actual == predicted) /
```

```python
                                             len(actual))) * 100, 2)

def train(X_test, y_test, epochs, alpha, lr):

    global test_loss
    test_loss = []

    w, b = initialize_weights()
    for ep in range(epochs):
        print('Epoch:', ep + 1)
        for index in range(len(X_test)):
            r_index = random.randint(0, len(X_test) - 1)

            ln_eqn = np.dot(X_test[r_index], w.T) + b
            error = y_test[r_index] - sigmoid(ln_eqn)

            dw = X_test[r_index] * error
            db = error

            w = w + (alpha * dw)
            b = b + (alpha * db)


        predicted, score = probability(X_test, w, b)
        loss = logloss(y_test, score)
        test_loss.append(loss)
        te_acc = find_accuracy(y_test, predicted)
        print('Accuracy',te_acc)


    return w, b
```

In [98]:
```python
lr   = 0.0001
alpha = 0.0001
N = len(Updated_testCV)
epochs = 10

w, b = train(Updated_testCV, test_CV, epochs, alpha, lr)
```

```
Epoch: 1
Accuracy 71.12
Epoch: 2
Accuracy 71.12
Epoch: 3
Accuracy 71.12
Epoch: 4
Accuracy 71.12
Epoch: 5
Accuracy 71.12
Epoch: 6
Accuracy 71.12
Epoch: 7
Accuracy 71.12
Epoch: 8
Accuracy 71.12
Epoch: 9
Accuracy 71.12
Epoch: 10
Accuracy 71.12
```

In [99]:
```python
w
```

Out[99]: array([0.11633561])

In [100]:
```python
b
```

Out[100]: array([-0.15534377])

In [148]:
```python
# Calibrate probablities
prob_values_calibrate = []
for (test_value) in custom_dec_val:
    #print(test_value)
    kernel = np.exp((-(w * (test_value+b))))
    prob = (1/(1+kernel))
    #print(prob)
    prob_values_calibrate.append((prob))
```

In [125]:
```python
len(prob_values_calibrate)
```

Out[125]: 800

In [152]:
```python
df_prob_calibriate = pd.DataFrame(data={'Prob_score_1':prob_values_calibrate,'Actual_testCV':test_CV})
```

In [154]:
```python
df_prob_calibriate['Prob_score_1'] = df_prob_calibriate['Prob_score_1'].str[0]
```

In [158]:
```python
df_prob_calibriate.head(5)
```

Out[158]:

|   | Prob_score_1 | Actual_testCV |
|---|---|---|
| 0 | 0.520982 | 1 |
| 1 | 0.462975 | 0 |
| 2 | 0.393492 | 0 |
| 3 | 0.477130 | 0 |
| 4 | 0.460889 | 0 |

**The empirical results show that after calibration gives best probablities which matches actual test_cv**

In [ ]: