

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

α_i

Check the documentation for better understanding of these attributes:

<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 $1 / (1 + \exp(\text{decision_value} * \text{probA_} + \text{probB_}))$ 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 $\text{sign}(wx + b)$, 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 α_i we get, we will find the value of $\text{sign}(\sum_{i=1}^n (y_i \alpha_i K(x_i, x_q)) + \text{intercept})$, 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>

TASK 2

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 boundary values f_{cv} on the X_{cv} data i.e. $f_{cv} = \text{decision_function}(X_{cv})$ **you need to implement this decision_function()**

In [19]:

```
import numpy as np
import pandas as pd
from sklearn.datasets import make_classification
import numpy as np
from numpy import linalg
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split
from math import exp
import math
```

In [20]:

```
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=
5)
```

In [27]:

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.40,
random_state=15)
X_test, X_cv, y_test, y_cv=train_test_split(X_test, y_test, test_size=.50,
random_state=15)
```

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 [29]:

```

# you can write your code here
clf = SVC(gamma=0.001, C=100)
clf.fit(X_train,y_train)
support_vectors=clf.support_vectors_
intercept=clf.intercept_
dual_coe=clf.dual_coef_

def kernal(mat1, mat2, gamma):

    k=exp(-gamma*np.sum((mat1-mat2)**2))

    return k

gamma=0.001
def decision_function(gamma,X_cv,dual_coe,intercept,support_vectors):
    result_dec_fn=[]
    for j in X_cv:
        tmp=0
        for i,k in zip(support_vectors,dual_coe[0]):
            Ker=kernal(i,j,gamma)
            tmp+=(k*Ker)
        tmp=tmp+intercept
        result_dec_fn.append(tmp)
    return result_dec_fn

result_dec_funcn=decision_function(gamma,X_cv,dual_coe,intercept,support_vectors)
result_list=list(i[0] for i in result_dec_funcn)
print("output of decision funcation ",result_list[0:10])
F_cv =[]
for i in result_list:
    if i>0:
        F_cv.append(1)
    else:
        F_cv.append(0)

print("output of model ",clf.decision_function(X_cv)[0:10])

output of decision funcation [-2.15980217966647, -2.6315643052289905, -3.610567063
094037, 1.9254014334365632, -0.9839950024528423, -1.84103624132939, -3.047397298728
9804, -0.9108054741150926, -1.0738614096121768, -2.8494875264428563]
output of model [-2.15980218 -2.63156431 -3.61056706 1.92540143 -0.983995 -1.84
103624
-3.0473973 -0.91080547 -1.07386141 -2.84948753]

```

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:

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

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:

$$\text{argmin} \{ - \sum u_i \log(n_i) + (1 - u_i) \log(1 - n_i) \} \quad (2)$$

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.

```
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 be 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(Xtest)`, W and b will be learned as mentioned in the above step

```

no_of_plus=np.count_nonzero(y_train==1)
no_of_minus=np.count_nonzero(y_train==0)

y_plus=(no_of_plus+1)/(no_of_minus+2)
y_minus=1/(no_of_minus+2)

```

In [31]:

```

def predict(y_cv):
    y_cv_predict=np.where(y_cv==0,y_minus,y_plus)
    return y_cv_predict

y_cv_pred=predict(y_cv)

F_test=clf.decision_function(X_test)

```

In [33]:

```

def initialize_weights(dim):
    w=np.zeros_like(dim)
    b=0
    return w,b

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

def logloss(y_true,y_pred):

    loss =0
    for i in range(len(y_true)):
        temp=y_true[i]*math.log10(y_pred[i])+(1-y_true[i])*math.log10(1-y_pred[i])
        loss+=temp
    loss=(-1*loss)/len(y_true)
    return loss

def gradient_dw(x,y,w,b,alpha,N):
    '''In this function, we will compute the gradient w.r.to w'''
    dw=x*(y-sigmoid(np.dot(w,x)+b)) - (alpha/N)*w
    return dw

def gradient_db(x,y,w,b):
    '''In this function, we will compute gradient w.r.to b'''
    db=y-(sigmoid(np.dot(w,x)+b))
    return db

def train(X_train,y_train,X_test,y_test,epochs,alpha,eta0):
    w,b = initialize_weights(X_train[0])
    train_loss = []
    test_loss = []
    for e in range(epochs):

        for x,y in zip(X_train,y_train):
            dw = gradient_dw(x,y,w,b,alpha,N)
            db = gradient_db(x,y,w,b)
            w = w + (eta0 * dw)
            b = b + (eta0 * db)

        train_pred = []

        for i in X_train:
            y_pred= sigmoid(np.dot(w,i) + b)

```

```

train_pred.append(y_pred)
loss1=logloss(y_train, train_pred)
train_loss.append(loss1)

```

```

return w,b,train_loss

```

```

alpha=0.0001
eta0=0.0001
N=len(result_list)
epochs=100
w,b,train_loss=train(result_list,y_cv_pred ,F_test,y_test,epochs,alpha,eta0)
print("W={}  intercept={} ".format(w,b))

print("train_loss =",train_loss)

```

```

W=0.7836275055532742  intercept=-1.1641575400972337
train_loss = [0.2738898730674267, 0.25292270156343827, 0.23653510496818148, 0.22353
8334974367, 0.21306885628487032, 0.20450422190302345, 0.1973944243909419, 0.1914112
6585906384, 0.18631259588648674, 0.18191745822119545, 0.1780888604132711, 0.1747217
6924067384, 0.17173467282292865, 0.16906358533578106, 0.16665773729347866, 0.164476
44049690363, 0.16248678059572427, 0.16066189939558828, 0.1589797022027556, 0.157421
87492274887, 0.15597312933760984, 0.15462061821151918, 0.15335347804529983, 0.15216
246867761588, 0.15103968701810322, 0.14997833800423616, 0.1489725500831522, 0.14801
7225599711103, 0.14710791874588794, 0.14624073541959998, 0.145412250611909, 0.144619
4399029328, 0.1438596223794855, 0.14313041284987219, 0.1424296816658115, 0.14175552
079947107, 0.14110621508807586, 0.14048021776675423, 0.1398761295750996, 0.13929268
085409002, 0.13872871615493573, 0.13818318096579313, 0.13765511023041202, 0.1371436
1838808004, 0.1366478907092743, 0.13616717573831205, 0.13570077868457778, 0.1352480
556288902, 0.13480840843226496, 0.13438128025149873, 0.1339661515803394, 0.13356253
67469777, 0.1331699808086555, 0.13278805679263767, 0.13241636323994632, 0.132054522
0142925, 0.1317021763437721, 0.1313589890672654, 0.13102464106119846, 0.13069882982
551373, 0.13038126821043153, 0.13007168326792545, 0.12976981521386893, 0.1294754164
8853514, 0.1291882509046684, 0.1289080928736239, 0.12863472670123563, 0.12836794594
603904, 0.12810755283335007, 0.12785335771944092, 0.12760517860072793, 0.1273628406
6344375, 0.12712617586978112, 0.12689502257694074, 0.12666922518590154, 0.126448633
81707616, 0.12623310401032153, 0.1260224964470405, 0.12581667669234725, 0.125615514
95548254, 0.1254188858668472, 0.12522666827019263, 0.12503874502865037, 0.124855002
84341938, 0.1246753320840374, 0.12449962662928052, 0.12432778371781236, 0.124159703
80780333, 0.12399529044480304, 0.1238344501372222, 0.12367709223883644, 0.123523128
83778344, 0.12337247465156713, 0.1232250469276289, 0.12308076534908659, 0.122939551
94527487, 0.12280133100675049, 0.12266602900446225, 0.12253357451280165, 0.12240389
813627972]

```

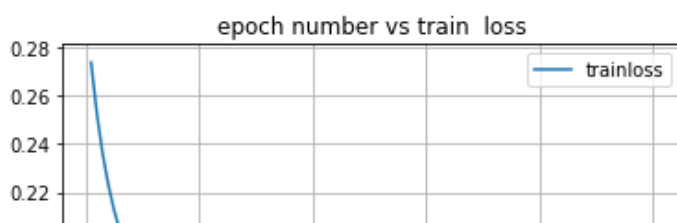
In [34]:

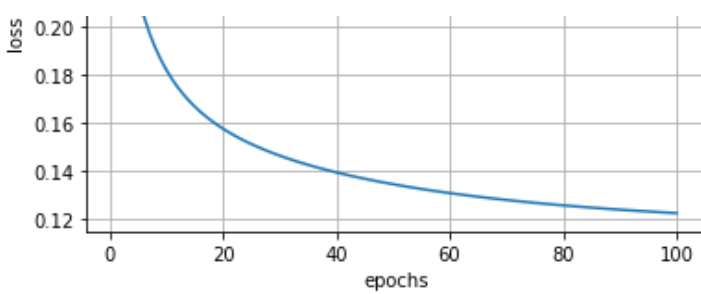
```

import matplotlib.pyplot as plt

epochs=list(range(1,101))
plt.plot(epochs, train_loss, label="trainloss")
plt.legend()
plt.xlabel("epochs")
plt.ylabel("loss")
plt.title("epoch number vs train loss")
plt.grid()
plt.show()

```





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 hyperparameter tuning part, but interested students can try that

If any one wants to try other calibration algorithm isotonic 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=1MzmA7QaP58RDzocBORBmRiWfI7Co_VJ7
3. 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)