# 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

| Attributes: | |
| --- | --- |
| | **support_** : *array-like, shape = [n_SV]*<br>Indices of support vectors. |
| | **support_vectors_** : *array-like, shape = [n_SV, n_features]*<br>Support vectors. |
| | **n_support_** : *array-like, dtype=int32, shape = [n_class]*<br>Number of support vectors for each class. |
| | **dual_coef_** : *array, shape = [n_class-1, n_SV]*<br>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]*<br>Weights assigned to the features (coefficients in the primal problem). This is only available in the case of a linear kernel.<br><br>coef_ is a readonly property derived from `dual_coef_` and `support_vectors_`. |

Saved successfully!                                    ✕

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]*<br>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

## ▾ 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()

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

Saved successfully! ✕

```
from sklearn.model_selection import train_test_split


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)


X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=0)
X_train,X_valid,y_train,y_valid=train_test_split(X_train,y_train,test_size=0.2,random_state=0
```

## ▾ 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 $(\sum_{i=1}^{all\ the\ support\ vectors}(y_i\alpha_i K(x_i, x_q)) + intercept)$, 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)

```python
from functools import partial

def kernel(a,b,gamma):
  return np.exp(-gamma*np.square(a-b).sum())

def decision_function(query,dual_coef,support_vectors,intercept,gamma):
  total=0
  for i,coef in enumerate(dual_coef.ravel()):
    total+=coef*kernel(support_vectors[i],query,gamma)

  total+=intercept
  return 0 if total<0 else 1

def predict(X):
  predictions=[]
  for query in X:
    predictions.append(
        decision_function(query,dual_coef,support_vectors,intercept,gamma))
  return predictions
```

Saved successfully!                              ✕

```python
gamma=0.001
svc=SVC(gamma=gamma,C=100)
svc.fit(X_train,y_train)
dual_coef,support_vectors,intercept=svc.dual_coef_,svc.support_vectors_,svc.intercept_
```

```python
(preds==svc_preds).all
```

```
<function ndarray.all>
```

```python
%%time
preds=predict(X)
svc_preds=svc.predict(X)
```

```
CPU times: user 20.5 s, sys: 33.4 ms, total: 20.5 s
Wall time: 20.5 s
```

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

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:

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

where

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

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 good posterior probabilities. This, back, since the same set can be used for model and parameter selection.

Saved successfully! ×

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

$$y_+ = \frac{N_+ + 1}{N_+ + 2}; \; y_- = \frac{1}{N_- + 2} \tag{4}$$

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

Check this [PDF](#)

▾ 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**

Saved successfully! ✕

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=1MzmA7QaP58RDzocB0RBmRiWfl7Co_VJ7

3. https://drive.google.com/open?id=133odBinMOIVb_rh_GQxxsyMRyW-Zts7a

4. https://stat.fandom.com/wiki/Isotonic_regression#Pool_Adjacent_Violators_Algorithm

```python
def sigmoid(f,A,B):
  return 1/(1+np.exp(-A*f+B))
```

```python
preds=svc.predict(X_valid)
y_pos,y_neg=preds[preds==1],preds[preds==0]
n_pos,n_neg=len(y_pos),len(y_neg)
```

```
  y_pos=(n_pos+1)/(n_pos+2)
  y_neg=1/(n_neg+2)



def gradient(x,y,w,b):
  return (x*(y-sigmoid(x,W,b)))


%%time

epoch=500000
eta0=0.0001
alpha=0.0001

#Lets initialize the weights
W,b=0,0

#lets train the model
for _ in range(epoch):

  #creating a random sample from training set
  random_index=np.random.randint(len(y_valid))
  f_rand,y_rand=preds[random_index],y_valid[random_index]

  #updating weights and biases with Gradients
  if y_rand==1:
    W+=eta0*gradient(f_rand,y_pos,W,b)
    b+=eta0*gradient(f_rand,y_pos,W,b)
```

Saved successfully!                    ✕

```
    W+=eta0*gradient(f_rand,y_neg,W,b)
    b+=eta0*gradient(f_rand,y_neg,W,b)

  #lets regularize the weights
  W-=eta0*alpha*np.square(W).sum()/len(preds)

    CPU times: user 10.9 s, sys: 29.3 ms, total: 10.9 s
    Wall time: 11.1 s


test_preds=svc.predict(X_test)
probabilities=sigmoid(test_preds,W,b)
```

✓ 0s completed at 2:48 PM ● ✕

Saved successfully! ✕