# DS-GA 1003 HW6

Yiyan Chen

April 2018

# 1 Reformulations of Multiclass Hinge Loss

## 1.1 Multiclass setting review

## 1.2 Two versions of multiclass hinge loss (or generalized hinge loss)

### 1.2.1

$$
\begin{aligned}
l_2(h, (x_i, y_i)) &= max_{y \in Y}[\Delta(y_i, y) + h(x_i, y) - h(x_i, y_i)] \\
&= max(\Delta(y_i, y_i) + h(x_i, y_i) - h(x_i, y_i), max_{y \in Y - y_i}[\Delta(y_i, y) + h(x_i, y) - h(x_i, y_i)]) \\
&= max(0, max_{y \in Y - y_i}[\Delta(y_i, y) + h(x_i, y) - h(x_i, y_i)]) \\
&= max(0, max_{y \in Y - y_i}(\Delta(y_i, y) - m_{i,y}(h)) \\
&= max_{y \in Y - y_i}(max[0, \Delta(y_i, y), m_{i,y}(h)]) \\
&= l_1(h, (x_i, y_i))
\end{aligned}
$$

### 1.2.2

(a):

$$
\begin{aligned}
l_1(h, (x_i, y_i)) &= max_{y \in Y - y_i}(max[0, \Delta(y_i, y) + h(x_i, y) - h(x_i, y_i)]) \\
&= max_{y \in Y - y_i}(max[0]) \\
&= 0
\end{aligned}
$$

$$
\begin{aligned}
l_2(h, (x_i, y_i)) &= max(\Delta(y_i, y_i) + h(x_i, y_i) - h(x_i, y_i), max_{y \in Y - y_i}(\Delta(y_i, y) + h(x_i, y) - h(x_i, y_i))) \\
&= max(0, max_{y \in Y - y_i}(< 0)) \\
&= 0
\end{aligned}
$$

(b):

$$h(x_i, y) \leq h(x_i, y_i) - \Delta(y_i, y)$$
$$\max h(x_i, y) = h(x_i, y_i) - \Delta(y_i, y)$$
$$\max h(x_i, y) = h(x_i, y_i) - 0$$
$$\max h(x_i, y) = h(x_i, y_i)$$
$$argmax_{y \in Y} h(x_i, y_i) = y_i$$

# 2 SGD for Multiclass Liner SVM

## 2.1 Optional

The first part of the function: $\lambda||w||^2$ is convex since its derivative $2\lambda w > 0$. The second part of the function is the sum of convex functions. Hence, two parts combined is the convex function.

## 2.2

Set g as the subgradient:

$$J(w) = J_1(w) + J_2(w)$$

$$J_1(w) = \lambda ||w||^2, J_2(w) = \frac{1}{n} \sum_{i=1}^{n} max_{y \in Y}[\Delta(y_i, y) + < w, \Psi(x_i, y) - \Psi(x_i, y_i) >]$$

$$J_2(w_1) = J_2(w_2) + g(w_1 - w_2)$$

$$\sum_{i=1}^{n}[\Delta(y_i, \hat{y_i} + < w_1, \Psi(x_i, y) - \Psi(x_i, y_i) >] = \sum_{i=1}^{n}[\Delta(y_i, \hat{y_i} + < w_2, \Psi(x_i, \hat{y}) - \Psi(x_i, y_i) >] + ng(w_1 - w_2)$$

$$g = \frac{\sum_{i=1}^{n} < w_1 - w_2, \Psi(x_i, \hat{y}) - \Psi(x_i, y_i) >}{n(w_1 - w_2)}$$

$$= \frac{1}{n} \sum_{i=1}^{n} \frac{(w_1 - w_2)^T(\Psi(x_i, \hat{y}) - \Psi(x_i, y_i))}{w_1 - w_2}$$

$$= \frac{1}{n} \sum_{i=1}^{n} [\Psi(x_i, \hat{y}) - \Psi(x_i, y_i)]$$

$$g = 2\lambda w + \frac{1}{n} \sum_{i=1}^{n} [\Psi(x_i, \hat{y}) - \Psi(x_i, y_i)]$$

## 2.3

The stochastic subgradient iteration takes one point $(x_i, y_i)$ at a time:

$$w \longleftarrow w - \eta g$$
$$\longleftarrow w - \eta(2\lambda w + [\Psi(x_i, \hat{y}) - \Psi(x_i, y_i)])$$

## 2.4

Minibatch subgradient iteration takes $(x_i, y_y) \ldots (x_{i+m-1}, y_{i+m-1})$:

$$w \longleftarrow w - \eta g$$

$$\longleftarrow w - \eta(2\lambda w + \frac{1}{m} \sum_{i=j}^{j+m-1} [\Psi(x_i, \hat{y}) - \Psi(x_i, y_i)])$$

4

# 3 [Optional]

# 4 Muliclass Classification - Implementation

## 4.1 One-vs-All

```python
from sklearn.base import BaseEstimator, ClassifierMixin, clone

class OneVsAllClassifier(BaseEstimator, ClassifierMixin):
    """
    One-vs-all classifier
    We assume that the classes will be the integers 0,..,(n_classes-1).
    We assume that the estimator provided to the class, after fitting, has a "decision_func
    returns the score for the positive class.
    """
    def __init__(self, estimator, n_classes):
        """
        Constructed with the number of classes and an estimator (e.g. an
        SVM estimator from sklearn)
        @param estimator : binary base classifier used
        @param n_classes : number of classes
        """
        self.n_classes = n_classes
        self.estimators = [clone(estimator) for _ in range(n_classes)]
        self.fitted = False

    def fit(self, X, y=None):
        """
        This should fit one classifier for each class.
        self.estimators[i] should be fit on class i vs rest
        @param X: array-like, shape = [n_samples,n_features], input data
        @param y: array-like, shape = [n_samples,] class labels
        @return returns self
        """
        #Your code goes here
        for i in range(self.n_classes):
            y_i = np.where(y==i, 1,0)
            self.estimators[i].fit(X,y_i)
        self.fitted = True
        return self


    def decision_function(self, X):
        """
        Returns the score of each input for each class. Assumes
```
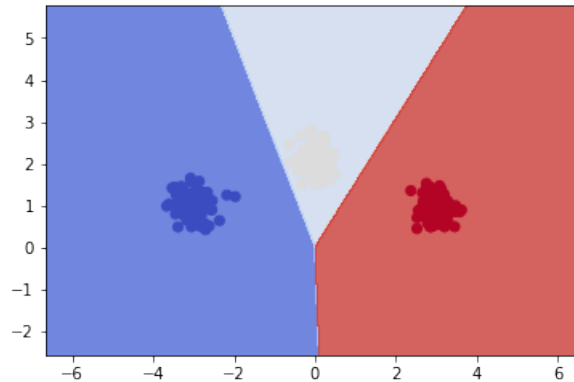
```python
40          that the given estimator also implements the decision_function method (which sklearn
41          and that fit has been called.
42          @param X : array-like, shape = [n_samples, n_features] input data
43          @return array-like, shape = [n_samples, n_classes]
44          """
45          if not self.fitted:
46              raise RuntimeError("You must train classifer before predicting data.")
47
48          if not hasattr(self.estimators[0], "decision_function"):
49              raise AttributeError(
50                  "Base estimator doesn't have a decision_function attribute.")
51
52          #Replace the following return statement with your code
53          if len(self.estimators) == 1:
54              score = self.estimator[0].decision_function(X)
55          else:
56              score = np.zeros([X.shape[0], self.n_classes])
57              for i in range(self.n_classes):
58                  score.T[i] = self.estimators[i].decision_function(X)
59          return score

61      def predict(self, X):
62          """
63          Predict the class with the highest score.
64          @param X: array-like, shape = [n_samples,n_features] input data
65          @returns array-like, shape = [n_samples,] the predicted classes for each input
66          """
67          #Replace the following return statement with your code
68          score = self.decision_function(X)
69          pred = np.zeros(X.shape[0])
70          for i in range(X.shape[0]):
71              zipped = zip(score[i], range(self.n_classes))
72              pred[i] = sorted(zipped, key=itemgetter(0))[-1][1]
73          return pred
```

```
Coeffs 0
[[-1.05852747 -0.90296521]]
Coeffs 1
[[ 0.22117096 -0.38900908]]
Coeffs 2
[[ 0.89162796 -0.82467394]]

array([[100,   0,   0],
       [  0, 100,   0],
       [  0,   0, 100]])
```

## 4.2 Multiclass SVM

```python
1  def zeroOne(y,a) :
2      '''
3      Computes the zero-one loss.
4      @param y: output class
5      @param a: predicted class
6      @return 1 if different, 0 if same
7      '''
8      return int(y != a)
9
10 def featureMap(X,y,num_classes) :
11     '''
12     Computes the class-sensitive features.
13     @param X: array-like, shape = [n_samples,n_inFeatures] or [n_inFeatures,], input featur
14     @param y: a target class (in range 0,..,num_classes-1)
15     @return array-like, shape = [n_samples,n_outFeatures], the class sensitive features for
16     '''
17     #The following line handles X being a 1d-array or a 2d-array
18     num_samples, num_inFeatures = (1,X.shape[0]) if len(X.shape) == 1 else (X.shape[0],X.sha
19     #your code goes here, and replaces following return
20     #x = (x1, x2), classes = (0,1,2)-> (x, 0) = (x1,x2, 0, 0, 0, 0)
21     n_outFeatures = num_classes*num_inFeatures
22     rn = np.zeros([num_samples, n_outFeatures])
23     y=np.array([y])
24     if num_samples == 1:
25         X = X.reshape((1, -1))
26     for i in range(num_samples):
27         if len(y) > 1:
28             y_i =y[i]
29         else:
```

```
30              y_i = y
31              rn[i][int(y_i)*num_inFeatures:int(y_i)*num_inFeatures+num_inFeatures] = X[i]
32      return rn
33
34  def sgd(X, y, num_outFeatures, subgd, eta = 0.001, T = 10000):
35      '''
36      Runs subgradient descent, and outputs resulting parameter vector.
37      @param X: array-like, shape = [n_samples,n_features], input training data
38      @param y: array-like, shape = [n_samples,], class labels
39      @param num_outFeatures: number of class-sensitive features
40      @param subgd: function taking x,y and giving subgradient of objective
41      @param eta: learning rate for SGD
42      @param T: maximum number of iterations
43      @return: vector of weights
44      '''
45      num_samples = X.shape[0]
46      #your code goes here and replaces following return statement
47      w=np.zeros(num_outFeatures)
48      for t in range(T):
49          i = np.random.randint(num_samples)
50          w = w-eta*subgd(X[i], y[i], w)
51      return w
52
53
54  class MulticlassSVM(BaseEstimator, ClassifierMixin):
55      '''
56      Implements a Multiclass SVM estimator.
57      '''
58      def __init__(self, num_outFeatures, lam=1.0, num_classes=3, Delta=zeroOne, Psi=featureMa
59          '''
60          Creates a MulticlassSVM estimator.
61          @param num_outFeatures: number of class-sensitive features produced by Psi
62          @param lam: l2 regularization parameter
63          @param num_classes: number of classes (assumed numbered 0,..,num_classes-1)
64          @param Delta: class-sensitive loss function taking two arguments (i.e., target marg
65          @param Psi: class-sensitive feature map taking two arguments
66          '''
67          self.num_outFeatures = num_outFeatures
68          self.lam = lam
69          self.num_classes = num_classes
70          self.Delta = Delta
71          self.Psi = lambda X,y : Psi(X,y,num_classes)
72          self.fitted = False
73
74      def subgradient(self,x,y,w):
75          '''
```

```python
76              Computes the subgradient at a given data point x,y
77              @param x: sample input
78              @param y: sample class
79              @param w: parameter vector
80              @return returns subgradient vector at given x,y,w
81              '''
82              #Your code goes here and replaces the following return statement
83              y_hat = 0
84              w = w.reshape((1,-1))
85              for i in range(self.num_classes):
86                  cal = self.Delta(y, i)+np.dot(w, (self.Psi(x,i)-self.Psi(x,y)).T)
87                  if cal>y_hat:
88                      y_hat = i
89              g = 2*self.lam*w+self.Psi(x,y_hat)-self.Psi(x,y)
90              return g


        def fit(self,X,y,eta=0.001,T=10000):
94              '''
95              Fits multiclass SVM
96              @param X: array-like, shape = [num_samples,num_inFeatures], input data
97              @param y: array-like, shape = [num_samples,], input classes
98              @param eta: learning rate for SGD
99              @param T: maximum number of iterations
100             @return returns self
101             '''
102             self.coef_ = sgd(X,y,self.num_outFeatures,self.subgradient,eta,T)
103             self.fitted = True
104             return self

106         def decision_function(self, X):
107             '''
108             Returns the score on each input for each class. Assumes
109             that fit has been called.
110             @param X : array-like, shape = [n_samples, n_inFeatures]
111             @return array-like, shape = [n_samples, n_classes] giving scores for each sample,cla
112             '''
113             if not self.fitted:
114                 raise RuntimeError("You must train classifer before predicting data.")

116             #Your code goes here and replaces following return statement
117             score = np.zeros([X.shape[0], self.num_classes])
118             for sample in range(X.shape[0]):
119                 for cls in range(self.num_classes):
120                     score[sample][cls] = np.dot(self.coef_, (self.Psi(X[sample], cls)).T)
121             return score
```
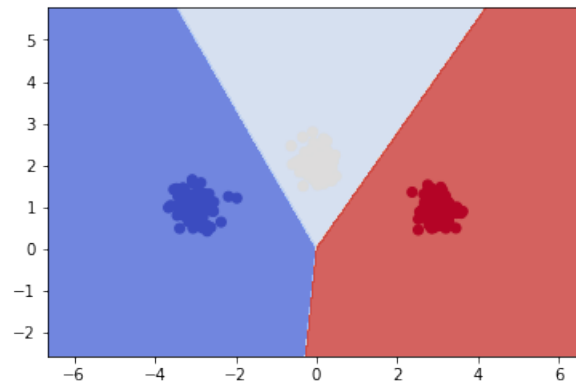
```
122
123    def predict(self, X):
124        '''
125        Predict the class with the highest score.
126        @param X: array-like, shape = [n_samples, n_inFeatures], input data to predict
127        @return array-like, shape = [n_samples,], class labels predicted for each data poin
128        '''
129        #Your code goes here and replaces following return statement
130        pred = np.zeros(X.shape[0])
131        score = self.decision_function(X)
132        for i in range(X.shape[0]):
133            zipped = zip(score[i], range(self.num_classes))
134            pred[i] = sorted(zipped, key=itemgetter(0))[-1][1]
135        return pred
```

```
w:
[[-0.45890788 -0.0553265  -0.01807324  0.20763369  0.47698112 -0.15230719]]

array([[100,   0,   0],
       [  0, 100,   0],
       [  0,   0, 100]])
```

# 5 [Optional]

# 6 [Optional]

# 7 Gradient Boosting Machines

## 7.1

$$(-g_m)_i = -\frac{\partial}{\partial f_{m-1}(x_i)} l(y_i, f_{m-1}(x_i))$$

$$= (y_i - f_{m-1}(x_i))$$

$$h_m = argmin_{h \in F} \sum_{i=1}^{n} [(-g_m)_i - h(x_i)]^2$$

$$= argmin_{h \in F} \sum_{i=1}^{n} [(y_i - f_{m-1}(x_i)) - h(x_i)]^2$$

## 7.2

$$\frac{\partial ln(1 + e^{-yf(x)})}{\partial f(x)} = \frac{-ye^{-yf(x)}}{1 + e^{-yf(x)}}$$

$$(-g_m)_i = \frac{y_i}{e^{y_i f(x_i)} + 1}$$

$$h_m = argmin_{h \in F} \sum_{i=1}^{n} [(-g_m)_i - h(x_i)]^2$$

$$= argmin_{h \in F} \sum_{i=1}^{n} [\frac{y_i}{e^{y_i f(x_i)} + 1} - h(x_i)]^2$$

# 8 Gradient Boosting Implementation

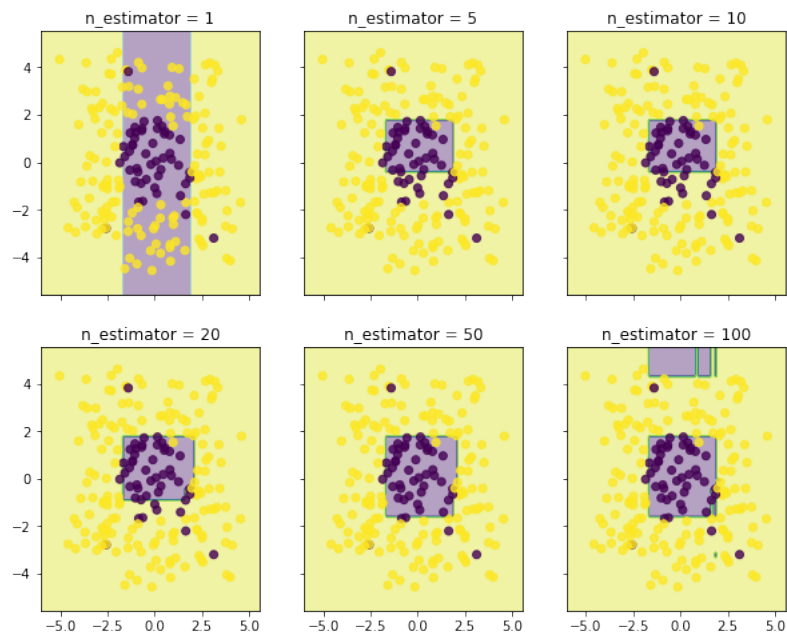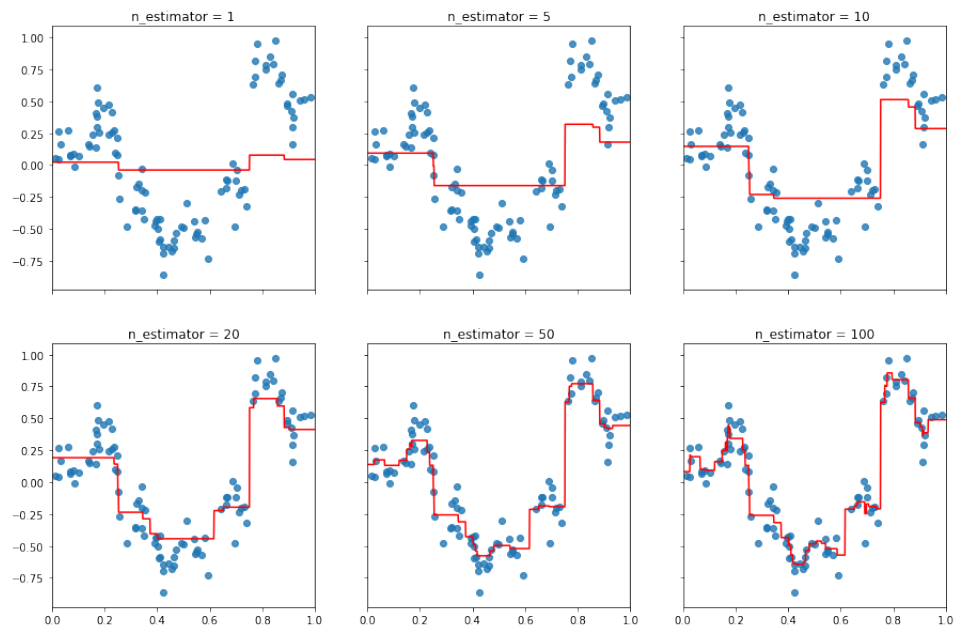## 8.1

```python
class gradient_boosting():
    '''
    Gradient Boosting regressor class
    :method fit: fitting model
    '''
    def __init__(self, n_estimator, pseudo_residual_func, learning_rate=0.1, min_sample=5, n
        '''
        Initialize gradient boosting class

        :param n_estimator: number of estimators (i.e. number of rounds of gradient boosting
        :pseudo_residual_func: function used for computing pseudo-residual
        :param learning_rate: step size of gradient descent
        '''
        self.n_estimator = n_estimator
        self.pseudo_residual_func = pseudo_residual_func
        self.learning_rate = learning_rate
        self.min_sample = min_sample
        self.max_depth = max_depth


    def fit(self, train_data, train_target):
        '''
        Fit gradient boosting model
        '''
        # Your code goes here
        #n_estimator - M, n - train_data.shape[0]
        fm = np.zeros(train_data.shape[0])
#          print(train_target.shape)
        train_target = train_target.squeeze()
#          print(train_target.shape)
        self.hm = []
        for i in range(self.n_estimator):
            rgs = DecisionTreeRegressor(min_samples_split=self.min_sample,
            max_depth=self.max_depth)
            rgs.fit(train_data, self.pseudo_residual_func(train_target, fm))
            self.hm.append(rgs)
            h = rgs.predict(train_data)
            fm+=self.learning_rate*h
        return self

    def predict(self, test_data):
```

```
42          '''
43          Predict value
44          '''
45          # Your code goes here
46          pred = np.zeros(test_data.shape[0])
47          for i in range(self.n_estimator):
48              rgs = self.hm[i]
49              pev = rgs.predict(test_data)
50              pred+= self.learning_rate*pev
51          return pred
```

## 8.2   [Optional]

```
1  def pseudo_residual_logistic(train_target, train_predict):
2      '''
3      Compute the pseudo-residual based on current predicted value.
4      '''
5      y= train_target
6      f = train_predict
7      return y/(np.exp(y*f)+1)
```