# DS-GA 1003: Machine Learning and Computational Statistics Homework 6: Generalized Hinge Loss and Multiclass SVM

**Due: Monday, April 11, 2016, at 6pm (Submit via NYU Classes)**

**Instructions**: Your answers to the questions below, including plots and mathematical work, should be submitted as a single file, either HTML or PDF. You may include your code inline or submit it as a separate file. You may either scan hand-written work or, preferably, write your answers using software that typesets mathematics (e.g. LATEX, LYX, or MathJax via iPython).

## 1   Introduction

This is an entirely written problem set, and relatively short. The goal of this problem set is to get more comfortable with the multiclass hinge loss and multiclass SVM. In several problems below, you are asked to justify that certain things are convex functions. For these problems, you may use any of the rules about convex functions described in our notes on Convex Optimization (`https://davidrosenberg.github.io/mlcourse/Notes/convex-optimization.pdf`) or in the Boyd and Vandenberghe book. In particular, you will need to make frequent use of the following result: If $f_1, \ldots, f_m : \mathbf{R}^n \to \mathbf{R}$ are convex, then their pointwise maximum

$$f(x) = \max \{f_1(x), \ldots, f_m(x)\}$$

is also convex.

## 2   Convex Surrogate Loss Functions

It's common in machine learning that the loss function we really care about leads to optimization problems that are not computationally tractable. The 0/1 loss for binary classification is one such example[1]. Since we have better machinery for minimizing convex functions, a standard approach is to find a **convex surrogate loss function.** A convex surrogate loss function is a convex function that is an upper bound for the loss function of interest[2]. If we can make the upper bound small,

---

[1]Interestingly, if our hypothesis space is linear classifiers and we are in the "realizable" case, which means that there is some hypothesis that achieves 0 loss (with the 0/1 loss), then we can efficiently find a good hypothesis using linear programming. This is not difficult to see: each data point gives a single linear constraint, and we are looking for a vector that satisfies the constraints for each data point.

[2]At this level of generality, you might be wondering: "A convex function of WHAT?". For binary classification, we usually are talking about a convex function of the margin. But to solve our machine learning optimization problems, we will eventually need our loss function to be a convex function of some $w \in \mathbf{R}^d$ that parameterizes our hypothesis space. It'll be clear in what follows what we're talking about.

then the loss we care about will also be small[3]. Below we will show that the multiclass hinge loss based on a class-sensitive loss $\Delta$ is a convex surrogate for the multiclass loss function $\Delta$, when we have a linear hypothesis space. We'll start with a special case, that the hinge loss is a convex surrogate for the 0/1 loss.

## 2.1 Hinge loss is a convex surrogate for 0/1 loss

1. Let $f : \mathcal{X} \to \mathbf{R}$ be a classification score function for binary classification.

   (a) For any example $(x, y) \in \mathcal{X} \times \{-1, 1\}$, show that

   $$1(y \neq \text{sign}(f(x)) \leq \max\{0, 1 - yf(x)\},$$

   where $\text{sign}(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases}$.

   SOLUTION: If $y = \text{sign}(f(x))$, then the LHS is 0 and the RHS is always at least 0. So the inequality holds. Otherwise, the LHS is 1 and $yf(x) \leq 0$, which means $1 - yf(x) \geq 1$. So the inequality again holds.

   (b) Show that the hinge loss $\max\{0, 1 - m\}$ is a convex function of the margin $m$.
   SOLUTION: $1 - m$ is an affine function and $0$ is a constant function. Both are convex, so their pointwise maximum is convex.

   (c) Suppose our prediction score functions are given by $f_w(x) = w^T x$. The hinge loss of $f_w$ on any example $(x, y)$ is then $\max\{0, 1 - yw^T x\}$. Show that this is a convex function of $w$.
   SOLUTION: $1 - yw^T x$ is an affine function of $w$, so by the same argument as for (b), the expression is a convex function of $w$.

## 2.2 Multiclass Hinge Loss

Consider the multiclass output space $\mathcal{Y} = \{1, \ldots, k\}$. Suppose we have a base hypothesis space $\mathcal{H} = \{h : \mathcal{X} \times \mathcal{Y} \to \mathbf{R}\}$ from which we select a compatibility score function. Then our final multiclass hypothesis space is $\mathcal{F} = \{f(x) = \arg\max_{y \in \mathcal{Y}} h(x, y) \mid h \in \mathcal{H}\}$. Since functions in $\mathcal{F}$ map into $\mathcal{Y}$, our action space $\mathcal{A}$ and output space $\mathcal{Y}$ are the same. Suppose we have a class-sensitive loss function $\Delta : \mathcal{Y} \times \mathcal{A} \to \mathbf{R}$. Even though $\mathcal{Y} = \mathcal{A}$, we write $\mathcal{Y} \times \mathcal{A}$ to indicate that the true class goes in the first argument of the function, while the prediction (i.e. the action) goes in the second slot. We do this because we don't assume that $\Delta(y, y') = \Delta(y', y)$. It would not be unusual to have this asymmetry in practice. For example, false alarms may be much less costly than no alarm when indeed something is going wrong.

Our ultimate goal would be to find $f \in \mathcal{F}$ minimizing the empirical cost-sensitive loss:

$$\min_{f \in \mathcal{F}} \sum_{i=1}^{n} \Delta(y_i, f(x_i)).$$

[3]This is actually fairly weak motivation for a convex surrogate. Much better motivation comes from the more advanced theory of **classification calibrated** loss functions. See Bartlett et al's paper "Convexity, Classification, and Risk Bounds." http://www.eecs.berkeley.edu/~wainwrig/stat241b/bartlettetal.pdf

Since binary classification with $0/1$ loss is intractable and is a special case of this formulation, we know that this more general formulation must also be computationally intractable. Thus we are looking for a convex surrogate loss function.

1. Suppose we have chosen an $h \in \mathcal{H}$, from which we get $f(x) = \arg\max_{y \in \mathcal{Y}} h(x, y)$. Justify that for any $x \in \mathcal{X}$ and $y \in \mathcal{Y}$, we have

$$h(x, y) \leq h(x, f(x)).$$

SOLUTION: Follows directly from the definition of $f(x)$.

2. Justify the following two inequalities:

$$\begin{aligned} \Delta\left(y, f(x)\right) &\leq \Delta\left(y, f(x)\right) + h(x, f(x)) - h(x, y) \\ &\leq \max_{y' \in \mathcal{Y}} \left[\Delta\left(y, y'\right)\right) + h(x, y') - h(x, y)\right] \end{aligned}$$

The RHS of the last expression is called the **generalized hinge loss:**

$$\ell\left(h, (x, y)\right) = \max_{y' \in \mathcal{Y}} \left[\Delta\left(y, y'\right)\right) + h(x, y') - h(x, y)\right].$$

We have shown that for any $x \in \mathcal{X}, y \in \mathcal{Y}, h \in \mathcal{H}$ we have

$$\ell\left(h, (x, y)\right) \geq \Delta(y, f(x)),$$

where, as usual, $f(x) = \arg\max_{y \in \mathcal{Y}} h(x, y)$. [You should think about why we cannot write the generalized hinge loss as $\ell\left(f, (x, y)\right)$.]

SOLUTION: First inequality is trivial from Part (1), since it implies $h(x, f(x)) - h(x, y) \geq 0$. Second inequality is also trivial, since we are replacing $f(x)$ with something that potentially makes the expression larger.

3. We now introduce a specific base hypothesis space $\mathcal{H}$ of linear functions. Consider a class-sensitive feature mapping $\Psi : \mathcal{X} \times \mathcal{Y} \to \mathbf{R}^d$, and $\mathcal{H} = \left\{h_w\left(x, y\right) = \langle w, \Psi(x, y) \rangle \mid w \in \mathbf{R}^d\right\}$. Show that we can write the generalized hinge loss for $h_w(x, y)$ on example $(x_i, y_i)$ as

$$\ell\left(h_w, (x_i, y_i)\right) = \max_{y \in \mathcal{Y}} \left[\Delta\left(y_i, y\right)\right) + \langle w, \Psi(x_i, y) - \Psi(x_i, y_i) \rangle\right].$$

SOLUTION: We just use linearity of inner product to write this with a single inner product in the maximum, rather than 2.

4. We will now show that the generalized hinge loss $\ell\left(h_w, (x_i, y_i)\right)$ is a convex function of $w$. Justify each of the following steps.

   (a) The expression $\Delta(y_i, y) + \langle w, \Psi(x_i, y) - \Psi(x_i, y_i) \rangle$ is an affine function of $w$.
       SOLUTION: With respect to $w$, $\Delta(y_i, y)$ is a constant scalar and $\Psi(x_i, y) - \Psi(x_i, y_i)$ is a constant vector. Thus we have an affine function of $w$.

   (b) The expression $\max_{y \in \mathcal{Y}} \left[\Delta\left(y_i, y\right)\right) + \langle w, \Psi(x_i, y) - \Psi(x_i, y_i) \rangle\right]$ is a convex function of $w$.
       SOLUTION: For each $y \in \mathcal{Y}$, we have an affine function of $w$, which isi convex. The pointwise maximum of convex functions is convex, and affine functions are convex.

5. Conclude that $\ell\left(h_w, (x_i, y_i)\right)$ is a convex surrogate for $\Delta(y_i, f_w(x_i))$.
   SOLUTION: We have shown that $\ell$ is a convex function of $w$ and gives an upper bound to the class-sensitive loss function $\Delta$.

# 3    Hinge Loss is a Special Case of Generalized Hinge Loss

Let $\mathcal{Y} = \{-1, 1\}$. Let $\Delta(y, \hat{y}) = 1(y \neq \hat{y})$. If $g(x)$ is the score function in our binary classification setting, then define our compatibility function as

$$
\begin{aligned}
h(x, 1) &= g(x)/2 \\
h(x, -1) &= -g(x)/2.
\end{aligned}
$$

Show that for this choice of $h$, the multiclass hinge loss reduces to hinge loss:

$$
\ell(h, (x, y)) = \max_{y' \in \mathcal{Y}} [\Delta(y, y')) + h(x, y') - h(x, y)] = \max\{0, 1 - yg(x)\}
$$

SOLUTION: We have

$$
\begin{aligned}
\ell(h, (x, y)) &= \max_{y' \in \mathcal{Y}} [\Delta(y, y')) + h(x, y') - h(x, y)] \\
&= \max\{\Delta(y, y), [\Delta(y, -y)) + h(x, -y) - h(x, y)]\} \\
&= \max\left\{0, \left[1 + \frac{1}{2}\begin{cases} -g(x) - g(x) & \text{for } y = 1 \\ g(x) - [-g(x)] & \text{for } y = -1 \end{cases}\right]\right\} \\
&= \max\left\{0, \left[1 + \begin{cases} -g(x) & \text{for } y = 1 \\ g(x) & \text{for } y = -1 \end{cases}\right]\right\} \\
&= \max\{0, 1 - yg(x)\}
\end{aligned}
$$

# 4    Another Formulation of Generalized Hinge Loss

In lecture we defined the **margin** of the compatibility score function $h$ on the $i$th example $(x_i, y_i)$ for class $y$ as

$$
m_{i,y}(h) = h(x_i, y_i) - h(x_i, y),
$$

and the loss on an individual example $(x_i, y_i)$ to be:

$$
\max_y \left[(\Delta(y_i, y) - m_{i,y}(h)))_+\right].
$$

Here we investigate whether this is just an instance of the generalized hinge loss $\ell(h, (x, y)))$ defined above.

1. Show that $\ell(h, (x_i, y_i)) = \max_{y' \in \mathcal{Y}} [\Delta(y_i, y')) - m_{i,y'}(h)]$.
   SOLUTION:

   $$
   \begin{aligned}
   \ell(h, (x_i, y_i)) &= \max_{y' \in \mathcal{Y}} [\Delta(y_i, y')) + h(x_i, y') - h(x_i, y_i)] \\
   &= \max_{y' \in \mathcal{Y}} [\Delta(y_i, y')) - m_{i,y'}(h)]
   \end{aligned}
   $$

2. Suppose $\Delta(y, y') \geq 0$ for all $y, y' \in \mathcal{Y}$. Show that for any example $(x_i, y_i)$ and any score function $h$, the multiclass hinge loss we gave in lecture and the generalized hinge loss presened above are equivalent, in the sense that

   $$
   \max_{y \in \mathcal{Y}} \left[(\Delta(y_i, y) - m_{i,y}(h)))_+\right] = \max_{y \in \mathcal{Y}} (\Delta(y_i, y) - m_{i,y}(h))).
   $$

4

(Hint: This is easy by piecing together other results we have already attained regarding the relationship between $\ell$ and $\Delta$.)

SOLUTION: By the assumption and by the fact that $\ell$ gives an upper bound to $\Delta$, and finally using the result of the previous problem, we get

$$0 \leq \Delta(y_i, f(x_i)) \leq \ell\left(h, (x_i, y_i)\right) = \max_{y \in \mathcal{Y}}\left(\Delta(y_i, y) - m_{i,y}(h)\right).$$

Since the maximum is nonnegative, the result is the same as if we had

$$\max_{y \in \mathcal{Y}}\left[\left(\Delta(y_i, y) - m_{i,y}(h)\right)_+\right]$$

3. In the context of the generalized hinge loss, $\Delta(y, y')$ is like the "target margin" between the score for true class $y$ and the score for class $y'$. Suppose that our prediction function $f$ gets the correct class on $x_i$. That is, $f(x_i) = \arg\max_{y' \in \mathcal{Y}} h(x_i, y') = y_i$. Furthermore, assume that all of our target margins are reached or exceeded. That is

$$m_{i,y}(h) = h(x_i, y_i) - h(x_i, y) \geq \Delta(y_i, y),$$

for all $y \neq y_i$. Show that $\ell\left(h, (x_i, y_i)\right) = 0$ if we assume that $\Delta(y, y) = 0$ for all $y \in \mathcal{Y}$.

SOLUTION:

$$
\begin{aligned}
\ell\left(h, (x_i, y_i)\right) &= \max_{y' \in \mathcal{Y}}\left[\Delta\left(y_i, y'\right) + h(x_i, y') - h(x_i, y_i)\right] \\[2mm]
&= \max\left\{\underbrace{\Delta\left(y_i, y_i\right)}_{\text{case:}y'=y_i}, \max_{y' \neq y_i}\left[\Delta\left(y_i, y'\right) - m_{i,y'}(h)\right]\right\} \\[2mm]
&= \max\left\{0, \max_{y' \neq y_i}\left[\underbrace{\Delta\left(y_i, y'\right) - m_{i,y'}(h)}_{\leq 0}\right]\right\} \\[2mm]
&= 0.
\end{aligned}
$$

# 5  SGD for Multiclass SVM

Suppose our output space and our action space are given as follows: $\mathcal{Y} = \mathcal{A} = \{1, \ldots, k\}$. Given a non-negative class-sensitive loss function $\Delta : \mathcal{Y} \times \mathcal{A} \to \mathbf{R}^{\geq 0}$ and a class-sensitive feature mapping $\Psi : \mathcal{X} \times \mathcal{Y} \to \mathbf{R}^d$. Our prediction function is $f : \mathcal{X} \to \mathcal{Y}$ is given by

$$f_w(x) = \arg\max_{y \in \mathcal{Y}} \langle w, \Psi(x, y)\rangle$$

1. For a training set $(x_1, y_1), \ldots (x_n, y_n)$, let $J(w)$ be the $\ell_2$-regularized empirical risk function for the multiclass hinge loss. We can write this as

$$J(w) = \lambda\|w\|^2 + \frac{1}{n}\sum_{i=1}^{n}\max_{y \in \mathcal{Y}}\left[\Delta\left(y_i, y\right)\right) + \langle w, \Psi(x_i, y) - \Psi(x_i, y_i)\rangle\right].$$

We will now show that that $J(w)$ is a convex function of $w$. Justify each of the following steps. As we've shown it in a previous problem, you may use the fact that $w \mapsto \max_{y \in \mathcal{Y}} [\Delta(y_i, y)) + \langle w, \Psi(x_i, y) - \Psi(x_i, y_i) \rangle]$ is a convex function.

(a) $\frac{1}{n} \sum_{i=1}^{n} \max_{y \in \mathcal{Y}} [\Delta(y_i, y)) + \langle w, \Psi(x_i, y) - \Psi(x_i, y_i) \rangle]$ is a convex function of $w$.
SOLUTION: Nonnegative combination of convex functions is convex.

(b) $\|w\|^2$ is a convex function of $w$.
SOLUTION: Norms in $\mathbf{R}^d$ are convex.

(c) $J(w)$ is a convex function of $w$.
SOLUTION: We have a nonnegative combination of convex functions, so it's convex.

2. Since $J(w)$ is convex, it has a subgradient at every point. Give an expression for a subgradient of $J(w)$. You may use any standard results about subgradients, including the result from an earlier homework about subgradients of the pointwise maxima of functions. (Hint: It may be helpful to refer to $\hat{y} = \arg\max_{y \in \mathcal{Y}} [\Delta(y_i, y)) + \langle w, \Psi(x_i, y) - \Psi(x_i, y_i) \rangle]$.)
SOLUTION: A subgradient of $J(w)$ at $w$ is given by the following expression for $g \in \mathbf{R}^d$:

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

3. Give an expression the stochastic subgradient based on the point $(x_i, y_i)$.
SOLUTION:
$$g_{\text{SGD}} = 2\lambda w + (\Psi(x_i, \hat{y}) - \Psi(x_i, y_i)).$$

4. Give an expression for a minibatch subgradient, based on the points $(x_i, y_i), \ldots, (x_{i+m-1}, y_{i+m-1})$.
SOLUTION:
$$g_{\text{SGD}} = 2\lambda w + \frac{1}{m} \sum_{j=i}^{i+m-1} (\Psi(x_j, \hat{y}) - \Psi(x_j, y_j))$$