

Machine Learning and Computational Statistics

Homework 5: Generalized Hinge Loss and Multiclass SVM

Due: Tuesday, April 11, 2017, at 10pm (Submit via Gradescope)

Instructions: Your answers to the questions below, including plots and mathematical work, should be submitted as a single PDF file. It's preferred that you write your answers using software that typesets mathematics (e.g. L^AT_EX, L^AT_EX, or MathJax via iPython), though if you need to you may scan handwritten work. You may find the [minted](#) package convenient for including source code in your L^AT_EX document. If you are using L^AT_EX, then the [listings](#) package tends to work better.

1 Introduction

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 functions are convex. 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 $f_1, \dots, f_m : \mathbf{R}^n \rightarrow \mathbf{R}$ are convex, then their pointwise maximum

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

is also convex.

2 Convex Surrogate Loss Functions

It's common in machine learning that the loss functions we really care about lead to optimization problems that are not computationally tractable. The 0/1 loss for binary classification is one such example¹. 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 interest². If we can make the upper bound small,

¹Interestingly, 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.

²At 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 small³. Below we will show that the multiclass hinge loss based on a class-sensitive loss Δ is a convex surrogate for the multiclass loss function Δ , 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 $f : \mathcal{X} \rightarrow \mathbf{R}$ be a classification score function for binary classification.

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

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

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

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 Generalized Hinge Loss

Consider the multiclass output space $\mathcal{Y} = \{1, \dots, k\}$. Suppose we have a base hypothesis space $\mathcal{H} = \{h : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbf{R}\}$ 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. Nevertheless, we will write our class-sensitive loss function as $\Delta : \mathcal{Y} \times \mathcal{A} \rightarrow \mathbf{R}$, even though $\mathcal{Y} = \mathcal{A}$. We do this to indicate that the true class goes in the first slot of the function, while the prediction (i.e. the action) goes in the second slot. This is important because we do not 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:

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

³This 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 both intractable and 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 the decision function $f(x) = \arg \max_{y \in \mathcal{Y}} h(x, y)$. 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:

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

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

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

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

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

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 $\ell(f, (x, y))$.]

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 $\Psi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbf{R}^d$, and $\mathcal{H} = \{h_w(x, y) = \langle w, \Psi(x, y) \rangle \mid w \in \mathbf{R}^d\}$. Show that we can write the generalized hinge loss for $h_w(x, y)$ on example (x_i, y_i) as

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

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 $\ell(h_w, (x_i, y_i))$ 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 $\max_{y \in \mathcal{Y}} [\Delta(y_i, y) + \langle w, \Psi(x_i, y) - \Psi(x_i, y_i) \rangle]$ is a convex function of w .
SOLUTION: For each $y \in \mathcal{Y}$, we have an affine function of w , which is convex. The pointwise maximum of convex functions is convex, and affine functions are convex.
5. Conclude that $\ell(h_w, (x_i, y_i))$ is a convex surrogate for $\Delta(y_i, f_w(x_i))$.
SOLUTION: We have shown that ℓ is a convex function of w and gives an upper bound to the class-sensitive loss function Δ .

3 SGD for Multiclass SVM

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

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

1. For a training set $(x_1, y_1), \dots, (x_n, y_n)$, let $J(w)$ be the ℓ_2 -regularized empirical risk function for the multiclass hinge loss. We can write this as

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

for some $\lambda > 0$. We will now show 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 $\hat{y}_i = \arg \max_{y \in \mathcal{Y}} [\Delta(y_i, y) + \langle w, \Psi(x_i, y) - \Psi(x_i, y_i) \rangle]$.)

SOLUTION: A subgradient of $J(w)$ at w is given by the following expression for $g \in \mathbf{R}^d$:

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

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 $(x_i, y_i), \dots, (x_{i+m-1}, y_{i+m-1})$.

SOLUTION:

$$g_{\text{SGD}} = 2\lambda w + \frac{1}{m} \sum_{j=i}^{i+m-1} (\Psi(x_j, \hat{y}) - \Psi(x_j, y_j))$$

4 [OPTIONAL] 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 $\ell(h, (x_i, y_i)) = \max_{y' \in \mathcal{Y}} [\Delta(y_i, y') - m_{i,y'}(h)]$. (In other words, it looks just like loss above, but without the positive part.)

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 presented above are equivalent, in the sense that

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

(Hint: This is easy by piecing together other results we have already attained regarding the relationship between ℓ and Δ .)

SOLUTION: By the assumption and by the fact that ℓ gives an upper bound to Δ , and finally using the result of the previous problem, we get

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

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

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

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 $y \neq y_i$. It seems like in this case, we should have 0 loss. This is almost the case. Show that $\ell(h, (x_i, y_i)) = 0$ if we assume that $\Delta(y, y) = 0$ for all $y \in \mathcal{Y}$.

SOLUTION:

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

5 [OPTIONAL] 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

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

6 Multiclass Classification - Implementation

In this problem we will work on a simple three-class classification example, similar to the one [given in lecture](#). The data is generated and plotted for you in the skeleton code.

6.1 One-vs-All (also known as One-vs-Rest)

In this problem we will implement one-vs-all multiclass classification. Our approach will assume we have a binary base classifier that returns a score, and we will predict the class that has the highest score.

1. Complete the class `OneVsAllClassifier` in the skeleton code. Following the `OneVsAllClassifier` code is a cell that extracts the results of the fit and plots the decision region. Include these results in your submission.
2. [Optional] Normalize the vectors corresponding to each of the linear SVM classifiers so that they have unit norm. Evaluate the results and plot the decision regions for these normalized vectors.
3. [Optional] You may notice that every time you run the cell that fits the models and plots the decision regions, you get slightly different results. This is more pronounced when C is larger, e.g. $C = 200$. Investigate and propose an explanation for this. You may use any means necessary, including google searching and asking other people (just like in real life). [It's generally good to investigate things that look odd – you'll almost always learn something, and sometimes you'll uncover a more serious underlying problem.]

6.2 Multiclass SVM

In this question, we will implement stochastic subgradient descent for the linear multiclass SVM described in lecture and in this problem set. We will use the class-sensitive feature mapping approach with the “multivector construction”, as described in our [multiclass classification lecture](#) and in SSBD Section 17.2.1.

1. Complete the skeleton code for multiclass SVM. Following the multiclass SVM implementation, we have included another block of test code. Make sure to include the results from these tests in your assignment, along with your code.

7 [Optional] Audio Classification

In this problem, we will work on the urban sound dataset [URBANSOUND8K](#) from the Center for Urban Science and Progress (CUSP) at NYU. We will first extract features from raw audio data using the [LibROSA](#) package, and then we will train multiclass classifiers to classify the sounds into 10 sound classes. URBANSOUND8K dataset contains 8732 labeled sound excerpts. For this problem, you may use the file `UrbanSound8K.csv` to randomly sample 2000 examples for training and 2000 examples for validation.

1. In LibROSA, there are many functions for visualizing audio waves and spectra, such as `display.waveplot()` and `display.specshow()`. Load a random audio file from each class as a floating point time series with `librosa.load()`, and plot their waves and [linear-frequency power spectrogram](#). If you are interested, you can also play the audio in the notebook with functions `display()` and `Audio()` in `IPython.display`.

2. **Mel-frequency cepstral coefficients (MFCC)** are a commonly used feature for sound processing. We will use MFCC and its first and second differences (like discrete derivatives) as our features for classification. First, use function `feature.mfcc()` from LibROSA to extract MFCC features from each audio sample. (The first MFCC coefficient is typically discarded in sound analysis, but you do not need to. You can test whether this helps in the optional problem below.) Next, use function `feature.delta()` to calculate the first and second differences of MFCC. Finally, combine these features and normalize each feature to zero mean and unit variance.
3. Train a linear multiclass SVM on your 2000 example training set. Evaluate your results on the validation set in terms of 0/1 error and generate a confusion table. Compare the results to a one-vs-all classifier using a binary linear SVM as the base classifier. For each model, may use your code from the previous problem, or you may use another implementation (e.g. from sklearn).
4. [More Optional] Compare results to any other multiclass classification methods of your choice.
5. [More Optional] Try different feature sets and see how they affect performance.