

CS 475 Machine Learning: Homework 4  
Online Learning, SVMs and Kernels  
Due: Wednesday October 24, 2012, 12:00pm  
100 Points Total      Version 1.2

## 1 Programming (50 points)

In this assignment, you will implement a kernelized version of logistic regression. Your implementation will train the model using a gradient ascent algorithm.

### 1.1 Review: Logistic Regression

The logistic regression model is used to model binary classification data. Logistic regression is a special case of generalized linear regression where the labels  $Y$  are modeled as a linear combination of the data  $X$ , but in a transformed space specified by  $g$ , sometimes called the “link function”:

$$P(y = 1|\mathbf{x}, \mathbf{w}) = g(\mathbf{w} \cdot \mathbf{x}) \quad (1)$$

This “link function” allows you to model inherently non-linear data with a linear model. In the case of logistic regression, the link function is the logistic function:

$$g(z) = \frac{1}{1 + e^{-z}} \quad (2)$$

### 1.2 Kernel Logistic Regression

The above model relies on a linear function of weight parameters and instances:  $\mathbf{w} \cdot \mathbf{x}$ . However, some data can be highly non-linear, which is impossible to learn using a linear classifier. In class, we learned about the kernel trick, which uses the dual of a linear classifier for kernel learning. Kernels project data into high dimensional spaces in which a linear separator may exist. The decision boundary learned by the linear classifier is linear in the high dimensional space but non-linear in the original feature space.

A kernel is a function which maps the input data into a new space using a basis function and computes the inner product between two basis functions. A kernel function is defined as:

$$K(x, x') = (\phi(x) \cdot \phi(x')^T)$$

In this assignment you will implement a dual version of the logistic regression model so that you can learn with a kernel. In class, we derived the updates for the dual perceptron, and we learned that the linear function of weights  $\mathbf{w}$  and instances  $\mathbf{x}$  can be rewritten:

$$\mathbf{w} \cdot \mathbf{x} = \sum_{i=1}^N \alpha_i y_i K(\mathbf{x}_i, \mathbf{x})$$

This dual formulation also holds in logistic regression. Under the kernel logistic regression model, we have:

$$P(y = 1|\mathbf{x}, \alpha) = g\left(\sum_{i=1}^N \alpha_i K(\mathbf{x}_i, \mathbf{x})\right) \quad (3)$$

$$P(y = 0|\mathbf{x}, \alpha) = 1 - g\left(\sum_{i=1}^N \alpha_i K(\mathbf{x}_i, \mathbf{x})\right) \quad (4)$$

Note that the  $y_i$  terms from the Perceptron formula are not needed here. Instead of learning  $\mathbf{w}$ , you will learn  $\alpha_i$  for each example  $\mathbf{x}_i$ . The next subsection describes the learning procedure you will use.

### 1.2.1 Training with Gradient Ascent

In this assignment, we will solve for the parameters  $\alpha$  in our logistic regression model using gradient ascent to find the maximum likelihood estimate.

Gradient ascent is an optimization technique that is both very simple and powerful. This optimization algorithm works by taking the gradient of the objective function and taking steps in directions where the gradient is positive, which increases the objective function.<sup>1</sup> Under certain conditions, gradient ascent is guaranteed to asymptotically converge to a local maximum of the objective function (or a global maximum if the objective is convex). This strategy is often referred to as “hill climbing.”

Given a likelihood function  $\ell(Y, X, \alpha)$ , we need to compute the gradient  $\nabla \ell(Y, X, \alpha)$ . The partial derivative with respect to the parameter  $\alpha_k$  is:

$$\frac{\partial \ell}{\partial \alpha_k} = \sum_{i=1}^N y_i g\left(-\sum_{j=1}^N \alpha_j K(x_j, x_i)\right) (K(x_i, x_k)) + (1 - y_i) g\left(\sum_{j=1}^N \alpha_j K(x_j, x_i)\right) (-K(x_i, x_k)) \quad (5)$$

You can see how this is derived in the appendix of this handout.

In gradient descent, the update rule is  $\alpha'_j = \alpha_j + \eta \frac{\partial}{\partial \alpha_j}$  for all values of  $j$  (one per training instance), where  $\eta$  is called the *learning rate*. Choosing  $\eta$  intelligently is a non-trivial task, and there are many ways to choose it in the literature. In this assignment, we will use a constant scalar  $\eta$ . See section 1.3.2 for more details.

This update is repeated for a given number of training iterations (see 1.3.3).

### 1.2.2 Making Predictions

Given an instance  $\mathbf{x}_i$ , you can predict a label  $y_i$  as the  $y$  value with higher probability under the model. You should choose a prediction  $\hat{y}_i$  according to:

$$\hat{y}_i = \arg \max_y P(y_i = y|\mathbf{x}_i, \alpha) \quad (6)$$

---

<sup>1</sup>In machine learning, we can either maximize or minimize an objective. In this assignment, we are *maximizing* an objective function (likelihood) as part of maximum likelihood estimation. One might also want to *minimize* an objective function that we think of as a *loss* or *error*. In this case, the algorithm is called gradient **descent**. You can trivially convert a maximization problem into a minimization problem by negating the objective function, so the terms gradient descent and gradient ascent are sometimes used interchangeably.

In other words, if  $g(\sum_{j=1}^N \alpha_j K(\mathbf{x}_j, \mathbf{x}_i)) \geq 0.5$ , you should predict  $y_i$  to be 1, otherwise you predict  $y_i$  to be 0.

### 1.2.3 Kernel Functions

You will implement three kernel functions for use with kernel logistic regression. Each version of the kernel logistic regression algorithm will be its own class, but it is highly recommended that most of the code be in a parent class: `KernelLogisticRegression` to avoid code duplication. Each subclass will overload and implement the necessary functions for their specific kernels.

1. Linear Kernel:  $K(x, x') = (x \cdot x'^T)$

Implement this algorithm as the `LinearKernelLogisticRegression` class. The command line argument for the `algorithm` parameter should be `logistic_regression_linear_kernel`.

2. Polynomial Kernel:  $K(x, x') = (1 + (xx'^T))^d$

Implement this algorithm as the `PolynomialKernelLogisticRegression` class. The command line argument for the `algorithm` parameter should be `logistic_regression_polynomial_kernel`. The parameter  $d$  is explained in section 1.3.1.

3. Gaussian Kernel:  $K(x, x') = \exp(-||x - x'||^2/2\sigma^2)$

Implement this algorithm as the `GaussianKernelLogisticRegression` class. The command line argument for the `algorithm` parameter should be `logistic_regression_gaussian_kernel`. The parameter  $\sigma$  is explained in section 1.3.1.

Note that the naive implementation may compute  $K(x_i, x_j)$  many times during training. To improve training efficiency you should store (cache) computed values of  $K(x_i, x_j)$  in a matrix,  $G_{ij} = K(x_i, x_j)$ .  $G$  is called a Gram Matrix.

## 1.3 Implementation Details

### 1.3.1 Kernel Parameters

The parameters  $d$  for the polynomial kernel and  $\sigma$  for the Gaussian kernel are set using command line options.

For the polynomial kernel parameter  $d$ , add this command line option by adding the following code block to the `createCommandLineOptions` method of `Classify`.

```
registerOption("polynomial_kernel_exponent", "double", true, "The exponent of the polynomial kernel.");
```

You can then read the value from the command line by adding the following to the main method of `Classify`:

```
double polynomial_kernel_exponent = 2;
if (CommandLineUtilities.hasArg("polynomial_kernel_exponent"))
    polynomial_kernel_exponent = CommandLineUtilities.getOptionValueAsFloat("polynomial_kernel_exponent");
```

For the Gaussian kernel parameter  $\sigma$ , add this command line option by adding the following code block to the `createCommandLineOptions` method of `Classify`.

```
registerOption("gaussian_kernel_sigma", "double", true, "The sigma of the Gaussian kernel.");
```

You can then read the value from the command line by adding the following to the main method of `Classify`:

```
double gaussian_kernel_sigma = 1;
if (CommandLineUtilities.hasArg("gaussian_kernel_sigma"))
    gaussian_kernel_sigma = CommandLineUtilities.getOptionValueAsFloat("gaussian_kernel_sigma");
```

The default value for  $d$  should be 2 and the default value for  $\sigma$  should be 1.

### 1.3.2 Learning Rate

Your default value for  $\eta$  should be 0.01. You *must* add a command line argument to allow this value to be adjusted via the command line.

Add this command line option by adding the following code to the `createCommandLineOptions` method of `Classify`.

```
registerOption("gradient_ascent_learning_rate", "double", true, "The learning rate for logistic regression.");
```

Be sure to add the option name exactly as it appears above. You can then read the value from the command line by adding the following to the main method of `Classify`:

```
double gradient_ascent_learning_rate = 0.01;
if (CommandLineUtilities.hasArg("gradient_ascent_learning_rate"))
    gradient_ascent_learning_rate =
        CommandLineUtilities.getOptionValueAsFloat("gradient_ascent_learning_rate");
```

### 1.3.3 Number of Training Iterations

We will define the number of times each algorithm iterates over all of the data by the parameter `gradient_ascent_training_iterations`. You *must* define a new command line option for this parameter. Use a default value of 5 for this parameter.

You can add this option by adding the following code to the `createCommandLineOptions` method of `Classify`.

```
registerOption("gradient_ascent_training_iterations", "int", true, "The number of training iterations.");
```

You can then read the value from the command line by adding the following to the main method of `Classify`:

```
int gradient_ascent_training_iterations = 5;
if (CommandLineUtilities.hasArg("gradient_ascent_training_iterations"))
    gradient_ascent_training_iterations =
        CommandLineUtilities.getOptionValueAsInt("gradient_ascent_training_iterations");
```

You should not change the order of examples. You must iterate over examples exactly as they appear in the data file, i.e. as provided by the data loader.

### 1.3.4 Initialization

You should initialize all  $\alpha$  values to 0.

### 1.3.5 Numeric Stability

For all numerical calculations involving floating point numbers, use the `double` type and NOT the `float` type to store values. This will help in achieving numerical precision.

## 1.4 Data Sets

We are providing a new synthetic data set for this assignment called `circle`. It may be helpful for testing your non-linear methods.

## 1.5 Appendix: Deriving the Gradient

In this section, we show how the gradient  $\nabla \ell(Y, X, \alpha)$  is derived. We begin by writing the conditional likelihood:

$$P(Y \mid \alpha, X) = \prod_{i=1}^N p(y_i \mid \alpha, \mathbf{x}_i) \quad (7)$$

Since  $y_i \in \{0, 1\}$ , we can write the conditional probability inside the product as:

$$P(Y \mid \alpha, X) = \prod_{i=1}^N p(y_i = 1 \mid \alpha, \mathbf{x}_i)^{y_i} \times (1 - p(y_i = 0 \mid \alpha, \mathbf{x}_i))^{1-y_i} \quad (8)$$

Note that one of these terms in the product will have an exponent of 0, and will evaluate to 1.

For ease of math and computation, we will take the log:

$$\ell(Y, X, \alpha) = \log P(Y \mid \alpha, X) = \sum_{i=1}^N y_i \log(p_{y_i=1}(\alpha, \mathbf{x}_i)) + (1 - y_i) \log(p_{y_i=0}(\alpha, \mathbf{x}_i)) \quad (9)$$

Plug in our logistic function  $g$  for the probability that  $y$  is 1:

$$\ell(Y, X, \alpha) = \sum_{i=1}^N y_i \log(g(\sum_{j=1}^N \alpha_j K(\mathbf{x}_j, \mathbf{x}_i))) + (1 - y_i) \log(1 - g(\sum_{j=1}^N \alpha_j K(\mathbf{x}_j, \mathbf{x}_i))) \quad (10)$$

To keep the notation clean, we will use the shorthand:

$$\kappa(\mathbf{x}_i) = \sum_{j=1}^N \alpha_j K(\mathbf{x}_j, \mathbf{x}_i) \quad (11)$$

With this notation, we have:

$$\ell(Y, X, \alpha) = \sum_{i=1}^N y_i \log(g(\kappa(\mathbf{x}_i))) + (1 - y_i) \log(1 - g(\kappa(\mathbf{x}_i))) \quad (12)$$

Recall that the link function,  $g$ , is the logistic function. It has the nice property  $1 - g(z) = g(-z)$ .

$$\ell(Y, X, \alpha) = \sum_{i=1}^N y_i \log(g(\kappa(\mathbf{x}_i))) + (1 - y_i) \log(g(-\kappa(\mathbf{x}_i))) \quad (13)$$

We can now use the chain rule to take the gradient with respect to  $\alpha$ :

$$\nabla \ell(Y, X, \alpha) = \sum_{i=1}^N y_i \frac{1}{g(\kappa(\mathbf{x}_i))} \nabla g(\kappa(\mathbf{x}_i)) + (1 - y_i) \frac{1}{g(-\kappa(\mathbf{x}_i))} \nabla g(-\kappa(\mathbf{x}_i)) \quad (14)$$

Since  $\frac{\partial}{\partial z}g(z) = g(z)(1 - g(z))$ :

$$\nabla \ell(Y, X, \alpha) = \sum_{i=1}^N y_i \frac{1}{g(\kappa(\mathbf{x}_i))} g(\kappa(\mathbf{x}_i))(1 - g(\kappa(\mathbf{x}_i))) \nabla(\kappa(\mathbf{x}_i)) \quad (15)$$

$$+ (1 - y_i) \frac{1}{g(-\kappa(\mathbf{x}_i))} g(-\kappa(\mathbf{x}_i))(1 - g(-\kappa(\mathbf{x}_i))) \nabla(-\kappa(\mathbf{x}_i)) \quad (16)$$

Simplify again using  $1 - g(z) = g(-z)$  and cancel terms

$$\nabla \ell(Y, X, \alpha) = \sum_{i=1}^N y_i g(-\kappa(\mathbf{x}_i)) \nabla(\kappa(\mathbf{x}_i)) + (1 - y_i) g(\kappa(\mathbf{x}_i)) \nabla(-\kappa(\mathbf{x}_i)) \quad (17)$$

You can now get the partial derivatives (components of the gradient) out of this gradient function by:

$$\frac{\partial \ell}{\partial \alpha_k} = \sum_{i=1}^N y_i g \left( - \sum_{j=1}^N \alpha_j K(x_j, x_i) \right) (K(x_i, x_k)) + (1 - y_i) g \left( \sum_{j=1}^N \alpha_j K(x_j, x_i) \right) (-K(x_i, x_k)) \quad (18)$$

because

$$\frac{\partial}{\partial \alpha_k} \kappa(\mathbf{x}_i) = \frac{\partial}{\partial \alpha_k} \left[ \sum_{j=1}^N \alpha_j K(\mathbf{x}_j, \mathbf{x}_i) \right] = K(\mathbf{x}_k, \mathbf{x}_i). \quad (19)$$

## 2 Analytical (50 points)

The following problems consider a standard binary classification setting: we are given  $n$  observations with  $m$  features,  $x_1, \dots, x_n \in \mathbb{R}^m$ .

**1) Overfitting (10 points)** SVMs using nonlinear kernels usually have two tuning parameters (regularization parameter  $C$  and kernel parameter  $\gamma$ ), which are usually determined by cross validation.

- (a) Suppose we use cross validation to determine the non-linear kernel parameter and slack variable for an SVM. We find that classifiers using parameters  $(c_1, \gamma_1)$  and  $(c_2, \gamma_2)$  achieve the same cross validation error, but  $(c_1, \gamma_1)$  leads to fewer support vectors than  $(c_2, \gamma_2)$ . Explain which set of parameters should we choose for the final model?
- (b) The optimization problem of linear SVMs can be in either primal or dual form. As we know, the primal form has  $m$  parameters to learn, while the dual form has  $n$  parameters to learn. If  $m \gg n$ , is it true that the dual form reduces over-fitting since it has fewer parameters? Explain.

**2) Hinge Loss (10 points)** Linear SVMs can be formulated in an unconstrained optimization problem

$$\min_{w, b} \sum_{i=1}^n H(y_i(w^T x_i)) + \lambda \|w\|_2^2, \quad (20)$$

where  $\lambda$  is the regularization parameter and  $H(a) = \max(1 - a, 0)$  is the well known hinge loss function. The hinge loss function can be viewed as a convex surrogate of the 0/1 loss function  $I(a \leq 0)$ .

- Prove that  $H(a)$  is a convex function of  $a$ .
- The function  $L(a) = \max(-a, 0)$  can also approximate the 0/1 loss function. What is the disadvantage of using this function instead?
- If  $H'(a) = \max(0.5 - a, 0)$ , please show that there exists  $\lambda'$  such that (21) is equivalent to (20)

$$\min_{w,b} \sum_{i=1}^n H'(y_i(w^T x_i)) + \lambda' \|w\|_2^2. \quad (21)$$

**3) Kernel Trick (12 points)** The kernel trick extends SVMs to handle with nonlinear data sets. However, an improper use of a kernel function can cause serious over-fitting. Consider the following kernels.

- Polynomial kernel:  $K(x, x') = (1 + (xx'^T))^d$ , where  $d \in \mathbb{N}$ . Does increasing  $d$  make over-fitting more or less likely?
- Gaussian kernel:  $K(x, x') = \exp(-\|x - x'\|^2 / 2\sigma^2)$ , where  $\sigma > 0$ . Does increasing  $\sigma$  make over-fitting more or less likely?

We say  $K$  is a kernel function, if there exists some transformation  $\phi : \mathbb{R}^m \rightarrow \mathbb{R}^{m'}$  such that  $K(x_i, x_{i'}) = \langle \phi(x_i), \phi(x_{i'}) \rangle$ .

- Let  $K_1$  and  $K_2$  be two kernel functions. Prove that  $K(x_i, x_{i'}) = K_1(x_i, x_{i'}) + K_2(x_i, x_{i'})$  is also a kernel function.

**4) Prediction using Kernel (8 points)** One of the differences between linear SVMs and kernel SVMs concerns computational complexity at prediction time.

- What is the computational complexity of prediction of a linear SVM in terms of the examples  $n$  and features  $m$ ?
- What is the computational complexity of prediction of a non-linear SVM in terms of the examples  $n$ , features  $m$  and support vectors  $s$ ?

**5) Stochastic Gradient Algorithm (10 points)** The stochastic gradient algorithm is a very powerful optimization tool to solve many online learning problems. Instead of computing the gradient over the entire data set before making an update, the stochastic gradient algorithm computes the gradient over a single example, then updates the parameters. By passing over the entire data set in this fashion we can converge to the optimal parameters.

A single iteration of stochastic gradient considers a single example. As a result, we know that stochastic gradient converges more slowly than gradient descent. While it takes many more iterations, each iteration is much faster, both in terms of memory and computation.

Consider a ridge regression problem with  $n$  samples:

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_2^2. \quad (22)$$

In each iteration, instead of using only one example, we randomly choose  $k$  out of  $n$  samples and obtain  $(x_{1'}, y_{1'}), \dots, (x_{k'}, y_{k'})$ .

- (a) What is the computational complexity of computing the gradient at each iteration (using  $k$  and  $n$ )?
- (b) What are the advantages/disadvantages of increasing  $k$  in terms of computational complexity (using  $k$  and  $n$ )? What is traded-off by increasing/decreasing  $k$ ?
- (c) What is an example of an algorithm we learned in class that uses stochastic gradient? Explain.

### 3 What to Submit

In each assignment you will submit two things.

1. **Code:** Your code as a zip file named `library.zip`. **You must submit source code (.java files)**. We will run your code using the exact command lines described above, so make sure it works ahead of time. Remember to submit all of the source code, including what we have provided to you.
2. **Writeup:** Your writeup as a **PDF file** (compiled from latex) containing answers to the analytical questions asked in the assignment. Make sure to include your name in the writeup PDF and use the provided latex template for your answers.

Make sure you name each of the files exactly as specified (`library.zip` and `writeup.pdf`).

To submit your assignment, visit the “Homework” section of the website (<http://www.cs475.org/>.)

### 4 Questions?

Remember to submit questions about the assignment to the appropriate group on the class discussion board: <http://bb.cs475.org>.