

# CPSC 340 Assignment 4 (due November 14)

## Linear Models Part 2

### 1 Logistic Regression with Sparse Regularization

If you run the function *example\_logistic*, it will:

1. Load a binary classification dataset containing a training and a validation set.
2. ‘Standardize’ the columns of  $X$  and add a bias variable.
3. Apply the same transformation to  $X_{\text{validate}}$ .
4. Fit a logistic regression model.
5. Report the number of features selected by the model (number of non-zero regression weights).
6. Report the error on the validation set.

Logistic regression does ok on this dataset, but it uses all the features (even though only the prime-numbered features are relevant) and the validation error is above the minimum achievable for this model (which is 1 percent). In this question, you will modify this demo to use different forms of regularization to improve on these aspects.

#### 1.1 L2-Regularization

Make a new function, *logRegL2*, that takes an input parameter  $\lambda$  and fits a logistic regression model with L2-regularization. Specifically, while *logReg* computes  $w$  by minimizing

$$f(w) = \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)),$$

your new function *logRegL2* should compute  $w$  by minimizing

$$f(w) = \sum_{i=1}^n [\log(1 + \exp(-y_i w^T x_i))] + \frac{\lambda}{2} \|w\|^2.$$

Hand in your updated code. Using this new code, report the number of non-zeroes and the validation error with  $\lambda = 1$ .

#### 1.2 L1-Regularization

Make a new function, *logRegL1*, that takes an input parameter  $\lambda$  and fits a logistic regression model with L1-regularization,

$$f(w) = \sum_{i=1}^n [\log(1 + \exp(-y_i w^T x_i))] + \lambda \|w\|_1.$$

Hand in your updated code. Using this new code, report the number of non-zeroes and the validation error with  $\lambda = 1$ .

Hint: the function `findMinL1` implements a proximal-gradient method to minimize the sum of a differentiable function  $g$  and  $\lambda\|w\|_1$ ,

$$f(w) = g(w) + \lambda\|w\|_1.$$

This function has a similar interface to `findMin`, except that you (a) only provide the code to compute the function/gradient of the differentiable part  $g$  and (b) need to provide the value  $\lambda$ .

### 1.3 L0-Regularization

The function `logRegL0` contains part of the code needed to implement the *forward selection* algorithm, which approximates the solution with L0-regularization,

$$f(w) = \sum_{i=1}^n [\log(1 + \exp(-y_i w^T x_i))] + \lambda\|w\|_0.$$

The ‘for’ loop in this function is missing the part where we fit the model using the subset `ind_new`, then compute the score and updates the `minScore/minInd`. Modify the ‘for’ loop in this code so that it fits the model using only the features `ind_new`, computes the score above using these features, and updates the `minScore/minInd` variables. Hand in your updated code. Using this new code, report the number of non-zeroes and the validation error with  $\lambda = 1$ .

Note that the code differs a bit from what we discussed in class, since we assume that the first feature is the bias variable and assume that the bias variable is always included. Also, note that for this particular case using the L0-norm with  $\lambda$  is equivalent to what is known as the Bayesian information criterion (BIC) for variable selection.

## 2 Convex Functions and MLE/MAP Loss Functions

This question gets you to explore two important concepts related to loss functions: the *convexity* of loss functions (since convex loss functions can be minimized with gradient descent) and the *probabilistic interpretation* of loss functions (since this allows us to define new loss functions when we encounter weird new situations).

### 2.1 Showing Convexity from Definitions

Show that the following functions are convex:

- |   |  |                           |
|---|--|---------------------------|
| 1. Quadratic                                | $f(w) = aw^2 + bw$   | $w \in \mathbb{R}, a > 0$ |
| 2. Negative logarithm                       | $f(w) = -\log(aw)$   | $w > 0$                   |
| 3. Regularized regression (arbitrary norms) | $f(w) = \ Xw - y\ _p + \lambda\ w\ _q$   | $p \geq 1, q \geq 1$      |
| 4. Logistic regression                      | $f(w) = \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$                                       | $w \in \mathbb{R}^d$      |
| 5. Support vector regression                | $f(w) = \sum_{i=1}^N \max\{0,  w^T x_i - y_i  - \epsilon\} + \frac{\lambda}{2}\ w\ _2^2$ |                           |

Hint: for the first two you can use the second-derivative test. For the last 3 you’ll have to use some of the results in class regarding how combining convex functions can yield convex functions (see Lecture 17).

## 2.2 MAP Estimation

In class, we considered MAP estimation in a regression model where we assumed that:

- The likelihood  $p(y_i|x_i, w)$  is a normal distribution with a mean of  $w^T x_i$  and a variance of 1.
- The prior for each variable  $j$ ,  $p(w_j)$ , is a normal distribution with a mean of zero and a variance of  $\lambda^{-1}$ .

Under these assumptions, we showed that this leads to the standard L2-regularized least squares objective function:

$$f(w) = \frac{1}{2} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|^2.$$

For each of the alternate assumptions below, show how the loss function would change (simplifying as much as possible):

1. We use a zero-mean Laplace prior for each variable with a scale parameter of  $\lambda^{-1}$ , so that

$$p(w_j) = \frac{\lambda}{2} \exp(-\lambda|w_j|).$$

2. We use a Laplace likelihood with a mean of  $w^T x_i$  and a scale of 1, so that

$$p(y_i|x_i, w) = \frac{1}{2} \exp(-|w^T x_i - y_i|).$$

3. We use a Gaussian likelihood where each datapoint where the variance is  $\sigma^2$  instead of 1,

$$p(y_i|x_i, w) = \frac{1}{\sqrt{2\sigma^2\pi}} \exp\left(-\frac{(w^T x_i - y_i)^2}{2\sigma^2}\right).$$

4. We use a Gaussian likelihood where each datapoint has its own variance  $\sigma_i^2$ ,

$$p(y_i|x_i, w) = \frac{1}{\sqrt{2\sigma_i^2\pi}} \exp\left(-\frac{(w^T x_i - y_i)^2}{2\sigma_i^2}\right).$$

5. We use a (very robust) student  $t$  likelihood with a mean of  $w^T x_i$  and a degree of freedom of  $\nu$ ,

$$p(y_i|x_i, w) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi}\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{(w^T x_i - y_i)^2}{\nu}\right)^{-\frac{\nu+1}{2}},$$

where  $\Gamma$  is the “gamma” function (which is always non-negative).

Why is loss coming from the student  $t$  distribution “very robust”?

## 3 Multi-Class Logistic

The function `example_multiClass` loads a multi-class classification dataset and fits a ‘one-vs-all’ classification model using least squares, then reports the validation error and shows a plot of the data/classifier. The performance on the validation set is ok, but could be much better. For example, this classifier never even predicts that examples will be in classes 1 or 5.

### 3.1 One-vs-all Logistic Regression

Using the squared error on this problem hurts performance because it has ‘bad errors’ (the model gets penalized if it classifies examples ‘too correctly’). Write a new function, *logLinearClassifier*, that replaces the squared loss in the one-vs-all model with the logistic loss. [Hand in the code and report the validation error.](#)

### 3.2 Softmax Classification

Using a one-vs-all classifier hurts performance because the classifiers are fit independently, so there is no attempt to calibrate the columns of the matrix  $W$ . An alternative to this independent model is to use the softmax probability,

$$p(y_i|W, x_i) = \frac{\exp(w_{y_i}^T x_i)}{\sum_{c=1}^k \exp(w_c^T x_i)}.$$

Here  $c$  is a possible label and  $w_c$  is column  $c$  of  $W$ . Similarly,  $y_i$  is the training label,  $w_{y_i}$  is column  $y_i$  of  $W$ , and in this setting we are assuming a discrete label  $y_i \in \{1, 2, 3, 4, 5\}$ . Before we move on to implementing the softmax classifier, let’s do a simple example:

Consider the dataset below, which has 10 training examples and 2 features:

$$X = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 0 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 2 \\ 2 \\ 3 \\ 3 \\ 3 \\ 3 \end{bmatrix}.$$

Suppose that you want to classify the following test example:

$$\hat{x} = \begin{bmatrix} 1 & 1 \end{bmatrix}.$$

Suppose we fit a multi-class linear classifier using the softmax loss, and we obtain the following weight matrix:

$$W = \begin{bmatrix} +2 & +2 & +3 \\ -1 & +2 & -1 \end{bmatrix}$$

[Under this model, what class label would we assign to the test example? \(Show your work.\)](#)

### 3.3 Softmax Loss

The loss function corresponding to the negative logarithm of the softmax probability is given by

$$f(W) = \sum_{i=1}^n \left[ -w_{y_i}^T x_i + \log \left( \sum_{c'=1}^k \exp(w_{c'}^T x_i) \right) \right].$$

Derive the derivative of this loss function with respect to a particular element  $W_{jc}$ . Try to simplify the derivative as much as possible (but you can express the result in summation notation).

Hint: for the gradient you can use  $x_{ij}$  to refer to element  $j$  of example  $i$ . You can use an ‘indicator’ function,  $I(y_i = c)$ , which is 1 when  $y_i = c$  and is 0 otherwise. Note that you can use the definition of the softmax probability to simplify the derivative.

### 3.4 Softmax Classifier

Make a new function, *softmaxClassifier*, which fits  $W$  using the softmax loss from the previous section instead of fitting  $k$  independent classifiers. [Hand in the code and report the validation error.](#)

Hint: you may want to use the *autoGrad* function from A3 to check that your gradient code is correct.