

CS 273A: Machine Learning: Fall 2018

Homework 3

Due Date: **Thursday, November 1, 2018**

Sameer Singh

<http://sameersingh.org/courses/gml/fa18/>

The submission for this homework, as for others so far, should be **one stand-alone PDF file** containing all of the relevant code, figures, and any text explaining your results. Although you will be primarily filling missing sections in a Python file, include all the relevant sections you have written as answers to the appropriate question.

Points: This homework adds up to a total of **100 points** with 10 points of extra credit, as follows:

Problem 1: Logistic Regression 75 (+10) points

Problem 2: Shattering 20 points

Statement of Collaboration 5 points

Be sure to re-download the data and replace the Python code; even though it is mostly the same, we may make changes to the data or fix bugs for each homework.

Problem 1: Logistic Regression (75+10 points)

In this problem, we'll build a logistic regression classifier and train it on separable and non-separable data. Since it will be specialized to binary classification, we've named the class `logisticClassify2`.

We'll start by building two binary classification problems, one separable and the other not:

```
1 iris = np.genfromtxt("data/iris.txt", delimiter=None)
2 X, Y = iris[:,0:2], iris[:, -1] # get first two features & target
3 X, Y = ml.shuffleData(X, Y)    # reorder randomly (important later)
4 X, _ = rescale(X)              # works much better on rescaled data
5
6 XA, YA = X[Y<2, :], Y[Y<2]    # get class 0 vs 1
7 XB, YB = X[Y>0, :], Y[Y>0]    # get class 1 vs 2
```

For this problem, we are focused on the learning algorithm, rather than performance — so, we will not bother creating training and validation splits; just use all your data for training.

Note: The code uses numpy's `permute` to iterate over data randomly; should avoid issues due to the default order of the data (by class). Similarly, rescaling and centering the data may help speed up convergence as well.

1. Show the two classes in a scatter plot (one for each data set) and verify that one data set is linearly separable while the other is not. (5 points)
2. Write (fill in) the function `plotBoundary` in `logisticClassify2.py` to compute the points on the decision boundary. In particular, you only need to make sure `x2b` is set correctly using `self.theta`. This will plot the data & boundary quickly, which is useful for visualizing the model during training. To demo your function plot the decision boundary corresponding to the classifier

$$\text{sign}(.5 - .25x_1 + 1x_2)$$

along with the A data, and again with the B data; these fixed parameters will look like an OK classifier on one data set, but a poor classifier on the other.

You can create a “blank” learner and set the weights by:

```
1 import mltools as ml
2 from logisticClassify2 import *
3
4 learner = logisticClassify2(); # create "blank" learner
5 learner.classes = np.unique(YA) # define class labels using YA or YB
```

```

6  wts = np.array([theta0,theta1,theta2]); # TODO: fill in values
7  learner.theta = wts;                  # set the learner's parameters

```

Include the lines of code you added to the function, and the two generated plots. (10 points)

- Complete the `logisticClassify2.predict` function to make predictions for your classifier. Verify that your function works by computing & reporting the error rate of the classifier in the previous part on both data sets A and B. (The error rate on one should be ≈ 0.0505 , and higher on the other.) Note that, in the code, the two classes are stored in the variable `self.classes`, with the first entry being the “negative” class (or class 0), and the second entry being the “positive” class, so you want to have different learner objects for each dataset, and you use `learner.err` directly.

Include the function definition and the two computed errors. (10 points)

- Verify that your predict code matches your boundary plot by using `plotClassify2D` with your manually constructed learner on the two data sets. This will call `predict` on a dense grid of points, and you should find that the resulting decision boundary matches the one you computed analytically. (5 points)
- In the provided code, we first transform the classes in the data Y into YY , with canonical labels for the two classes: “class 0” (negative) and “class 1” (positive). In our notation, let $r^{(j)} = x^{(j)} \cdot \theta^T$ be the linear response of the classifier, and σ is the standard logistic function

$$\sigma(r) = (1 + \exp(-r))^{-1}.$$

The logistic negative log likelihood loss for a single data point j is then

$$J_j(\theta) = -y^{(j)} \log \sigma(x^{(j)} \theta^T) - (1 - y^{(j)}) \log(1 - \sigma(x^{(j)} \theta^T))$$

where $y^{(j)}$ is either 0 or 1. Derive the gradient of the negative log likelihood J_j for logistic regression, and give it in your report. (You will need this in your gradient descent code for the next part.)

Provide the gradient equations for $\frac{\partial}{\partial \theta_0} J_j$, $\frac{\partial}{\partial \theta_1} J_j$, and $\frac{\partial}{\partial \theta_2} J_j$. (10 points)

- Complete `train` function to perform stochastic gradient descent on the logistic loss function. This will require that you fill in:
 - computing the surrogate loss function at each epoch ($J = \frac{1}{m} \sum J_j$, from the previous part);
 - computing the response ($r^{(j)}$) and gradient associated with each data point $x^{(j)}, y^{(j)}$;
 - a stopping criterion consisting of two conditions (stop when either you have reached `stopEpochs` epochs or J has not changed by more than `stopTol` since the last epoch).

Include the complete implementation of `train`. (25 points)

- Run `train` for your logistic regression classifier on both data sets (A and B). Describe your parameter choices for each dataset (stepsize, etc.) and include plots showing the convergence of the surrogate loss and error rate (e.g., the loss values as a function of epoch during gradient descent), and the final converged classifier with the data (the included `train` function does that for you already). (10 points)

Note: Debugging machine learning algorithms can be quite challenging, since the results of the algorithm are highly data-dependent, and often somewhat randomized (initialization, etc.). I suggest starting with an extremely small step size and verifying both that the learner’s prediction evolves slowly in the correct direction, and that the objective function J decreases monotonically. If that works, go to larger step sizes to observe the behavior. I often manually step through the code — for example by pausing after each parameter update using `raw_input()` (Python 2.7) or `input()` (Python 3) — so that I can examine its behavior. You can also (of course) use a more sophisticated debugger.

Note on plotting: The code generates plots as the algorithm runs, so you can see its behavior over time; this is done with `pyplot.draw()`. Run your code either interactively or as a script to see these display over time; refer to discussion notebooks on how to plot the loss over time in Jupyter.

- Extra Credit (10 points):** Add an L2 regularization term ($+\alpha \sum_i \theta_i^2$) to your surrogate loss function, and update the gradient and your code to reflect this addition. Try re-running your learner with some regularization (e.g. $\alpha = 2$) and see how different the resulting parameters are. Find a value of α that gives noticeably different results & explain them.

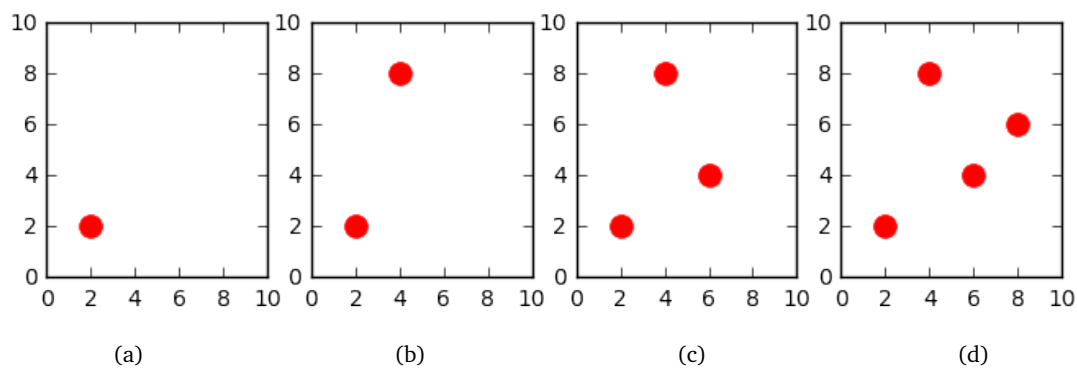


Figure 1: Four datasets to test whether they can be *shattered* by a given classifier, i.e. can the classifier exactly separate their all possible binary colorings. **No three data points are on a line.**

Problem 2: Shattering and VC Dimension (20 points)

Consider the data points in Figure 1 which have two real-valued features x_1, x_2 . We are also giving a few learners below. For the learners below, $T[z]$ is the sign threshold function, $T[z] = +1$ for $z \geq 0$ and $T[z] = -1$ for $z < 0$. The learner parameters a, b, c, \dots are real-valued scalars, and each data point has two real-valued features x_1, x_2 .

Which of the four datasets can be shattered by each learner? Give a brief explanation/justification and use your results to guess the VC dimension of the classifier (you do not have to give a formal proof, just your reasoning).

1. $T(a + bx_1)$ (5 points)
2. $T((a * b)x_1 + (c/a)x_2)$ (5 points)
3. $T((x_1 - a)^2 + (x_2 - b)^2 + c)$ (5 points)
4. $T(a + bx_1 + cx_2) \times T(d + bx_1 + cx_2)$ (5 points)

Hint: The two equations are two parallel lines.

Statement of Collaboration (5 points)

It is **mandatory** to include a *Statement of Collaboration* in each submission, with respect to the guidelines below. Include the names of everyone involved in the discussions (especially in-person ones), and what was discussed.

All students are required to follow the academic honesty guidelines posted on the course website. For programming assignments, in particular, I encourage the students to organize (perhaps using Campuswire) to discuss the task descriptions, requirements, bugs in my code, and the relevant technical content *before* they start working on it. However, you should not discuss the specific solutions, and, as a guiding principle, you are not allowed to take anything written or drawn away from these discussions (i.e. no photographs of the blackboard, written notes, referring to Campuswire, etc.). Especially *after* you have started working on the assignment, try to restrict the discussion to Campuswire as much as possible, so that there is no doubt as to the extent of your collaboration.

Acknowledgements

This homework is adapted (with minor changes) from one made available by Alex Ihler's machine learning course. Thanks, Alex!