

6.867: Homework 2

Anonymous authors

Abstract—This is Machine Learning homework 2. Topics include Logistic Regression, SVM. All work has been done in Python.

I. LOGISTIC REGRESSION (LR)

In this section, we consider Logistic Regression (LR) for binary classification of several 2D datasets, and compare the effect of different model parameters and regularization.

A. Part 1

[todo]

B. Part 2

An alternative regularization (L_1) was then compared against L_2 from the previous section. In Fig. 1, the elements of the weight vector are compared. The top row shows L_1 regularization for each of the four datasets, and the bottom row shows L_2 . In general, increased regularization parameter, λ , causes smaller magnitude weight vector, since the blue and green points (small λ) are further from zero than the magenta and yellow points (large λ). This is true for both L_1 and L_2 regularization. Especially clear in the 2nd and 4th columns, L_1 regularization creates a more sparse weight vector, as all elements are zero (yellow) for the highly regularized case, whereas with the same regularization constant with a L_2 penalty has some non-zero elements.

In addition to the weight vector, the decision boundary is affected by λ and regularization method. In Fig. 2, the decision boundaries for various λ values are shown for same dataset with L_1 on the left and L_2 on the right. As λ increases, the boundary is placed in a position that will cause training errors ($\lambda > 1$). The boundaries from the two regularization strategies are slightly different, but not by very much. In Fig. 3, the same concept is shown across all four datasets, with the top row using L_1 and bottom row using L_2 . Green, blue, and yellow lines are decision boundaries with low λ , and red magenta, cyan are high λ . Again, as λ increases, training error increases. The first dataset (left) is linearly separable, the middle two have some class overlap, and the right dataset is not even close. Accordingly, the decision boundary does not appear in the range of the data on the rightmost plot.

Finally, the test-set classification error is also compared across λ values and regularization methods in Fig. 4. When λ is very high for L_1 regularization (blue), test accuracy suffers, since model is too sparse. L_2 reg. has almost constant accuracy throughout.

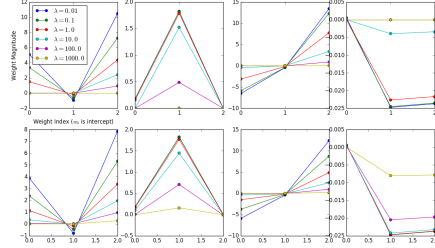


Fig. 1: Weight vectors of L_1 regularization (top row) and L_2 reg. (bottom row) for four datasets. Each column is the model for the same dataset. Within each subplot, several values of λ are shown. For high lambda, weights are smaller in magnitude for all plots. L_1 reg causes sparser weight vector (more zero elements) than L_2 .

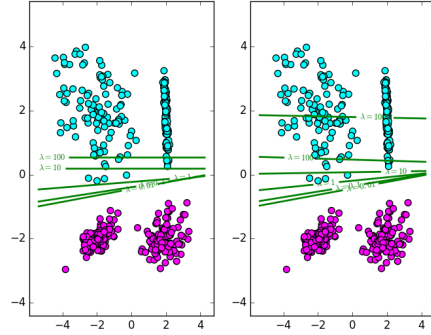


Fig. 2: LR on one dataset with L_1 reg on left, L_2 on right. Green lines show decision boundaries for various λ s. As λ increases, training error increases.

C. Part 3

To pick the best regularizer and λ , the LR weights are trained with many values of λ for each regularizer, using the training set. Then, the generalization is evaluated on the validation set, by measuring the classification accuracy. A model that generalizes well will have high accuracy on data not seen during training. In cases where the validation accuracy is identical for multiple values of λ , the higher λ is chosen, because this corresponds to lower model complexity (more regularized).

Dataset	Regularizer	λ	Train Acc. (%)	Val. Acc. (%)
1	L_2	10	100.0	100.0
2	L_2	10	100.0	100.0
3	L_1	10	100.0	100.0
4	L_2	10	100.0	100.0

TABLE I: Accuracy of LR on four datasets.

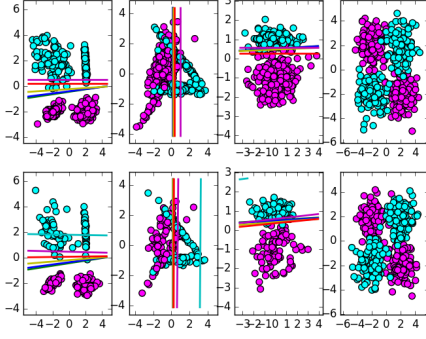


Fig. 3: LR on four datasets, with the top row using L_1 and bottom row using L_2 . Green, blue, and yellow lines are decision boundaries with low λ , and red magenta, cyan are high λ . As λ increases, training error increases. The first dataset (left) is linearly separable, the middle two have some class overlap, and the right dataset is not even close. Accordingly, the decision boundary does not appear in the range of the data on the rightmost plot.

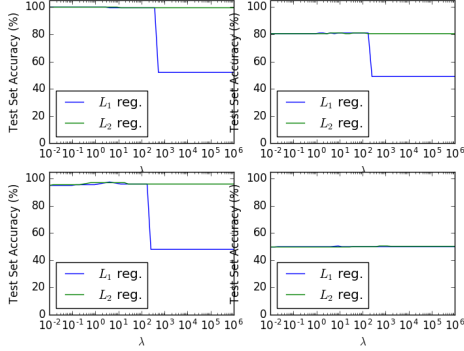


Fig. 4: For each dataset, the test accuracy is plotted across a wide range of λ values. When λ is very high for L_1 regularization (blue), test accuracy suffers, since model is too sparse. L_2 reg has almost constant accuracy throughout.

II. SUPPORT VECTOR MACHINE (SVM)

In this section, we consider the Support Vector Machine (SVM) for binary classification of several 2D datasets.

A. Part 1

First, we implemented the dual form of a linear SVM with slack variables.

[dual svm eqns]

I first converted the input data into a standard format, then we used the cvxopt Python package to execute the quadratic programming to find the desired α values.

[data format P, etc.]

With the simple 4 element dataset $\{((2, 2), +1), ((2, 3), +1), ((0, -1), -1), ((-3, -2), -1)\}$, the algorithm outputs $(2, 2)$ and $(0, -1)$ as the support vectors, meaning they lie on the boundaries of the margin, seen in Fig. 5. The weight vector has elements $w = [0.308, 0.462]$ and bias, $b = -0.538$. The training error

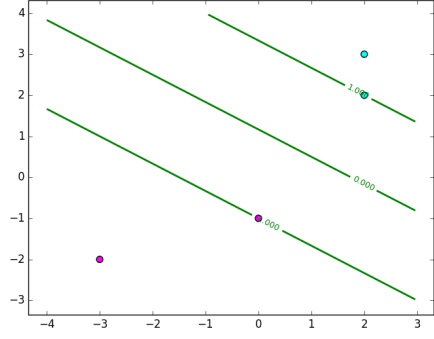


Fig. 5: Simple 4-element dataset (positive elements in cyan, negative in magenta) is separated by SVM. The three lines are the decision boundary ($\hat{y} = 0.0$) and positive and negative margin boundaries ($\hat{y} = \pm 1.0$) and two elements (support vectors) lie directly on boundaries of margin. Training error is zero.

is zero for this simple dataset. We use regularization term $C = 1$ for all results in this and the next subsection.

B. Part 2

Next, we used SVM to classify data from four much larger 2D datasets. The weights and bias of the decision boundary are generated using the training data, shown on the top row of ??, and then tested on the validation set shown in the bottom row of the same figure. Each column represents one dataset, presumably from the same distribution but with slightly different data. The leftmost dataset is easiest to linearly separate, whereas the middle two datasets have significant overlap between classes in this feature space. The rightmost dataset is not well-suited for a linear classifier because the data seems to have four distinct clusters at opposite ends of a rectangle.

The classification accuracy on the validation set matches the intuition from the plots. The first and third datasets are classified most accurately, while the fourth dataset is essentially equivalent to a random coinflip. For three of the four datasets, the training accuracy is at least as high as the validation accuracy. This is expected, because the model typically performs worse on data that was not seen during testing. In these cases, though, the accuracy difference is very minor.

Dataset	Training Accuracy (%)	Training Accuracy (%)
1	100.0	100.0
2	82.75	84.0
3	98.75	96.5
4	51.0	48.5

TABLE II: Accuracy of linear SVM on four datasets in ??.

C. Part 3

[todo]

