# Machine Learning: Assignment #2
## Winter 2018

**Due:** Dec 20th, 23:59:59 CST (UTC +8).

## 1. A Walk Through Linear Models

In this problem, you will implement a whole bunch of linear classifiers and compare their performance and properties.

We assume our target function is $f : [-1, 1] \times [-1, 1] \rightarrow \{-1, +1\}$. In each run, choose a random line in the plane as your target function $f$, where one side of the line maps to $+1$ and the other maps to -1.

Skeleton code and MATLAB/Python functions including *mkdata* and *plotdata* are given for your convenience, see the comments in the code for more details. What you need to do is implementing each algorithm and write scripts to play with them. Your algorithm implementations should be able to handle data in arbitrary dimension. For most of the questions below, you should repeat the experiments for 1000 times and take the average. See *run.m/run.ipynb* for a script example.

Hint: Do not forget to add the bias constant term!

(a) **Perceptron**

   Implement Perceptron learning algorithm (in *perceptron.m/perceptron.py*), then answer the following questions.

   (i) What is the training error rate and expected testing error rate (we can generate a large set of testing points to estimate test error), when the size of training set is 10 and 100 respectly?

   (ii) What is the average number of iterations before your algorithm converges when the size of training set is 10 and 100 respectively ?

   (iii) What will happen if the training data is not linearly separable (Use *mkdata(N, 'noisy')* to generate non-linearly separable data) ?

(b) **Linear Regression**

   Implement Linear Regression (in *linear_regression.m/linear_regression.py*), then answer the following questions.

   Note that we use Linear Regression here to classify data points. This technique is called Linear Regression on indicator response matrix[1].

   (i) What is the training error rate and expected testing error rate if the size of training set is 100 ?

   (ii) What is the training error rate and expected testing error rate if the training data is noisy and not linearly separable (nTrain = 100) ?

---

[1] `http://books.google.com/books?id=VRzITwgNV2UC&pg=PA81`

   (iii) Run Linear Regression on dataset *poly_train.mat*. What is the training error rate? What is the testing error rate on *poly_test.mat* ?

   (iv) Do the following data transformation

$$(1, x_1, x_2) \rightarrow (1, x_1, x_2, x_1 x_2, x_1^2, x_2^2).$$

on dataset *poly_train.mat*. After this transformation, what is the training error rate? What is the testing error rate on *poly_test.mat* ?

(c) **Logistic Regression**

Implement Logistic Regression (in *logistic.m/logistic.py*), then answer the following questions.

Remember the form of Logistic Regression

$$p(y|\mathbf{x}) = \frac{1}{1 + e^{-\theta^T x}}.$$

Let us assume

$$P(y = 1|x; \theta) = h_\theta(x), \quad P(y = 0|x; \theta) = 1 - h_\theta(x),$$

then

$$P(y|x; \theta) = (h_\theta(x))^y (1 - h_\theta(x))^{1-y}.$$

Now you can learn your logistic model using MLE. Not like Linear Regression, we have no closed-form solution to solve the MLE problem. So you should use gradient descent to maximize the log-likelihood function iteratively.

Derive the gradient descent update rule, choose a proper learning rate, and repeat the experiment for 100 times to take average.

   (i) What is the training error rate and expected testing error rate if the size of training set is 100 ?

   (ii) What is the training error rate and expected testing error rate if the training data is noisy and not linearly separable (nTrain = 100) ?

(d) **Support Vector Machine**

Implement Support Vector Machine (in *svm.m/svm.py*) without introducing slack variables , then answer the following questions.

Hint: For MATLAB users, using MATALB built-in function *quadprog*, you should be able to implement SVM in less than 10 lines of code. For Python users, using *scipy.optimize*, you should be able to implement in 20 lines of code.

   (i) What is the training error rate and expected testing error rate if the size of training set is 30 ? (Use *plotdata* to plot your learnt model, hope the graph will help you understand SVM as a maximum margin classifier.)

   (ii) What is the training error rate and expected testing error rate if the size of training set is 100 ?

   (iii) For the case nTrain = 100, what is average number of support vectors in your trained SVM models?

## 2. Regularization and Cross-Validation

We now learnt the undesirable consequence of overfitting, just as this quotation says:

> If you torture the data for long enough, it will confess to anything.

Regularization and Cross-Validation[2] are the two most common used methods to control overfitting. In this problem, we will use these two techniques together to relieve the overfitting problem.

Cross-Validation is a standard technique for adjusting regularization parameters of learning models. We first split the dataset into training set and validation set. Then we compute the validation error on validation set using different values for regularization parameters. Finally, we pick the parameter with the lowest validation error and use it for training our learning model. To reduce variability, we can run validation multiple rounds using different dataset partitions, and the validation results are averaged over all the rounds.

In the following questions, we will use leave-one-out cross-validation to choose regularization parameters. As the name suggests, leave-one-out cross-validation (LOOCV) involves using a single observation from the original sample as the validation data, and the remaining observations as the training data. This is repeated such that each observation in the sample is used once as the validation data.

You may use *validation.m/validation.ipynb* as a skeleton.

(a) Implement Ridge Regression (in *ridge.m/ridge.py*), and use LOOCV to tune the regularization parameter $\lambda$.
Train your algorithm on *digit_train.mat* and test on *digit_test.mat*.

Hint: Since the dataset is a set of raw images (hand-written digits) without preprocessing, you should do feature normalization[3] to make the feature have zero-mean and unit-variance. You can visualize the dataset using *show_digit.m*.

  (i) What is the $\lambda$ chosen by LOOCV ?

  (ii) What is $\sum_{i=1}^{m} \omega_i^2$ with and without regularization (let $\lambda = 0$) ?

  (iii) What's training/testing error rate with and without regularization ?

(b) Implement Logistic Regression with regularization (in *logistic_r.m/logistic_r.py*), and use LOOCV to tune the regularization parameter $\lambda$.

Train your algorithm on *digit_train.mat* and test on *digit_test.mat*. Then report the training/testing error rate with and without regularization (let $\lambda = 0$). You should also report the $\lambda$ chosen by LOOCV.

---

[2]https://en.wikipedia.org/wiki/Cross-validation_(statistics)
[3]https://en.wikipedia.org/wiki/Feature_scaling

## 3.   Bias Variance Trade-off

Let us review the bias-variance decomposition first.

   The intuition behind it is straight-forward: if the model is too simple, the learnt function is biased and does not fit the data. If the model is too complex then it is very sensitive to small changes in the data.

   If we were able to sample a dataset $D$ infinite many times, we will learn different $g(x)$ for each time, and get an expected hypothesis $\bar{g}(x)$. So bias means the difference between the truth and what you expect to learn. It measures how well our model can approximate the truth at best.

   However, it is impossible to sample the training dataset multiple time, so variance means the difference between what you learn from a particular dataset and what you expect to learn.

   Now please answer the following questions:

(a) True of False:

   (i)   If a learning algorithm is suffering from high bias, adding more training examples will improve the test error significantly.

   (ii)  We always prefer models with high variance (over those with high bias) as they will able to better fit the training set.

   (iii) A model with more parameters is more prone to overfitting and typically has higher variance.

   (iv)  Introducing regularization to the model always results in equal or better performance on the training set.

   (v)   Using a very large value of regularization parameter $\lambda$ cannot hurt the performance of your hypothesis.

---

   Please submit your homework report to at `http://assignment.zjulearning.org:8081` in pdf format, with all your code in a zip archive.