# EE 240: Pattern Recognition and Machine Learning
# Homework 2

**Due date:** May 3, 2018

**Description:** These questions will explore multiple topics in logistic regression, perceptron learning, and support vector machines.

**Reading assignment:** ESL: Ch. 2, 3, 4, & 12, AML: Ch. 3 & 8

---

**Homework and lab assignment submission policy:**

All homework and lab assignments must be submitted online via https://iLearn.ucr.edu.

Homework solutions should be written and submitted individually, but discussions among students are encouraged.

All assignments should be submitted by the due date. There will be 25% penalty per day for late assignments. No grade will be given to homework submitted 3 days after the due date.

---

**H2.1** Logistic regression: Suppose we have labeled data $(\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathbb{R}^d$, $y \in \{0, 1\}$. We want to compute $\mathbf{w}$ that minimizes the following log-likelihood function:

$$\operatorname*{minimize}_{\mathbf{w}} \sum_{i=1}^{N} y_i \log(1 + \exp(-\mathbf{w}^T \mathbf{x}_i)) + (1 - y_i) \log(1 + \exp(\mathbf{w}^T \mathbf{x}_i)). \tag{1}$$

Let us denote the objective of this function as $J(\mathbf{w})$. (Note that we could define $y_i \in \{-1, 1\}$ and rewrite the objective in (1) as $J(\mathbf{w}) = \sum_{i=1}^{N} \log(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i))$.)

(a) Show the the following log-logistic function is convex: **(10 pts)**

$$f(\mathbf{w}) = \log(1 + \exp(-\mathbf{w}^T \mathbf{x})).$$

You may assume that $\mathbf{x}$ is a scalar and show that the second derivative of $f(\mathbf{w})$ w.r.t. $\mathbf{w}$ is always positive.

(b) You can prove that the expression for $\nabla_{\mathbf{w}} J = \sum_{i=1}^{N} (y_i - h_{\mathbf{w}}(\mathbf{x}_i)) \mathbf{x}_i$, where $h_{\mathbf{w}}(\mathbf{x}) = \dfrac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$ (e.g., see cs229 notes 1). Write a Python script for solving (1) using gradient ascent algorithm with a fixed step size $\eta$. **(10 pts)**

(c) Implement a solution for Exercise 1B described at the following link using Python: MNIST logistic regression **(30 pts)**
*See the last question for more information about downloading MNIST data in Python.*

**H2.2** Perceptron learning algorithm (PLA) and kernel extension: In this problem you will implement PLA to build a simple binary classification system. Suppose we have labeled data $(\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathbb{R}^{d+1}$ with $\mathbf{x}_i(1) = 1$ and $y \in \{-1, +1\}$. Let us define $h_{\mathbf{w}}(\mathbf{x}_i) = \operatorname{sign}\left(\mathbf{w}^T \mathbf{x}_i\right)$. We want to compute $\mathbf{w}$ using PLA. **(50 pts)**

(a) Data Processing: Generate two examples of 2D *linearly separable* dataset with $N = 100$ samples each. (To do this, you will first generate a weight and bias vector, $\mathbf{w}$, and then assign $\pm 1$ labels to your data samples as $y_i = h_{\mathbf{w}}(x_i)$.) Let us call the two datasets "Data1" and "Data2". For Data1, randomly select 80% of the samples for training and the remaining 20% for testing on Data1 (80/20). For Data2, randomly select 30% of the samples for training and the remaining 70% for testing (30/70).

(b) Implementation: Write a script for PLA in Python by initializing $\mathbf{w} = 0$ and using the following update rule:

$$\text{for } i=1,\ldots,N$$

$$\mathbf{w} = \mathbf{w} + \frac{1}{2}(y_i - h_{\mathbf{w}}(\mathbf{x}_i))\mathbf{x}_i.$$

(c) Training: Use your PLA implementation to train two linear classifiers for Data1 and Data2 using the respective training samples. Plot the training samples, test samples, and separation lines for both examples in two separate plots.

(d) Testing: Using the test samples in each example, compute precision, recall, and f1-score for your classifiers.

Precision is defined as $\dfrac{tp}{tp + fp}$; $tp$ and $fp$ denote the number of *true* and *false* positives.

Recall is defined as $\dfrac{tp}{tp + fn}$; $fn$ denotes the number of *false* negatives.

f1-score is defined as $2\dfrac{precision \times recall}{precision + recall}$; it can be viewed as a measure of accuracy and ranges between 0 (worst) and 1 (best).

(e) Kernel perceptron: Note that we can define the weight vector as $\mathbf{w} = \sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i$. We can view this as the dual perceptron in which we initialize $\alpha_i = 0$ for $i = 1, \ldots, N$ and update them as follows.
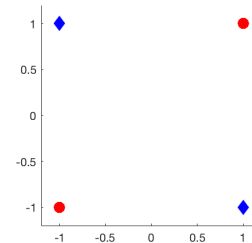
$$\text{for } j=1,\ldots,N$$

$$\hat{y} = \text{sign}\left(\sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_j\right)$$

$$\text{if } \hat{y} \neq y_j, \text{ update } \alpha_j = \alpha_j + 1$$

Implement the kernel perceptron and demonstrate that the classifier you learn is identical to the one in step (b).

**H2.3** Use a second degree polynomial kernel to classify four samples in 2D space, distributed in an XOR pattern (i.e., points on opposite diagonals of a rectangle have same labels, as shown in the table and figure below.) **(40 pts)**

| $y$ | $\mathbf{x}(1)$ | $\mathbf{x}(2)$ |
|-----|-----|-----|
| 1 | -1 | -1 |
| -1 | -1 | 1 |
| -1 | 1 | -1 |
| 1 | 1 | 1 |



**Table 1:** Table for class labels and sample values in an XOR pattern

**Figure 1:** Blue diamonds and red circles denote samples for two classes

(a) Map 2D feature vectors $\mathbf{x} = \begin{bmatrix} \mathbf{x}(1) & \mathbf{x}(2) \end{bmatrix}^T$ using kernel function $K(\mathbf{x}_1, \mathbf{x}_2) = (1 + \mathbf{x}_1^T \mathbf{x}_2)^2$, which results in 6D feature vectors: $\phi(\mathbf{x}) = \begin{bmatrix} 1 & \sqrt{2}\mathbf{x}(1) & \sqrt{2}\mathbf{x}(2) & \mathbf{x}(1)^2 & \mathbf{x}(2)^2 & \sqrt{2}\mathbf{x}(1)\mathbf{x}(2) \end{bmatrix}^T$[1]. Show that the features vectors are linearly separable in the 5D space.

---

[1]You can also use kernel function $K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^T \mathbf{x}_2)^2$, which results in 3D feature vectors: $\phi(\mathbf{x}) = \begin{bmatrix} \mathbf{x}(1)^2 & \mathbf{x}(2)^2 & \sqrt{2}\mathbf{x}(1)\mathbf{x}(2) \end{bmatrix}^T$.

(b) The dual objective for SVM can be written as

$$L(\alpha) = -\frac{1}{2}\alpha^T G\alpha + \sum_{i=1}^{4} \alpha_i,$$

where $G$ is a $4 \times 4$ matrix with $G(i,j) = y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$. Compute $G$ matrix and $\alpha$ that maximize $L(\alpha)$.

(c) Show that $\alpha$ computed in the previous step satisfies the following constrains: $\alpha_i \geq 0$ and $\sum_{i=1}^{4} \alpha_i y_i = 0$.

(d) Compute the SVM weight vector $\mathbf{w} = \sum_i \alpha_i y_i \phi(\mathbf{x}_i)$ and an expression for $h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \sum_i \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x})$.

**H2.4** Support vector machines (SVM): **(60 pts)**
In this problem you will use SVMs to build a simple text classification system. To begin, you need to get the MNIST dataset. You can do this in Python using the commands

```
from sklearn.datasets import fetch_mldata
mnist = fetch_mldata('MNIST original')
```

This downloads the dataset and stores it in a default location ($\sim$/scikit learn data/). You can also use the optional parameter data `home=path` in `fetch mldata` to direct the dataset to a custom path if you wish. The first time you run this command, the data will be downloaded from `mldata.org`, but once it has been downloaded this will simply load the dataset into the variable `mnist`.

This data set contains 70,000 handwritten digits, each of size $28 \times 28$ pixels[2]. You should begin by putting this data into a more convenient form by setting

```
X = mnist.data
y = mnist.target
```

Each row of `X` corresponds to one image. You can use the following to plot the $j^{th}$ image:

```
import matplotlib.pyplot as plt
plt.title('The jth image is a {label}'.format(label=int(y[j])))
plt.imshow(X[j].reshape((28,28)), cmap='gray')
plt.show()
```

In this problem you will build a classifier to classify between the (sometimes very similar) images of the digits "4" and "9". To get just this data, you can use

```
X4 = X[y==4,:]
X9 = X[y==9,:]
```

There are a little under 7,000 examples from each class. You should begin by using this data to form three distinct datasets: the first 4,000 points from each class will be used for designing the classifier (this is the *training* set), and the remainder will be used to evaluate the performance of our classifier (this is the *testing* set).

Since SVMs involve tuning a parameter $C$, you will also need to set this parameter. We will do this in a principled way using the so-called "holdout method". This means that you take the training set and divide it into two parts: you use the first to fit the classifier, and the second—which we call the *holdout* set—to gauge performance for a given value of $C$. You will want to decide on a finite set of

---

[2]You can learn more about this dataset at http://yann.lecun.com/exdb/mnist/.

"grid points" on which to test $C$ (I would suggest a logarithmic grid of values to test both $C \ll 1$ and $C \gg 1$). For each value of $C$, you will train an SVM on the first part of the set, and then compute the error rate on the holdout set. In the end, you can then choose the value of $C$ that results in a classifier that gives the smallest error on the holdout set

*Note:* Once you have selected $C$, you may then retrain on all the entire training set (including the holdout set), but you should never use the testing set until every parameter in your algorithm is set. These are used exclusively for computing the final test error.

To train the SVM , you can use the built in solver from scikit-learn. As an example, to train a linear SVM on the full dataset with a value of $C = 1$, we would use

```
from sklearn import svm
clf= svm.SVC(C=1.0,kernel='linear')
clf.fit(X,y)
```

Once you have trained the classifier, you can calculate the probability of error for the resulting classifier on the full dataset via

```
Pe = 1 - clf.score(X,y)
```

(a) Train an SVM using the kernels $k(u, v) = (u^T v + 1)^p$, $p = 1, 2$, that is, the inhomogeneous linear and quadratic kernel. To do this, you will want to set the parameters `kernel='poly'` and `degree=1` or `degree=2`.

For each kernel, report the best value of $C$, the test error, and the number of data points that are support vectors (returned via `clf.support_vectors_`). Turn in your code.

(b) Repeat the above using the radial basis function kernel $k(u, v) = e^{-\gamma \|u-v\|^2}$ (by setting the parameters `kernel='rbf'`, `gamma=gamma`. You will now need to determine the best value for both $C$ and $\gamma$. Report the best value of $C$ and $\gamma$, the test error, and the number of support vectors.

(c) For each kernel, turn in a $4 \times 4$ subplot showing images of the 16 support vectors that violate the margin by the greatest amount (these are, in a sense, the "hardest" examples to classify), and explain how these are determined. Above each subplot, indicate the true label (4 or 9).

To help get started with this, you can use the following code:

```
f, axarr = plt.subplots(4, 4)
axarr[0, 0].imshow(X[j].reshape((28,28)), cmap='gray')
axarr[0, 0].set_title('{label}'.format(label=int(y[j])))
...
plt.show()
```

You may find some help from this blog post: http://www.eric-kim.net/eric-kim-net/posts/1/kernel_trick.html

**Maximum points: 200**