

UNIVERSITY OF MICHIGAN
Department of Electrical Engineering and Computer Science
EECS 445 Introduction to Machine Learning
Winter 2019

Homework 2, Due: Tue. 02/26 at 11:59pm

Submission: Please upload your completed assignment by 11:59pm on Feb. 26, 2019 to Gradescope.

1 Support Vectors [8 pts]

1.1 SVM Primal [4 pts]

Suppose we are looking for a maximum hard margin (no slack variables) linear classifier *through the origin*, i.e., $b = 0$ with binary labels in $\{-1, 1\}$. In other words, we minimize $\frac{\|\bar{\theta}\|^2}{2}$ subject to $y^{(i)}\bar{\theta} \cdot \bar{x}^{(i)} \geq 1$, $i = 1, \dots, n$

- a) (2pts) Given a single training vector $\bar{x} = [a_1, a_2]^T$ with label $y = -1$, **what** is the $\bar{\theta}^*$ that satisfies the above constrained minimization?

Hint: Try thinking about this problem geometrically.

One way to solve this is to use Lagrange multipliers. We move the constraint into the objective function and introduce a Lagrange multiplier to get the following optimization:

$$\max_{\lambda} \min_{\bar{\theta}} \frac{\|\bar{\theta}\|^2}{2} + \lambda(1 + \bar{\theta} \cdot \bar{x})$$

$$\lambda \geq 0$$

Note that we've substituted $y = -1$. For minimizing wrt $\bar{\theta}$, we take the derivative of the above wrt $\bar{\theta}$ and set it equal to zero to get:

$$\bar{\theta} + \lambda \bar{x} = 0$$

We also note that since there is only a single point \bar{x} , it must be a support vector. Recall that $1 = y(\bar{\theta}^* \cdot \bar{x} + b)$ for any support vector \bar{x} that lies on the margin, and since here we are using a hard-margin, every support vector lies on the margin! This gives us another equation:

$$1 + \bar{\theta} \cdot \bar{x} = 0$$

We then have the following system of equations:

$$\theta_1 + \lambda a_1 = 0$$

$$\theta_2 + \lambda a_2 = 0$$

$$1 + \theta_1 a_1 + \theta_2 a_2 = 0$$

We can solve this linear system of equations to get the following: $\bar{\theta}^* = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} = -\frac{1}{a_1^2 + a_2^2} [a_1, a_2]^T$

We can also derive the solution geometrically. Since we are looking at a hard margin linear classifier through the origin, we know that our point $[a_1, a_2]$ is on the margin and at a distance $\frac{1}{\|\bar{\theta}\|}$ away from the origin. This gives us:

$$\sqrt{a_1^2 + a_2^2} = \frac{1}{\|\bar{\theta}\|}$$

$$\bar{\theta} \cdot \bar{\theta} = \frac{1}{a_1^2 + a_2^2}$$

$$-\bar{\theta} \cdot \bar{\theta} = \frac{-1}{a_1^2 + a_2^2}$$

We know that $\bar{\theta} \cdot \bar{x} = -1$, because $y = -1$. Thus:

$$-\bar{\theta} \cdot \bar{\theta} = \frac{\bar{\theta} \cdot \bar{x}}{a_1^2 + a_2^2}$$

$$\bar{\theta}^* = \frac{-\bar{x}}{a_1^2 + a_2^2} = \frac{-[a_1, a_2]^T}{a_1^2 + a_2^2} = -\frac{1}{a_1^2 + a_2^2} [a_1, a_2]^T$$

Another method is to solve the dual directly.

$$\max_{\alpha} \min_{\bar{\theta}} \frac{1}{2} \|\bar{\theta}\|^2 + \alpha(1 - y\bar{\theta} \cdot \bar{x})$$

subject to $\alpha \geq 0$

As we derived in class, this simplifies to:

$$\max_{\alpha} \alpha - \frac{1}{2} \alpha^2 \bar{x} \cdot \bar{x}$$

We can solve for α by taking the gradient:

$$1 - \alpha \bar{x} \cdot \bar{x} = 0$$

$$\alpha = \frac{1}{\bar{x} \cdot \bar{x}} = \frac{1}{a_1^2 + a_2^2}$$

We also showed in class that $\bar{\theta}^* = \sum_{i=1}^n \alpha_i y^{(i)} \bar{x}^{(i)}$, so:

$$\bar{\theta}^* = \alpha y \bar{x} = -\frac{[a_1, a_2]^T}{a_1^2 + a_2^2} = -\frac{1}{a_1^2 + a_2^2} [a_1, a_2]^T$$

- b) (1pt) Suppose we have two training examples, $\bar{x}^{(1)} = [-2, -1]^T$ and $\bar{x}^{(2)} = [-1, -1]^T$ with labels $y^{(1)} = 1$ and $y^{(2)} = -1$. **What** is $\bar{\theta}^*$ in this case, and **what** is the margin γ to the support vector?

In this case, we have two points of different labels. Note also that the angle between the two points is less than ninety degrees. Intuitively, what limits “pushing” the boundary around the margin will be these two points. This means that each point must be a support vector. Again, since we are using a hard margin, these points will be on the boundary. This gives us two equations:

$$1 - y^{(1)} \bar{\theta} \cdot \bar{x}^{(1)} = 0$$

$$1 - y^{(2)} \bar{\theta} \cdot \bar{x}^{(2)} = 0$$

We can solve the linear system:

$$1 - \bar{\theta} \cdot [-2, -1]^T = 0$$

$$1 + \bar{\theta} \cdot [-1, -1]^T = 0$$

$$-2\bar{\theta}_1 - \bar{\theta}_2 = 1$$

$$\bar{\theta}_1 + \bar{\theta}_2 = 1$$

$$\bar{\theta}^* = [-2, 3]^T, \gamma = \frac{1}{\|\bar{\theta}^*\|} = \frac{1}{\sqrt{13}}$$

- c) (1pt) **How** would the classifier and the margin in the previous question change if the offset parameter b were allowed to be non-zero? What are $(\bar{\theta}^*, b^*)$ and γ in this case?

In this example, we again have two support vectors and two equations (as above), but we have three unknowns! How do we solve this problem? Well, we could use Lagrange multipliers to get additional equations, but there is an easier way: by observation! Observe that since the data only varies along the horizontal axis, the decision boundary will be vertical. Thus the norm of the decision boundary, which is at a right angle to the boundary, will be vertical. This gives us three equations:

$$1 - y^{(1)}(\bar{\theta} \cdot \bar{x}^{(1)} + b) = 0$$

$$1 - y^{(2)}(\bar{\theta} \cdot \bar{x}^{(2)} + b) = 0$$

$$\theta_2 = 0$$

We can solve this system of equations:

$$1 - (\bar{\theta} \cdot [-2, -1]^T + b) = 0$$

$$1 + (\bar{\theta} \cdot [-1, -1]^T + b) = 0$$

$$-2\bar{\theta}_1 + b = 1$$

$$\bar{\theta}_1 - b = 1$$

$$\bar{\theta}^* = [-2, 0]^T, b^* = -3, \gamma = \frac{1}{\sqrt{4}} = \frac{1}{2}$$

1.2 SVM Dual [4 pts]

Consider solving the dual SVM optimization problem in the simple case of only two 2-dimensional training examples $\bar{x}^{(1)}, \bar{x}^{(2)} \in \mathbb{R}^2$ with labels $y^{(1)} = 1$ and $y^{(2)} = -1$. For this problem, we assume in addition that $\|\bar{x}^{(1)}\| = \|\bar{x}^{(2)}\| = 1$ and that the kernel function is $K(\bar{x}, \bar{x}') = \bar{x} \cdot \bar{x}'$ (linear kernel). Note that we are including the offset parameter in this case. The goal is to maximize

$$\sum_{i=1}^2 \alpha_i - \frac{1}{2} \sum_{i=1}^2 \sum_{j=1}^2 \alpha_i \alpha_j y^{(i)} y^{(j)} K(\bar{x}^{(i)}, \bar{x}^{(j)})$$

subject to $\alpha_i \geq 0$ and $\sum_{i=1}^2 \alpha_i y^{(i)} = 0$. Express your answers in terms of $\bar{x}^{(1)}$ and $\bar{x}^{(2)}$.

- a) (1pt) **What** are the resulting optimal values for the Lagrange multipliers α_1^* and α_2^* ? (Hint: use the constraints first, write the optimization problem in terms of α_1^* alone). Make sure that you use the condition on the norm of the training examples.

We can expand the objective function:

$$\max_{\alpha_1, \alpha_2} \alpha_1 + \alpha_2 - \frac{1}{2} \left(\alpha_1^2 y^{(1)2} \|\bar{x}^{(1)}\|^2 + \alpha_2^2 y^{(2)2} \|\bar{x}^{(2)}\|^2 + 2\alpha_1 \alpha_2 y^{(1)} y^{(2)} \bar{x}^{(1)} \cdot \bar{x}^{(2)} \right)$$

Noting that $y^{(i)2} = 1$ and $\|\bar{x}^{(i)}\| = 1$, we write:

$$\max_{\alpha_1, \alpha_2} \alpha_1 + \alpha_2 - \frac{1}{2} \left(\alpha_1^2 + \alpha_2^2 - 2\alpha_1 \alpha_2 \bar{x}^{(1)} \cdot \bar{x}^{(2)} \right)$$

$$\text{subject to } \alpha_1 \geq 0, \alpha_2 \geq 0, \alpha_1 = \alpha_2$$

Where the constraints come from expanding out the terms and substituting in the label values. Since $\alpha_1 = \alpha_2$, we can again re-write our objective:

$$\max_{\alpha_1} 2\alpha_1 - \frac{1}{2} \left(2\alpha_1^2 - 2\alpha_1^2 \bar{x}^{(1)} \cdot \bar{x}^{(2)} \right)$$

We can take the derivative and set it equal to 0:

$$2 - 2\alpha_1 + 2\alpha_1 \bar{x}^{(1)} \cdot \bar{x}^{(2)} = 0$$

We solve to get:

$$\alpha_1 = \alpha_2 = \frac{1}{1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)}} = \frac{2}{\|\bar{x}^{(1)} - \bar{x}^{(2)}\|^2}$$

- b) (1pt) Using (a), write the optimal parameter values of $\bar{\theta}^*$ in terms of the training examples only ($\bar{x}^{(1)}$ and $\bar{x}^{(2)}$)

Using the fact that $\bar{\theta}^* = \sum_{i=1}^2 \alpha_i y^{(i)} \bar{x}^{(i)}$, we get:

$$\begin{aligned} \bar{\theta}^* &= \frac{\bar{x}^{(1)}}{1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)}} - \frac{\bar{x}^{(2)}}{1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)}} = \frac{\bar{x}^{(1)} - \bar{x}^{(2)}}{1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)}} \\ \bar{\theta}^* &= \frac{2\bar{x}^{(1)}}{\|\bar{x}^{(1)} - \bar{x}^{(2)}\|^2} - \frac{2\bar{x}^{(2)}}{\|\bar{x}^{(1)} - \bar{x}^{(2)}\|^2} = \frac{2(\bar{x}^{(1)} - \bar{x}^{(2)})}{\|\bar{x}^{(1)} - \bar{x}^{(2)}\|^2} \end{aligned}$$

These solutions are equivalent:

$$\begin{aligned} \frac{2(\bar{x}^{(1)} - \bar{x}^{(2)})}{\|\bar{x}^{(1)} - \bar{x}^{(2)}\|^2} &= \frac{2(\bar{x}^{(1)} - \bar{x}^{(2)})}{(\bar{x}^{(1)} - \bar{x}^{(2)}) \cdot (\bar{x}^{(1)} - \bar{x}^{(2)})} = \frac{2(\bar{x}^{(1)} - \bar{x}^{(2)})}{\bar{x}^{(1)} \cdot \bar{x}^{(1)} - 2\bar{x}^{(1)} \cdot \bar{x}^{(2)} + \bar{x}^{(2)} \cdot \bar{x}^{(2)}} \\ &= \frac{2(\bar{x}^{(1)} - \bar{x}^{(2)})}{2(1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)})} = \frac{\bar{x}^{(1)} - \bar{x}^{(2)}}{1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)}} \end{aligned}$$

- c) (1pt) **What** is the offset parameter b^* ?

Since the norm of the points are fixed at one, we know that all the points that we considering together form a circle of radius one. Intuitively, the best way to separate any two of the points on the circle is a line between the origin and the midpoint of the two points (thus there should be no offset). Now, let's consider this algebraically, Note that in this case, we are using a hard-margin SVM, so all points that are support vectors lie on the margin. Thus, we know that for any support vector \bar{x} and its associated label y , we have $1 = y(\bar{\theta}^* \cdot \bar{x} + b)$. We can solve this equation using either $\bar{x}^{(1)}, y^{(1)}$ or $\bar{x}^{(2)}, y^{(2)}$ to get:

$$b = 0$$

d) (1pt) **What** is the distance (margin γ) from the decision boundary to the support vector?

$$\begin{aligned} \gamma &= \frac{1}{\|\bar{\theta}^*\|} = \frac{1}{\left\| \frac{\bar{x}^{(1)} - \bar{x}^{(2)}}{1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)}} \right\|} = \frac{\|1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)}\|}{\|\bar{x}^{(1)} - \bar{x}^{(2)}\|} = \frac{1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)}}{\sqrt{(\bar{x}^{(1)} - \bar{x}^{(2)}) \cdot (\bar{x}^{(1)} - \bar{x}^{(2)})}} \\ &= \frac{1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)}}{\sqrt{1 - 2\bar{x}^{(1)} \cdot \bar{x}^{(2)} + 1}} = \frac{1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)}}{\sqrt{2(1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)})}} = \frac{\sqrt{1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)}}}{\sqrt{2}} = \frac{\sqrt{2}\sqrt{1 - \bar{x}^{(1)} \cdot \bar{x}^{(2)}}}{2} \\ &= \frac{\sqrt{2 - 2\bar{x}^{(1)} \cdot \bar{x}^{(2)}}}{2} = \frac{\|\bar{x}^{(1)} - \bar{x}^{(2)}\|}{2} \\ \gamma &= \frac{1}{\|\bar{\theta}^*\|} = \frac{1}{\left\| \frac{2(\bar{x}^{(1)} - \bar{x}^{(2)})}{\|\bar{x}^{(1)} - \bar{x}^{(2)}\|^2} \right\|} = \frac{\|\bar{x}^{(1)} - \bar{x}^{(2)}\|^2}{2\|\bar{x}^{(1)} - \bar{x}^{(2)}\|} = \frac{\|\bar{x}^{(1)} - \bar{x}^{(2)}\|}{2} \end{aligned}$$

2 Soft Margin SVM Dual [2 pts]

Consider the dual formulation of a soft-margin SVM with regularization, where C is a regularization hyper-parameter.

$$\begin{aligned} &\underset{\bar{\alpha}}{\text{maximize}} && \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y^{(i)} y^{(j)} \bar{x}^{(i)} \cdot \bar{x}^{(j)} \\ &\text{subject to} && \sum_{i=1}^N \alpha_i y^{(i)} = 0, \\ &&& 0 \leq \alpha_i \leq C, \forall i = 1, \dots, N. \end{aligned}$$

i	$\bar{x}^{(i)}$	$y^{(i)}$	α_i
1	[-2.78,-2.08]	-1	0
2	[-5.20,1.99]	-1	0
3	[-2.69,4.26]	-1	0.05
4	[-2.25,-3.96]	-1	0
5	[-2.59,-3.36]	-1	0
6	[-1.17,1.44]	-1	0.05
7	[-1.94,-3.68]	-1	0.0159
8	[-1.19,-2.99]	-1	0.05
9	[-2.14,-0.22]	-1	0.05
10	[-3.12,-1.99]	-1	0

i	$x^{(i)}$	$y^{(i)}$	α_i
11	[0.68, 1.76]	1	0.05
12	[2.38, 0.52]	1	0
13	[-1.49, 1.65]	1	0.05
14	[1.74, 0.22]	1	0.05
15	[-0.13, -1.77]	1	0.05
16	[1.59, 5.35]	1	0
17	[1.65, 2.17]	1	0.0159
18	[2.56, -4.04]	1	0
19	[2.75, 0.73]	1	0
20	[2.34, 2.92]	1	0

- a) (1pt) Given the table above, which points are the support vectors?

All points s.t. $\alpha_i \neq 0$, including both points exactly on the margin and those inside the margins: 3, 6, 7, 8, 9, 11, 13, 14, 15, 17

- b) (1pt) Assuming the regularization hyperparameter $C = 0.05$, given the α_i for each of the points above, what are the equations to compute the optimal parameter values $\bar{\theta}^*$ and b^* ? What are the resulting $\bar{\theta}^*$ and b^* for the data above?

$$\bar{\theta}^* = \sum_{i=1}^n \alpha_i y^{(i)} \bar{x}^{(i)}$$

$$\bar{\theta}^* = 0.05 * -1 * [-2.69, 4.26] + 0.05 * -1 * [-1.17, 1.44] + 0.0159 * -1 * [-1.94, -3.68] + 0.05 * -1 * [-1.19, -2.99] + 0.5 * -1 * [-2.14, -0.22] + 0.5 * [0.68, 1.76] + 0.5 * [-1.49, 1.65] + 0.5 * [1.74, 0.22] + 0.5 * [-0.13, -1.77] + 0.0159 * [1.65, 2.17] = [0.456581, 0.061515]$$

For points on the margin, we know that they satisfy $y^{(i)}(\bar{\theta} \cdot \bar{x}^{(i)} + b) = 1$. Thus, we can use points on the margin to directly solve for b^* . We know points on the margin have $0 < \alpha < C$, so points 7 and 17 are on the margin.

Acceptable answers for offset calculation:

Using first point on the margin (i=7): $b^* = y^{(7)} - \bar{\theta} \cdot \bar{x}^{(7)} = (-1) - (-1.11214234) = 0.11214234$

Using second point on the margin (i=17): $b^* = y^{(17)} - \bar{\theta} \cdot \bar{x}^{(17)} = (1) - (0.8868462) = 0.1131538$

Using average of both: $b^* = \frac{1}{2}(0.11214234 + 0.1131538) = 0.11264807$

3 Kernels [8 pts]

3.1 From feature mapping to kernel [2 pts]

We have two kernels, $K_1(\bar{x}, \bar{z}) = \bar{x} \cdot \bar{z}$ and $K_2(\bar{x}, \bar{z}) = ?$, where $\bar{x}, \bar{z} \in \mathbb{R}^2$. The first kernel is defined for us, but the second kernel is unknown. We can define a third kernel, which is the product of the first two: $K_3(\bar{x}, \bar{z}) = K_1(\bar{x}, \bar{z})K_2(\bar{x}, \bar{z})$. The feature mapping corresponding to $K_3(\bar{x}, \bar{z})$ is given:

$$\phi(\bar{x}) = [x_1^3, \sqrt{2}x_1^2x_2, x_1x_2^2, x_1^2x_2, \sqrt{2}x_1x_2^2, x_2^3, \sqrt{6}x_1^2, 2\sqrt{3}x_1x_2, \sqrt{6}x_2^2, 3x_1, 3x_2]$$

In other words, $K_3(\bar{x}, \bar{z}) = \phi(\bar{x})\phi(\bar{z})$.

Hint: Group the first six terms, the next three, and the last two terms together, and simplify from there.

(a) What is $K_2(\bar{x}, \bar{z})$?

$$\begin{aligned}
 K_3(\bar{x}, \bar{z}) &= \phi(\bar{x})\phi(\bar{z}) \\
 &= (x_1^3 z_1^3 + 2x_1^2 x_2 z_1^2 z_2 + x_1 x_2^2 z_1 z_2^2 + x_1^2 x_2 z_1^2 z_2 + 2x_1 x_2^2 z_1 z_2^2 + x_2^3 z_2^3) + (6x_1^2 z_1^2 + 12x_1 x_2 z_1 z_2 + 6x_2^2 z_2^2) + \\
 &\quad (9x_1 z_1 + 9x_2 z_2) \\
 &= (x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + x_2^2 z_2^2)(\bar{x} \cdot \bar{z}) + (6\bar{x} \cdot \bar{z})(\bar{x} \cdot \bar{z}) + 9(\bar{x} \cdot \bar{z}) \\
 &= ((\bar{x} \cdot \bar{z})^2 + 6\bar{x} \cdot \bar{z} + 9)(\bar{x} \cdot \bar{z}) \\
 &= (\bar{x} \cdot \bar{z} + 3)^2 (\bar{x} \cdot \bar{z}) \\
 K_2(\bar{x}, \bar{z}) &= (\bar{x} \cdot \bar{z} + 3)^2.
 \end{aligned}$$

3.2 Kernelizing logistic regression [6 pts]

Here is the algorithm of SGD for logistic regression.

```

1 Initialization:  $\bar{\theta}^{(0)} = \bar{0}$ 
2 Repeat until convergence:
3   for  $i = 1, \dots, N$ :
4      $\bar{\theta}^{(k+1)} \leftarrow \bar{\theta}^{(k)} + \eta \sigma(-y^{(i)} \bar{\theta}^{(k)} \cdot \bar{x}^{(i)}) y^{(i)} \bar{x}^{(i)}$ 
5      $k \leftarrow k + 1$ 
6  $\bar{\theta} = \bar{\theta}^{(k)}$ 
7 Calculate probability that  $\bar{x}$  is classified with label  $y = +1$ :
8  $h(\bar{x}) = \sigma(\bar{\theta} \cdot \bar{x}) = \frac{1}{1 + \exp(-\bar{\theta} \cdot \bar{x})}$ 

```

We would like to kernelize this algorithm. This can be done by mapping each example $\bar{x} \in \mathbb{R}^d$ into a new feature vector $\phi(\bar{x}) \in \mathbb{R}^{d'}$ (in many cases, $d' > d$) and reformulating the algorithm such that it only uses the kernel function $k(\bar{x}, \bar{x}')$ corresponding to $\phi(\bar{x}) \cdot \phi(\bar{x}')$. This will require a number of changes. In the algorithm above, we need to calculate $\bar{\theta} \cdot \bar{x}$. However, once we introduce a feature mapping $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$, $\bar{\theta}$ will lie in a d' -dimensional space, where d' may be very large. We want a new formulation that does not explicitly depend on $\bar{\theta}$.

(a) (2pts) We will derive a new version of the above algorithm, by first reformulating $\bar{\theta}$ as a combination of training examples $\bar{x}^{(i)}$ for $i = 1, \dots, N$, i.e. $\bar{\theta} = \sum_{i=1}^N \alpha_i \bar{x}^{(i)}$.

Now consider a dataset with 3 data points, i.e. $N = 3$. Suppose we update the parameters for 9 times before convergence. These assumptions are only for part (a).

We DO NOT shuffle the points after each epoch. **Write down** the expression for α_i for $i = 1, 2, 3$ as a function of η , $\bar{x}^{(i)}$, $y^{(i)}$, and $\bar{\theta}^{(k)}$ for $k = 0, \dots, 8$.

$$\begin{aligned}
 \alpha_1 &= \eta \sigma(-y^{(1)} \bar{\theta}^{(0)} \cdot \bar{x}^{(1)}) y^{(1)} + \eta \sigma(-y^{(1)} \bar{\theta}^{(3)} \cdot \bar{x}^{(1)}) y^{(1)} + \eta \sigma(-y^{(1)} \bar{\theta}^{(6)} \cdot \bar{x}^{(1)}) y^{(1)} \\
 \alpha_2 &= \eta \sigma(-y^{(2)} \bar{\theta}^{(1)} \cdot \bar{x}^{(2)}) y^{(2)} + \eta \sigma(-y^{(2)} \bar{\theta}^{(4)} \cdot \bar{x}^{(2)}) y^{(2)} + \eta \sigma(-y^{(2)} \bar{\theta}^{(7)} \cdot \bar{x}^{(2)}) y^{(2)} \\
 \alpha_3 &= \eta \sigma(-y^{(3)} \bar{\theta}^{(2)} \cdot \bar{x}^{(3)}) y^{(3)} + \eta \sigma(-y^{(3)} \bar{\theta}^{(5)} \cdot \bar{x}^{(3)}) y^{(3)} + \eta \sigma(-y^{(3)} \bar{\theta}^{(8)} \cdot \bar{x}^{(3)}) y^{(3)}
 \end{aligned}$$

- (b) These new parameters, $\bar{\alpha}$, will help us replace $\bar{\theta}$ in the algorithm. We can turn the to-be-trained parameters from $\bar{\theta} \in \mathbb{R}^d$ to $\bar{\alpha} \in \mathbb{R}^N$. Next, we will come up with the appropriate update for $\bar{\alpha}$ as outlined in the algorithm below.

```

1 Initialization:  $\bar{\alpha}^{(0)} = \bar{0}$ 
2 Repeat until convergence:
3   for  $i = 1, \dots, N$ :
4     <Update all  $N$  elements in  $\bar{\alpha}^{(k+1)}$  in terms of  $\eta, \bar{\alpha}^{(k)}, y^{(i)}, \bar{x}^{(i)}$ >
5      $k \leftarrow k + 1$ 
6    $\bar{\alpha} = \bar{\alpha}^{(k)}$ 
7 Calculate probability that  $\bar{x}$  is classified with  $y = +1$ :
8 <Update compute  $h(\bar{x})$  in terms of  $\bar{\alpha}, \bar{x}^{(j)}$  for  $j = 1, \dots, N$ >

```

Consider the following questions:

- (i) (1pt) First, consider the parameter update (**line 4**). We seek to update the parameters of our model. Based on the relationship between $\bar{\alpha}$ and $\bar{\theta}$, think about how many elements in $\bar{\alpha}$ need updating. **How** should we update the parameters? You may write more than one line of pseudocode.
- (ii) (1pt) When we finish updating the parameters, we need to calculate the probability that a new example \bar{x} is classified as $y = +1$ (**line 8**). **How** can we express this probability in terms of $\bar{\alpha}$?

From the relationship between $\bar{\alpha}$ and $\bar{\theta}$ in part (a), we see that each $\bar{\theta}^{(k)}$ only affects the value of one element of $\bar{\alpha}$. Thus, when we are updating $\bar{\theta}$ in the original algorithm, we only have to modify the value of one element of $\bar{\alpha}$.

```

1 Initialization:  $\bar{\alpha}^{(0)} = \bar{0}$ 
2 Repeat until convergence:
3   for  $i = 1, \dots, N$ : (1pt)
4      $\bar{\alpha}^{(k+1)} \leftarrow \bar{\alpha}^{(k)}$ 
5      $\alpha_i^{(k+1)} \leftarrow \alpha_i^{(k)} + \eta \sigma(-y^{(i)} \sum_{j=1}^N \alpha_j^{(k)} \bar{x}^{(j)} \cdot \bar{x}^{(i)}) y^{(i)}$ 
6      $k \leftarrow k + 1$ 
7    $\bar{\alpha} = \bar{\alpha}^{(k)}$ 
8 Calculate probability that  $\bar{x}$  is classified with  $y = +1$ :
9  $h(\bar{x}) = \sigma(\sum_{j=1}^N \alpha_j \bar{x}^{(j)} \cdot \bar{x})$  (1pt)

```

- (c) (2pts) Now we can kernelize the algorithm from part (b): map each example \bar{x} into a feature vector $\phi(\bar{x})$; reformulate the algorithm such that it only uses the kernel function $k(\bar{x}, \bar{x}')$ corresponding to $\phi(\bar{x}) \cdot \phi(\bar{x}')$. **Complete** the algorithm below to kernelize the logistic regression:

```

1 Initialization:  $\bar{\alpha}^{(0)} = \bar{0}$ 
2 Repeat until convergence:
3   for  $i = 1, \dots, N$ :
4     # You can write more than one line here.
5     _____ (1pt)

```

```

6 |            $k \leftarrow k + 1$ 
7 |  $\bar{\alpha} = \bar{\alpha}^{(k)}$ 
8 | Calculate probability that  $\bar{x}$  is classified with  $y = +1$ :
9 |  $h(\bar{x}) =$  _____ (1pt)

```

```

1 | Initialization:  $\bar{\alpha}^{(0)} = \bar{0}$ 
2 | Repeat until convergence:
3 |   for  $i = 1, \dots, N$ : (1pt)
4 |      $\bar{\alpha}^{(k+1)} \leftarrow \bar{\alpha}^{(k)}$ 
5 |      $\alpha_i^{(k+1)} \leftarrow \alpha_i^{(k)} + \eta \sigma(-y^{(i)} \sum_{j=1}^N \alpha_j^{(k)} k(\bar{x}^{(j)}, \bar{x}^{(i)})) y^{(i)}$ 
6 |      $k \leftarrow k + 1$ 
7 |  $\bar{\alpha} = \bar{\alpha}^{(k)}$ 
8 | Calculate probability that  $\bar{x}$  is classified with  $y = +1$ :
9 |  $h(\bar{x}) = \sigma(\sum_{j=1}^N \alpha_j k(\bar{x}^{(j)}, \bar{x}))$  (1pt)

```

4 Entropy [2 pts]

Feature			Classification
Punctuality	Boarding Efficiency	Quality of Service	Satisfied with the flight
Delayed	Decent	Good	Yes
On Time	Decent	Poor	No
Delayed	Slow	Good	No
Delayed	Fast	Good	Yes
On Time	Slow	Great	No
Delayed	Fast	Poor	Yes
On Time	Decent	Good	Yes
On Time	Decent	Great	Yes
On Time	Slow	Poor	No

- a) (1pt) Consider the table above, which maps variables pertaining to a customer's experience to whether or not they were satisfied with their flight. We consider a categorical representation (as opposed to an ordinal representation). **Which feature(s)** result in the highest conditional entropy $H(Y|X)$, where Y is the outcome and X is one of the features (Punctuality, Boarding Efficiency, and Quality of Service)?

$$H(Y|\text{Punctuality}) = \frac{4}{9} \left(-\frac{1}{4} \log \frac{1}{4} - \frac{3}{4} \log \frac{3}{4} \right) + \frac{5}{9} \left(-\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} \right) = 0.9000$$

$$H(Y|\text{Boarding Efficiency}) = \frac{2}{9} \times 0 + \frac{3}{9} \times 0 + \frac{4}{9} \left(-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \right) = 0.3606$$

$$H(Y|\text{Quality Of Service}) = \frac{2}{9} \left(-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right) + \frac{3}{9} \left(-\frac{1}{3} \log \frac{1}{3} - \frac{2}{3} \log \frac{2}{3} \right) + \frac{4}{9} \left(-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \right) = 0.8889$$

The feature “Punctuality” has the highest conditional entropy of the three features.

- b) (1pt) **Calculate** the information gain $IG(X, Y)$ for each of the features and the outcome. **Which feature(s)** result in the highest information gain?

Our current entropy is $H(Y) = -\frac{5}{9} \log \frac{5}{9} - \frac{4}{9} \log \frac{4}{9} = 0.9911$. Thus:

$$IG(\text{Punctuality}, Y) = 0.9911 - 0.9000 = 0.0911$$

$$IG(\text{Boarding Efficiency}, Y) = 0.9911 - 0.3606 = 0.6305$$

$$IG(\text{Quality of Service}, Y) = 0.9911 - 0.8889 = 0.1022$$

The feature “Boarding Efficiency” gives the highest information gain.

5 Ensemble Methods [5 pts]

We will study the performance of two ensemble methods on the very popular MNIST dataset consisting of handwritten digits. This dataset contains 70000 samples (each a 28×28 grayscale image having 784 features) of handwritten digits, classified into 10 classes (0 through 9). Here, we will consider a subset of the data pertaining to four classes, which you can fetch using the provided `load_mnist(classes)` function. Please be aware that due to the relatively large dataset size, the code will take a longer time to run.

Within `HW2_ensemble.py`, the following functions have been implemented for you:

- `load_mnist(classes)`
- `get_avg_performance(X, y, m_vals, nsplits=50)`
- `plot_data(bagging_scores, random_forest_scores, m_range)`

It also contains the following function declarations that you will implement:

- `bagging_ensemble(X_train, y_train, X_test, y_test, n_clf=10)`
- `random_forest(X_train, y_train, X_test, y_test, m, n_clf=10)`

- a) (2pts) **Implement** `random_forest(X_train, y_train, X_test, y_test, m, n_clf=10)` based on the specification provided in the skeleton code. Random forest consists of `n_clf`-many decision trees where each decision tree is trained independently on a bootstrap sample of the training

data (for this problem, `n_clf=10`). For each node, we randomly select m features as candidates for splitting on.

Here, the final prediction of the bagging classifier is determined by a majority vote of these `n_clf` decision trees. In the case of ties, randomly sample among the plurality classes (i.e. the classes that are tied) to choose a label.

You should use the `sklearn.tree.DecisionTreeClassifier` class. Set `criterion='entropy'` to avoid the default setting. Also, see the `max_features` parameter within this class.

Note: Do not set the `max_depth` parameter. Remember that you are free to use helper functions to keep your code organized.

Implementations will vary

- b) (1pt) **Implement** `bagging_ensemble(X_train, y_train, X_test, y_test, n_clf=10)` based on the specification provided in the skeleton code. Like random forest, a bagging ensemble classifier consists of `n_clf`-many decision trees where each decision tree is trained independently on a bootstrap sample of the training data. However, all features are considered at every split. Again, the final prediction is determined by a majority vote of these `n_clf` decision trees.

Implementations will vary

- c) (2pts) Now, we will compare the performance of these ensemble classifiers using `get_avg_performance()`. Measure the median performance across 50 random splits of the digits dataset into training (80%) and test (20%) sets, and **plot** the returned performance for each algorithm over the range of m values specified in the skeleton code (use the `plot_data()` function). See Figure 1 for an example of the plot (yours may differ due to randomness, but the general trend should be the same). How does the average test performance of the two methods compare as we vary m (the size of the randomly selected subset of features)?

See the plot in Figure 1. Random forest initially does better than bagging: the best performance for random forest is achieved around $m = 56$, and steadily declines as m increases. This makes intuitive sense: as m increases, it lets each tree choose more of the same features, because the feature choice is less random. Thus, each tree becomes less independent, and the performance decreases. When, in the limit, random forest uses all available features, it is no different from bagging, which is why random forest tends towards bagging in this plot.

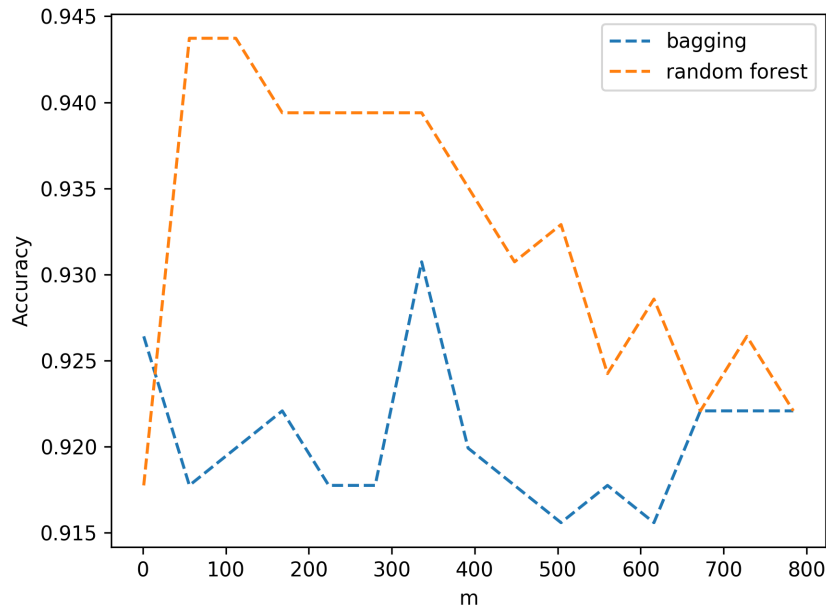


Figure 1: Plot of accuracy of random forest, bagging ensemble methods for $m = 1, 56, 112, \dots, 784$

6 Boosting [5 pts]

Boosting combines a set of weak classifiers into a stronger ensemble classifier $h_M(\bar{x}) = \sum_{j=1}^M \alpha_j h(\bar{x}; \bar{\theta}_j)$, where $\alpha_j \geq 0$ are votes allocated to each base classifier (“decision stump”) $h(\bar{x}; \bar{\theta}_j) = \text{sign}(\theta_1^{(j)} x_k + \theta_0^{(j)})$ described by a parameter vector $\bar{\theta}_j = \{k, \theta_1^{(j)}, \theta_0^{(j)}\}$ that encodes the co-ordinate, direction and location information. Finding a jointly optimum solution of $\bar{\theta}_j$ and α_j for all j is a hard problem and therefore, we take an iterative approach exemplified by the adaptive boosting (AdaBoost) algorithm:

```

Set  $W_0(i) = \frac{1}{n}$  for  $i = 1 \dots n$ 
For  $m = 1$  to  $M$  do:
    Find  $h(\bar{x}; \bar{\theta}_m)$  that minimizes the weighted training error  $\epsilon_m$ :
         $\epsilon_m = \sum_{i=1}^n W_{m-1}(i) [y^{(i)} \neq h(\bar{x}^{(i)}; \bar{\theta}_m)]$ 

    Given  $\bar{\theta}_m$ , compute  $\alpha_m$  that minimizes weighted training loss:
         $\alpha_m = \frac{1}{2} \log\left(\frac{1-\epsilon_m}{\epsilon_m}\right)$ 

    Update weights on all training examples
    For  $i = 1$  to  $n$  do:
         $W_m(i) = c_m W_{m-1}(i) \exp\{-y^{(i)} \alpha_m h(\bar{x}^{(i)}; \bar{\theta}_m)\}$ 

```

Consider the labeled training points in Figure 2, where the \bullet 's and \times 's denote negative and positive labels, respectively. We wish to apply Adaboost with decision stumps to solve the classification problem. In each boosting iteration, we select the stump that minimizes the weighted training error, breaking ties arbitrarily.

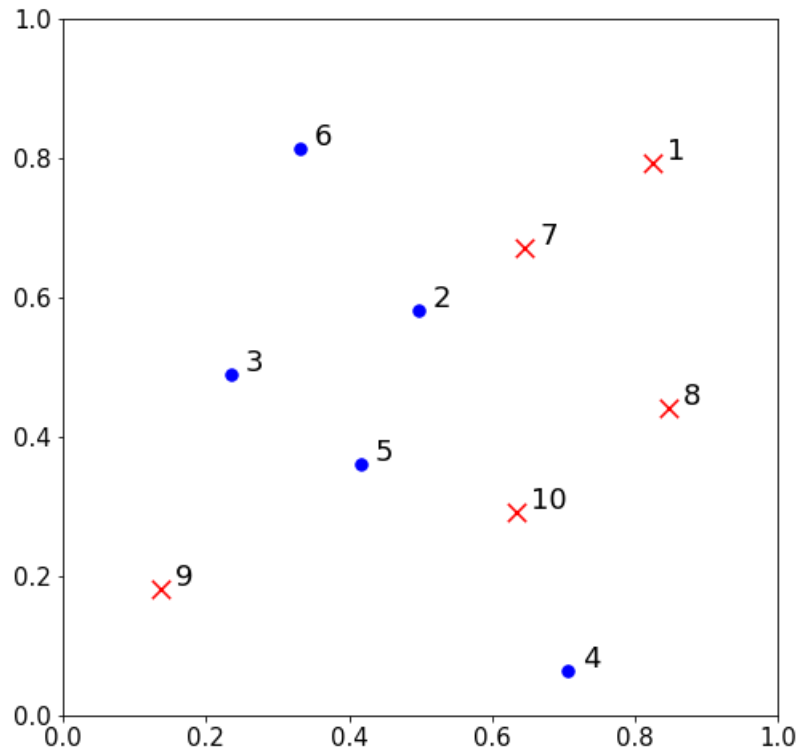


Figure 2: •-negative points and ×-positive points.

- a) (1pt) In Figure 2, **draw** the decision boundary corresponding to the first decision stump that the boosting algorithm would choose. **Label** this boundary as (1), and also **indicate** the $+/ -$ sides of the decision boundary. (Note: we have provided the data in the skeleton code associated with this problem, so that you can recreate the plot yourself if needed.)

Any vertical line between point 2 and point 10 that classifies the right side as positive and the left side as negative.

- b) (1pt) In the same Figure 2, **circle** the point(s) that have the largest weight after the first boosting iteration.

Points 4 and 9.

- c) (1pt) **What is the weighted error** of the first decision stump after the first boosting iteration, i.e., after the points have been re-weighted?

$$\begin{aligned}
 W_0(i) &= \frac{1}{10} \\
 \epsilon_1 &= \frac{2}{10} = \frac{1}{5} \\
 \alpha_1 &= \frac{1}{2} \log\left(\frac{1-\frac{1}{5}}{\frac{1}{5}}\right) = \frac{1}{2} \log(4) \\
 \text{For } i \neq 4, 9: W_1(i) &= c_1 W_0(i) \exp\left\{-\frac{1}{2} \log(4)\right\} = c_1 W_0(i) \exp\left\{\log\left(\frac{1}{2}\right)\right\} = \frac{1}{2} \frac{1}{10} c_1 = \frac{1}{20} c_1 \\
 \text{For } i = 4, 9: W_1(i) &= c_1 W_0(i) \exp\left\{\frac{1}{2} \log(4)\right\} = c_1 W_0(i) \exp\{\log(2)\} = c_1 \frac{1}{10} (2) = \frac{1}{5} c_1 \\
 c_1 &= \left(8 * \frac{1}{20} + 2 * \frac{1}{5} = \frac{4}{5}\right)^{-1} = \left(\frac{4}{5}\right)^{-1} = \frac{5}{4} \\
 \text{For } i \neq 4, 9: W_1(i) &= \frac{1}{20} c_1 = \frac{1}{16} \\
 \text{For } i = 4, 9: W_1(i) &= \frac{1}{5} c_1 = \frac{1}{4} \\
 \text{Weighted error of first decision stump:} \\
 \text{First decision stump misclassifies points 4 and 9, so } \epsilon &= \frac{1}{4} + \frac{1}{4} = \frac{1}{2} = 0.5
 \end{aligned}$$

- d) (1pt) **Draw** the decision boundary corresponding to the second decision stump, again in Figure 2, and **label** it as (2) also **indicating** the $+/ -$ sides of the boundary.

There are two valid choices here, both of which achieve an error 0.25. The choices are: either a vertical line between points 3 and 9 which classifies the left side as positive and the right direction as negative, or a horizontal line between points 4 and 9 that classifies the top side as positive and the bottom side as negative.

- e) (1pt) Will any of the points be misclassified by the combined classifier after the two boosting iterations? **Provide a brief justification**, no calculations are necessary (the point will be awarded for the justification, not whether your y/n answer is correct).

Yes. The first decision stump has a higher weight than the second. Since the first decision stump does not correctly classify points 4 and 9, these two points will not be correctly classified by the combined classifier.

REMEMBER Submit your completed assignment by 11:59pm on Feb. 26, 2018 to Gradescope.