**1.** Consider a binary linear classification where the data points are two dimensional, i.e.,
$\mathbf{x} \in \mathbb{R}^2$ and the labels $y \in \{-1, 1\}$. The training set consists of the following points:

$$\mathbf{x}_1 = (1, 0)^T, \qquad y_1 = +1$$
$$\mathbf{x}_2 = (-1, 0)^T, \qquad y_2 = -1$$
$$\mathbf{x}_3 = (1, d)^T, \qquad y_3 = +1$$

Assume that $d > 0$ is some constant. **For the perceptron algorithm please treat the points that fall exactly on the decision boundary as mistakes and do the update accordingly. Assume that the initial weight is the all-zero vector.**

1 mark

(**a**) Is the data-set linearly separable? Sketch the points in the $x_1 - x_2$ plane and show the labels.

6 marks

(**b**) Run the perceptron algorithm in the following order $(1, 2, 3), (1, 2, 3), \ldots$ until it converges. Show the output of the perceptron algorithm in each step, sketch the resulting decision boundary, and find the resulting margin (distance of the decision boundary to the nearest training point). s
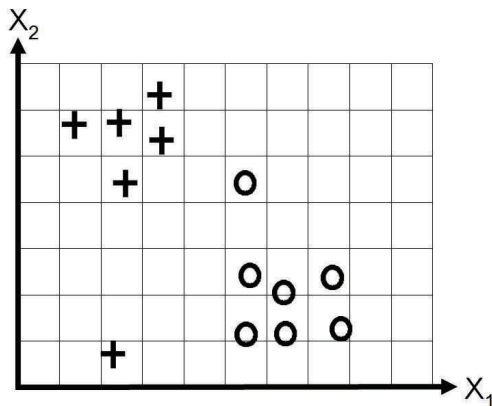
7 marks      (**c**) Repeat part (b) when the order is $(3, 2, 1), (3, 2, 1, ), \ldots$. Also comment on what the decision boundary approaches when $d \to \infty$

1 mark      (**d**) Between parts (b) and (c) which solution will be preferred? Briefly explain.

**2**. Consider a binary linear classification problem where $\mathbf{x} \in \mathbb{R}^2$ and $y \in \{-1, +1\}$. We illustrate the training dataset below. The '+' label refers to $y = +1$ and the 'o' label refers to $y = -1$. We would like to construct a classifier $h_{\mathbf{w}}(\mathbf{x}) = \text{sign}(w_0 + w_1 x_1 + w_2 x_2)$ where $\text{sign}(\cdot)$ is the *sign* function as discussed in class.



In the figure above, the adjacent vertical (and horizontal) lines are 1 unit apart from each other. Assume that the training points are above are $(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_N, y_N)$ (with $N = 13$). We consider the classification loss

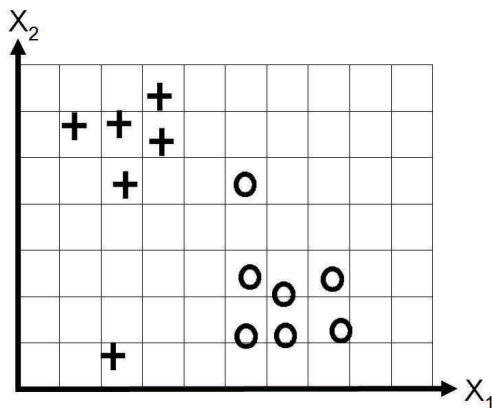$$L(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^{N} \mathbb{I}(y_i \neq h_{\mathbf{w}}(\mathbf{x}_i))$$

where $\mathbb{I}$ denotes the indicator function.

**(a)** Draw a decision boundary in the figure above that achieves zero training error.

**(b)** Suppose that we attempt to minimize the following loss function over $\mathbf{w}$ : $J(\mathbf{w}) = L(\mathbf{w}) + \lambda w_0^2$, where $\lambda = 10^7$ is a huge constant. Sketch a possible decision boundary in the figure below. How many points are mis-classified.
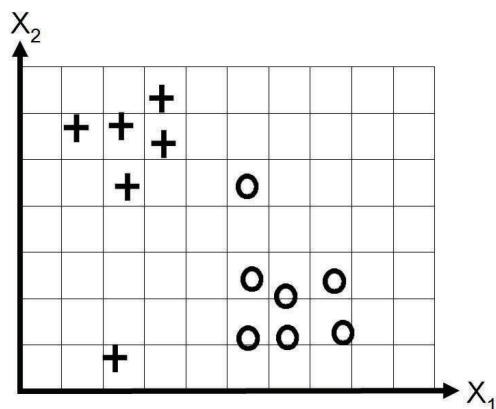
2 marks

(**c**) Suppose that we attempt to minimize the following loss function over $\mathbf{w}$ : $J(\mathbf{w}) = L(\mathbf{w}) + \lambda w_1^2$, where $\lambda = 10^7$ is a huge constant. Sketch a possible decision boundary in the figure below. How many points are mis-classified.

2 marks

(**d**) Suppose that we attempt to minimize the following loss function over $\mathbf{w}$ : $J(\mathbf{w}) = L(\mathbf{w}) + \lambda w_2^2$, where $\lambda = 10^7$ is a huge constant. Sketch a possible decision boundary in the figure below. How many points are mis-classified.
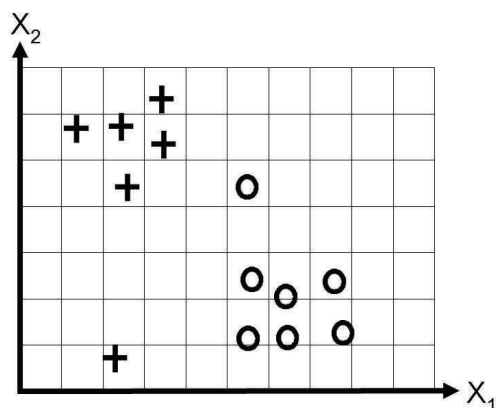
25 marks  **3**. Suppose that a data vector $\mathbf{y} \in \mathbb{R}^n$ and a matrix $\mathbf{A} \in \mathbb{R}^{n \times p}$ are given. Assume that $n > p$ and that $\mathbf{A}^T \mathbf{A}$ is an $p \times p$ invertible matrix. We would like to approximate $\mathbf{y}$ using a vector of the form $\hat{\mathbf{y}}_{\mathbf{w}} = \mathbf{A}\mathbf{w}$.

2 marks  **(a)** Let $\mathbf{w}_{\mathrm{LS}} = \arg\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{y}_{\mathbf{w}}\|^2$, be the least square estimate of $\mathbf{y}$. Provide an expression for $\mathbf{w}_{\mathrm{LS}}$ and the associated $\hat{\mathbf{y}}_{\mathrm{LS}} = A\mathbf{w}_{\mathrm{LS}}$. No calculation is required in this part.

4 marks  **(b)** Show that for any $\mathbf{w} \in \mathbb{R}^p$ the following identity holds:

$$\|\mathbf{y} - \hat{\mathbf{y}}_{\mathbf{w}}\|^2 = \|\mathbf{y} - \hat{\mathbf{y}}_{\mathrm{LS}}\|^2 + \|\hat{\mathbf{y}}_{\mathbf{w}} - \hat{\mathbf{y}}_{\mathrm{LS}}\|^2$$

No calculation is needed for this part, just a clear geometric argument and an accompanying figure that uses properties of $\hat{\mathbf{y}}_{\mathrm{LS}}$ is sufficient.

6 marks  **(c)** **For parts (c) and (d) of this question assume that $\mathbf{A}^T\mathbf{A}$ is a diagonal matrix and $\lambda > 0$ is a constant.** Let $a_1, a_2, \ldots, a_p$ be the elements on the diagonal of $\mathbf{A}^T\mathbf{A}$, assumed to be non-zero. Let $\mathbf{w}^\star = \arg\min_{\mathbf{w}} \left\{ \|\mathbf{y} - \hat{\mathbf{y}}_{\mathbf{w}}\|^2 + \lambda\|\mathbf{w}\|^2 \right\}$. Using part (b) show that

$$w_i^\star = \arg\min_w \left\{ a_i (w - w_{\mathrm{LS},i})^2 + \lambda w^2 \right\}, \quad i = 1, 2, \ldots, p$$

where $w_i^\star$ and $w_{\mathrm{LS},i}$ denote the $i$-th elemenet of $\mathbf{w}^\star$ and $\mathbf{w}_{\mathrm{LS}}$ respectively. Hence express $w_i^\star$ in terms of $w_{\mathrm{LS},i}$, $a_i$ and $\lambda$.

4 marks

(**d**) Suppose we wish to compute:

$$\mathbf{w}^\star = \arg\min_{\mathbf{w}} \left\{ \|\mathbf{y} - \hat{\mathbf{y}}_{\mathbf{w}}\|^2 + \lambda \sum_{i=1}^{p} |w_i| \right\}$$

where $|\cdot|$ is the absolute value. Using the method analogous to part (c) express $w_i^\star$ (the $i$-th component of $\mathbf{w}^\star$) in terms of $w_{\mathrm{LS},i}$, $a_i$ and $\lambda$. *Hint: Consider the function $f(x) = (x - c)^2 + \eta|x|$, where $\eta > 0$. The minimizing value of $f(x)$, say $x_0$, can be expressed as follows: If $c > \eta/2$ then $x_0 = c - \eta/2$, if $c < -\eta/2$ then $x_0 = c + \eta/2$ and if $|c| \le \eta/2$ then $x_0 = 0$.*

4 marks

(**e**) Sketch $w_i^\star$ as a function of $w_{\mathrm{LS},i}$ for $\lambda = 1$ and $a_i = 1$ for both parts (c) and (d). Also explain qualitatively how does your answer in the two cases differ for large values of $\lambda$?

(**f**) Suppose that we have $N$ training examples: $(\mathbf{x}_1, y_1), \ldots (\mathbf{x}_N, y_N)$ where $\mathbf{x}_i \in \mathbb{R}^p$ and we wish to minimize the following loss function:

$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_i^T \mathbf{w} - y_i)^2 + \frac{\lambda}{N} \sum_{j=1}^{p} |w_j|,$$

where $\mathbf{w} = (w_1, \ldots, w_p)^T$. Provide a SGD update rule for minimizing $J(\mathbf{w})$. In each step assume that one example is selected at random in each update, and the gradient from the regularization term is added in each step.

**4.** Consider a logistic regression binary classification model that given $\mathbf{x} \in \mathbb{R}^2$ outputs:

$$P_{\mathbf{w}}(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(w_1 x_1 + w_2 x_2)}}, \qquad P_{\mathbf{w}}(y = -1|\mathbf{x}) = \frac{1}{1 + e^{(w_1 x_1 + w_2 x_2)}}.$$

Assume that the bias term in the model is set to zero for simplicity in this problem. Suppose that we have only two training points:

$$\mathbf{x}_1 = (1, 1), \quad y_1 = 1$$
$$\mathbf{x}_2 = (1, 0), \quad y_2 = -1$$

We intend to select $\mathbf{w}$ that minimizes the following regularized loss function:

$$J(\mathbf{w}) = -\sum_{i=1}^{2} \log P_{\mathbf{w}}(y_i|\mathbf{x}_i) + \lambda \|\mathbf{w}\|^2$$

**(a)** Suppose that $\lambda = 0$. What will the optimal choice of $\mathbf{w}$ be? What will the resulting value of $J(\mathbf{w})$ be?

**(b)** Suppose that $\lambda$ is a large constant such that it only suffices to consider $\|\mathbf{w}\| \ll 1$ when minimizing $J(\mathbf{w})$. In this case we can approximate

$$\log(1 + e^{-y_i \cdot \mathbf{w}^T \mathbf{x}_i}) \approx \log 2 - \frac{1}{2} y_i \cdot \mathbf{w}^T \mathbf{x}_i$$

Assuming that the above approximation is exact find $\mathbf{w}$ that minimizes $J(\mathbf{w})$.