## 1.1 Problem: Overfitting and Regularized Logisitc Regression (5 pts.)

1) (2 pts.) Plot the sigmoid function $1/(1 + \exp^{-wX})$ for increasing weights $w \in \{1, 5, 100\}$ with $X \in \mathbb{R}$. A qualitative sketch will suffice. Utilize these plots to explain why large weights can lead to overfitting in logistic regression.

2) (3 pts.) To mitigate overfitting, it is preferable to have smaller weights. To accomplish this, rather than utilizing maximum conditional likelihood estimation M(C)LE for logistic regression:
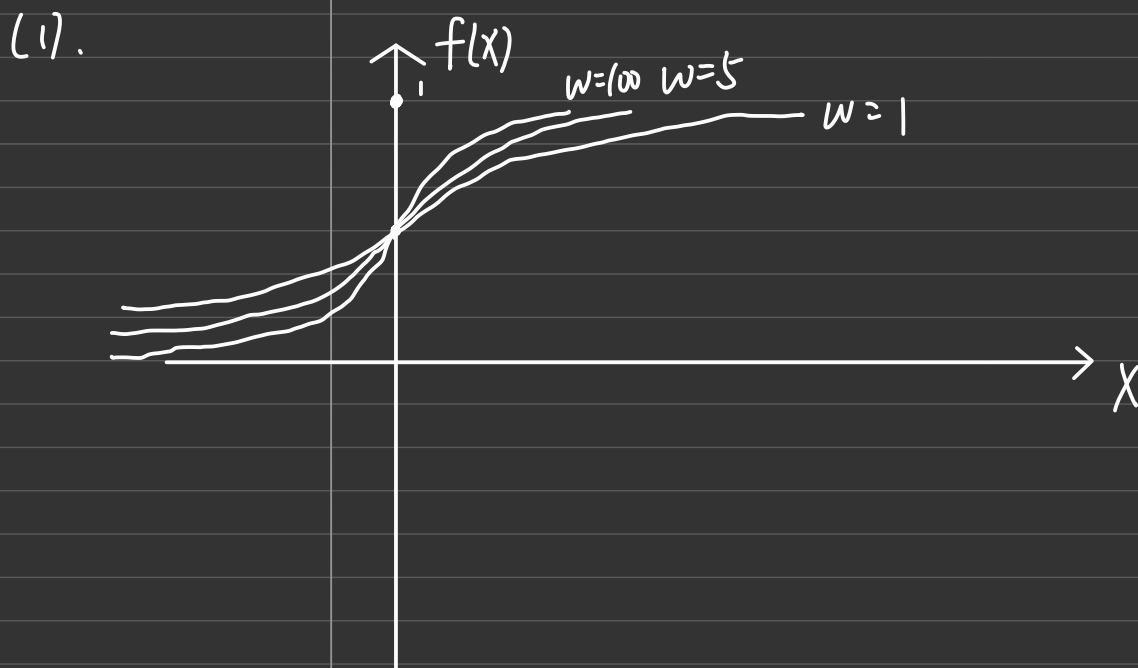
$$\max_{w_0, \cdots, w_d} \prod_{i=1}^{n} P(Y_i | X_i, w_0, \cdots, w_d), \tag{1}$$

we can consider maximum conditional a posterior M(C)AP estimation:

$$\max_{w_0, \cdots, w_d} \prod_{i=1}^{n} P(Y_i | X_i, w_0, \cdots, w_d) P(w_0, \cdots, w_d), \tag{2}$$

where $P(w_0, \cdots, w_d)$ is a prior on the weights.

Given a standard Gaussian prior $\mathcal{N}(0, \mathbf{I})$ for the weight vector $w$, please derive the **gradient ascent** update rules for the weights and explain why M(C)AP can address overfitting issue.

(1).



Larger weights makes $f(x)$ larger with same $X$, so the penalty is larger when a sample is misclassified. As a result, more samples are forced to be classified to the "right" class even though they are very close to another class. So it leads to overfitting.

(2).

$$\max_{w_0, \cdots, w_d} \prod_{i=1}^{n} P(Y_i | X_i, w_0, \cdots, w_d) P(w_0, \cdots, w_d)$$

$$\Longleftrightarrow \max_{w_0, \cdots, w_d} \sum_{i=1}^{n} \log(P(Y_i | X_i, \cdots, w_d)) + n \log P(w_0, \cdots, w_d)$$

$$\Longrightarrow \max_{w_0 \cdots w_d} \sum_{i=1}^{n} Y_i \log f_w(x_i) + (1 - Y_i) \log(1 - f_w(x_i)) - \frac{n}{2} w^T w$$

$\nabla J$ denotes gradient of $J(w) = y_i \log f_w(x_i) + (1-y_i) \log(1-f_w(x_i))$
$$- \frac{\eta}{2} w^T w$$

$w_{update} = w + \alpha \nabla J$, where $\alpha$ is learning rate.

As you can see, in $J(w)$ we get $-\frac{\eta}{2} w^T w$, since we need to maximize $J(w)$, so we would make $w$ small and thus to mitigate overfitting.

2. (1).  $P(y_c^i = 1 | X^i; W)$

$$= \frac{\exp(W_c^T X^i)}{\sum_{c'} \exp(W_{c'}^T X^i)}$$

$$= \frac{\exp((W_c - W_k)^T X^i)}{\sum_{c'} \exp((W_{c'} - W_k)^T X^i)} \qquad \left(\text{divide } \exp(W_k X^i)\right)$$

$$= \frac{\exp((W_{c'} - W_k)^T X^i)}{1 + \sum_{c'=1}^{k-1} \exp((W_{c'} - W_k)^T X^i)}$$

$$= \begin{cases} \dfrac{\exp(W_c^T X^i)}{1 + \sum_{c'=1}^{k-1} \exp(W_{c'}^T X^i)} & , \text{ if } c < k. \\[4mm] \dfrac{1}{1 + \sum_{c'=1}^{k-1} \exp(W_{c'}^T X^i)} & , \text{ if } c = k \end{cases} \qquad \left(\text{set } W_k^T = 0\right)$$

(2). $L(W) = \ln \prod_{j=1}^{n} P(y_c^j | X^j, W)$

$$= \ln \prod_{j=1}^{n} \sum_{c'} \frac{\exp(W_c^T X^j)}{\sum_{c'} \exp(W_{c'}^T X^j)} \cdot y_c \qquad \left(\begin{array}{l}\text{Only need the right class} \\ \text{probability. so multiply } y_c\end{array}\right)$$

$$= \sum_{j=1}^{n} \sum_{c=1}^{k} \left[ y_c^j (W_c^T X^j) - y_c^j \ln\left(\sum_{c'} \exp(W_{c'}^T X^j)\right)\right]$$

$$\left(y_c = 1 \text{ when } j = C, \text{ else } 0\right)$$

(3). $\nabla_{W_c} \sum_{c=1}^{k} y_c^j \ln\left(\sum_{c'} \exp(W_{c'}^T X^j)\right)$

$$= \left(\frac{\nabla_{W_c}\left(\exp(W_c^T X^j)\right)}{\sum_{c'} \exp(W_c^T X^j)}\right)$$

$$= X^j \frac{\exp(W_{c'}^T X^j)}{\sum_{c'} \exp(W_{c'}^T X^j)} \qquad \left(\text{consider fixed } j, \text{ most values } 0\right)$$

(4). $\nabla_{W_c} L(w) = \sum_{j=1}^{n} y_c^j x^j + \sum_{j=1}^{n} x^j \dfrac{\exp(W_c^T x^j)}{\sum_{c'} \exp(W_c^T x^j)}$

$= \sum_{j=1}^{n} x^j (y_c^j - P(y_c^j = 1 \mid x^j; w)$

(5).

$W_{\text{update}} = W + \eta \cdot \nabla_{W_c} L(w)$

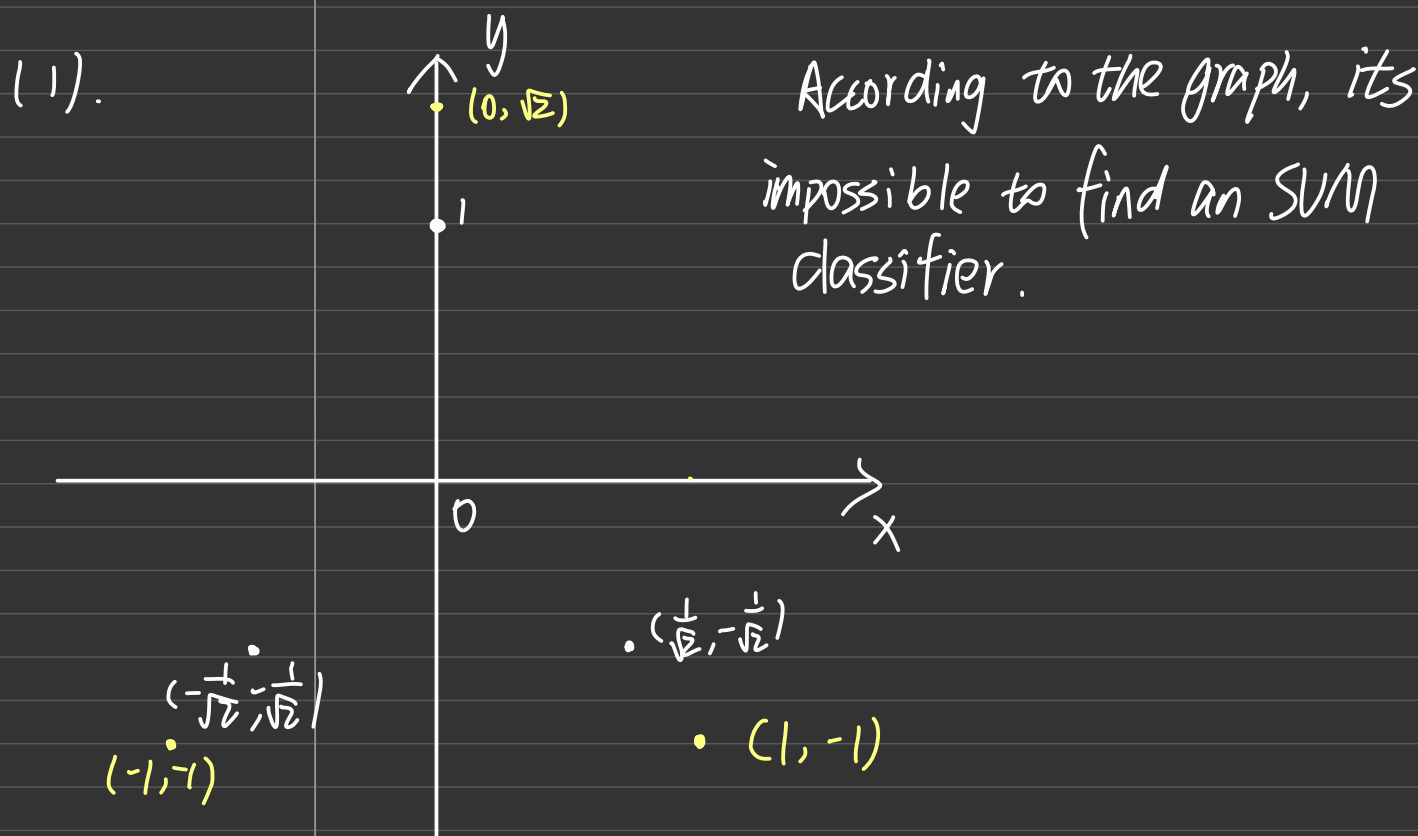$= W + \eta \cdot \sum_{j=1}^{n} x^j (y_c^j - P(y_c^j = 1 \mid x^j; w))$

Given a binary data set:
Class -1: $\{(0,1), (\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}), (-\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}})\}$. Class +1: $\{(0, \sqrt{2}), (1, -1), (-1, -1)\}$.
1) (3 pts.)  Could you find an SVM classifier (without slack variable) for this dataset?  Please explain your reasoning, possibly with the help of a plot sketch.
2) (3 pts.)  Use SVM by expanding the original feature vector $x = [x_1, x_2]$ to $x = [x_1^2, x_2^2]$, find the svm of this given data set and predict the label of $(-\frac{1}{2}, \sqrt{2})$.

(1).

According to the graph, its impossible to find an SVM classifier.



(2). Class $-1 : \{(0,1), (\frac{1}{2}, \frac{1}{2})  (\frac{1}{2}, \frac{1}{2})\}$
Class $1 : \{(0,2), (1,1), (1,1)\}$



The optimization problem is
$$\min \frac{1}{2}\|w\|^2$$
$$\text{s.t.} \quad (w_1 + b) \cdot (-1) \geqslant 1$$
$$(\frac{1}{2}w_0 + \frac{1}{2}w_1 + b) \cdot (-1) \geqslant 1$$
$$(2w_1 + b) \cdot 1 \geqslant 1$$
$$(w_0 + w_1 + b) \cdot 1 \geqslant 1$$

$$L(w, b, \alpha) = \frac{1}{2}(w_1^2 + w_0^2) + \alpha_1(1 + w_1 + b) + \alpha_2(\frac{1}{2}w_0 + \frac{1}{2}w_1 + b)$$
$$+ \alpha_3(1 - 2w_1 - b) + \alpha_4(1 - w_0 - w_1 - b)$$

$$\frac{\partial L}{\partial W} = \begin{pmatrix} W_0 + \frac{1}{2}d_2 - d_4 \\ W_1 + d_1 + \frac{1}{2}d_2 - 2d_3 - d_4 \end{pmatrix} = 0$$

$$\frac{\partial L}{\partial b} = d_1 + d_2 - d_3 - d_4 = 0$$

$d_i \geq 0$, for $i = 1, 2, 3, 4$

$$\left\{ \begin{array}{l} 1 + W_1 + b \leq 0 \\ \frac{1}{2}W_0 + \frac{1}{2}W_1 + b + 1 \leq 0 \\ 1 - 2W_1 - b \geq 0 \\ 1 - W_0 - W_1 - b \geq 0 \end{array} \right.$$

$d_1 (1 + W_1 + b) = 0$

$d_2 (\frac{1}{2}W_0 + \frac{1}{2}W_1 + b + 1) = 0$

$d_3 (1 - 2W_1 - b) = 0$

$d_4 (1 - W_0 - W_1 - b) = 0$

$$\Rightarrow \left\{ \begin{array}{l} W_0 = 2 \\ W_1 = 2 \\ b = -3 \end{array} \right.$$

} KKT Condition

$\frac{1}{2}$ tp.

The lable of $(-\frac{1}{2}, \sqrt{2})$ should be $+1$.

1.4  (1),   $\min\limits_{W,s} \frac{1}{2}\|W\|^2 + C\sum\limits_{i=1}^{n}(S_i + t_i)$

s.t.
$$y_i - f(X_i) - \varepsilon \leq S_i \quad, i=1, 2, 3\cdots, n$$
$$f(X_i) - y_i - \varepsilon \leq t_i \quad, i=1, 2, 3, \cdots, n$$
$$S_i \geq 0, t_i \geq 0, \quad i=1, 2, \cdots, n$$

Where  $f(X) = WX$

(2)  $L(W, t, s, \partial, \beta, r, \theta)$
$$= \frac{1}{2}\|W\|^2 + C\sum\limits_{i=1}^{n}S_i + \sum\limits_{i=1}^{n}(y_i - f(X_i) - \varepsilon - S_i)\partial_i + \beta_i(y_i + f(X_i) - \varepsilon - S_i)$$
$$- \sum\limits_{i=1}^{n}S_i r_i - \sum\limits_{i=1}^{n}\theta_i t_i$$

(3).

$g(\partial, \beta, r, \theta) = \min\limits_{X \in R^n} L(\omega, b, s, \partial, \beta, r, \theta)$

$\frac{\partial g}{\partial W} = 0 \Rightarrow \quad W = \sum\limits_{i=1}^{n}(\partial_i - \beta_i)X_i$

and       $\partial_i(\varepsilon + S_i - y_i + f(X_i)) = 0$

$\beta_i(\varepsilon + t_i + y_i - f(X_i)) = 0$

$0 \leq \partial_i \quad, i=1, 2, \cdots, n$
$0 \leq \beta_i \quad, i=1, 2, \cdots, n$

so dual  form  is

$\max g(\partial, \beta, r, \theta)$
$\Updownarrow$

$\min \frac{1}{2}\sum\limits_{i=1}^{n}\sum\limits_{j=1}^{n}(\partial_i - \beta_i)(\partial_j - \beta_j)X_j^T X_i$
$+ \sum\limits_{i=1}^{n}((\varepsilon - y_i)\cdot\partial_i + (\varepsilon + y_i)\cdot\beta_i)$

subject to  $\partial_i \geq 0, \beta_i \geq 0$

$i=1, 2, \cdots, n$

(4). Yes.

(5). The support vectors are points lie on the margin or in the $\varepsilon$ boundary around the trained decision boundary.

(6). $g(x) = \sum_{i=1}^{n} (d_i - \frac{\beta}{2}) x_i x$

(7). Yes. substitue $X$ with $k(X)$.

(8). The decision area will get larger, which means the error that can be toleranted gets larger if $\varepsilon$ is smaller, and the model can be more complex and easier to be overfitting.

If $\varepsilon$ is larger, the model gets less complex, more values can be fitted in the area without large penalty. However, it may leads to loss of precision.

(9). If $C$ is smaller, it may result in a larger margin and a simpler model, and the model may be more generalized.

If $C$ is larger, it may result in a smaller margin and a complex model, which means it provides higher precision but may lead to overfitting.

1.5 First iteration:

$$H(D|A) = \sum_{i=1}^{n} \frac{|D_i|}{D} H(D_i)$$

$$= \sum_{i=1}^{n} \frac{|D_i|}{|D|} \sum_{k=1}^{k} \frac{|D_{ik}|}{|D_i|} \log_2 \frac{|D_{ik}|}{|D_i|}$$

$$H(D) = -\frac{6}{15} \log_2 \frac{6}{15} - \frac{9}{15} \log_2 \frac{9}{15}$$

$$g_R(D,A) = \frac{H(D) - H(D|A)}{H_A(D)}$$

Age: $H(D|A) = \frac{1}{3} \cdot \frac{3}{5} \log_2 \frac{2}{5} - \frac{1}{3} \cdot \frac{2}{5} \log_2 \frac{2}{5} - \frac{1}{3} \log_2 \frac{3}{5} - \frac{1}{3} \log_2 \frac{2}{5} -$

$$\frac{1}{3} \cdot \frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{3} \cdot \frac{1}{5} \log_2 \frac{1}{5}$$

$$= -\frac{2}{5} \log_2 \frac{3}{5} - \frac{4}{15} \log_2 \frac{2}{5} - \frac{4}{15} \log_2 \frac{4}{5} - \frac{1}{15} \log_2 \frac{1}{5}$$

$$H_A(D) = -\frac{3}{5} \log_2 \frac{1}{5}$$

$$g_R(D, Age) = \frac{H(D) - H(D|A)}{H_A(D)} = 0.05$$

$$g_R(D, work) = 0.35$$

$$g_R(D, house) = 0.43$$

$$g_R(D, credit) = 0.21$$

$$A_g = house$$

$$D_1 = \{1, 2, 5, 6, 7, 13, 14, 15\}$$

$$D_2 = \{4, 8, 9, 10, 11, 12\} \qquad A = \{Age, Work, Credit\}$$

For $D_1$. $g_R(D_1, Age) = 0.16$

$$g_R(D, work) = 1$$

so work is chosen

$$D_1 \longrightarrow \begin{cases} D_3 \{1, 2, 5, 6, 7, 15\} \\ D_4 \{3, 13, 14\} \end{cases}$$

For $D_2$: belongs to Yes,
   $D_3$: belongs to No
   $D_4$: belongs to Yes. Done.

1.2, 3, 4, 5, 6, 7, 8, 9, 10,
11, 12, 13, 14, 15          $D$
house

$D_1$
1, 2, 5, 6, 7, 15

$D_2$: Yes
4, 8, 9, 10, 11, 12

Work

$D_3$: No
1, 2, 5, 6, 7, 15

$D_4$: Yes
3, 13, 14