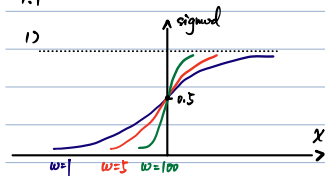


Written Problems

1.1

1)



As the weight increases, the slope near $x=0$ become steeper. This indicates an intensive change in the sigmoid when x changes. Hence there is high possibility that the model would regard the noise as critic characteristics, leading to overfitting in many cases.

2) For logistic regression, $P(Y_i | X_i, w_0, w_1, \dots, w_n) = \sigma_i^{y_i} (1 - \sigma_i)^{1-y_i}$ where $\sigma_i = (1 + e^{-w^T X_i})^{-1}$

For $w \sim \mathcal{N}(0, I)$, $P(w_0, w_1, \dots, w_n) = (2\pi)^{-\frac{n}{2}} \exp(-\frac{1}{2} \sum w_j^2)$

Hence the log-likelihood $\ln \prod_{i=1}^n P(Y_i | X_i, w) P(w)$

$$= \ln \prod_{i=1}^n \left(\frac{1}{1 + e^{-w^T X_i}} \right)^{y_i} \left(1 - \frac{1}{1 + e^{-w^T X_i}} \right)^{1-y_i} (2\pi)^{-\frac{n}{2}} \exp(-\frac{1}{2} \sum w_j^2)$$

$$= \sum_{i=1}^n y_i \ln \left(\frac{1}{1 + e^{-w^T X_i}} \right) + (1 - y_i) \ln \left(1 - \frac{1}{1 + e^{-w^T X_i}} \right) - \frac{n}{2} \ln(2\pi) - \frac{1}{2} \sum w_j^2$$

$$= -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \sum w_j^2 + \sum_{i=1}^n y_i \ln \left(\frac{1}{1 + e^{-w^T X_i}} \right) + (1 - y_i) \ln \left(1 - \frac{1}{1 + e^{-w^T X_i}} \right)$$

Drop those irrelevant with the weights. $J(w) = -\frac{1}{n} \sum_{i=1}^n y_i \ln \left(\frac{1}{1 + e^{-w^T X_i}} \right) - (1 - y_i) \ln \left(1 - \frac{1}{1 + e^{-w^T X_i}} \right)$

The gradient $\nabla_w J(w) = \frac{1}{n} \sum_{i=1}^n (y_i - \frac{1}{1 + e^{-w^T X_i}}) X_i - w$

Hence we have the update rule:

(1) Initialize the vector w appropriately;

(2) Update w s.t. $w := w + \alpha \nabla_w J(w)$, where $\nabla_w J(w) = \frac{1}{n} \sum_{i=1}^n (y_i - \frac{1}{1 + e^{-w^T X_i}}) X_i - w$

(3) If $\|\nabla_w J(w)\| \leq \epsilon$ or reaching the maximum # of iterations, stop.

Otherwise, execute (2) and (3) repeatedly.

(1C) AP addresses overfitting by introducing a penalty term as a prior, working as a regularization to the term of weights.

1.2

$$1) P(Y_i = 1 | X_i; w) = \frac{\exp(w^T X_i)}{\sum_{c=1}^2 \exp(w^T X_i^c)} = \frac{\exp(w^T X_i^1 - w^T X_i^2) \exp(w^T X_i^2)}{\sum_{c=1}^2 \exp(w^T X_i^c - w^T X_i^2) \exp(w^T X_i^2)}$$

$$= \frac{\exp(w^T X_i^1 - w^T X_i^2)}{\sum_{c=1}^2 \exp(w^T X_i^c - w^T X_i^2)} = \frac{\exp(w^T X_i^1 - w^T X_i^2)}{\sum_{c=1}^2 \exp(w^T X_i^c - w^T X_i^2) + 1}$$

In the alternative expression, $w_c = w_1 - w_2$, $\forall c=1, 2$

$$P(Y_i = 1 | X_i; w) = \frac{\exp(w^T X_i^1)}{\sum_{c=1}^2 \exp(w^T X_i^c) + 1}$$

$$P(Y_i = 2 | X_i; w) = 1 - \frac{\exp(w^T X_i^1)}{\sum_{c=1}^2 \exp(w^T X_i^c) + 1} = \frac{1}{\sum_{c=1}^2 \exp(w^T X_i^c) + 1}$$

$$2) L(w) = \ln \prod_{j=1}^n P(Y_j^1 | X_j^1, w)$$

$$= \sum_{j=1}^n \ln \frac{\exp(w^T X_j^1)}{\sum_{c=1}^2 \exp(w^T X_j^c) + 1}$$

with only the one with a true label $Y_j^1 = 1$.

$$= \sum_{j=1}^n \ln \frac{\exp(w^T X_j^1)}{\sum_{c=1}^2 \exp(w^T X_j^c) + 1}$$

hence the logarithm could move in.

$$= \sum_{j=1}^n \ln \frac{\exp(w^T X_j^1)}{\sum_{c=1}^2 \exp(w^T X_j^c) + 1}$$

$$= \sum_{j=1}^n \ln \frac{\exp(w^T X_j^1)}{\sum_{c=1}^2 \exp(w^T X_j^c) + 1}$$

$$= \sum_{j=1}^n \ln \frac{\exp(w^T X_j^1)}{\sum_{c=1}^2 \exp(w^T X_j^c) + 1}$$

$$3) \nabla_w L(w) = \frac{\partial \sum_{j=1}^n \ln \frac{\exp(w^T X_j^1)}{\sum_{c=1}^2 \exp(w^T X_j^c) + 1}}{\partial w} = \frac{\partial \ln \left(\sum_{c=1}^2 \exp(w^T X_j^c) + 1 \right)}{\partial w} \cdot \frac{\partial \exp(w^T X_j^1)}{\partial w}$$

$$= \frac{X_j \exp(w^T X_j^1)}{\sum_{c=1}^2 \exp(w^T X_j^c) + 1}$$

$$4) \frac{\partial \sum_{j=1}^n y_j^1 w^T X_j^1}{\partial w} = \sum_{j=1}^n X_j^1 y_j^1, \quad y_j^1 = 1 \text{ only for the correct one.}$$

Combining with $\nabla_w L(w)$ yields $\nabla_w L(w) = \sum_{j=1}^n X_j (y_j^1 - P(Y_j^1 = 1 | X_j^1; w))$

5) Hence we have the update rule:

(1) Initialize the vector w appropriately;

(2) Update w s.t. $w := w + \alpha \nabla_w L(w)$, where $\nabla_w L(w) = \sum_{j=1}^n X_j (y_j^1 - P(Y_j^1 = 1 | X_j^1; w))$

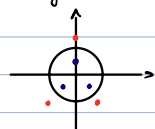
(3) If $\|\nabla_w L(w)\| \leq \epsilon$ or reaching the maximum # of iterations, stop.

Otherwise, execute (2) and (3) repeatedly.

1.3

1) Recall $\|x\|_2 = \sqrt{x_1^2 + x_2^2}$, and in this case $\|x_1\|_2 = 1$, $\|x_2\|_2 = \sqrt{2}$

To find the largest possible geometric margin, it is intuitive to construct a circle right in the middle.



=> The appropriate SVM classifier could be $\|x\|_2 = \frac{1+\sqrt{2}}{2}$

2) In the $[x_1^2, x_2^2]$ space, we have $\begin{cases} \text{Class -1: } \{(0,1), (\frac{1}{2}, \frac{1}{2})\} \\ \text{Class +1: } \{(1,2), (1,1)\} \end{cases}$

$$\min \frac{1}{2} \|w\|^2 \text{ s.t. } \forall i, y_i (w^T X_i + b) \geq 1$$

By KKT, $L(w, b, \alpha) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^m \alpha_i [1 - y_i (w^T X_i + b)]$

$$\frac{\partial L}{\partial w} = 0, \quad w = \sum_{i=1}^m \alpha_i y_i X_i$$

$$\frac{\partial L}{\partial b} = 0, \quad \sum_{i=1}^m \alpha_i y_i = 0$$

$$\alpha_i \geq 0, \quad 1 - y_i (w^T X_i + b) \leq 0$$

$$\alpha_i (1 - y_i (w^T X_i + b)) = 0$$

Hence the dual: $\max \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j X_i^T X_j$ s.t. $\forall i, \sum_{i=1}^m \alpha_i y_i = 0, \alpha_i \geq 0$

i.e. In this setting, $\max f(\alpha_1, \alpha_2, \alpha_3, \alpha_4)$ s.t. $-\alpha_1 - \alpha_2 + \alpha_3 + \alpha_4 = 0$

$$\text{where } f(\alpha_1, \alpha_2, \alpha_3, \alpha_4) = \alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 - \frac{1}{2} \alpha_1^2 - \frac{1}{2} \alpha_2^2 - 2\alpha_3^2 - \alpha_4^2 - \frac{1}{2} \alpha_1 \alpha_2 + 2\alpha_1 \alpha_3 + \alpha_1 \alpha_4 + \alpha_2 \alpha_3 + \alpha_2 \alpha_4 - 2\alpha_3 \alpha_4$$

$$\text{Thus } w = \begin{bmatrix} 2 \\ 3 \end{bmatrix}, \quad b = -3$$

$$\text{For } (-\frac{1}{2}, \sqrt{2}), \quad 2 \times (-\frac{1}{2})^2 + 2 \times (\sqrt{2})^2 - 3 = \frac{3}{2} > 0$$

=> It is predicted as class +1.

1.4

$$1) \min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max(0, |y_i - w^T X_i| - \epsilon)$$

Converted into:

$$\min \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i$$

$$\text{s.t. } |y_i - w^T X_i| - \epsilon \leq \xi_i$$

$$\xi_i \geq 0$$

Further converted into:

$$\min \frac{1}{2} w^T w + C \sum_{i=1}^n (\xi_i^+ + \xi_i^-)$$

$$\text{s.t. } y_i - w^T X_i - \epsilon \leq \xi_i^+$$

$$w^T X_i - y_i - \epsilon \leq \xi_i^-$$

$$\xi_i^+, \xi_i^- \geq 0$$

2) Define the lagrangian $L(w, \xi^+, \xi^-, \alpha^+, \alpha^-, \mu^+, \mu^-)$

$$= \frac{1}{2} w^T w + C \sum_{i=1}^n (\xi_i^+ + \xi_i^-) + \sum_i \alpha_i^+ (y_i - w^T X_i - \epsilon - \xi_i^+)$$

$$+ \sum_i \alpha_i^- (w^T X_i - y_i - \epsilon - \xi_i^-) - \sum_i \mu_i^+ \xi_i^+ - \sum_i \mu_i^- \xi_i^-$$

3) By KKT, we have conditions as follows:

$$\text{Stationary: } \frac{\partial L}{\partial w} = w - \sum_i \alpha_i^+ X_i + \sum_i \alpha_i^- X_i = 0$$

$$\frac{\partial L}{\partial \xi_i^+} = C - \alpha_i^+ - \mu_i^+ = 0 \quad \frac{\partial L}{\partial \xi_i^-} = C - \alpha_i^- - \mu_i^- = 0$$

$$\text{Feasibility: } \alpha_i^+, \alpha_i^- \geq 0, \quad \xi_i^+, \xi_i^- \geq 0, \quad \mu_i^+, \mu_i^- \geq 0$$

$$\text{Complementary: } \alpha_i^+ (y_i - w^T X_i - \epsilon - \xi_i^+) = 0 \quad \alpha_i^- (w^T X_i - y_i - \epsilon - \xi_i^-) = 0$$

$$\mu_i^+ \xi_i^+ = 0 \quad \mu_i^- \xi_i^- = 0$$

$$\text{Hence } L(w) = \frac{1}{2} w^T w + C \sum_{i=1}^n (\xi_i^+ + \xi_i^-) + \sum_i \alpha_i^+ (y_i - w^T X_i - \epsilon - \xi_i^+)$$

$$+ \sum_i \alpha_i^- (w^T X_i - y_i - \epsilon - \xi_i^-) - \sum_i \mu_i^+ \xi_i^+ - \sum_i \mu_i^- \xi_i^-$$

$$= \frac{1}{2} w^T w + C \sum_{i=1}^n (\xi_i^+ + \xi_i^-) + \sum_i (\alpha_i^+ - \alpha_i^-) y_i - \sum_i (\alpha_i^+ - \alpha_i^-) \epsilon - \sum_i (\alpha_i^+ - \alpha_i^-) (y_i - \epsilon)$$

$$- \sum_i (\alpha_i^+ - \alpha_i^-) \epsilon - \sum_i \alpha_i^+ \xi_i^+ - \sum_i \alpha_i^- \xi_i^- - \sum_i \mu_i^+ \xi_i^+ - \sum_i \mu_i^- \xi_i^-$$

$$= -\frac{1}{2} \sum_{i,j} (\alpha_i^+ - \alpha_i^-) (\alpha_j^+ - \alpha_j^-) X_i^T X_j + \sum_i [(y_i - \epsilon) \alpha_i^+ - (y_i + \epsilon) \alpha_i^-]$$

The dual could be written as

$$\max -\frac{1}{2} \sum_{i,j} (\alpha_i^+ - \alpha_i^-) (\alpha_j^+ - \alpha_j^-) X_i^T X_j + \sum_i [(y_i - \epsilon) \alpha_i^+ - (y_i + \epsilon) \alpha_i^-]$$

$$\text{s.t. } 0 \leq \alpha_i^+ \leq C$$

$$0 \leq \alpha_i^- \leq C$$

4) Yes, obviously.

5) $(\alpha^+ = 0, \alpha^- \neq 0)$ or $(\alpha^+ \neq 0, \alpha^- = 0)$.

6) The prediction should be $\hat{y} = w^T X_i$, w as the trained weights.

7) Yes. replacing this $X_i^T X_j$ by $\phi(x_i, x_j)$.

8) $\epsilon \downarrow \Rightarrow$ Width of boundary $\downarrow \Rightarrow$ Cost \uparrow

$\epsilon \uparrow \Rightarrow$ Width of boundary $\uparrow \Rightarrow$ Lose some feature but cost \downarrow

9) $C \uparrow \Rightarrow$ More penalty on outliers instead of regularization \Rightarrow More inclusive.

$C \downarrow \Rightarrow$ More penalty on regularization \Rightarrow The hyperplane is flatter.

1.5

$$H(D) = -\frac{1}{15} \log_2 \frac{1}{15} - \frac{9}{15} \log_2 \frac{9}{15} = 0.971$$

$$\text{Age: } H(D| \text{Age} = \text{young}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

$$H(D| \text{Age} = \text{middle}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

$$H(D| \text{Age} = \text{old}) = -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} = 0.722$$

$$g(D, A) = H(D) - H(D|A) = 0.971 - 0.971 \times \frac{1}{3} - 0.971 \times \frac{1}{3} - 0.722 \times \frac{1}{3} = 0.083$$

$$H_A(D) = -\frac{1}{15} \log_2 \frac{1}{15} - \frac{5}{15} \log_2 \frac{5}{15} - \frac{5}{15} \log_2 \frac{5}{15} = 1.585$$

$$g_R(D, A) = \frac{g(D, A)}{H_A(D)} = \frac{0.083}{1.585} = 0.052$$

$$\text{Work: } H(D| \text{Work} = \text{yes}) = -\log_2 1 = 0$$

$$H(D| \text{Work} = \text{no}) = -\frac{6}{10} \log_2 \frac{6}{10} - \frac{4}{10} \log_2 \frac{4}{10} = 0.971$$

$$g(D, A) = H(D) - H(D|A) = 0.971 - \frac{2}{3} \times 0.971 = 0.324$$

$$H_A(D) = -\frac{10}{15} \log_2 \frac{10}{15} - \frac{5}{15} \log_2 \frac{5}{15} = 0.918$$

$$g_R(D, A) = \frac{g(D, A)}{H_A(D)} = 0.353$$

Similarly, for "House" $g_R(D, A) = 0.433$; for "Credit" $g_R(D, A) = 0.232$.

Hence we first divide it using "House". The ones with "House" labelled as yes all, so we continue to divide the ones with "House" labelled as no.

In the next stage, we similarly have $H(D) = 0.918$

The "Age" $g_R(D, A) = 0.165$, the "Work" $g_R(D, A) = 1$, and the "Credit" $g_R(D, A) = 0.341$

Hence we next divide it using "Work".

As a result,

