

Homework 3

2021年10月9日 星期六 下午3:06

Problem 1

1. First, we plug in $y=0$ in $P[y|x; w]$:

$$\begin{aligned} P[y|x; w] &= P[y=1|x; w]^0 \cdot P[y=0|x; w]^1 \\ &= P[y=0|x; w] \end{aligned}$$

Then, we plug in $y=1$ in $P[y|x; w]$:

$$\begin{aligned} P[y|x; w] &= P[y=1|x; w]^1 \cdot P[y=0|x; w]^0 \\ &= P[y=1|x; w] \end{aligned}$$

So $P[y|x; w] = P[y=1|x; w]^y \cdot P[y=0|x; w]^{1-y}$

2.

$$\begin{aligned} \log \left(\prod_{i=1}^n P[y_i|x_i; w] \right) &= \sum_{i=1}^n \log P[y_i|x_i; w] \\ &= \sum_{i=1}^n \log \left[\frac{1}{1+e^{-w^T x_i}} \right]^{y_i} \cdot \left(\frac{1}{1+e^{-w^T x_i}} \right)^{1-y_i} \\ &= \sum_{i=1}^n y_i \log \left[\frac{1}{1+e^{-w^T x_i}} \right] + (1-y_i) \log \left[\frac{1}{1+e^{-w^T x_i}} \right] \\ &= -\sum_{i=1}^n y_i \log (1+e^{-w^T x_i}) + (1-y_i) \log (1+e^{-w^T x_i}) \end{aligned}$$

$$\arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n y_i \log (1+e^{-w^T x_i}) + (1-y_i) \log (1+e^{-w^T x_i})$$

3.

$$\begin{aligned} w_{t+1} &= w_t - \eta_t \left(\sum_{i=1}^n (w_t^T x_i - y_i) x_i + 2\lambda w_t \right) \\ \frac{\partial}{\partial w} &= \sum_{i=1}^n y_i \frac{-x_i e^{-w^T x_i}}{1+e^{-w^T x_i}} + \sum_{i=1}^n (1-y_i) \frac{x_i e^{w^T x_i}}{1+e^{w^T x_i}} \\ &= \sum_{i=1}^n y_i \frac{-x_i (e^{-w^T x_i} + 1 - 1)}{1+e^{-w^T x_i}} + \sum_{i=1}^n (1-y_i) \frac{x_i (e^{w^T x_i} + 1 - 1)}{1+e^{w^T x_i}} \\ &= \sum_{i=1}^n -x_i y_i (1 - \sigma(w^T x_i)) + \sum_{i=1}^n (1-y_i) x_i (1 - \sigma(w^T x_i)) \\ &= \sum_{i=1}^n -x_i y_i \sigma(-w^T x_i) + \sum_{i=1}^n -x_i y_i \sigma(w^T x_i) + x_i \sigma(w^T x_i) \\ &= \sum_{i=1}^n -x_i y_i \sigma(-w^T x_i) - x_i y_i \sigma(w^T x_i) + x_i \sigma(w^T x_i) \\ &= \sum_{i=1}^n -x_i y_i [\sigma(-w^T x_i) + \sigma(w^T x_i)] + x_i \sigma(w^T x_i) \\ &= \sum_{i=1}^n -x_i y_i \cdot (1 + x_i \sigma(w^T x_i)) \\ &= \sum_{i=1}^n [x_i (\sigma(w^T x_i) - y_i)] \end{aligned}$$

correctly classified:

We know that when $y_i = 1$, the confidence should be $\sigma(w^T x_i)$. So the difference between $\sigma(w^T x_i)$ and y_i will be small. That means the weight of the x_i is small, also the summation will be small. Therefore, the contribution of the correctly classified example to updated solution in every iteration of GD misclassified example:

To the contrary, if that example is misclassified, the confidence should be far from $\sigma(w^T x_i)$. So the difference will be significant. And the summation will be large in the end. It will lead bigger contribution than correctly classified example.

Problem 2

1. Attributes and their values:

Color: Yellow, Green

Size: Small, Large

Shape: Round, Irregular

$$\begin{aligned} H(Y|X) &= \sum_{x \in \text{values}(X)} p(x) H(Y|x) \\ H(X) &= -\sum_{x \in \text{values}(X)} p(x) \log_2(p(x)) \\ &= -[\log_2 p + (1-p) \log_2(1-p)] \end{aligned}$$

Target label: {Yes, No}

splitting for mushroom data

$$H(Y) = -\left[\frac{9}{16} \log_2 \left(\frac{9}{16}\right) + \frac{7}{16} \log_2 \left(\frac{7}{16}\right)\right] = 0.989$$

$$H(Y|Color) = -\left[\frac{3}{16} \log_2 \left(\frac{3}{16}\right) + \frac{13}{16} \log_2 \left(\frac{13}{16}\right)\right] = 0.961$$

$$H(Y|Green) = -\left[\frac{1}{8} \log_2 \left(\frac{1}{8}\right) + \frac{3}{8} \log_2 \left(\frac{3}{8}\right)\right] = 0.918$$

$$InfoGain(Color) = H(Y) - H(Y|Color)$$

$$= 0.989 - \left[\frac{13}{16} H(Y|Color) + \frac{3}{16} H(Y|Green)\right]$$

$$= 0.989 - \left[\frac{13}{16} (0.961) + \frac{3}{16} (0.918)\right]$$

$$= 0.036$$

$$H(Y|Size) = -\left[\frac{6}{16} \log_2 \left(\frac{6}{16}\right) + \frac{10}{16} \log_2 \left(\frac{10}{16}\right)\right] = 0.811$$

$$H(Y|Large) = -\left[\frac{3}{8} \log_2 \left(\frac{3}{8}\right) + \frac{5}{8} \log_2 \left(\frac{5}{8}\right)\right] = 0.954$$

$$InfoGain(Size) = H(Y) - H(Y|Size)$$

$$= 0.989 - \left[\frac{8}{16} (0.811) + \frac{8}{16} (0.954)\right]$$

$$= 0.107 \quad \textcircled{1}$$

$$H(Y|Shape) = -\left[\frac{6}{16} \log_2 \left(\frac{6}{16}\right) + \frac{10}{16} \log_2 \left(\frac{10}{16}\right)\right] = 0.811$$

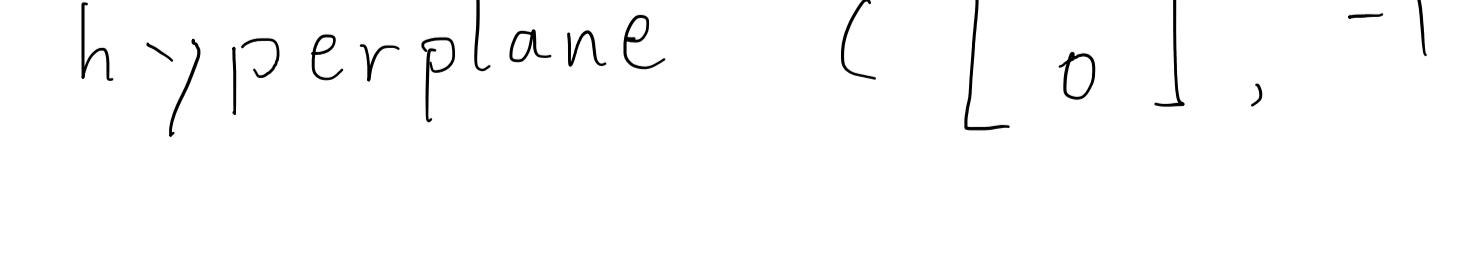
$$H(Y|Irregular) = -\left[\frac{3}{8} \log_2 \left(\frac{3}{8}\right) + \frac{5}{8} \log_2 \left(\frac{5}{8}\right)\right] = 0.954$$

$$InfoGain(Shape) = H(Y) - H(Y|Shape)$$

$$= 0.989 - \left[\frac{12}{16} (0.811) + \frac{4}{16} (0.954)\right] = 0.036$$

"Size" would the algorithm choose to use for the root of tree

2.



Problem 3

There is a possibility of overfitting. During the design of decision trees, overfitting occurs when all samples from the training set are perfectly fitted into the tree. Therefore, it should branches with sparse data.

Problem 4

$$\vec{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$\begin{cases} (-1)(w_1(0) + w_2(0) + b) \geq 1 \\ (-1)(w_1(0) + w_2(-1) + b) \geq 1 \\ (+1)(w_1(-2) + w_2(0) + b) \geq 1 \end{cases}$$

$$\begin{cases} -b \geq 1 & \textcircled{1} \\ w_2 - b \geq 1 & \textcircled{2} \\ -2w_1 + b \geq 1 & \textcircled{3} \end{cases}$$

$$\textcircled{1} + \textcircled{3} \longrightarrow -2w_1 \geq 2 \quad w_1 \leq -1$$

$$\textcircled{2} + \textcircled{3} \longrightarrow w_2 - b \geq 2 \quad -2w_1 + w_2 \geq 2$$

$$To \quad \min_{w, b} \frac{1}{2} \|w\|_2^2, \quad w_1 = -1, w_2 = 0$$

$$\begin{cases} -b \geq 1 \\ 0 - b \geq 1 \\ 2 - b \geq 1 \end{cases} \quad \begin{cases} -b \geq 1 \\ 2 + b \geq 1 \\ b \leq -1 \end{cases} \quad \begin{cases} b \leq -1 \\ b \geq -1 \end{cases}$$

$$\therefore b_* = -1$$

$$\text{Optimal hyperplane } \left(\begin{bmatrix} -1 \\ 0 \end{bmatrix}, -1 \right)$$

margin:

$$w^T x + b = 0$$

$$(0, 0)$$

$$(-2, 0)$$

$$(0, -1)$$

$$x_1$$

$$x_2$$

$$y$$

$$d(x_i, w, b)$$

$$x_i \in S$$

$$= 1$$

$$margin(w, b, S) = \min_{x_i \in S} d(x_i, w, b)$$

$$= 1$$

$$w_1 = -1$$

$$w_2 = 0$$

$$b = -1$$

$$w^T x + b = 0$$

$$(-1, 0)$$

$$(0, 1)$$

$$x_1$$

$$x_2$$

$$y$$

$$d(x_i, w, b)$$

$$x_i \in S$$

$$= 1$$

$$margin(w, b, S) = \min_{x_i \in S} d(x_i, w, b)$$

$$= 1$$

$$w_1 = -1$$

$$w_2 = 0$$

$$b = -1$$

$$w^T x + b = 0$$

$$(-1, 0)$$

$$(0, 1)$$

$$x_1$$

$$x_2$$

$$y$$

$$d(x_i, w, b)$$

$$x_i \in S$$

$$= 1$$

$$margin(w, b, S) = \min_{x_i \in S} d(x_i, w, b)$$

$$= 1$$

$$w_1 = -1$$

$$w_2 = 0$$

$$b = -1$$

$$w^T x + b = 0$$

$$(-1, 0)$$

$$(0, 1)$$

$$x_1$$

$$x_2$$

$$y$$

$$d(x_i, w, b)$$

$$x_i \in S$$

$$= 1$$

$$margin(w, b, S) = \min_{x_i \in S} d(x_i, w, b)$$

$$= 1$$

$$w_1 = -1$$

$$w_2 = 0$$

$$b = -1$$

$$w^T x + b = 0$$

$$(-1, 0)$$

$$(0, 1)$$