

$$\begin{aligned} J(w) &= -\sum_{n=1}^N [y_n (\log(1 + e^{-w^T x_n})) + (1 - y_n) (\log(e^{-w^T x_n}) - \log(1 + e^{-w^T x_n}))] \\ &= -\sum_{n=1}^N [y_n \log(1 + e^{-w^T x_n}) - w^T x_n - \log(1 + e^{-w^T x_n}) + y_n w^T x_n] \\ &= \sum_{n=1}^N [-w^T x_n - \log(1 + e^{-w^T x_n}) + y_n w^T x_n] \end{aligned}$$

let $w^T = [w_1, w_2]$,

the prediction is $y_i = \text{sgn}(w^T x_i + b)$

$$w^T x + b = w_1 x_1 + w_2 x_2 + b$$

$$\begin{cases} w_1 + w_2 + b \geq 0 \\ w_1 - w_2 + b \leq 0 \end{cases}$$

$$\begin{cases} w_1 + w_2 + b \geq 0 \\ -w_1 + w_2 + b \leq 0 \end{cases}$$

$$\begin{cases} w_1 + w_2 + b \geq 0 \\ -w_1 - w_2 + b \leq 0 \end{cases}$$

∴ one solution is $\begin{cases} w^T = [1, 1] \\ b = -1 \end{cases}$

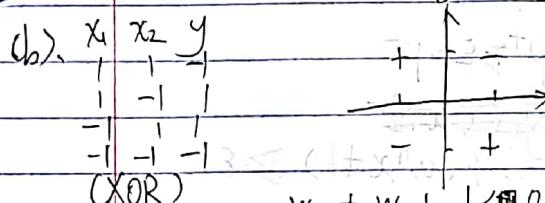
$$\frac{\partial J(w)}{\partial w_j} = \frac{\partial}{\partial w_j} \left[\sum_{n=1}^N [-x_{n,j} - \frac{x_{n,j} e^{-w^T x_n}}{1 + e^{-w^T x_n}} + y_n x_{n,j}] \right] = \frac{\partial}{\partial w_j} \left[\sum_{n=1}^N \frac{-x_{n,j} - x_{n,j} e^{-w^T x_n} + x_{n,j} e^{-w^T x_n}}{1 + e^{-w^T x_n}} + y_n x_{n,j} \right] = \sum_{n=1}^N \left[\frac{x_{n,j}}{1 + e^{-w^T x_n}} - y_n x_{n,j} \right]$$

another valid linear model is

$$w^T = [1, 2]$$

$$\begin{cases} b = -2 \end{cases}$$

$$= \sum_{n=1}^N \left[\frac{1}{1 + e^{-w^T x_n}} - y_n \right] x_{n,j}$$



$$3. (a) \frac{\partial J}{\partial w_0} = \sum_{n=1}^N 2\alpha_n (w_0 + w_1 x_{n,1} - y_n) \cdot 1$$

$$\frac{\partial J}{\partial w_1} = \sum_{n=1}^N 2\alpha_n (w_0 + w_1 x_{n,1} - y_n) x_{n,1}$$

$$\therefore \nabla J = \left[\sum_{n=1}^N 2\alpha_n (w_0 + w_1 x_{n,1} - y_n), \sum_{n=1}^N 2\alpha_n (y_n x_{n,1}) \right]$$

$$\begin{array}{l} \textcircled{1} + \textcircled{2} \Rightarrow 2w_2 \quad 2b < 0 \\ \textcircled{3} + \textcircled{4} \Rightarrow -2b > 0 \end{array}$$

$$\therefore b < 0$$

$$\textcircled{2} + \textcircled{3} \Rightarrow 2b > 0 \Rightarrow b \geq 0$$

∴ There is no solution.

No such linear model exists because the data is not linearly separable.

$$\begin{aligned} \textcircled{1} &\Leftrightarrow \sum_{n=1}^N 2\alpha_n w_0 + \sum_{n=1}^N 2\alpha_n w_1 x_{n,1} - \sum_{n=1}^N 2\alpha_n y_n = 0 \\ \sum_{n=1}^N 2\alpha_n w_0 &= -\sum_{n=1}^N 2\alpha_n w_1 x_{n,1} + \sum_{n=1}^N 2\alpha_n y_n \\ &= \sum_{n=1}^N \alpha_n w_1 x_{n,1} + \sum_{n=1}^N \alpha_n y_n \end{aligned}$$

$$2. (a) \sigma(w^T x_n) = \frac{1}{1 + e^{-w^T x_n}}, \therefore 1 - \frac{1}{1 + e^{-w^T x_n}} = \frac{e^{-w^T x_n}}{1 + e^{-w^T x_n}}$$

$$\begin{aligned} \textcircled{2} &\Leftrightarrow \sum_{n=1}^N 2\alpha_n w_0 + \sum_{n=1}^N 2\alpha_n w_1 x_{n,1}^2 - \sum_{n=1}^N 2\alpha_n w_1 x_{n,1} y_n = 0 \\ \therefore \sum_{n=1}^N \alpha_n w_0 x_{n,1} &= w_0 \sum_{n=1}^N \alpha_n x_{n,1} \end{aligned}$$

$$J(w) = -\sum_{n=1}^N [y_n (\log(1) - \log(1 + e^{-w^T x_n}) + (1 - y_n) (\log(e^{-w^T x_n}) - \log(1 + e^{-w^T x_n})))]$$

$$\therefore J(w) = -\sum_{n=1}^N [y_n \log(1 + e^{-w^T x_n}) - w^T x_n + y_n w^T x_n]$$

$w_1 - \sum_{n=1}^N \alpha_n x_{n,1} + \sum_{n=1}^N \alpha_n y_n \cdot \sum_{n=1}^N \alpha_n x_{n,1} + w_1 \sum_{n=1}^N \alpha_n x_{n,1}^2$ shift the plane of $(w^T x + b)$ by ε :

$$p^+ - \varepsilon > p^- - \varepsilon$$

$$= \sum_{n=1}^N 2\alpha_n y_n x_{n,1}$$

$$-w_1 \sum_{n=1}^N \alpha_n x_{n,1} + w_1 \sum_{n=1}^N \alpha_n x_{n,1}^2 = \sum_{n=1}^N \alpha_n y_n p^+ + \sum_{n=1}^N \alpha_n y_n \cdot \sum_{n=1}^N \alpha_n x_{n,1}^2 \text{ for } (w^T x + (b - \varepsilon))$$

$$p_{\text{new}} = \min_{i \in \text{positive } y_i} (w^T x + (b - \varepsilon))$$

$$-w_1 \left(\sum_{n=1}^N \alpha_n x_{n,1}^2 + \sum_{n=1}^N \alpha_n x_{n,1}^2 \right) = \text{RHS}$$

$$= p^+ - \varepsilon$$

$$= \frac{p^+ - p^-}{2}$$

$$p_{\text{new}} = \max_{i \in \text{negative } y_i} (w^T x + (b - \varepsilon))$$

$$= p^- - \varepsilon$$

$$= \frac{p^- + p^+}{2}$$

$$p_{\text{new}} \geq 0, p_{\text{new}} < 0$$

$$p^+ - \varepsilon > 0 > p^- - \varepsilon$$

$$w_0 = \sum_{n=1}^N \alpha_n \left(\frac{\sum_{n=1}^N \alpha_n y_n x_{n,1}}{\sum_{n=1}^N \alpha_n} + \frac{\sum_{n=1}^N \alpha_n x_{n,1}^2}{\sum_{n=1}^N \alpha_n} / \sum_{n=1}^N \alpha_n y_n \right) \quad \therefore p^+ - \varepsilon > 0 > p^- - \varepsilon$$

$$p^+ - \varepsilon > p^-$$

$$y_i (w^T x + b)$$

$$\therefore y_i (w^T x + b) \geq \varepsilon$$

$$\frac{y_i (w^T x + b)}{\varepsilon} \geq 1$$

$$\text{define } w^T = \frac{w^T}{\varepsilon}, b_1 = \frac{b}{\varepsilon}$$

$$\text{then } y_i (w^T x + b_1) \geq 1 - 0$$

$$\therefore \text{optimal } \delta = 0$$

(b). If there is an optimal solution with $\delta = 0$,
 $y_i (w^T x_i + b) \geq 1 - 0 = 1$

4. (a). define $p^+ = \min_{i \in \text{positive } y_i} (w^T x_i + b)$,

$$p^- = \max_{i \in \text{negative } y_i} (w^T x_i + b)$$

$\therefore D$ is linearly separable and satisfies (1),

$$\text{by (1), } p^+ \geq 0 > p^-$$

let $\varepsilon > 0$ be the solution of $\frac{|p^+ - \varepsilon|}{\|w\|} = \frac{|p^- - \varepsilon|}{\|w\|}$

$$\Leftrightarrow \frac{p^+ - \varepsilon}{\|w\|} = \frac{-p^- + \varepsilon}{\|w\|}$$

$$\Rightarrow \varepsilon = \frac{p^+ + p^-}{2}$$

when $y_i = 1, w^T x_i + b \geq 1 > 0$

when $y_i = -1, w^T x_i + b \leq -1 < 0$

$y_i \geq 1 \text{ if } w^T x_i + b \geq 0$

$y_i \leq -1 \text{ if } w^T x_i + b < 0$

it satisfies condition (1)

$\therefore D$ is linearly separable

(c), when $0 < \delta < 1$,

$$y_i = 1 \Rightarrow w^T x_i + b \geq 1 - \delta > 0$$

$$y_i = -1 \Rightarrow w^T x_i + b < -(1 - \delta) < 0$$

condition (1) is still satisfied.

when $\delta \geq 1$, $1 - \delta \leq 0$

$$y_i = 1 \Rightarrow w^T x_i + b \geq 1 - \delta \leq 0$$

$$y_i = -1 \Rightarrow w^T x_i + b \leq -(1 - \delta) \geq 0$$

(1) may not be satisfied

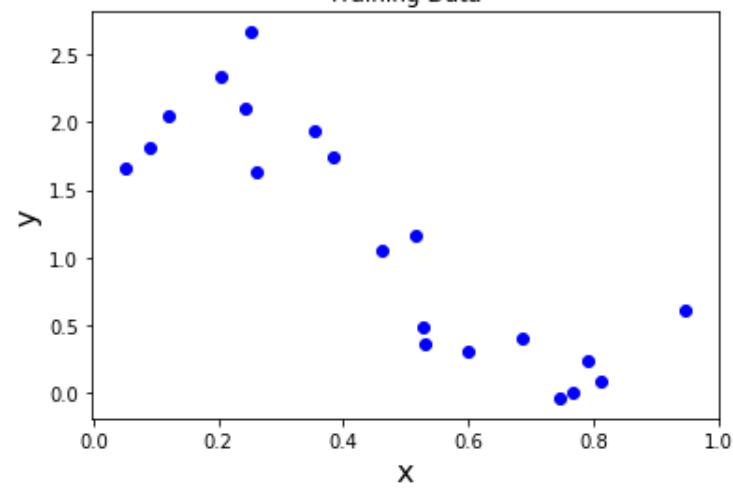
\therefore It is uncertain.

when $\delta \in (0, 1)$, D is still linearly

separable; when $\delta \geq 1$, D may or
may not be linearly separable.

5(a)

Training Data



(d). The optimal solution is: $\delta = 0$,

$$\vec{w}^T = \vec{0}, b = 0$$

So the problem is that this formulation
may be satisfied even when there
is no separating hyperplane.

(e). Let $w^T = [w_1 \ w_2 \ w_3]$,

then (2) \Rightarrow min δ

$$\text{subject to } y_i(w_1x_1 + w_2x_2 + w_3x_3 + b) \geq 1 - \delta, \forall (x_i, y_i) \in D, \delta \geq 0$$

$$\delta > 0$$

$$\therefore w_1 + w_2 + w_3 + b \geq 1 - \delta \quad ①$$

$$-(w_1 + w_2 + w_3 + b) \geq 1 - \delta \quad ②$$

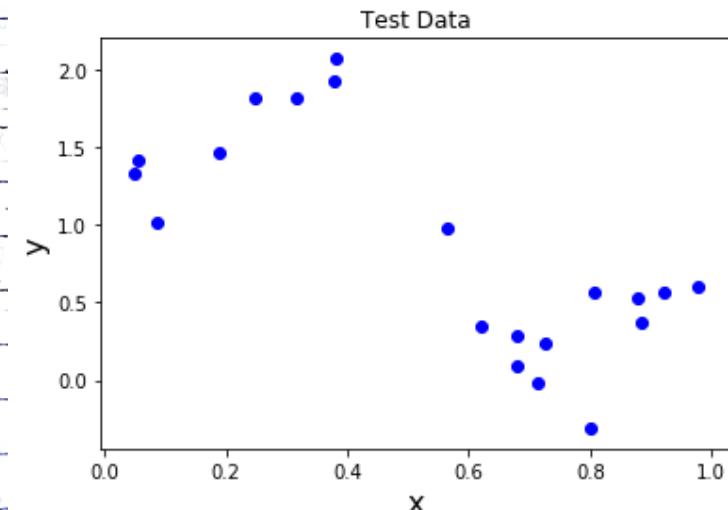
$$② \Leftrightarrow w_1 + w_2 + w_3 - b \geq 1 - \delta \quad ③$$

$$\text{when } \delta = 0, \quad ① + ③ \Rightarrow 2(w_1 + w_2 + w_3) \geq 2$$

$$w_1 + w_2 + w_3 \geq 1$$

$$\therefore w_1 + w_2 + w_3 \geq 1 \pm b$$

$$\therefore w_1 + w_2 + w_3 \geq 1 + |b|$$



The optimal solutions are $\vec{w} = [w_1 \ w_2 \ w_3]$ and b satisfying $w_1 + w_2 + w_3 \geq 1 + |b|$ at $\delta = 0$.

I observe that the distributions of training data and test data are similar. However, in the training data, there is a clear negative linear correlation, but it is hard to find in the test data. Therefore, my

educated guess is that the linear regression (c) implemented in regression.py.

from the training data may not generalize well on the test data.

(b). In regression.py

(c). In regression.py

(d). Implemented in regression.py

The closed-form solution is:

$$[w_0, w_1] = [2.44640709, -2.81635359], \text{ cost} = 3.91258$$

almost the same as the results of GD when convergence is reached.

Both the $[w_0, w_1]$ and the cost are

Its time ≈ 0.0 , so it is much faster than GD.

(f). Implemented in regression.py

time	cost	coefficients (w_0, w_1)
0.26563	3.91258	[2.44640709, -2.81635359]

with # iterations = 1356

η	time	#iterations	cost	coefficients (w_0, w_1)
0.0001	9.0966780185699463	10000	4.0863970	[2.27044798, -2.46064834]
0.001	1.48087	7020	3.91258	[2.44640709, -2.81635359]
0.01	0.27312	764	3.91258	[2.44640709, -2.81635359]
0.0407	2.10207	10000	2.71092e+3	[-9.407093e+18, -4.65929895e+18]

For $\eta = 10^{-4}$ and $\eta = 0.0407$, GD does (g). In regression.py,

not converge in 10000 steps. The reason (h). Implemented in regression.py

is that 10^{-4} is too small while 0.0407

Using RMSE can get rid of the influence is too large, so for 10^{-4} it is converging of the size of data on the error, so too slowly while for 0.0407 it is unstable. errors of different models can be compared

From $\eta = 10^{-4}$, $\eta = 10^{-3}$ to $\eta = 10^{-2}$, the

It normalizes the error.

number of iterations needed to the error and the time used decrease.

(i). As shown in the following plot, degrees of 4, 5 and 6 best fit the data. These values minimize the test error, and the training error is also very low.

The coefficients of $\eta = 10^{-3}$ and $\eta = 10^{-2}$

When $m \leq 3$, both training error and test error are larger, indicating underfitting.

are largely the same, since in these

When $m \geq 8$, the training error continue to decrease while the test error increases dramatically, indicating overfitting.

to cases the iterations has finished within

10000 steps. The coefficients of $\eta = 10^{-4}$ is

close to them, and it is

less accurate because the iterations was

terminated before convergence. The

coefficients of $\eta = 0.0407$ does not make sense.

RMSE versus m

