

## Contribution Report

Discussion and Solve Problems for Questions 1, 2, and 3:  
Equal Contribution by both partners

Solutions Write-up for Questions 1, 2, and 3:  
David Wu

Proof-reading for Questions 1, 2, and 3:  
Shu Wang

Let  $u, v \in \mathbb{R}^n$ ,  $c \in \mathbb{R}$ , and  $y$  be any linear regression predictor.

$$\begin{aligned}\text{Then: } y(u+v) &= w^T (u+v) \\ &= \sum_i w_i (u_i + v_i) \\ &= \sum_i w_i (u_i + v_i) \\ &= \sum_i w_i u_i + w_i v_i \\ &= (w_1 u_1 + w_1 v_1) + \dots + (w_n u_n + w_n v_n)\end{aligned}$$

Rearranging the above gives:

$$\begin{aligned}y(u+v) &= (w_1 u_1 + w_2 u_2 + \dots + w_n u_n) + (w_1 v_1 + w_2 v_2 + \dots + w_n v_n) \\ &= \sum_i w_i u_i + \sum_i w_i v_i \\ &= w^T u + w^T v \\ &= y(u) + y(v).\end{aligned}$$

Similarly, we have:

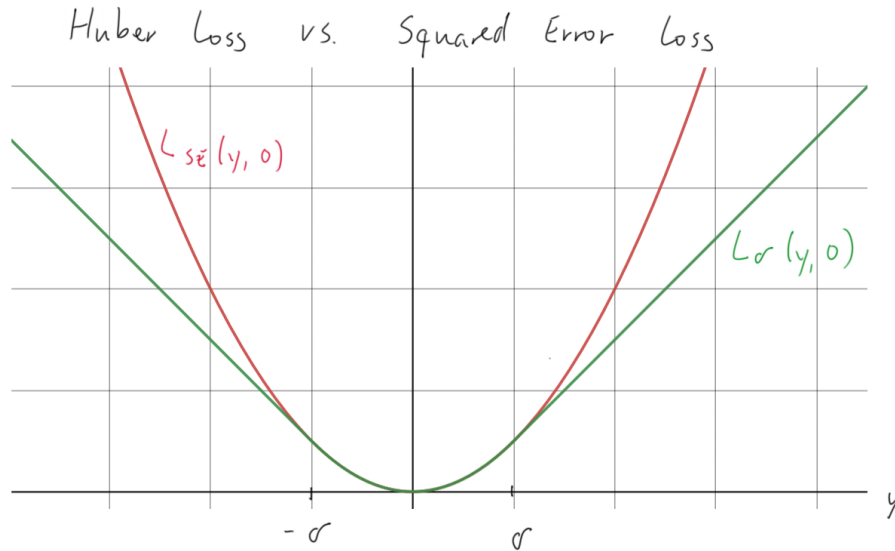
$$\begin{aligned}y(cu) &= w^T (cu) \\ &= \sum_i w_i (c u_i) \\ &= w_1 c u_1 + w_2 c u_2 + \dots + w_n c u_n\end{aligned}$$

Rearranging the above:

$$\begin{aligned}y(cu) &= c w_1 u_1 + c w_2 u_2 + \dots + c w_n u_n \\ &= c (w_1 u_1 + \dots + w_n u_n) \\ &= c \sum_i w_i u_i \\ &= c w^T u = c y(u)\end{aligned}$$

Therefore all linear regression predictors are linear functions.

2. a.



Consider  $N$  data points with no outliers and  $y^*$  be the optimal predictor. Now consider additional outlier points  $(x^{(i)}, t^{(i)})$ . Since these points come from a different conditional distribution, we expect that  $y^*(x^{(i)}) - t^{(i)}$  to be large, meaning  $|y^*(x^{(i)}) - t^{(i)}| > \delta$ .

From the graphs,  $L_{SE}$  assigns a larger loss to  $y^*(x^{(i)}) - t^{(i)}$  than  $L_{\delta}$ . For a modelling fitting all the points ( $N + \text{outliers}$ ),  $y^*$  is no longer the optimal predictor, but since we are trying to minimize average loss, a small value for  $L(y^*(x^{(i)}) - t^{(i)})$  terms mean that we have to adjust  $y^*$  less. Therefore, the predictor under  $L_{\delta}$  is affected by outliers less when compared to  $L_{SE}$ .

2.b.

$$H_{\sigma}(a) = \begin{cases} \sigma(-a - \frac{1}{2}\sigma) & a < -\sigma \\ \frac{1}{2}a^2 & -\sigma \leq a \leq \sigma \\ \sigma(a - \frac{1}{2}\sigma) & a > \sigma \end{cases} \Rightarrow H'_{\sigma}(a) = \begin{cases} -\sigma & a < -\sigma \\ a & |a| \leq \sigma \\ \sigma & a > \sigma \end{cases}$$

$$\vec{y} = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{pmatrix} = \begin{pmatrix} x^{(1)T} w \\ \vdots \\ x^{(n)T} w \end{pmatrix} = Xw.$$

$$\frac{\partial \hat{R}}{\partial \vec{y}} = \left( \frac{\partial \hat{R}}{\partial y^{(1)}} \quad \dots \quad \frac{\partial \hat{R}}{\partial y^{(n)}} \right)^T. \text{ Since only the } j\text{th term of } \hat{R} = \frac{1}{N} \sum_{i=1}^N L_{\sigma}(y^{(i)}, t^{(i)})$$

$$\text{depends on } y^{(j)}, \quad 0 \leq j \leq N, \quad \frac{\partial \hat{R}}{\partial y^{(j)}} = \frac{1}{N} H'_{\sigma}(y^{(j)} - t^{(j)})$$

$$\text{In vector notation } \boxed{\frac{\partial \hat{R}}{\partial \vec{y}} = \frac{1}{N} H'_{\sigma}(\vec{y} - \vec{t})}$$

$$\text{Using the chain rule } \frac{\partial \hat{R}}{\partial w_j} = \frac{1}{N} \sum_{i=1}^N L'_{\sigma}(w^T x^{(i)}, t) \cdot x_j^{(i)}, \text{ where } x_j^{(i)} \text{ is}$$

the  $j$ th element of  $x^{(i)}$ , and  $L'$  is derivative with respect to  $y^{(i)}$

$$\text{The above can be written as } \frac{\partial \hat{R}}{\partial w_j} = \frac{1}{N} \begin{pmatrix} x_j^{(1)} & x_j^{(2)} & \dots & x_j^{(N)} \end{pmatrix} \begin{bmatrix} L'_{\sigma}(w^T x^{(1)}, t) \\ \vdots \\ L'_{\sigma}(w^T x^{(n)}, t) \end{bmatrix}$$

$$\text{Therefore } \boxed{\frac{\partial \hat{R}}{\partial \vec{w}} = \frac{1}{N} X^T H'_{\sigma}(Xw - \vec{t})}$$

c. See python file

d. See python file

e. The training set loss is calculated by  $\frac{1}{2N} \sum (\hat{y} - t)^2$ , and when training the model we find weights that minimize  $\frac{1}{N} \sum H_\delta(w^T x^{(i)} - t)$ .

However, for larger  $\delta$ , Huber loss becomes closer to squared error loss (see picture in part a), and the training model chooses weights to minimize a function that is closer to  $\frac{1}{2N} \sum (\hat{y} - t)^2$ . Therefore, we expect the training set loss to decrease for increasing  $\delta$ .

Validation loss does not always increase, because the model generalizes poorly. For larger  $\delta$ , the model fits squared errors loss better on the training set but also gets affected by the outliers more. The optimal  $\delta$  needs to balance between these two factors.

#### Appendix: Output from python

```
linear regression validation loss: 28.44202896790706
delta: 0.1, valid. squared error loss: 4.909228730093678, train squared error loss: 38.041444425291125
delta: 0.5, valid. squared error loss: 4.018972761905316, train squared error loss: 38.311166891032116
delta: 1, valid. squared error loss: 4.200164017924187, train squared error loss: 37.993439863363456
delta: 5, valid. squared error loss: 23.636105018763637, train squared error loss: 28.84756778475285
delta: 10, valid. squared error loss: 26.420467729491516, train squared error loss: 27.79829743150903
```



3.a.  $\begin{array}{c|c|c} + & - & + \\ -1 & 1 & 3 \end{array} x^{(i)}$  Suppose some weights perfectly classifies all points. Then it classifies  $x^{(i)} = -1$  and  $x^{(i)} = 3$  as 1.

From lecture, we know that the set of points that have the same prediction is a convex set, so the point  $\frac{1}{2}(-1) + \frac{1}{2}(3) = -0.5 + 1.5 = 1$  must be classified with label 1. However, if  $x^{(i)} = 1$ ,  $t^{(i)} = 0$ , so we have a contradiction and the data is not linearly separable.

b.  $\begin{array}{c|c|c} \psi_1 & \psi_2 & t^{(i)} \\ -1 & 1 & 1 \\ 1 & 1 & 0 \\ 3 & 9 & 1 \end{array}$  we want  $w_1 \psi_1 + w_2 \psi_2 \geq 0$  only for points with  $t^{(i)} = 1$ . This gives:

$$-w_1 + w_2 \geq 0 \quad \text{--- (1)} \quad \quad 3w_1 + 9w_2 \geq 0 \quad \text{--- (3)}$$

$$w_1 + w_2 < 0 \quad \text{--- (2)}$$

A possible solution is  $w_1 = -2$  and  $w_2 = 1$ . Then:

(1) becomes  $-(-2) + 1 = 2 + 1 \geq 0$

(2) becomes  $-2 + 1 = -1 < 0$

(3) becomes  $3(-2) + 9(1) = -6 + 9 \geq 0$ .

All 3 points are classified correctly.