

# COMP 3105 — Fall 2025 — Assignment 1

Pascal Law 101318299 | Ayaan Aleem 101307254

## Question 1:

Subquestion b1.

$$\underset{1 \times (d+1)}{\mathbf{c}^T} \underset{(d+1) \times 1}{\mathbf{u}} = \delta$$

Since  $\mathbf{u} = \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix}$ , we can make all the weights in  $\mathbf{w}$  zero, and just have  $\delta$  in the last entry of  $\mathbf{u}$ .

As such, we can set:

$$\mathbf{c} = \underset{(d+1) \times 1}{\begin{bmatrix} \mathbf{0}_{1 \times d} & 1 \end{bmatrix}^T}$$

Subquestion b2.

We can break

$$\underset{1 \times (d+1)}{G^1} \underset{1 \times 1}{\begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix}} \leq \underset{1 \times 1}{\mathbf{h}^1}$$

into two parts, and it becomes:

$$\underset{1 \times d}{G^{11}} \underset{1 \times 1}{G^{12}} \underset{1 \times 1}{\begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix}} \leq \underset{1 \times 1}{\mathbf{h}^1}$$

As such,

$$\underset{1 \times d}{G^{11}} \mathbf{w} + \underset{1 \times 1}{G^{12}} \delta \leq \underset{1 \times 1}{\mathbf{h}^1}$$

Thus, given that the first part of the linear programming problem is:

$$-\delta \leq 0$$

we can let:

$$\underset{1 \times d}{G^{11}} = \mathbf{0}_{1 \times d}$$

$$\underset{1 \times 1}{G^{12}} = -1$$

$$\underset{1 \times (d+1)}{G^1} = \underset{1 \times (d+1)}{\begin{bmatrix} \mathbf{0}_{1 \times d} & -1 \end{bmatrix}}$$

$$\underset{1 \times 1}{\mathbf{h}^1} = 0$$

Subquestion b3.

We can break

$$G_{n \times (d+1)}^2 \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix} \leq \mathbf{h}_{n \times 1}^2$$

into two parts, and it becomes:

$$G_{n \times d}^{21} G_{n \times 1}^{22} \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix} \leq \mathbf{h}_{n \times 1}^2$$

As such,

$$G_{n \times d}^{21} \mathbf{w} + G_{n \times 1}^{22} \delta \leq \mathbf{h}_{n \times 1}^2$$

Thus, given that the second part of the linear programming problem is:

$$X\mathbf{w} - \delta \cdot \mathbf{1}_n \leq \mathbf{y}$$

we can let:

$$G_{n \times d}^{21} = X$$

$$G_{n \times 1}^{22} = -\mathbf{1}_{n \times 1}$$

$$G_{n \times (d+1)}^2 = \begin{bmatrix} X & -\mathbf{1}_{n \times 1} \end{bmatrix}_{n \times (d+1)}$$

$$\mathbf{h}_{n \times 1}^2 = \mathbf{y}$$

Subquestion b4.

We can break

$$G_{n \times (d+1)}^3 \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix} \leq \mathbf{h}_{n \times 1}^3$$

into two parts, and it becomes:

$$G_{n \times d}^{31} G_{n \times 1}^{32} \begin{bmatrix} \mathbf{w} \\ \delta \end{bmatrix} \leq \mathbf{h}_{n \times 1}^3$$

As such,

$$G_{n \times d}^{31} \mathbf{w} + G_{n \times 1}^{32} \delta \leq \mathbf{h}_{n \times 1}^3$$

Thus, given that the third part of the linear programming problem is:

$$-X\mathbf{w} - \delta \cdot \mathbf{1}_n \leq -\mathbf{y}$$

we can let:

$$G_{n \times d}^{31} = -X$$

$$\begin{aligned}
G_{n \times 1}^{32} &= -\mathbf{1}_{n \times 1} \\
G_{n \times (d+1)}^3 &= \begin{bmatrix} -X & -\mathbf{1}_{n \times 1} \end{bmatrix}_{n \times (d+1)} \\
\mathbf{h}_{n \times 1}^3 &= -\mathbf{y}
\end{aligned}$$

Subquestion c1.

Table 1: Different training losses for different models

<b>Model</b>	$L_2$ loss	$L_\infty$ loss
$L_2$ model	0.09557022	1.58819212
$L_\infty$ model	0.23302085	0.86917599

Table 2: Different test losses for different models

<b>Model</b>	$L_2$ loss	$L_\infty$ loss
$L_2$ model	0.05291075	1.03171424
$L_\infty$ model	0.29564879	2.1535331

Subquestion c2.

Analysis of result:

When we compare our training data with our test data, we find that our  $L_2$  loss has a smaller difference when compared to our  $L_\infty$  loss, and similarly our  $L_2$  model also has a smaller difference when compared to our  $L_\infty$  model.

This means that our  $L_2$  model using  $L_2$  loss has the smallest difference and so is the most accurate prediction of the data. On the other hand, the  $L_\infty$  model has the biggest difference which means that  $L_\infty$  model using  $L_\infty$  loss is the most inaccurate prediction of the data.

The reason that the  $L_\infty$  loss is so inaccurate is because of the difference of size between the test data and the training data. We only generate 30 samples for the training while we generate 1000 samples for the test data. This means that our  $L_\infty$  loss is modelled on a very limited set of training data, and so it is heavily affected by outliers in our much larger test data.

This is in contrast with our  $L_2$  loss, which is not as affected by outliers, which makes it much more accurate when the training data sample is small.

When considering the models, the  $L_2$  model is more accurate because it has a smaller tolerance, which means it doesn't chase outliers like our  $L_\infty$  model.

These factors lead to the conclusion that our  $L_2$  loss and  $L_2$  model are more accurate than the  $L_\infty$  counterparts for this set of training and test data.

Subquestion d2.

Table 3: Different training losses for different models

<b>Model</b>	$L_2$ loss	$L_\infty$ loss
$L_2$ model	53.18368801	32.21643524
$L_\infty$ model	66.09571725	26.39169568

Table 4: Different test losses for different models

<b>Model</b>	$L_2$ loss	$L_\infty$ loss
$L_2$ model	55.20252554	33.40449987
$L_\infty$ model	68.04960819	36.9655386

**Question 2:**

Subquestion c1.

Table 5: Training accuracies with different hyper-parameters

<b>m</b>	<b>Train Accuracy</b>	<b>dim1</b>	<b>Train Accuracy</b>	<b>dim2</b>	<b>Train Accuracy</b>
10	0.9725	1	0.8467	1	0.9266
50	0.9252	2	0.9262	2	0.9275
100	0.9217	4	0.98285	4	0.9255
200	0.923125	8	0.99935	8	0.9349

Table 6: Testing accuracies with different hyper-parameters

<b>m</b>	<b>Test Accuracy</b>	<b>dim1</b>	<b>Test Accuracy</b>	<b>dim2</b>	<b>Test Accuracy</b>
10	0.882065	1	0.836885	1	0.91926
50	0.911655	2	0.91675	2	0.915435
100	0.91711	4	0.968755	4	0.915865
200	0.91836	8	0.996295	8	0.90904

Subquestion c2.

Our  $m$  represents the number of data points in the training data sample. As our  $m$  increases, our training accuracy decreases, but our test accuracy increases. This is because when  $m$  is low, we have very few data points to train our optimal parameter which causes us to create a line that is unsuitable for a large number of data points, like the amount we have in the test data. On the other hand, as our  $m$  increases, we have a much larger sample to train our optimal parameter on, which helps to create a line that is more suitable for our test data points. However, this may not necessarily be completely accurate for the set of training data points due to the higher amount of randomly distributed points that may skew the model.

Both for training and test data, as  $\text{dim1}$  increases our accuracy becomes almost perfect. This is because  $\text{dim1}$  is an expression of the number of useful dimensions; when we have a higher  $\text{dim1}$  in training, our info on the data is also higher which allows us to predict labels with a higher degree of accuracy. Similarly, when our test data has a higher  $\text{dim1}$ , our accuracy is again higher because of more useful information on the data points.

However, as our  $\text{dim2}$  increases, it has opposite effects in training and test data; in training data, it marginally increases the accuracy whereas in test data, the accuracy drops slowly but steadily. This is because  $\text{dim2}$  is an expression of the number of noise dimensions which add to the variance. When we train using a larger  $\text{dim2}$ , it doesn't really increase our accuracy by much, however when we use test data with larger  $\text{dim2}$ , the increase in variance negatively affects our accuracy.

## References:

- Asked on Discord group about an error message related to `scipy.optimize.minimize` method.
- Checked Discord group messages about the error messages which others have encountered (especially `scipy.special.expi()` method and how to suppress warnings).
- Asked ChatGPT to explain `dim1` and `dim2` meanings in Q2.
- Asked ChatGPT to explain Python syntax as in `for i, m in enumerate((10, 50, 100, 200)):`
- Googled `numpy`, `pandas`, `scipy`, `autograd`, etc. documentations.