

## Q1

|                               | Larger C                 | Smaller C                |
|-------------------------------|--------------------------|--------------------------|
| Model capacity (large/small?) | Large                    | Low                      |
| Overfitting/Underfitting?     | Overfitting              | Underfitting             |
| Bias variance (how/low?)      | Low bias & High variance | High bias & Low variance |

## Q2

$$\begin{aligned}
p(w|\mathcal{D}) &\propto p(w)p(w|\mathcal{D}) \\
&= p(w) \prod_{i=1}^m p(t^{(i)}|x^{(i)}, w) \\
&= \frac{1}{\sqrt{2\pi}\sigma_w} \exp\left\{-\frac{1}{2\sigma_w^2}w^2\right\} \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma_\epsilon} \exp\left\{-\frac{1}{2\sigma_\epsilon^2}(wx^{(i)} - t^{(i)})^2\right\} \\
&= \frac{1}{\sqrt{2\pi}\sigma_w} \exp\left\{-\frac{1}{2\sigma_w^2}w^2\right\} \left(\frac{1}{\sqrt{2\pi}\sigma_\epsilon}\right)^m \exp\left\{-\frac{1}{2\sigma_\epsilon^2} \sum_{i=1}^m (wx^{(i)} - t^{(i)})^2\right\} \\
&\propto \exp\left\{-\frac{1}{2}\left(\frac{w^2}{\sigma_w^2} + \frac{\sum_{i=1}^m (wx^{(i)} - t^{(i)})^2}{\sigma_\epsilon^2}\right)\right\} \\
&\quad \text{(Dropping constant term inside exp.)} \\
&= \exp\left\{-\frac{1}{2}\left(\frac{w^2}{\sigma_w^2} + \frac{\sum_{i=1}^m (x^{(i)})^2 w^2 - 2 \sum_{i=1}^m t^{(i)} x^{(i)} w + \sum_{i=1}^m (t^{(i)})^2}{\sigma_\epsilon^2}\right)\right\} \\
&\quad \text{(Rearrange.)} \\
&\propto \exp\left\{-\frac{1}{2\sigma_w^2\sigma_\epsilon^2} \left( (\sigma_\epsilon^2 + \sigma_w^2 \sum_{i=1}^m (x^{(i)})^2) w^2 - 2\sigma_w^2 \sum_{i=1}^m t^{(i)} x^{(i)} w \right)\right\} \\
&\quad \text{(Dropping constant term inside exp.)} \\
&\propto \exp\left\{-\frac{\sigma_\epsilon^2 + \sigma_w^2 \sum_{i=1}^m (x^{(i)})^2}{2\sigma_w^2\sigma_\epsilon^2} \left( w^2 - 2\frac{\sigma_w^2 \sum_{i=1}^m t^{(i)} x^{(i)}}{\sigma_\epsilon^2 + \sigma_w^2 \sum_{i=1}^m (x^{(i)})^2} w \right)\right\} \\
&\quad \text{(Dropping constant term inside exp.)} \\
&\propto \exp\left\{-\frac{\sigma_\epsilon^2 + \sigma_w^2 \sum_{i=1}^m (x^{(i)})^2}{2\sigma_w^2\sigma_\epsilon^2} \left( w - \frac{\sigma_w^2 \sum_{i=1}^m t^{(i)} x^{(i)}}{\sigma_\epsilon^2 + \sigma_w^2 \sum_{i=1}^m (x^{(i)})^2} \right)^2\right\} \\
&\quad \text{(Dropping constant term inside exp.)}
\end{aligned}$$

We can see that  $p(w|\mathcal{D})$  is also of the form of a Gaussian distribution, where

$$\sigma_{post}^2 = \frac{\sigma_w^2\sigma_\epsilon^2}{\sigma_\epsilon^2 + \sigma_w^2 \sum_{i=1}^m (x^{(i)})^2} \text{ and } \mu_{post} = \frac{\sigma_w^2 \sum_{i=1}^m t^{(i)} x^{(i)}}{\sigma_\epsilon^2 + \sigma_w^2 \sum_{i=1}^m (x^{(i)})^2}$$

### Q3

We know that  $p(\mathbf{w}|\mathcal{D}) \propto p(\mathbf{w})p(\mathcal{D}|\mathbf{w})$ .

For maximization, it is easy to see that

$$\operatorname{argmax}_{\mathbf{w}} \log p(\mathbf{w}|\mathcal{D}) \propto \operatorname{argmax}_{\mathbf{w}} \left( \log p(\mathbf{w}) + \log p(\mathcal{D}|\mathbf{w}) \right)$$

Following the lecture notes given a dataset  $\mathcal{D}$  and a linear regression model, we know that

$$p(\mathcal{D}|\mathbf{w}) = \prod_{i=1}^m \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{1}{2\sigma^2} (t^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 \right\}$$

and

$$\log p(\mathcal{D}|\mathbf{w}) = -\frac{1}{2\sigma^2} \sum_{i=1}^m (t^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 + \text{const} \quad (1)$$

Assume each  $w_i$  in  $\mathbf{w}$  is drawn from  $\text{Laplace}(0, b)$ , denoted as

$$p(w_i) = \frac{1}{2b} \exp \left\{ -\frac{|w_i|}{b} \right\}$$

We have

$$p(\mathbf{w}) = \prod_{i=1}^m \frac{1}{2b} \exp \left\{ -\frac{|w_i|}{b} \right\}$$

and

$$\begin{aligned} \log p(\mathbf{w}) &= -\frac{1}{b} \sum_{i=1}^m |w_i| + \text{const} \\ &= -\frac{1}{b} \|\mathbf{w}\|_1 + \text{const} \end{aligned} \quad (2)$$

From (1) and (2), we have

$$\begin{aligned} &\log p(\mathbf{w}) + \log p(\mathcal{D}|\mathbf{w}) \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^m (t^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 - \frac{1}{b} \|\mathbf{w}\|_1 + \text{const} \end{aligned}$$

Therefore, the Laplace prior gives us an linear regression model with  $L_1$  regularization.