

2b)

$$AA_{i,j}^T = \sum_{j=0}^M \sum_{k=0}^N A_{i,k} A_{j,k}$$

$$\begin{aligned}
 3a) \quad l(w_0, w) &= \sum [t^{(n)} - y(x^{(n)})]^2 \\
 &= \sum [y(x^{(n)}) - t^{(n)}]^2 \\
 &= \sum [y^{(n)} - t^{(n)}]^2
 \end{aligned}$$

$y^{(n)}$ & $t^{(n)}$ can be written as vectors $\vec{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$, $\vec{t} = \begin{pmatrix} t_1 \\ \vdots \\ t_n \end{pmatrix}$

Also, the summation can be written as the magnitude.

Thus,

$$\sum [y^{(n)} - t^{(n)}]^2 = \|\vec{y} - \vec{t}\|^2, \text{ as required } \square$$

$$3b) \quad y(x) = w_0 + \sum w_m z_m$$

$\sum w_m z_m$ can be written as Zw , where Z is a feature matrix with m features and n columns.

w can be written as a column vector, $w = \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix}$

Zw would produce an $m \times 1$ matrix.

$$\text{so, } \sum w_m z_m = Zw$$

To add w_0 to Zw we must turn w_0 into a column matrix with m rows. We can do this by multiplying w_0 by a column vector of 1's.

Thus,

$$y = w_0 \vec{1} + Zw, \text{ as required } \square$$

$$\begin{aligned}
 3c) \quad l(w_0, w) &= \sum [t^{(n)} - y(x^{(n)})]^2 = \sum [y^{(n)} - t^{(n)}]^2 \\
 \frac{\partial l(w_0, w)}{\partial w} &= 2 \sum [y^{(n)} - t^{(n)}] \cdot \frac{\partial}{\partial w} [y^{(n)} - t^{(n)}] \\
 &= 2(y - t) \cdot \frac{\partial}{\partial w} [w_0 + \sum w_m z_m - t^{(n)}] \\
 &= 2(y - t) \cdot \frac{\partial}{\partial w} (Zw) \\
 &= 2(y - t) \cdot Z^T \\
 &= 2
 \end{aligned}$$

$$3d) l(\omega_0, \omega) = \sum [y^{(n)} - t^{(n)}]^2$$

$$\frac{dl(\omega_0, \omega)}{d\omega_0} = 2 \sum [y^{(n)} - t^{(n)}] \cdot \frac{d}{d\omega_0} (y^{(n)} - t^{(n)})$$

$$= 2 \sum [y^{(n)} - t^{(n)}] \frac{d}{d\omega_0} [\omega_0 + \sum z_m \omega_m] \quad \text{Note: } \frac{d}{d\omega_0} [\omega_0 + \sum z_m \omega_m]$$

$$= 2 \sum [y^{(n)} - t^{(n)}] \cdot 1$$

$$= 1$$

$$= 2 \cdot \vec{1}^T (y - t)$$

4e)

The testing error is larger than the training error because the model is poor in general, due to the lack of basis functions. Although it is poor, it fits the training data much better than the test data.

4f)

The model begins to slowly overfit the training data, thus resulting in a large difference between the training error and testing error.

4g)

Just by looking at the graph, we can see that the model is highly overfit. The model goes through every point of the training data, instead of generalizing the data. The large delta between the training and testing error is a quantification of the figure.

5d)

It can be seen that the smallest gamma fails to regularize the model enough, and hence we see many inflection points. On the other hand, the largest gamma value regularizes the model too much and doesn't have enough variability and thus turns the model into a parabola, which is inaccurate to our data. The optimal gamma regularizes the model just enough to take the general structure of how the points are laid out.

$$6a) \tilde{L}(\omega_0, \omega) = L(\omega_0, \omega) + \gamma \sum \omega_j^2$$

$$\frac{\partial \tilde{L}(\omega_0, \omega)}{\partial \omega} = \frac{\partial L}{\partial \omega} + \frac{\partial}{\partial \omega} \left(\gamma \sum_{j=1}^n \omega_j^2 \right)$$

$$= 2Z^T(y-t) + 2\gamma \sum_{j=1}^n \omega_j$$

$$= 2 \left[Z^T(y-t) + \gamma \sum_{j=1}^n \omega_j \right]$$

$$6b) \frac{\partial \tilde{L}(\omega_0, \omega)}{\partial \omega_0} = \frac{\partial L}{\partial \omega_0} + \frac{\partial}{\partial \omega_0} \left(\gamma \sum_{j=1}^n \omega_j^2 \right)$$

$$= \frac{\partial L}{\partial \omega_0}$$

$$= 2\tilde{I}^T(y-t)$$