

Question 2

2(b)

The equation for computing AA^T without computing A^T is:

$$AA^T_{ij} = \sum_k A_{ik} A_{jk}$$

2(d)

myfun performs MN^2 floating-point multiplications and $MN(N-1)$ floating-point additions.

Question 3

3(a)

Question 3(a)

Prove that $\ell(w_0, w) = \|y - t\|^2$,

where $t = [t^{(1)}, t^{(2)}, \dots, t^{(n)}]$
 $y = [y^{(1)}, y^{(2)}, \dots, y^{(n)}]$

Note: Definition of magnitude of a vector

$$\|v\|^2 = \sum_n v_n^2$$
$$\|y - t\|^2 = \sum_n [y^{(n)} - t^{(n)}]^2$$

we can reorder $[y^{(n)} - t^{(n)}]^2$
b/c $x^2 = (-x)^2 \Rightarrow (a-b)^2 = [-1(a-b)]^2 = (b-a)^2$

$$= \sum_n [t^{(n)} - y^{(n)}]^2$$
$$= \ell(w_0, w) \quad (\text{def. of loss function (4)}) \quad \blacksquare$$

3(b)

Question 3(b)

Prove that $y = w_0 \vec{1} + Z w$,

where y and w are treated as column vectors

$$w_0 \vec{1} + Z w = \begin{bmatrix} \phi_1(x^{(1)}) & \phi_2(x^{(1)}) & \dots & \phi_M(x^{(1)}) \\ \phi_1(x^{(2)}) & \phi_2(x^{(2)}) & \dots & \phi_M(x^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(x^{(N)}) & \phi_2(x^{(N)}) & \dots & \phi_M(x^{(N)}) \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_M \end{bmatrix} + w_0 \vec{1}$$

$$= \begin{bmatrix} w_1 \phi_1(x^{(1)}) + w_2 \phi_2(x^{(1)}) + \dots + w_M \phi_M(x^{(1)}) \\ w_1 \phi_1(x^{(2)}) + w_2 \phi_2(x^{(2)}) + \dots + w_M \phi_M(x^{(2)}) \\ \vdots \\ w_1 \phi_1(x^{(N)}) + w_2 \phi_2(x^{(N)}) + \dots + w_M \phi_M(x^{(N)}) \end{bmatrix} + w_0 \vec{1}$$

let a_i represent row i of $Z w$ where $i \in [1, N]$

$$= \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{bmatrix} + \begin{bmatrix} w_0 \\ w_0 \\ \vdots \\ w_0 \end{bmatrix} = \begin{bmatrix} w_0 + \sum_{m=1}^M w_m z_m \\ w_0 + \sum_{m=1}^M w_m z_m \\ \vdots \\ w_0 + \sum_{m=1}^M w_m z_m \end{bmatrix} \quad (\text{by eqn. \#5})$$

\leftarrow where $x = x^{(1)}$
 \leftarrow where $x = x^{(2)}$
 \vdots
 \leftarrow where $x = x^{(N)}$

$$= \begin{bmatrix} y(x^{(1)}) \\ y(x^{(2)}) \\ \vdots \\ y(x^{(N)}) \end{bmatrix} = y$$

$$\text{Note: } y(x) = w_0 + \sum_{m=1}^M w_m z_m \quad (5)$$

3(c)

Question 3(c)Prove that $\frac{d\ell(w, u)}{du} = 2Z^T(y - t)$

$$\begin{aligned}
\frac{d\ell(w, u)}{du} &= \frac{d \sum_{n=1}^N [y(x^{(n)}) - t^{(n)}]^2}{du} \\
&= \sum_{n=1}^N \frac{d [y(x^{(n)}) - t^{(n)}]^2}{du} \\
&= \sum_{n=1}^N 2(y(x^{(n)}) - t^{(n)}) \cdot \frac{d(y(x^{(n)}) - t^{(n)})}{du} \\
&= 2 \sum_{n=1}^N (y(x^{(n)}) - t^{(n)}) \cdot \frac{d(\sum_{m=1}^M w_m z_m - t^{(n)})}{du} \\
&= 2 \sum_{n=1}^N (y(x^{(n)}) - t^{(n)}) \cdot \phi_m(x^{(n)}) \\
&= 2 \sum_{n=1}^N (y(x^{(n)}) - t^{(n)}) \cdot z_{nm} \\
&= 2 \begin{bmatrix} \sum_{n=1}^N (y(x^{(n)}) - t^{(n)}) z_{n0} \\ \vdots \\ \sum_{n=1}^N (y(x^{(n)}) - t^{(n)}) z_{nM} \end{bmatrix} \\
&= 2 \begin{bmatrix} \sum_{n=1}^N (y(x^{(n)}) - t^{(n)}) \phi_0(x^{(n)}) \\ \vdots \\ \sum_{n=1}^N (y(x^{(n)}) - t^{(n)}) \phi_M(x^{(n)}) \end{bmatrix} \\
\text{Note: } \phi_0(x) &= 1 \\
&= 2 \begin{bmatrix} 1 & \dots & 1 \\ \vdots & & \vdots \\ \phi_M(x^{(1)}) & \dots & \phi_M(x^{(N)}) \end{bmatrix} \cdot \begin{bmatrix} y(x^{(1)}) - t^{(1)} \\ \vdots \\ y(x^{(N)}) - t^{(N)} \end{bmatrix} \\
&= 2Z^T(y - t) \quad \blacksquare
\end{aligned}$$

3(d)

Question 3(d)

Prove that $\frac{d\mathcal{L}(w_0, w)}{dw_0} = 2\vec{1}^T(y-t)$

$$\begin{aligned}\frac{d\mathcal{L}(w_0, w)}{dw_0} &= \frac{d \sum_{n=1}^N [y(x^{(n)}) - t^{(n)}]^2}{dw_0} \\&= \sum_{n=1}^N \frac{d[y(x^{(n)}) - t^{(n)}]^2}{dw_0} \\&= \sum_{n=1}^N 2(y(x^{(n)}) - t^{(n)}) \cdot \frac{d(y(x^{(n)}) - t^{(n)})}{dw_0} \\&= 2 \sum_{n=1}^N (y(x^{(n)}) - t^{(n)}) \cdot \frac{d\left(\sum_{m=0}^M w_m z_m - t^{(n)}\right)}{dw_0} \\&= 2 \sum_{n=1}^N (y(x^{(n)}) - t^{(n)}) \cdot \phi_0(x^{(n)}) \\&= 2 \sum_{n=1}^N (y(x^{(n)}) - t^{(n)}) \cdot z_{n0} \\&= 2 \sum_{n=1}^N (y(x^{(n)}) - t^{(n)}) \phi_0(x^{(n)})\end{aligned}$$

Note: $\phi_0(x) = 1$

$$\begin{aligned}&= 2 \begin{bmatrix} 1 & \dots & 1 \end{bmatrix} \cdot \begin{bmatrix} y(x^{(1)}) - t^{(1)} \\ \vdots \\ y(x^{(N)}) - t^{(N)} \end{bmatrix} \\&= 2\vec{1}^T(y-t)\end{aligned}$$

Question 4

4(e)

The test error is larger than the training error because the model was trained to fit the points in the training data to minimize specifically this training error. In other words, a fitted model adapts to the training data while the test data is new and it will have more difficulty predicting it.

4(f)

The model is starting to overfit the training data, which is evident by the training error being much lower than the test error. In other words, the model has adapted to the training data specifically, so much that it does not generalize well for other inputs.

4(g)

The model is now a perfect fit for the training data, which is a bad thing if you want to use it to predict anything other than the training data. The training error is very low because the plot of the model passes through every point in the training data. The test error is very high because the model is perfectly adapted to the training data and cannot predict the new inputs of the test data.

5(d)

The curve with the weights for the smallest gamma looks like a wavelet. While it has many weights that are very negative, there are also many weights that are very large so it does not limit the effective model complexity much. This means that many weights decay to zero but many are also scaled up by a large amount. This leads to the corresponding plot overfitting because its model complexity is not limited.

The curve with the optimal weights looks similar to the curve in 5(b), which has 19 basis functions and a gamma of 10^{-9} . The optimal choice of gamma decays the weights that are not supported by the data and effectively limits the model complexity. This leads to the corresponding plot fitting nicely.

The curve with the weights for the largest gamma looks like a parabola that has been flipped over the x-axis. Its weights are very small and this limits the effective model complexity to being linear like. This leads to the corresponding plot underfitting.

Question 6

6(a)

Question 6(a)

Derive $\frac{d\tilde{\ell}(w, u)}{dw}$

$$= \frac{d\ell(w, u)}{dw} + \frac{d\left(\gamma \sum_{j=1}^M w_j^2\right)}{dw} \quad (\text{by (6)})$$

$$= 2Z^T(y-t) + \gamma 2w$$

by 3(c)
and
by Aside

Aside

$$\frac{d\left(\gamma \sum_{j=1}^M w_j^2\right)}{dw_j} = \gamma 2w_j$$

$$\Rightarrow \frac{d\left(\gamma \sum_{j=1}^M w_j^2\right)}{dw} = \gamma 2\vec{w}$$

6(b)

Question 6(b)Derive $\frac{d \tilde{\ell}(w_0, u)}{dw_0}$

$$= \frac{d \ell(w_0, u)}{dw_0} + \frac{d \left(\gamma \sum_{j=1}^M w_j^2 \right)}{dw_0} \quad (\text{by (6)})$$

$$= 2 \vec{1}^T (y - t) + \gamma \frac{d \sum_{j=1}^M w_j^2}{dw_0} \quad (\text{by 3(a)})$$

$$= 2 \vec{1}^T (y - t) + 0 \quad (\text{by Aside})$$

$$\gamma \frac{d \sum_{j=1}^M w_j^2}{dw_0} \quad \text{Aside:}$$

$$= \gamma \frac{d (w_1 + w_2 + \dots + w_M)}{dw_0}$$

$$= \gamma (0)$$

$$= 0$$