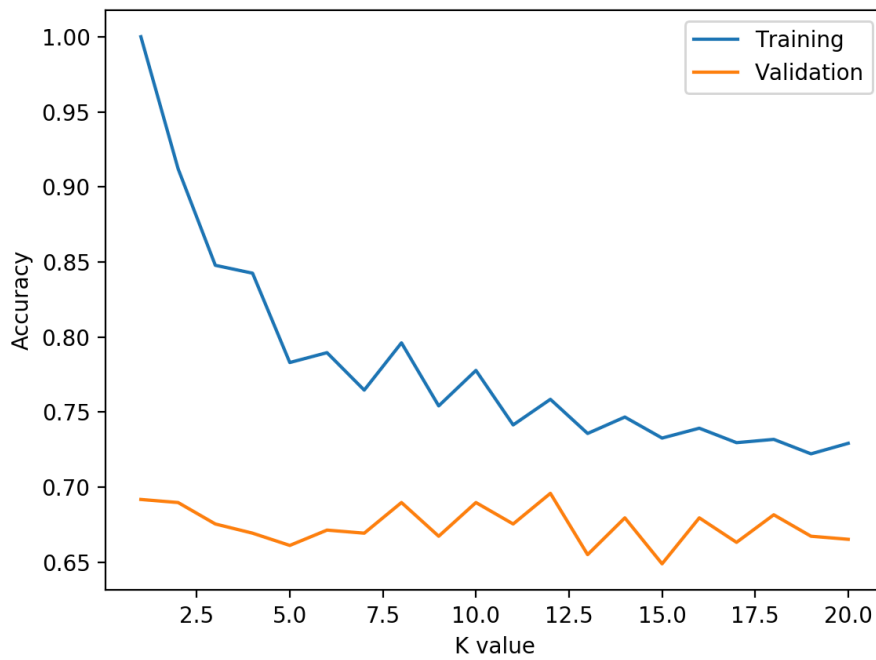
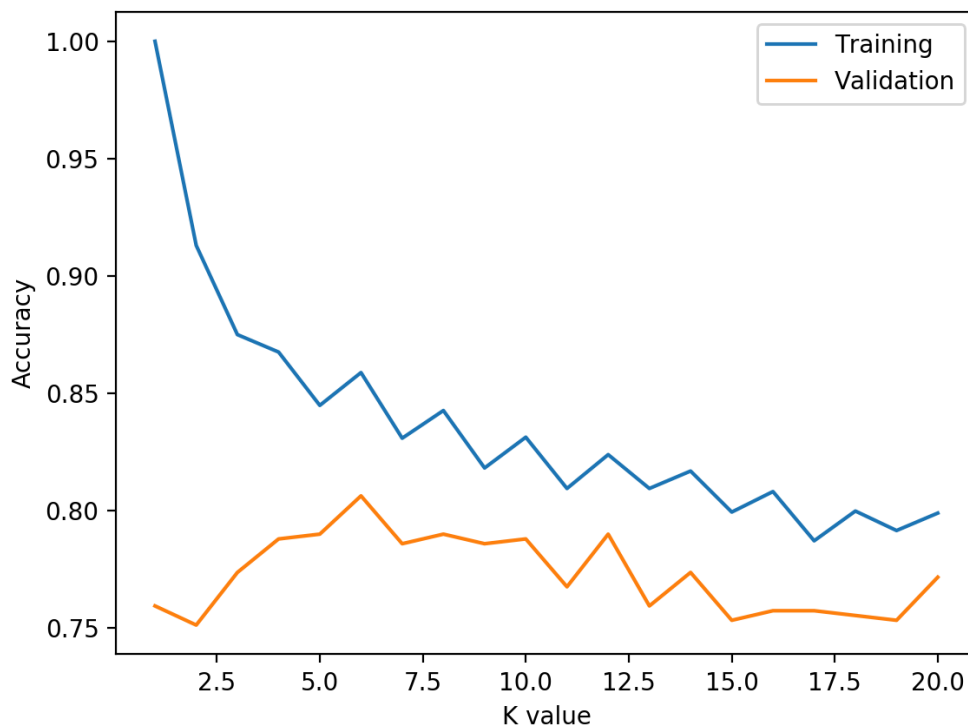


Question 1.
(b).



The k value of the best model is: 12. Its accuracy on the test set is 0.6306122448979592.

(c).



Claim: the cosine similarity performs better than Euclidean metric in text classification problems.

Consider the example dataset ['cat', 'bulldozer', 'cat cat cat'],

which is represented by the following matrix after vectorization.

$$A = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 3 & 0 \end{pmatrix}$$

Euclidean metric:

$$\text{Distance between 'cat', 'bulldozer'} = \sqrt{(1-0)^2 + (0-1)^2} = \sqrt{2}$$

$$\text{Distance between 'cat cat cat', 'cat'} = \sqrt{(3-1)^2 + 0^2} = 2$$

Hence 'cat' and 'bulldozer' are closer,
which is obviously wrong.

Cosine Similarity:

$$\cos(\text{angle between 'cat', 'bulldozer'}) = 0$$

$$\cos(\text{angle between 'cat', 'cat cat cat'}) = 1$$

The result is correct.

Since the distance between words is large because of the # of it appears, in this case, the cosine similarity can still give a smaller angle between them.

Question 2.

(a).

$$w_j \leftarrow (1 - \alpha \beta_j) w_j - \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_j^{(i)}$$

$$b \leftarrow \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)})$$

Because we subtract $\alpha \beta_j^T \vec{w}$ every time when we update \vec{w} .

Details:

$$\begin{aligned}\frac{\partial J_{\text{reg}}^{\beta}}{\partial \omega_j} &= \frac{\partial}{\partial \omega_j} \left[\frac{1}{2N} \sum_{i=1}^N \left(\sum_{k=1}^D \omega_k x_k^{(i)} + b - t^{(i)} \right)^2 \right] + \frac{\partial}{\partial \omega_j} \left[\frac{1}{2} \sum_{i=1}^D \beta_i \omega_i^2 \right] \\ &= \frac{1}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) \cdot \frac{\partial}{\partial \omega_j} [\omega_j x_j^{(i)}] + \frac{\partial}{\partial \omega_j} \left[\frac{1}{2} \beta_j \omega_j^2 \right] \\ &= \frac{1}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) \cdot x_j^{(i)} + \beta_j \omega_j\end{aligned}$$

$$\begin{aligned}\text{so } \omega_j &\leftarrow \omega_j - \alpha \frac{\partial J_{\text{reg}}^{\beta}}{\partial \omega_j} \\ &\leftarrow (1 - \alpha \beta_j) \omega_j - \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) \cdot x_j^{(i)}\end{aligned}$$

$$\frac{\partial J_{\text{reg}}^{\beta}}{\partial b} = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)})$$

$$\begin{aligned}b &\leftarrow b - \alpha \frac{\partial J_{\text{reg}}^{\beta}}{\partial b} \\ &\leftarrow b - \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)})\end{aligned}$$

(b).

By (a).

$$\begin{aligned}\frac{\partial J_{\text{reg}}^{\beta}}{\partial \omega_j} &= \frac{1}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) \cdot x_j^{(i)} + \beta_j \omega_j \\ &= \frac{1}{N} \sum_{i=1}^N \left(\sum_{j'=1}^D x_{j'}^{(i)} \omega_{j'} - t^{(i)} \right) x_j^{(i)} + \beta_j \omega_j \\ &= \left(\sum_{j'=1}^D \frac{1}{N} \sum_{i=1}^N x_{j'}^{(i)} x_j^{(i)} \omega_{j'} \right) + \beta_j \omega_j - \frac{1}{N} \sum_{i=1}^N x_j^{(i)} \cdot t^{(i)} = 0\end{aligned}$$

so

$$A_{jj'} = \begin{cases} \frac{1}{N} \sum_{i=1}^N (x_j^{(i)} x_{j'}^{(i)}) + \beta_j & \text{if } j=j' \\ \frac{1}{N} \sum_{i=1}^N x_j^{(i)} x_{j'}^{(i)} & \text{if } j \neq j' \end{cases}$$

$$c_j = \frac{1}{N} \sum_{i=1}^N x_j^{(i)} \cdot t^{(i)}$$

(c). By (b). Clearly:

$$A = \frac{1}{N} X^T X + \text{diag}(\beta_1, \beta_2, \dots, \beta_D)$$

$$\vec{c} = \frac{1}{N} X^T \vec{t}$$

Since $A\vec{w} - \vec{c} = 0$

$$\vec{w} = A^{-1} \vec{c}$$

$$= \left[\frac{1}{N} X^T X + \text{diag}(\beta_1, \dots, \beta_D) \right]^{-1} \cdot \frac{1}{N} X^T \vec{t}$$

$$\vec{w} = \left[X^T X + N \text{diag}(\beta_1, \dots, \beta_D) \right]^{-1} X^T \vec{t}$$

Question 3.

$$y = X\vec{w} + b \cdot \vec{1}$$

$$\frac{\partial J}{\partial y} = \frac{\partial}{\partial y} \left[\frac{1}{N} \sum_{i=1}^N 1 - \cos(y^{(i)} - t^{(i)}) \right]$$

$$= \frac{1}{N} \sin(\vec{y} - \vec{t})$$

$$\frac{\partial J}{\partial \vec{w}} = \frac{\partial J}{\partial \vec{y}} \frac{\partial \vec{y}}{\partial \vec{w}},$$

$$\text{by } \frac{\partial y^{(i)}}{\partial w_j} = x_j^{(i)},$$

$$= \frac{1}{N} X^T \cdot \text{sm}(\vec{y} - \vec{t})$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial \vec{y}} \frac{\partial \vec{y}}{\partial b}$$

$$\text{by } \frac{\partial y^{(i)}}{\partial b} = 1,$$

$$= \frac{1}{N} [\text{sm}(\vec{y} - \vec{t})]^T \vec{1}$$

Question 4.

(b).

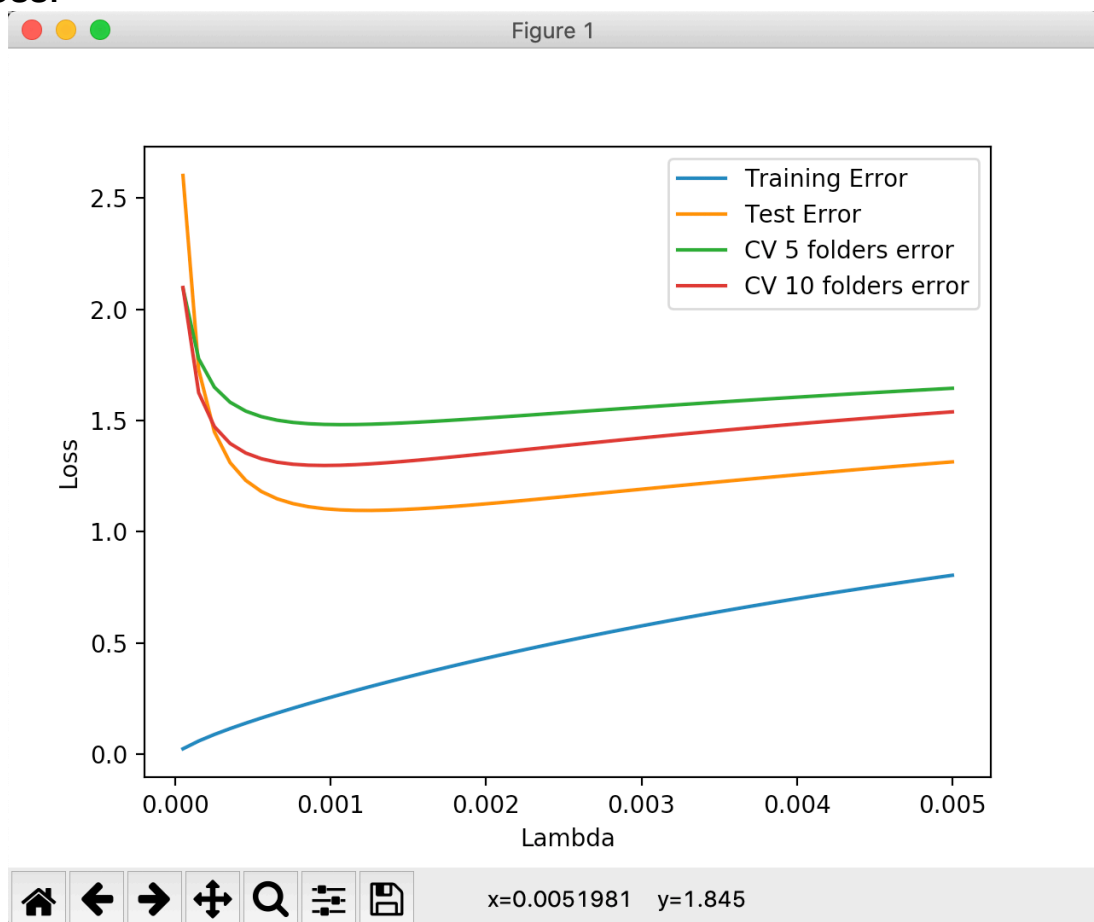
Order:

Firstly, `shuffle_data(data) + split_data(data, num_folders, fold)` to preprocess the dataset.

Then train our model by calling `train_model(data, lambda)`.

Next, predict the result of the training set by `predict(data, model)` and use `loss()` to compute the loss.

(c) & (d).



```
Best lambda by 5-fold CV: 0.0009591836734693879
Best lambda by 10-fold CV: 0.001060204081632653
```

The training error increases as the coefficient of L^2 penalty becoming greater.
The cross validation errors reach to their local minimum near $\lambda = 0.001$, which is the proposed parameter.

Notice the initial loss is very big if we set the penalty to 0 (i.e. $\lambda=0$), that's how regularization works.