

1.(a) Word-sequence kernel function is the kernel function could deal with this task. Word-sequence kernel function is a upgrade version of string kernel that propose the use of sequences of words rather than characters. It inherits the advantage of string kernel that finite sequences of symbols that need not be of the same length. Thus, Word-sequence kernel could deal with the condition that each document is an arbitrary length sequence of words taken from a vocabulary W .

The formula of Word-sequence Kernel is:

$$F_n(x, y) = \sum_{i=1}^n \mu^{(i)} \hat{K}_i(x, y).$$

$$\phi_u(x) = \sum_{i: x_i \neq u} \lambda_i x_i$$

For computing the similarity between pairs of documents d_i and d_j , this kern function will find the more similar two sequences a and b are, the higher the value of a sequence kernel $F_n(a, b)$ will be. Also, this kernel function can convert the tuples of word counts into a sequence. Therefore, Word-sequence kernel function, K_{ij} can computes the similarity between pairs of documents d_i and d_j .

Word-sequence Kernel

(b) We can use Multiple kernel learning (function) to solve this problem. Multiple kernel learning uses a predefined set of kernels and learn an optimal linear or non-linear combination of kernels. So, for this case, we can make a combination of image kernel and string kernel, to handle this problem.

The formula of Multiple kernel is $K' = \sum_{i=1}^h B_i K_i$, where B is a vector of coefficients for each kernel.

$$K(x, y) = \sum_{s=a}^a \sum_{t=b}^b w(s, t) f(x-s, y-t).$$

String Kernel.

If $K(x, y) = \varphi(x) \cdot \varphi(y)$ with φ mapping the arguments into an inner product space.

For this task, we may save the image kernel to deal with the image data. After this step, we can use the string kernel to handle the text data. Instead of creating a new kernel, multiple kernel algorithms can be used to combine kernels already established for each individual data source. So, Multiple Kernel function is the kernel function to solve this problem. To calculate the similarity, multiple kernel function will compare the similarity between two combinations of kernel functions. So, This should be a multiple kernel function.

Multiple Kernel function

(Image Kernel + String Kernel)

2. the original Stochastic Gradient descent function is

$$w_i \leftarrow -\lambda w_i + c(1-\lambda) y_i x_{ij} I[y_i > 1] I[y_i (w \cdot x_i) + b < 1]$$
$$b \leftarrow b + c(1-\lambda) y_i I[y_i (w \cdot x_i) + b < 1]$$

For this case, we have two definitions for each class.

$$w_{j0} \leftarrow -\lambda w_{j0} + c_0(1-\lambda) y_i x_{ij} I[y_i (w \cdot x_i) + b_0 < 1]$$

$$b_0 \leftarrow b_0 + c_0(1-\lambda) y_i I[y_i (w \cdot x_i) + b_0 < 1]$$

$$w_{i1} \leftarrow -\lambda w_{i1} + c_1(1-\lambda) y_i x_{ij} I[y_i (w \cdot x_i) + b_1 < 1]$$

$$b_1 \leftarrow b_1 + c_1(1-\lambda) y_i I[y_i (w \cdot x_i) + b_1 < 1]$$

In original formula, both classes used same w and b , but right now the costs of misclassification of the two classes are unequal, that's why each class has their own function.

To Balance the unequal cost, we can adjust bias and weight to make the margin becomes closer. Since we have three distinct margin. (one for overall Margin, one for c_0 , and one for c_1). We can see the gap between each margin, and adjust the bias and weight to make the cost of misclassification for each class are close to each other. Finally, the cost will be same after many times of adjusting bias and weight, and we should use the overall margin to justify the cost. Above its how to make the cost of misclassification of two classes become equal.

(Adjust the bias and weights separately for each class.)

3. (a) The choice of the activation function Z_{jp} satisfy the requirements for universal function approximation theorem. $Z_{jp} = \frac{1}{\delta^2 + h_{jp}}$. Where δ is a constant and $h_{jp} = \sum_i w_{ji} x_{ip}$. The requirements for UFAT is a (1) non-constant (2) non-linear (3) monotone (4) continuous function. For (1), Z_{jp} is clearly a non-constant because the value of Z_{jp} is always changing depending on δ and h_{jp} . For (2), Obviously, Z_{jp} is not a linear function because there are δ^2 and h_{jp}^2 . For (3), Z_{jp} is a monotone function as well because it depend on δ^2 and h_{jp}^2 , that without change a straight line immediately. For (4), This is obviously a continuous function. Thus, the activation function $Z_{jp} = \frac{1}{\delta^2 + h_{jp}}$ satisfy all the requirements for UFAT.

Yes, it satisfies the requirements of UFAT

$$(b) E_a = \frac{1}{2} \sum_{p=1}^P (d_p - o_{ip})^2 = \frac{1}{2} \sum_{p=1}^P (d_p - \sum_j w_{ji} \frac{1}{\delta^2 + (\sum_i w_{ij} x_{ip})^2})^2$$

To update w_{ji} , we need the update equation for input to hidden units.

$$\begin{aligned} \frac{\partial E_p}{\partial w_{ji}} &= \sum_{p=1}^P \frac{\partial E_p}{\partial o_{ip}} \frac{\partial o_{ip}}{\partial w_{ji}} = \sum_{p=1}^P \frac{\partial E_p}{\partial o_{ip}} \frac{\partial o_{ip}}{\partial Z_{jp}} \cdot \frac{\partial Z_{jp}}{\partial h_{jp}} \cdot \frac{\partial h_{jp}}{\partial w_{ji}} \\ &= \sum_{p=1}^P \frac{\partial}{\partial o_{ip}} \left[\frac{1}{2} \sum_{p=1}^P (d_p - o_{ip})^2 \right] (u_j) (Z_{jp}) (1 - Z_{jp}) (x_{ip}) \\ &= - \sum_{p=1}^P (d_p - o_{ip}) (u_j) (Z_{jp}) (1 - Z_{jp}) (x_{ip}) \\ &= - \underbrace{\left(\sum_{p=1}^P \delta_p (u_j) \frac{1}{\delta^2 + (\sum_i w_{ij} x_{ip})^2} \right)}_{\delta_{ip}} (1 - \frac{1}{\delta^2 + (\sum_i w_{ij} x_{ip})^2}) (x_{ip}) \\ &= - \delta_{ip} x_{ip} \end{aligned}$$

$w_{ji} \leftarrow w_{ji} + n \delta_{ip} x_{ip}$

Above is the update equation for w_{ji} that minimize E_a .

Now see the update equation for u_j that Hidden-to-output.

$$\frac{\partial E_p}{\partial u_j} = \frac{\partial E_p}{\partial h_{jp}} \cdot \frac{\partial h_{jp}}{\partial u_j} \quad \frac{\partial h_{jp}}{\partial u_j} = Z_{jp}$$

$$\frac{\partial E_p}{\partial h_{jp}} = \frac{\partial E_p}{\partial o_{ip}} \cdot \frac{\partial o_{ip}}{\partial h_{jp}} = -(d_p - o_{ip})(1)$$

$$u_j \leftarrow u_j - n \frac{\partial E_p}{\partial u_j} = u_j + (d_p - o_{ip}) Z_{jp} = u_j + \delta_{ip} Z_{jp}$$

$u_j \leftarrow u_j + \delta_p \left(\frac{1}{\sigma^2 (\sum_i w_{ij} x_{ip})^2} \right)$, where δ can be found from the previous part that update equation for w_{ji} .

Above is the update equation for u_j that minimize E_a .

4. (a) Since $E_b = \sum_{p=1}^P \sum_{i=0}^N \left(\frac{\partial E_a}{\partial x_{ip}} \right)^2$, the error function $E = \lambda E_a + (1-\lambda) E_b$ where $0 \leq \lambda \leq 1$ will never be decreased. To the impact of minimizing E_b on the sensitivity of the network output to relatively small amounts of noise in the input sample, the sensitivity to the noise will be small, which means the network output is insensitive to noise because E_b makes the weight not that accurate for real results. In other words, E_b makes the increase of the robustness of this neural network. Overall, the sensitivity to small amounts of noise is small because of the impact of E_b .

The tendency of the network to over-fit the training data will be decreased (slow) as well because of a similar reason that E_b makes the weight not that accurate for desired output. In other words, it generalizes the capability of the network, and it slows down the tendency of the network to over-fit the training data (avoid the risk of over-fitting).

Above are the impact of minimizing E_b on the sensitivity of the network and the tendency of the network to over-fit the training data.

$$(b) E = \lambda E_a + (1-\lambda) E_b \text{ where } 0 \leq \lambda \leq 1.$$

$$E_a = \frac{1}{2} \sum_{p=1}^P (d_p - o_p)^2 = \frac{1}{2} \sum_{p=1}^P \left[d_p - \sum_j u_j \frac{1}{(b^2 + (\sum_i w_{ij} x_{ip})^2)} \right]^2$$

$$E_b = \sum_{p=1}^P \sum_{i=0}^N \left(\frac{\partial E_a}{\partial x_{ip}} \right)^2$$

We will calculate the error function separately. For E_a , we have already derive the update equations in 3(b), so here only complete the update equation for E_b .

$$\frac{\partial E_a}{\partial x_{ip}} = \frac{\sum_{p=1}^P (d_p - o_p)}{\partial x_{ip}} = \sum_{p=1}^P \left[d_p - \sum_j u_j \frac{1}{(b^2 + (\sum_i w_{ij} x_{ip})^2)} \right]$$

$$E_b = \sum_{p=1}^P \sum_{i=0}^N \left[d_p - \sum_j u_j \frac{1}{(b^2 + (\sum_i w_{ij} x_{ip})^2)} \right]^2$$

$$\frac{\partial E_b}{\partial w_{ji}} = -2\lambda \left[d_p - \sum_j w_{ji} \frac{1}{(G^2 + \sum_i (w_{ij}))} \right]$$

Thus $w_{ji} \leftarrow w_{ji} + \gamma_{dp} [d_p - 2\lambda \left[d_p - \sum_j w_{ji} \frac{1}{(G^2 + \sum_i (w_{ij}))} \right]]$

and w_j is

$$w_j \leftarrow w_j - \frac{\partial E_b}{\partial w_j} = w_j - \lambda \left[d_p - \sum_i w_{ji} \frac{1}{(G^2 + \sum_i (w_{ij}))} \right] + \beta$$

Above is the update equation for the parameters w_j and w_j for E_b .

To minimize E , we can simply add the update equations of E_a and E_b together, and calculate the final error. However, the logic behind this is very similar to the 3(b), the update equations of E_a .

5.(a) Since the samples are approximately linearly separable, we can use Perceptron learning algorithm to classify this task. Perceptron is an algorithm that attempts to fix all errors encountered in the training set, and divides the data into two classes. For the data that is non-linearly separable, we can use linear kernel to make them become linearly separable. However, I think Perceptron learning algorithm is the best choice for this task.

(b) Naive-Bayes classifier is the best choice for this task, since the features are likely to be independent. Naive Bayes classifier is a simple probabilistic classifier based on applying Bayes' theorem with strong independence assumptions between the features. Thus, NB fits this requirement. Additionally, the dataset has a limited number of training examples, however NB can deal with this situation too.

Naive Bayes

(c) I recommend neural network is the best algorithm for this task. First, neural networks can do that the features are numeric; Second, the robustness of neural networks is good, and introduce small amounts of noise in the weight update, can help neural network reduces the risk of over-fitting. Thus, I recommend neural network for this task.

Neural Network

(d) To solve this problem, we can chose the combination of sparse regularization, universal function approximation theorem, and finally, neural network as our learning algorithm. Since the number of features far exceeds the number of training example, we should make the features become sparsely regularized, and select a couple of the most important features based on our prior knowledge of features. After this, the number of features becomes small, and we can use UFACT to compute all the eligible

functions into neural network classifier. Thus, we will use neural network to do normal learning since we have already known prior changes to guide the selection of features to be used to train a neural network function approximator.

Sparse regularization + UFAT + Neural Network

(e) I recommended neural network to do this task. First, we can divided the data into three subgroups based on their types (namely, Text, Images). Since neural network will have a weight for each input, thus, we donot need to worry about the weight problem. Finally, neural networks will combine them into a Net input function. So neural network is the algorithm to handle the composed of three types of data.

Neural Network