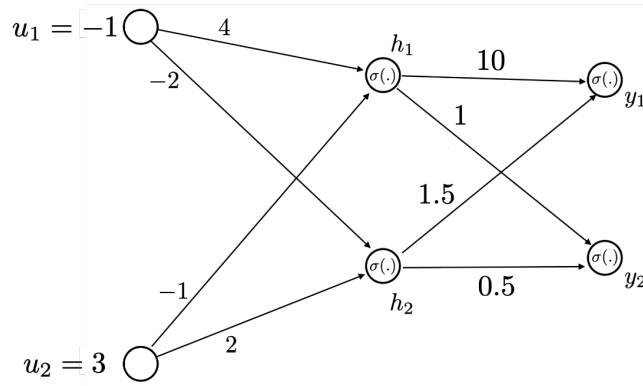


# CS/ECE/ME 532

## Unit 6 Practice Problems

1. **Neural Network Basic.** Consider the following neural network.



**Figure 1:** A Neural Network Architecture

- a) How many hidden layers are in this network? How many neurons are in this network? How many outputs does this network have? What is  $\sigma(\cdot)$  called?

**SOLUTION:** One hidden layer/6 neurons including the input neurons/2 outputs/activation function

- b) Assume that the neural network shown above uses an activation called *Leaky ReLU*, defined as follows:

$$\sigma(z) = \begin{cases} z & \text{if } z > 0 \\ 0.01z & \text{otherwise.} \end{cases}$$

Note that Leaky ReLU behaves differently from ReLU when the input value is less than 0. The values labeled  $h_1, h_2, y_1, y_2$  are the outputs of the corresponding nodes. Find the numerical values of  $h_1, h_2, y_1, y_2$  for an input feature  $u_1 = -1, u_2 = 3$ .

**SOLUTION:**

$$h_1 = \sigma(-4 - 3) = \sigma(-7) = -0.07$$

$$h_2 = \sigma(2 + 6) = \sigma(8) = 8$$

$$y_1 = \sigma(-0.7 + 12) = \sigma(11.3) = 11.3$$

$$y_2 = \sigma(-0.07 + 4) = \sigma(3.93) = 3.93.$$

- 2. SGD to learn the weights for a single neuron** We use the single neuron shown in the figure for classification. Here  $x_j^i$  is the  $j$ -th feature in the  $i$ -th training sample and the output is  $\hat{y}^i = \sigma\left(\sum_{j=0}^P w_j x_j^i\right)$ , where the activation function is *ReLU6*, defined as follows:

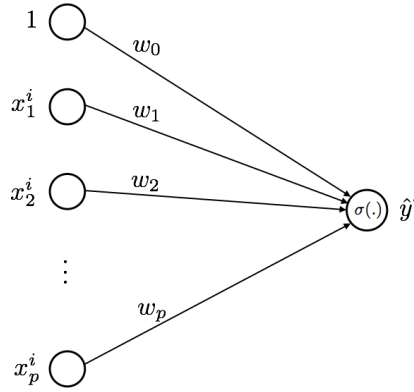
$$\sigma(z) = \min(\max(0, z), 6).$$

Note that ReLU6 behaves differently from ReLU when the input value is larger than 6.

Now suppose we use ridge regression for the loss function

$$f^i(\mathbf{w}) = \frac{1}{2}(\hat{y}^i - y^i)^2 + \lambda \sum_{j=0}^P w_j^2.$$

Derive the gradient for the update step,  $\nabla f^{it}(\mathbf{w}^{(t)})$  and write the update equation for  $\mathbf{w}^{(t+1)}$ .



**Figure 2:** A Neuron

**SOLUTION:** Let  $x_0^i = 1$  for all  $i$ . Note that

$$\frac{df^{it}(\mathbf{w}^{(t)})}{dw_k} = (\hat{y}^{it} - y^{it}) \frac{d\hat{y}^{it}}{dw_k} + 2\lambda w_k.$$

To compute  $\frac{d\hat{y}^{it}}{dw_k}$ , we use the chain rule:

$$\frac{d\hat{y}^{it}}{dw_k} = \frac{d\hat{y}^{it}}{dz^{it}} \frac{dz^{it}}{dw_k},$$

where  $z^{it} = \sum_{j=0}^P w_j x_j^{it}$ . By the definition of ReLU6, it is clear  $\frac{d\hat{y}^{it}}{dz^{it}} = \mathbb{1}\{0 \leq z^{it} \leq 6\} = \mathbb{1}\{0 \leq \sum_{j=0}^P w_j x_j^{it} \leq 6\}$ , where  $\mathbb{1}(\text{condition}) := \begin{cases} 1 & \text{if condition is true} \\ 0 & \text{otherwise} \end{cases}$ . Also, it is clear  $\frac{dz^{it}}{dw_k} = x_k^{it}$ . Thus,

$$\frac{d\hat{y}^{it}}{dw_k} = x_k^{it} \cdot \mathbb{1}\{0 \leq \sum_{j=0}^P w_j x_j^{it} \leq 6\}.$$

Define  $\delta^{it} = (\hat{y}^{it} - y^{it}) \cdot \mathbb{1}\{0 \leq \sum_{j=0}^P w_j x_j^{it} \leq 6\}$  so we can write

$$\nabla f^{it}(\mathbf{w}^{(t)}) = \delta^{it} \mathbf{x}^{it} + 2\lambda \mathbf{w}^{(t)}$$

and thus the update is

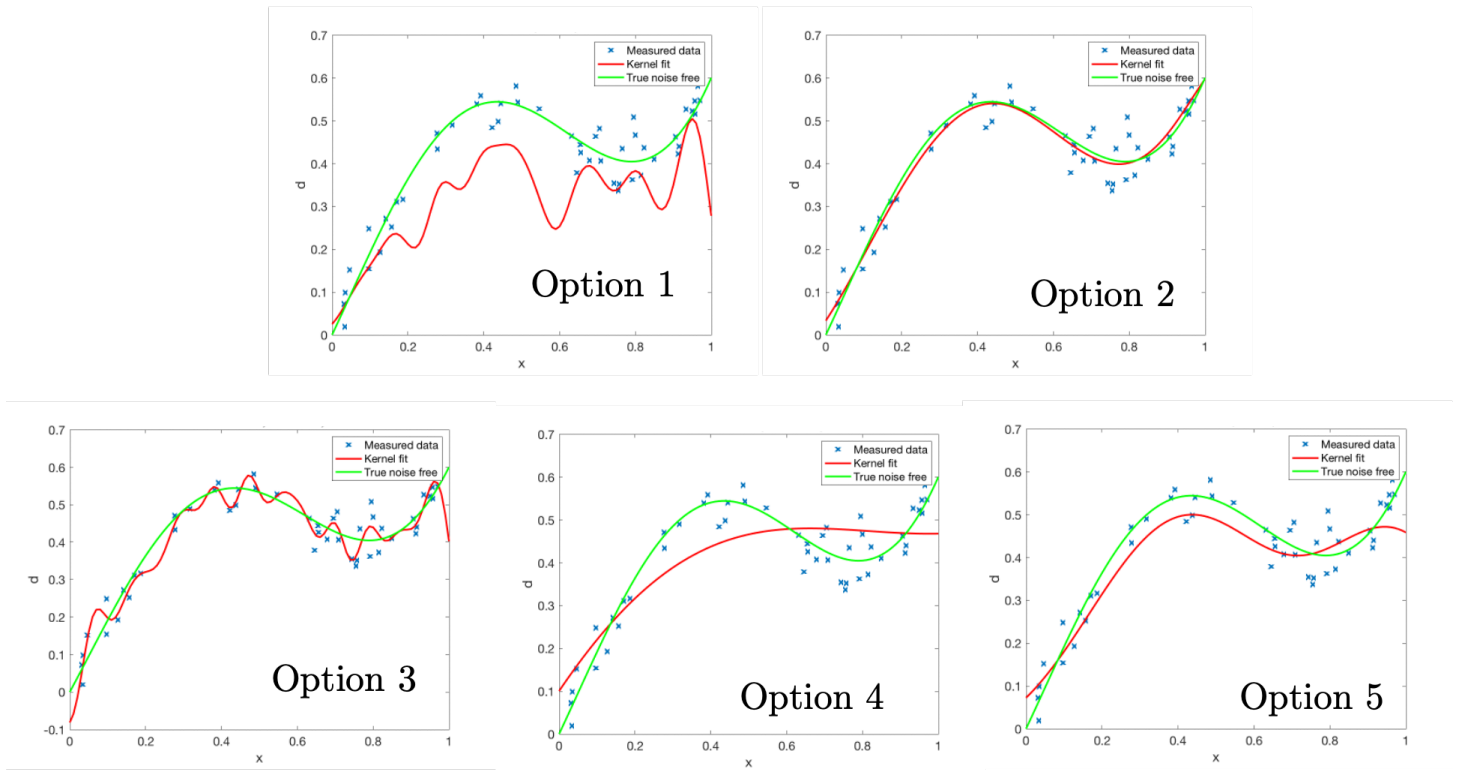
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \alpha_t \delta^{it} \mathbf{x}^{it} - 2\alpha_t \lambda \mathbf{w}^{(t)}$$

**3. Kernel Regression.** Consider the Gaussian kernel regression with  $\ell_2$  regularization. The first hyperparameter is  $\sigma$  for the Gaussian kernel  $K(\mathbf{u}, \mathbf{v}) = \exp(-\|\mathbf{u} - \mathbf{v}\|_2^2 / (2\sigma^2))$ . The second hyperparameter is  $\lambda$  which is the weight for the  $\ell_2$  regularization term.

Now, consider the following five choices:

- a)  $\lambda = 0.01, \sigma = 0.04$
- b)  $\lambda = 0.01, \sigma = 0.2$
- c)  $\lambda = 0.01, \sigma = 1$
- d)  $\lambda = 1, \sigma = 0.04$
- e)  $\lambda = 1, \sigma = 0.2$

Find the correct mapping between each of these choices and each of the following options (each row corresponds one choice).



**Figure 3:** Kernel regression results with five different hyperparameters

**SOLUTION:** Let us first consider  $\lambda$ .  $\lambda$  is associated with the 2-norm penalty in the ridge regression problem, so as  $\lambda$  increases, the 2-norm penalty receives increasing emphasis. Hence  $\alpha$  for the  $\lambda = 1$  case will have smaller norm than for the  $\lambda = 0.01$  case. This results in underestimation (and a slightly smoother solution) when  $\lambda = 1$ . From this, one can see Option 1 and Option 5 correspond to  $\lambda = 1$ , and Option 2, 3, 4 correspond to  $\lambda = 0.01$ .

Let's now consider  $\sigma$ , which determines the width of the kernels. Larger values of  $\sigma$  result in wider kernels, which will lead to a smoother solution. Smaller values of  $\sigma$  results in narrower kernels, fully overfitting the measured data. Between Option 1 and Option 5, Option 5's fit is smoother, so

Option 5 corresponds to case e) and Option 1 corresponds to case d). Among those with  $\lambda = 0.01$ , Option 4 has the smoothest fit, and Option 3 fully overfits the data, and Option 2 is something in between them. Thus, Option 4 corresponds to case c), Option 3 corresponds to case a), and Option 2 corresponds to case b).

4. **Kernel Classification.** Consider the kernel classification with the squared error loss using the Gaussian kernel  $K(\mathbf{u}, \mathbf{v}) = \exp(-\|\mathbf{u} - \mathbf{v}\|_2^2 / (2\sigma^2))$  for the following dataset.

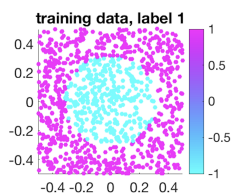


Figure 4: Binary classification dataset

Now, consider the following three choices:

- a)  $\sigma = 5$
- b)  $\sigma = 0.05$
- c)  $\sigma = 0.005$

Find the correct mapping between each of these choices and each of the following options (each row corresponds one choice).

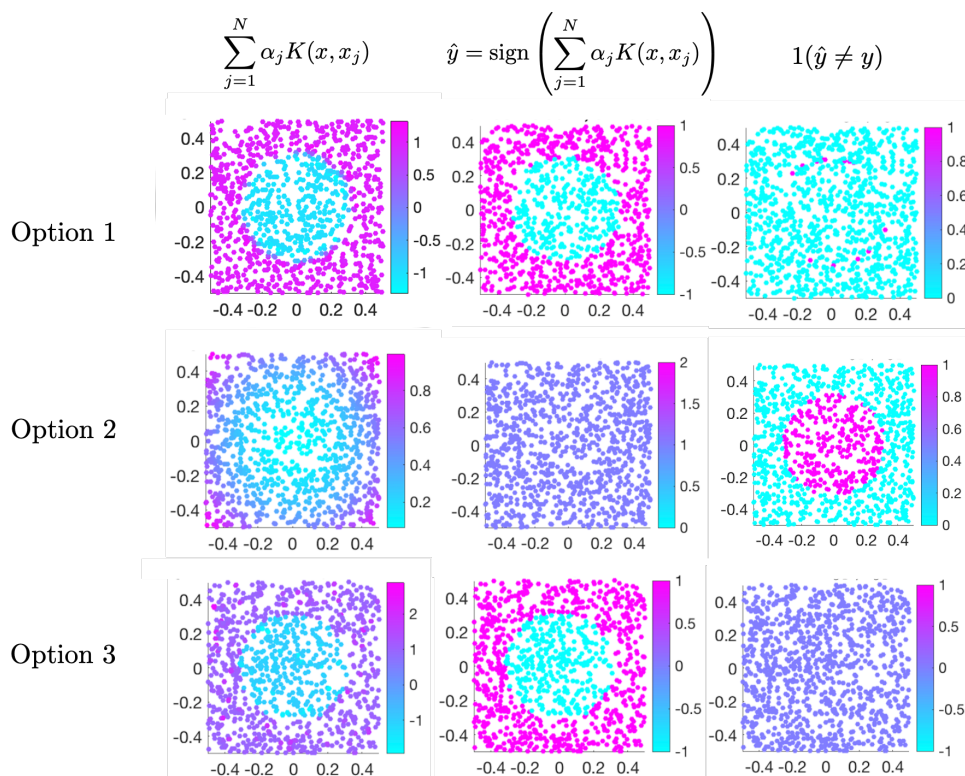


Figure 5: Kernel classification results with three different options

**SOLUTION:**  $\sigma$  controls the smoothness of the decision boundary by controlling the “width” of the kernels. Wide kernels (large  $\sigma$ ) result in smooth decision boundaries. Narrow kernels (small  $\sigma$ ) result in more complex decision boundaries, making a fewer mistakes on train set. Thus, a) corresponds to Option 2, b) corresponds to Option 1, and c) corresponds to Option 3.

**5. Kernel Classification.** When using a Gaussian kernel, is there a downside to choosing a very small value for  $\sigma$ ?

**SOLUTION:** Overfitting is very likely. Each of the kernels tends to contribute to a single measurement ( $\mathbf{K}$  is nearly diagonal), which results in perfect classification, but would also result in very poor performance generalizing to data not used to train the classifier.

- 6. Kernel SVM.** You use a kernel-based support vector machine for binary classification with labels  $d^i = \{+1, -1\}$ . Given training features and labels  $(\mathbf{x}^i, d^i), i = 1, 2, \dots, N$  you use a kernel  $K(\mathbf{u}, \mathbf{v})$  and design the classifier weights  $\boldsymbol{\alpha}$  as

$$\hat{\boldsymbol{\alpha}} = \arg \min_{\boldsymbol{\alpha}} \sum_{i=1}^N \left( 1 - d^i \sum_{j=1}^N \alpha_j K(\mathbf{x}^i, \mathbf{x}^j) \right)_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j K(\mathbf{x}^i, \mathbf{x}^j)$$

- a) Assume the optimization problem has been solved to obtain the weights  $\boldsymbol{\alpha}$ . Express the classification procedure for a measured feature  $\mathbf{x}$ .
- b) Suppose  $N = 1000$  and  $\alpha_i = 0, i = 1, 2, \dots, 99, 102, 103, \dots, 1000$ . Identify the support vectors and write the classification procedure in terms of the support vectors.

**SOLUTION:**

- a) Kernel regression gives  $\widehat{d(\mathbf{x})} = \sum_{i=1}^N \alpha_i K(\mathbf{x}, \mathbf{x}^i)$ . Therefore, the classification procedure is

$$\text{sign} \left\{ \sum_{i=1}^N \alpha_i K(\mathbf{x}, \mathbf{x}^i) \right\}$$

- b) The only non zero  $\alpha_i$  are  $\alpha_{100}$  and  $\alpha_{101}$ . Therefore,  $\mathbf{x}^{100}$  and  $\mathbf{x}^{101}$  are the support vectors. This means that the classification procedure is

$$\text{sign} \{ \alpha_{100} K(\mathbf{x}, \mathbf{x}^{100}) + \alpha_{101} K(\mathbf{x}, \mathbf{x}^{101}) \}$$