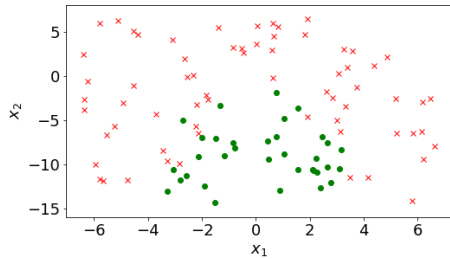


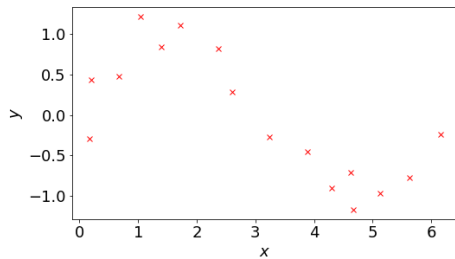
- Machine learning contains a large number of matrix multiplications. Now we need to calculate the product of three matrices A, B, C . Suppose the dimensions of A, B, C are $m \times n, n \times p, p \times q$ respectively, where $m < n < p < q$. Which of the following calculations is correct and also the most efficient.
 - $(AB)C$
 - $(AC)B$
 - $A(BC)$
 - $A(CB)$
- In Gradient Descent, what will happen if the learning rate is too small?
 - The model may not converge
 - The model may converge fast
 - The model may converge slowly
 - The model may converge properly
- In Gradient Descent, what will happen if the learning rate is too large?
 - The model may not converge
 - The model may converge fast
 - The model may converge slowly
 - The model may converge properly
- Given a linear function $\mathbf{y} = \mathbf{w}\mathbf{x} + \mathbf{b}$ and an activation function $\mathbf{y} = \sigma(\mathbf{x})$, what is a neuron?
 - $\mathbf{w}\sigma(\mathbf{x}) + \mathbf{b}$
 - $\sigma(\mathbf{w}\mathbf{x} + \mathbf{b})$
 - $\mathbf{w}\mathbf{x} + \mathbf{b} + \sigma(\mathbf{x})$
 - $\sigma(\mathbf{w}\mathbf{x}) + \mathbf{b}$
- Suppose $\mathbf{z} \in \mathbb{R}^D$ is the output of a model. Which of the following functions should you use to get the prediction $\hat{\mathbf{y}}$ if you want to do linear regression?
 - $\hat{\mathbf{y}} = \mathbf{z}$, where $D = 1$
 - $\hat{\mathbf{y}} = \sigma(\mathbf{z})$, where $\sigma(\cdot)$ is the Sigmoid function and $D = 1$
 - $\hat{\mathbf{y}} = \sigma(\mathbf{z})$, where $\sigma(\cdot)$ is the Softmax function and $D = 3$
 - $\hat{\mathbf{y}} = \mathbf{w}\mathbf{z} + \mathbf{b}$ and $D = 3$
- Suppose $\mathbf{z} \in \mathbb{R}^D$ is the output of a model. Which of the following functions should you use to get the prediction $\hat{\mathbf{y}}$ if you want to do 2-class classification?
 - $\hat{\mathbf{y}} = \mathbf{z}$, where $D = 1$
 - $\hat{\mathbf{y}} = \sigma(\mathbf{z})$, where $\sigma(\cdot)$ is the Sigmoid function and $D = 1$
 - $\hat{\mathbf{y}} = \sigma(\mathbf{z})$, where $\sigma(\cdot)$ is the Softmax function and $D = 3$
 - $\hat{\mathbf{y}} = \mathbf{w}\mathbf{z} + \mathbf{b}$ and $D = 3$
- Suppose $\mathbf{z} \in \mathbb{R}^D$ is the output of a model. Which of the following functions should you use to get the prediction $\hat{\mathbf{y}}$ if you want to do N -class classification?
 - $\hat{\mathbf{y}} = \mathbf{z}$, where $D = 1$
 - $\hat{\mathbf{y}} = \sigma(\mathbf{z})$, where $\sigma(\cdot)$ is the Sigmoid function and $D = 1$
 - $\hat{\mathbf{y}} = \sigma(\mathbf{z})$, where $\sigma(\cdot)$ is the Softmax function and $D = N$
 - $\hat{\mathbf{y}} = \mathbf{w}\mathbf{z} + \mathbf{b}$ and $D = N$
- In neural networks, which of the following parameters is not the hyperparameter you need to set.
 - The number of layers L
 - The number of neurons in the l layer $D^{[l]}$
 - The bias in the in the l layer $\mathbf{b}^{[l]}$

- D. The learning rate λ
9. Which of the following loss functions is the loss function for linear regression.
- $(\hat{y} - y)^2/2$
 - $-y \log \hat{y} - (1 - y) \log(1 - \hat{y})$
 - $-\mathbf{y}^T \log \hat{\mathbf{y}}$
 - Any one of the above functions
10. Which of the following loss functions is the loss function for logistic regression.
- $(\hat{y} - y)^2/2$
 - $-y \log \hat{y} - (1 - y) \log(1 - \hat{y})$
 - $-\mathbf{y}^T \log \hat{\mathbf{y}}$
 - Any one of the above functions
11. Which of the following loss functions is the loss function for Softmax regression.
- $(\hat{y} - y)^2/2$
 - $-y \log \hat{y} - (1 - y) \log(1 - \hat{y})$
 - $-\mathbf{y}^T \log \hat{\mathbf{y}}$
 - Any one of the above functions
12. Which of the following loss functions is the loss function for neural networks.
- $(\hat{y} - y)^2/2$
 - $-y \log \hat{y} - (1 - y) \log(1 - \hat{y})$
 - $-\mathbf{y}^T \log \hat{\mathbf{y}}$
 - Any one of the above functions
13. What is a linear function in 2-dimensional space.
- A line
 - A plane
 - A hyperplane
 - A point
14. What is a linear function in 3-dimensional space.
- A line
 - A plane
 - A hyperplane
 - A point
15. What is a linear function in 4-dimensional space.
- A line
 - A plane
 - A hyperplane
 - A point
16. In Perceptron, how to determine whether sample \mathbf{x} is mis-classified.
- See if $y(\mathbf{w}^T \mathbf{x} + b) < 0$
 - See if $(\mathbf{w}^T \mathbf{x} + b) < 0$
 - See if $y < 0$
 - See if $(\mathbf{w}^T \mathbf{x} + b) > 0$

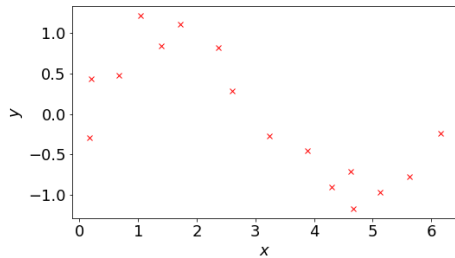
17. Given the following training set, which of the following answers is correct to train a classification model $\hat{y} = \sigma(w_0 + w_1x_1 + w_2x_2)$, where $\sigma(\cdot)$ is the Sigmoid function.



- A. The model has a high bias and a low variance
 - B. The model will be overfitting
 - C. The model has a low bias and a low variance
 - D. The model has a low bias and a high variance
18. Given the following training set, which of the following answers is correct to train a regression model $\hat{y} = w_0 + w_1x$.



- A. The model has a high bias and a high variance
 - B. The model will be underfitting
 - C. The model has a low bias and a low variance
 - D. The model has a low bias and a high variance
19. Given the following training set, which of the following answers is correct to train a regression model $\hat{y} = w_0 + w_1x + w_2x + w_3x + w_4x + w_5x + w_6x + w_7x + w_8x + w_9x$.



- A. The model has a high bias and a high variance
 - B. The model will be underfitting
 - C. The model has a low bias and a low variance
 - D. The model has a low bias and a high variance
20. Given a cost function with L^2 regularization $\mathcal{J}(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$, which of the following answers is correct when $\lambda = 0$?
- A. It is equivalent to the cost function $\mathcal{J}(\mathbf{w})$ without any regularization
 - B. It leads to a result that every $w_i \approx 0$
 - C. It successfully solves the problem of overfitting
 - D. It brings overfitting to the model

21. Given a cost function with L^2 regularization $\mathcal{J}(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$, which of the following answers is correct when $\lambda \rightarrow \infty$?
- A. It is equivalent to the cost function $\mathcal{J}(\mathbf{w})$ without any regularization
 - B. It leads to a result that every $w_i \approx 0$
 - C. It successfully solves the problem of underfitting
 - D. It brings overfitting to the model
22. Suppose the dataset contains two positive samples $\mathbf{x}^{(1)} = [2, 2]^T$ and $\mathbf{x}^{(2)} = [2, 4]^T$, and two negative samples $\mathbf{x}^{(3)} = [0, 0]^T$ and $\mathbf{x}^{(4)} = [-1, 0]^T$. Please calculate the SVM decision hyperplane