

# Final Review

Junyang Deng

January 27, 2023

1. For a hyperplane,  $g(x) = w^T x + w_0$ , please describe its meanings in regression, artificial neural networks, and support vector machine.

## SOLUTION

- (a) In regression, a hyperplane  $g(x)$  minimizes the difference between predicted  $y$  and true  $y$  (measured by MSE). Most of the time, the optimal hyperplane can be solved analytically by  $w = (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{D}^T \mathbf{r}$  where  $w$  is the hyperplane,  $\mathbf{D}$  is a matrix of  $x_1, \dots, x_n$ , and  $\mathbf{r}$  is the class label. When an analytical solution is unavailable, we can use gradient descent to search for the best plane. With the hyperplane, people can predict new entry with known variables.
  - (b) In ANN, each node (or called perceptron) includes a hyperplane. Each perceptron has inputs either from the environment or outputs of other perceptrons. For each input, the node can calculate a weighted sum of inputs by  $y = \sum_{j=1}^d w_j x_j + w_0$ ,  $w_0$  is the intercept value to make the model more general. This hyperplane divides the input space into two, which allows the perceptron to separate classes.
  - (c) In support vector machine, the hyperplane is solved by maximizing the margin between classes. To be specific, the task is to find  $\min \frac{1}{2} \|w\|^2$  subject to  $r^t (w^T x^t + w_0) \geq +1, \forall t$ . The resulting  $w$  can be used to construct hyperplane.
2. What is the “gradient descent” optimization algorithm? Describe its usages in machine learning approaches.

## SOLUTION

The “gradient descent” optimization algorithm is first construct a loss function  $E(w)$ , which is a differentiable function of a vector of variables. Then, for each parameter  $w$ ,

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}, \forall i$$
$$w_i = w_i + \Delta w_i$$

where  $\eta$  is called the learning factor, and determines how much to move in that direction. The use of a good value for  $\eta$  is critical; if it is too small, the convergence may be too slow, and a large value may cause oscillations and even divergence.

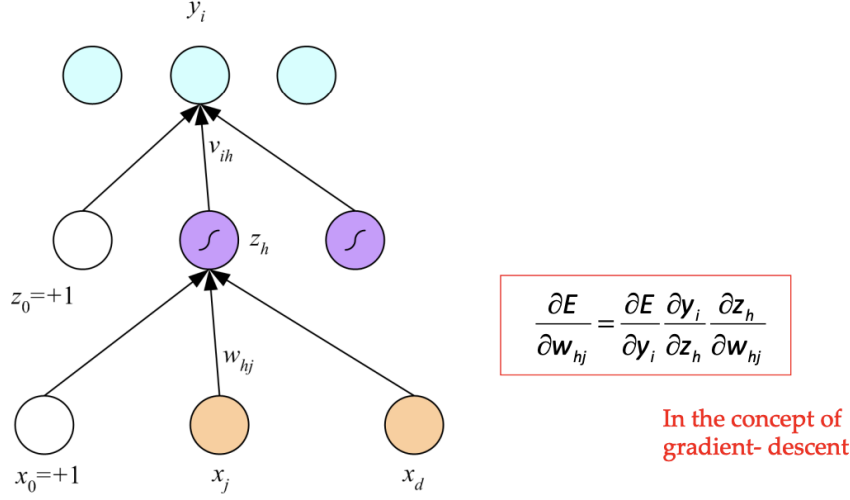
Gradient descent can be applied in many machine learning approaches, including linear regression, logistic regression and neural network.

3. Describe the backpropagation algorithms in multilayer perceptrons.

## SOLUTION

When we have a set of parameters, we can make predictions about classes or conduct regression based

on the given information. The predicted value may be different from true value, and such difference is measured by error function. In order to use this error to update parameters, we use the backpropagation algorithm. For a nonlinear regression, For example, with  $\frac{\partial E}{\partial w_{hj}} = \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial z_h} \frac{\partial z_h}{\partial w_{hj}}$ , error E propagates back to  $w$ . The general update rule is Update = LearningFactor  $\times$  (DesiredOutput - ActualOutput)  $\times$  Input.



$$\begin{aligned}
 \Delta w_{hj} &= -\eta \frac{\partial E}{\partial w_{hj}} \\
 &= -\eta \sum_t \frac{\partial E^t}{\partial y^t} \frac{\partial y^t}{\partial z_h^t} \frac{\partial z_h^t}{\partial w_{hj}} \\
 &= -\eta \sum_t \underbrace{-(r^t - y^t)}_{\partial E^t / \partial y^t} \underbrace{v_h}_{\partial y^t / \partial z_h^t} \underbrace{z_h^t (1 - z_h^t) x_j^t}_{\partial z_h^t / \partial w_{hj}} \\
 &= \eta \sum_t (r^t - y^t) v_h z_h^t (1 - z_h^t) x_j^t
 \end{aligned}$$

4. Describe “Kernel Trick” in support vector machine and present at least three different kernel functions.

#### SOLUTION

Kernel trick allows people to solve nonlinear problems without transformation. The idea is to replace the inner product of basis functions,  $\phi(\mathbf{x}^t)^T \phi(\mathbf{x}^s)$ , by a kernel function,  $K(\mathbf{x}^t, \mathbf{x}^s)$ , between instances in the original input space. Instead of mapping two instances  $\mathbf{x}^t$  and  $\mathbf{x}^s$  to the  $z$ -space and doing a dot product there, we directly apply the kernel function in the original space.

There are many different kernel functions:

- Dot product kernel:  $K(\mathbf{x}^t, \mathbf{x}) = \mathbf{x}^T \mathbf{x}^t$
- Polynomial (with degree q):  $K(\mathbf{x}^t, \mathbf{x}) = (\mathbf{x}^T \mathbf{x}^t)^q$
- Radial-basis function:  $K(\mathbf{x}^t, \mathbf{x}) = \exp \left[ -\frac{\|\mathbf{x}^t - \mathbf{x}\|^2}{2s^2} \right]$
- Mahalanobis kernel:  $K(\mathbf{x}^t, \mathbf{x}) = \exp \left[ -\frac{1}{2} (\mathbf{x}^t - \mathbf{x})^T \mathbf{S}^{-1} (\mathbf{x}^t - \mathbf{x}) \right]$
- Sigmoid function:  $K(\mathbf{x}^t, \mathbf{x}) = \tanh(2\mathbf{x}^T \mathbf{x}^t + 1)$

5. Classes are in non-linear space; how to make them linear in a new space. Propose methods and apply them when using linear discrimination and SVM.

**SOLUTION**

Transformation can map non-linear separable classes into linear separable ones. Specifically, we can define a set of nonlinear basis functions  $\Phi_{ij}$  to map the original data to a new space where the function can be written in a linear form.

After transformation, linear discrimination and SVM can be applied directly in the new space.

6. Describe “sigmoid function,  $\text{sigmoid}(x)$ ”, which is applied to logistic discrimination using a hyperplane. (Hint: using the posteriors of  $C1$  and  $C2$ , given  $x$ ).

**SOLUTION**