

## Function Approximation and RBF Neural Networks

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-1

## What is function approximation ?

Consider two sets  $F \subseteq R^m$  and  $L \subseteq R^l$ , and a mapping  $f$  from  $F$  to  $L$ . Suppose that for a set  $\Omega \subset F$ , a mapping  $g$  is given such that

$$g(x) \approx f(x), \quad \text{for } x \in \Omega$$

The so called function approximation problem is to find a mapping  $\hat{f}$  satisfying

$$\|\hat{f}(x) - f(x)\| \leq \varepsilon, \quad \text{for } x \in F$$

where  $\varepsilon > 0$  is the tolerance, and  $\|\bullet\|$  can be any norm.

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-2

## Some practical considerations

Since the original mapping is unknown, it cannot be used to evaluate the quality of a solution. In practice, we should find  $\hat{f}$  such that

$$\|\hat{f}(x) - g(x)\| < \varepsilon, \quad \text{for } x \in \Omega$$

In general, function approximation is an ill - posed problem. A solution can be good for the training set  $\Omega$ , but bad for  $x \in F - \Omega$ . To get solutions that generalize well, some regularization constraints are usually used.

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-3

## Polynomial based approximation

One model for function approximation is polynomial. Taylor expansion is one example. If the function is smooth enough, the Taylor expansion is given as follows :

$$f(x) = f(a) + \frac{f'(a)}{1!}(x-a) + \cdots + \frac{f^{(n-1)}(a)}{(n-1)!}(x-a)^{n-1} + R_n$$

where  $R_n$  is the reminder or the approximation error. That is, we can get an approximated value of  $f(x)$  using information at only one point.

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-4

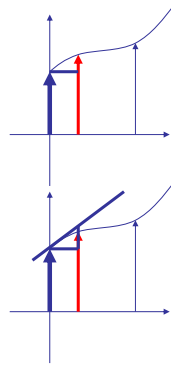
## 0-th order and 1-st order approximation

⇒ Nearest neighbor approximation

Replace  $x$  with the nearest example  $a$ , and approximate  $f(x)$  with  $f(a)$ .

⇒ Linear approximation

If the first order derivative exists, we can approximate  $f(x)$  with  $f(a) + f'(a)(x-a)$ .



Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-5

## General polynomial approximation

In general, a function  $f(x)$  can be approximated by

$$f(x) = \sum_{i=0}^{n-1} a_i x^i + \varepsilon$$

The coefficients can be found using the training examples in  $\Omega$ . Methods include

- 1) Solving a simultaneous linear equation.
- 2) Solving a quadratic optimization problem in case there are noises.

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-6

## Approximation based on Fourier transformation

Analysis :

$$F(\omega) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(x) e^{-j\omega x} dx = \langle f, \exp(-j\omega x) \rangle$$

Synthesis (Approximation)

$$f(x) = \frac{1}{\sqrt{2\pi}} \int_a^b F(\omega) e^{j\omega x} d\omega = \langle F(\omega), \exp(j\omega x) \rangle$$

where  $-\infty \leq a < b \leq \infty$

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-7

## Basis function approximation

In fact, both polynomial and transformation based approximations can be formulated in a more general form as follows :

$$f(x) = \sum_{k=0}^N a_k \phi_k(x) + \varepsilon$$

where  $\phi_k(x)$  is called a basis function.

Quizzes : What are the basis functions in the approximation methods discussed so far ?

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-8

## Function approximation with local basis functions

One example of local basis function is the Gaussian function given as follows :

$$\varphi_k(\mathbf{x}) = \exp\left(-\frac{1}{2\sigma_k^2}\|\mathbf{x} - \mathbf{x}_k\|^2\right)$$

where  $\mathbf{x}_i$  and  $\sigma_i$  are, respectively, the center and the variance of the  $k$  - th basis function. This kind of basis function is called radial basis function (RBF).

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-9

## RBF neural networks

- A RBF neural network is a three layer neural network
  - Input neurons: same as the MLP
  - Output neurons: linear combinations
  - Hidden neurons: basis functions
- For a one-output network, the output is given by

$$o = \sum_{k=1}^N w_k \varphi(\|\mathbf{x} - \mathbf{x}_i\|)$$

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-10

## Training of RBF neural networks

- Provide the training patterns with known function values.
- Use the training patterns as the centers of the radial basis functions.
- Find the weights of the output neuron by solving simultaneous linear equation.

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-11

## Problems to be solved

- Too many hidden neurons if the number of training patterns is large
  - Selection of important training data.
  - Re-location of the centers to improve the performance of the network.
  - Optimizing the variances of each RBF.
- For noisy data, the approximation is ill-posed
  - Regularization is needed.

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-12

## How to regularize ?

- The basic idea of regularization is to add a penalty in the error function.
- Minimizing the error along with the penalty can find a solution as smooth as possible.
- The cost function is given by

$$E(\hat{f}) = \sum_{i=1}^P (d_i - \hat{f}(x_i))^2 + \lambda \|P\hat{f}\|^2$$

where  $\lambda > 0$  is the regularization parameter, and  $P$  is a linear differential operator.

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-13

## A generalized RBF neural network

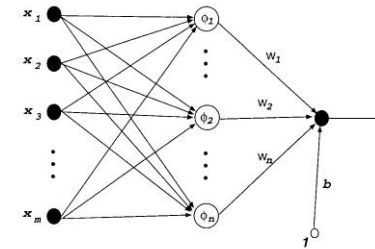


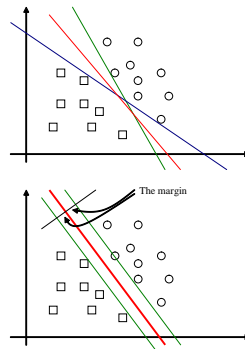
Figure 1: Radial basis function neural network

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-14

## Support vector machine

- Support vector machine (SVM) is an algorithm for pattern classification.
- Proposed by V. Vapnik at AT&A in 1995.
- The first important feature of SVM is to make decisions with the maximum margin.



Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-15

## Problem formulation (1)

Consider a two-class problem. Suppose that we have  $l$  linearly separable examples  $[(x_1, y_1), (x_2, y_2), \dots, (x_l, y_l)]$ , where  $y_i \in \{-1, 1\}$  is the teacher signal of  $x_i \in D \subseteq \mathbb{R}^n$ . We want to find a hyperplane to divide all examples into two classes, so that

- 1) The class labels are the same as the teacher signals.
- 2) The nearest examples of both sides are equal distance.

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-16

## Problem formulation (2)

Suppose that the hyperplane to find is

$$\mathbf{H} : \mathbf{w}'\mathbf{x} + b = 0$$

The distance between an example  $\mathbf{x}_i$  and  $\mathbf{H}$  is given by

$$\text{distance}(\mathbf{H}, \mathbf{x}_i) = \frac{|\mathbf{w}'\mathbf{x}_i + b|}{\|\mathbf{w}\|}$$

$\mathbf{H}$  can be scaled so that for the nearest example we have

$$\mathbf{w}'\mathbf{x} + b = \pm 1 \quad \text{or} \quad y(\mathbf{w}'\mathbf{x} + b) = 1$$

Thus the distance between  $\mathbf{H}$  and the nearest examples is  $1/\|\mathbf{w}\|$ , and the margin is  $2/\|\mathbf{w}\|$ .

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-17

## Problem formulation (3)

For an example outside the margin, we have

$$\mathbf{w}'\mathbf{x} + b > 1 \quad \text{or} \quad \mathbf{w}'\mathbf{x} + b < -1$$

$$\text{or} \quad y(\mathbf{w}'\mathbf{x} + b) > 1$$

Now the problem can be formulated as follows :

$$\begin{aligned} \min \quad & \|\mathbf{w}\|^2 / 2 \\ \text{s.t.} \quad & y_i(\mathbf{w}'\mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, l \end{aligned}$$

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-18

## Problem formulation (4)

The problem can be solved by optimizing the following Lagrangian function

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^l \alpha_i [y_i(\mathbf{w}'\mathbf{x}_i + b) - 1]$$

where  $\alpha_i \geq 0$  are the Lagrange multipliers. This is called the primal Lagrangian. The dual form can be found by differentiating  $L$  with respect to  $\mathbf{w}$  and  $b$ .

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-19

## Problem formulation (5)

The maximum margin classification problem can be formulated as a quadratic optimization problem

$$\begin{aligned} \max \quad & W(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l y_i y_j \alpha_i \alpha_j \mathbf{x}_i' \mathbf{x}_j \\ \text{s.t.} \quad & \sum_{i=1}^l y_i \alpha_i = 0; \quad \alpha_i \geq 0, \quad i = 1, 2, \dots, l \end{aligned}$$

If  $\alpha^*$  is the optimal solution, the optimal  $\mathbf{w}$  and  $b$  are given by

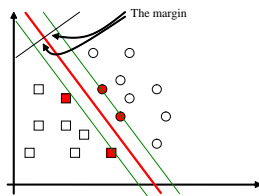
$$\mathbf{w} = \sum_{i=1}^l y_i \alpha_i^* \mathbf{x}_i, \quad b = -\frac{1}{2} \mathbf{w}'(\mathbf{x}_{\text{nearest}}^{+1} + \mathbf{x}_{\text{nearest}}^{-1})$$

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-20

## Some remarks

- From the optimal solution we can see that the weight vector  $w$  is the linear combination of all training examples.
- Not all  $\alpha_i$  take non-zero values.
- If  $\alpha_i$  is non-zero, the corresponding training example  $x_i$  is called a **support vector**.
- Support vectors are the examples nearest to the hyperplane  $H$ .
- Only support vectors are useful for making decisions.



Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-21

## How to make the decision?

The linear discriminant function is defined as

$$f(x) = \text{sgn}(w^t x + b)$$

$$= \text{sgn}\left(\sum_{i=1}^l y_i \alpha_i (x_i^t x) + b\right)$$

For any un-known example  $x$ , it is classified to class +1 if  $f(x) = 1$ ; or to class -1 if  $f(x) = -1$ .

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-22

## Support vector machine with soft margin

- So far we have assumed that all data are linearly separable.
- In practice, most problems are not linearly separable.
- The original problem can be relaxed by allowing some classification errors.
- That is, some data points can be inside the margin, or equivalently, the constraints can take the following form:

$$y_i (w^t x + b) \geq 1 - \xi_i, \quad i = 1, 2, \dots, l, \quad \text{with } \xi_i \geq 0$$

- This problem can be solved in the same way as before.

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-23

## Support vector machine with soft margin

- In the dual form optimization problem, the constraints should be modified as

$$0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, l, \quad \text{and} \quad \sum_{i=1}^l \alpha_i y_i = 0$$

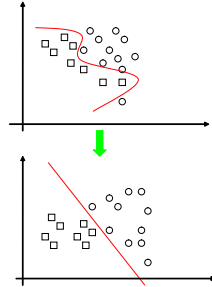
where  $C$  is the upper bound (to be chosen by the user) on the Lagrange multipliers  $\alpha_i$ .

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-24

## Non-linear support vector machine

- For non-linear problems, soft margin alone is not enough.
- Another important feature of SVM is to use a non-linear mapping.
- All data are first mapped from a low dimensional space to a high dimensional space.
- All data will become linearly separable in the mapped space, if the dimensionality of that space is high enough.



Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-25

## Basic considerations

- However, if we try to find the optimal hyperplane in the high dimensional space, the computational cost will become very large.
- This problem can be avoided if we introduce the concept of kernel function.
- A function  $k(x,y)$  is a kernel function if it can be represented as

$$k(x,y) = \phi(x)^T \phi(y) = \langle \phi(x), \phi(y) \rangle$$

- Where  $x$  and  $y$  are  $n$ -dimensional vectors, and  $\phi$  is the function for mapping  $x$  to the high dimensional space.

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-26

## Examples of kernel functions

### 1) Polynomial kernel

$$k(x, y) = (\langle x, y \rangle + c)^d$$

### 2) Gaussian kernel

$$k(x, y) = \frac{1}{c} e^{-\|x-y\|^2}$$

### 3) Sigmoid kernel

$$k(x, y) = \tanh[c \langle x, y \rangle + \theta]$$

where  $c$ ,  $d$ , and  $\theta$  are parameters.

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-27

## Optimal hyperplane in the mapped space

In the mapped space, the best hyperplane can be found by solving the following quadratic optimization problem

$$\begin{aligned} \max W(\alpha) &= \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l y_i y_j \alpha_i \alpha_j \phi^T(x_i) \phi(x_j) \\ &= \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l y_i y_j \alpha_i \alpha_j k(x_i, x_j) \end{aligned}$$

$$s.t. \quad \sum_{i=1}^l y_i \alpha_i = 0; \quad \alpha_i \geq 0, \quad i = 1, 2, \dots, l$$

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-28

## How to make the decision?

The linear discriminant function in the mapped space is defined as

$$\begin{aligned} f(x) &= \text{sgn}(w^t \phi(x) + b) \\ &= \text{sgn}\left(\sum_{i=1}^l y_i \alpha_i (\phi^t(x_i) \phi(x) + b)\right) \\ &= \text{sgn}\left(\sum_{i=1}^l y_i \alpha_i k(x_i, x) + b\right) \end{aligned}$$

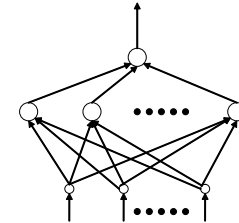
For any unknown example  $x$ , it is classified to class +1 if  $f(x) = 1$ ; or to class -1 if  $f(x) = -1$ .

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-29

## SVM is also a neural network!

- The original form of SVM is for 2-class problems, so there is only one output.
- For  $n$ -dimensional data, there are  $n$  inputs.
- The number of actually used hidden neurons is the number of support vectors.



Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-30

## Relation between SVM and RBF neural network

- SVM provides an efficient way for selecting the patterns to be used in an RBF-NN.
- However, SVM is not able to fine-tune the positions of the RBF centers and the widths of the basis functions.
- On the other hand, we can use different kernels in SVM for solving different problems.
- So far, SVM is known as one of the best method for pattern classification/recognition.

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-31

## Team Project VI

- Try to find some free program on the internet for designing SVM.
- Down-load at least two databases from the UCI machine learning repository, and design the SVMs using the program.
- Compare SVMs and the multilayer feedforward neural networks (MLPs) trained using BP, and make some conclusions about the accuracy, computing cost, system size, etc.

Produced by Qiangfu Zhao (Sine 1997), All rights reserved ©

Lecture 7-32