

Neural Networks

Other ANNs



André C P L F de Carvalho
Instituto de Ciências Matemáticas e de Computação
Universidade de São Paulo, Brazil

1

Topics

- Radial Basis Function neural networks
- Hybrid training
- Comparison with MLP neural networks
- Large margin classifiers
- Support Vector Machines
- Limitations
- Alternatives to overcome limitations

2

RBF networks

- So far, the activation function used by all ANNs receive
 - The internal product between the input and weight vectors
- Some multi-layer networks use activation functions that receive different values
 - E.g.: the distance between the input and weight vectors
 - Radial Basis Function (RBF) networks

3

André de Carvalho - LABIC/USP

3

RBF networks


$$\phi(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \boldsymbol{\mu}_j\|^2}{2\sigma_j^2}\right)$$

RBF networks typically uses a single hidden layer

4

André de Carvalho - LABIC/USP

4




RBF networks

- Two layers
 - First layer
 - Non-linear activation functions
 - Radial basis functions
 - Second layer
 - Non-linear or linear activation functions
 - Similar to Adaline

André de Carvalho - LABIC/USP 5

5



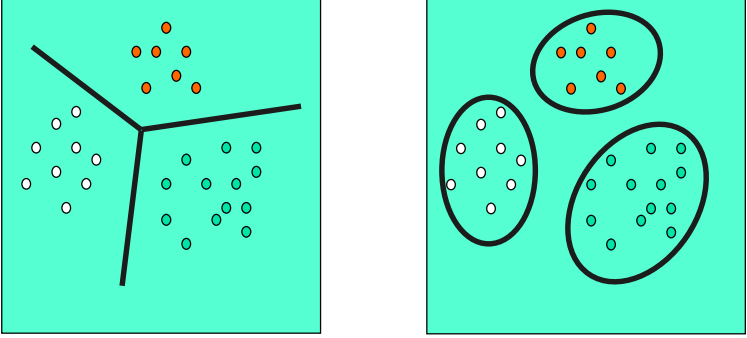
RBF networks

- RBF χ MLP networks
 - MLP use hyper-planes to partition the input space (hidden layer)
 - Defined by functions like $f(\sum w_i x_i)$
 - RBF use hyper-ellipsoids to partition the input space (hidden layer)
 - Defined by functions like $\phi(\|x_i - \mu_i\|)$
 - Distance between the input vector and the centre of a cluster

André de Carvalho - LABIC/USP 6

6

Decision boundaries



MLP

RBF

André de Carvalho - LABIC/USP

7

7


RBF networks

- Each node in the hidden layer computes a radial basis function
 - Centre
 - Defines the cluster prototype
 - Width
 - Defines the area covered by the cluster
- Can be much faster than MLP networks

André de Carvalho - LABIC/USP

8

8




RBF networks

- Total input
 - $u = \|x_i - \mu_i\|$ (hidden layer)
 - $u = \sum w_i \phi_i(\|x_i - \mu_i\|)$ (output layer)
- Distance measure
 - Usually, Euclidean distance

$$\|x_i - \mu_i\| = \sqrt{\sum_{i=1}^N (x_i - y_i)^2}$$

André de Carvalho - LABIC/USP 9

9



RBF networks


- Hidden layer activation functions
 - Non-linear
 - Value either increases or decreases when the total input moves away from the cluster centre
 - Typical functions:
 - Gaussian

$$\phi(v) = \exp\left(-\frac{v^2}{2\sigma^2}\right)$$

$v = \|x - \mu\|$
 x : input vector
 μ : radial function center
 σ : radial function width

André de Carvalho - LABIC/USP 10

10




RBF networks

- Usually employs hybrid learning
 - Unsupervised
 - Create clusters
 - K-means
 - Supervised
 - Least Mean Square, LMS
 - Singular Value Decomposition, SVD

André de Carvalho - LABIC/USP

11

11



RBF networks

- Main parameters to be defined:
 - Number of centres
 - Centres position
 - Centres width
 - Activation functions

André de Carvalho - LABIC/USP

12

12

Large Margin classifiers

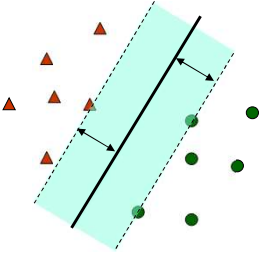
- Maximize the separation margin between different classes
 - Support Vector Machines (SVMs)
 - Boosting
- Higher generalization capacity
- Based on the statistical learning theory - Vapnik and Chervonenkis (1968)

18/12/2019 André de Carvalho 13

13

Support Vector Machines

- SVM looks for a hyperplane with maximum margin
 - Originally employed for linearly separable data



18/12/2019 André de Carvalho 14

14

Support Vector Machines

ANNs

SVMs

© André de Carvalho - ICMC/USP 15

15

Support Vector Machines

Maximum Margin

Optimal separation hiperplane

$$\text{sign}(h(x)) = \begin{cases} +1 & \text{if } w \cdot x + b \geq 1 \\ -1 & \text{if } w \cdot x + b \leq -1 \end{cases}$$

$$y_i \times (w \cdot x_i + b) \geq 1$$

$$\text{Margin} = \frac{2}{\|w\|^2}$$

© André de Carvalho - ICMC/USP 16

16

Slack variables

- Further increase the margins

© André de Carvalho - ICMC/USP

17

Linearly separable problems

- SVMs perform well for linearly separable problems
 - However, in its original format, cannot handle nonlinearly separable problems
- Some datasets require more complex than linear decision borders
 - For them, the Cover theorem can be used

© André de Carvalho - ICMC/USP

18

Teorema de Cover

Datasets that are nonlinearly separable in a given space can be transformed to another space in which, with high probability, they become linearly separable

- Conditions:
 - Transformation is nonlinear
 - Dimension of the new space is high enough

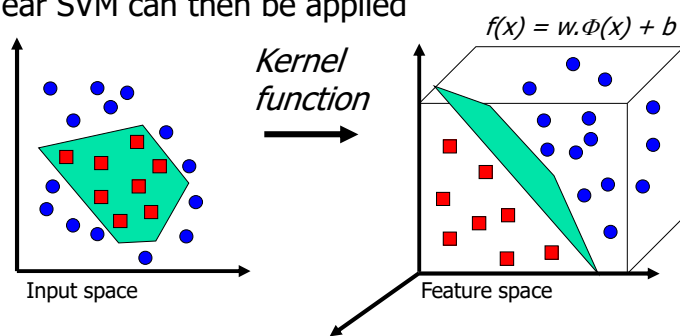
© André de Carvalho - ICMC/USP

19

19

Support Vector Machines


- Generalization for nonlinear problems
 - Mapping original data to higher dimensional space
 - Linear SVM can then be applied



© André de Carvalho - ICMC/USP

20

20




Example

- Suppose dataset with 2 predictive attributes
- Define 3 location points in the original set
- Use these points to transform the 2 original attributes into 3 new attributes
 - E.g. Distance between each example x_i and each of the 3 location points

© André de Carvalho - ICMC/USP 21

21



Kernel functions

- Several
 - Gaussian
 - Polynomial
 - Linear
 - Sigmoid
 - For specific applications
- Follow Mercer theorem conditions

© André de Carvalho - ICMC/USP 22

22

Multiclass classification

- SVMs can induce only binary classifiers
 - Other ML algorithms have the same limitation
- A large number of real problems has more than 2 classes
 - Multiclass strategies are necessary
 - Decompositional strategies are often used

© André de Carvalho - ICMC/USP 23

23

Code matrix

AAA

	f_1	f_2	f_3	f_4	f_5	f_6
1	+1	+1	+1	0	0	0
2	-1	0	0	+1	+1	0
3	0	-1	0	-1	0	+1
4	0	0	-1	0	-1	-1

↓ (1)x(2) ↓ (1)x(3) ↓ (2)x(3) ↓ (2)x(4) ↓ (3)x(4)

OAA

	f_1	f_2	f_3	f_4
1	+1	-1	-1	-1
2	-1	+1	-1	-1
3	-1	-1	+1	-1
4	-1	-1	-1	+1

↓ (1)x(2,3,4) ↓ (2)x(1,3,4) ↓ (3)x(1,2,4) ↓ (4)x(1,2,3)


ECOC

	f_1	f_2	f_3	f_4	f_5	f_6	f_7
1	+1	+1	+1	+1	+1	+1	+1
2	-1	-1	-1	-1	+1	+1	+1
3	-1	-1	+1	+1	-1	-1	+1
4	-1	+1	-1	+1	-1	+1	-1

↓ (1,2,3)x(4) ↓ (1,2,4)x(3) ↓ (1,2,3,4)x(2) ↓ (1,3,4)x(2,4) ↓ (1,4)x(2,3) ↓ (1)x(2,3,4)

24

24



Quiz 1

- What are limitations of the original SVM?
 - A) Only works for linearly separable problems
 - B) Minimize separation margin
 - C) Cannot deal with multiclass classification problems
 - D) Only works with one of 3 kernel functions

© André de Carvalho - ICMC/USP 25

25



Next: ANN applications

© André de Carvalho - ICMC/USP 26

26