

## 1 Basic Information

2 How to choose an activation function

3 Author: H. N. Mhaskar (CSLA), C. A. Micchelli (IBM. Watson) & Anne Schroeder

4 Proceeding NIPS'93 Proceedings of the 6th International Conference on Neural Information Processing Systems,  
5 pages 319-326

6 Citation: 39

## 7 Summary

8 This paper focuses on the complexity of estimating real-valued functions using neural networks. Complexity means  
9 approximating the number of neurons required to achieve a certain level of estimation accuracy. Previous work have  
10 proved that single hidden layer networks with any non-linear activation function can estimate any continuous function  
11 to any desired degree of accuracy. However, people still do not know how to construct (how many neurons to use)  
12 such networks given a known class of functions before this paper was published. Barron's theory only gave the results  
13 on sigmoidal activation function estimation for functions under certain conditions. Other work proved that functions  
14 with a bounded gradient with  $s$  variables require  $O(n^s)$  neurons. In this paper, the authors focused on approximating  
15 functions with a certain number of smooth derivatives, which expand the functions that can be estimated. They  
16 investigate the number of neurons requires given activation functions. In addition to examining sigmoid activation  
17 function, they also consider more classes of functions. They gave a method to construct such networks with radial  
18 basis activation functions, including Gaussian, squashing and sigmoidal functions. They also showed that the smoother  
19 the activation function, the better the rate of approximation. This paper provides an overview for choosing activation  
20 functions and expands the deal-able functions to a larger class.

## 21 Strong points

### 22 1. Originality.

23 The idea is not new according to the related work stated by the authors. This paper only examines more classes of  
24 activation functions based on previous work. However, the idea of using Fourier coefficients and infinitely many times  
25 continuously differentiable function to estimate is interesting and useful.

### 26 2. Technical detail.

27 The paper is theoretically well-founded on the following points. The authors clearly stated the notations and back-  
28 ground knowledge. The authors mentioned that this method is also valid for  $L^p$  - norms approximation, instead of  
29 only being valid for uniform approximation.

### 30 3. Writing

31 I like the way the authors decompose the problem. They added conditions so that the problem can be reduced to a  
32 simplified version but without loss of generality. The conversion made the proof procedure simpler and easier for  
33 readers to understand.

## 34 Weak points

### 35 1. Experiment/Example.

36 This is a purely theoretical work. Although experiments may not be easy to implement for theoretical work, I think  
37 more concrete examples could be made. For Table 1, most results are remarkable. However, when I am reading the  
38 table, I care more about the comparison of the most common functions we are using today, including  $ReLU(x^{k=1})$  if  
39  $x > 0, 0$ , otherwise). However, the paper only considered  $k > 1$ , *sigmoid*, *tanh* and so on. Unfortunately, famous  
40 activations such as *ReLU* and *tanh* are not examined in this table. Computing the result for a certain activation  
41 function requires much mathematics foundations and is not easy for readers to accomplish in a short time. Additionally,  
42 the authors also did not explicitly prove the relation between continuity and complexity to estimate, which is stated as  
43 a conclusion in the abstract.

### 44 2. Technical detail.

45 The paper mentioned some results with  $l$ , which is the number of hidden layers. However, the authors did not mention  
46 how to decide the number of neurons used in each layer. This makes the statements unclear.

### 47 3. Writing

48 Most core results have been published in the authors' previous work. They could have explained the most impor-  
49 tant conclusion instead of referring the reader to their other work. A brief explanation could help the readers better  
50 understand their results.