

Neural Network Learning: Theoretical Foundations

Chapter 24 ~ 26

Martin Anthony and Peter L. Bartlett

Presenter: Yongchan Kwon

Department of Statistics, Seoul National University, Seoul, Korea

ykwon0407@snu.ac.kr

Reviews; pp.312-313

Definition (RP)

The class of decision problems that can be solved by a polynomial-time randomized algorithm is denoted by RP.

Definition (H-FIT)

Instance: $z \in (\mathbf{R}^n \times \{0,1\})^m$ and an integer k between 1 and m .

Question: Is there $h \in H_n$ such that $\hat{e}r_z(h) \leq k/m$?

where H_n is a class of a binary function on n -dimensional inputs.

Reviews; pp.312-313

Theorem (23.7)

Let $H = \cup_n H_n$ be a graded binary function class. If there is an efficient learning algorithm for H , then there is a polynomial time randomized algorithm for H -FIT; in other words, H -FIT is in RP .

Theorem (23.8)

Suppose $RP \neq NP$ and that H is a graded class of binary functions. If H -FIT is NP -hard, then there is no efficient learning algorithm for H .

Contents

- 1 Ch. 24: The Boolean Perceptron
- 2 Ch. 25: Hardness Results for Feed-Forward Networks
- 3 Ch. 26: Constructive Learning Algorithms for Two-Layer Networks

Ch. 24: The Boolean Perceptron

Learning is Hard for the Simple Perceptron

Definition (BP-FIT)

Instance: $z \in (\{0,1\}^n \times \{0,1\})^m$ and an integer k between 1 and m .

Question: Is there $h \in BP_n$ such that $\hat{e}_z(h) \leq k/m$?

where BP_n is the set of boolean function from $\{0,1\}^n$ to $\{0,1\}$ computed by the boolean perceptron, and $BP = \cup_n BP_n$.

Definition (Simple perceptron)

A simple perceptron is a function $f : \mathbf{R}^n \rightarrow \{0,1\}$ of the form

$$f(x) = \begin{cases} 0, & \text{if } w^T x - \theta < 0 \\ 1, & \text{if } w^T x - \theta \geq 0 \end{cases}$$

for input vector $x \in \mathbf{R}^n$, $w \in \mathbf{R}^n$, and $\theta \in \mathbf{R}$.

Learning is Hard for the Simple Perceptron

Theorem (24.2)

BP-FIT is NP-hard.

Key idea: The problem is at least as hard as a well-known NP-hard problem in the field of graph theory.

Vertex cover problem [NP-hard]

A vertex cover of the graph is a set U of vertices such that for each edge (i, j) of the graph, at least one of the vertices i, j belongs to U .

Instance: A graph $G = (V, E)$ and an integer $k \leq |V|$

Question: Is there a vertex cover $U \subset V$ such that $|U| \leq k$?

Corollary (24.3)

If $RP \neq NP$, then there is no efficient learning algorithm for BP.

Learning is Easy for Fixed Fan-In Perceptrons

- The previous theorem shows that learning the simple perceptron is difficult. We consider simpler perceptrons in which the number of non-zero weights is constrained.

Definition (fan-in)

A simple perceptron with weights $w \in \mathbf{R}^n$ and threshold $\theta \in \mathbf{R}$ has fan-in k if the number of non-zero components of w is no more than k .

Pseudocode for the *Splitting* procedure

```

argument: Training set,  $S = \{x_1, \dots, x_m\} \subset \mathbb{R}^n$ 
returns: Set of weights and thresholds,  $W = \{(w, \theta)\}$ 

function Splitting( $S$ )
   $W := \emptyset$ 
   $P := \emptyset$ 
  for all  $t_1 < \dots < t_k$  from  $\{1, \dots, n\}$ 
    for all  $l$  from  $\{1, \dots, k+1\}$ 
      for all  $r_1 < \dots < r_l$  from  $\{1, \dots, m\}$ 
        for all  $\alpha_1, \dots, \alpha_l$  from  $\{\pm 1\}$ 
          if there is a solution  $(w, \theta)$  to the system
            of linear equations
               $x_{r_i} \cdot w + \theta = \alpha_i \quad i = 1, \dots, l$ 
            satisfying
               $\{i : w_i \neq 0\} = \{t_1, t_2, \dots, t_k\}$ 
          then
             $S' := \{x \in S : w \cdot x - \theta < 0\}$ 
             $S'' := \{x \in S : w \cdot x - \theta \geq 0\}$ 
            if  $\{S', S''\} \notin P$ 
              then
                 $W := W \cup \{(w, \theta)\}$ 
                 $P := P \cup \{S', S''\}$ 
            endif
          endif
        endfor
      endfor
    endfor
  endfor
  return  $W$ 
end

```

Learning is Easy for Fixed Fan-In Perceptrons

Theorem (24.4)

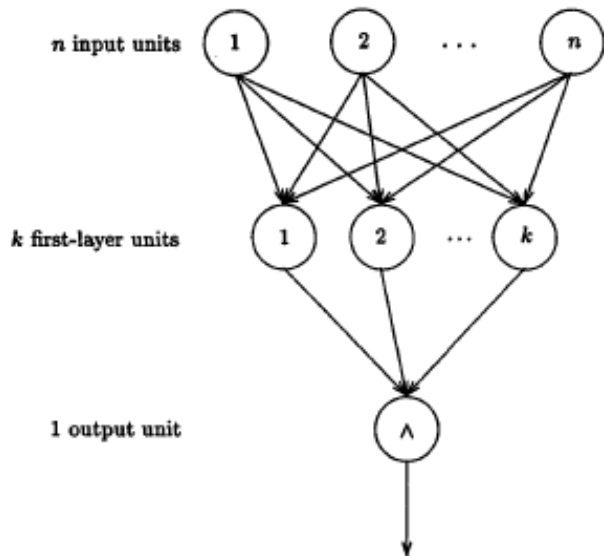
The procedure **Splitting** returns all dichotomies of its arguments $S \subset \mathbf{R}^n$ that can be computed by some simple perceptron with fan-in no more than k . For $|S| = m$, it takes time $O(n^{2k} 2^k m^{2k+3})$

Corollary (24.5)

For fixed k , define the graded class $H^k = \cup_n H_n^k$, where H_n^k is the class of simple perceptrons defined on \mathbf{R}^n with fan-in no more than k . The class H^k is efficiently learnable.

Ch. 25: Hardness Results for Feed-Forward Networks

Linear Threshold Networks with Binary Inputs



Linear Threshold Networks with Binary Inputs

- Let $N_{\wedge, n}^k$ be a neural network on n binary inputs and $k + 1$ linear threshold units. Further, we only consider $N_{\wedge, n}^k$ has two layers of computation units, the first consisting of k linear threshold units.
- The output unit is also a linear threshold unit, with a connection of fixed weight 1 from each of the other k threshold units.
- Consider the graded space $N_{\wedge}^k = \cup_n N_{\wedge, n}^k$.

N_{\wedge}^k – CONSISTENCY

Instance: $z \in (\{0, 1\}^n \times \{0, 1\})^m$

Question: Is there $h \in N_{\wedge, n}^k$ such that $\hat{e}r_z(h) = 0$?

Linear Threshold Networks with Binary Inputs

Corollary (25.2)

Let $k \leq 3$ be any fixed integer. Then, $N_{\wedge}^k - \text{CONSISTENCY}$ is NP-hard.

Key idea: Again, the problem is at least as hard as a well-known NP-hard problem in the field of graph theory.

$k - \text{colouring}$ [NP-hard]

A $k - \text{colouring}$ of G is a function $\chi : V \rightarrow \{1, 2, \dots, k\}$ with the property that whenever $(i, j) \in E$, then $\chi(i) \neq \chi(j)$.

Instance: A graph G

Question: Does G have a $k - \text{colouring}$?

Corollary (25.3)

Unless $RP = NP$, there is no efficient learning algorithm for the graded class $H = \cup_n H_n$, where H_n is the set of functions computable by $N_{\wedge, n}^k$.

Linear Threshold Networks with Real Inputs

- The result of the previous section is limited, since it shows that learning is difficult for a rather unusual network class. But...

Theorem (25.4)

Unless $RP = NP$, there is no efficient learning algorithm for the graded class $H = \cup_n H_n$, where H_n is the set of functions computable by N_n^k , a network with n real inputs.

- Similar results are obtained for sigmoid networks. (chapter 25.4).

Ch. 26: Constructive Learning Algorithms for Two-Layer Networks

Real Estimation with Convex Combinations of Basis Functions

- We consider learning algorithms for classes F of real valued functions that can be expressed as convex combinations of functions from some class G of basis functions.
- Some boosting and neural networks classes are example of F under some constraints

Real Estimation with Convex Combinations of Basis Functions

Theorem (26.1)

Let V be a vector space with an inner product, and let $\|f\| = \sqrt{(f, f)}$ be the induced norm on V . Suppose that $G \subset V$ and that, for some $B > 0$, $\|g\| \leq B$ for all $g \in G$. Fix $f \in V$, $k \in \mathbf{N}$, and $c \geq B^2$, and define $\hat{f}_0 = 0$. Then for $i = 1, \dots, k$, choose $g_i \in G$ such that

$$\|f - \hat{f}_i\|^2 \leq \inf_{g \in G} \|f - ((1 - \alpha_i)\hat{f}_{i-1} + \alpha_i g)\|^2 + e_i,$$

where $\alpha_i = 2/(i + 1)$, $e_i \leq 4(c - B^2)/(i + 1)^2$, and $\hat{f}_i = (1 - \alpha_i)\hat{f}_{i-1} + \alpha_i g_i$. Then,

$$\|f - \hat{f}_k\|^2 < \inf_{\hat{f} \in \text{co}(G)} \|f - \hat{f}\|^2 + \frac{4c}{k}.$$

Real Estimation with Convex Combinations of Basis Functions

Note that $\|f - ((1 - \alpha_i)\hat{f}_{i-1} + \alpha_i g)\|^2 = \alpha_i^2 \|\tilde{f} - g\|^2$, where $\tilde{f} = (f - (1 - \alpha_i)\hat{f}_{i-1})/\alpha_i$. This suggests using an approximate-SEM algorithm for a class G to approximately minimize sample error over the class $\text{co}(G)$.

Corollary (26.2)

Suppose that $G = \cup_n G_n$ is a graded class of real-valued functions that map to some bounded interval, and L is an efficient approximate-SEM algorithm for G , with running time $O(p(m, n, 1/e))$ for some polynomial p . Then, the algorithm Construct_L can be used as the basis of an efficient approximate-SEM algorithm for $\text{co}(G) = \cup_n \text{co}(G_n)$, and this algorithm has running time $O(p(m, n, 1/e)/e)$.

Pseudocode for the *Construct* procedure

arguments: Training set, $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \subset (X \times \mathbb{R})^m$
 Number of iterations, k
 Bound, B , on range of functions in G

returns: Convex combination of functions from G , $\hat{f}_k = \sum_{i=1}^k \gamma_k g_i$

function $\text{Construct}_L(S, k, B)$

```

 $\hat{f}_0 := 0$ 
for  $i := 1$  to  $k$ 
   $\alpha_i := 2/(i + 1)$ 
  for  $j := 1$  to  $m$ 
     $\tilde{y}_j = (1/\alpha_i) (y_j - (1 - \alpha_i)\hat{f}_{i-1}(x_j))$ 
  end for
   $\tilde{S} = \{(x_1, \tilde{y}_1), \dots, (x_m, \tilde{y}_m)\}$ 
   $g_i := L(\tilde{S}, B^2)$ 
   $\hat{f}_i := (1 - \alpha_i)\hat{f}_{i-1} + \alpha_i g_i$ 
endfor
return  $\hat{f}_k$ 
end

```

Fig. 26.1. Pseudocode for the Construct_L algorithm. (L is an approximate-SEM algorithm for $G \subset [-B, B]^X$.)

Real Estimation with Convex Combinations of Basis Functions

Let $G = BH \cup -BH$, with $H = \{\text{sgn}(w^T x + w_0) : w \in \mathbf{R}^n, w_0 \in \mathbf{R}\}$, and $BH = \{B \times h : h \in H\}$. Then, $F = \text{co}(G)$ is the class of two-layer networks with linear threshold units in the first layer and a linear output unit.

Theorem (26.6)

Let H_n^k be the set of k fan-in linear threshold functions, and let $F_n^k = \text{co}(BH_n^k \cup -BH_n^k)$. Then, the algorithm Construct, based on the algorithm Splitting, is an efficient learning algorithm for the graded class $F^k = \cup_n F_n^k$.