

구간분석을 통한 심층신경망의 안정성검증의 한가지 방법

리철진, 최창일

경애하는 김정은동지께서는 다음과 같이 말씀하시였다.

《과학기술부문에서 첨단돌파전을 힘있게 벌려야 하겠습니다.》

선행연구[1]에서는 비선형활성함수를 가지는 신경망의 안정성을 정의하고 한계 모형 검사도구를 리용하여 안정성을 검증하기 위한 일반적인 한가지 방법을 제기하였다. 선행연구[2]에서는 화상분류를 진행하는 신경망에 대하여 주어진 화상과 주어진 잡음에 관한 안정성을 정의하고 술어론리를 리용하여 검증하는 방법을 제기하였다. 선행연구[3]에서는 일정한 조건을 만족시키고 ReLU활성함수를 리용하는 신경망에 대하여 선행연구[2]에서 정의된 안정성을 검증하는 방법과 도구를 제기하였다. 선행연구[4]에서는 주어진 입력모임에 관한 심층신경망의 안정성문제를 출력모임과 불안정모임의 사립을 구하는 문제로 귀착시키는 방법을 제기하였다.

논문에서는 활성함수에 대한 제한조건이 선행연구[3]에서보다 약화된 심층신경망의 안정성을 선행연구[4]에서 제기한 방법을 갱신하여 검증하는 한가지 방법을 제기하였다.

심층신경망 N 을 다음과 같은 4원조 (L, W, b, Act) 로 표시한다.

① 층들의 모임 $L: \{L_k \mid L_k: n_k \text{ 개의 마디를 가진다. } 0 \leq k \leq l\}$

② 무게결수행렬들의 모임

$$W: \{W^{(k)} = [w_1^{(k)}, \dots, w_{n_k}^{(k)}]^T, w_i^{(k)} \in \mathbf{R}^{n_{k-1}}, 1 \leq i \leq n_k, 1 \leq k \leq l\}$$

③ 편향벡토르들의 모임

$$b: \{b^{(k)} = [b_1^{(k)}, \dots, b_{n_k}^{(k)}]^T, b_i^{(k)} \in \mathbf{R}, 1 \leq i \leq n_k, 1 \leq k \leq l\}$$

④ 활성함수모임 $Act: \{\phi_k, 1 \leq k \leq l\}$

심층신경망 $N = (L, W, b, Act)$ 의 층 $L_k (1 \leq k \leq l)$ 가 층 L_{k-1} 로부터 받아들이는 벡토르 $\mathbf{x}^{(k-1)}$ 에 대하여 이 층의 출력벡토르 $\mathbf{y}^{(k)}$ 는 다음과 같이 표시된다.

$$\mathbf{y}^{(k)} = \phi_k(W^{(k)}\mathbf{x}^{(k-1)} + \mathbf{b}^{(k)})$$

특히 심층신경망 N 의 출력층 L_l 의 출력벡토르 $\mathbf{y}^{(l)}$ 은 다음과 같이 표시된다.

$$\mathbf{y}^{(l)} = \Phi(\mathbf{x}^{(0)}), \Phi(\mathbf{x}^{(0)}) := \hat{\phi}_l \circ \dots \circ \hat{\phi}_1(\mathbf{x}^{(0)}), \hat{\phi}_k := \phi_k(W^{(k)}\mathbf{x}^{(k-1)} + \mathbf{b}^{(k)}) \quad (1)$$

정의 1 [4] 심층신경망 $N = (L, W, b, Act)$ 과 주어진 입력벡토르들의 유한모임 $X \subseteq \mathbf{R}^{n_0}$ 에 대하여 식 (1)에 의해 얻어지는 모임

$$Y = \{\mathbf{y}^{(l)} \in \mathbf{R}^{n_l} \mid \mathbf{y}^{(l)} = \Phi(\mathbf{x}^{(0)}), \mathbf{x}^{(0)} \in X\}$$

를 X 의 출력모임이라고 부른다. 또한 Φ 를 N 의 출력함수라고 부른다.

심층신경망의 검증에서는 특정한 입력벡토르들의 모임 $X \subseteq \mathbf{R}^{n_0}$ 에 대하여 그것의 출력모임의 원소들이 취해서는 안되는 모임 $S_X \subseteq \mathbf{R}^{n_l}$ 을 X 에 관한 불안정모임이라고 부른다.

실례로 softmax활성함수를 리용하여 0부터 9까지의 필기체수자화상을 분류하는 심층 신경망모형이 있다고 하자. 이 신경망에 대하여 수자 2를 나타내는 유한개의 화상들을 입력모임으로 할 때 출력벡토르에서 2에 대응하는 성분값이 0.5보다 작지 않으면 2로 정확히 분류된다.

이로부터 이 입력모임의 불안정모임 S_X 는 $S_X = \{y = [y_0, \dots, y_9] \in \mathbf{R}^{10} \mid y_2 < 0.5\}$ 로 볼 수 있다.

일반적으로 불안정모임은 심층신경망의 구조와 입력모임에 따라 결정된다.

정의 2 신경망 $N = (L, W, b, Act)$ 와 입력벡토르들의 유한모임 $X \subseteq \mathbf{R}^{n_0}$, X 의 불안정모임 S_X 와 출력모임 Y 에 대하여 $Y \cap S_X = \emptyset$ 이면 신경망 N 은 입력 X 에 대하여 안정하다고 말한다.

정의 3 \mathbf{R} 우에서의 닫힌구간전부의 모임을 \mathbf{IR} 로 표시하자. 함수 $\phi: \mathbf{R} \rightarrow \mathbf{R}$ 에 대하여 $[\phi]([x, x]) = \phi(x)$, $x \in \mathbf{R}$ 를 만족시키는 함수 $[\phi]: \mathbf{IR} \rightarrow \mathbf{IR}$ 를 ϕ 의 구간확장이라고 부른다. 마찬가지로 함수 $\Phi: \mathbf{R}^n \rightarrow \mathbf{R}^m$ ($n, m \in \mathbf{N}$)에 대하여

$$[\Phi]([x_1, x_1] \times \dots \times [x_n, x_n]) = \Phi(\mathbf{x}), \quad \mathbf{x} = [x_1, \dots, x_n] \in \mathbf{R}^n$$

이 성립하는 함수 $[\Phi]: \mathbf{IR}^n \rightarrow \mathbf{IR}^m$ 을 Φ 의 구간확장이라고 부른다.

n 차원벡토르들의 유한모임 $X \subseteq \mathbf{R}^n$ 에 대하여 X_i 를 X 의 원소들의 i 째 성분들의 모임이라고 할 때 닫힌구간 $[\min(X_i), \max(X_i)]$, $i=1, \dots, n$ 들의 직적을 $[X]$ 로 표시한다.

$[X] = [\underline{x}_1, \bar{x}_1] \times \dots \times [\underline{x}_n, \bar{x}_n] \in \mathbf{IR}^n$ 에 대하여 $w([X]) = \max_{i=1, \dots, n} (\bar{x}_i - \underline{x}_i)$ 를 $[X]$ 의 너비라고 부른다.

정의 4 함수 $\Phi: \mathbf{R}^n \rightarrow \mathbf{R}^m$ ($n, m \in \mathbf{N}$)의 구간확장 $[\Phi]: \mathbf{IR}^n \rightarrow \mathbf{IR}^m$ 이 임의의 $[X_1], [X_2] \in \mathbf{IR}^n$ 에 대하여

$$[X_1] \subseteq [X_2] \Rightarrow [\Phi]([X_1]) \subseteq [\Phi]([X_2])$$

를 만족시키면 $[\Phi]$ 는 포함단조라고 부른다.

정의 4로부터 다음의 보조정리를 얻는다.

보조정리 1 함수 $\Phi: \mathbf{R}^n \rightarrow \mathbf{R}^m$ ($n, m \in \mathbf{N}$)의 구간확장 $[\Phi]: \mathbf{IR}^n \rightarrow \mathbf{IR}^m$ 이 포함단조일 때 임의의 n 차원벡토르들의 유한모임 X 에 대하여 다음의 성질이 만족된다.

$$\forall \mathbf{x} \in [X] \Rightarrow \Phi(\mathbf{x}) \in [\Phi]([X]) \quad (2)$$

증명 $[\Phi]$ 는 함수 Φ 의 구간확장이므로 $\forall \mathbf{x} = [x_1, \dots, x_n] \in [X]$ 에 대하여

$$\Phi(\mathbf{x}) = [\Phi]([x_1, x_1] \times \dots \times [x_n, x_n])$$

이 성립한다. 또한 $[x_1, x_1] \times \dots \times [x_n, x_n] \subset [X]$ 이며 $[\Phi]$ 는 포함단조인 구간확장이므로

$$[\Phi]([x_1, x_1] \times \dots \times [x_n, x_n]) \subseteq [\Phi]([X])$$

가 성립하며 따라서 식 (2)가 만족된다.(증명끝)

보조정리 1로부터 다음의 보조정리를 얻는다.

보조정리 2 함수 $\Phi: \mathbf{R}^n \rightarrow \mathbf{R}^m$ ($n, m \in \mathbf{N}$)의 구간확장 $[\Phi]: \mathbf{IR}^n \rightarrow \mathbf{IR}^m$ 이 포함단조일 때 임의의 n 차원벡토르들의 유한모임 X 에 대하여 다음의 성질이 만족된다.

$$\Phi([X]) \subset [\Phi]([X]) \quad (3)$$

증명 보조정리 1로부터 $\forall \mathbf{x} \in [X]$ 에 대하여 $\Phi(\mathbf{x}) \in [\Phi]([X])$ 가 성립하며 따라서 식 (3)이 만족된다.(증명끝)

보조정리 2로부터 신경망 $N=(L, W, b, Act)$ 의 입력벡토르들의 유한모임 X 와 그것의 불안정모임 S_X 에 대하여 출력함수 Φ 의 구간확장 $[\Phi]$ 가 포함단조이고 $[\Phi]([X]) \cap S_X = \emptyset$ 이면 N 은 X 에서 안정하다는것을 알수 있다.

출력함수 Φ 의 포함단조인 구간확장 $[\Phi]$ 를 구하기 위하여 신경망의 활성화함수에 대한 한가지 가정을 제기한다.

가정 신경망 $N=(L, W, b, Act)$ 의 활성화함수 $\phi_k(k=1, \dots, l)$ 는 단조증가함수이다. 즉

$$\forall x_1, x_2 \in \mathbf{R}, x_1 \leq x_2 \Rightarrow \phi_k(x_1) \leq \phi_k(x_2)$$

심층신경망에서 리용되는 ReLU함수, tanh함수, 로지스틱함수, sigmoid함수 등 대부분의 활성화함수들은 모두 가정을 만족시킨다.

가정을 만족시키는 활성화함수를 리용하는 신경망에 대하여 출력함수의 포함단조인 구간확장을 구성하는 한가지 방법을 제기한다.

정리 가정을 만족시키는 신경망 $N=(L, W, b, Act)$ 와 N 의 입력벡토르들의 유한모임 X 에 대하여 다음과 같이 구성되는 함수 $[\Phi]$ 는 출력함수 Φ 의 구간확장이며 포함단조이다.

$$[\Phi]([X]) = [\hat{\phi}_l] \circ \dots \circ [\hat{\phi}_1]([X]) \quad (4)$$

여기서

$$[\hat{\phi}_k]([X^{\{k\}}]) := [\phi_k](W^{\{k\}}[X^{\{k\}}] + b^{\{k\}}) = [Y^{\{k\}}] = [\underline{y}_1^{\{k\}}, \bar{y}_1^{\{k\}}] \times \dots \times [\underline{y}_{n_k}^{\{k\}}, \bar{y}_{n_k}^{\{k\}}] \quad (5)$$

$$\begin{aligned} \underline{y}_i^{\{k\}} &:= \phi_k \left(\sum_{j=1}^{n_{k-1}} \underline{p}_{ij}^{\{k\}} + b_i^{\{k\}} \right), \quad \underline{p}_{ij}^{\{k\}} := \begin{cases} w_{ij}^{\{k\}} \underline{x}_j^{\{k\}}, & w_{ij}^{\{k\}} \geq 0 \\ w_{ij}^{\{k\}} \bar{x}_j^{\{k\}}, & w_{ij}^{\{k\}} < 0 \end{cases} \\ \bar{y}_i^{\{k\}} &:= \phi_k \left(\sum_{j=1}^{n_{k-1}} \bar{p}_{ij}^{\{k\}} + b_i^{\{k\}} \right), \quad \bar{p}_{ij}^{\{k\}} := \begin{cases} w_{ij}^{\{k\}} \bar{x}_j^{\{k\}}, & w_{ij}^{\{k\}} \geq 0 \\ w_{ij}^{\{k\}} \underline{x}_j^{\{k\}}, & w_{ij}^{\{k\}} < 0 \end{cases} \end{aligned} \quad (i=1, \dots, n_k) \quad (6)$$

이며 $[X^{\{k\}}] = [Y^{\{k-1\}}]$, $[X^{\{1\}}] = [X]$ 이다.

증명 먼저 $[\Phi]$ 가 출력함수 Φ 의 구간확장이라는것을 밝히자.

임의의 입력벡토르 $\mathbf{x} = [x_1, \dots, x_{n_0}] \in \mathbf{R}^{n_0}$ 에 대하여 $[x_1, x_1] \times \dots \times [x_{n_0}, x_{n_0}]$ 을 $[\mathbf{x}]$ 로 표시하면

$$[\hat{\phi}_1]([\mathbf{x}]) = [\underline{y}_1^{\{1\}}, \bar{y}_1^{\{1\}}] \times \dots \times [\underline{y}_{n_1}^{\{1\}}, \bar{y}_{n_1}^{\{1\}}]$$

$$\underline{y}_i^{\{1\}} = \bar{y}_i^{\{1\}} = \phi \left(\sum_{j=1}^{n_0} w_{ij}^{\{1\}} x_j + b_i^{\{1\}} \right)$$

로서 $[\hat{\phi}_1]([x_1, x_1] \times \dots \times [x_{n_0}, x_{n_0}])$ 은 $\hat{\phi}_1(\mathbf{x})$ 와 같다. 마찬가지로

$$[\Phi]([\mathbf{x}]) = [\hat{\phi}_l] \circ \dots \circ [\hat{\phi}_1]([\mathbf{x}]) = \hat{\phi}_l \circ \dots \circ \hat{\phi}_1(\mathbf{x}) = \Phi(\mathbf{x})$$

가 성립하며 따라서 $[\Phi]$ 는 출력함수 Φ 의 구간확장이다.

다음으로 $[\Phi]$ 가 포함단조이라는것을 밝히자.

$[X_1] \subseteq [X_2]$ 인 $[X_m] = [\underline{x}_{m,1}, \bar{x}_{m,1}] \times \cdots \times [\underline{x}_{m,n_0}, \bar{x}_{m,n_0}]$, $m=1, 2$ 가 있다고 하자. 그러면

$$\underline{x}_{2,j} \leq \underline{x}_{1,j} \leq \bar{x}_{1,j} \leq \bar{x}_{2,j}, \quad j=1, \dots, n_0 \quad (7)$$

이 성립한다. $[X_m]$, $m=1, 2$ 에 $[\hat{\phi}]$ 을 적용한 경우를 보기로 하자.

$$[\hat{\phi}](X_m) = [\underline{y}_{m,1}^{\{\}} , \bar{y}_{m,1}^{\{\}}] \times \cdots \times [\underline{y}_{m,n_1}^{\{\}} , \bar{y}_{m,n_1}^{\{\}}], \quad m=1, 2$$

$$\begin{aligned} \underline{y}_{m,i}^{\{\}} &= \phi_1 \left(\sum_{j=1}^{n_0} \underline{p}_{m,ij}^{\{\}} + b_i^{\{\}} \right), \quad \underline{p}_{m,ij}^{\{\}} = \begin{cases} w_{ij}^{\{\}} \underline{x}_{m,j}, & w_{ij}^{\{\}} \geq 0 \\ w_{ij}^{\{\}} \bar{x}_{m,j}, & w_{ij}^{\{\}} < 0 \end{cases} \\ \bar{y}_{m,i}^{\{\}} &= \phi_1 \left(\sum_{j=1}^{n_0} \bar{p}_{m,ij}^{\{\}} + b_i^{\{\}} \right), \quad \bar{p}_{m,ij}^{\{\}} = \begin{cases} w_{ij}^{\{\}} \bar{x}_{m,j}, & w_{ij}^{\{\}} \geq 0 \\ w_{ij}^{\{\}} \underline{x}_{m,j}, & w_{ij}^{\{\}} < 0 \end{cases} \end{aligned} \quad (i=1, \dots, n_1)$$

식 (7)로부터 $\underline{p}_{1,ij}^{\{\}} \geq \underline{p}_{2,ij}^{\{\}}, \bar{p}_{1,ij}^{\{\}} \geq \bar{p}_{2,ij}^{\{\}} (i=1, \dots, n_1, j=1, \dots, n_0)$ 이 성립한다는것을 쉽게 알수 있다. 또한 가정에 의하여 활성함수 $\phi_k (k=1, \dots, l)$ 는 단조증가함수이므로

$$\begin{aligned} \underline{y}_{1,i}^{\{\}} &= \phi_1 \left(\sum_{j=1}^{n_0} \underline{p}_{1,ij}^{\{\}} + b_i^{\{\}} \right) \geq \phi_1 \left(\sum_{j=1}^{n_0} \underline{p}_{2,ij}^{\{\}} + b_i^{\{\}} \right) = \underline{y}_{2,i}^{\{\}} \\ \bar{y}_{1,i}^{\{\}} &= \phi_1 \left(\sum_{j=1}^{n_0} \bar{p}_{1,ij}^{\{\}} + b_i^{\{\}} \right) \leq \phi_1 \left(\sum_{j=1}^{n_0} \bar{p}_{2,ij}^{\{\}} + b_i^{\{\}} \right) = \bar{y}_{2,i}^{\{\}} \end{aligned}$$

이 성립하며 따라서 $[\underline{y}_{1,i}^{\{\}}, \bar{y}_{1,i}^{\{\}}] \subseteq [\underline{y}_{2,i}^{\{\}}, \bar{y}_{2,i}^{\{\}}], i=1, \dots, n_1$ 로서 $[\hat{\phi}](X_1) \subseteq [\hat{\phi}](X_2)$ 도 성립한다. 즉 $[\hat{\phi}]$ 은 포함단조이다.

마찬가지로 $[\hat{\phi}] \circ \cdots \circ [\hat{\phi}](X_1) \subseteq [\hat{\phi}] \circ \cdots \circ [\hat{\phi}](X_2)$ 도 성립하며 $[\Phi]$ 가 포함단조라는것이 증명된다.(증명끝)

정리 1과 보조정리 2로부터 다음과 같은 따름을 얻게 된다.

따름 가정을 만족시키는 신경망 $N=(L, \mathbf{W}, \mathbf{b}, Act)$ 와 N 의 입력벡토르들의 유한모임 X , 그것의 불안정모임 S_X , 식 (4)–(6)과 같이 구성된 출력함수 Φ 의 구간확장 $[\Phi]$ 에 대하여 $[\Phi](X) \cap S_X = \emptyset$ 이면 N 은 X 에 대하여 안정하다.

우의 따름에 기초하여 가정을 만족시키는 신경망 $N=(L, \mathbf{W}, \mathbf{b}, Act)$ 와 N 의 입력벡토르들의 유한모임 X 와 그것의 불안정모임 S_X 가 주어졌을 때 안정성을 검증하는 알고리즘을 제기한다.

알고리즘

입력: 가정을 만족시키는 신경망 $N=(L, \mathbf{W}, \mathbf{b}, Act)$ 와 입력모임 $X \subseteq \mathbf{R}^{n_0}$, X 의 불안정모임 $S_X \subseteq \mathbf{R}^{n_l}$, 턱값 $\varepsilon > 0$

출력: 모임 *Unsafe*

단계 1: $M \leftarrow \{X\}$, *Unsafe* $\leftarrow \{\}$

단계 2: M 이 빈모임이면 단계 7로 넘어간다. 아니면 단계 3으로 넘어간다.

단계 3: M 에서 한 원소 $[X]$ 를 선택한다. 동시에 M 에서 그 원소를 제거한다.

단계 4: $[\Phi](X) \cap S_X$ 가 빈모임이면 단계 2로 넘어간다. 아니면 단계 5로 넘어간다.

단계 5: $w([X]) > \varepsilon$ 이면 단계 6으로 넘어간다. 아니면 *Unsafe* 에 $[X]$ 를 추가하고 단계 2로 넘어간다.

단계 6: $[X] = [x_1, \bar{x}_1] \times \cdots \times [x_{n_0}, \bar{x}_{n_0}]$ 에서 $\bar{x}_i - x_i = w([X])$ 인 한 닫힌구간을 $[x_i, x_i + w([X])/2]$, $[x_i + w([X])/2, \bar{x}_i]$ 로 분할하여 $[X_1], [X_2]$ 를 구성하고 모임 M 에 추가한 다음 단계 2로 넘어간다.

단계 7: *Unsafe* 를 출력한다.

모임 *Unsafe* 가 빈모임이면 신경망 N 은 X 에 대하여 안정하다.

참 고 문 헌

- [1] L. Pulina et al.; AI Communications, 25, 2, 117, 2012.
- [2] X. Huang et al.; CAV, LNCS, 10426, 3, 2017.
- [3] G. Katz et al.; CAV, LNCS, 10426, 97, 2017.
- [4] W. Xiang et al.; IEEE Transactions on Neural Network and Learning Systems, 22, 5777, 2018.

주체109(2020)년 9월 5일 원고접수

Safety Verification of Deep Neural Network Using Interval Analysis

Ri Chol Jin, Choe Chang Il

A method based on interval analysis for safety verification of deep neural network is presented in this paper.

Keywords: deep neural network, safety verification, interval analysis