

클라우드 네이티브 방식의 분산 엣지 클라우드를 위한 연합학습

구자빈, 김종원

광주과학기술원

pseineshi@gist.ac.kr, jongwon@gist.ac.kr

Federated Learning for
Cloud-native Distributed Edge Clouds

Jabin Koo, JongWon Kim

Gwangju Institute of Science and Technology

요 약

본 논문은 클라우드 네이티브 방식의 분산 엣지 클라우드에서의 전통적인 중앙화된 기계학습 방법의 현실적인 한계들을 제시하고 해당 한계들을 극복할 수 있는 방법인 연합학습의 모의구현을 위해 클라우드 네이티브 방식의 분산 엣지 클라우드를 가정한 연합학습 클러스터 환경을 구축하고 해당 클러스터 환경 위에서 연합학습의 실행을 보인다.

1. 서 론

5G 네트워크와 사물인터넷 시대의 출현에 의해 클라우드에 연결된 지능형 단말들의 규모가 급속하게 커졌고, 이로 인해 연산 또는 저장을 데이터의 근원에 더 가까운 곳으로 나눠서 처리하는 분산 엣지 클라우드 기술이 주목받게 되었다. 엣지 클라우드 기술은 공공 안전감시, 교통제어, 생활 IoT 등의 다양한 분야에서 기계학습과 긴밀하게 연계되고 있으며 기계학습에 유용한 대량의 데이터를 실시간으로 뽑아낸다.

그러나 이렇게 엣지에서 생성된 데이터를 전통적인 중앙화된 기계학습을 이용해 학습하는 것은 큰 한계들을 동반한다. 엣지에 연결된 단말들의 수가 증가하고 분산된 엣지의 규모가 커짐에 따라 실시간으로 생성되는 모든 데이터를 중앙 클라우드로 전송하고 저장, 처리하는 모든 과정은 큰 부하가 된다. 또한 안전감시, 생활 IoT, 모바일 데이터 등 엣지에서 생성되는 많은 데이터가 프라이버시에 민감한 정보를 포함할 가능성이 커 데이터를 중앙에 전송하고 저장하여 기계학습 모델훈련을 진행하는 것은 개인의 프라이버시에 큰 위협을 초래한다.

이에 따라 해당 한계들을 극복할 수 있는 연합학습[1]이 분산 엣지 클라우드에서의 기계학습 방법론으로서 주목받고 있다. 분산된 데이터를 개별 클러스터 내부에서 학습 후 데이터가 아닌 그레디언트나 모델만을 중앙으로 전송하여 기계학습 모델을 학습하는 연합학습 방법은 프라이버시에 민감한 데이터를 외부로 공유하지 않으며 효율적으로 연산 및 저장 부하를 분산하고 전송량마저 큰 폭으로 줄일 수 있다.

본 논문에서는 특별히 클라우드 네이티브 방식의 분산 엣지 클라우드에서의 연합학습을 모의 구현하였다. 클라우드 네이티브 방식의 분산 엣지 클라우드는 분산된 엣지 클라우드로 확장성, 통일성 그리고 회복성 있는 운영을 제공하며 특히 클라우드 네이티브 방식의

확장성과 회복성은 연합학습에서의 데이터 편향 그리고 개별 지점의 실패로 인한 낙오자 문제를 최소화하는데 유용할 수 있다.

II. 클라우드 네이티브 방식의 분산 엣지 클라우드를 가정한 연합학습 클러스터 환경 구축

클라우드 네이티브 방식의 분산 엣지 환경에서의 연합학습 구현을 위해 그림 1과 같이 분산 엣지 클라우드를 가정한 쿠버네티스 기반의 멀티 클러스터 환경을 구축하였다.

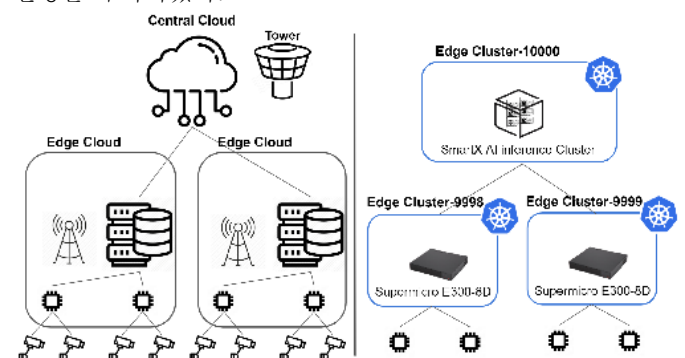


그림 1. 클라우드 네이티브 방식의 분산 엣지 클라우드 환경을 가정한 멀티 클러스터 구조

해당 환경은 SmartX AI inference cluster[2]와 두 개의 단일노드로 구성된 쿠버네티스 클러스터로 구성되어있으며, 각 클러스터는 10000, 9999, 9998의 번호를 가진다. 모든 클러스터는 각자의 고유 데이터를 가지고 연합학습 모델 학습에 협력하는 분산 엣지 클라우드의 역할을 수행하며 클러스터 10000은 모델 학습을 지시하고, 그레디언트를 병합하여 기계학습 모델을 갱신하고 재배포하는 중앙 클라우드의 역할 또한 수행한다.

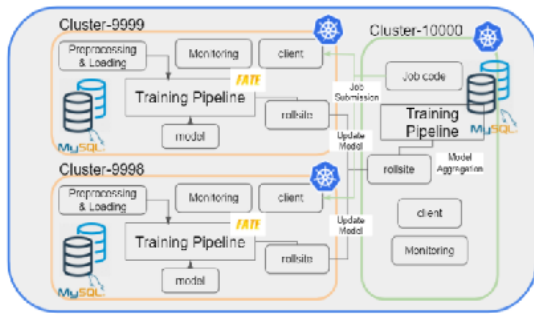


그림 2. 클라우드 네이티브 방식의 연합학습 클러스터 모듈 구조

각 쿠버네티스 클러스터 위에는 WeBank의 FATE(Federated AI Technology Enabler)[3] 프레임워크와 Federated AI Ecosystem 기반의 클러스터 모듈을 배포하여 그림 2와 같이 클라우드 네이티브 방식으로 연합학습 구조를 구축하였다.

각 클러스터는 독립적으로 데이터를 관리하며 인그레스를 통해 소통용 rollsite와 연합학습 monitoring site 만이 외부에 노출되며 각 클러스터 간에는 암호화된 작업 명령, 그래디언트와 모델만을 공유한다.

III. 클라우드 네이티브 방식의 연합학습 클러스터를 활용한 연합학습

본 논문에서는 클라우드 네이티브 방식으로 구성된 연합학습 클러스터의 검증에 위해 간단한 연합 신경망 모델 학습 실험 과정을 구현하였다.

해당 실험에서는 수평적으로 각 클러스터에 균등하게 나누어진 Breast Cancer Wisconsin Diagnostic Data Set[4]을 이용하여 악성 유방종양을 분류하는 하나의 신경망 모델을 연합하여 학습하였다.

학습은 그래디언트를 하나의 클러스터에 모아 그래디언트 평균화를 통해 신경망 모델을 갱신하고 각 클러스터에서의 학습을 위해 재배포하는 FedAVG[5] 방법을 통해 진행되었으며, 각 클러스터의 학습 과정은 그림 3과 같이 웹 인터페이스를 통해 모니터링된다.

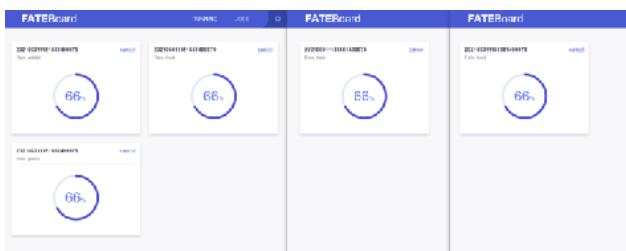


그림 3. 실증구조에서의 연합학습 과정 중 모니터링 창.

해당 실험의 학습에 이용된 신경망은 2계층 밀집 신경망으로 각 클러스터에서 full batch size로 2회씩 학습한 그래디언트를 총 10회의 모델병합을 통해 학습하였다. 학습된 신경망 모델은 0.9736842의 정확도를 달성하였으며 학습시간은 27.0269초가 소요되었다.

위 실험을 통해 본 논문은 구성된 클라우드 네이티브 방식의 연합학습 클러스터 환경에서의 엣지 클러스터간 연합학습의 실행을 보였다. 그러나 본 실험은 균등하고 Non-IID 하지 않은 분산 데이터를 이용해 얻은

신경망을 학습한다는 분산 엣지 클라우드에서 다소 현실적이지 않은 상황을 가정하여 한정된 검증을 보인다.

IV. 결론 및 추가 계획

본 논문에서는 클라우드 네이티브 방식의 분산 엣지 클라우드에서의 기존 기계학습 방법의 한계를 제시하였으며, 해당 한계들을 해결할 수 있는 연합학습의 모의구현을 위해 클라우드 네이티브 방식의 분산 엣지 클라우드를 가정한 연합학습 클러스터를 구축하고 해당 클러스터 위에서의 연합학습의 실행을 보였다.

그러나 본 논문에서는 다소 현실적이지 못한 실험상황을 가정하여 실제 분산 엣지 클라우드를 가정한 클라우드 네이티브 방식의 연합학습의 구동을 완전히 검증하지 못했다. 따라서 기존 연합학습 클러스터를 이용해 Real-World Image Dataset for Federated learning[6]의 더욱 현실적인 데이터와 현실적인 신경망 모델을 학습하는 분산 엣지 클라우드 모의실험을 진행하며 클라우드 네이티브 방식의 연합학습의 확장성과 회복성에 대한 추가 검증과 해당 특성들로 인한 낙오자 문제의 완화 가능성에 대한 검증 실험을 추가 진행할 예정이다.

ACKNOWLEDGMENT

본 논문은 2019년도 정부(과학기술정보통신부)의 재원으로 정보통신기획평가원의 지원을 받아 수행된 연구임 (No. 2019-0-01842, 인공지능대학원지원(광주과학기술원))

참 고 문 헌

- [1] Reza, S and Vitaly, S. "Privacy-Preserving Deep Learning." In Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communication Security - CCS '15, pages 1310-1320, Denver, Colorado, USA, 2015, ACM Press
- [2] 권진철, 김남곤, 김중원 (2019). ICT 응용 서비스 지능화를 지원하는 클라우드-네이티브 기반 SmartX AI 컴퓨팅 클러스터 설계 및 검증. 정보과학회 컴퓨팅의 실제 논문지. 25(12). 571-584
- [3] An Industrial Grade Federated Learning Framework. Available online: <https://fate.fedai.org/>.
- [4] Breast Cancer Wisconsin Diagnostic Data Set. Available online: <https://www.kaggle.com/uciml/breast-cancer-wisconsin-data>
- [5] Jakub, K. Brendan, M and Daniel, R. "Federated Optimization: Distributed Optimization Beyond the Datacenter." arXiv:1511.03575[cs, math], Nov. 2015
- [6] Jiahuan, L. "Real-World Image Datasets for Federated Learning," arXiv:1910.11089, Oct. 2019.