

學 士 學 位 論 文

# 프라이버시 보호 딥러닝 서비스 개발

충남대학교

공과대학 컴퓨터공학과

이 상 화

조 승 현

김 수 민

김 주 희

지도교수 임 성 수

2021 年 2 月

# 프라이버시 보호 딥러닝 서비스 개발

지도교수 임 성 수

이 논문을 공학사학위  
청구논문으로 제출함

2020 年 11 月

충 남 대 학 교

공과대학 컴퓨터공학과

201402392 이 상 화

201402433 조 승 현

201704145 김 수 민

201704145 김 주 희

## 목 차

I. 서론 .....	1
II. 연구 내용 .....	2
1. 차등 프라이버시(Differential Privacy) .....	2
2. 연합학습(Federated Learning) .....	4
III. 설계 .....	7
1. 실험도구 .....	7
2. Data Set .....	8
3. 실험 내용 .....	11
IV. 실험 결과 .....	14
1. 차등 프라이버시 실험 결과 .....	14
2. 연합학습 실험 결과 .....	17
V. 정리 .....	20
참고문헌 .....	21

## 표목차

〈표 1〉 비식별화 기술의 재식별 가능성 .....	2
〈표 2〉 개발 도구 .....	7
〈표 3〉 차등 프라이버시 적용 전과 후의 상품 추천 결과 .....	16

## 그림목차

<그림 1> 차등 프라이버시를 통한 데이터 변조 과정 .....	3
<그림 2> 연합학습 과정 .....	4
<그림 3> 각 디바이스에서 학습 수행 .....	5
<그림 4> 학습된 파라미터를 중앙으로 전송 .....	5
<그림 5> 업데이트된 모델 디바이스로 재전송 .....	5
<그림 6> 전송 가중치에 차등 프라이버시를 적용한 연합학습 .....	6
<그림 7> 개인신용정보 데이터 관계 .....	8
<그림 8> 차등 프라이버시 적용을 위해 가공한 개인신용정보 데이터 .....	9
<그림 9> 4차원 데이터 분포 그래프 .....	9
<그림 10> 연합학습 적용을 위해 가공한 개인신용정보 데이터 .....	10
<그림 11> 모델의 추천 과정 .....	11
<그림 12> 본 연구에서 실험 해 볼 연합학습 과정 .....	12
<그림 13> 본 연구에서 진행 할 차등프라이버시 적용 된 연합학습 과정 ...	12
<그림 14> 전처리 및 정규화 결과 데이터 .....	14
<그림 15> 전처리 및 정규화 결과 데이터의 4차원 분포 그래프 .....	14
<그림 16> K-means Clustering 결과 .....	15
<그림 17> 차등 프라이버시 적용 K-means Clustering 결과 .....	15
<그림 18> 딥러닝 모델의 loss와 accuracy .....	17
<그림 19> 연합학습 모델의 loss와 accuracy .....	18
<그림 20> 차등 프라이버시가 적용된 연합학습 모델의 loss와 accuracy ...	19
<그림 21> epoch 200까지 수행한 차등 프라이버시 적용 연합학습 모델의 loss와 accuracy .....	19

# I. 서론

우리는 다양하고 수많은 데이터들로 둘러싸인 일상을 살아가고 있다. SNS(Social Network Service) 그리고 각종 온라인 플랫폼의 발달과 그 수요의 증가로 다양한 데이터의 확보가 가능해진 데이터 전성기 시대가 도래하게 되었으며, 대량의 데이터와 컴퓨팅 파워를 필요로 하는 인공지능 기술의 발전 역시 빠르게 이루어지고 있다. 빅데이터와 인공지능 기술의 발전으로 각종 업계에서는 개인화에 대응하며 인공지능을 적용한 개인화 서비스를 적용시키고 있다.

그러나 이러한 인공지능 기술의 발전은 대량의 데이터를 바탕으로 이루어지기 때문에 자연스레 데이터 주체자들의 프라이버시 문제로도 이어지고 있다. 서비스를 사용하기 위해 사용자들은 자신의 개인 정보를 밝혀야 하며 이러한 정보들이 데이터로 활용되고 개인의 식별 정보로 재사용되면서 프라이버시 침해가 발생하게 되는 것이다. 서비스를 제공하는 기업의 입장에서는 사용자의 정보가 안전하게 활용될 것임을 약속하지만 악의를 가진 공격자가 해킹 공격을 가하여 프라이버시가 침해되는 위험성을 배제할 수 없다.

본 논문에서는 빅데이터 시대의 프라이버시 문제를 해결하기 위한 방법에 대해 연구한 내용을 다룰 것이며 연구한 내용을 바탕으로 프라이버시 침해 문제를 해결할 수 있는 동시에 사용자의 특성과 상품 유형을 고려한 맞춤형 서비스를 개발하는 것을 목적으로 한다.

2장에서는 두 가지 프라이버시 문제 해결방법을 다루고 있다. 첫 번째로 프라이버시 비식별 처리 기법인 차등 프라이버시(Differential Privacy)에 대한 내용이 서술되어 있으며, 두 번째로는 분산형 머신 러닝 접근법인 연합학습(Federated Learning)에 대한 내용이 서술되어 있다. 마지막으로 연합학습 기술에서 서버 가중치만으로도 데이터 복원을 할 수 있다는 위험성을 인지하였기에 가중치에 차등프라이버시를 적용하여 모델 역추적을 통한 원 데이터 예측을 방지하고자 했다. 3장과 4장에서는 두 가지의 실험 내용이 서술되어 있는데 첫 번째 실험에서는 차등 프라이버시 적용 전/후의 클러스터링의 결과를 비교·분석하였으며, 두 번째 실험에서는 차등프라이버시 적용 전/후의 연합학습에 대한 성능을 비교하고 결과를 분석하였다.

## II. 연구 내용

### 1. 차등 프라이버시 (Differential Privacy)

차등 프라이버시(Differential Privacy)는 데이터에 수학적 노이즈를 추가하는 개인정보 비식별화 기술 중 하나이다. 비식별화란 식별 가능한 데이터의 수정을 통해 특정 개인 식별이 불가능하도록 하는 것을 의미한다. 비식별화 기술은 크게 데이터의 신뢰성을 임의로 낮춤으로써 특정 데이터와 개인 간의 강한 연결성을 제거하는 방법인 무작위화 방법과 데이터 값을 보편적 범위나 의미로 변경하여 특정 개인을 식별하지 못하도록 하는 방법인 일반화 방법으로 나누어지는데, 이 중 차등 프라이버시는 무작위화 방식에 속한다. 데이터에 비식별화 기술을 적용하게 되었을 때 발생하는 문제점으로는 재식별 가능성을 들 수 있다. 재식별 가능성을 판별하는 개별화 가능성, 연결 가능성, 추론 가능성 이 3가지 기준으로 비식별화 기술들을 살펴볼 수 있다.

구분	개별화 가능성	연결 가능성	추론 가능성
k-익명성	X	O	O
L-다양성	X	O	△
차등 프라이버시	△	△	△
해싱/토큰화	O	O	△

〈표 1〉 비식별화 기술의 재식별 가능성 [1]

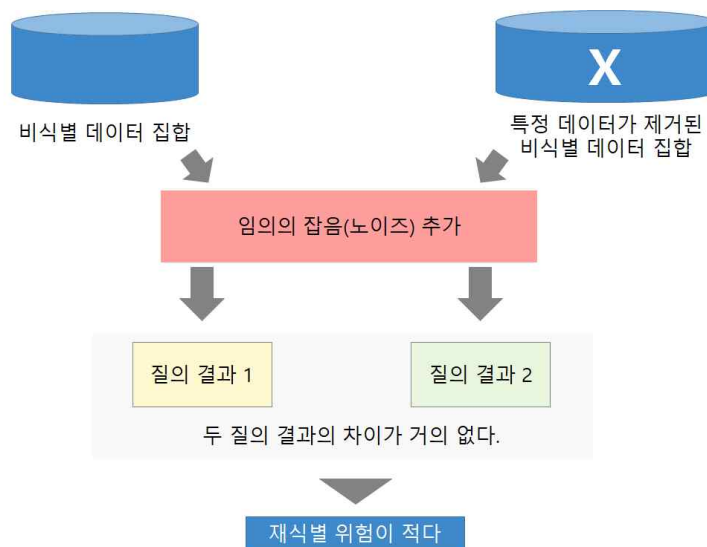
위의 표 1과 같이 차등 프라이버시의 경우에는 다른 비식별화 기술들과 비교해보았을 때 3가지 위험성에 대한 재식별 가능성이 가장 낮음을 확인할 수 있다. 즉 이전에 연구되었던 k-익명성, L-다양성의 경우 공격자가 모델의 가정 수준 이상의 배경지식을 가지고 있는 경우 개인을 식별하여 연결할 가능성이 있다는 한계를 보이지만 차등 프라이버시는 특정 개인의 존재 여부와 상관없이 비슷한 통계가 도출될 수 있도록 변조하는 방법을 통해 공격자의 배경지식과 상관없이 프라이버시를 보호할 수 있다는 점을 눈여겨 볼 만 하다.

수학적으로 차등 프라이버시는 다음 식 (1)과 같이 정의할 수 있다.

$$\Pr[f(D_1) \in S] \leq e^\epsilon \times \Pr[f(D_2) \in S] \dots (1)$$

$D_1, D_2$ 는 레코드 하나만 다른 두 개의 데이터베이스 이고,  $f$ 는 임의의 랜덤함수 그리고 데이터를 입력으로 하는 질의 결과가  $S \in \text{range}(f)$  라고 할 때, 위의 수식을 만족하면 차등 프라이버시를 제공한다고 정의할 수 있다. (이 때  $\epsilon$ 은 양의 실수이다.)

이러한 차등 프라이버시 기술은 프라이버시 뿐만 아니라 데이터의 유용성 수준을 노이즈 삽입 수준으로 동시에 고려한다는 점에서 큰 장점을 갖는다. 노이즈 삽입 수준은 차등 프라이버시의 파라미터  $\epsilon$  값으로 결정할 수 있다.  $\epsilon$  값이 0에 가까워질수록 서로 다른 입력 값에서 동일한 통계 결과가 나올 확률이 비슷해져 개인의 프라이버시 보호 수준이 강해지고,  $\epsilon$  값이 커질수록 서로 다른 입력 값에 대한 동일한 통계 결과가 나올 확률의 차가 커지면서 프라이버시 보호 수준은 약해지지만 데이터를 덜 변조하게 되면서 통계 결과의 정확성은 높아진다. 따라서 우리는 상황에 맞는 적절한 파라미터 값  $\epsilon$ 을 선정해야 한다.

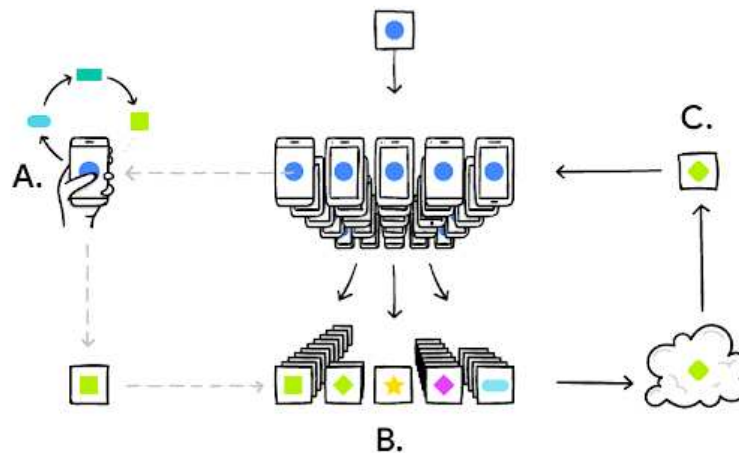


<그림 1> 차등 프라이버시를 통한 데이터 변조 과정



## 2. 연합 학습 (Federated Learning)

연합학습(Federated Learning)은 머신러닝을 중앙 클라우드가 아닌 사용자 개별 디바이스에서 스스로 데이터를 처리하고 업데이트시키는 구글에서 제안된 새로운 기법이다.

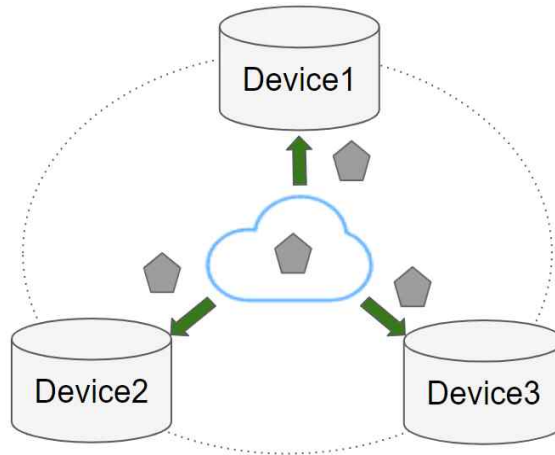


<그림 2> 연합학습 과정 [9]

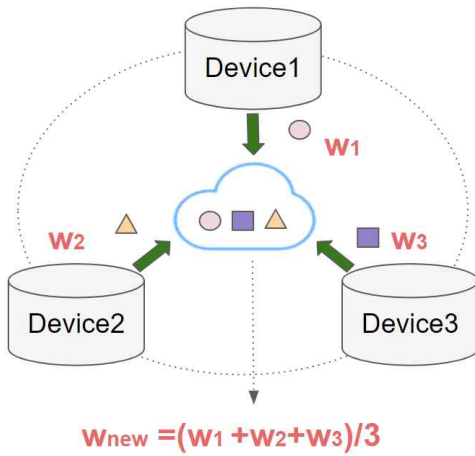
기존의 학습 방법으로는 중앙으로 데이터를 모아 학습하게 되는데, 이러한 방식이 프라이버시 침해 문제를 발생시키기도 한다. 데이터 비식별화 과정을 거쳐 전송한다고 하더라도 민감한 개인 정보가 이동하는 것은 누구라도 내키지 않을 것이다.

이러한 문제를 해결할 수 있는 방법 중 하나가 바로 연합학습이다. 앞에서 정의한 바와 같이 연합학습은 데이터를 중앙에 모으지 않고, 각 사용자의 디바이스에서 학습한 모델을 중앙으로 취합하여 학습을 수행하기 때문에 데이터의 정보보호 및 프라이버시 보호가 가능해진다.

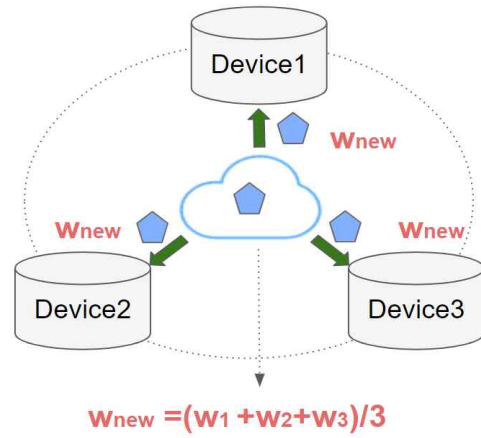
연합학습 수행 과정을 단계별로 살펴보면 다음과 같다. 먼저 <그림 3>와 같이 모델을 각 디바이스로 전송하여 사용자의 디바이스 내에서 학습을 수행하게 된다. 이때 사용자의 데이터는 외부로 유출되지 않고 온전히 사용자의 디바이스 내에서만 존재하며 그것을 바탕으로 학습이 이루어지기 때문에 프라이버시를 완벽하게 보존할 수 있다.



<그림 3> 각 디바이스에서 학습 수행



<그림 4> 학습된 파라미터를 중앙으로 전송

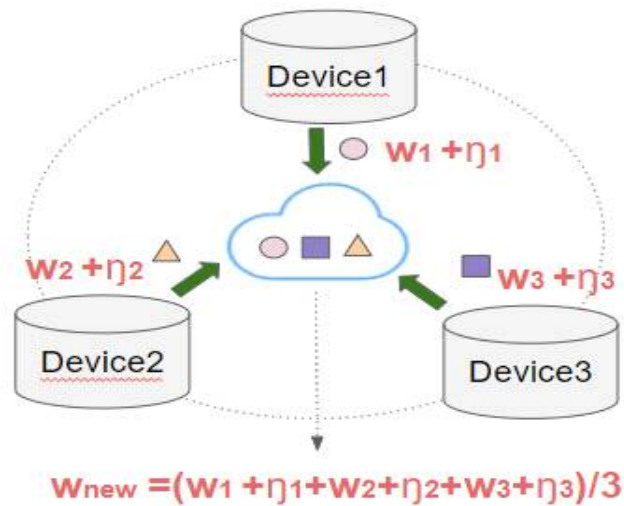


<그림 5> 업데이트된 모델 디바이스로 재전송

사용자의 디바이스 내에서 학습이 이루어지고 난 후에는 <그림 4>과 같이 각 디바이스의 학습된 파라미터는 다시 중앙 서버로 보내져 통합된다. 통합은 식(2)에 의해 이루어진다.

$$w_{\neq w} = (w_1 + w_2 + w_3 + \dots + w_n) / n \dots (2)$$

$w$ 는 weight를 나타내며,  $n$ 은 디바이스의 수를 나타낸다. 평균화 알고리즘을 통해 단순히 평균화 작업을 수행하게 되며, 해당 파라미터를 업데이트하게 된다. 이렇게 업데이트된 모델을 그림 4와 같이 각 디바이스로 재전송 해주게 되면, 최종적으로 각각의 디바이스에서 사용자는 직접적으로 데이터를 공유하지 않았지만 통일된 하나의 모델을 학습하는 효과를 얻을 수 있게 된다.



<그림 6> 전송 가중치에 차등 프라이버시를 적용한 연합학습

나아가, Federated Learning with Differential Privacy: Algorithms and Performance Analysis(Kang Wei et al. (2019))에서 정보 유출을 효과적으로 방지하기 위해 차등 프라이버시 (DP) 개념을 기반으로 집계 전 클라이언트 측에서 매개 변수에 인공 노이즈를 추가하는, 즉 모델 집계 전 노이즈 FL (NbAFL) 프레임 워크를 제안한다. 해당 과정을 수행함으로써 중간 통계를 통한 프라이버시 노출을 방지하면서 연합학습을 수행할 수 있게 된다.

### Ⅲ. 설 계

#### 1. 실험 도구

구분		항목
HW	컴퓨팅 시스템	CPU : Intel Core i5 8 <sup>th</sup> Gen RAM : 16GB
	I/O	마우스, 키보드, 모니터 등
SW	OS	Window 10
	개발 언어	python
	IDE	Jupyter Notebook
	Library	차등 프라이버시: diffprivlib-IBM 연합학습: shfl

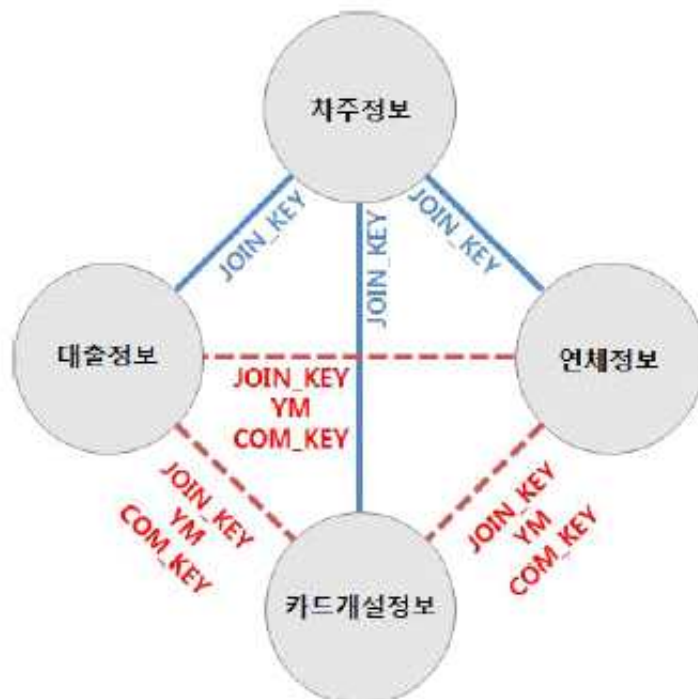
<표 2> 개발 도구

<표 2> 는 본 실험에 사용된 개발환경에 대한 내용이다.

## 2. Data Set

### 1) 한국 신용 정보원 빅데이터 센터 데이터

한국 신용 정보원 빅데이터 센터에서 제공받은 데이터를 사용하였다. <그림 7>를 참고하여 각 실험에 맞게 여러 테이블을 조인하고 불필요한 속성을 제거하였다.



<그림 7> 개인신용정보 데이터 관계 [8]

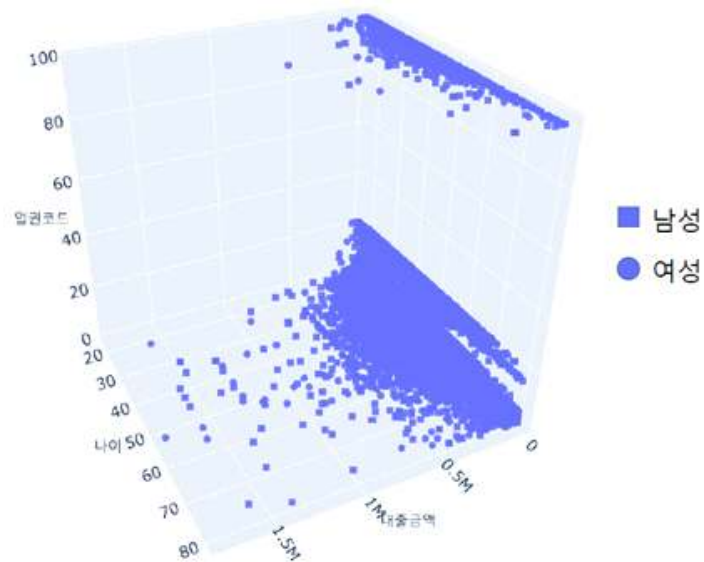
### 2) 차등 프라이버시 적용 전/후의 클러스터링

데이터 레코드 각각은 개인을 의미하고 총 5가지 속성을 가진다. 대출상품은 모델 학습에는 사용되지 않고 추후 대출상품 추천 시 해당 클러스터의 대출상품 빈도를 계산하여 추천 리스트를 출력한다.

아래의 <그림 8>은 데이터를 시각화하여 그래프로 나타낸 것이다.

	나이	성별	업권코드	대출금액	대출상품
0	50	1	3	50000	230
1	49	2	1	3500	200
2	49	2	8	5000	100
3	37	2	1	20000	240
4	29	1	21	3100	100
...	...	...	...	...	...
57672	61	2	17	6400	500
57673	50	1	1	30000	100
57674	48	1	3	30000	100
57675	48	1	3	170000	230
57676	77	2	1	37000	270

<그림 8> 차등 프라이버시 적용을 위해 가공한 개인신용정보 데이터



<그림 9> 4차원 데이터 분포 그래프

### 3) 차등 프라이버시 적용 전/후의 연합학습

클러스터링 실험에서 사용한 데이터와 동일하게 한국신용정보원 빅데이터 센터에서 제공받은 데이터 셋을 정제하여 사용하였다. 차주 정보 테이블, 대출 정보 테이블, 연체 정보 테이블, 카드 개설 정보 테이블을 조인하고 불필요한 속성을 제거한 <그림 10>로 연합학습을 수행하였다. 데이터 레코드는 사용자 개개인을 의미하고 총 10가지 속성을 가진다.

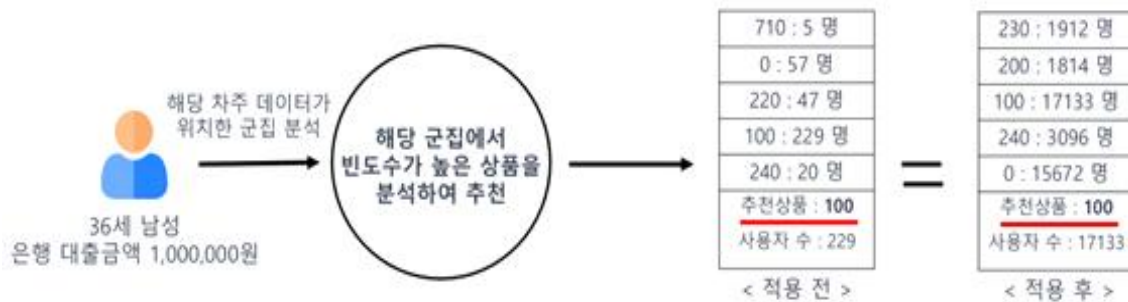
	나이	성별	업권코드	대출금액	연체유형코드	연체사유코드	등록사유코드	연체등록금액	개설사유코드	카드유형코드	대출상품
0	26	0	1	19000.0	0.0	0	0	0.0	0	0	6
1	62	1	0	0.0	0.0	0	0	0.0	1	1	0
2	62	1	0	0.0	0.0	0	0	0.0	1	1	0
3	62	1	0	0.0	0.0	0	0	0.0	1	2	0
4	62	1	0	0.0	0.0	0	0	0.0	1	1	0
...	...	...	...	...	...	...	...	...	...	...	...
6060	41	0	0	0.0	0.0	0	0	0.0	1	1	0
6061	79	0	0	0.0	0.0	0	0	0.0	1	1	0
6062	43	0	0	0.0	0.0	0	0	0.0	1	1	0
6063	36	0	0	0.0	0.0	0	0	0.0	1	1	0
6064	34	1	0	0.0	0.0	0	0	0.0	1	1	0

<그림 10> 연합학습 적용을 위해 가공한 개인신용정보 데이터

### 3. 실험 내용

#### 1) 차등 프라이버시 적용 전/후의 클러스터링

차등 프라이버시가 개인정보를 잘 보호하면서 모델의 성능에 얼마나 영향을 끼치는지 알아보기 위해 차등 프라이버시 적용 전과 후의 클러스터링을 비교해보는 실험을 진행하였다. 먼저 추천을 위해 훈련 데이터를 클러스터링(군집화) 한다. 클러스터링은 비슷한 속성을 가진 사람들을 그룹으로 묶어주는 방법이다. 본 문서에서는 클러스터링의 대표로 k-means 방식을 사용한다. k-means 알고리즘을 사용한 이유는 효율적이고 해석이 쉽기 때문이다. 특히 군집의 수를 자동화하여 결정해줄 수 있어 더욱 편리하다. 데이터를 학습에 적합한 형태로 전처리 후 동일한 훈련 데이터로 k-means와 차등 프라이버시를 적용한 k-means를 사용하여 모델을 각각 학습한다. k-means로 각 데이터의 군집을 나눈 결과를 이용해 추천 과정을 진행하였다. 추천 결과로 상품 모델을 반환하는 recommend 함수를 생성하였다. recommend 함수는 클러스터 군집에서 가장 많이 사용되는 상품을 추천해줄도록 동작한다. 동일한 테스트 데이터를 사용해 차등 프라이버시가 적용되지 않은 모델과 차등 프라이버시가 적용된 모델, 두 모델로부터 클러스터를 각각 예측한다. 예측한 recommend 함수를 수행하여 차등 프라이버시 적용 전 모델과 적용 후 모델의 추천 결과를 비교·분석한다.



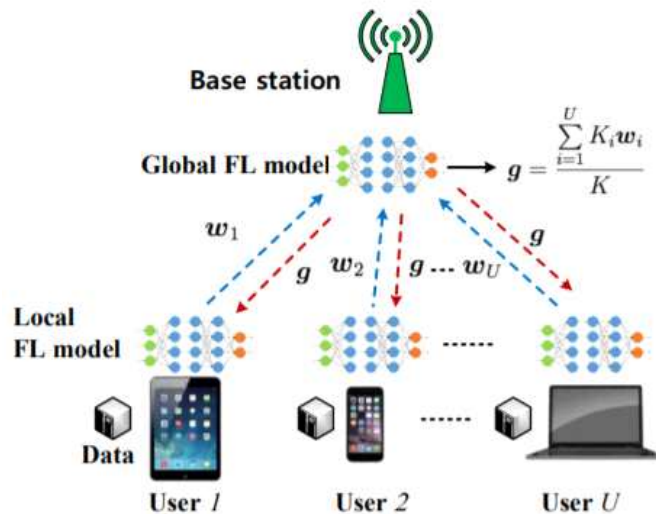
<그림 11> 모델의 추천 과정



## 2) 차등 프라이버시 적용 전/후의 연합학습

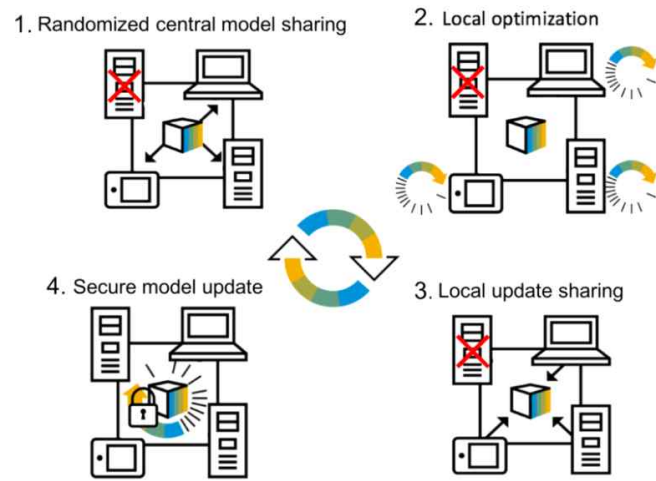
본 논문의 최종 목표인 연합학습에 차등프라이버시를 적용해 프라이버시를 보호하기 위해 연합학습을 구현하고 차등프라이버시를 적용한 연합학습과 성능을 비교해보았다. 연합학습은 오픈소스 머신러닝 플랫폼인 TensorFlow를 사용했다. 연합학습을 위해 정제된 데이터로 연합학습을 진행하여 성능을 측정하고 각 단말의 학습 파라미터를 넘겨주는 과정에 차등 프라이버시를 적용하여 연합학습을 진행한 후 성능을 측정하여 차등 프라이버시 적용 전과 후의 연합학습 성능을 비교하고 결과를 분석한다.

차등 프라이버시 적용 전 과정은 <그림 12>를 참고하며 진행하였다. 해당 알고리즘을 구현하기 위해서 Sherpa. ai의 shfl 라이브러리를 사용한다.



<그림 12> 본 연구에서 실험 해 볼 연합학습 과정 [6]

차등 프라이버시 적용 후 과정은 <그림 13>를 참고하며 진행하였다. 첫 번째 단계에서 중앙 모델에서 클라이언트로의 모델 전송 과정을 거치며 두 번째 단계에서 클라이언트에서 모델의 최적화 과정을 진행하며 세 번째 단계에서 중앙 서버로의 모델을 낸다. 마지막 네 번째 단계에서는 모델 업데이트를 차등 프라이버시 적용 전과 다르게 안전하게 진행하는데 이 과정을 클라이언트에서 차등프라이버시가 적용된 모델 가중치를 서버에 전송 시킬 때 그 값을 이용하여 모델을 업데이트 하는 과정을 거치게 된다. 차등 프라이버시 적용 전과 마찬가지로 해당 과정은 Sherpa.ai의 shfl 라이브러리를 사용해 진행한다.



<그림 13> 본 연구에서 진행 할 차등프라이버시 적용 된  
연합학습 과정 [3]

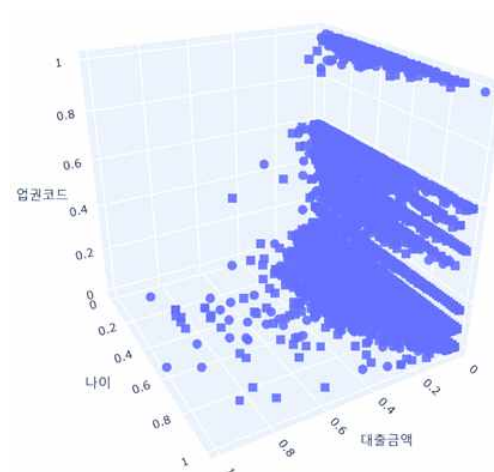
## IV. 실험 결과

### 1. 차등 프라이버시 실험 결과

모델 학습에 앞서 데이터를 학습에 적합한 형태로 전처리하고 정규화 하였다.

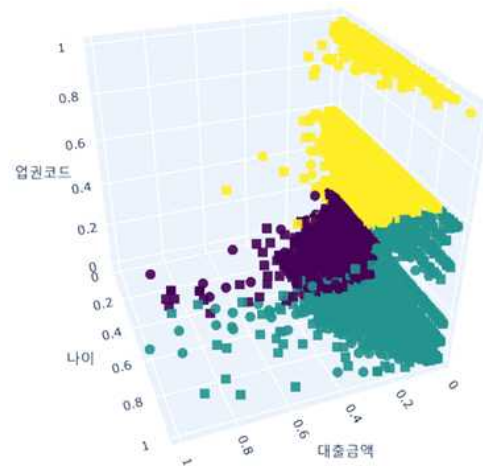
	나이	성별	업권코드	대출금액	대출상환
0	0.515625	0.0	0.083333	0.027772	230
1	0.500000	0.1	0.000000	0.001939	200
2	0.500000	0.1	0.333333	0.002772	100
3	0.312500	0.1	0.000000	0.011106	240
4	0.187500	0.0	0.500000	0.001717	100
...	...	...	...	...	...
57672	0.687500	0.1	0.416667	0.003550	500
57673	0.515625	0.0	0.000000	0.016661	100
57674	0.484375	0.0	0.083333	0.016661	100
57675	0.484375	0.0	0.083333	0.094439	230
57676	0.937500	0.1	0.000000	0.020550	270

<그림 14> 전처리 및 정규화 결과 데이터



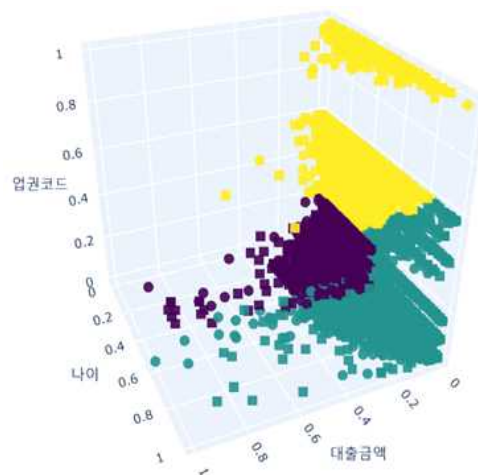
<그림 15> 전처리 및 정규화 결과 데이터의 4차원 분포 그래프

Clustering 기법으로 K-means를 사용하였다. 이때 elbow 함수를 이용해 적절한 k 값을 얻었다.



<그림 16> K-means Clustering 결과

K-means 알고리즘에 차등 프라이버시를 적용하여 군집화를 진행하였다. 이때 차등 프라이버시의 noise의 정도를 결정하는 epsilon 값으로 0.7을 주었다.



<그림 17> 차등 프라이버시 적용  
K-means Clustering 결과

<표 3>은 테스트 데이터에 대한 차등 프라이버시 적용 전과 후의 상품 추천 결과이다.

차등프라이버시 적용 전		차등프라이버시 적용 후	
대출상품 코드	사용자 수	대출상품 코드	사용자 수
0	6917명	0	6361명
100	6084명	100	5729명
220	3819명	220	3574명
230	1240명	230	1155명
240	985명	240	914명
추천 상품: 0		추천 상품: 0	

<표 3> 차등 프라이버시 적용 전과 후의 상품 추천 결과

차등 프라이버시 적용 전과 후의 추천 결과를 비교해보면 상품 사용자의 수에서 약간의 차이가 발생했지만 상품 추천의 우선순위는 동일함을 알 수 있다.

위의 실험을 통해 차등 프라이버시 적용은 민감한 정보를 보호하여 안전성을 높이는 동시에 적용 전/후의 결과 값에 영향을 주지 않아 데이터의 유용성을 훼손시키지 않음을 확인할 수 있다.

## 2. 연합학습 실험 결과

### 1) 딥러닝 모델

연합학습 모델, 차등프라이버시가 적용된 연합학습 모델과의 성능 비교를 위해 keras 라이브러리의 Sequential 모델을 사용해 딥러닝 모델을 생성하였다. 훈련 데이터로 100Epoch 만큼 훈련을 하고 10Epoch마다의 Loss, Accuracy를, 검증 데이터를 통하여 10Epoch마다의 Val\_loss와 Val\_accuracy를 출력하고 값을 확인하였다. 100Epoch을 진행하고 딥러닝 모델의 최종 Loss는 0.286813, 최종 Accuracy는 0.893757임을 <그림 18>에서 확인 할 수 있다.

	loss	accuracy	val_loss	val_accuracy
10	0.428391	0.859835	0.545459	0.835165
20	0.383477	0.871143	0.541318	0.846154
30	0.359497	0.874912	0.578755	0.847802
40	0.338095	0.884570	0.610851	0.846154
50	0.326082	0.882921	0.649346	0.853846
60	0.317531	0.886926	0.669423	0.852747
70	0.310580	0.887397	0.707870	0.846154
80	0.302665	0.887632	0.749481	0.850000
90	0.290680	0.895877	0.781243	0.843407
100	0.286813	0.893757	0.832614	0.845604

<그림 18> 딥러닝 모델의 loss와 accuracy

### 2) 연합학습 모델

연합학습 모델에서는 분산된 클라이언트에서 모델 파라미터만을 가져와 중앙 서버에서 평균화 방식 등을 이용해 학습하게 된다. 해당 실험은 딥러닝 모델과 동일한 상황에서 수행되도록 하였다. 연합학습 모델에서 중앙 모델의 가중치 갱신을 의미하는 Round를 딥러닝 모델의 Epoch과 같은 의미로 판단하여 딥러닝 모델에서 수행한 100Epoch을 연합학습 모델에서는 100Rounds로 수행하였다. 딥러닝 모델에서의 실험과 마찬가지로 10Rounds마다의 Loss, Accuracy를 확인하였고 최종 Loss는 0.761186, 최종 Accuracy는 0.805441임을 <그림 19>에서 확인할 수 있다.

	loss	accuracy
10	2.135458	0.655400
20	1.439673	0.657873
30	1.094696	0.676010
40	0.973483	0.700742
50	0.905005	0.777411
60	0.856158	0.790602
70	0.818512	0.800495
80	0.792516	0.803792
90	0.774349	0.804617
100	0.761186	0.805441

<그림 19> 연합학습 모델의  
loss와 accuracy

### 3) 차등 프라이버시 적용 연합학습 모델

차등 프라이버시 적용 연합학습 모델에서는 클라이언트에서 중앙 서버로 전송하는 가중치에 차등프라이버시를 적용하여 학습을 진행하게 된다. 해당 실험은 딥러닝, 연합학습 모델 실험과 마찬가지로 10Rounds마다 Loss, Accuracy를 출력하였고 학습을 완료하였다. 100Rounds 학습 후, 최종 Loss는 1.556252이고 최종 Accuracy는 0.769992임을 <그림 20>에서 확인할 수 있다. 전송 가중치에 차등 프라이버시 적용으로 인하여 감소된 Accuracy는 Rounds의 수를 늘려 향상시킬 수 있다는 것을 <그림 21>의 실험을 통해 확인 할 수 있다.

	loss	accuracy
10	1.967384	0.624897
20	1.765698	0.661171
30	1.767950	0.638087
40	1.477620	0.662819
50	1.027212	0.684254
60	1.037482	0.691674
70	0.917180	0.799670
80	0.925493	0.801319
90	1.362558	0.773289
100	1.556252	0.769992

<그림 20> 차등 프라이버시가 적용된 연합학습 모델의  
loss와 accuracy

	loss	accuracy
110	1.189343	0.789777
120	1.292304	0.767519
130	1.045923	0.792251
140	0.897034	0.787304
150	0.932219	0.801319
160	1.279750	0.772465
170	0.944983	0.800495
180	0.890658	0.790602
190	0.888937	0.801319
200	0.789528	0.807090

<그림 21> epoch 200까지 수행한 차등 프라이버시 적용  
연합학습 모델의 loss와 accuracy



## V. 정 리

본 논문에서는 데이터 비식별화 기술 알고리즘과 연합학습 적용을 통한 프라이버시 보호 실험을 수행하였다. 기존 비식별화 기술인 K-익명성, L-다양성, 해싱/토큰화 등 보다 재식별 위험성이 낮은 차등 프라이버시를 제안하였다. 데이터에 적절한 임의의 noise를 추가하여 안전성을 높이면서도 데이터 활용에 있어서 결과 값에는 영향을 미치지 않기 때문에 데이터의 유용성을 훼손시키지 않을 수 있었으며, 또한 분산형 인공지능 학습인 연합학습을 통해 중앙으로 데이터를 직접적으로 공유하지 않아 프라이버시 보호라는 실험의 목적에 도달할 수 있었다.

마지막으로 이를 종합하여 연합학습 수행 과정에서 중앙 클라우드로 전달되는 사용자의 각 학습 파라미터에 차등 프라이버시를 적용시킴으로써 데이터의 유용성을 보존하는 동시에 민감한 데이터는 보호하며, 보호된 데이터를 바탕으로 사용자에게는 적절한 상품을 추천해 줄 수 있는 결과를 얻을 수 있었다. 이는 향후 다른 여러 분야의 민감한 데이터를 활용할 때에도 적용할 수 있을 것으로 보인다.

## 참고문헌

- [1] C. Dwork, “Differential privacy,” Automata, Languages and Programming, Lecture Notes in Computer Science vol. 4052, M. Bugliesi et al., editors, Heidelberg: Springer, pp. 1-12, 2006.
- [2] A. Hard, K. Rao, R. Mathews, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, “Federated learning for mobile keyboard prediction” , 2018. [Online], Available: arXiv:1811.03604
- [3] Brendan McMahan and Daniel Ramage, “Federated Learning: Collaborative Machine Learning without Centralized Training Data”
- [4] Kang Wei, Jun Li, Ming Ding, Chuan Ma, Howard H. Yang, Farokhi Farhad, Shi Jin, Tony Q. S. Quek, H. Vincent Poor, “Federated Learning with Differential Privacy: Algorithms and Performance Analysis” , 2019.11.1.
- [5] ARTICLE 29 DATA PROTECTION WORKING PARTY, “Opinion on Anonymisation Techniques” , 2014.4.10.
- [6] Mingzhe Chen, Zhaohui Yang, Walid Saad, Changchuan Yin, H. Vincent Poor, Shuguang Cui, “A Joint Learning and Communications Framework for Federated Learning over Wireless Networks” , 2019.09.17., arXiv:1909.07972
- [7] N. Rodriguez-Barroso et al., “Federated learning and differential privacy: Software tools analysis, the Sherpa.ai FL framework and methodological guidelines for preserving data privacy,” Information Fusion, vol. In press, 2020.
- [8] 한국신용정보원(www.kcredit.or.kr), 개인신용정보 표본DB 매뉴얼
- [9] ai.googleblog.com