# Federated Learning

**2025. 03. 07**

**Korea University**

**D**ata **M**ining & **Q**uality **A**nalytics Lab.

**최지형**

Data Mining
Quality Analytics

# 발표자 소개



❖ **최지형 (Jihyung Choi)**
- 고려대학교 산업경영공학과 대학원 재학
- Data Mining & Quality Analytics Lab. (김성범 교수님)
- M.S Student (2024.09 ~ Present)

❖ **Research Interest**
- Federated Learning
- Fine-tuning Foundation Models
- Agent AI

❖ **Contact**
- jibro@korea.ac.kr

Data Mining
Quality Analytics

# Introduction

❖ **Federated Learning (FL)**

*Q.* 데이터가 부족한가?



"우리 집 고양이 츄르를 좋아해."

# Introduction

❖ **Federated Learning (FL)**

*Q.* 데이터가 부족한가?



"우리 집 고양이 츄르를 좋아해"

데이터 생성

# Introduction

❖ **Federated Learning (FL)**

*Q.* **데이터가 부족한가?**

# Introduction

❖ **Federated Learning (FL)**

*Q.* 데이터가 부족한가?

# Introduction

❖ **Federated Learning (FL)**

*Q.* 데이터가 부족한가?

*A.* 데이터 활용 방법이 부족하다!

# Introduction

❖ **Federated Learning (FL)**

# Introduction

❖ **Federated Learning (FL)**

# Introduction

❖ **Federated Learning (FL)**

➢ **데이터가 분산된 환경에, 데이터 공유 없이 모델 학습!!**

# Introduction

Federated Learning vs. Distributed Learning

❖ **Distributed Learning:** 데이터를 여러 연산 장치로 나누어 학습

➢ **데이터가 분산된 환경에서 학습한다는 점에서,** FL과 동일

| 데이터가 분산된 환경에서 학습 | |
|---|---|
| 단일 디바이스 | 여러 디바이스 |
| 프라이버시 문제 X | 프라이버시 문제 O |

**Distributed Learning**

GPU 0　GPU 1　GPU 2　GPU 3

**Federated Learning**

Data Mining
Quality Analytics

# Introduction

❖ **Distributed Learning:** 데이터를 여러 연산 장치로 나누어 학습

➢ **데이터가 분산된 환경에서 학습한다는 점에서,** FL과 동일

| 데이터가 분산된 환경에서 학습 | |
|:---:|:---:|
| 단일 디바이스 | 여러 디바이스 |
| 프라이버시 문제 X | 프라이버시 문제 O |

Distributed Learning

G P U 0  G P U 1  G P U 2  G P U 3

Federated Learning

# Introduction

❖ **Distributed Learning:** 데이터를 여러 연산 장치로 나누어 학습

➢ **데이터가 분산된 환경에서 학습한다는 점에서,** FL과 동일

| 데이터가 분산된 환경에서 학습 | |
|---|---|
| 단일 디바이스 | 여러 디바이스 |
| 프라이버시 문제 X | 프라이버시 문제 O |

| communication cost |
|---|
| partial participation |

**Distributed Learning**

| G P U 0 | G P U 1 | G P U 2 | G P U 3 |

**Federated Learning**

Data Mining
Quality Analytics

# Introduction

❖ **Distributed Learning:** 데이터를 여러 연산 장치로 나누어 학습

➢ **데이터가 분산된 환경에서 학습한다는 점에서,** FL과 동일

| 데이터가 분산된 환경에서 학습 | |
|---|---|
| 단일 디바이스 | 여러 디바이스 |
| 프라이버시 문제 X | 프라이버시 문제 O |

communication cost

partial participation

**Distributed Learning**

G
P
U
0

G
P
U
1

G
P
U
2

G
P
U
3

**Federated Learning**

On   Off   On   Off

Data Mining
Quality Analytics

# Introduction

❖ **Distributed Learning:** 데이터를 여러 연산 장치로 나누어 학습

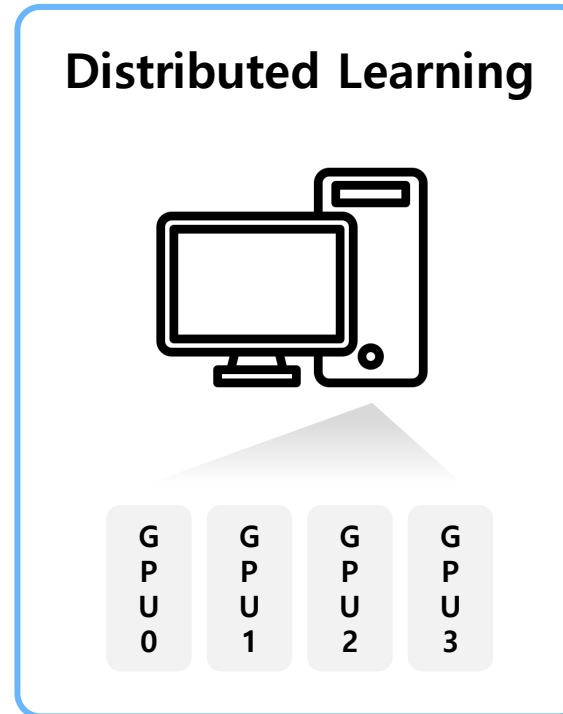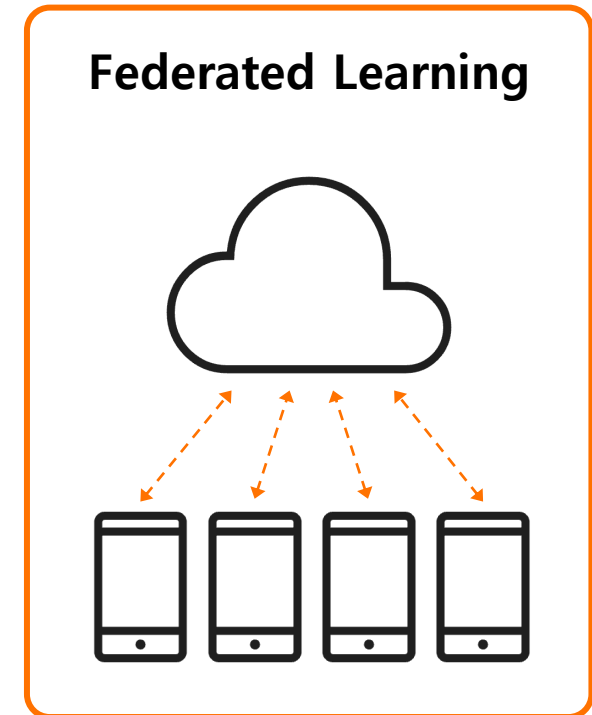➢ **데이터가 분산된 환경에서 학습한다는 점에서,** FL과 동일

# Introduction

Federated Learning vs. Distributed Learning

❖ **Distributed Learning:** 데이터를 여러 연산 장치로 나누어 학습
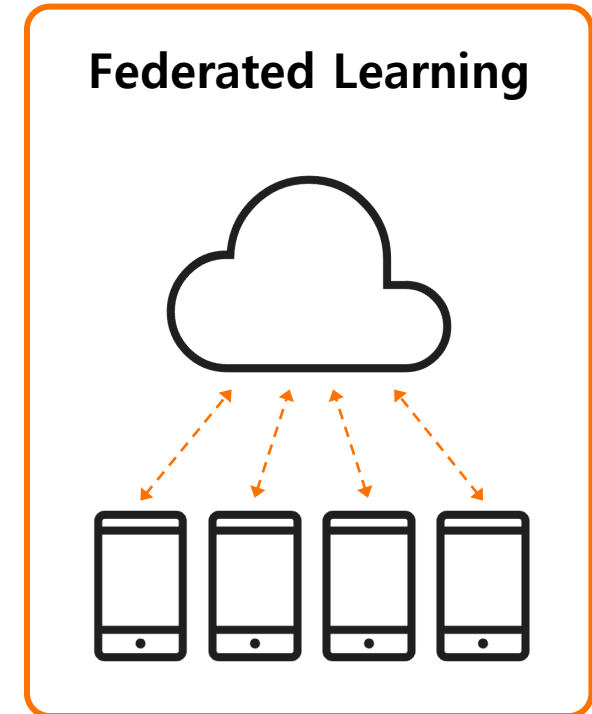
  ➢ **데이터가 분산된 환경에서 학습한다는 점에서,** FL과 동일

| 데이터가 분산된 환경에서 학습 | |
|---|---|
| 단일 디바이스 | 여러 디바이스 |
| 프라이버시 문제 X | 프라이버시 문제 O |

communication cost

partial participation

**Distributed Learning**

G P U 0  G P U 1  G P U 2  G P U 3

**Federated Learning**

각 디바이스에서 데이터 생성

# Introduction

❖ **Distributed Learning:** 데이터를 여러 연산 장치로 나누어 학습

➢ **데이터가 분산된 환경에서 학습한다는 점에서, FL과 동일**

| 데이터가 분산된 환경에서 학습 | |
|---|---|
| 단일 디바이스 | 여러 디바이스 |
| 프라이버시 문제 X | 프라이버시 문제 O |

communication cost

partial participation

Non-I.I.D

데이터 불균형

**Distributed Learning**

| G P U 0 | G P U 1 | G P U 2 | G P U 3 |
|---|---|---|---|

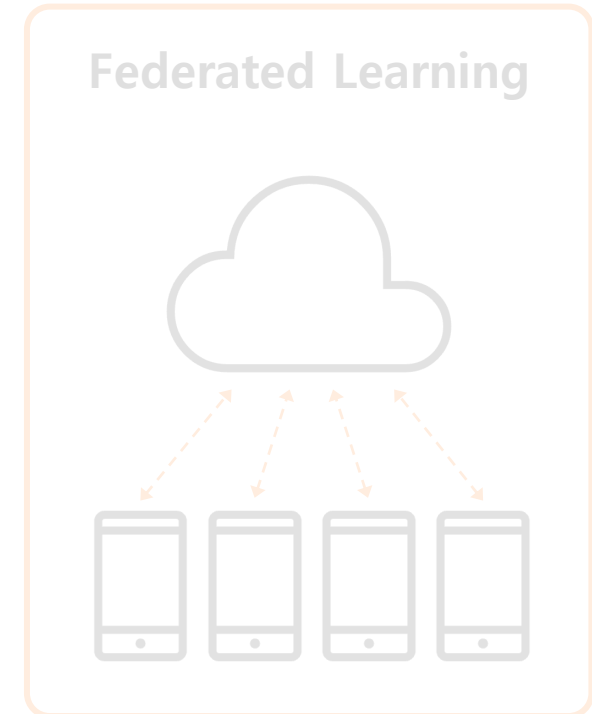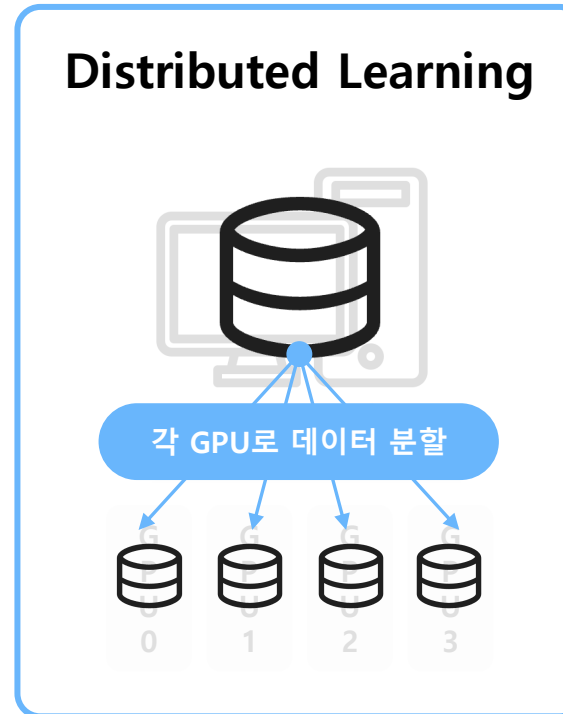**Federated Learning**

각 디바이스에서 데이터 생성

# Introduction

Federated Learning vs. Distributed Learning

❖ **Distributed Learning:** 데이터를 여러 연산 장치로 나누어 학습

➤ **데이터가 분산된 환경에서 학습한다는 점에서,** FL과 동일

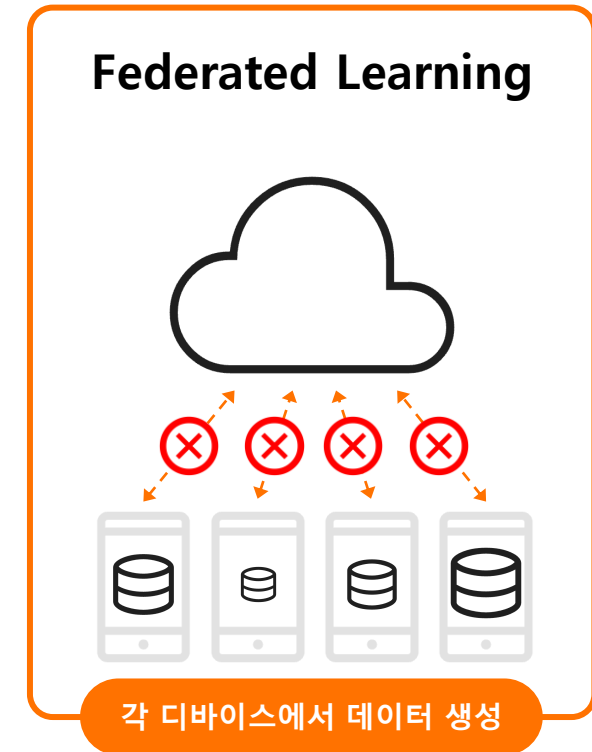| 데이터가 분산된 환경에서 학습 | |
|---|---|
| 단일 디바이스 | 여러 디바이스 |
| 프라이버시 문제 X | 프라이버시 문제 O |

FL이 해결해야 하는
4가지 문제
- communication cost
- partial participation
- Non-I.I.D
- 데이터 불균형

**Distributed Learning**

G P U 0  G P U 1  G P U 2  G P U 3

**Federated Learning**

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Federated learning 개념과 federated averaging 알고리즘 제안**

  ➤ AISTATS'17

  ➤ 피인용 21686회 (2025년 3월 기준)

McMahan, H. B., Moore, E., Ramage, D., Hampson, S., & Agüera y Arcas, B. (2017). Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, JMLR: W&CP, 54.

Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Distributed Optimization**

   ➢ **전체 데이터에 대한 손실값, $f(w)$ 최소화**

*Goal* $\qquad$ $w^* \triangleq \min_{w} f(w)$



$w$       $D$

*Data size: $n$*

*Total Loss:* $f(w) := \frac{1}{n} \sum f_i(w)$

| GPU: 1 | GPU: 2 | GPU: 3 | ... | GPU: K |

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Distributed Optimization**

➢ **전체 데이터에 대한 손실값, $f(w)$ 최소화**

*Goal* $\quad\bigg|\quad w^* \triangleq \min_w f(w)$

Data size: $n$

Total Loss: $f(w) \coloneqq \frac{1}{n}\sum f_i(w)$

**Data Partitioning**

| GPU: 1 | GPU: 2 | GPU: 3 | ... | GPU: K |

$D_1$　$D_2$　$D_3$　　$D_K$

Data size: $\frac{n}{K}$

Partial Loss:  $F_k(w) = \frac{K}{n}\sum_{i \in D_k} f_i(w)$

# Federated Averaging (FedAvg)

Communication-Efficient Learning of Deep Networks from Decentralized Data

❖ **Distributed Optimization**

➢ **전체 데이터에 대한 손실값, $f(w)$ 최소화**

$$Goal \quad \Big| \quad w^* \triangleq \min_{w} f(w)$$



Data size: $n$
Total Loss: $f(w) := \frac{1}{n} \sum f_i(w)$

**Data Partitioning**

| GPU: 1 | GPU: 2 | GPU: 3 | ... | GPU: K |

$D_1$    $D_2$    $D_3$    $D_K$

Data size: $\frac{n}{K}$
Partial Loss: $F_k(w) = \frac{K}{n} \sum_{i \in D_k} f_i(w)$

$$f(w) = \frac{1}{K} \sum_{k=1}^{K} F_k(w)$$
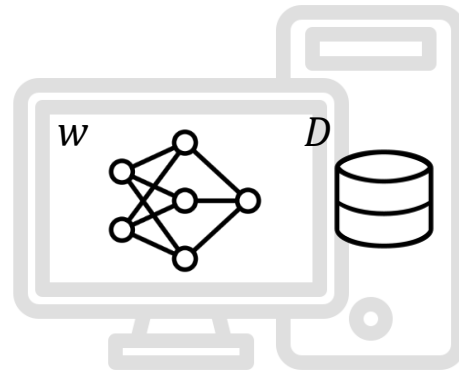
Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Distributed Optimization**

➢ 전체 데이터에 대한 손실함수 값, $f(w)$ 최소화

*Goal* $\quad\bigg|\quad$ $w^* \triangleq \min\limits_{w} f(w)$

**Gradient Descent**

$$w_{t+1} = w_t - \eta \underbrace{\nabla f(w_t)}_{\text{Total Loss Gradient}}$$

$$= w_t - \eta \nabla \{\frac{1}{K} \sum_{k=1}^{K} F_k(w_t)\}$$

$$= w_t - \frac{1}{K} \sum_{k=1}^{K} \eta \nabla F_k(w_t)$$

$$= \frac{1}{K} \sum_{k=1}^{K} (w_t - \eta \nabla F_k(w_t))$$
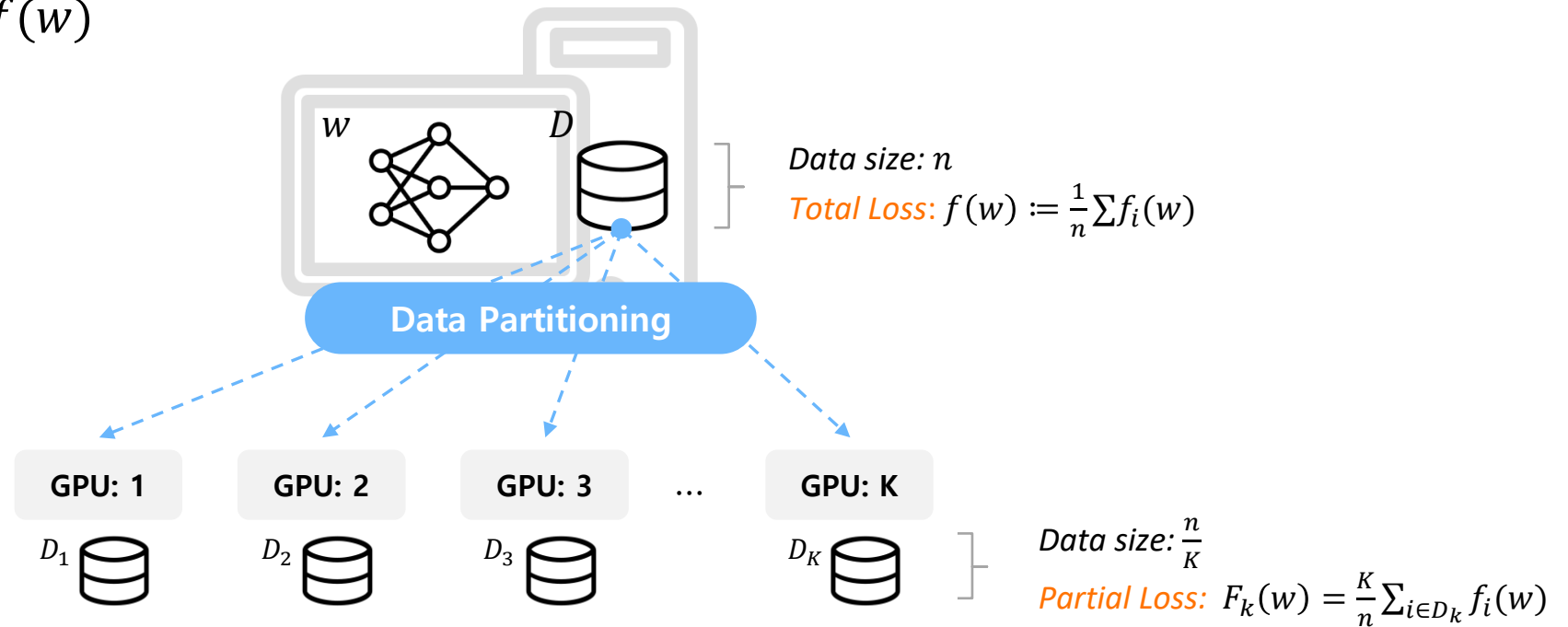
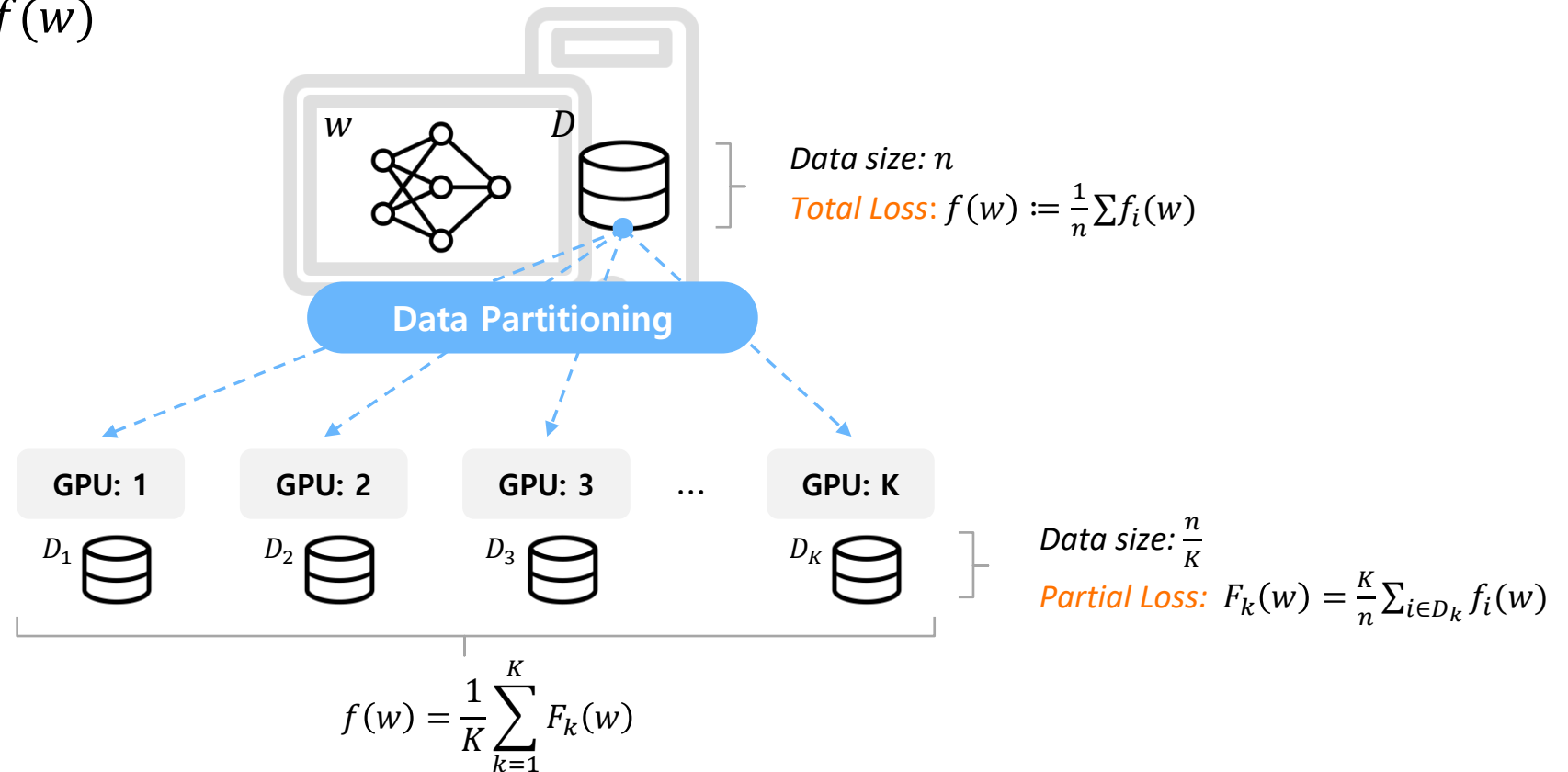$$= \frac{1}{K} \sum_{k=1}^{K} w_{t+1}^{k}$$

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Distributed Optimization**

➢ **전체 데이터에 대한 손실함수 값, $f(w)$ 최소화**

*Goal*

$$w^* \triangleq \min_w f(w)$$

**Gradient Descent**

$$w_{t+1} = w_t - \eta \boxed{\nabla f(w_t)} \text{ Total Loss Gradient}$$

$$= w_t - \eta \nabla \{ \frac{1}{K} \sum_{k=1}^{K} F_k(w_t) \}$$

$$f(w) = \frac{1}{K} \sum_{k=1}^{K} F_k(w)$$

$$= w_t - \frac{1}{K} \sum_{k=1}^{K} \eta \nabla F_k(w_t)$$

$$= \frac{1}{K} \sum_{k=1}^{K} (w_t - \eta \nabla F_k(w_t))$$

$$= \frac{1}{K} \sum_{k=1}^{K} w_{t+1}^k$$

Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Distributed Optimization**

➢ 전체 데이터에 대한 손실함수 값, $f(w)$ 최소화

*Goal* $\quad\bigg|\quad$ $w^* \triangleq \min_{w} f(w)$ **Gradient Descent**

$$w_{t+1} = w_t - \eta \boxed{\nabla f(w_t)} \text{ \textit{Total Loss Gradient}}$$

$$= w_t - \eta \nabla\{\frac{1}{K}\sum_{k=1}^{K} F_k(w_t)\}$$

$$f(w) = \frac{1}{K}\sum_{k=1}^{K} F_k(w)$$

$$= w_t - \frac{1}{K}\sum_{k=1}^{K} \eta \nabla F_k(w_t)$$

$$= \frac{1}{K}\sum_{k=1}^{K} (w_t - \eta \nabla F_k(w_t))$$

$$= \frac{1}{K}\sum_{k=1}^{K} w_{t+1}^k$$

Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Distributed Optimization**

➢ 전체 데이터에 대한 손실함수 값, $f(w)$ 최소화

*Goal*  |  $w^* \triangleq \min_{w} f(w)$

**Gradient Descent**

$$w_{t+1} = w_t - \eta \boxed{\nabla f(w_t)} \text{ \textit{Total Loss Gradient}}$$

$$= w_t - \eta \nabla \{ \frac{1}{K} \sum_{k=1}^{K} F_k(w_t) \}$$

$$f(w) = \frac{1}{K} \sum_{k=1}^{K} F_k(w)$$

$$= w_t - \frac{1}{K} \sum_{k=1}^{K} \eta \nabla F_k(w_t)$$

$$= \frac{1}{K} \sum_{k=1}^{K} (w_t - \eta \nabla F_k(w_t))$$

$$= \frac{1}{K} \sum_{k=1}^{K} w_{t+1}^k$$

Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Distributed Optimization**

➢ 전체 데이터에 대한 손실함수 값, $f(w)$ 최소화

*Goal* | $w^* \triangleq \min_{w} f(w)$

**Gradient Descent**

$$w_{t+1} = w_t - \eta \boxed{\nabla f(w_t)} \; \textit{Total Loss Gradient}$$

$$f(w) = \frac{1}{K} \sum_{k=1}^{K} F_k(w)$$

$$= w_t - \eta \nabla \{ \frac{1}{K} \sum_{k=1}^{K} F_k(w) \}$$

$$= w_t - \frac{1}{K} \sum_{k=1}^{K} \eta \nabla F_k(w)$$

$$= \frac{1}{K} \sum_{k=1}^{K} \boxed{(w_t - \eta \nabla F_k(w_t))}$$

*New Gradient Descent???*

$$= \frac{1}{K} \sum_{k=1}^{K} w_{t+1}^k$$

Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Distributed Optimization**

➢ 전체 데이터에 대한 손실함수 값, $f(w)$ 최소화

*Goal* $\quad\Big|\quad w^* \triangleq \min_{w} f(w)$

**Gradient Descent**

$$w_{t+1} = w_t - \eta \underbrace{\nabla f(w_t)}_{\textit{Total Loss Gradient}}$$

$$= w_t - \eta \nabla \{\frac{1}{K} \sum_{k=1}^{K} F_k(w)\}$$

$$= w_t - \frac{1}{K} \sum_{k=1}^{K} \eta \nabla F_k(w)$$

*Partial Loss Gradient*

$$= \frac{1}{K} \sum_{k=1}^{K} (w_t - \eta \nabla F_k(w_t))$$

$$= \frac{1}{K} \sum_{k=1}^{K} w_{t+1}^{k}$$

$$f(w) = \frac{1}{K} \sum_{k=1}^{K} F_k(w)$$

Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Distributed Optimization**

➤ 전체 데이터에 대한 손실함수 값, $f(w)$ 최소화

*Goal*

$$w^* \triangleq \min_w f(w)$$

**Gradient Descent**

$$w_{t+1} = w_t - \eta \nabla f(w_t)_{\text{Total Loss Gradient}}$$

$D$ : $w_t \rightarrow w_{t+1}$

$$= w_t - \eta \nabla \{ \frac{1}{K} \sum_{k=1}^{K} F_k(w) \}$$

$$= w_t - \frac{1}{K} \sum_{k=1}^{K} \eta \nabla F_k(w)$$

*Partial Loss Gradient*

$$= \frac{1}{K} \sum_{k=1}^{K} (w_t - \eta \nabla F_k(w_t))$$

$D_k$ : $w_t \rightarrow w_{t+1}^k$    *Local Update*

$$= \frac{1}{K} \sum_{k=1}^{K} w_{t+1}^k$$

Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Distributed Optimization**

➢ 전체 데이터에 대한 손실함수 값, $f(w)$ 최소화

*Goal* $\qquad w^* \triangleq \min_w f(w)$

**Gradient Descent**

$$w_{t+1} = w_t - \eta \nabla f(w_t) \quad \text{\textit{Total Loss Gradient}}$$

$D$ : $w_t \to w_{t+1}$

$$= w_t - \eta \nabla \{ \frac{1}{K} \sum_{k=1}^{K} F_k(w) \}$$

$$= w_t - \frac{1}{K} \sum_{k=1}^{K} \eta \nabla F_k(w)$$

*Partial Loss Gradient*

$$= \frac{1}{K} \sum_{k=1}^{K} (w_t - \eta \nabla F_k(w_t))$$

$D_k$ : $w_t \to w_{t+1}^k$ *Local Update*

$$= \frac{1}{K} \sum_{k=1}^{K} w_{t+1}^k$$

**글로벌 업데이트 = 로컬 업데이트 수행 후 평균 !!**

Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

Communication-Efficient Learning of Deep Networks from Decentralized Data

❖ **Distributed Optimization**

➢ **전체 데이터에 대한 손실함수 값, $f(w)$ 최소화**
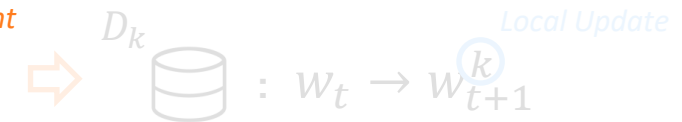
*Goal* $\quad w^* \triangleq \min_w f(w)$

**Gradient Descent**

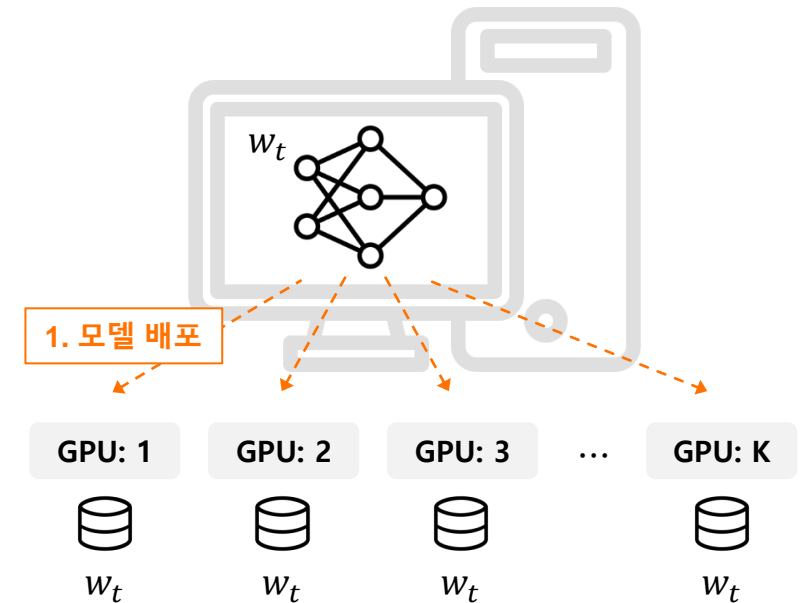$$w_{t+1} = w_t - \eta \underbrace{\nabla f(w_t)}_{\textit{Total Loss Gradient}}$$

$$= w_t - \eta \nabla \{ \frac{1}{K} \sum_{k=1}^{K} F_k(w) \}$$

$$= w_t - \frac{1}{K} \sum_{k=1}^{K} \eta \nabla F_k(w)$$

$$= \frac{1}{K} \sum_{k=1}^{K} (w_t - \eta \nabla F_k(w_t))$$

*Partial Loss Gradient*

*Local Update*

$$= \frac{1}{K} \sum_{k=1}^{K} w_{t+1}^k$$

$w_t$

**1. 모델 배포**

| GPU: 1 | GPU: 2 | GPU: 3 | ... | GPU: K |

$w_t \qquad w_t \qquad w_t \qquad\qquad w_t$

Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Distributed Optimization**

➢ **전체 데이터에 대한 손실함수 값, $f(w)$ 최소화**

*Goal* $\quad\Big|\quad w^* \triangleq \min_w f(w)$

**Gradient Descent**

$$w_{t+1} = w_t - \eta \nabla f(w_t)$$ *Total Loss Gradient*

$$= w_t - \eta \nabla \{ \frac{1}{K} \sum_{k=1}^{K} F_k(w) \}$$

$$= w_t - \frac{1}{K} \sum_{k=1}^{K} \eta \nabla F_k(w)$$

*Partial Loss Gradient*

$$= \frac{1}{K} \sum_{k=1}^{K} (w_t - \eta \nabla F_k(w_t))$$

*Local Update*

$$= \frac{1}{K} \sum_{k=1}^{K} w_{t+1}^k$$



$w_t$

**1. 모델 배포**

| GPU: 1 | GPU: 2 | GPU: 3 | ⋯ | GPU: K |

$w_t \qquad w_t \qquad w_t \qquad\qquad w_t$

**2. Local update 1회**

$w_{t+1}^1 \qquad w_{t+1}^2 \qquad w_{t+1}^3 \qquad w_{t+1}^k$

Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Distributed Optimization**

  ➢ **전체 데이터에 대한 손실함수 값, $f(w)$ 최소화**

*Goal* $\qquad w^* \triangleq \min_{w} f(w)$

**Gradient Descent**

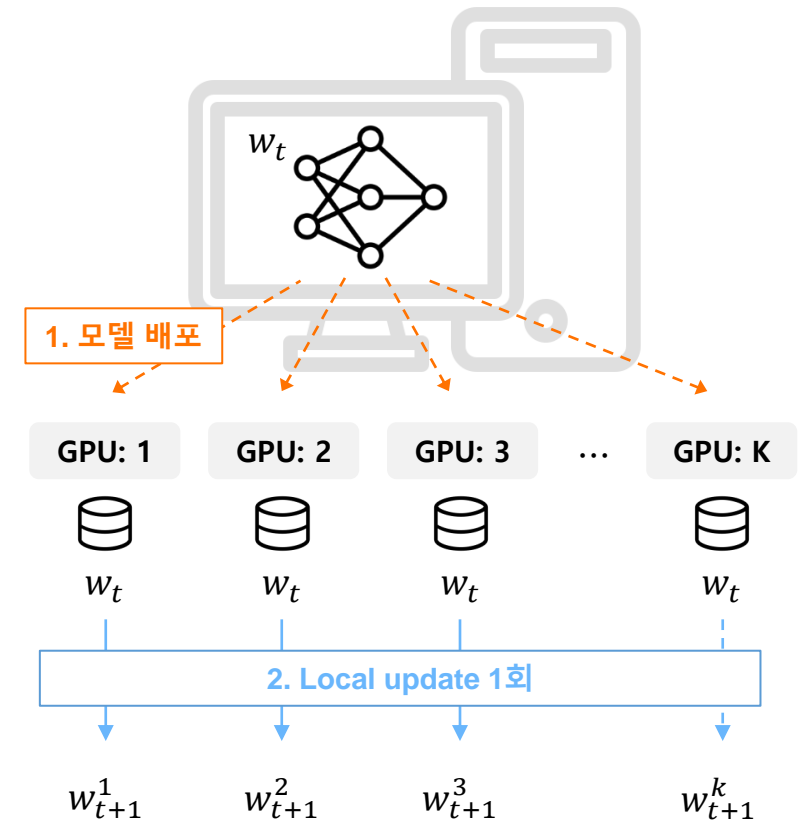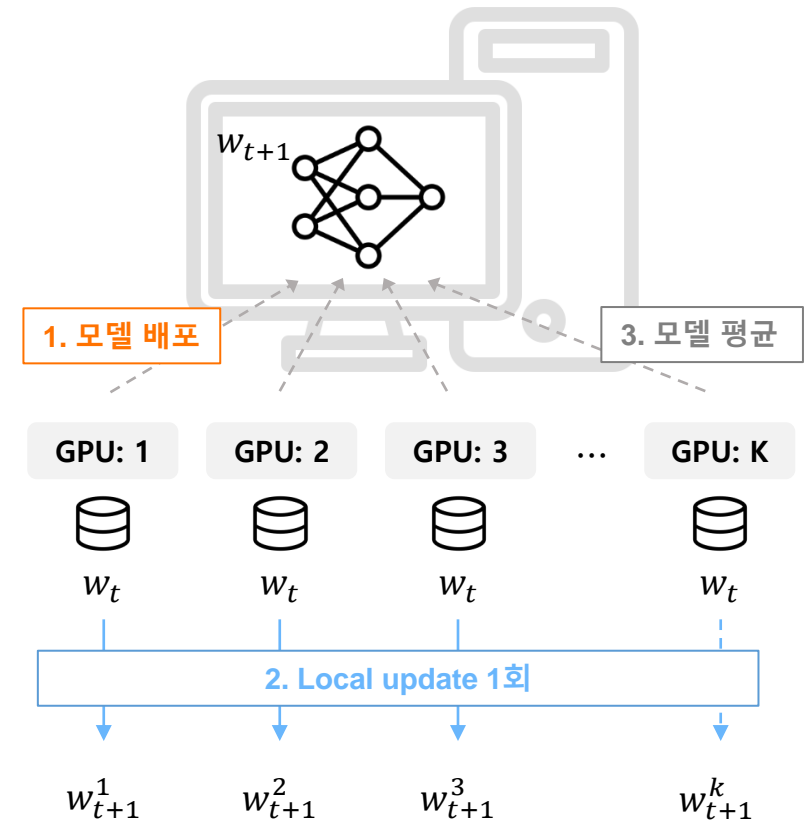$$w_{t+1} = w_t - \eta \nabla f(w_t) \quad \textit{Total Loss Gradient}$$

$$= w_t - \eta \nabla \{\frac{1}{K} \sum_{k=1}^{K} F_k(w)\}$$

$$= w_t - \frac{1}{K} \sum_{k=1}^{K} \eta \nabla F_k(w)$$

$$= \frac{1}{K} \sum_{k=1}^{K} (w_t - \eta \nabla F_k(w_t)) \quad \begin{array}{l} \textit{Partial Loss Gradient} \\ \textit{Local Update} \end{array}$$

$$= \frac{1}{K} \sum_{k=1}^{K} w_{t+1}^k$$



1. 모델 배포

3. 모델 평균

| GPU: 1 | GPU: 2 | GPU: 3 | ... | GPU: K |

$w_t \quad w_t \quad w_t \quad w_t$

2. Local update 1회

$w_{t+1}^1 \quad w_{t+1}^2 \quad w_{t+1}^3 \quad w_{t+1}^k$

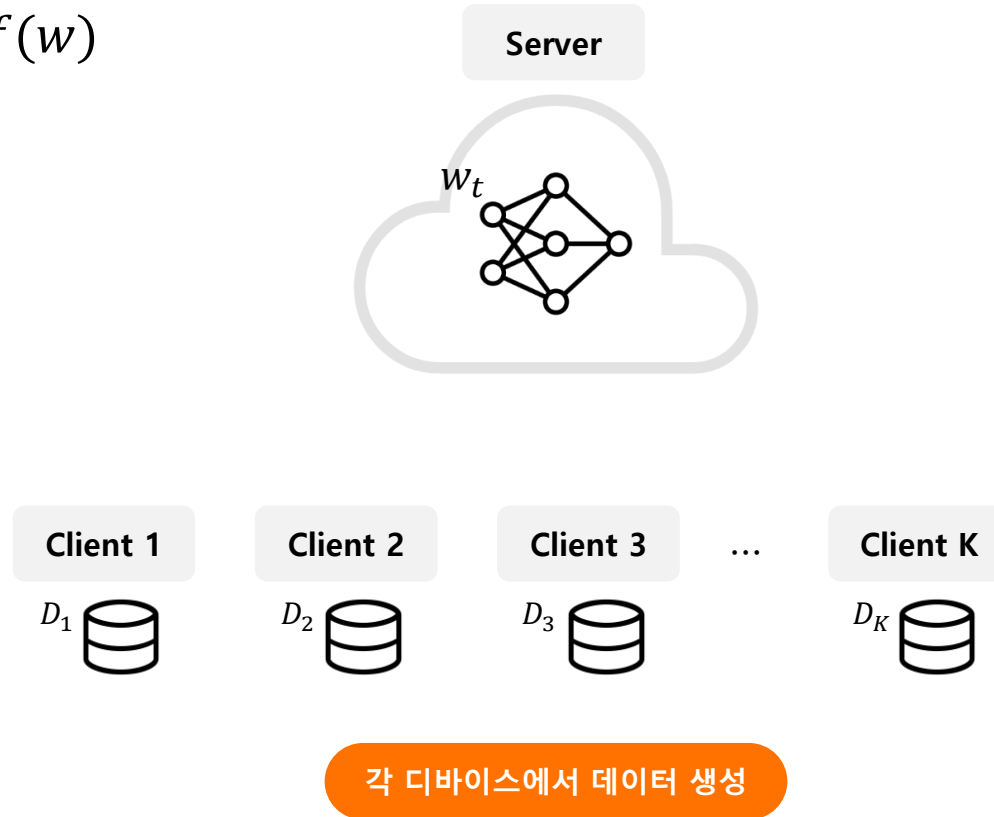Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **FedSGD**

➤ **Distributed Optimization에 partial participation과 데이터 불균형 반영**

$Goal$ $\quad \Big| \quad$ $w^* \triangleq \min_{w} f(w)$

Server

$w_t$

Client 1    Client 2    Client 3   ...   Client K

$D_1$       $D_2$       $D_3$       $D_K$

**각 디바이스에서 데이터 생성**

communication cost

partial participation

Non-I.I.D

데이터 불균형
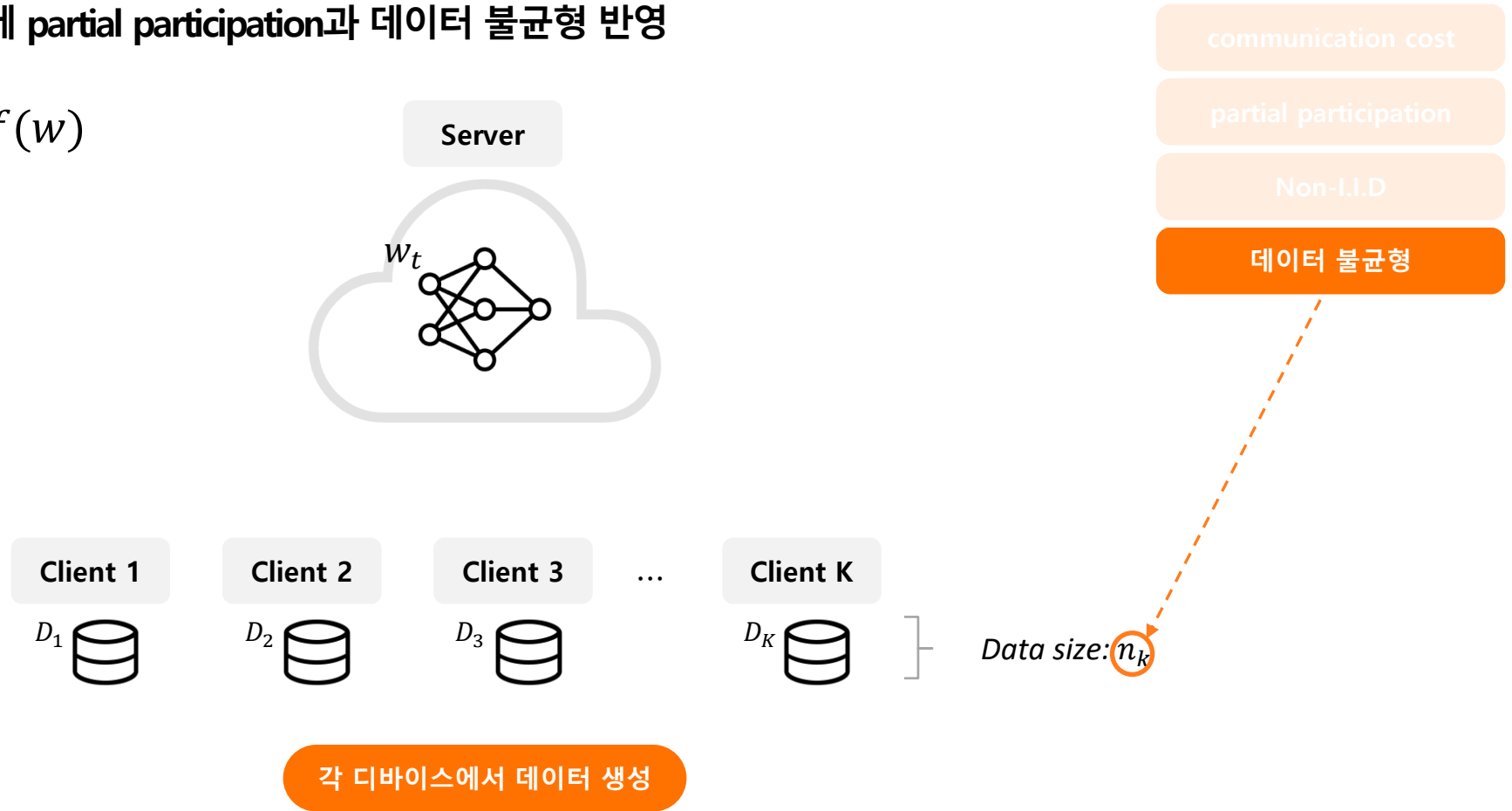
Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **FedSGD**

➢ **Distributed Optimization에 partial participation과 데이터 불균형 반영**

*Goal* | $w^* \triangleq \min_{w} f(w)$

**Server**

$w_t$

**Client 1**    **Client 2**    **Client 3**    …    **Client K**

$D_1$      $D_2$      $D_3$      $D_K$      *Data size:* $n_k$

각 디바이스에서 데이터 생성

communication cost
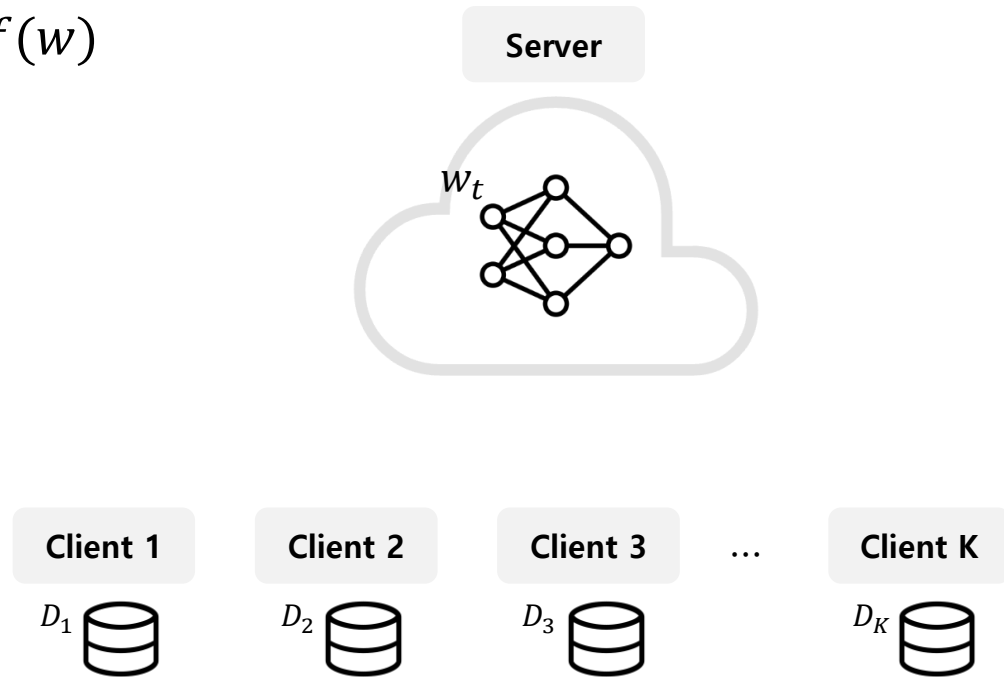
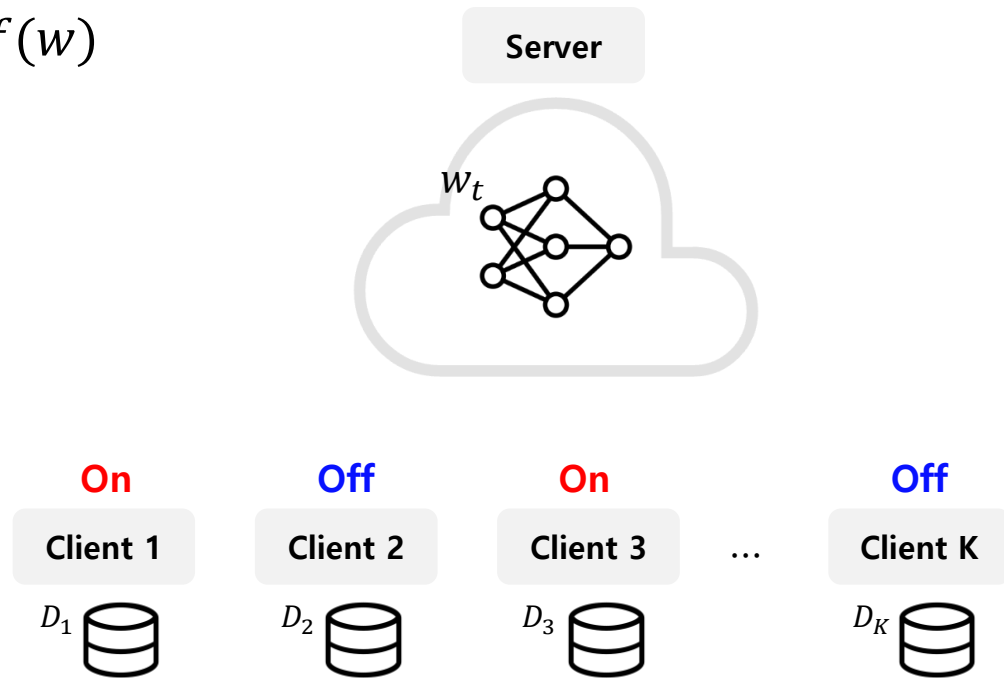partial participation

Non-I.I.D

데이터 불균형

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **FedSGD**

➢ **Distributed Optimization에 partial participation과 데이터 불균형 반영**

*Goal* $\quad\bigg|\quad w^* \triangleq \min_w f(w)$

**Server**

$w_t$



| communication cost |
| partial participation |
| Non-I.I.D |
| **데이터 불균형** |

**Client 1**   **Client 2**   **Client 3**   …   **Client K**

$D_1$   $D_2$   $D_3$   $D_K$

*Data size:* $n_k$
*Partial Loss:* $F_k(w) = \frac{1}{n_k}\sum_{i \in D_k} f_i(w)$
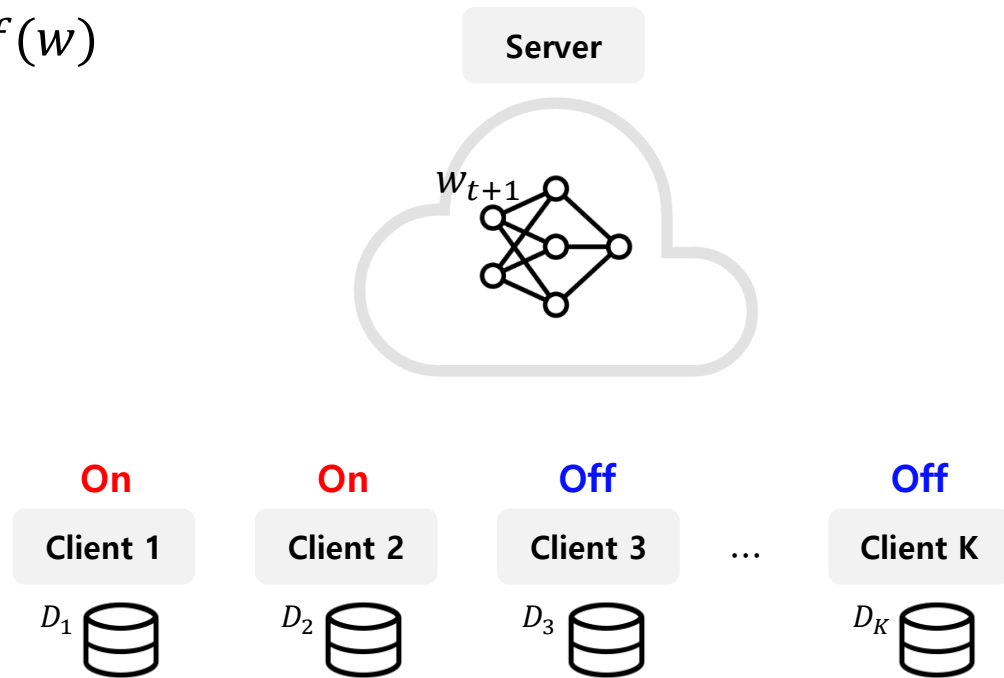
**각 디바이스에서 데이터 생성**

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **FedSGD**

➢ **Distributed Optimization에 partial participation과 데이터 불균형 반영**

$$Goal \quad \bigg| \quad w^* \triangleq \min_w f(w)$$

Server

$w_t$

| communication cost |
| --- |
| **partial participation** |
| Non-I.I.D |
| **데이터 불균형** |

**On**      **Off**      **On**      **Off**

Client 1    Client 2    Client 3   …   Client K

$D_1$     $D_2$     $D_3$     $D_K$

*Data size:* $n_k$

*Partial Loss:* $F_k(w) = \frac{1}{n_k} \sum_{i \in D_k} f_i(w)$
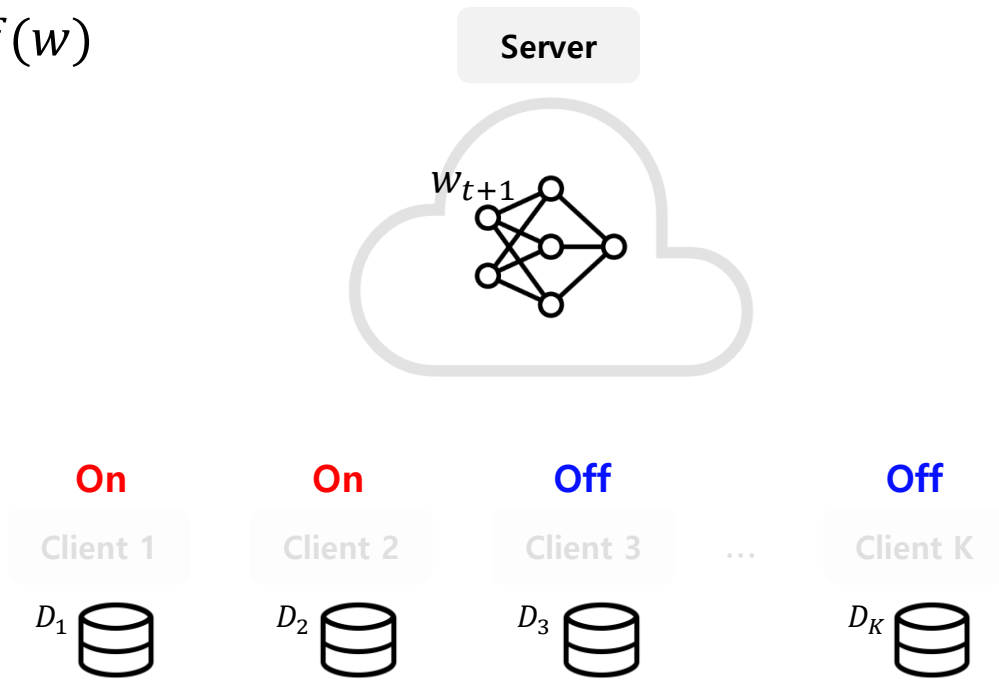
Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **FedSGD**

➢ **Distributed Optimization에 partial participation과 데이터 불균형 반영**

$$Goal \quad \Big| \quad w^* \triangleq \min_{w} f(w)$$

**Server**

$w_{t+1}$

| communication cost |
| partial participation |
| Non-I.I.D |
| 데이터 불균형 |

**On** **On** **Off** **Off**

**Client 1**   **Client 2**   **Client 3**   …   **Client K**

$D_1$   $D_2$   $D_3$   $D_K$

Data size: $n_k$

Partial Loss:  $F_k(w) = \frac{1}{n_k}\sum_{i \in D_k} f_i(w)$
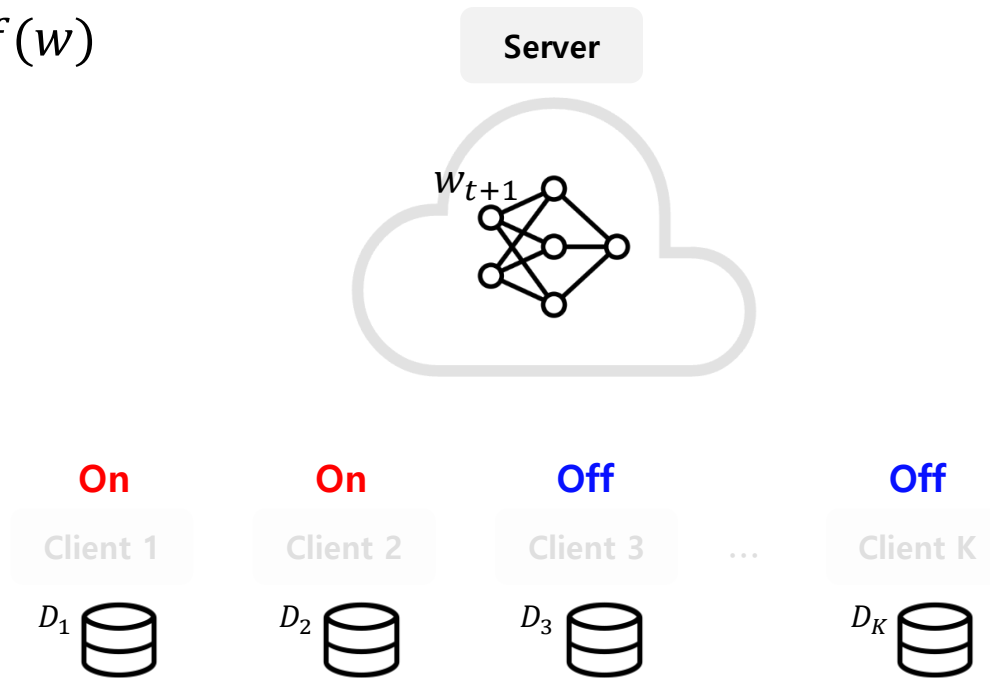
# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **FedSGD**

➢ **Distributed Optimization에 partial participation과 데이터 불균형 반영**

$$Goal \quad \Big| \quad w^* \triangleq \min_{w} f(w)$$

**Server**

$$w_{t+1}$$

| communication cost |
|---|
| **partial participation** |
| Non-I.I.D |
| **데이터 불균형** |

| **On** | **On** | **Off** | | **Off** |
|---|---|---|---|---|
| Client 1 | Client 2 | Client 3 | … | Client K |

$D_1$    $D_2$    $D_3$    $D_K$

*Data size:* $n_k$

*Partial Loss:* $F_k(w) = \frac{1}{n_k} \sum_{i \in D_k} f_i(w)$

**매 시점마다, 전체 데이터 중 일부 데이터만을 가지고 학습?**
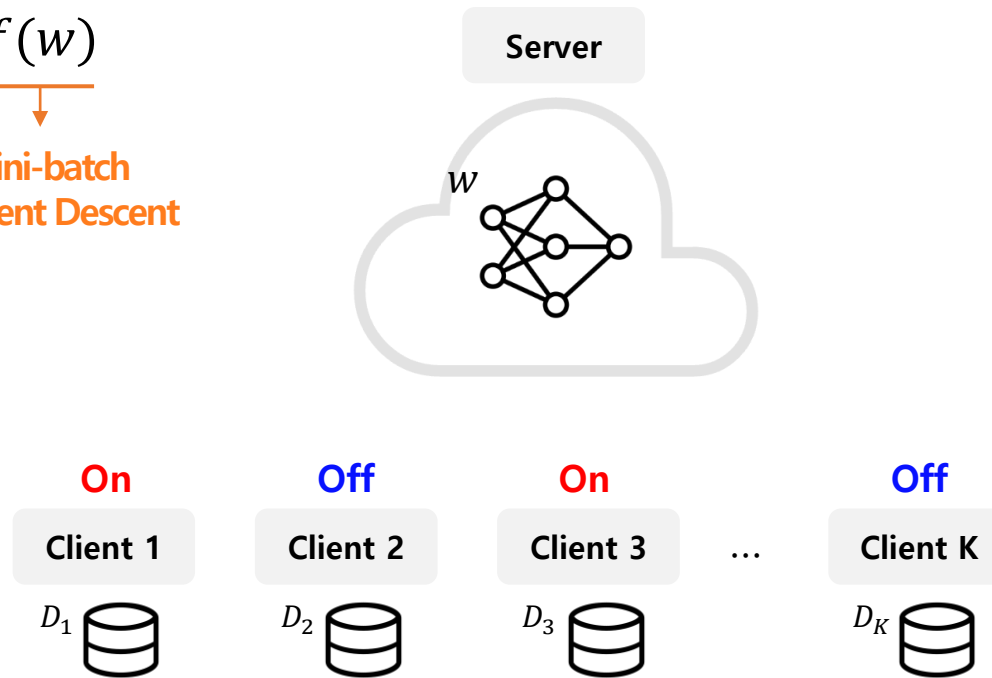
Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **FedSGD**

  ➢ **Distributed Optimization에 partial participation과 데이터 불균형 반영**



*Goal* | $w^* \triangleq \min_{w} f(w)$

**Server**

$w_{t+1}$

communication cost

**partial participation**

Non-I.I.D

**데이터 불균형**

**On**   **On**   **Off**   **Off**

Client 1   Client 2   Client 3   …   Client K

$D_1$   $D_2$   $D_3$   $D_K$

Data size: $n_k$

Partial Loss: $F_k(w) = \frac{1}{n_k} \sum_{i \in D_k} f_i(w)$

**매 시점마다, 전체 데이터 중 일부 데이터만을 가지고 학습?**

*Mini-Batch !!*

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **FedSGD**

➢ **Distributed Optimization에 partial participation과 데이터 불균형 반영**

*Goal* $\quad\Big|\quad w^* \triangleq \min_{w} f(w)$

**Mini-batch Gradient Descent**

**Server**

$w$

communication cost

**partial participation**

Non-I.I.D

**데이터 불균형**

| On | Off | On | | Off |
|---|---|---|---|---|
| **Client 1** | **Client 2** | **Client 3** | … | **Client K** |

$D_1$ $\quad\quad D_2$ $\quad\quad D_3$ $\quad\quad\quad D_K$

*Data size:* $n_k$

*Mini-batch size:* $m = \sum_{k \in S_t} n_k$

*Partial Loss:* $F_k(w) = \frac{1}{n_k} \sum_{i \in D_k} f_i(w)$

$$f_B(w) = \frac{1}{m} \sum_i f_i(w) = \frac{1}{m} \sum_{k \in S_t} n_k F_k(w) = \sum_{k \in S_t} \frac{n_k}{m} F_k(w)$$

Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **FedSGD**

➢ **Distributed Optimization에 partial participation과 데이터 불균형 반영**

*Goal* $\quad \Big| \quad w^* \triangleq \min_w f(w)$

**Mini-batch Gradient Descent**

$$w_{t+1} = w_t - \eta \nabla f_B(w_t) \quad \text{\textit{Mini-batch Loss Gradient}}$$

$$= w_t - \eta \nabla \Big\{ \sum_{k \in S_t} \frac{n_k}{m} F_k(w_t) \Big\}$$

$$= w_t - \sum_{k \in S_t} \frac{n_k}{m} \eta \nabla F_k(w_t)$$

*Partial Loss Gradient*

$$= \sum_{k \in S_t} \frac{n_k}{m} \cdot (w_t - \eta \nabla F_k(w_t))$$

*Local Update*

$$= \sum_{k \in S_t} \frac{n_k}{m} w_{t+1}^k$$

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **FedSGD**

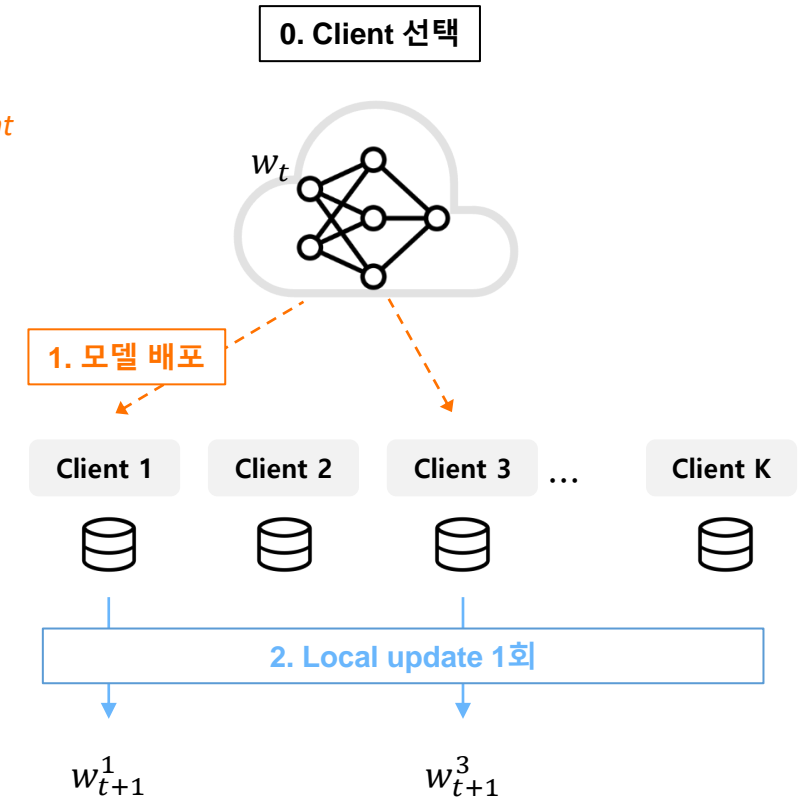➢ **Distributed Optimization에 partial participation과 데이터 불균형 반영**

*Goal* $\quad\Big|\quad w^* \triangleq \min_w f(w)$

**Mini-batch Gradient Descent**

$$w_{t+1} = w_t - \eta \nabla f_B(w_t) \quad \text{\textit{Mini-batch Loss Gradient}}$$

$$= w_t - \eta \nabla \{ \sum_{k \in S_t} \frac{n_k}{m} F_k(w_t) \}$$

$$= w_t - \sum_{k \in S_t} \frac{n_k}{m} \eta \nabla F_k(w_t)$$

*Partial Loss Gradient*

$$= \sum_{k \in S_t} \frac{n_k}{m} \cdot (w_t - \eta \nabla F_k(w_t))$$

*Local Update*

$$= \sum_{k \in S_t} \frac{n_k}{m} w_{t+1}^k$$

0. Client 선택

$w_t$

✔ Client 1  Client 2  ✔ Client 3  …  Client K

# Federated Averaging (FedAvg)

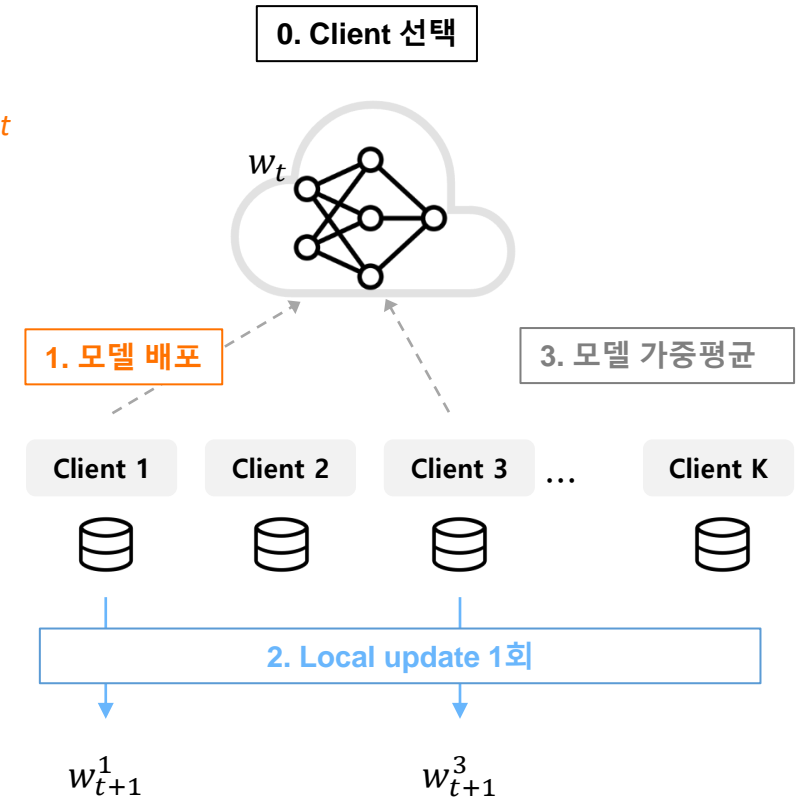**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **FedSGD**

➢ **Distributed Optimization에 partial participation과 데이터 불균형 반영**

*Goal*  $\quad w^* \triangleq \min_w f(w)$

**Mini-batch Gradient Descent**

$$w_{t+1} = w_t - \eta \nabla f_B(w_t) \quad \textit{Mini-batch Loss Gradient}$$

$$= w_t - \eta \nabla \{ \sum_{k \in S_t} \frac{n_k}{m} F_k(w_t) \}$$

$$= w_t - \sum_{k \in S_t} \frac{n_k}{m} \eta \nabla F_k(w_t)$$

*Partial Loss Gradient*

$$= \sum_{k \in S_t} \frac{n_k}{m} \cdot (w_t - \eta \nabla F_k(w_t))$$

*Local Update*

$$= \sum_{k \in S_t} \frac{n_k}{m} w_{t+1}^k$$

**0. Client 선택**

$w_t$

**1. 모델 배포**

| Client 1 | Client 2 | Client 3 | ... | Client K |

**2. Local update 1회**

$w_{t+1}^1 \qquad w_{t+1}^3$

# Federated Averaging (FedAvg)

Communication-Efficient Learning of Deep Networks from Decentralized Data

❖ **FedSGD**

➢ **Distributed Optimization에 partial participation과 데이터 불균형 반영**

*Goal* | $w^* \triangleq \min_w f(w)$

**Mini-batch Gradient Descent**

$$w_{t+1} = w_t - \eta \nabla f_B(w_t) \quad \text{\textit{Mini-batch Loss Gradient}}$$

$$= w_t - \eta \nabla \{ \sum_{k \in S_t} \frac{n_k}{m} F_k(w_t) \}$$

$$= w_t - \sum_{k \in S_t} \frac{n_k}{m} \eta \nabla F_k(w_t)$$

*Partial Loss Gradient*

$$= \sum_{k \in S_t} \frac{n_k}{m} \cdot (w_t - \eta \nabla F_k(w_t))$$

*Local Update*

$$= \sum_{k \in S_t} \frac{n_k}{m} w_{t+1}^k$$

**0. Client 선택**

$w_t$

**1. 모델 배포**

**3. 모델 가중평균**

| Client 1 | Client 2 | Client 3 | ... | Client K |

**2. Local update 1회**

$w_{t+1}^1$ $w_{t+1}^3$

Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Federated Averaging (FedAvg)**

➢ **FedSGD:** Local update 당 communication 1회

➢ **FedAvg:** Local update $E$ 회당 communication 1회

➢ 동일 local update 수 기준, communication 횟수 $1/E$로 감소

| communication cost |
| :---: |

| partial participation |
| :---: |

| Non-I.I.D |
| :---: |

| 데이터 불균형 |
| :---: |

0. Client 선택

$w_{t+1}$

1. 모델 배포

3. 모델 가중평균

Client 1    Client 2    Client 3    …    Client K

2. Local update E회

$w_{t+1}^1$          $w_{t+1}^3$

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Federated Averaging (FedAvg)**

➢ **FedSGD:** Local update 당 communication 1회

➢ **FedAvg:** Local update $E$ 회당 communication 1회

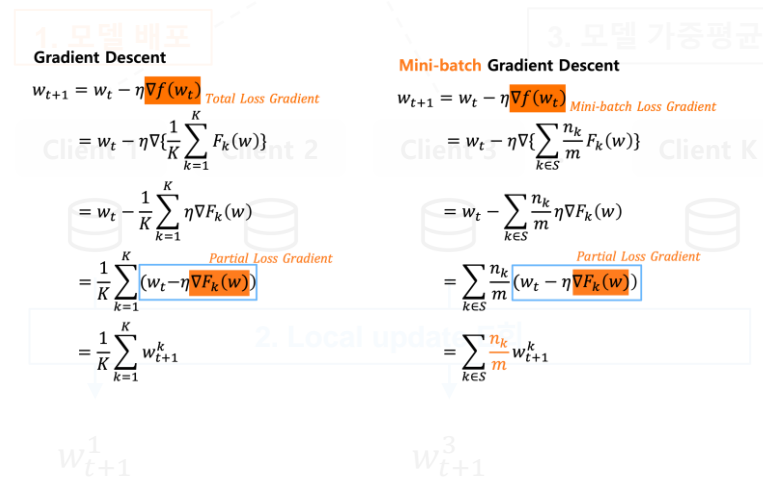➢ 동일 local update 수 기준, communication 횟수 **1/$E$**로 감소

**communication cost**

**partial participation**

Non-I.I.D

데이터 불균형

0. Client 선택

$w_{t+1}$

---

**Algorithm 1** FederatedAveraging. The $K$ clients are indexed by $k$; $B$ is the local minibatch size, $E$ is the number of local epochs, and $\eta$ is the learning rate.

---

**Server executes:**
  initialize $w_0$
  **for** each round $t = 1, 2, \ldots$ **do**
    $m \leftarrow \max(C \cdot K, 1)$
    $S_t \leftarrow$ (random set of $m$ clients)
    **for** each client $k \in S_t$ **in parallel do**
      $w_{t+1}^k \leftarrow$ ClientUpdate$(k, w_t)$
    $m_t \leftarrow \sum_{k \in S_t} n_k$
    $w_{t+1} \leftarrow \sum_{k \in S_t} \frac{n_k}{m_t} w_{t+1}^k$   // *Erratum*[4]

**ClientUpdate**$(k, w)$:   // *Run on client $k$*
  $\mathcal{B} \leftarrow$ (split $\mathcal{P}_k$ into batches of size $B$)
  **for** each local epoch $i$ from 1 to $E$ **do**
    **for** batch $b \in \mathcal{B}$ **do**
      $w \leftarrow w - \eta \nabla \ell(w; b)$
  return $w$ to server

---

Client 2    Client 3    Client K

- $C$: 매 round마다 참여할 클라이언트 비율 (Partial Participation)
- $E$: 매 round마다 학습하는 local epoch 수
- $B$: 매 local epoch마다 학습에 사용하는 local mini-batch 크기

2. Local update E회

$w_{t+1}^1$      $w_{t+1}^3$

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Federated Averaging (FedAvg)**

➢ **FedSGD:** Local update 당 communication 1회

➢ **FedAvg:** Local update $E$ 회당 communication 1회

➢ 동일 local update 수 기준, communication 횟수 $1/E$로 감소

0. Client 선택

**communication cost**

**partial participation**

Non-I.I.D

데이터 불균형

**1. Distributed Optimization이나**
**2. FedSGD와는 달리**
**'Local update 평균'이**
**'Global update'와 동일하다는 이론적 보장 X**

**Gradient Descent**

$$w_{t+1} = w_t - \eta \nabla f(w_t) \quad \textit{Total Loss Gradient}$$
$$= w_t - \eta \nabla \{ \frac{1}{K} \sum_{k=1}^{K} F_k(w) \}$$
$$= w_t - \frac{1}{K} \sum_{k=1}^{K} \eta \nabla F_k(w)$$
$$= \frac{1}{K} \sum_{k=1}^{K} (w_t - \eta \nabla F_k(w)) \quad \textit{Partial Loss Gradient}$$
$$= \frac{1}{K} \sum_{k=1}^{K} w_{t+1}^k$$

**Mini-batch Gradient Descent**

$$w_{t+1} = w_t - \eta \nabla f(w_t) \quad \textit{Mini-batch Loss Gradient}$$
$$= w_t - \eta \nabla \{ \sum_{k \in S} \frac{n_k}{m} F_k(w) \}$$
$$= w_t - \sum_{k \in S} \frac{n_k}{m} \eta \nabla F_k(w)$$
$$= \sum_{k \in S} \frac{n_k}{m} (w_t - \eta \nabla F_k(w)) \quad \textit{Partial Loss Gradient}$$
$$= \sum_{k \in S} \frac{n_k}{m} w_{t+1}^k$$

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Federated Averaging (FedAvg)**

➢ **FedSGD:** Local update 당 communication 1회

➢ **FedAvg:** Local update $E$ 회당 communication 1회

➢ 동일 local update 수 기준, communication 횟수 $1/E$ 로 감소

> communication cost
> partial participation
> Non-I.I.D
> 데이터 불균형

0. Client 선택

**1. Distributed Optimization**이나
**2. FedSGD**와는 달리
'**Local update 평균**'이
'**Global update**'와 동일하다는 **이론적 보장 X**

1. 모델 배포      3. 모델 가중평균

Client 1   Client 2   Client 3   …   Client K

**따라서, FedAvg에 대한 수렴성 증명 필요!**
**On the Convergence of FedAvg on Non-IID Data (ICLR'20)**

2. Local update E회

$w_{t+1}^1$      $w_{t+1}^3$

# Federated Averaging (FedAvg)

Communication-Efficient Learning of Deep Networks from Decentralized Data

❖ **Experiments**

➤ **하이퍼 파라미터 $C$, $E$, $B$를 어떻게 설정해야 하는가?**

- $C$: 매 round마다 참여할 클라이언트 비율 (Partial Participation)

- $E$: 매 round마다 학습하는 local epoch 수

- $B$: 매 local epoch마다 학습에 사용하는 local mini-batch 크기

**Test accuracy 97% 도달 round**

**Test accuracy 99% 도달 round**

| 2NN | —— IID —— | | —— NON-IID —— | |
|---|---|---|---|---|
| $C$ | $B = \infty$ | $B = 10$ | $B = \infty$ | $B = 10$ |
| 0.0 | 1455 | 316 | 4278 | 3275 |
| 0.1 | 1474 (1.0×) | 87 (3.6×) | 1796 (2.4×) | 664 (4.9×) |
| 0.2 | 1658 (0.9×) | 77 (4.1×) | 1528 (2.8×) | 619 (5.3×) |
| 0.5 | — (—) | 75 (4.2×) | — (—) | 443 (7.4×) |
| 1.0 | — (—) | 70 (4.5×) | — (—) | 380 (8.6×) |
| CNN, $E = 5$ | | | | |
| 0.0 | 387 | 50 | 1181 | 956 |
| 0.1 | 339 (1.1×) | 18 (2.8×) | 1100 (1.1×) | 206 (4.6×) |
| 0.2 | 337 (1.1×) | 18 (2.8×) | 978 (1.2×) | 200 (4.8×) |
| 0.5 | 164 (2.4×) | 18 (2.8×) | 1067 (1.1×) | 261 (3.7×) |
| 1.0 | 246 (1.6×) | 16 (3.1×) | — (—) | 97 (9.9×) |

**클라이언트 당 label 2개씩만 할당 (Label Shift)**

**MNIST 데이터셋 분류 실험**

Data Mining
Quality Analytics

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Experiments**

➢ **하이퍼 파라미터 $C, E, B$를 어떻게 설정해야 하는가?**

- $C$: 매 round마다 참여할 클라이언트 비율 (Partial Participation)
- $E$: 매 round마다 학습하는 local epoch 수
- $B$: 매 local epoch마다 학습에 사용하는 local mini-batch 크기

$$Local\ Computation \propto \frac{E}{B}$$



**MNIST 데이터셋 분류 실험**

**Shakespeare Next-token prediction 실험**

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Experiments**

➢ **하이퍼 파라미터 $C, E, B$를 어떻게 설정해야 하는가?**

- $C$: 매 round마다 참여할 클라이언트 비율 (Partial Participation)

- $E$: 매 round마다 학습하는 local epoch 수

- $B$: 매 local epoch마다 학습에 사용하는 local mini-batch 크기

$$Local\ Computation \propto \frac{E}{B}$$

**Local computation**이 증가($E$ ↓, $B$ ↑)할수록, 빠르게 수렴

MNIST 데이터셋 분류 실험                    Shakespeare Next-token prediction 실험

# Federated Averaging (FedAvg)

**Communication-Efficient Learning of Deep Networks from Decentralized Data**

❖ **Discussion**

➢ **의의: Federated Learning과 이를 위한 FedAvg 알고리즘 제안**

➢ **한계점:**

- 개선의 여지가 많은 알고리즘 (e.g. Non-IID 조건에 대한 고려 X) ➡ **Federated Optimization in Heterogeneous Networks (MLSys'20)**

- 실험적으로 잘 작동함을 보였으나, 엄밀한 설명 부족 ➡ **On the Convergence of FedAvg on Non-IID Data (ICLR'20)**

# FedProx

**Federated Optimization in Heterogeneous Networks**

❖ **FedAvg 한계점 개선한 알고리즘 제안 및 convergence analysis 제공**

➢ MLSys'20

➢ 피인용 6028회 (2025년 3월 기준)

Li, T., Sahu, A. K., Zaheer, M., Sanjabi, M., Talwalkar, A., & Smith, V. (2020). Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems(MLSys)*, 2, 429–450.

# FedProx

**Federated Optimization in Heterogeneous Networks**

❖ **FedAvg 한계점**

➢ **클라이언트 간 이질성에 취약**

- **통계적 이질성(Non-IID)**: 데이터 분포 차이가 클 경우, 수렴하지 못하고 발산 (Client Drift)
- **시스템적 이질성**: 학습 속도가 느린 디바이스 존재 (Straggler) → 동일 Epoch 강제는 비효율적

| [ 통계적 이질성 ] | [ 시스템적 이질성 ] |
| --- | --- |



**Local 학습 시 Global model에서
지나치게 멀어지는 것을 방지**

**디바이스 사양에 따라 Epoch 수 조절**

[1] https://ar5iv.labs.arxiv.org/html/2103.00710
[2] https://www.researchgate.net/figure/Stragglers-impact-on-FL-performance-In-synchronous-FL-all-clients-wait-for-the_fig1_372163244

Data Mining
Quality Analytics

# FedProx

❖ **Proximal Term & $\gamma_t^k$-inexact Solution**

➢ **두 이질성에 대한 보완책**

- **Proximal Term**: 로컬 모델과 글로벌 모델 간 차이를 억제 → 통계적 이질성 ↓

- $\gamma_t^k$**-inexact Solution** : 학습 속도가 느린 디바이스는 다소 부정확한 모델이라도 반환할 수 있도록 허락 → 시스템적 이질성↓

**[ Proximal Term ]**                                                  **[ $\gamma_t^k$-inexact Solution ]**

$$\min_{w} h_k(w; w_t) = F_k(w) + \frac{\mu}{2} \|w - w_t\|^2$$

**Proximal Term**

Local 학습 시 Global model에서
지나치게 멀어지는 것을 방지                        디바이스 사양에 따라 Epoch 수 조절

Data Mining
Quality Analytics

56

# FedProx

❖ **Proximal Term & $\gamma_t^k$-inexact Solution**

➢ **두 이질성에 대한 보완책**

- **Proximal Term**: 로컬 모델과 글로벌 모델 간 차이를 억제 → 통계적 이질성 ↓
- **$\gamma_t^k$-inexact Solution** : 학습 속도가 느린 디바이스는 다소 부정확한 모델이라도 반환할 수 있도록 허락 → 시스템적 이질성↓

**[ Proximal Term ]**　　　　　　　　　　　　　　　　**[ $\gamma_t^k$-inexact Solution ]**

$$\min_w h_k(w; w_t) = F_k(w) + \frac{\mu}{2}\|w - w_t\|^2$$

**Proximal Term**

$$\nabla h_k(w; w_t) = \nabla F_k(w) + \boxed{\mu}(w - w_t)$$

**페널티 파라미터**

**Local 학습 시 Global model에서
지나치게 멀어지는 것을 방지**

**디바이스 사양에 따라 Epoch 수 조절**

# FedProx

❖ **Proximal Term & $\gamma_t^k$-inexact Solution**

➤ **두 이질성에 대한 보완책**

- **Proximal Term**: 로컬 모델과 글로벌 모델 간 차이를 억제 → 통계적 이질성 ↓

- $\gamma_t^k$**-inexact Solution** : 학습 속도가 느린 디바이스는 다소 부정확한 모델이라도 반환할 수 있도록 허락 → 시스템적 이질성↓

**[ Proximal Term ]**

$$\min_{w} h_k(w; w_t) = F_k(w) + \frac{\mu}{2}\|w - w_t\|^2$$

**Proximal Term**

$$\nabla h_k(w; w_t) = \nabla F_k(w) + \boxed{\mu}(w - w_t)$$

**페널티 파라미터**

**Local 학습 시 Global model에서
지나치게 멀어지는 것을 방지**

**[$\gamma_t^k$-inexact Solution ]**

$$\|\nabla h_k(w^*; w_t)\| \leq \gamma_t^k \|\nabla h_k(w_t; w_t)\|$$

$$\|\nabla F_k(w^*) + \mu(w^* - w_t)\| \leq \gamma_t^k \|\nabla F_k(w_t)\|$$

**디바이스 사양에 따라 Epoch 수 조절**

# FedProx

❖ **Proximal Term & $\gamma_t^k$-inexact Solution**

➢ **두 이질성에 대한 보완책**

- **Proximal Term**: 로컬 모델과 글로벌 모델 간 차이를 억제 → 통계적 이질성 ↓

- **$\gamma_t^k$-inexact Solution** : 학습 속도가 느린 디바이스는 다소 부정확한 모델이라도 반환할 수 있도록 허락 → 시스템적 이질성↓

**[ Proximal Term ]**

$$\min_w h_k(w; w_t) = F_k(w) + \frac{\mu}{2}\|w - w_t\|^2$$

**Proximal Term**

$$\nabla h_k(w; w_t) = \nabla F_k(w) + \boxed{\mu}(w - w_t)$$

**페널티 파라미터**

**Local 학습 시 Global model에서
지나치게 멀어지는 것을 방지**

**[ $\gamma_t^k$-inexact Solution ]**

$$\|\nabla h_k(w^*; w_t)\| \leq \gamma_t^k \|\nabla h_k(w_t; w_t)\|$$

$$\|\nabla F_k(w^*) + \mu(w^* - w_t)\| \leq \gamma_t^k \|\nabla F_k(w_t)\|$$

**부정확한 로컬 모델이 글로벌
모델에 악영향을 주지는 않을까?**

**디바이스 사양에 따라 Epoch 수 조절**

# FedProx

Federated Optimization in Heterogeneous Networks

❖ **Proximal Term & $\gamma_t^k$-inexact Solution**

➤ **두 이질성에 대한 보완책**

- **Proximal Term**: 로컬 모델과 글로벌 모델 간 차이를 억제 → 통계적 이질성 ↓
- $\gamma_t^k$**-inexact Solution** : 학습 속도가 느린 디바이스는 다소 부정확한 모델이라도 반환할 수 있도록 허락 → 시스템적 이질성↓

**[ Proximal Term ]**

**[ $\gamma_t^k$-inexact Solution ]**

$$\min_{w} h_k(w; w_t) = F_k(w) + \frac{\mu}{2}\|w - w_t\|^2$$

**Proximal Term**

$$\nabla h_k(w; w_t) = \nabla F_k(w) + \mu(w - w_t)$$

**페널티 파라미터**

$$\|\nabla h_k(w^*; w_t)\| \leq \gamma_t^k \|\nabla h_k(w_t; w_t)\|$$

$$\|\nabla F_k(w^*) + \mu(w^* - w_t)\| \leq \gamma_t^k \|\nabla F_k(w_t)\|$$

**부정확한 로컬 모델이 글로벌 모델에 악영향을 주지는 않을까?**

**Local 학습 시 Global model에서 지나치게 멀어지는 것을 방지**

**디바이스 사양에 따라 Epoch 수 조절**

# FedProx

Federated Optimization in Heterogeneous Networks

❖ **FedProx 알고리즘**

➢ **FedAvg 알고리즘에서** $E$ **(local train epoch 수) 대체** → $\mu$, $\gamma$

---

**Algorithm 1** Federated Averaging (FedAvg)

**Input:** $K, T, \eta, E, w^0, N, p_k, k = 1, \cdots, N$
**for** $t = 0, \cdots, T-1$ **do**
    Server selects a subset $S_t$ of $K$ devices at random (each device $k$ is chosen with probability $p_k$)
    Server sends $w^t$ to all chosen devices
    Each device $k \in S_t$ updates $w^t$ for $E$ epochs of SGD on $F_k$ with step-size $\eta$ to obtain $w_k^{t+1}$
    Each device $k \in S_t$ sends $w_k^{t+1}$ back to the server
    Server aggregates the $w$'s as $w^{t+1} = \frac{1}{K}\sum_{k \in S_t} w_k^{t+1}$
**end for**

---

**Algorithm 2** FedProx (Proposed Framework)

**Input:** $K, T, \mu, \gamma, w^0, N, p_k, k = 1, \cdots, N$
**for** $t = 0, \cdots, T-1$ **do**
    Server selects a subset $S_t$ of $K$ devices at random (each device $k$ is chosen with probability $p_k$)
    Server sends $w^t$ to all chosen devices
    Each chosen device $k \in S_t$ finds a $w_k^{t+1}$ which is a $\gamma_k^t$-inexact minimizer of: $w_k^{t+1} \approx \arg\min_w h_k(w; w^t) = F_k(w) + \frac{\mu}{2}\|w - w^t\|^2$
    Each device $k \in S_t$ sends $w_k^{t+1}$ back to the server
    Server aggregates the $w$'s as $w^{t+1} = \frac{1}{K}\sum_{k \in S_t} w_k^{t+1}$
**end for**

# FedProx

❖ **Convergence Analysis**

➢ **FedProx 사용 시, total Loss $f(w)$가 최솟값 $f^*$으로 수렴함을 증명**

*Definition.* B-local dissimilarity

$$\mathbb{E}_k \left[ \|\nabla F_k(w)\|^2 \right] = \|\nabla f(w)\|^2 B(w)^2$$

Data Mining
Quality Analytics

# FedProx

❖ **Convergence Analysis**

➢ **FedProx 사용 시, total Loss $f(w)$가 최솟값 $f^*$으로 수렴함을 증명**

*Definition.* B-local dissimilarity

$$\mathbb{E}_k \left[ \|\nabla F_k(w)\|^2 \right] = \|\nabla f(w)\|^2 B(w)^2$$

$$\mathbb{E}_k \left[ \|\nabla F_k(w)\|^2 \right] = \|\mathbb{E}_k \left[ \nabla F_k(w) \right] \|^2 B(w)^2$$

분산과 유사!

$$\nabla F_1(w) = 1 \qquad \nabla F_2(w) = 1 \qquad \nabla F_3(w) = 1$$

$$1 = B(w)^2$$

$$\nabla F_1(w) = -3 \qquad \nabla F_2(w) = 3 \qquad \nabla F_3(w) = 6$$

$$1.5 = B(w)^2$$

# FedProx

❖ **Convergence Analysis**

➤ **FedProx 사용 시, total Loss $f(w)$가 최솟값 $f^*$으로 수렴함을 증명**

*Definition.* B-local dissimilarity

$$\mathbb{E}_k \left[ \|\nabla F_k(w)\|^2 \right] = \|\nabla f(w)\|^2 B(w)^2$$

*Assumption.* Bounded dissimilarity

$$B(w) \leq B$$

# FedProx

Federated Optimization in Heterogeneous Networks

❖ **Convergence Analysis**

  ➢ **FedProx 사용 시, total Loss $f(w)$가 최솟값 $f^*$으로 수렴함을 증명**

*Definition.* B-local dissimilarity

$$\mathbb{E}_k \left[ \|\nabla F_k(w)\|^2 \right] = \|\nabla f(w)\|^2 B(w)^2$$

*Assumption.* Bounded dissimilarity

$$B(w) \leq B$$

*Theorem 1.* Non-convex FedProx convergence

$$\mathbb{E}_{S_t} \left[ f(w_{t+1}) \right] \leq f(w_t) - \rho \|\nabla f(w_t)\|^2$$

*Theorem 2. Convergence rate*

$B$, $\mu$, $\gamma$에 대한 함수
$\rho > 0$일 때, 수렴

$$f(w_0) - f^* =: \Delta$$

$$T := O \left( \frac{\Delta}{\rho \epsilon} \right)$$

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \|\nabla f(w_t)\|^2 \right] \leq \epsilon$$

Data Mining
Quality Analytics

# FedProx

Federated Optimization in Heterogeneous Networks

❖ **Convergence Analysis**

➢ **FedProx 사용 시, total Loss $f(w)$가 최솟값 $f^*$으로 수렴함을 증명**

*Definition.* B-local dissimilarity

$$\mathbb{E}_k \left[ \|\nabla F_k(w)\|^2 \right] = \|\nabla f(w)\|^2 B(w)^2$$

*Assumption.* Bounded dissimilarity

$$B(w) \leq B$$

*Theorem 1.* Non-convex FedProx convergence

$$\mathbb{E}_{S_t} \left[ f(w_{t+1}) \right] \leq f(w_t) - \rho \|\nabla f(w_t)\|^2$$

*Theorem 2.* Convergence rate

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \|\nabla f(w_t)\|^2 \right] \leq \epsilon$$

$$f(w_0) - f^* =: \Delta$$

$$T := O \left( \frac{\Delta}{\rho \epsilon} \right)$$
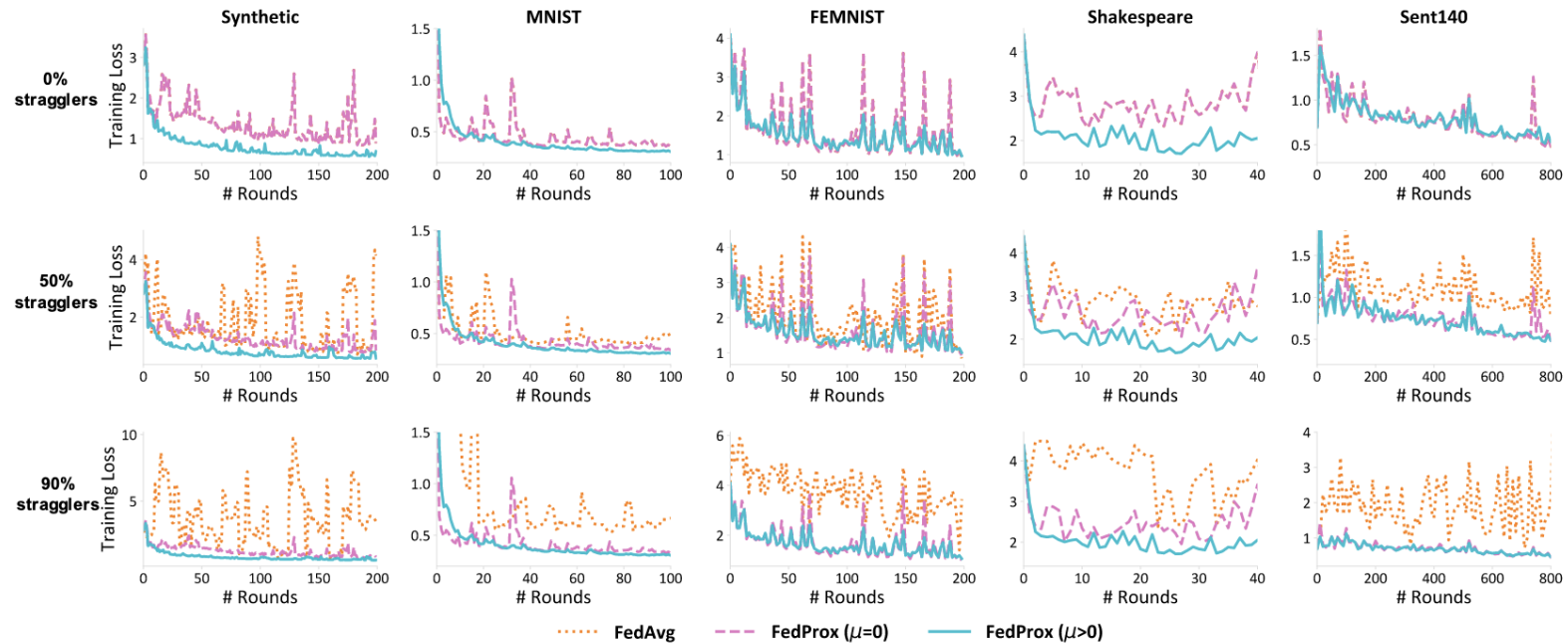
Data Mining
Quality Analytics

# FedProx

Federated Optimization in Heterogeneous Networks

❖ **Experiments**

  ➢ **시스템적 이질성**

   • 매 round마다 0%, 50%, 혹은 90%의 straggler가 발생
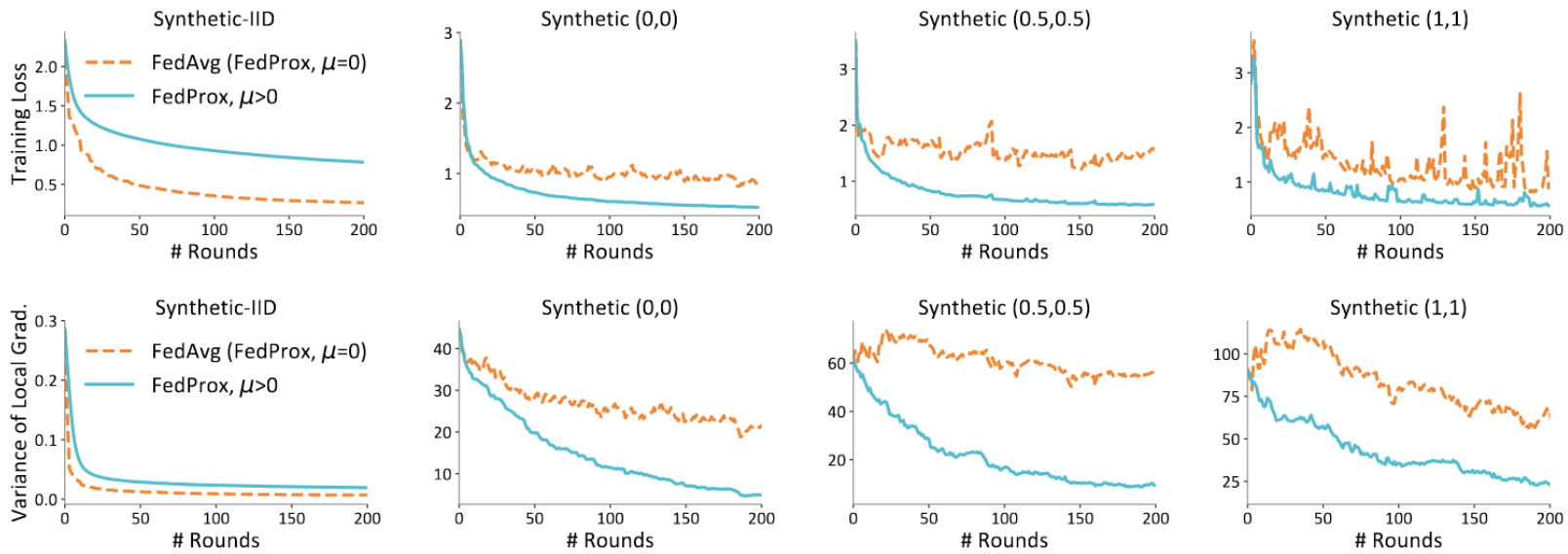
   • FedAvg는 Straggler를 FL에서 배제하도록 설정

# FedProx

Federated Optimization in Heterogeneous Networks

❖ **Experiments**

➢ **통계적 이질성**

- Straggler는 존재하지 않는다고 가정
- 클라이언트 간 통계적 이질성을 달리 하며 실험 진행

# Convergence of FedAvg

On the Convergence of FedAvg on Non-IID Data

❖ **FedAvg에 대한 convergence analysis 제공**

  ➢ ICLR'20

  ➢ 피인용 2832회 (2025년 3월 기준)

# Convergence of FedAvg

**On the Convergence of FedAvg on Non-IID Data**

❖ **증명 과정**

➢ Full participation에 대해 증명 후, partial participation으로 확장

➢ Partial participation의 sampling과 averaging 전략은 아래와 같음

**2023년 개정 전 알고리즘**

*Client selection*

**FedAvg 제안 논문**

| Paper | Sampling | Averaging | Convergence rate |
|---|---|---|---|
| McMahan et al. (2017) | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \notin \mathcal{S}_t} p_k \mathbf{w}_t + \sum_{k \in \mathcal{S}_t} p_k \mathbf{w}_t^k$ | - |
| Sahu et al. (2018) | $\mathcal{S}_t \sim \mathcal{W}(N, K, \mathbf{p})$ | $\frac{1}{K} \sum_{k \in \mathcal{S}_t} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})$ [5] |
| Ours | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \in \mathcal{S}_t} p_k \frac{N}{K} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})$ [6] |

*Sampling and Averaging Schemes*

Data Mining
Quality Analytics

# Convergence of FedAvg

❖ **증명 과정**

➢ Full participation에 대해 증명 후, partial participation으로 확장

➢ Partial participation의 sampling과 averaging 전략은 아래와 같음

|  | **개정 전** | **개정 후** |
|---|---|---|
| **FedAvg 제안 논문** | $w_{t+1} = \sum_{k \notin S_t} \dfrac{n_k}{n} w_t + \sum_{k \in S_t} \dfrac{n_k}{n} w_{t+1}^k$ | $w_{t+1} = \sum_{k \in S_t} \dfrac{n_k}{m_t} w_{t+1}^k$ |

# Convergence of FedAvg

**On the Convergence of FedAvg on Non-IID Data**

❖ **증명 과정**

➢ Full participation에 대해 증명 후, partial participation으로 확장

➢ Partial participation의 sampling과 averaging 전략은 아래와 같음

**2023년 개정 전 알고리즘**

*Client selection*

| Paper | Sampling | Averaging | Convergence rate |
|---|---|---|---|
| McMahan et al. (2017) | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \notin \mathcal{S}_t} p_k \mathbf{w}_t + \sum_{k \in \mathcal{S}_t} p_k \mathbf{w}_t^k$ | - |
| Sahu et al. (2018) | $\mathcal{S}_t \sim \mathcal{W}(N, K, \mathbf{p})$ | $\frac{1}{K} \sum_{k \in \mathcal{S}_t} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})^5$ |
| Ours | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \in \mathcal{S}_t} p_k \frac{N}{K} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})^6$ |

**FedAvg 제안 논문** (McMahan et al. (2017), Sahu et al. (2018))

**수렴성 보장 실패**

*Sampling and Averaging Schemes*

# Convergence of FedAvg

On the Convergence of FedAvg on Non-IID Data

❖ 증명 과정

➢ Full participation에 대해 증명 후, partial participation으로 확장

➢ Partial participation의 sampling과 averaging 전략은 아래와 같음

_Client selection_

| Paper | Sampling | Averaging | Convergence rate |
|---|---|---|---|
| McMahan et al. (2017) | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \notin \mathcal{S}_t} p_k \mathbf{w}_t + \sum_{k \in \mathcal{S}_t} p_k \mathbf{w}_t^k$ | - |
| Sahu et al. (2018) | $\mathcal{S}_t \sim \mathcal{W}(N, K, \mathbf{p})$ | $\frac{1}{K} \sum_{k \in \mathcal{S}_t} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})^5$ |
| Ours | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \in \mathcal{S}_t} p_k \frac{N}{K} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})^6$ |

**FedProx 제안 논문**

_Sampling and Averaging Schemes_

# Convergence of FedAvg

On the Convergence of FedAvg on Non-IID Data

❖ **증명 과정**

➢ Full participation에 대해 증명 후, partial participation으로 확장

➢ Partial participation의 sampling과 averaging 전략은 아래와 같음

*Client selection*

**수렴성 보장을 위한 변형**

| Paper | Sampling | Averaging | Convergence rate |
|---|---|---|---|
| McMahan et al. (2017) | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \notin \mathcal{S}_t} p_k \mathbf{w}_t + \sum_{k \in \mathcal{S}_t} p_k \mathbf{w}_t^k$ | - |
| Sahu et al. (2018) | $\mathcal{S}_t \sim \mathcal{W}(N, K, \mathbf{p})$ | $\frac{1}{K} \sum_{k \in \mathcal{S}_t} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})$[5] |
| Ours | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \in \mathcal{S}_t} p_k \frac{N}{K} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})$[6] |

*Sampling and Averaging Schemes*

$$\sum_{k \in \mathcal{S}_t} \frac{n_k}{n} \neq 1$$

# Convergence of FedAvg

**On the Convergence of FedAvg on Non-IID Data**

❖ **증명 과정**

➢ Full participation에 대해 증명 후, partial participation으로 확장

➢ Partial participation의 sampling과 averaging 전략은 아래와 같음

*Client selection*

**수렴성 보장을 위한 변형**

| Paper | Sampling | Averaging | Convergence rate |
|---|---|---|---|
| McMahan et al. (2017) | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \notin \mathcal{S}_t} p_k \mathbf{w}_t + \sum_{k \in \mathcal{S}_t} p_k \mathbf{w}_t^k$ | - |
| Sahu et al. (2018) | $\mathcal{S}_t \sim \mathcal{W}(N, K, \mathbf{p})$ | $\frac{1}{K} \sum_{k \in \mathcal{S}_t} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})$[5] |
| Ours | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \in \mathcal{S}_t} p_k \frac{N}{K} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})$[6] |

*Sampling and Averaging Schemes*

$$\mathbb{E}\left[\frac{N}{K} \sum_{k \in \mathcal{S}_t} \frac{n_k}{n}\right] = 1$$

$N$: 전체 클라이언트 수

$K$: 한 라운드에 클라이언트 수 제한

# Convergence of FedAvg

**On the Convergence of FedAvg on Non-IID Data**

❖ **가설**

➢ FedProx와는 달리, convexity에 대한 가정 존재

➢ DNN 적용에 한계 (Non-convex)

## Assumptions

**Assumption 1:** $L$-smoothness

$$F_k(v) - F_k(w) \leq (v-w)^T \nabla F_k(w) + \frac{L}{2}\|v-w\|_2^2.$$

**Assumption 2:** $\mu$-strong **convexity**

$$F_k(v) - F_k(w) \geq (v-w)^T \nabla F_k(w) + \frac{\mu}{2}\|v-w\|_2^2.$$

**Assumption 3: Bounded Variance**

$$\mathbb{E}[\|\nabla F_k(w_k^t, \xi_k^t) - \nabla F_k(w_k^t)\|^2] \leq \sigma_k^2.$$

**Assumption 4: Uniformly Bounded Squared Expectation**

$$\mathbb{E}[\|\nabla F_k(w_k^t, \xi_k^t)\|^2] \leq G^2$$

# Convergence of FedAvg

On the Convergence of FedAvg on Non-IID Data

❖ **Full Participation**

  ➢ 최종 시점 모델에 대한 손실 값이 이론적 최솟값에 수렴함을 증명

  ➢ $T$: 전체 local update 수

  ➢ $E$: Communication 1회 당 local update 수

$$\mathbb{E}\left[F(w_T)\right] - F^* \leq \frac{\kappa}{\gamma + T - 1}) \left( \frac{2B}{\mu} + \frac{\mu\gamma}{2}\mathbb{E}\|w_1 - w^*\|^2 \right),$$

$$B = \sum_{k=1}^{N} p_k^2 \sigma_k^2 + 6L\Gamma + 8(E-1)^2 G^2$$

# Convergence of FedAvg

On the Convergence of FedAvg on Non-IID Data

❖ **Full Participation**

➢ 최종 시점 모델에 대한 손실 값이 이론적 최솟값에 수렴함을 증명

➢ $T$: 전체 local update 수

➢ $E$: Communication 1회 당 local update 수

*T가 증가함에 따라 0으로 수렴*

이론적 최솟값

$$\mathbb{E}\left[F(w_T)\right] - F^* \leq \boxed{\frac{\kappa}{\gamma + T - 1}) \left(\frac{2B}{\mu} + \frac{\mu\gamma}{2}\mathbb{E}\|w_1 - w^*\|^2\right)},$$

T 시점 모델 손실 값

$$B = \sum_{k=1}^{N} p_k^2 \sigma_k^2 + 6L\Gamma + 8(E-1)^2 G^2$$

# Convergence of FedAvg

On the Convergence of FedAvg on Non-IID Data

❖ **Partial Participation**

➢ 두 가지 전략에 대해서 증명

➢ Scheme 2에 대해서는, **데이터 균형**을 가정 (비현실적 가정)

| Paper | Sampling | Averaging | Convergence rate |
|---|---|---|---|
| McMahan et al. (2017) | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \notin \mathcal{S}_t} p_k \mathbf{w}_t + \sum_{k \in \mathcal{S}_t} p_k \mathbf{w}_t^k$ | - |
| Sahu et al. (2018) | $\mathcal{S}_t \sim \mathcal{W}(N, K, \mathbf{p})$ | $\frac{1}{K} \sum_{k \in \mathcal{S}_t} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})^5$ |
| Ours | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \in \mathcal{S}_t} p_k \frac{N}{K} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})^6$ |

*Scheme 1*
*Scheme 2*

$$\mathbb{E}\left[F(w_T)\right] - F^* \leq \frac{\kappa}{\gamma + T - 1}\left(\frac{2(B + \boxed{C})}{\mu} + \frac{\mu\gamma}{2}\mathbb{E}\|w_1 - w^*\|^2\right),$$

*Scheme 1*

*Scheme 2 + Balanced Data Assumption*

$$C = \frac{4}{K}E^2G^2$$

$$C = \frac{N - K}{N - 1}\frac{4}{K}E^2G^2$$

Data Mining
Quality Analytics

# Convergence of FedAvg

On the Convergence of FedAvg on Non-IID Data

❖ **Partial Participation**

➢ 두 가지 전략에 대해서 증명

➢ Scheme 2에 대해서는, **데이터 균형**을 가정 (비현실적 가정)

|  | Paper | Sampling | Averaging | Convergence rate |
|---|---|---|---|---|
| *Scheme 1* | McMahan et al. (2017) | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \notin \mathcal{S}_t} p_k \mathbf{w}_t + \sum_{k \in \mathcal{S}_t} p_k \mathbf{w}_t^k$ | - |
| *Scheme 2* | Sahu et al. (2018) | $\mathcal{S}_t \sim \mathcal{W}(N, K, \mathbf{p})$ | $\frac{1}{K} \sum_{k \in \mathcal{S}_t} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})$ [5] |
|  | Ours | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \in \mathcal{S}_t} p_k \frac{N}{K} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})$ [6] |

$$\mathbb{E}\left[F(w_T)\right] - F^* \le \frac{\kappa}{\gamma + T - 1}) \left( \frac{2(B + \boxed{C})}{\mu} + \frac{\mu\gamma}{2} \mathbb{E}\|w_1 - w^*\|^2 \right),$$

**_Scheme 1_ – N개의 디바이스 중 K개를 임의로 Sampling 할 수 있어야 함**

**샘플링 된 디바이스 중 straggler 가 존재한다면 비효율적**

Data Mining
Quality Analytics

# Convergence of FedAvg

On the Convergence of FedAvg on Non-IID Data

❖ **Partial Participation**

➢ 두 가지 전략에 대해서 증명

➢ Scheme 2에 대해서는, **데이터 균형**을 가정 (비현실적 가정)

*Scheme 1*
*Scheme 2*

| Paper | Sampling | Averaging | Convergence rate |
|---|---|---|---|
| McMahan et al. (2017) | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \notin \mathcal{S}_t} p_k \mathbf{w}_t + \sum_{k \in \mathcal{S}_t} p_k \mathbf{w}_t^k$ | - |
| Sahu et al. (2018) | $\mathcal{S}_t \sim \mathcal{W}(N, K, \mathbf{p})$ | $\frac{1}{K} \sum_{k \in \mathcal{S}_t} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})$[5] |
| Ours | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \in \mathcal{S}_t} p_k \frac{N}{K} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})$[6] |

$$\mathbb{E}\left[F(w_T)\right] - F^* \leq \frac{\kappa}{\gamma + T - 1}) \left(\frac{2(B + \boxed{C})}{\mu} + \frac{\mu\gamma}{2} \mathbb{E}\|w_1 - w^*\|^2\right),$$

*Scheme 2* – N개의 디바이스 중 먼저 학습이 끝난 K개를 선택

**Straggler**가 존재한다면 *Scheme 1*의 현실적 대안이 될 수 있음

Data Mining
Quality Analytics

# Convergence of FedAvg

On the Convergence of FedAvg on Non-IID Data

❖ **Partial Participation**

➢ 두 가지 전략에 대해서 증명

➢ Scheme 2에 대해서는, **데이터 균형**을 가정 (비현실적 가정)

Scheme 1
Scheme 2

| Paper | Sampling | Averaging | Convergence rate |
|---|---|---|---|
| McMahan et al. (2017) | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \notin \mathcal{S}_t} p_k \mathbf{w}_t + \sum_{k \in \mathcal{S}_t} p_k \mathbf{w}_t^k$ | - |
| Sahu et al. (2018) | $\mathcal{S}_t \sim \mathcal{W}(N, K, \mathbf{p})$ | $\frac{1}{K} \sum_{k \in \mathcal{S}_t} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})$ [5] |
| Ours | $\mathcal{S}_t \sim \mathcal{U}(N, K)$ | $\sum_{k \in \mathcal{S}_t} p_k \frac{N}{K} \mathbf{w}_t^k$ | $\mathcal{O}(\frac{1}{T})$ [6] |

$$\mathbb{E}\left[F(w_T)\right] - F^* \leq \frac{\kappa}{\gamma + T - 1}) \left( \frac{2(B + \boxed{C})}{\mu} + \frac{\mu \gamma}{2} \mathbb{E}\|w_1 - w^*\|^2 \right),$$

*Scheme T − Balanced data assumption 완화*

$$\tilde{F}_k(w) := \underline{p_k N F_k(w)} \qquad F(w) = \sum_{k=1}^{N} p_k F_k(w) = \boxed{\frac{1}{N}} \sum_{k=1}^{N} \tilde{F}_k(w)$$

*Scaling*

$$\nu := N \cdot \max_k p_k \qquad \varsigma := N \cdot \min_k p_k$$

$$\tilde{L} := \nu L$$
$$\tilde{\mu} := \varsigma \mu$$
$$\tilde{\sigma}_k := \sqrt{\nu} \sigma_k$$
$$\tilde{G} := \sqrt{\varsigma} G$$

Data Mining
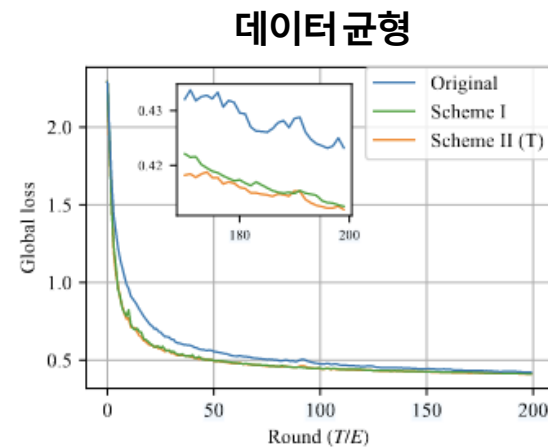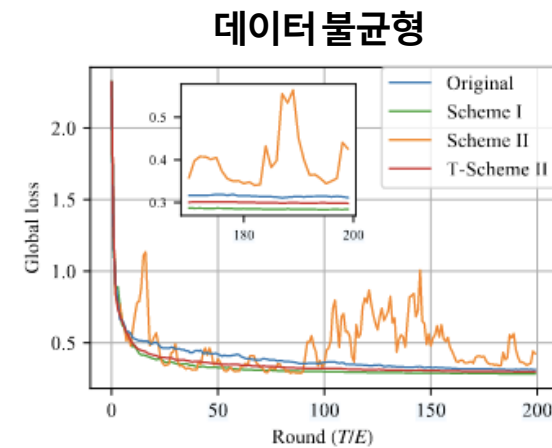Quality Analytics

# Convergence of FedAvg

**On the Convergence of FedAvg on Non-IID Data**

❖ **Experiments (1)**

➢ 데이터 균형: *Scheme* 관계 없이 안정적으로 수렴

➢ 데이터 불균형: *Scheme 1*과 *Scheme T*가 안정적으로 수렴



데이터 균형          데이터 불균형

(c) Different schemes      (d) Different schemes

# Convergence of FedAvg
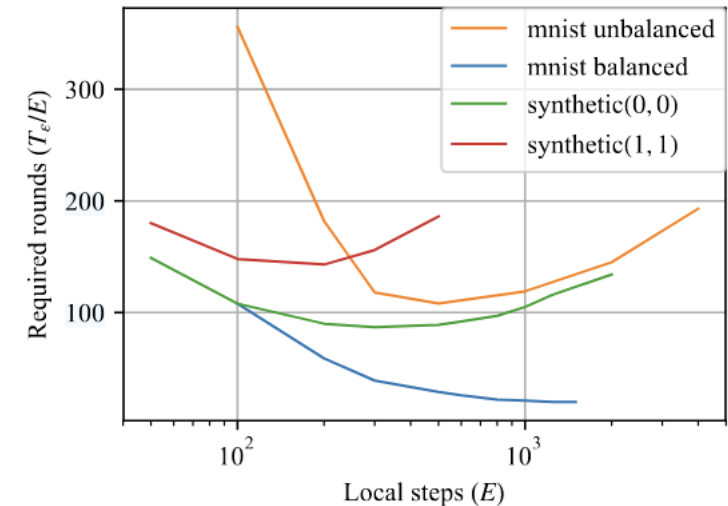
❖ **Hyperparameter $E$의 영향**

➢ $T_\epsilon$: $\epsilon$만큼의 정확도를 달성하는데 필요한 local update 수

➢ $\dfrac{T_\epsilon}{E}$: $\epsilon$만큼의 정확도를 달성하는데 필요한 **communication round** 수

$$\frac{T_\epsilon}{E} \propto \left(1 + \frac{1}{K}\right) EG^2 + \frac{\sum_{k=1}^{N} p_k^2 \sigma_k^2 + L\Gamma + \kappa G^2}{E} + G^2$$

**Communication round**를 결정하는 것은 $E$

$E \downarrow$ : *Local model* *underfitting*

$E \uparrow$ : *Local model* *overfitting*



(a) The impact of $E$