FedNLP: Federated Learning for Natural Language Processing

Xiang Ren, USC CS Mahdi Soltanolkotabi, USC ECE



amazon

Prof. Xiang Ren (PI)

- Expertise: natural language processing, explainable AI, continual learning
- Make NLP systems trustworthy, cheaper to build, easier to maintain
- Create benchmark datasets for a range of NLP tasks



Commonsense reasoning



Human-in-the-loop learning, Explanability & Robustness



Cross-task Generalization

Prof. Mahdi Soltanolkotabi



- Expertise: Foundations of AI & ML, (non)convex optimization, high-dimensional statistics and probability
- Developed some of the first guarantees for deep learning





Generalization (how neural networks predict)

• Extensive contributions to distributed/federated learning

Bill Yuchen Lin, PhD candidate @ USC



- (Bill) Yuchen Lin is a Ph.D. candidate in USC CS, working with Prof. Xiang Ren at the Intelligence and Knowledge Discovery Research Lab
- Develop intelligent systems that demonstrate a deep understanding of the world with common-sense knowledge and reasoning ability
- Research interests: information extraction, knowledge graphs, logical reasoning, graph neural networks, explanations, robustness

Chaoyang He, PhD candidate @ USC



advised by Professor Salman Avestimehr, Professor Mahdi Soltanolkotabi, and Professor Murali Annavaram (USC). He also works closely with researchers/engineers at Google, Facebook, Amazon, and Tencent. Previously, He was an R&D Team Manager and Staff Software Engineer at Tencent (2014-2018), a Team Leader and Senior Software Engineer at Baidu (2012-2014), and a Software Engineer at Huawei (2011-2012).

• Chaoyang He is a Ph.D. Candidate in the CS department at the University of Southern California. He is

- His research focuses on distributed/federated machine learning algorithms, systems, and applications (<u>https://FedML.ai</u>, <u>https://DistML.ai</u>).
- Chaoyang He has received a number of awards in academia and industry, including Amazon ML
 Fellowship (2021-2022), Qualcomm Innovation Fellowship (2021-2022), Tencent Outstanding Staff
 Award (2015-2016), WeChat Special Award for Innovation (2016), Baidu LBS Group Star Awards (2013),
 and Huawei Golden Network Award (2012).

A Surprisingly "Simple" Recipe for Modern NLP



A Surprisingly "Simple" Recipe for Modern NLP



The "pre-train then fine-tune" paradigm for NLP



Randomly masked A quick [MASK] fox jumps over the [MASK] dog



WIKIPEDIA The Free Encyclopedia





Fine-tuning a pretrained transformer

What happen in workplace..



Disadvantages of Data-Centralized Learning

- User Privacy Concerns
- Data Sharing Regulation Laws
- High Cost of Transferring Raw Data
- Expensive Computation of Centralized Training



Federated Learning (FL)



Taking *FedAvg* as an example FL method

the local data is never exposed to others!

Federated Learning for NLP: Promises and Challenges

Background: Current FL research mainly focus on testing methods on *toy* datasets or computer vision (image classification).

Our goal: Provide a universal platform for benchmarking and developing FL methods for various NLP tasks.



Different NLP tasks have distinct task formulations.

Interface between Transformer LMs and FL

- Most existing FL studies are customized for computer vision tasks/datasets, with specialized model architectures
- Modern NLP are primarily based on pre-training & fine-tuning Transformer LMs.
- No existing FL framework connecting Transformer LMs with FL methods.



Creating non-IID datasets for FL in NLP is an open problem

- Current NLP datasets are mainly collected for *centralized* learning, and thus does not have a *natural*, *non-IID* partition
- Existing FL datasets are mainly for computer vision tasks such as object detection
- Ideal yet not accessible: private user data in large companies; but we cannot make them public for community use





A pool of clients where each client has a relatively unique data distribution.

What happen in workplace..





Three ways to create non-IID data partitions.

Task	Txt.Cls.	Seq.Tag.	QA	Seq2Seq
Dataset	20News	Onto.	MRQA	Giga.
# Training # Test # Labels	11.3k 7.5k 20	50k 5k 37*	53.9k 3k N/A	10k 2k N/A
Metrics	Acc.	F-1	F-1	ROUGE

We select four typical datasets for each formulation.

A showcase of non-IIDness on 20news dataset using *Dirichlet Allocation Methods*



- Label distribution: darker color
 higher probability for a client's data assigned with a certain label
- Smaller alpha
 more distinct label distributions between the clients (non-IID)
- When alpha=100
 uniform label

 distribution for all clients

Our Proposed FedNLP framework

- We support a wide range of FL methods such as FedAvg, FedOpt, FedNova, FedProx, etc.
- We implement the interface between these FL methods to Transformer LMs.
 - We support fine-tuning the full model as well as part of the models (last few layers).
 - All task formulations.

Algorithm 1: FEDOPT (Reddi et al. 2020)): A Generic FedAvg Algorithm

Input: Initial model $\boldsymbol{x}^{(0)}$, CLIENTOPT, SERVEROPT 1 for $t \in \{0, 1, \dots, T-1\}$ do Sample a subset $\mathcal{S}^{(t)}$ of clients 2 for client $i \in \mathcal{S}^{(t)}$ in parallel do On each client. 3 Initialize local model $m{x}_i^{(t,0)} = m{x}^{(t)}$ Download from server. 4 for $k = 0, ..., \tau_i - 1$ do 5 Compute local stochastic gradient $g_i(\boldsymbol{x}_i^{(t,k)})$ 6 Perform local update $\boldsymbol{x}_{i}^{(t,k+1)} =$ 7 CLIENTOPT $(\boldsymbol{x}_{i}^{(t,k)}, q_{i}(\boldsymbol{x}_{i}^{(t,k)}), \eta, t)$ Compute local model changes 8 $\Delta_i^{(t)} = oldsymbol{x}_i^{(t, au_i)} - oldsymbol{x}_i^{(t,0)}$ Upload to server. Aggregate local changes 9 $\Delta^{(t)} = \sum_{i \in \mathcal{S}^{(t)}} p_i \Delta_i^{(t)} / \sum_{i \in \mathcal{S}^{(t)}} p_i$ Update global model 10 $\boldsymbol{x}^{(t+1)} = \operatorname{ServerOpt}(\boldsymbol{x}^{(t)}, -\Delta^{(t)}, \eta_s, t)$ On the server.





Q1: How do popular FL methods perform over different NLP tasks?



Y-axis: evaluation metric of the task. Higher the better. X-axis: # of iterations in the algorithm

- FedOpt outperforms the other two FL methods in the first three tasks
- FedProx and FedAvg are comparable with each other.

Q2: How do different non-IID partitions of influence FL performance?

Figure 5: Testing FedOPT with DistilBERT for 20News under different data partition strategies.

- Non-IID data partition (smaller alpha) creates more challenges for FL methods to perform
- Non-IID data partition (smaller alpha) also makes the FL algorithm less stable
- Uniform and quantity-skew partitions are less challenging to learn

Q3: How does freezing of Transformers influence the FL performance?

Q4: Are compact model DistilBERT adequate for FL+NLP?

- BERT-base is 2x larger than DistilBERT

- DistilBERT is a more cost-effective choice.
- It's reasonable to do experiments with
 DistilBERT as the curve is similar to BERT-base.

Figure 7: FedOPT for 20News with different LMs.

Task 1 [heterogenous FedNLP]: Current FL methods focus on the case where all local models are of the same architecture and model size. This is inflexible and can cause problems when users have different devices.

How should we perform FL when clients are using different BERT-style architecture?

Task 2 [privacy]: Concerns about Transformer LMs that can memorize private information.

Can we quantitively measure such data leak? Can we design methods to prevent this?