

Federated Learning White Paper V1.0

WeBank, Shenzhen, China

WeBank AI Group

2018.9

Content

1. Introduction	5
1.1 Background.....	5
1.2 The GDPR and New Challenge of AI	6
1.3 Federated Learning a Feasible Solution	6
2 Federated Learning	8
2.1 An Overview of Federated Learning.....	8
2.2 Definition of Federated Learning.....	8
3 A Categorization of Federated Learning.....	10
3.1 Horizontally Federated Learning	11
3.2 Vertically Federated Learning	11
3.3 Federated Transfer Learning	11
3.4 Architecture for a Federated Learning System	12
4 Related Work.....	14
4.1 Federated Learning vs Differential Privacy	14
4.2 Federated Learning vs Distributed Machine Learning.....	14
4.3 Federated Learning vs Federated Database System	15
5 Applications of Federated Learning.....	16
5.1 Intelligent Marketing.....	16
5.2 Intelligent Diagnosis.....	17
5.3 Federated Learning and Industry Data Alliance	17
6 Development Roadmap of Federated Learning.....	19
6.1 Establishing Federated Learning Domestic and Global Standards.....	19
6.2 Building Federated Learning Usecase in Vertical Market	19
6.3 Forming Federated Learning Industrial Data Alliance	20
7 Conclusions and Prospects.....	21
Reference	22

1. Introduction

1.1 Background

2016 is the year when AI came of age. With AlphaGo defeating two top human Go players in succession, we have truly seen the huge potential of AI, and began to expect more complex, more cutting-edge AI technology in driverless cars, medical care, finance, etc. The field is showing its strengths. However, when we look back at the development of AI, it is inevitable that the development of AI has experienced two troughs and three peaks. The two troughs of AI were due to the lack of algorithms, computing power and data. Now the AI driven by a big data environment has entered the third golden development period.

The success of AI relies on the availability of big data. Deep learning systems that recognize images require tens of millions of training images to reach top performance. This is true not only in computer vision, but speech recognition, question answering chatbots, and large-scale recommendation and prediction systems that empower e-commerce systems. A typical example is AlphaGo in 2016, which used a total of 300,000 games as training data and achieved excellent results. With AlphaGo's success, people naturally hope that the big data-driven AI like AlphaGo will be realized in all aspects of life soon. However, the real situation is very disappointing: with the exception of few industries, most fields have only limited data or poor quality data, making the realization of AI technology difficult. The majority organizations and applications only have 'small data', as data collection is often costly if not impossible today. That is the case with many medical applications such as diagnosis, drug design, and health care. Many of these datasets are scattered across different organizations, departments, and businesses. These data may look like isolated islands on a vast ocean, and we may refer to this problem as 'small-data problem'. At the same time, it is also hard to break the barriers between data sources. In general, the data required by AI involves multiple fields. For example, in an AI-driven product recommendation service, the product seller has information about the product, data of the user's purchase, but no data of the user's purchasing ability and payment habits. In most industries, data exists in the form of isolated islands. Due to industry competition, privacy security, and complicated administrative procedures, even data integration between different departments of the same company faces heavy resistance. It is almost impossible to integrate the data scattered around the country and institutions, or the cost is prohibited.

1.2 The GDPR and New Challenge of AI

On the other hand, with the advancement of big data, the emphasis on data privacy and security has become a worldwide trend. Every leak of public data will cause great concern to the media and the public. For example, the recent data breach of Facebook has caused a wide range of protests. At the same time, countries are strengthening the protection of data security and privacy. Take the General Data Protection Regulation (GDPR)^[11], which was enforced by the European Union on May 25, 2018, for example. GDPR aims to protect users' personal privacy and data security. It requires business to use clear and plain language for its user agreement and grants users the "right to be forgotten", that is, users can have their personal data deleted or withdrawn. The GDPR has nearly banned all kinds of autonomous activities in collecting, transferring and using user data. Which means, it is no longer acceptable to simply collect sources of data and integrate them in one location without user permission. Also, many normal operations in the big data domain, such as merging user data from various source parties for building an AI model without any user agreement, are to be considered illegal in the new regulatory framework. The GDPR brings a fundamental shift in the protection of data and privacy, shaping the way how businesses operate; companies will face serious monetary fines for the violation of the regulation.

Similarly, China's Cyber Security Law^[12] and the General Principles of the Civil Law^[13], implemented in 2017, also pointed out that internet business must not disclose, tamper with, or destroy the personal information they collect, and when conducting data transactions with third parties, they need to ensure that the proposed contract clearly defines the scope of the data to be traded and the data protection obligations. The establishment of these regulations to various degrees poses new challenges to the traditional data processing of AI.

In the field of AI, the traditional data processing model often involves in one party collecting and transferring data to another party for processing, cleaning and modeling, and finally selling the model to a third party. However, as the above regulations and monitoring become stricter and more substantial, it is possible to break the law by leaving the collector or the user unclear about the specific use of the model. Our data is already in the form of isolated islands. A straightforward solution is to collect all the data to one place for processing. However, it is now illegal to do so because the law does not allow businesses to arbitrarily consolidate data. How to legally solve the problem of isolated data islands is a major challenge for AI scholars and practitioners, because the big data dilemma is likely to lead to the next AI winter.

1.3 Federated Learning a Feasible Solution

It is thus argued that for AI to be a genuinely successful and transforming technology, there need to be efforts on two fronts to address the challenges of the small-data problem and data privacy problem. However, traditional methods for solving this dilemma of big data have run into bottlenecks. Simply exchanging data between two

companies is not allowed by many regulations including GDPR. First, the user is the owner of the original data, and the company cannot exchange data without the user's approval. Second, the usage of models can't be changed until the user approves it. Therefore, many attempts at exchanging data in the past, such as Data Exchanges, also require drastic changes to be compliant. At the same time, the data owned by commercial companies often has huge potential value. Two organizations and even two departments of the same organization must consider the interests of exchanging data. Under this premise, one department often choose not to consolidate data with other departments, resulting in data appearing in isolated islands even in the same company.

We propose to shift the focus of research to how to solve the big data dilemma, that is, the problem of isolated data islands. We believe that the focus of AI in the next step will shift from the AI-based algorithm to the algorithm-oriented big data architecture that guarantees security and privacy. Here we present a possible solution called federated learning^[14-15]. The federative learning framework intends to make industries effectively and accurately use data across organizations while meeting the privacy, security and regulatory requirements, in addition to building more flexible and powerful models to enable business cooperation by using data collectively but without data exchange directly.

Federated learning is a system that:

- Data distributed located in each data entities, with no privacy revealing and no compliance violation.
- Multiple data parties build a virtual shared model under a data federation system, gaining mutual benefit from the system.
- Under such a federal mechanism, the identity and status of each participant are the same.
- This virtual model has the same, or nearly the same performance as the model that built by putting all data together.

Federated learning permit learning to be done while multiple data sets stay put – no data exchanges are needed on the raw data to protect privacy and secrecy, providing a feasible solution to the date isolation problem.

2 Federated Learning

2.1 An Overview of Federated Learning

What is Federated Learning? Suppose there are two companies A and B with different data. For example, Company A has user profile data; Company B has product feature data and label data. According to the above GDPR guidelines, the two companies cannot rudely combine the data of both parties because the original providers of the data, their respective users, did not agree to do so. Suppose that each party builds a learning model for a classification or prediction task respectively, and these tasks are already recognized by their respective users when the data is obtained. Now the question is how to build higher quality models for both A and B. However, because the data is incomplete (for example, A lacks label data, and B lacks feature data), or data is insufficient (the amount of data is insufficient to build a good model), the models at each end may not be established or the results may not be satisfactory.

Therefore, the purpose of federal learning is to solve this problem: it aims at building a model across organizations while individual data of each organization stay in their local environment, and model parameters are exchanged under encryption mechanism in a federated system. That is, a virtual shared model is built without violating the data privacy regulations. This virtual model has the same performance as the model that you build by putting all data together. But when building a virtual model, the data itself does not move, nor does it reveal privacy or affect data specifications. In this way, the model built serve only local tasks in their respective regions. Under such a federal mechanism, the identity and status of each participant are the same, and the federal system helps everyone establish a "common wealth" strategy, which is why this system is called "federated learning."

The above examples illustrate the basic ideas of federated learning. In the following, we provide a formal definition for federated learning, as well as its categorizations based on the distribution characteristics of the island data. Finally, the workflow and system architecture of the federated learning system are described.

2.2 Definition of Federated Learning

Define multiple data owners $F_i, i=1...N$ who all wish to train a machine learning model by consolidating their respective data D_i . A conventional method is to put all data together and use $D=\{D_i, i=1...N\}$ to train a model M_{sum} . However, this solution is not possible to implement due to legal issues such as privacy and data security. To solve this problem, we propose federal learning. Federated Learning is a learning process in which data owners collaboratively train a model M_{FED} and in the process any data owner F_i

does not expose its data D_i . In addition, the performance of M_{FED}, V_{FED} should be very close to the performance of M_{SUM}, V_{SUM} . That is,

$$|V_{FED} - V_{SUM}| < \delta, \delta \text{ is bounded.}$$

Federated learning is first proposed to deal with the pain points of financial institutions, especially private commercial banks like WeBank. A use case is detecting multi-party borrowing, which always been a headache in the banking industry, especially in the Internet finance industry. Multi-party borrowing refers to the return of a bad user to another lending institution after borrowing from a financial institution. A large number of such illegal actions will cause the entire financial system to collapse. To find such users, the traditional approach is that financial institutions go to a central database to query user information, and each organization must upload all their users, but this is equivalent to exposing all important user privacy and data security of financial institutions, which is not allowed under GDPR. Under federated learning, there is no need to establish a central database, and any financial institution participating in federated learning can use the federated mechanism to issue new user queries to other agencies within the federation. Other agencies only need to answer questions about local lending without knowing the specific information of the user. This can not only protect the privacy and data integrity of existing users in various financial institutions, but also solve the important issue of querying multi-party lending.

3 A Categorization of Federated Learning

The above definition of federated learning does not discuss how to specifically design and implement federated learning. In practice, the island data has different distribution characteristics. According to these characteristics, we can propose a corresponding federated learning framework. Below, we will classify federal learning based on the feature and sample ID distribution characteristics of the island data.

Considering that there are multiple data owners, the data set D_i held by each data owner can be represented by a matrix. Each row of the matrix represents a user, and each column represents a user feature. At the same time, some data sets may also contain label data. If you want to build a predictive model of user behavior, you must have label data. We can call the user feature X and the labels Y . For example, in the financial field, the user's credit is the label Y that needs to be predicted; in the marketing field, the label is the user's purchase desire Y ; in the education field, Y is the degree of knowledge of the student. The user feature X plus the label Y constitutes the complete training dataset (X, Y) . However, in reality, it is common that the users of the various data sets are not identical, or the user characteristics are not identical. Specifically, taking federated learning with two data owners as an example, the data distribution can be divided into the following three cases:

- The overlap of features (X_1, X_2, \dots) is large, whereas the overlap of users (U_1, U_2, \dots) is small;
- The overlap of users (U_1, U_2, \dots) is large, whereas the overlap of features (X_1, X_2, \dots) is small;
- The overlap of users (U_1, U_2, \dots) and the overlap of features (X_1, X_2, \dots) are both small.

In order to provide solutions for the above three scenarios, we classify federated learning into horizontally federated learning, vertically federated learning and federated transfer learning (shown in Figure 1).

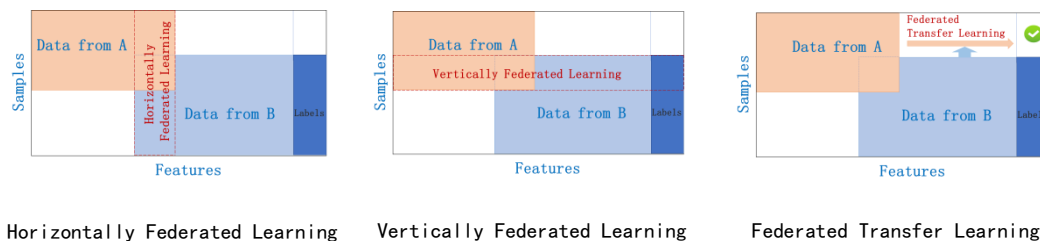


Figure 1 The categorization of federated learning

3.1 Horizontally Federated Learning

In the scenarios that two data sets share the same feature space but different in samples, a federated learning system is called horizontal federated learning. For example, two regional banks may have very different user groups from their respective regions, and the intersection of users is very small. However, their business is very similar, so the recorded user features are the same. In this case, a horizontal federated learning model can be built. In 2017, Google proposed a horizontal federated learning solution for Android phone model updates^[6-7]: A single user using an Android phone constantly updates the model parameters locally and uploads the parameters to the Android cloud, thus jointly training the centralized model together with other data owners. A secure aggregation scheme to protect the privacy of aggregated user updates under their federated learning framework is also introduced.

3.2 Vertically Federated Learning

Vertically federated learning is applicable to the cases that two data sets share the same users but differ in feature space. For example, consider two different companies in the same city, one is a bank, and the other is an e-commerce. Their user base is likely to contain most of the residents of the area, so the size of common users is large. However, since the bank records the user's revenue and expenditure behavior and credit rating, and the e-commerce retains the user's browsing and purchasing history, their user features are very different. Vertically federated learning is the process of aggregating these different features in an encrypted state and computing the training loss and gradients in a privacy-preserving manner to build a model with both data collaboratively. At present, machine learning models such as logistic regression models, tree structure models and neural network-based models have all been proved being able to incorporate into this federated system.

3.3 Federated Transfer Learning

Federated Transfer Learning applies to the scenarios that the two data sets differ not only in samples but also in feature space. In this case, transfer learning^[9] techniques can be applied to overcome the lack of data or labels. Consider two institutions, one is a bank located in China, and the other is an e-commerce company located in the United States. Due to geographical restrictions, the user groups of the two institutions have a small intersection. On the other hand, due to the different businesses, only a small part of the data features of the two companies overlap. In this case, in order to carry out effective federated learning, it is necessary to introduce transfer learning to solve the problem of small data size and weak supervision, thereby improving the performance of the model.

3.4 Architecture for a Federated Learning System

In this section, we use the vertically federated learning as an example to introduce the architecture of the federated learning system and to explain the detailed process of how it works.

First, let's take the scenario of two data owners (i.e, companies A and B) as an example to introduce the architecture of the federated learning system, which can be extended to scenarios involving multiple data owners. Suppose that companies A and B want to jointly train a machine learning model, and their business systems each have their own data. In addition, Company B also has label data that the model needs to predict. For data privacy and security reasons, A and B cannot directly exchange data. At this point, the model can be built using the federated learning system, which consists of two parts, as shown in Figure 2a.

Part 1: Encrypted entity alignment. Since the user groups of the two companies are not the same, the system uses the encryption-based user ID alignment technology to confirm the common users of both parties without A and B exposing their respective data, and the system does not expose users that do not overlap with each other.

Part 2 : Encrypted model training. After determining the common entities, we can use these common entities' data to train the machine learning model. In order to ensure the confidentiality of the data during the training process, it is necessary to use a third-party collaborator C for encryption. Taking the linear regression model as an example, the training process can be divided into the following four steps (as shown in Figure 2b):

- Step ①: collaborator C creates encryption pairs, send public key to A and B;
- Step ②: A and B encrypt and exchange the intermediate results for gradient and loss calculations;
- Step ③: A and B computes encrypted gradients respectively, and B also computes encrypted loss; A and B send encrypted values to C.
- Step ④: C decrypts and send the decrypted gradients and loss back to A and B; A and B update the model parameters accordingly.

Iterations through the above steps continue until the loss function converges, thus completing the entire training process. During entity alignment and model training, the data of A and B are kept locally, and the data interaction in training does not lead to data privacy leakage. Therefore, the two parties achieve training a common model cooperatively with the help of federated learning.

Part 3: Incentives Mechanism. A major characteristic of federated learning is that it solves the problem of why different organizations need to jointly build a model. After the model is built, the performance of the model will be manifested in the actual applications and recorded in a permanent data recording mechanism (such as blockchain). Organizations that provide more data will be better off, and the model's effectiveness depends on the data provider's contribution to the system. The effectiveness of these models are distributed to parties based on

federated mechanisms and continue to motivate more organizations to join the data federation.

The implementation of the above three steps not only considers the privacy protection and effectiveness of collaboratively-modeling among multiple organizations, but also considers how to reward organizations that contribute more data, and how to implement incentives with a consensus mechanism. Therefore, federated learning is a "closed loop" learning mechanism.

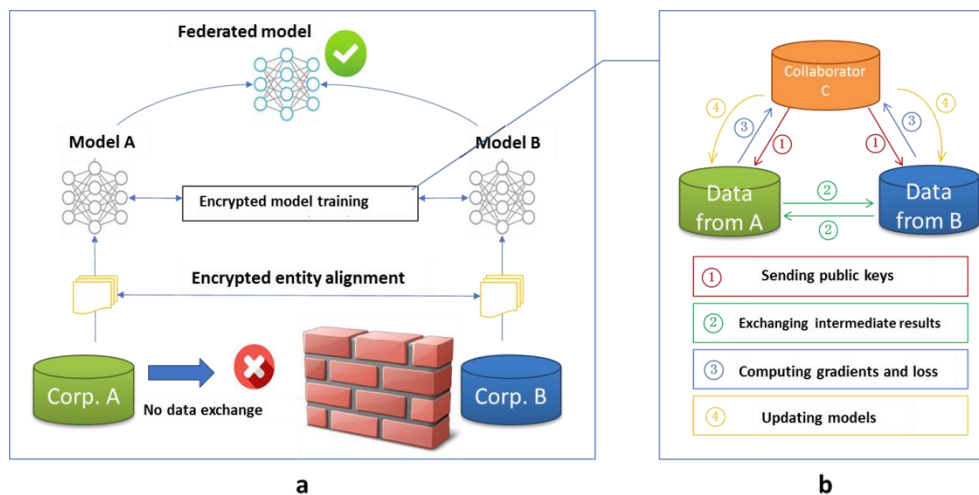


Figure 2 Architecture for a federated learning system

4 Related Work

As a novel technology, federated learning has some originality while drawing on some mature technologies. Below we explain the relationship between federated learning and other related concepts from multiple perspectives.

4.1 Federated Learning vs Differential Privacy

Federated Learning protects the privacy of user data in a very different way that Differential Privacy^[1], k-Anonymity^[2] or 1-Diversity^[3] does for data privacy protection. First, federated learning protects user data privacy through parameter exchange under the encryption mechanism. The encryption scheme includes homomorphic encryption^[10]. Unlike differential privacy protection, the data and the model itself are not transmitted, nor can they be guessed by the other party's data. Therefore, there is no possibility of leakage at the raw data level, nor does it violate stricter data protection laws such as GDPR. The diversification of methods like differential privacy, k-anonymity, and 1-Diversity, involve in adding noise to the data, or using generalization methods to obscure certain sensitive attributes until the third party cannot distinguish the individual, thereby making the data impossible to be restored to protect user privacy. However, the root of these methods are still transmitting raw data, therefore there is a possibility of potential attack, and the protection of data privacy may not be applicable under stricter data protection regulations such as GDPR. Correspondingly, federated learning is a more powerful tool of protecting user data privacy.

4.2 Federated Learning vs Distributed Machine Learning

Horizontally federated learning at first sight is somewhat similar to Distributed Machine Learning. Distributed machine learning covers many aspects, including distributed storage of training data, distributed operation of computing tasks, distributed distribution of model results, etc. Parameter Server^[4] is a typical element in distributed machine learning. As a tool to accelerate the training process, the parameter server stores data on distributed working nodes, allocates data and computing resources through a central scheduling node, so as to train the model more efficiently. For horizontally federated learning, the working node represents the data owner. It has full autonomy for the local data, and can decide when and how to join the federated learning. In the parameter server, the central node always takes the control, so federated learning is faced with a more complex learning environment. Secondly, federated learning emphasizes the data privacy protection of the data owner during the model training process. Effective measures to protect data privacy can better cope with the increasingly stringent data privacy and data security regulatory environment in the future.

4.3 Federated Learning vs Federated Database System

Federated Database System^[5] is a system that integrates multiple database units and manages the integrated system as a whole. It is proposed to achieve interoperability with multiple independent databases. The federated database system often uses distributed storage for database units, and in practice the data in each database unit is heterogeneous. Therefore, it has many similarities with federated learning in terms of the type and storage of data. However, the federated database system does not involve any privacy protection mechanism in the process of interacting with each other, and all database units are completely visible to the management system. In addition, the focus of the federated database system is on the basic operations of data including inserting, deleting, searching, and merging, etc., while the purpose of federated learning is to establish a joint model for each data owner under the premise of protecting data privacy, so that the various values and laws the data contain serve us better.

5 Applications of Federated Learning

5.1 Intelligent Marketing

As a modeling method that could ensure data security, federated learning has a promising future in the financial and marketing industry, where raw data could not be aggregated brutally for training models in consideration of intellectual property, data privacy and data security problems. Therefore, it is expected to train a federated model without sharing data, which could be fulfilled by federated learning.

Take the intelligent marketing as an example. The purpose of intelligent marketing is to provide personalized services such as products recommendation for clients via the machine learning techniques. Data features involved in this task mainly include the purchasing power and the preference of the clients, as well as the characteristics of products. In real applications, these three types of features could be distributed in different companies. For example, the purchasing power of a people could be inferred by her bank savings and her preference for different products is shown on her social network, while the characteristics of products are recorded on an e-shop. Actually, we are facing two problems. First, in order to ensure the data privacy and security, it is hard to break the barrier of data between the bank, the social network and the e-shop. As a result, the data could not be aggregated directly. Second, the data stored in the three companies is usually heterogeneous, and traditional machine learning methods would not work on heterogeneous data. These problems are not solved efficiently till now, which impede the development of artificial intelligence techniques in many fields.

Fortunately, the arising of federated learning brings hope to solve these problems. Imagine that we build a united model for the data from the three companies through federated learning and transfer learning in the intelligent marketing task. We will have the following benefits. First, we could train a federated model without exporting the data from either company, which is guaranteed by the characteristics of federated learning. Such method could not only protect the data privacy but also provide personalized services, further realize the purpose of benefiting together. Meanwhile, we could tackle the data heterogeneity problem by the ideas of transfer learning. Transfer learning aims to learn knowledge from data and transfer the knowledge to deal with other tasks, which could break the limitations of traditional artificial intelligence techniques. It is believed that federated learning would play an important role in establishing a cross-enterprise, cross-field and cross-data ecosphere for artificial intelligence.

5.2 Intelligent Diagnosis

Intelligent diagnosis is a popular topic that combines medicine and artificial intelligence together. However, existing intelligent diagnosis systems are far from real “intelligence”. In this part, we will discuss the shortages of the current intelligent diagnosis systems with an example of the IBM Watson, and propose a conception that could overcome the shortcomings with the help of federated learning.

IBM’s supercomputer system Watson is one of the most famous applications in the field of intelligent diagnosis. In the medicine industry, Watson is used for automated medical diagnosis, especially for cancer by medical institutions from China, America and many other countries. However, Watson is suffering from people’s doubts recently due to an exposed file that reveals a misdiagnosis which is possible to cause a death. So why would Watson make such a misdiagnosis? We find out that the training data used by Watson should have contained the features of diseases, gene sequences, medical reports, examination results and academic papers. But in reality, there are no stable data sources and the majority of data is facing the problem of missing labels. Someone assumes that it would take 10 years for 10,000 experts to gather an effective dataset. The insufficiency of data and labels results in the bad performance of machine learning models, which becomes the bottleneck of intelligent diagnosis.

People would ask, “How to break through the bottleneck, then?” We assume that medical institutions all over the world unite together by sharing their data, and then they would possess a dataset large enough to train a model far better than before. Combining federated learning with transfer learning is vital to achieve this goal. The feasibility lies in the two following reasons. First, the data from medical institutions must be sensitive to privacy and security problems and brutal data exchange would be infeasible, while federated learning allows training models without sharing data directly. Second, the problem of missing labels is severe, and transfer learning could be applied to fill the missing labels, so that the dataset could be enlarged and the performance of the model could be improved largely. Therefore, federated transfer learning would play an important role in the development of intelligent diagnosis systems. If all the medical institutions could establish a data alliance together in the future, the medical level of human would step up to a new stage.

5.3 Federated Learning and Industry Data Alliance

Federated learning is not only a technology standard but also a business model. When people realize the effects of big data, the first thought that occurs to them is to aggregate the data together, compute the models through a remote processor and then download the results for further use. Cloud computing comes into being under such demands. However, with the increasing importance of data privacy and data security and a closer relationship between a company’s profits and its data, the cloud computing model has been challenged.

However, the business model of federated learning has provided a new paradigm for

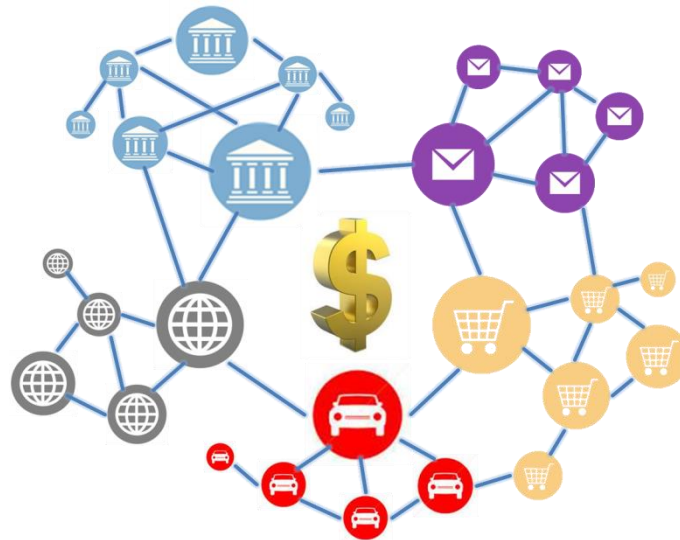


Figure 4 Data alliance allocates the benefits on blockchain

applications of big data. When the isolated data occupied by each institution fails to produce an ideal model, the mechanism of federated learning makes it possible for institutions and enterprises to share a united model without data exchange. Furthermore, federated learning could make equitable rules for profits allocation with the help of consensus mechanism from blockchain techniques. The data possessors, regardless of the scale of data they have, will be motivated to join in the data alliance and make their own profits. We believe that the establishment of the business model for data alliance and the technical mechanism for federated learning should be carried out together. We would also make standards for federated learning in various fields to put it into use as soon as possible.

6 Development Roadmap of Federated Learning

Considering the regulatory requirement, industrial pain points and the application scenario of the AI and big data community, we suggest the development roadmap of federated learning should be: 1) establishing federated learning domestic and global standards; 2) building federated learning usecase in vertical market; 3) forming federated learning industrial data alliance.

6.1 Establishing Federated Learning Domestic and Global Standards

More effort has been put in establishing AI related standards in recent years. International Organization for Standardization had found technical program subcommittee focusing on artificial intelligence standards (ISO/IEC JTC 1/SC 42) in October 2017. The US and Germany have program proposals for AI terminology and model guideline. In China, with the highly support from the SAC and MIIT, the National Artificial Intelligence Standardization Group was founded in January 2018, with members including famous industrial enterprise and academic institutes, aiming to promote the standard system of the AI area.

Studying and establishing federated learning domestic (e.g., AIOSS standard) and global standards (e.g., IEEE standard), and formulating architectural framework and application guideline of federated learning, will greatly facilitate industry cooperation. In the meantime, such standard and guideline will help different data entities make full use of data that scattered and isolated in different organizations. By addressing the privacy and security issues, establishing federated learning standard will set up a paradigm framework of mutual benefits win-win cooperation among data entities and industrial alliance

6.2 Building Federated Learning Usecase in Vertical Market

The application and commercialization of federated learning in industry will bring a revolutionary change to the current landscape. The application scenarios of federated learning in the vertical market can be categorized to intra-business market and inter-business market. The intra-business market refers to the scenario that two parties are both in the same or similar industry, for example, business cooperation between banks or financial institutes. In this scenario, data owned by two parties have almost the same attributes and features, but usually with different user IDs. Applying horizontally

federated learning in this scenario will build a better collective model in a way of increasing sample size as if all the samples from the two data sources are put together. The inter-business market refers to the scenario that two parties are in different industry, for example, business cooperation between a bank and e-commerce company. In this scenario, data owned by two parties share some overlapping user IDs, but with very different attributes and features, say, bank has users' income and transaction behavior and e-commerce company has users' shopping behavior. Applying vertically federated learning in this scenario will build a better collective model in a way of increasing the number of features as if all the features from the two data sources are put together. Both of two application scenarios will yield better models than each data party building models using its own data.

Promoting application of federated learning in vertical market, especially inter-business market, will help building a new business paradigm and ecosystem based on the framework.

6.3 Forming Federated Learning Industrial Data Alliance

The above mentioned new business paradigm based on federated learning is best to operate within and supported by an industrial data alliance. The alliance may have N entities, by joining the alliance, entities can cooperate with each other using data under federated learning framework. Companies and organizations are encouraged to join the alliance. Alliance will have clear incentive mechanism. Members in the alliance enjoy rights and interests, and also fulfill responsibilities. The alliance may use blockchain to build consensus of all parties, record each party's contribution, and award parties that yielding outstanding contribution. With consensus and technical support, we can design industrial data alliance in many vertical markets, for example, financial industry can form a financial data alliance, while medical industry can form a medical data alliance.

7 Conclusions and Prospects

In recent years, the isolation of data and the emphasis on data privacy are becoming the next challenges for artificial intelligence, but federated learning has brought us new hope. It could establish a united model for multiple enterprises while the local data is protected, so that enterprises could win together taking the data security as premise. This article generally introduces the basic concept, architecture and techniques of federated learning, and discusses its potential in various applications. It is expected that in the near future, federated learning would break the barriers between industries and establish a community where data and knowledge could be shared together with safety, and the benefits would be fairly distributed according to the contribution of each participant. The bonus of artificial intelligence would finally be brought to every corner of our lives.

Reference

- [1] Dwork C. Differential privacy: A survey of results[C]//International Conference on Theory and Applications of Models of Computation. Springer, Berlin, Heidelberg, 2008: 1-19.
- [2] Sweeney L. k-anonymity: A model for protecting privacy[J]. International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, 2002, 10(05): 557-570.
- [3] Li N, Li T, Venkatasubramanian S. t-closeness: Privacy beyond k-anonymity and l-diversity[C]//Data Engineering, 2007. ICDE 2007. IEEE 23rd International Conference on. IEEE, 2007: 106-115.
- [4] Ho Q, Cipar J, Cui H, et al. More effective distributed ml via a stale synchronous parallel parameter server[C]//Advances in neural information processing systems. 2013: 1223-1231.
- [5] Sheth A P, Larson J A. Federated database systems for managing distributed, heterogeneous, and autonomous databases[J]. ACM Computing Surveys (CSUR), 1990, 22(3): 183-236.
- [6] Konečný J, McMahan H B, Yu F X, et al. Federated learning: Strategies for improving communication efficiency[J]. arXiv preprint arXiv:1610.05492, 2016.
- [7] McMahan H B, Moore E, Ramage D, et al. Communication-efficient learning of deep networks from decentralized data[J]. arXiv preprint arXiv:1602.05629, 2016.
- [8] Hardy S, Henecka W, Ivey-Law H, et al. Private federated learning on vertically partitioned data via entity resolution and additively homomorphic encryption[J]. arXiv preprint arXiv:1711.10677, 2017.
- [9] Pan S J, Yang Q. A survey on transfer learning[J]. IEEE Transactions on knowledge and data engineering, 2010, 22(10): 1345-1359.
- [10] Hesamifard E, Takabi H, Ghasemi M. CryptoDL: Deep Neural Networks over Encrypted Data[J]. arXiv preprint arXiv:1711.05189, 2017.
- [11] <https://www.eugdpr.org>
- [12] http://www.xinhuanet.com/politics/2016-11/07/c_1119867015.htm
- [13] http://www.npc.gov.cn/npc/xinwen/2017-03/15/content_2018907.htm
- [14] <https://zhuanlan.zhihu.com/p/42646278> 杨强：GDPR 对 AI 的挑战和基于联邦迁移学习的对策
- [15] <https://zhuanlan.zhihu.com/p/41052548> 机器之心专访杨强教授：联邦迁移学习与金融领域的 AI 落地