

RRCM: A Fairness Framework for Federated Learning

Jianyi Zhang^{1,2,*}, Wenxin Wang¹, Zhi Sun¹, Zixiao Xiang¹ Yuyang Han¹,

¹Beijing Electronic Science and Technology Institute, Beijing 100070, China

²University of Louisiana at Lafayette, Louisiana 70503, US.

*Corresponding author

zjy@besti.edu.cn

Abstract

In recent years, federated learning still suffers from unfair incentive mechanisms in industrial systems. This situation will result in data owners in the system that may no longer actively contribute higher weighted local models to the central server. In this paper, we propose a method to achieve fairness in federated learning processing by setting a reputation system, reward-punishment and cost-interest compensation mechanism. Moreover, we have introduced a method of compensating for costs and interest to rationalize the commercialization process for the federal system. Experiments on fairness, accuracy, and compensation trends on benchmark datasets show that the proposed method can achieve higher fairness than traditional frameworks.

1 Introduction

With the promulgation of the data protection regulations GDPR[Yang *et al.*, 2019] and CCPA[Truex *et al.*, 2020], federated learning (FL) has received extensive attention from all walks of life because of its ability to protect data information through the transfer of models[Hard *et al.*, 2018]. For example, WeBank officially open-sourced the world’s first industrial-grade FL framework FATE in 2019[Kholod *et al.*, 2021]. Although the research on communication routing and backdoor defense of federation technology continues to mature in recent years, the incentive mechanism may become a shortcoming that restricts its future development[Li *et al.*, 2021]. In a typical commercialized scenario of federated learning, the central coalition makes multiple iterations of the contributions uploaded by each alliance to form a global model. The global model can be used as a commercial product to obtain revenue. All clients obtain rewards from the central alliance by sharing data resources. Since the data owners in the alliance are independent communities of interest, they are rational and selfish. When the rewards obtained by the data owners do not match their contributions (the incentives are not fair), independent members of the alliance may pursue the maximization of short-term benefits. This situation puts the cooperation of the federal system at risk. However, the existing algorithms do not reasonably address the rational

distribution of incentives in current federated learning[Khan *et al.*, 2020]. Hence, one of the most significant technical challenges in the federal system is how to design a proper incentive mechanism to keep the system’s fairness.

Nowadays, the FL incentive mechanism’s reward methods mainly include income reward and gradient reward.: income reward and gradient reward. The former provides biased information and financial tips. For example, Xu Xinyi’s team rewards each participant with biased information[Xu and Lyu, 2020]. However, the introduction of unbiased information may lead to the impact of the overall system fairness. From the perspective of economics and game theory, Tu Xuezheng’s team solves the distribution problem of incentive mechanisms through economic rewards[Tu *et al.*, 2021]. However, this method introduces other variables, which increase the communication burden of the system and communication loss. The latter mainly relies on each round of participants to obtain gradient model formation optimizations commensurate with their contributions. However, some of the literature does not discuss the Non-IID problem of FL[Xinyi *et al.*, 2021]. The federated system assigning different models to each participant will result in the participants producing different types of data items and characteristic attributes. Therefore, the central unit cannot simply adopt the Fed-avg aggregation mode in the second iteration. Besides, most literature does not consider the drawbacks of gradient reward. For example, assigning a weight with less similarity to a participant with a lower contribution will worsen the next round of the global model[Lyu *et al.*, 2020]. This situation will eventually lead to lower weights for participants with higher contributions in subsequent gradient assignments than expected.

In addition, since there is no protection mechanism, there is a risk of being attacked by adversaries in the multi-party cooperation under the traditional FL framework. Opportunists can influence the global model of the central server by uploading uncorrelated gradients or gradients with low contributions. Therefore, the federal system needs to design penalties to prevent abnormal participants from joining[Kotsogianis and Schwager, 2006].

Finally, traditional incentives only consider the cost of joining a federated system for the technical challenge of fairness in federated learning. For example, Ye’s team proposes that model training and commercialization will take time, resulting in delays in federal system compensation[[Yu *et al.*,

2020a)]. However, it does not take into account that costs will accrue interest over time, and existing dividend schemes have not yet reasonably addressed temporary mismatches. The participant will not join the federation system if the total cost and interest are more significant than the respective benefits. At this point, the actual benefits of data owners in the federal system should include costs, profits, and incentives. Therefore, this paper makes the proposed framework more reasonable by adding a compensation mechanism.

Our contributions can be summarized as follows:

- We propose a *Reputation, Reward-punishment, and Cost-interest Mechanism* (RRCM) framework to achieve fairness of the federated learning incentive mechanism.
- RRCM iteratively calculates participants' contributions through a reputation system and assigns differentiated rewards to each participant according to different performances.
- Experiments on benchmark datasets show that our framework can achieve high fairness and satisfied results. Moreover, by introducing the cost-profit model, the incentive mechanism of federated learning becomes more reasonable.

To the best of our knowledge, this paper is the first to combine the dynamics of reputation systems, reward-punishment measures, and cost-interest mechanisms. It provides a novel framework for the federated federation to achieve a more rational distribution of FL incentives.

The remaining chapters of this paper are as follows: "Related Work" reviews the fairness standards and incentive mechanisms in the existing literature to provide the basis for the research in this paper; "Fairness Definition" describes the different fairness measures and the cooperative fairness of this paper. "RRCM Framework" details the design of each module and the inter-module relationships; "Experiments" include The data set settings and experiments are compared, and it is concluded that the RRCM framework proposed in this paper is more reasonable. Finally, this paper concludes with "Conclusion" which is used to discuss the incentive mechanism of federated learning in this paper.

2 Related Work

In this section, we review the literature on incentives for FL in to link our research with existing research.

Existing research will classify federated learning incentive mechanisms into five categories: Stackelberg game, auction, contract theory, Shapley value, and reputation system[Zeng *et al.*, 2021]. The Stackelberg game[Xiao *et al.*, 2020] is often used during the sale or purchase phase. Sarikaya *et al.*[Sarikaya and Ercetin, 2020] used the Stackelberg game model to incentivize the CPU supply of multiple workers to reduce the budget of the FL major league with fully synchronized SGD's local training time. An auction[Le *et al.*, 2020] is a mathematical tool for pricing, task assignment, and node selection. Zeng's team proposed a federated learning lightweight multi-dimensional incentive scheme Fmore based on procurement auction in the mobile edge computing scenario[Zeng *et al.*, 2020]. Contract theory[Kang *et*

al., 2019] is how participants construct and develop optimal agreements in the case of conflicting interests and unequal information levels. A contract menu is provided to participants on a public procurement timing server, and each participant proactively selects a different option without informing participants of private costs. Shapley values[Wang, 2019] derived from cooperative game theory are widely adopted for contribution evaluation and profit distribution in FL. The benefit distribution of alliance members based on the Shapley value reflects the contribution of each member to the overall goal of the alliance, which can avoid egalitarianism in distribution; In the paper[Song *et al.*, 2019], Wang *et al.* adopted a variant of the Shapley group value to measure the utility of a subset of features. They merge some private features into a standard joint part and then compute the Shapley group value of this collaborative feature in the case of two participants. The reputation system mechanism[ur Rehman *et al.*, 2020] is a common method of FL incentives. Teacher Yang Qiang's team conducts research on incentive fairness in this way. For example, the paper[Yu *et al.*, 2020b] forms a fairer incentive method by adding the reputation dynamic and regret models.

Benefit-sharing schemes can be divided into equal gain, marginal contribution, and marginal loss[Stark *et al.*, 2015]. Equal benefit means that the user-generated by the central server of the federated system is equally distributed among the data contributors participating in the system federation. Marginal contribution implies that the benefit of the participants in the system alliance is the overall increase in the benefit value of data contributors when they join the system alliance. Marginal loss means that the participants' profit in the system alliance is the overall loss value when the data contributor leaves the system alliance.

To sum up, the federated learning incentive mechanism can combine the reputation system and benefit-sharing schemes. For example, the contribution of allies is calculated by reputation trust to allocate different rewards. In addition, the incentive mechanism of federated learning can be improved by introducing discrimination rate, reward rate, punishment measures, etc.

3 Fairness Definition

A reasonable federated learning incentive mechanism needs to be fair to every participant. Fairness classification at different stages is depicted in Table 1. The primary representative of the early fairness mechanism is egalitarianism[Mohri *et al.*, 2019], and differentiated allies get the same incentives for training iterations in the system.

Today, different federated learning fairness articles have different standards for fairness measurement. According to formalizing fairness, Gajane and Pechenizkiy divide fairness into individual fairness, group fairness, unconscious fairness, preference-based fairness, and counterfactual fairness[Gajane, 2017]. Individual fairness means that if a pair of individuals have similar attributes, the federated learning algorithm should input similar probabilities[Binns, 2020]. Group fairness implies that a specific attribute should present the same possibility among different groups through the FL algorithm[Binns, 2020]. Unconscious fairness means that in-

Method	Fairness classification	Reference
Equal incentive	Egalitarianism	[Mohri <i>et al.</i> , 2019]
Formalization fairness	Individual fairness	[Binns, 2020]
	Group fairness	[Binns, 2020]
	Unconscious fairness	[Hardt <i>et al.</i> , 2016]
	Preference-based fairness	[Zafar <i>et al.</i> , 2017]
	Counterfactual fairness	[Kusner <i>et al.</i> , 2017]
Training process	Cognitive representation fairness	[Shin, 2020]
	Algorithm modeling fairness	
	Decision evaluation fairness	
Algorithm level	Contribution fairness	[Shi <i>et al.</i> , 2021]
	Regret distribution fairness	
	Expectation fairness	
Reputation, Reward, and Cost	Contribution difference	Ours

Table 1: Thematic taxonomy of machine learning fairness.

dividuals with the same type of attributes (protected attributes and general attributes) appear similar decisions in the federated learning process[Hardt *et al.*, 2016]. Preference-based fairness indicates that when multiple choices are given among different groups, individuals in the group spontaneously choose decisions that are beneficial to their development[Zafar *et al.*, 2017]. Counterfactual fairness means that the results of the protected data in the real world are consistent with the predicted results in the counterfactual world[Kusner *et al.*, 2017].

According to the training process of machine learning pre-processing, processing and post-processing, Shin divides fairness into three parts: cognitive representation fairness, algorithm modeling fairness, and decision evaluation fairness[Shin, 2020].

In terms of federated learning incentive mechanisms, most studies divide the fairness of incentives into contribution fairness, regret distribution fairness, and expectation fairness[Shi *et al.*, 2021]. Contribution fairness means that the data owner’s benefit must be positively related to its assistance. Regret distribution fairness minimizes the difference in regret and temporal regret among data owners. Desired right refers to minimizing fluctuations in data owner regret and temporal regret.

The RRCM framework proposed in this paper includes three mechanisms: reputation system, reward-punishment measure, and cost-interest. The fairness of its incentive mechanism mainly includes determining the size of the benefits based on the participants’ contributions. In addition to this qualitative relationship, we also consider the relationship between the data owner’s grant and reward described by the Pearson correlation coefficient, which is used to quantitatively represent the cooperative fairness of the federated learning incentive mechanism.

Definition 1 (Cooperative Fairness of Federated Learning Incentive Mechanism). Assuming that the actual contribution of the participants is a group of α , and the rewards obtained by them are distributed to a group of σ , the cooperation fairness of the FL incentive mechanism can be expressed as $\rho_p(\alpha, \sigma)$. $\rho_p(\cdot, \cdot)$ is Pearson correlation coefficient[Mu *et al.*, 2018]. The larger $\rho_p(\cdot, \cdot)$ is, the more cooperative fairness the

RRCM framework proposed in this paper is.

4 RRCM Framework

This section will introduce three mechanisms of a reputation system, reward-punishment measure, and cost-interest in the FL system. In this way, a federated learning incentive optimization based on the *Reputation, Reward-punishment system, and Cost-interest Mechanism* (RRCM) framework is formed. The core principle we follow is that the rewards obtained by each participant from the central alliance are related to their contribution[Nishio *et al.*, 2020].

4.1 Overall Framework

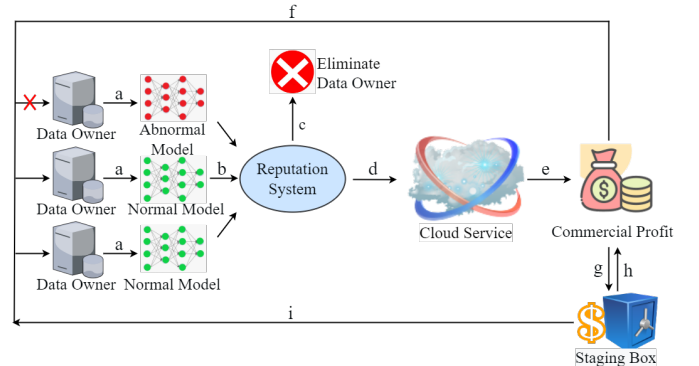


Figure 1: FL incentive mechanism RRCM framework.

The overall training process of the FL incentive mechanism RRCM framework is shown in Figure 1.

First, Data owners train local data to become local models. Local models can be divided into two categories: normal and abnormal models (a). Data owners need to pass the threshold detection of the reputation system when passing local models (b). And then, RRCM framework federal system eliminates data owners with abnormal reputations (c). Reputable local models are uploaded to the cloud server after many iterations of training to form the global model (d). Then, cloud server trades with external systems through commercial activities to

generate profits (e). A portion of the commercialization profit is compensated to reputable data owners (f). And a portion of the commercialized profits will be temporarily held in a staging box (g). Then funds from the staging box can be used to increase the cost of commercializing the investment (h). At the end of the federated system training, the staging box will return the remaining funds to the reputable data owner (i). For example, suppose that in an industry alliance system, the normal data model aggregated by local enterprises reaches the cloud server through the reputation system and participates in multiple aggregations of the central alliance. Each time a round of aggregation is completed, the Federal Center will give the enterprise part of the commercialization benefits. If the participant data models remain in good shape throughout the iterations, the Federal Center will return the remaining commercialization proceeds to individual enterprises.

The Federated learning RRCM framework incentive mechanism includes the reputation system, reward-punishment measures, and cost-interest. They are independent of each other yet connected. Cost-benefit helps the federal system more reasonably compensate data owners for costs and benefits through incentives and penalties. Benefit distribution can be rewarded and punished through the reputation system. Among them, penalties include excluding low-contributors from the RRCM framework federal system. Second, when the reputation of the data owner falls below the reputation threshold, the staging box will no longer compensate abnormal participants with staging earnings.

4.2 Reputation System

The reputation system [Gupta *et al.*, 2003] is a redirected self-feedback mechanism, which can reflect the state of its credit through the collaborative approval of related parties. It aims to show the influence of the reputation of the participants on the decision of the system. This paper adopts the reputation system as the criterion. Then, according to different standings, the central server gives each participant a profit commensurate with the contribution through different reputations. According to the study of cosine similarity representing gradient quality [Cao *et al.*, 2020]: $\cos(u, v) = \langle u, v \rangle / (\|u\| \times \|v\|)$, The contribution of each participant in this paper is represented by the cosine similarity of the local weight and the center weight: $\alpha_i^{(t)} = \cos(\Delta\omega_i^{(t)}, \Delta\omega)$. In the initial stage of the federation system, we set the same value (reputation threshold A) for each participant's initial reputation. Assuming that there is a certain positive relationship between the temporary reputation \bar{r} and the contribution degree α in this round ($\bar{r} \simeq \alpha$), \bar{r} can be equal to $\cos(\Delta\omega_i^{(t)}, \Delta\omega)$. The actual reputation can be obtained from the historical reputation and temporary reputation of this round. The reputation for each round is calculated as follows:

$$r_i^{(t)} = \beta r_i^{(t-1)} + (1 - \beta) \bar{r}_i^{(t)} \quad (1)$$

Where β is a settable weight coefficient. $r_i^{(t-1)}$ is the reputation value of the previous round. $\bar{r}_i^{(t)}$ is the temporary reputation of the current round. As a result, this framework assigns different incentives to different data owners through reputation and contribution links.

4.3 Reward-punishment Measure

Reward-punishment measures are divided into two parts: reward and punishment. Rewards are mainly given incentives to data owners through the commercialization of the federal system. By definition 1, the framework proposed in this paper can allocate incentives to various participants more reasonably. In this way, participants who contribute more can get more rewards.

Punishments are mainly done by setting a reputation threshold of A . The central federation removes data owners that fall below the reputation threshold each round from the federation system. This measure prevents low-contribution participants (such as free riders or hostile participants) from undermining the results of joint training of the system.

In addition, the federal system will temporarily store part of the incentives. When the federated learning training is over, the central server will return the temporary incentives to the reputable participants. Participants with lousy reputations will not refund the stash of incentives.

4.4 Cost-interest Mechanism

In the traditional FL system, the participation of different allies in joint training needs to be paid to the central alliance in advance. These fees are mainly used for the continuous reproduction process of the federal system. For example, data owners build local models and upload them to a central coalition. The jointly trained global model can also benefit from transactions with companies outside the alliance. However, model aggregation and commercialization will take time, which will result in the central collaboration needing to accumulate enough budget to reimburse the participating parties for the franchise cost.

The existing federated learning incentive mechanism solves the temporary mismatch between partner fees and incentives by researching incentive reward sharing schemes [Yu *et al.*, 2020a]. However, this method ignores the role of interest. The entire process takes time, from initial joining the federation system to commercializing the federal system. So over time, the central collaboration needs to repay the cost of each participant and consider the interest generated by the compensation cost.

The federated system can allow the appropriate data owners to join the scheme by requiring parties wishing to join the federation to pay the desired membership fee in advance. In the compensation process, the benefits of the alliance system first repay the cost-interest of the participants. Assuming that C_i is the cost contributed by the i -th participant to the federation, the repayment process of the i -th participant is as follows:

$$C_i \rightarrow S_i^t + \sum_{t=1}^t u_i^t (1 + \gamma) \quad (2)$$

Where S_i^t is the part of the cost of the i -th participant remaining in the central alliance in round t . u_i^t is the cost compensated to participant i by the central alliance in the t -th round. γ is the cost rate, which can be set by parameters. $u_i^t (1 + \gamma)$ represents the total return transferred to participant i in round t .

4.5 The Implementation of the RRCM

The particular implementation of RRCM in Algorithm 1 is as follows:

Algorithm 1 Reputation, Reward-punishment and Cost-interest Mechanism (RRCM)

- 1: **Input:** investment cost of each participant joining the alliance C_i , federal system incentive u_i^t , federal system rate γ , reputation threshold A .
 - 2: **Participant** i
 - 3: Download the allocated gradient $\nabla w_i^{(t-1)}$, the allocated reward $\sigma_i^t, \sigma_i^t \in T_t$
 - 4: **if** $\sum_{t=1}^t \sigma_i^t < \sum_{t=1}^t u_i^t(1 + \gamma)$ **then**
 - 5: This stage is to repay the cost
 - 6: **else**
 - 7: This stage is the actual benefit
 - 8: **end if**
 - 9: Local training $\Delta w_i^{(t)}$
 - 10: Upload local gradients $\Delta w_i^{(t)}$ to server
 - 11: **Server**
 - 12: Aggregation: $\Delta w^{(t)} = \sum_{i=1}^N \psi_i \Delta w_i^{(t)}$
 - 13: $\alpha_t = \text{cov}(\Delta w_i^{(t)}, \Delta w^{(t)})$
 - 14: **for** $i \in R$ **do**
 - 15: $\tilde{r}_i^t = \rho_p(\alpha_i^t, \sigma_i^t)$
 - 16: $r_i^{(t)} = \beta r_i^{(t-1)} + (1 - \beta)\tilde{r}_i^t$
 - 17: **if** $r_i^{(t)} < A$ **then**
 - 18: $R = R \setminus \{i\}$ Remove too low reputations
 - 19: $T_{t+1} = T_t - \sum_{i=1}^N \sum_{t=1}^t \sigma_i^t + S_i$
 - 20: **end if**
 - 21: **end for**
-

In Algorithm 1, there are two punishment measures in the RRCM framework: the first is to directly remove the participants whose reputation is lower than the reputation threshold from the federation, thus ensuring the accuracy of the training gradient aggregation of the federated system. The second is the system to increase the cost-interest compensation mechanism. The partnership will keep a portion of the cost and store it on a central server. Suppose the participant's reputation $r_i^{(t)}$ is always more significant than the reputation threshold in the complete training. The central server will return the reserved cost to the participant when the participant exits the federation system. If a participant's reputation $r_i^{(t)}$ is less than the reputation threshold, the cost compensation for the remaining storage will not be returned to this participant. This part of the money can be used for more commercialization of the system or more compensation for good players. $T_{t+1} = T_t - \sum_{i=1}^N \sum_{t=1}^t \sigma_i^t + S_i$ represents the total revenue process of the central alliance.

5 Experiment

5.1 Dataset

We selected three datasets, MNIST[Lecun *et al.*, 1998], CIFAR-10[Zhao *et al.*, 2018], and Movie Review (MR)[Pang

and Lee, 2005] to complete the control of this experiment. MNIST is a handwritten image classification dataset that includes 55,000 training data and 10,000 testing data. CIFAR-10 is a color image classification dataset with 50,000 training data and 10,000 testing data. MR is a sentiment binary classification dataset that contains 25,000 training movie reviews and 25,000 testing movie reviews.

In terms of standard IID, we choose a uniform cut of the dataset and denote it as UNI; in terms of Non-IID, considering the heterogeneity of the data, we randomly distribute among 5, 10, 20 participants according to the power-law Split 3000, 6000, 12000 MNIST samples and record it as POW[McMahan *et al.*, 2017].

5.2 Experimental Settings

We use three metrics as evaluation criteria for this experiment: accuracy, fairness and compensation trend. The accuracy is obtained by comparing the output results of the federated system with the test set. The RRCM framework proposed in this paper uses the FedAvg algorithm combined with the reputation system, reward-punishment measures, and cost-benefit mechanisms. Therefore, this experimental framework focuses on comparing with FedAvg in terms of accuracy. Fairness is quantitatively represented by cooperative fairness in Definition 1. The larger the Pearson coefficient ($\rho_p(\alpha, \sigma)$) of contribution and incentive, the more fair the federated learning framework is. For the FedAvg framework[McMahan *et al.*, 2017], the RRCM proposed in this experiment is also compared with two fairness standard frameworks, q-FFL[Li *et al.*, 2019] and CFFL[Lyu *et al.*, 2020]. In addition, we also explored the accuracy of q-FFL and CFFL in one hundred epochs of training. The compensation trend mainly compares the reward trend of the incentive mechanism in the three schemes of the incentive mechanism with no cost, cost and cost-interest, to determine the superiority of our proposed framework.

Referring to the relevant literature on reputation incentives for FL, we set the reputation threshold as $A = 1/(3N)$. According to the salary distribution principle, we set the storage cost of the central alliance as $S = 1/(10T)$. The federal center stores part of the cost to prevent data owners from providing models with lower similarity.

5.3 Experimental Results

Accuracy comparison. Table 2 lists the accuracy of different participants in the case of UNI and POW through RRCM and FedAvg. According to experimental data, The accuracy of RRCM is higher than that of the traditional fairness framework, and it is roughly the same as the accuracy of FedAvg. On the one hand, the higher accuracy of RRCM may be because q-FFL and CFFL are the frameworks that mainly address fairness, and the improvement of fairness will inevitably affect the accuracy. On the other hand, the accuracy rates of the two frameworks are roughly similar, mainly because the allocation of participants in the RRCM framework is based on the FedAvg algorithm, so their accuracy rates are not much different.

Fairness comparison. Table 3 lists the values of different cooperative fairness of other numbers of participants under

Framework	MNIST				CIFAR-10		MR
	10		20		10		5
Data Split	UNI	POW	UNI	POW	UNI	POW	POW
FedAvg	93	92	93	92	48	47	50
q-FFL	85	32	90	51	42	38	15
CFFL	90	85	90	88	40	46	42
RRCM	93	92	93	91	48	47	50

Table 2: Accuracy[%] comparison of commonly used frameworks.

the MNIST and cifar10 datasets. The Pearson coefficient can calculate the collaborative fairness value. According to the content in the table, RRCM is better than FedAvg, q-FFL, and CFFL in data training fairness on three datasets and two cuts. The scheme proposed in this paper can give data owners with higher contributions better incentives.

Framework	MNIST				CIFAR-10		MR
	10		20		10		5
Data Split	UNI	POW	UNI	POW	UNI	POW	POW
FedAvg	-30.2	77.3	3.8	-3.6	-38.6	40.1	22.1
q-FFL	-44.7	39.1	-22.0	38.7	-17.6	49.7	54.3
CFFL	83.5	91.8	82.5	94.6	72.5	76.3	92.8
RRCM	84.6	96.5	87.2	97.8	81.2	85.4	93.4

Table 3: Fairness[%] comparison of commonly used frameworks.

Compensation trend. As shown in Figure 2, it is a simulation graph of the compensation trend of the federated learning incentive mechanism in three cases. The left side represents the degree of compensation, and the right means the degree of incentive. Min n represents the minimum number of rounds for iterative aggregation of the federated system. According to the illustration, participants in the no-cost scheme do not need the central server to compensate for the cost but directly gain incentives from the alliance. The federal system will reimburse participants for costs in the cost scheme before providing incentives. In the cost-interest scheme, the federal system pays the costs and interest caused by the costs before distributing the rewards. Therefore, at the beginning of training, the cost and interest scheme does not directly reward each participant but first compensates each participant for the sum of part of the cost and interest. Furthermore, the cost-interest scheme does not gain as much compensation as the cost-interest scheme by preventing participants from contributing lower similarity weights in subsequent training. But at the end of the training, the central server will reimburse the reputable participants for the installment cost and interest.

In summary, according to fairness and accuracy, despite the accuracy of RRCM and FAV, the fairness of RRCM is higher, so the framework proposed in this paper works better. According to compensation trends, introducing cost-interest into an interest in this program can make the federal system more realistic. Therefore, the RRCM incentive mechanism proposed in this paper is superior and reasonable compared with the traditional framework.

6 Conclusions

This paper proposes a federated learning incentive optimization based on the reputation system, reward-punishment measures, and cost-interest (RRCM). It improves the cooperative fairness in FL accordingly. At the same time, the cost-

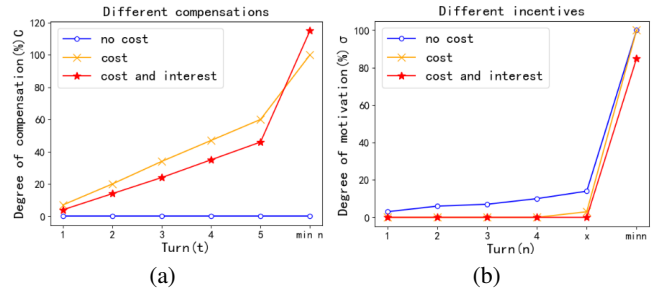


Figure 2: Compensation trend simulation graph. (a) is the simulation diagram of compensation in three ways. (b) is the compensation excitation map under the mode.

interest generated by the data owner joining the federation and the incentives obtained by the participants are positively related to their contribution degrees. According to the experiments, our proposed scheme can ensure lossless accuracy and improve fairness. Therefore, the RRCM framework presented in this paper has more advantages. In terms of reward-punishment measures, this paper proposes that the cost of punishment can be used as a reward for the participants, but it has no practical application. Subsequent experiments can further improve the reward method. It is hoped that this framework will be more optimized and perfected in the future and applied in actual enterprise alliances.

References

- [Binns, 2020] Reuben Binns. On the apparent conflict between individual and group fairness. *FAT* '20*, page 514–524, New York, NY, USA, 2020. Association for Computing Machinery.
- [Cao *et al.*, 2020] Xiaoyu Cao, Minghong Fang, Jia Liu, and Neil Zhenqiang Gong. Fltrust: Byzantine-robust federated learning via trust bootstrapping. *CoRR*, abs/2012.13995, 2020.
- [Gajane, 2017] Pratik Gajane. On formalizing fairness in prediction with machine learning. *CoRR*, abs/1710.03184, 2017.
- [Gupta *et al.*, 2003] Minaxi Gupta, Paul Judge, and Mostafa Ammar. A reputation system for peer-to-peer networks. In *Proceedings of the 13th International Workshop on Network and Operating Systems Support for Digital Audio and Video*, NOSSDAV '03, page 144–152, New York, NY, USA, 2003. Association for Computing Machinery.
- [Hard *et al.*, 2018] Andrew Hard, Kanishka Rao, Rajiv Mathews, Swaroop Ramaswamy, Françoise Beaufays, Sean Augenstein, Hubert Eichner, Chloé Kiddon, and Daniel Ramage. Federated learning for mobile keyboard prediction, 2018.
- [Hardt *et al.*, 2016] Moritz Hardt, Eric Price, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information*

- Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [Kang *et al.*, 2019] Jiawen Kang, Zehui Xiong, Dusit Niyato, Han Yu, Ying-Chang Liang, and Dong In Kim. Incentive design for efficient federated learning in mobile networks: A contract theory approach. In *2019 IEEE VTS Asia Pacific Wireless Communications Symposium (AP-WCS)*, pages 1–5, 2019.
- [Khan *et al.*, 2020] Latif U. Khan, Shashi Raj Pandey, Nguyen H. Tran, Walid Saad, Zhu Han, Minh N. H. Nguyen, and Choong Seon Hong. Federated learning for edge networks: Resource optimization and incentive mechanism. *IEEE Communications Magazine*, 58(10):88–93, 2020.
- [Kholod *et al.*, 2021] Ivan Kholod, Evgeny Yanaki, Dmitry Fomichev, Evgeniy Shalugin, Evgenia Novikova, Evgeny Filippov, and Mats Nordlund. Open-source federated learning frameworks for iot: A comparative review and analysis. *Sensors*, 21(1), 2021.
- [Kotsogiannis and Schwager, 2006] Christos Kotsogiannis and Robert Schwager. On the incentives to experiment in federations. *Journal of Urban Economics*, 60(3):484–497, 2006.
- [Kusner *et al.*, 2017] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [Le *et al.*, 2020] Tra Huong Thi Le, Nguyen H. Tran, Yan Kyaw Tun, Zhu Han, and Choong Seon Hong. Auction based incentive design for efficient federated learning in cellular wireless networks. In *2020 IEEE Wireless Communications and Networking Conference (WCNC)*, pages 1–6, 2020.
- [Lecun *et al.*, 1998] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [Li *et al.*, 2019] Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. Fair resource allocation in federated learning, 2019.
- [Li *et al.*, 2021] Qinbin Li, Zeyi Wen, Zhaomin Wu, Sixu Hu, Naibo Wang, Yuan Li, Xu Liu, and Bingsheng He. A survey on federated learning systems: Vision, hype and reality for data privacy and protection. *IEEE Transactions on Knowledge and Data Engineering*, pages 1–1, 2021.
- [Lyu *et al.*, 2020] Lingjuan Lyu, Xinyi Xu, Qian Wang, and Han Yu. *Collaborative Fairness in Federated Learning*, pages 189–204. Springer International Publishing, Cham, 2020.
- [McMahan *et al.*, 2017] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguerre y Arcas. Communication-Efficient Learning of Deep Networks from Decentralized Data. In Aarti Singh and Jerry Zhu, editors, *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pages 1273–1282. PMLR, 20–22 Apr 2017.
- [Mohri *et al.*, 2019] Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 4615–4625. PMLR, 09–15 Jun 2019.
- [Mu *et al.*, 2018] Yashuang Mu, Xiaodong Liu, and Lidong Wang. A pearson’s correlation coefficient based decision tree and its parallel implementation. *Information Sciences*, 435:40–58, 2018.
- [Nishio *et al.*, 2020] Takayuki Nishio, Ryoichi Shinkuma, and Narayan B. Mandayam. Estimation of individual device contributions for incentivizing federated learning. In *2020 IEEE Globecom Workshops (GC Wkshps)*, pages 1–6, 2020.
- [Pang and Lee, 2005] Bo Pang and Lillian Lee. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. *CoRR*, abs/cs/0506075, 2005.
- [Sarikaya and Ercetin, 2020] Yunus Sarikaya and Ozgur Ercetin. Motivating workers in federated learning: A stackelberg game perspective. *IEEE Networking Letters*, 2(1):23–27, 2020.
- [Shi *et al.*, 2021] Yuxin Shi, Han Yu, and Cyril Leung. A survey of fairness-aware federated learning, 2021.
- [Shin, 2020] Donghee Shin. User perceptions of algorithmic decisions in the personalized ai system: perceptual evaluation of fairness, accountability, transparency, and explainability. *Journal of Broadcasting & Electronic Media*, 64(4):541–565, 2020.
- [Song *et al.*, 2019] Tianshu Song, Yongxin Tong, and Shuyue Wei. Profit allocation for federated learning. In *2019 IEEE International Conference on Big Data (Big Data)*, pages 2577–2586, 2019.
- [Stark *et al.*, 2015] Oded Stark, Marcin Jakubek, and Martyna Kobus. A bitter choice turned sweet: How acknowledging individuals’ concern at having a low relative income serves to align utilitarianism and egalitarianism. Discussion Papers 199354, University of Bonn, Center for Development Research (ZEF), 2015.
- [Truex *et al.*, 2020] Stacey Truex, Ling Liu, Ka-Ho Chow, Mehmet Emre GURSOY, and Wenqi Wei. Ldp-fed: Federated learning with local differential privacy. In *Proceedings of the Third ACM International Workshop on Edge Systems, Analytics and Networking*, EdgeSys ’20, page 61–66, New York, NY, USA, 2020. Association for Computing Machinery.
- [Tu *et al.*, 2021] Xuezhen Tu, Kun Zhu, Nguyen Cong Luong, Dusit Niyato, Yang Zhang, and Juan Li. Incentive mechanisms for federated learning: From economic and game theoretic perspective. *CoRR*, abs/2111.11850, 2021.

- [ur Rehman *et al.*, 2020] Muhammad Habib ur Rehman, Khaled Salah, Ernesto Damiani, and Davor Svetinovic. Towards blockchain-based reputation-aware federated learning. In *IEEE INFOCOM 2020 - IEEE Conference on Computer Communications Workshops (INFOCOM WK-SHPS)*, pages 183–188, 2020.
- [Wang, 2019] Guan Wang. Interpret federated learning with shapley values, 2019.
- [Xiao *et al.*, 2020] Guiliang Xiao, Mingjun Xiao, Guojia Gao, Sheng Zhang, Hui Zhao, and Xiang Zou. Incentive mechanism design for federated learning: A two-stage stackelberg game approach. In *2020 IEEE 26th International Conference on Parallel and Distributed Systems (ICPADS)*, pages 148–155, 2020.
- [Xinyi *et al.*, 2021] Xu Xinyi, Lingjuan Lyu, Xingjun Ma, Chenglin Miao, Chuan-Sheng Foo, and Bryan Kian Hsiang Low. Gradient driven rewards to guarantee fairness in collaborative machine learning. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
- [Xu and Lyu, 2020] Xinyi Xu and Lingjuan Lyu. Towards building a robust and fair federated learning system. *CoRR*, abs/2011.10464, 2020.
- [Yang *et al.*, 2019] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept and applications. *ACM Trans. Intell. Syst. Technol.*, 10(2), jan 2019.
- [Yu *et al.*, 2020a] Han Yu, Zelei Liu, Yang Liu, Tianjian Chen, Mingshu Cong, Xi Weng, Dusit Niyato, and Qiang Yang. *A Fairness-Aware Incentive Scheme for Federated Learning*, page 393–399. Association for Computing Machinery, New York, NY, USA, 2020.
- [Yu *et al.*, 2020b] Han Yu, Zelei Liu, Yang Liu, Tianjian Chen, Mingshu Cong, Xi Weng, Dusit Niyato, and Qiang Yang. A sustainable incentive scheme for federated learning. *IEEE Intelligent Systems*, 35(4):58–69, 2020.
- [Zafar *et al.*, 2017] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, Krishna P. Gummadi, and Adrian Weller. From parity to preference-based notions of fairness in classification, 2017.
- [Zeng *et al.*, 2020] Rongfei Zeng, Shixun Zhang, Jiaqi Wang, and Xiaowen Chu. Fmore: An incentive scheme of multi-dimensional auction for federated learning in mec. In *2020 IEEE 40th International Conference on Distributed Computing Systems (ICDCS)*, pages 278–288, 2020.
- [Zeng *et al.*, 2021] Rongfei Zeng, Chao Zeng, Xingwei Wang, Bo Li, and Xiaowen Chu. A comprehensive survey of incentive mechanism for federated learning, 2021.
- [Zhao *et al.*, 2018] Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated learning with non-iid data, 2018.

A Basic Formulas

Federated learning basic formulas. This paper adopts the standard model trained on the local dataset of the FL client: $\min\{F(\omega) := \sum_{i=1}^N \psi_i F_i(\omega)\}$. $F(\omega)$ is the gradient of the global model. $F_i(\omega)$ is the trained model of the local model. ψ_i is the weight of the i -th participant such that $\psi_i \geq 0$ and $\sum_{i=1}^N \psi_i = 1$. In the t -th round of updates, $\Delta\omega_i^{(t)} := \nabla F_i(\omega_{(t-1)})$ and $\Delta\omega := \sum_{i=1}^N \psi_i \Delta\omega_i^{(t)}$. Please refer to Table 1 for the meanings of the main symbols in this paper.

B Symbolic Meaning

The following table indicates the meaning of the symbols used in this paper.

Symbol	Meaning
N	number of participants
i	i -th participant
ψ	model weight
ω	model
C	total cost
S	the cost of temporary storage
u	compensation cost
γ	interest
R	good reputation set
A	reputation threshold
r	reputation
\tilde{r}	temporary reputation
σ	income distribution
α	contribution
t	number of rounds
T	total center revenue
Δ	upload
∇	download