

---

# Collusion-free Cross-feature Logistic Regression

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Federated learning is an emerging machine learning paradigm to learn a shared  
2 model from multiple independent participants while the data privacy is well pre-  
3 served. Much recent research is based on the *honest-but-curious* security model,  
4 which assumes that there's no collusion among the participants. In this paper, we  
5 take collusion into consideration and propose a framework for the cross-feature  
6 federated logistic regression on the *semi-honest-but-curious* security model. The  
7 proposed framework is secure enough even  $N - 1$  participants collude with each  
8 other. We further provide an implementation of the proposed framework based on  
9 the Ec-Elgamal encryption. Experimental results on 5 widely used datasets validate  
10 the effectiveness of our proposed framework.

## 11 1 Introduction

12 Federated learning (FL), recently proposed in [1], aims to learn a shared model from independent  
13 participants, where each participant has different data. The key is to ensure that the data only stays  
14 locally when training the federated model, thus data privacy is preserved. Specifically, if each  
15 participant holds disjoint feature subsets of the same samples, it is called cross-feature FL [2]<sup>1</sup>. The  
16 ability allows for privacy preserving in several widely used applications including but not limited to  
17 credit assessment, personalized recommendation and medical diagnosis.

18 Generally, there are three kinds of participants in cross-feature federation: an active party holding a  
19 feature subset and the labels, multiple passive parties holding only feature subsets and an authority  
20 party may be required to assist in the data exchange and the encryption during the training process.  
21 Note the authority party usually has privileges in the federation, e.g. generate the encryption key  
22 pairs [6], aggregating gradients [1]. It is difficult to assume that there will be no colluding among all  
23 the participants. For instance, the passive parties may collude to acquire the label information from  
24 the active party. It worth noting that since the authority party usually has a superior position in the  
25 federation, the risk of information leaking is especially high if the colluding participants include the  
26 authority party.

27 In this paper, we present a collusion-free version of cross-feature logistic regression. The intuition  
28 behind our method is as follows: on the one hand, we removed the superior authority party in  
29 the federation, on the other hand, we further make sure that if a ciphertext is to be decrypted, all  
30 participants need to participate in the decryption in order to complete the decryption process. Thus  
31 no ciphertext will be leaked even when collusion happens. The main contributions of this paper are  
32 listed as follows:

- 33 • We have analyzed the robustness to collusion of several popular cross-feature logistic  
34 regression methods.

---

<sup>1</sup>Cross-feature FL is also known as vertical FL in [6].

- 35 • We proposed a framework for collusion-free cross-feature logistic regression. The authority  
36 party is no longer necessary in the proposed framework. Furthermore, we proved that the  
37 method is robust to both the collusion among the active party and partial passive parties and  
38 the collusion among the passive parties.
- 39 • We showed that the Ec-Elgamal encryption is able to perform exchangeable distributed  
40 decryption and described an implementation of the proposed framework based on the  
41 Ec-Elgamal encryption.

## 42 2 Related Works

43 Much recent research of cross-feature logistic regression has been focused on optimizing the encryption  
44 methods or the logistic function approximations. Anti-collusion is merely mentioned. We try to  
45 supply the analysis of the collusion robustness of the modern popular cross-feature logistic regression  
46 methods in this section.

47 Recently, Hardy et al. [6] proposed a cross-feature logistic regression method. This method is  
48 based on the Taylor expansion of the logistic regression loss function and the additive homomorphic  
49 encryption. There are three participants: A, B and C, where A and B hold data and C generates the  
50 encryption key pairs. Suppose one of the intermediate variable  $y = w \times (x_A + x_B)$  is transmitted to  
51 C, in which  $x_A$  and  $x_B$  are encrypted and hold by A and B separately. When A colludes with C, then  
52 A could send the encrypted  $x_A$  to C, and C can decrypt and build equations to recover features  $x_B$   
53 belonging to B. Our method also utilized the Taylor expansion of the logistic regression loss function,  
54 but will not leak confidential data even if the participants collude.

55 Kim et al. [3] used an homomorphic encryption for approximate arithmetic and developed a least  
56 square approximation of the logistic function in privacy preserving logistic regression. The private  
57 data at multiple institutions are encrypted and sent to a public server, where the approximated logistic  
58 regression is performed totally in the encrypted format. This method enjoys high security, high  
59 accuracy and very low transmission expense when no collusion occurs. However, if we take collusion  
60 into consideration, all the confidential data may be leaked when one of the participants who holds the  
61 private key of encryption collude with the public server.

62 Mugunthan et al. [4] proposed a mechanism based on secure Multi-Party Computation and Differential  
63 Privacy. Aggregating logistic regression weights was taken as an example, where the collusion  
64 problem was resolved. Specifically, to achieve the secure Differential Privacy, each party receives  
65 its noise from the other parties, where the honest party would have contributed to the noise of the  
66 other malicious parties. The core idea of this method is very alike to our method, but we take another  
67 approach that instead of generate random noises in a distributed manner, we generate private keys  
68 of homomorphic encryption distributively. Additionally, this method aims to solve the problem of  
69 securely aggregating locally trained models instead of solving the logistic regression problem.

## 70 3 Preliminaries

### 71 3.1 The Security Model

72 In this paper, we assume that all the participants are *semi-honest-but-curious*:

- 73 • All the participants follow the protocol;
- 74 • They will infer as much information as they could from the information received;
- 75 • They may collude, leaking local and received information to the other participants.

76 Note this security model is different to the conventional *honest-but-curious* model in that a part of the  
77 participants may share the information to infer the confidential of the other participants.

### 78 3.2 Revisit Cross-feature Logistic Regression via Taylor Expansion

79 Logistic regression is a classic machine learning algorithm that can be used to distinguish two classes.  
80 It can calculate the posterior probability for each sample based on the input. Denote the input features

81 as  $\mathbf{x} \in \mathbb{R}^n$ , the model parameter as  $\boldsymbol{\theta} \in \mathbb{R}^n$  and the data label as  $\mathbf{y} \in \{-1, 1\}$ , the average logistic  
 82 loss is approximated via the Taylor series expansion as follows [6]:

$$L(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i \boldsymbol{\theta}^T \mathbf{x}_i}) \approx \log 2 - \frac{1}{2} \mathbf{y} \boldsymbol{\theta}^T \mathbf{x} + \frac{1}{8} (\boldsymbol{\theta}^T \mathbf{x})^2 \quad (1)$$

83 Note Equation 1 is a polynomial, making the use of the additive homomorphic encryption available.  
 84 Based on Equation 1, the gradient is calculated as follows:

$$\nabla L(\boldsymbol{\theta}) = \left( \frac{1}{4} \boldsymbol{\theta}^T \mathbf{x} - \frac{1}{2} \mathbf{y} \right) \mathbf{x} \quad (2)$$

85 Specifically, under the cross-feature setting, the calculation of the gradient can be distributed to each  
 86 participant. Denote  $k \in (A, B, C, \dots)$  as the indicator of a participant.  $\mathbf{x}_k$  is the features held by each  
 87 participant and  $\boldsymbol{\theta}_k$  is the distributed model parameters. Thus the input features  $\mathbf{x} = (\mathbf{x}_A | \mathbf{x}_B | \dots)$ ,  
 88 the model parameters  $\boldsymbol{\theta} = (\boldsymbol{\theta}_A | \boldsymbol{\theta}_B | \dots)$ , and the gradients  $\nabla L(\boldsymbol{\theta}) = (\nabla L(\boldsymbol{\theta}_A) | \nabla L(\boldsymbol{\theta}_B) | \dots)$ . The  
 89 gradient of loss on each local model parameters is calculated as follows:

$$\nabla L(\boldsymbol{\theta}_k) = \left( \sum_k \frac{1}{4} \boldsymbol{\theta}_k^T \mathbf{x}_k - \frac{1}{2} \mathbf{y} \right) \mathbf{x}_k \quad (3)$$

90 In this paper, we also build the method based on the Taylor expansion of the logistic loss as in [6].

### 91 3.3 The Ec-Elgamal Encryption

92 Koblitz et al. introduced the analog of Elgamal cryptosystem based on the elliptic curve cryptography  
 93 [8], that is the Ec-Elgamal encryption. We use the definition of elliptic curve on  $F_p$  proposed in [9].  
 94 Given the elliptic curve  $E(a, b)$  with a base point  $P$ , an integer  $s$  can be randomly generated as the  
 95 private key, and the product of  $s$  and  $P$  is the public key, namely  $Q$ . At encryption, the plaintext  $m$   
 96 is mapped to the point  $M$  on the elliptic curve and random number is chosen to generate a pair of  
 97 cipher text points  $(R, S)$  as depicted in Equation 4. The decryption from  $(R, S)$  is shown in Equation  
 98 5. It can be seen that the encryption and decryption results are the same point on the elliptic curve.

$$\text{Encryption} : R = rP, S = M + rQ \quad (4)$$

$$\text{Decryption} : M = S - kR \quad (5)$$

99 Note in Equation 4, the encryption of the same plaintext point will generate different ciphertext points  
 100 using random  $r$ . Therefore, EC-Elgamal is probabilistic.

101 Furthermore, Ec-Elgamal is also additive homomorphic: Assume there are two plaintext points  $M_1$   
 102 and  $M_2$ , with cipher text as  $(R_1, S_1) = (r_1P, M_1 + r_1Q)$  and  $(R_2, S_2) = (r_2P, M_2 + r_2Q)$ . The  
 103 sum of the ciphertexts:

$$(R_3, S_3) = (R_1 + R_2, S_1 + S_2) = ((r_1 + r_2)P, M_1 + M_2 + (r_1 + r_2)Q) \quad (6)$$

104 Decrypting  $(R_3, S_3)$  as depicted in Equation 5, we get

$$M_3 = M_1 + M_2 + (r_1 + r_2)Q - s(r_1 + r_2)P = M_1 + M_2 \quad (7)$$

105 Thus, Ec-Elgamal is additive homomorphic.

Table 1: One training iteration of our approach

	Participant A	Participant B	Participant C
Step1	Generate the encryption key pair $(Q, (s_A, s_B, s_C))$		
Step2	initialize model parameters $\theta_A$	initialize model parameters $\theta_B$	initialize model parameters $\theta_C$
Step3	compute $u_A$ , encrypt to $[u_A]_{ABC}$	compute $u_B$ , encrypt to $[u_B]_{ABC}$	compute $u_C$ , encrypt to $[u_C]_{ABC}$
Step4	Calculate $[w]_{ABC} = \sum_k [u_k]_{ABC}$ and $[\nabla L(\theta_k)]_{ABC} = [w]_{ABC} \mathbf{x}_k$		
Step5	send $[\nabla L(\theta_A)]_{ABC}$ to B, get $[\nabla L(\theta_C)]_{ABC}$ and partially decrypt it to $[\nabla L(\theta_C)]_{BC}$	send $[\nabla L(\theta_B)]_{ABC}$ to C, get $[\nabla L(\theta_A)]_{ABC}$ and partially decrypt it to $[\nabla L(\theta_A)]_{AC}$	send $[\nabla L(\theta_C)]_{ABC}$ to A, get $[\nabla L(\theta_B)]_{ABC}$ and partially decrypt it to $[\nabla L(\theta_B)]_{AB}$
Step6	send $[\nabla L(\theta_C)]_{BC}$ to B, get $[\nabla L(\theta_B)]_{AB}$ and partially decrypt it to $[\nabla L(\theta_B)]_B$	send $[\nabla L(\theta_A)]_{AC}$ to C, get $[\nabla L(\theta_C)]_{BC}$ and partially decrypt it to $[\nabla L(\theta_C)]_C$	send $[\nabla L(\theta_B)]_{AB}$ to A, get $[\nabla L(\theta_A)]_{AC}$ and partially decrypt it to $[\nabla L(\theta_A)]_A$
Step7	send $[\nabla L(\theta_B)]_B$ to B, get $[\nabla L(\theta_A)]_A$ and finally decrypt it to $\nabla L(\theta_A)$	send $[\nabla L(\theta_C)]_C$ to C, get $[\nabla L(\theta_B)]_B$ and finally decrypt it to $\nabla L(\theta_B)$	send $[\nabla L(\theta_A)]_A$ to A, get $[\nabla L(\theta_C)]_C$ and finally decrypt it to $\nabla L(\theta_C)$
Step8	update $\theta_A$	update $\theta_B$	update $\theta_C$
obtain	$\theta_A$	$\theta_B$	$\theta_C$

## 106 4 Collusion-Free Cross-Feature Logistic Regression

107 Anti-collusion has been widely discussed in the secure Multi-Party Computation community. Com-  
108 mon solutions to anti-collusion is based on sharing any intermediate results used for the calculation,  
109 and all or a majority of the participants are necessary to get the final result. In fact, the whole logistic  
110 regression training process or an iteration can be regarded as a Multi-Party Computation protocol.  
111 However, exact secret sharing requires heavy inter-participants communication. In this paper, we aim  
112 to imitate the anti-collusion solutions via homomorphic encryption in each logistic regression training  
113 iteration. The core idea is to ensure that all the intermediate results and the final calculated gradient  
114 can only be decrypted if all the participants are involved, and regardless of the order of decryption.  
115 We denote this property as *exchangeable distributed decryption*.

116 We detail the proposed framework based on the exchangeable distributed decryption in Section 4.1  
117 and an implementation based on the Ec-Elgamal encryption in Section 4.2.

### 118 4.1 The Proposed Framework

119 An encryption method is said to has the property of *exchangeable distributed decryption* if the private  
120 key is distributed to all the participants. Only when all the participants participate in the decryption  
121 using its own private key, a confidential is able to be decrypted.

122 The proposed framework is depicted in Table 1. Without losing generality, we assume that there are  
123 three participants, where participant A holds both a subset of features  $\mathbf{x}_k$  and the labels  $\mathbf{y}$ , while  
124 participant B and C only hold subsets of features  $\mathbf{x}_k$ . This framework can be easily extended to  $N$   
125 ( $N \geq 2$ ) participants. Note the confidential information includes but not limited to the input features,  
126 the groundtruth labels and the model parameters. The framework is decentralized, i.e. we do not need  
127 an authority party to assist the training in this framework. Detailed steps are shown as following:

- 128 • Step1: Generate the encryption key pair  $(Q, (s_A, s_B, s_C))$ , where  $Q$  is the public key,  
129  $k_i, i \in (A, B, C)$  is the private keys of each participant.
- 130 • Step2: All participants initialize the model parameters  $\theta_k$ .
- 131 • Step3: Each party calculates  $u_k = \frac{1}{4}\theta_k^T \mathbf{x}_k - \alpha \mathbf{y}$ , where  $\alpha$  is 0 except on Participant A it is  
132  $\frac{1}{2}$ . Then  $u_k$  is encrypted by  $Q$  to  $[u_k]_{ABC}$ .
- 133 • Step4: Each participant get  $[w]_{ABC} = \sum_k [u_k]_{ABC}$  and calculate gradient  $[\nabla L(\theta_k)]_{ABC} =$   
134  $[w]_{ABC} \mathbf{x}_k$
- 135 • Step5-7: Steps 5-7 aim to decrypt the gradients. At each step,  $[\nabla L(\theta_k)]$  will pass to the right  
136 neighbor and performs the partial decryption by it. After looping  $N$  times, each participant  
137  $k$  will get partially decrypted gradient  $[\nabla L(\theta_k)]_k$  and can decrypt it to the final plaintext by  
138 its own private key  $s_k$ .
- 139 • Step8: All the participants use the plaintext gradient to update their model parameters.

140 The predicting procedure is similar to the training procedure. Each participant calculates  $\theta_k^T \mathbf{x}_k$  and  
141 encrypt it by  $Q$  to  $[\theta_k^T \mathbf{x}_k]_{ABC}$ . Then  $\sum_k [\theta_k^T \mathbf{x}_k]_{ABC}$  is calculated. Assume Participant A needs the  
142 prediction result, the encrypted sum go through all the participants except A. Thus finally, A gets  
143  $\sum_k [\theta_k^T \mathbf{x}_k]_A$  and decrypt it to calculate the prediction result.

#### 144 4.2 Implementation via Ec-Elgamal

145 The key to implement the framework is to find a homomorphic encryption method that can perform  
146 the exchangeable distributed decryption. In this section, we provide a simple implementation based  
147 on the composition of the Ec-Elgamal encryption as follows:

- 148 • Step1: Each participant generates a pair of local public and private keys  $(Q_k, s_k)$  via  
149 conventional Ec-Elgamal as described in Section 3.3.
- 150 • Step2: Calculate the public key  $Q = \sum_k Q_k$

151 These two steps build  $(Q, (s_A, s_B, \dots))$ , where  $Q$  is publicly available. At meanwhile,  $s_k$  is  
152 distributed at each participant and is kept secretly. Suppose plaintext point  $M$  is encrypted to  $(R, S)$   
153 as Equation 4, then

$$S = M + rQ = M + r \sum_k Q_k \quad (8)$$

154 All  $s_k$  are necessary in decryption, and the order of applying  $s_k$  does not affect the decryption:

$$M = M + r \sum_k Q_k - r \sum_k s_k P = S - rP \sum_k s_k \quad (9)$$

155 where  $P$  is the base point and  $Q_k = s_k P$ . Thus the composition of Ec-Elgamal is able to accomplish  
156 exchangeable distributed decryption. Furthermore, obviously, encryption using  $Q$  is also additive  
157 homomorphic.

#### 158 4.3 Transmission and Security Analysis

159 Assume there are  $N$  participants. As shown in Table 1, in one training iteration, the inter-participant  
160 transmission occurs in 3 phases: 1. the generation of the encryption keys, 2. the inter-participant  
161 sum to calculate the encrypted intermediate variables  $[u_k]_{ABC}$  and 3. the distributed decryption of  
162 the encrypted gradients. Specifically, the information transmitted in phase 1 is the local public keys,  
163 which are not confidential. At phase 2, the transmitted information is confidential but encrypted  
164 by the public key  $Q$ . Since the decryption needs all the participants, so the confidential will not be  
165 leaked even when  $N - 1$  participants collude. Phase 3 includes  $N$  rounds of transmissions to decrypt  
166 the gradients. Note in the proposed framework, each participant  $k$  will never send out  $[\nabla L(\theta_k)]$  that  
167 has been decrypted by its own private key  $s_k$ . Therefore,  $\nabla L(\theta_k)$  will never be leaked even when  
168 the other  $N - 1$  participants collude.

Table 2: A brief description of datasets

Dataset	Rows	Cols	Is Binary
breast cancer	569	30	Yes
digits	1797	64	No
credit	30000	24	Yes
MNIST	70000	784	No
covertypes	581012	54	No

Table 3: Numerical results on public datasets

Dataset	Time	Batch Size	ACC		AUC	
			Our LR	Original LR	Our LR	Original LR
breast cancer	39m	128	0.93	0.98	0.98	0.99
digits	3h35m	128	0.83	0.92	0.91	0.97
credit	49m	4096	0.79	0.77	0.71	0.73
MNIST	32h42m	4096	0.85	0.9	0.93	0.96
covertypes	29h30m	4096	0.71	0.74	0.78	0.81

## 169 5 Experiments

### 170 5.1 Settings

171 **Benchmark datasets:** We evaluate the collusion-free cross-feature logistic regression on 5 widely  
 172 used datasets: breast cancer, digits, credit, MNIST and covertypes. Table 2 describes the overall  
 173 characteristics of these datasets, including their size and whether the labels are binarized. For datasets  
 174 with multi-category labels, we converted the labels into binary classification via odd vs. even. All  
 175 data features are standardized. To simulate the cross-feature scenario, we divided the features of  
 176 the data to 5 participants approximately evenly, and then ran the program on 5 independent servers.  
 177 For the MNIST dataset, there are already separated train sets and test sets. For the rest datasets, we  
 178 randomly split 20% from the total samples for testing to evaluate the performance.

179 **Evaluation metrics:** We utilize the classification accuracy (ACC) and area under the ROC curve  
 180 (AUC) to evaluate the proposed method.

181 **Hardware settings:** We use 5 servers to play 5 participants in all the experiments. Each server is  
 182 equipped with a Intel(R) Xeon(R) Gold 5218 CPU @ 2.30GHz and 125 GB memory.

### 183 5.2 Results

184 In the experiment, since we are not interested in the converging time, all the number of epochs are  
 185 fixed to 2 and the learning rate to 0.01. For the breast cancer and the digits datasets, we set the  
 186 mini-batch size to 128. For the credit, the MNIST and the covertypes datasets, the batch size is set  
 187 to 4096. The proposed method is compared with the original logistic regression implemented by  
 188 Sklearn [10].

189 Table 3 shows the experimental results of the both the proposed method and the original logistic  
 190 regression, including the training time and the performances on each dataset. As can be seen that the  
 191 training time of proposed method is linearly and positively correlated with the size of the data set  
 192 and the number of features owned by each participant. The prediction accuracy is reasonable even  
 193 after only 2 epochs. The performance is worse than the original LR, which is within the expectation,  
 194 because of the Taylor expansion of the logistic loss function.

## 195 6 Conclusions

196 In this paper, we have proposed a framework for collusion-free cross-feature logistic regression and  
 197 provided an implementation based on the composition of the Ec-Elgamal encryption. The privacy  
 198 can be protected during both training and prediction even when  $N - 1$  participants collude under the

199 *semi-honest-but-curious* security model. Experiments show that the proposed method have achieved  
200 reasonable performance on 5 classification datasets and the training time grows linearly with the size  
201 the input dataset.

## 202 **References**

- 203 [1] McMahan, H. B., Moore, E., Ramage, D., & y Arcas, B. A. (2016). Federated learning of deep  
204 networks using model averaging. CoRR abs/1602.05629 (2016). arXiv preprint arXiv:1602.05629.
- 205 [2] Li, H., Meng, D., & Li, X. (2020). Knowledge Federation: Hierarchy and Unification. arXiv  
206 preprint arXiv:2002.01647.
- 207 [3] Kim, M., Song, Y., Wang, S., Xia, Y., & Jiang, X. (2018). Secure logistic regression based on  
208 homomorphic encryption: Design and evaluation. JMIR medical informatics, 6(2), e19.
- 209 [4] Mugunthan, V., Polychroniadou, A., Byrd, D., & Balch, T. H. (2019) SMPAI: Secure Multi-Party  
210 Computation for Federated Learning.
- 211 [5] Xu, R., Baracaldo, N., Zhou, Y., Anwar, A., & Ludwig, H. (2019, November). Hybridalpha:  
212 An efficient approach for privacy-preserving federated learning. In Proceedings of the 12th ACM  
213 Workshop on Artificial Intelligence and Security (pp. 13-23).
- 214 [6] Hardy, S., Henecka, W., Ivey-Law, H., Nock, R., Patrini, G., Smith, G., & Thorne, B. (2017). Pri-  
215 vate federated learning on vertically partitioned data via entity resolution and additively homomorphic  
216 encryption. arXiv preprint arXiv:1711.10677.
- 217 [7] Aono, Y., Hayashi, T., Trieu Phong, L., & Wang, L. (2016). Scalable and secure logistic regression  
218 via homomorphic encryption. In Proceedings of the Sixth ACM Conference on Data and Application  
219 Security and Privacy (pp. 142-144).
- 220 [8] Koblitz, N. (1994). A course in number theory and cryptography (Vol. 114). Springer Science &  
221 Business Media.
- 222 [9] Koblitz, N., Menezes, A., & Vanstone, S. (2000). The state of elliptic curve cryptography. Designs,  
223 codes and cryptography, 19(2-3), 173-193.
- 224 [10] Pedregosa, Fabian and Varoquaux, Gaël and Gramfort, Alexandre and Michel, Vincent and  
225 Thirion, Bertrand and Grisel, Olivier and Blondel, Mathieu and Prettenhofer, Peter and Weiss, Ron  
226 and Dubourg, & Vincent and others. (2011). Journal of machine learning research, 12(Oct), (pp.  
227 2825-2830).