EFFICIENT NEWTON-TYPE FEDERATED LEARNING WITH NON-IID DATA

Anonymous authors

Paper under double-blind review

ABSTRACT

The mainstream federated learning algorithms only communicate the first-order information across the local devices, i.e., FedAvg and FedProx. However, only using first-order information, these methods are often inefficient and the impact of heterogeneous data is yet not precisely understood. This paper proposes an efficient federated Newton method (FedNewton), by sharing both first-order and second-order knowledge over heterogeneous data. In general kernel ridge regression setting, we derive the generalization bounds for FedNewton and obtain the minimax-optimal learning rates. For the first time, our results analytically quantify the impact of the number of local examples, the data heterogeneity and the model heterogeneity. Moreover, as long as the local sample size is not too small and data heterogeneity is moderate, the federated error in FedNewton decreases exponentially in terms of iterations. Extensive experimental results further validate our theoretical findings and illustrate the advantages of FedNewton over the first-order methods.

023 024 025

026 027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Owing to the great potential in privacy preservation and in lowering the computational costs, federated learning (FL) McMahan et al. (2017); Li et al. (2020a); Zhang et al. (2021) becomes a promising framework in processing large-scale tasks. However, federated learning is facing massive challenges from the heterogeneous data Zhao et al. (2018); Zhou et al. (2023); Ye et al. (2023), including both the data heterogeneity and the model heterogeneity. The data heterogeneity comes from that inputs across devices are usually sampled from heterogeneous distributions, while the model heterogeneity measures the response shift due to inconsistency between local models and the global model.

First-order approaches, including FedAvg McMahan et al. (2017) and FedProx Li et al. (2020a). share the first-order information rather than the data across devices and tolerate the heterogeneity in federated learning, while Newton-type FL methods Ghosh et al. (2020); Gupta et al. (2021); Sa-037 faryan et al. (2022); Islamov et al. (2023); Liu et al. (2023); Dal Fabbro et al. (2024); Li et al. (2023) utilized second-order information for updating federated model. To the best of our knowledge, most of existing learning guarantees for FL methods are derived in the context of optimization and fo-040 cused on in-sample predictive errors only, i.e., the convergence analysis (optimization) of first-order 041 FL [Li et al.] (2020b); [Karimireddy et al.] (2020); [Pathak & Wainwright] (2020); [Glasgow et al.] (2022) 042 and Newton-type FL Ghosh et al. (2020); Safaryan et al. (2022); Qian et al. (2022). However, 043 beyond the optimization, the generalization guarantees (out-sample predictive performance) are of 044 great practical and theoretical interests for FL. Despite recent efforts and progress on the generalization for first-order algorithms Mohri et al. (2019); Yagli et al. (2020); Su et al. (2021); Yuan et al. (2022), the generalization guarantees for Newton-type FL algorithms remain elusive, especially on 046 heterogeneous data and localized models. Therefore, a challenging problem in FL is how to quantify 047 the impact of heterogeneity from the generalization perspective? 048

In this paper, motivated by sharing second-order information, we propose a second-order federated
 optimization method, named FedNewton. It approximates the global predictor on the entire data
 by utilizing the global gradient and local Hessians, improving the predictive accuracy in an efficient
 communications framework. We then study the statistical properties of FedNewton, and derive the
 generalization bounds with the minimax optimal rates. We conclude with experiments on simulated
 data and publicly available tasks that complement our theoretical results, exhibiting the computa-

tional and statistical benefits of our approach. Due to the length limit, we leave the experiment part in the appendix. We summarize our contributions as below:

1) On the algorithmic front. We propose a fast second-order federated learning algorithm, which improves the approximation of the centralized model while only requiring similar computational and communication costs as the first-order methods. The convergence of FedNewton is exponentially fast and a few communications, for example, $t \le 2$, can approximate the global model well.

2) On the statistical front. To our best knowledge, in presence of both data heterogeneity and
 model heterogeneity, we present the optimal generalization guarantees for the first time. Our results
 further analytically quantify the impacts of the local sample size, the data heterogeneity, and the
 model heterogeneity. Especially, the federated error decreases exponentially fast in benign cases,
 i.e., a sufficient number of local examples and moderate data heterogeneity.

066 067

068

081

082

084

085 086 087

090 091

092

093 094

096

097 098 099

100 101

102 103

104

105

2 PROBLEM SETUP

In a standard framework of federated learning, there is a global parameter server and m local computational clients. On the *j*-th local machine $\forall j \in [m]$, the local data $\mathfrak{D}_j = \{(x_{ij}, y_{ij})\}_{i=1}^{|\mathfrak{D}_j|}$ is drawn from a local distribution ρ_j on the joint space $\mathcal{X} \times \mathcal{Y}$. The total sample $\mathcal{D} = \bigcup_{j=1}^m \mathfrak{D}_j$ is the disjoint union of local data and corresponds to a global distribution ρ . For any local devices $j, k \in [m]$ and $j \neq k$, data distributions are identical $\rho_j = \rho_k = \rho$ in the homogeneous setting (iid data), while data distributions are distinct $\rho_j \neq \rho_k$ in the heterogeneous case (non-iid data).

We base our analysis on the standard non-parametric regression setup and assume that the target solution f^* belongs to a reproducing kernel Hilbert space (RKHS) induced by a Mercer kernel $K: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$. Mercer's theorem guarantees the kernel function admits an implicit feature mapping $K(\boldsymbol{x}, \boldsymbol{x}') = \langle \phi(\boldsymbol{x}), \phi(\boldsymbol{x}') \rangle_K$ and the norm by $\|\cdot\|_K$. The predictor can be stated as $f_{\mathcal{D},\lambda}(\boldsymbol{x}) = \langle \boldsymbol{w}_{\mathcal{D},\lambda}, \phi(\boldsymbol{x}) \rangle$ where $\boldsymbol{w}_{\mathcal{D},\lambda}$ minimizes the objective on the entire data \mathcal{D}

$$\underset{\boldsymbol{w}\in\mathcal{H}_{K}}{\operatorname{arg\,min}}\left\{\frac{1}{2|\mathcal{D}|}\sum_{i=1}^{|\mathcal{D}|}\left(f(\boldsymbol{x}_{i})-y_{i}\right)^{2}+\frac{\lambda}{2}\|\boldsymbol{w}\|_{K}^{2}\right\},\tag{1}$$

where $(x_i, y_i) \in D$, and $\lambda > 0$ is the regularity parameter. The above regression problem, known as Kernel Ridge Regression (KRR), admits a closed-form solution

$$\boldsymbol{w}_{\mathcal{D},\lambda} = (\boldsymbol{\Phi}_{\mathcal{D}}^{\top} \boldsymbol{\Phi}_{\mathcal{D}} + \lambda I)^{-1} \boldsymbol{\Phi}_{\mathcal{D}}^{\top} \boldsymbol{y}_{\mathcal{D}},$$
(2)

where $\mathbf{\Phi}_{\mathcal{D}} = \frac{1}{\sqrt{|\mathcal{D}|}} \left[\phi(\boldsymbol{x}_1), \cdots, \phi(\boldsymbol{x}_{|\mathcal{D}|}) \right]^T \in \mathbb{R}^{|\mathcal{D}|} \times \mathcal{H}_K$ are feature mappings on the training set \mathcal{D} and $\boldsymbol{y}_{\mathcal{D}} = \frac{1}{\sqrt{|\mathcal{D}|}} \left(y_1, \cdots, y_{|\mathcal{D}|} \right)^\top$ are the corresponding labels.

By averaging the local models, the simplest federated method only communicates once, known as Distributed Kernel Ridge Regression (DKRR) with the closed-form solution

$$\bar{\boldsymbol{w}}_{\mathcal{D},\lambda} = \sum_{j=1}^{m} p_j (\boldsymbol{\Phi}_{\mathfrak{D}_j}^{\top} \boldsymbol{\Phi}_{\mathfrak{D}_j} + \lambda I)^{-1} \boldsymbol{\Phi}_{\mathfrak{D}_j}^{\top} \boldsymbol{y}_{\mathfrak{D}_j},$$

where p_j is the weight of the *j*-th local model, which is usually set $p_j = |\mathfrak{D}_j|/|\mathcal{D}|$. Note that, $\Phi_{\mathfrak{D}_j} = \frac{1}{\sqrt{|\mathfrak{D}_j|}} \left[\phi(\boldsymbol{x}_1), \cdots, \phi(\boldsymbol{x}_{|\mathfrak{D}|_j})\right]^T \in \mathbb{R}^{|\mathfrak{D}_j|} \times \mathcal{H}_K$ are local feature mappings and $\boldsymbol{y}_{\mathfrak{D}_j} = \frac{1}{\sqrt{|\mathfrak{D}_j|}} \left(y_1, \cdots, y_{|\mathfrak{D}_j|}\right)^{\top}$ are labels on the *j*-th local train set $\mathfrak{D}_j = \left\{(\boldsymbol{x}_{ij}, y_{ij})\right\}_{i=1}^{|\mathfrak{D}_j|}, \quad \forall j \in [m].$

The solution of KRR equation 2 can be rewritten in the Newton's method form $w_{D,\lambda} = w - H_{D,\lambda}^{-1} a_{D,\lambda}.$

$$\boldsymbol{w}_{\mathcal{D},\lambda} = \boldsymbol{w} - \boldsymbol{H}_{\mathcal{D},\lambda}^{-1} \boldsymbol{g}_{\mathcal{D},\lambda}.$$
(3)

where the gradient and Hessian matrix are defined as

106

$$\boldsymbol{g}_{\mathcal{D},\lambda} := (\boldsymbol{\Phi}_{\mathcal{D}}^{\top} \boldsymbol{\Phi}_{\mathcal{D}} + \lambda I) \boldsymbol{w} - \boldsymbol{\Phi}_{\mathcal{D}}^{\top} \boldsymbol{y}_{\mathcal{D}}$$

107
 $\boldsymbol{H}_{\mathcal{D},\lambda} := (\boldsymbol{\Phi}_{\mathcal{D}}^{\top} \boldsymbol{\Phi}_{\mathcal{D}} + \lambda I).$

Alg	orithm 1 Federated Learning with Newton Method (FedNewton)
Inp	ut: Local training data subset $\mathfrak{D}_j, \forall j \in [m]$. Feature mapping $\phi : \mathcal{X} \to \mathbb{R}^M$.
	Put. The global estimator $w_{\mathcal{D},\lambda}$.
1:	Local machines: Compute feature mapping $\Psi_{\mathfrak{D}_j}$, $\mathbf{h}_{\mathfrak{D}_j,\lambda} = (\Psi_{\mathfrak{D}_j}\Psi_{\mathfrak{D}_j} + \lambda I)$, $\mathbf{h}_{\mathfrak{D}_j,\lambda}$ and \mathbf{E}_j^{\top}
_	$\Psi_{\mathfrak{D}_j} y_{\mathfrak{D}_j}$ for any $j \in [m]$.
2:	Local machines: Initialize the local estimators by $w_{\mathfrak{D}_j,\lambda}^{\mathfrak{G}} = H_{\mathfrak{D}_j,\lambda}^{\mathfrak{G}} \Phi_{\mathfrak{D}_j}^{\mathfrak{G}} y_{\mathfrak{D}_j}$ and upload them
2.	to the global server (\uparrow).
3:	Global server: initialize the solution by $w_{\mathcal{D},\lambda} = \sum_{j=1} p_j w_{\hat{\mathfrak{D}}_j,\lambda}$, and send it to the local nodes
4:	(\downarrow). for $t = 1$ to T do
5:	Local machines: Compute local gradients $g_{\mathfrak{D}_i,\lambda}^{t-1} = H_{\mathfrak{D}_i,\lambda} \bar{w}_{\mathcal{D},\lambda}^{t-1} - \Phi_{\mathfrak{D}_i}^{\top} y_{\mathfrak{D}_i}$ and upload
	them to global server (\uparrow) .
6:	Global server: Compute the global gradient $g_{\mathfrak{D},\lambda}^{t-1} = \sum_{j=1}^{m} p_j g_{\mathfrak{D}_j,\lambda}^{t-1}$ and send it to local
	nodes (\downarrow) .
7:	Local machines: Compute the local updates $H_{\mathfrak{D}_j,\lambda}^{-1} g_{\mathfrak{D},\lambda}^{t-1}$ and upload it to the global server
	(↑).
8:	Global server: Update the global estimator $\bar{w}_{\mathcal{D},\lambda}^t = \bar{w}_{\mathcal{D},\lambda}^{t-1} - \sum_{j=1}^m p_j H_{\mathfrak{D}_j,\lambda}^{-1} g_{\mathfrak{D},\lambda}^{t-1}$ and
0	communicate it to local machines (\downarrow) .
9:	
Froi	n equation 3, the global gradient $g_{\mathcal{D},\lambda}$ and Hessian $H_{\mathcal{D},\lambda}$ is the key to achieving the cen-
trali	zed model $\mathbf{w}_{\mathcal{D},\lambda}$. Note that, since the fact $\mathbf{\Phi}_{\mathcal{D}}^{\top}\mathbf{\Phi}_{\mathcal{D}} = \sum_{i=1}^{m} p_{i}\mathbf{\Phi}_{\mathfrak{D}_{i}}^{\top}\mathbf{\Phi}_{\mathfrak{D}_{i}}$ for data partition
\mathcal{D} =	$= \bigcup_{i=1}^m \mathfrak{D}_i$, one can easily obtain the following property for the global gradient and global
Hes	sian.
Pro	position 1 (Partitonability). If the loss is squared loss, the global gradient and Hessian matrix
con.	sist of the local ones, i.e. $\boldsymbol{g}_{\mathcal{D},\lambda} = \sum_{j=1}^{m} p_j \boldsymbol{g}_{\mathfrak{D}_j,\lambda}$ and $\boldsymbol{H}_{\mathcal{D},\lambda} = \sum_{j=1}^{m} p_j \boldsymbol{H}_{\mathfrak{D}_j,\lambda}$.
Ren	nark 1 (Computation of local inverse Hessian). The compute of the inverse of local Hessians
$H_{\mathfrak{D}}^{-}$	$\mathcal{D}_{j,\lambda}^{1}$ is time consuming $\mathcal{O}(\mathfrak{D}_{j} M^{2}+M^{3})$, which is a common problem in second-order optimiza-
tion	Botton at al (2018) There are many classic work to reduce the time complexity of the inverse of

¹¹ $\mathfrak{D}_{j,\lambda}$ is time consuming $\mathcal{O}(|\mathcal{D}_j||\mathcal{W}| + |\mathcal{W}|)$, which is a common problem in second-order optimization Bottou et al. (2018). There are many classic work to reduce the time complexity of the inverse of Hessian, i.e. BFGS Broyden (1970), L-BFGS Liu & Nocedal (1989), inexact Newton Dembo et al. (1982), Gauss-Newton Schraudolph (2002) and Newton sketch Pilanci & Wainwright (2017). Those techniques can be used to improve the efficiency of FedNewton, but it is beyond the scope of this paper. We focus on theoretical novelties and leave further computational improvements in the future.

Remark 2 (Feature mapping instead of kernel methods). Without loss of generality, we assume the feature mappings are finite dimensional $\phi : \mathcal{X} \to \mathbb{R}^M$, which covers a wide range of generalized linear models, for example neural networks Neal (1995); Jacot et al. (2018), kernel methods Vapnik (2000), random features Rahimi & Recht (2007); Le et al. (2013); Yang et al. (2014), and random sketching Woodruff et al. (2014); Yang et al. (2017).

149 150

151 152

3 FEDERATED LEARNING WITH NEWTON METHOD

Motivated by recent gradient-based distributed learning Wang et al. (2018); Lin et al. (2020), we propose a Newton-type federated learning method to quantity the impact of data heterogeneity and model heterogeneity. Using Proposition 1, the exact Federated Newton's method communicate local Hessians $H_{\mathfrak{D}_{j},\lambda}$ for computing the global Hessian matrix equation 3 whose the communication complexity is $O(M^2)$, which is infeasible in federated learning. To reduce communication costs, we propose FedNewton that approximates the Newton's updates with the global gradient and local Hessian matrices, such that

160

161

$$\boldsymbol{H}_{\mathcal{D},\lambda}^{-1}\boldsymbol{g}_{\mathcal{D},\lambda} \approx \sum_{j=1}^{m} p_j \boldsymbol{H}_{\mathfrak{D}_j,\lambda}^{-1} \boldsymbol{g}_{\mathcal{D},\lambda}.$$
(4)



Figure 1: The computations and communications in the *t*-th iteration for FedNewton.

The global learner $\bar{f}^t_{\mathcal{D},\lambda}(x) = \langle \bar{w}^t_{\mathcal{D},\lambda}, \phi(x) \rangle$ is updated by

174 175 176

177 178 179

181

182

183

185 186

187

$$\bar{\boldsymbol{w}}_{\mathcal{D},\lambda}^{t} = \bar{\boldsymbol{w}}_{\mathcal{D},\lambda}^{t-1} - \sum_{j=1}^{m} p_{j} \boldsymbol{H}_{\mathfrak{D}_{j},\lambda}^{-1} \boldsymbol{g}_{\mathcal{D},\lambda}^{t-1},$$
(5)

where $\bar{\boldsymbol{w}}_{\mathcal{D},\lambda}^t$ is the model after t iterations and the global gradient is $\boldsymbol{g}_{\mathcal{D},\lambda}^{t-1} = \sum_{j=1}^m p_j \boldsymbol{g}_{\mathfrak{D}_j,\lambda}^{t-1}$ from Proposition [] The approximation error between equation [] and equation [] is analyzed in Section []. Without loss of generality, we present the details of FedNewton in Algorithm [] and Figure [], which includes two times communications as the first-order methods in per round. Note that, the algorithm uploads local Newton updates $\boldsymbol{H}_{\mathfrak{D}_j,\lambda}^{-1} \boldsymbol{g}_{\mathcal{D},\lambda}^{t-1} \in \mathbb{R}^M$ instead of local inverse Hessians $\boldsymbol{H}_{\mathfrak{D}_j,\lambda}^{-1} \in \mathbb{R}^{M \times M}$, reducing communication costs from $\boldsymbol{O}(M^2)$ to $\boldsymbol{O}(M)$.

Computational complexity analysis. With finite-dimensional feature mappings $\phi : \mathcal{X} \to \mathbb{R}^M$, we 188 compute time complexity, space complexity, and communication complexity of FedNewton. The 189 space complexity on the *j*-th local machine is $\mathcal{O}(|\mathfrak{D}_j|M+M^2)$ to store $\Phi_{\mathfrak{D}_j}, H_{\mathfrak{D}_j,\lambda}$ and $H_{\mathfrak{D}_j,\lambda}^{-1}$ 190 while the global server requires $\mathcal{O}(mM)$ space to store $g_{\mathfrak{D}_j,\lambda}$ and $H^{-1}_{\mathfrak{D}_j,\lambda}g_{\mathcal{D},\lambda}$. Before the iterations, 191 192 the computations of $H_{\mathfrak{D}_j,\lambda}$ and $H_{\mathfrak{D}_j,\lambda}^{-1}$ costs $O(|\mathfrak{D}_j|M^2+M^3)$ time. In each iteration, the local time 193 complexity is $\mathcal{O}(M^2)$ to compute local gradient $g_{\mathfrak{D}_j,\lambda}$ and local Newton update $H_{\mathfrak{D}_j,\lambda}^{-1} g_{\mathfrak{D},\lambda}$, while the time complexity on the global server is $\mathcal{O}(mM)$ to update the global gradient and estimator. 194 195 Therefore, the total time complexity is $\mathcal{O}(\max_{j \in [m]} |\mathfrak{D}_j| M^2 + M^3 + M^2t + mMt)$. 196

Remark 3 (Communication burdens). The per iteration communication costs of the proposed 197 FedNewton are 2 times as compared to the first-order FL algorithms, e.g. FedAvg and FedProx, but the number of iterations for FedNewton is much fewer. The total communication complex-199 ity is $\mathcal{O}(Mt)$, the same as most first-order Federated algorithms. Notably, from Theorem \overline{I} the 200 iteration complexity is a linear convergence $t = \Omega(\log(1/\epsilon))$ where ϵ is the federated error, i.e., 201 FedNewton converges exponentially to the global estimator equation 2, while first-order feder-202 ated algorithms requires a large number of communication rounds $t = \Omega(1/\epsilon)$ Su et al. (2021). 203 Therefore, FedNewton cannot reduce the communication complexity for once communication as 204 communication-efficient FL algorithms Sattler et al. (2019); Reisizadeh et al. (2020); Wu et al. 205 (2022), but it significantly reduces the number of communication rounds, e.g., FedNewton with 206 $t \leq 2$ achieves good predictive performance in Section 7

207 **Remark 4** (Beyond the squared loss). To quantify the impacts from local sample size, data het-208 erogeneity and model heterogeneity, we apply the squared loss for FedNewton because it admits 209 closed-form solutions and is convenient for the theoretical analysis. Nevertheless, the proposed al-210 gorithm FedNewton is not applies to a broad range of loss functions as long as they are twice differentiable to compute the gradient $g_{\mathfrak{D}_j,\lambda}^{t-1}$ and the Hessian matrix $H_{\mathfrak{D}_j,\lambda}$. If the Hessian is in-211 dependent from the weights, the compute of local Hessians can be out of the loop, e.g. ReLU and 212 213 the squared loss. However, if the Hessian is relevant to the weights, for example exponential loss functions and trigonometric loss functions, we should compute the local Hessians for all iterations, 214 causing huge computational burdens. For other type loss functions, the weights can be initialized as 215 $\bar{\boldsymbol{w}}_{\mathcal{D},\lambda}^{0} = \boldsymbol{0}.$

216 4 MAIN RESULTS

In this section, to explore the factors that affect performance, we derive the excess risk bounds for FedNewton in homogeneous settings and heterogeneous settings, respectively.

4.1 NOTATIONS AND ASSUMPTIONS

We consider a broader scenario for federated learning, where the local training sets contain both heterogenous inputs (covariate shift) $\mathfrak{D}_j \sim \rho_j$ and different responses (concept shift) $y_{\mathfrak{D}_j} \sim \rho_j(y|\mathbf{x})$. The concept shift is represented as

$$f^*(\boldsymbol{x}) = \int_{\mathcal{Y}} y d\rho(y|\boldsymbol{x}), \, \boldsymbol{x} \in \mathcal{X}, \qquad f_j^*(\boldsymbol{x}) = \int_{\mathcal{Y}} y d\rho_j(y|\boldsymbol{x}), \, \boldsymbol{x} \in \mathcal{X}, \, j \in [m], \tag{6}$$

where f_j^* is the underlying mechanism governing the true responses on the *j*-th worker. Give a $x \in \mathcal{X}$ and $j, k, \in [m]$, the responses may be different $f_j^*(x) \neq f_k^*(x)$ when $j \neq k$.

Definition 1 (Operators with feature mapping ϕ). Using the feature mapping $\phi : \mathcal{X} \to \mathcal{H}_K, \forall \beta \in \mathcal{H}_K$, the covariance operators $C, C_j, C_{\mathcal{D}}, C_{\mathfrak{D}_j} : \mathcal{H}_K \to \mathcal{H}_K$ are defined as

$$C\boldsymbol{\beta} = \int_{X} \langle \boldsymbol{\beta}, \phi(\boldsymbol{x}) \rangle \phi(\boldsymbol{x}) d\rho_{X}(\boldsymbol{x}), \qquad C_{\mathcal{D}}\boldsymbol{\beta} = \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \langle \boldsymbol{\beta}, \phi(\boldsymbol{x}_{i}) \rangle \phi(\boldsymbol{x}_{i}), \,\forall \, (\boldsymbol{x}_{i}, y_{i}) \in \mathcal{D},$$
$$C_{j}\boldsymbol{\beta} = \int_{X} \langle \boldsymbol{\beta}, \phi(\boldsymbol{x}) \rangle \phi(\boldsymbol{x}) d\rho_{j}(\boldsymbol{x}), \qquad C_{\mathfrak{D}_{j}}\boldsymbol{\beta} = \frac{1}{|\mathfrak{D}_{j}|} \sum_{i=1}^{|\mathfrak{D}_{j}|} \langle \boldsymbol{\beta}, \phi(\boldsymbol{x}_{i}) \rangle \phi(\boldsymbol{x}_{i}), \,\forall \, (\boldsymbol{x}_{i}, y_{i}) \in \mathfrak{D}_{j}.$$

> Note that, $C_{\mathcal{D}} = \mathbf{\Phi}_{\mathcal{D}}^{\top} \mathbf{\Phi}_{\mathcal{D}}, C_{\mathfrak{D}_j} = \mathbf{\Phi}_{\mathfrak{D}_j}^{\top} \mathbf{\Phi}_{\mathfrak{D}_j}$ are the empirical covariance operators on \mathcal{D} and \mathfrak{D}_j , while $C = \mathbb{E}_{\rho}[C_{\mathcal{D}}], C_j = \mathbb{E}_{\rho_j}[C_{\mathfrak{D}_j}]$ are their expected counterparts.

For the sake of readability, we provide some notations

$$\mathcal{P}_{\mathfrak{D}_j,\lambda} := \| (C_{\mathfrak{D}_j} + \lambda I)^{-1} (C_j + \lambda I) \|, \qquad \mathcal{R}_{\mathfrak{D}_j,\lambda} := \| (C_j + \lambda)^{-1} (C_j - C_{\mathfrak{D}_j}) \|,$$

$$\Delta_{\mathfrak{D}_j} := \| C - C_j \|, \qquad \qquad \Delta_{f_j} := \| f^* - f_j^* \|.$$

The quantities $\mathcal{P}_{\mathfrak{D}_j,\lambda}$ and $\mathcal{R}_{\mathfrak{D}_j,\lambda}$ measure the similarity between the expected covariance operator and its empirical counterpart. From contraction inequalities for self-adjoint operators, a larger number of local samples $|\mathfrak{D}_j|$ leads to smaller $\mathcal{P}_{\mathfrak{D}_j,\lambda}$ and $\mathcal{R}_{\mathfrak{D}_j,\lambda}$. Note that, $\Delta_{\mathfrak{D}_j}$ measures the data heterogeneity on the expected covariance operator, while Δ_{f_j} measures the model heterogeneity on the true regressions.

254 We let $||f||_2 = \sqrt{\langle f, f \rangle} = \sqrt{\int_X |f(x)|^2 d\mathbb{P}(x)}$ denote the $L^2(\mathbb{P})$ norm and $L^2(\mathbb{P}) = \{f : \mathcal{X} \to \mathbb{R} \mid ||f||_2^2 < \infty\}$. Throughout this paper, we assume the outputs are bounded $|y| \le B$ almost surely for some B > 0 and $\kappa := ||\phi(x)||_K < \infty$ for any $x \in \mathcal{X}$.

Assumption 1 (Federated capacity condition). For $\lambda \in (0, 1)$, we define the effective dimensions on the global distribution ρ and local distributions ρ_j , $\forall j \in [m]$ as

$$\mathcal{N}(\lambda) = Tr(C(C+\lambda I)^{-1}), \ \mathcal{N}_j(\lambda) = Tr(C_j(C_j+\lambda I)^{-1}).$$

Assume there exists Q > 0 and $\gamma \in [0, 1]$, such that

$$\max\left(\mathcal{N}(\lambda), \mathcal{N}_1(\lambda), \cdots, \mathcal{N}_m(\lambda)\right) \le Q^2 \lambda^{-\gamma}.$$

Assumption 2 (Source condition). Define the integral operators $L: L^2(\mathbb{P}) \to L^2(\mathbb{P})$,

$$(Lg)(\cdot) = \int_X \langle \phi(\cdot), \phi({oldsymbol x})
angle g({oldsymbol x}) d
ho_X({oldsymbol x}), \quad orall \, g \in L^2(\mathbb{P}).$$

Assume there exists R > 0, r > 0, such that $||L^{-r}f^*|| \le R$. where the operator L^r denotes the r-th power of L as a compact and positive operator.

270 Capacity condition and source condition are standard assumptions in the optimal statistical learning 271 for the KRR related literature Caponnetto & De Vito (2007); Smale & Zhou (2007); Rudi & Rosasco 272 (2017); Lin & Cevher (2020); Liu et al. (2021). The effective dimensions $\mathcal{N}(\lambda)$ and $\mathcal{N}_j(\lambda)$ measure 273 the capacities of the RKHS \mathcal{H}_K on the global distribution ρ and the local distributions ρ_i , $\forall j \in [m]$. 274 Here, we modify the conventional capacity condition for federated learning to impose 275 constraints on local estimators. Note that, for effective dimensions, it holds 1/2 \leq 276 $\max(\mathcal{N}(\lambda), \mathcal{N}_1(\lambda), \cdots, \mathcal{N}_m(\lambda)) \leq \kappa^2 \lambda^{-1}$ Rudi et al. (2015). Assumption 1 reflects the vari-277 ance of the estimator. A larger γ leads to a larger \mathcal{H}_K and $\gamma = 1$ corresponds to the capacity 278 independence case. Assumption 2 controls the bias of an estimator, which reflects the regularity of 279 the estimator. The bigger r leads to the stronger regularity of the regression and the easier learning problem. The general settings $(r = 1/2, \gamma = 1)$ lead to $O(1/\sqrt{|D|})$ convergence rates for KRR 281 related approaches.

4.2 ERROR DECOMPOSITION

Theorem 1. Let $f_{\mathcal{D},\lambda}, \bar{f}^t_{\mathcal{D},\lambda}, f^*$ be defined according to equation 2 equation 5 and equation 6. Then, the following error decomposition holds

$$\|\bar{f}_{\mathcal{D},\lambda}^{t} - f^{*}\| \leq \underbrace{\|\bar{f}_{\mathcal{D},\lambda}^{t} - f_{\mathcal{D},\lambda}\|}_{\text{federated error}} + \underbrace{\|f_{\mathcal{D},\lambda} - f^{*}\|}_{\text{centralized excess risk}}, \tag{7}$$

and the federated error for FedNewton is bounded by:

$$\|\bar{f}_{\mathcal{D},\lambda}^t - f_{\mathcal{D},\lambda}\|_2 \leq \Upsilon^t \left\| (C + \lambda I)^{1/2} (\bar{\boldsymbol{w}}_{\mathcal{D},\lambda}^0 - \boldsymbol{w}_{\mathcal{D},\lambda}) \right\|_K,$$

291 292 293

294

282

284

285

286 287

289

290

where $\Upsilon = \sum_{j=1}^{m} p_j \mathcal{P}_{\mathfrak{D}_j,\lambda} \left(2\mathcal{R}_{\mathfrak{D}_j,\lambda} + \frac{\Delta_{\mathfrak{D}_j}}{\lambda} \right) \left(1 + \frac{\Delta_{\mathfrak{D}_j}}{\lambda} \right).$

In the above theorem, we decompose the excess risk for FedNewton into two parts: the federated error $\|\bar{f}_{\mathcal{D},\lambda}^t - f_{\mathcal{D},\lambda}\|$ and the excess risk for the centralized KRR $\|f_{\mathcal{D},\lambda} - f^*\|$. Since the generalization analysis for $\|f_{\mathcal{D},\lambda} - f^*\|$ is standard Caponnetto & De Vito (2007); Smale & Zhou (2007), we focus on the federated error $\|\bar{f}_{\mathcal{D},\lambda}^t - f_{\mathcal{D},\lambda}\|$.

299 From Theorem Π we find that the value of Υ determines the effectiveness of multiple iterations. 300 If $\Upsilon \geq 1$, FedNewton with multiple communications is worse than oneshot federated learning 301 (DKRR). However, when $\Upsilon < 1$, the federated error decreases exponentially and the rate of convergence is referred to as *linear convergence* in the optimization literature Bottou et al. (2018). 302 The quantities $\mathcal{P}_{\mathfrak{D}_i,\lambda}$ and $\mathcal{R}_{\mathfrak{D}_i,\lambda}$ measure the similarity between $C_{\mathfrak{D}_i}$ and C_j where those quanti-303 ties decrease as the local sample size $|\mathfrak{D}_j|$ increases. Because Υ is proportional to $\mathcal{P}_{\mathfrak{D}_j,\lambda}, \mathcal{P}_{\mathfrak{D}_j,\lambda}$ 304 and $\Delta_{\mathfrak{D}_j}$, the *linear convergence* requires both a sufficient number of local examples $|\mathfrak{D}_j|$ and 305 moderate data heterogeneity $\Delta_{\mathfrak{D}_j}$. If t = 0, the above error bound degrades into that for DKRR 306 $\|\bar{f}_{\mathcal{D},\lambda} - f_{\mathcal{D},\lambda}\|_2 \le \left\| (C + \lambda I)^{1/2} (\bar{\boldsymbol{w}}_{\mathcal{D},\lambda}^0 - \boldsymbol{w}_{\mathcal{D},\lambda}) \right\|_{K}.$ 307

Theorem 2. Under Assumption 2 with a high probability $1 - \delta$, $\forall \delta \in (0, 1)$, the federated error can be bounded

$$\begin{split} \|\bar{f}_{\mathcal{D},\lambda}^t - f_{\mathcal{D},\lambda}\|_2 \lesssim &\Upsilon^t \sum_{j=1}^m p_j \sqrt{1 + \frac{\Delta_{\mathfrak{D}_j}}{\lambda}} \left(2\mathcal{R}_{\mathfrak{D}_j,\lambda} + \frac{(1 + \mathcal{R}_{\mathfrak{D}_j,\lambda})\Delta_{\mathfrak{D}_j}}{\lambda} \right) \\ & \left(\left(\frac{1}{|\mathfrak{D}_j|\sqrt{\lambda}} + \sqrt{\frac{\mathcal{N}(\lambda)}{|\mathfrak{D}_j|}} \right) \log \frac{2}{\delta} + \frac{\Delta_{\mathfrak{D}_j}}{\lambda} + \Delta_{f_j} \right). \end{split}$$

313 314

308

310 311 312

315

316 Theorem 2 illustrates the key factors that affect the federated error: the discrepancy between ex-317 pected and empirical covariance operators $\mathcal{R}_{\mathfrak{D}_i,\lambda}$, the covariate shift $\Delta_{\mathfrak{D}_i}$, and the model hetero-318 geneity Δ_{f_i} . The smaller these factors, the smaller the federated error. The federated error results 319 from three parts: distributed error $\frac{1}{\sqrt{\lambda}|\mathfrak{D}_j|} + \sqrt{\frac{\mathcal{N}(\lambda)}{|\mathfrak{D}_j|}}$, covariate shift $\Delta_{\mathfrak{D}_j}/\lambda$ and concept shift Δ_{f_j} . 320 Specifically, as the increase of local sample size, the distributed error decreases. However, the con-321 cept shifts Δ_{f_i} is a constant and it will dominate the federated error when model heterogeneity Δ_{f_i} 322 is large. In the case $\Upsilon < 1$, iterators can reduce the federated error, alleviating the entire federated 323 error term.

4.3 HOMOGENEOUS SETTING

324

325

326 327 328

330

Theorem 3. Let $\delta \in (0, 1/3]$, $\lambda = |\mathcal{D}|^{\frac{-1}{2r+\gamma}}$ and $2r + \gamma \ge 1$. Under Assumptions 2 if $\Delta_{\mathfrak{D}_j} = 0$ and $\Delta_{f_j} = 0$, with the probability at least $1 - 3\delta$, it holds

$$\|\bar{f}_{\mathcal{D},\lambda}^t - f^*\|_2 \lesssim \Upsilon^t \sum_{j=1}^m p_j \aleph_j \log^2 \frac{2}{\delta} + |\mathcal{D}|^{\frac{-r}{2r+\gamma}} \log \frac{2}{\delta}.$$

Here, \aleph_i and Υ have different values w.r.t local sample size

$$\aleph_{j} = \begin{cases} |\mathfrak{D}_{j}|^{-2} |\mathcal{D}|^{\frac{1.5}{2r+\gamma}}, & \text{if } |\mathfrak{D}_{j}| \lesssim |\mathcal{D}|^{\frac{1-\gamma}{2r+\gamma}} \\ |\mathfrak{D}_{j}|^{-1.5} |\mathcal{D}|^{\frac{1+0.5\gamma}{2r+\gamma}}, & \text{if } |\mathcal{D}|^{\frac{1-\gamma}{2r+\gamma}} \lesssim |\mathfrak{D}_{j}| \lesssim |\mathcal{D}|^{\frac{1}{2r+\gamma}} \\ |\mathfrak{D}_{j}|^{-1} |\mathcal{D}|^{\frac{1+\gamma}{4r+2\gamma}}, & \text{if } |\mathcal{D}|^{\frac{1}{2r+\gamma}} \lesssim |\mathfrak{D}_{j}| \lesssim |\mathcal{D}|^{\frac{2r+\gamma+1}{4r+2\gamma}} \\ |\mathcal{D}|^{\frac{-r}{2r+\gamma}}, & \text{if } |\mathfrak{D}_{j}| \gtrsim |\mathcal{D}|^{\frac{2r+\gamma+1}{4r+2\gamma}}, \end{cases}$$

and $\Upsilon = 2 \sum_{j=1}^{m} p_j \mathcal{P}_{\mathfrak{D}_j,\lambda} \mathcal{R}_{\mathfrak{D}_j,\lambda}$ holds

344 Note that, the second term in the above bound is from the centralized model $||f_{\mathcal{D},\lambda} - f^*||_2$, where 345 the learning rate $O(|\mathcal{D}|^{\frac{-r}{2r+\gamma}})$ is optimal in a minimax sense Caponnetto & De Vito (2007). The 346 performance of FedNewton in the homogeneous setting is only affected by the local sample size. 347 We discuss the above result in three parts. First, when the number of local examples is limited $|\mathfrak{D}_j| \leq 1$ 348 $|\mathcal{D}|^{\frac{1}{2n+\gamma}}$, in another word the number of local machines is larger than $m \gtrsim |\mathcal{D}|^{\frac{2r+\gamma-1}{2r+\gamma}}$, the federated error dominates the excess risk and fails to achieve the optimal rate, where the convergence rates 349 350 are slower than $\mathcal{O}(|\mathcal{D}|^{\frac{\gamma-1}{4r+2\gamma}})$. Meanwhile, when the number of local examples is limited, it leads 351 to $\Upsilon \geq 1$ and multiple communications hurt the performance. Second, when $|\mathcal{D}|^{\frac{1}{2r+\gamma}} \lesssim |\mathfrak{D}_i| \lesssim 1$ 352 353 $|\mathcal{D}|^{\frac{2r+\gamma+1}{4r+2\gamma}}$, although the convergence rates of federated error are still not the optimal, the iterator Υ 354 is smaller than one, leading to a linear convergence. As the increase of communications $t \to \infty$, the 355 centralized excess risk will dominate the error bound that achieves the optimal rate. Third, with a large number of local examples $|\mathfrak{D}_i| \gtrsim |\mathcal{D}|^{\frac{2r+\gamma+1}{4r+2\gamma}}$, even with insufficient communications $t \to 0$, 356 357 the error bound still achieves the optimal rate $O(|\mathcal{D}|^{\frac{-r}{2r+\gamma}})$. 358

Theorem 3 can be further simplified in some special cases. For example, we consider the general case $(r = 1/2, \gamma = 1)$, where r = 1/2 is equivalent to assuming $f^* \in \mathcal{H}_K$ and $\gamma = 1$ is the capacity independent case. The learning rate achieves $O(1/\sqrt{|\mathcal{D}|})$ when $|\mathfrak{D}_j| \gtrsim |\mathcal{D}|^{0.5}$ with multiple iterations or $|\mathfrak{D}_j| \gtrsim |\mathcal{D}|^{0.75}$ with only one communication.

Remark 5. The existing theoretical guarantees for DKRR Zhang et al. (2015); Guo et al. (2017); Lin & Cevher (2020) focused on how to achieve the optimal rate by a sufficient number of local examples (or lower the number of partitions), but they ignored the sub-optimal case that the local sample size is fixed and insufficient. However, in federated learning, the number of partitions is fixed and local examples are generated locally, such that sub-optimal cases are more general. Theorem Jillustrate that a sufficient number of local examples is crucial for both learning rates (in generalization) and convergence rate (in optimization).

Remark 6 (Finite dimensional case). In the proofs of theoretical findings, we consider the estimator 370 in RKHS with $w \in \mathcal{H}_K$. However, the finite-dimensional cases are more general, i.e. $w \in \mathbb{R}^M$ in 371 Algorithm 1 where the feature mappings are explicit and can be neural networks or random features 372 Rahimi & Recht (2007). With a simple modification of our proofs, one can derive similar results for 373 finite-dimensional cases. In particular, under same assumptions of Theorem 3 and $(r = 1/2, \gamma = 0)$, 374 then with high probability, $\|\bar{f}_{\mathcal{D},\lambda}^t - f^*\|_2 \lesssim |\mathfrak{D}_j|^{-2} |\mathcal{D}|^{1.5} + \sqrt{M/|\mathcal{D}|}$, provided that $|\mathcal{D}| \gtrsim M \log M$. 375 As shown in Rudi & Rosasco (2017), a large number of random features $M \gtrsim |\mathcal{D}|^{\frac{1+\gamma(2r-1)}{2r+\gamma}}$ can guarantee the optimal rates for $||f_{\mathcal{D},\lambda} - f^*||_2$, and thus we can also provide similar results as 376 377 Theorem 3

3783794.4 HETEROGENEOUS SETTING

Theorem 4. Let $\delta \in (0, 1/3]$, $\lambda = |\mathcal{D}|^{\frac{-1}{2r+\gamma}}$ and $2r + \gamma \ge 1$. Under Assumptions 1 2 with the probability at least $1 - 3\delta$, the excess risk bound for FedNewton holds

$$\|\bar{f}_{\mathcal{D},\lambda}^t - f^*\|_2 \lesssim \Upsilon^t \sum_{j=1}^m p_j \sqrt{1 + \frac{\Delta_{\mathfrak{D}_j}}{\lambda}} (\aleph_j + \Pi_j) \log^2 \frac{2}{\delta} + |\mathcal{D}|^{\frac{-r}{2r+\gamma}} \log \frac{2}{\delta}.$$

386

384

387 388 Here, $\Upsilon = \sum_{j=1}^{m} p_j \mathcal{P}_{\mathfrak{D}_j,\lambda}(2\mathcal{R}_{\mathfrak{D}_j,\lambda} + \frac{\Delta_{\mathfrak{D}_j}}{\lambda})(1 + \frac{\Delta_{\mathfrak{D}_j}}{\lambda}), \ \aleph_j \text{ is same to Theorem } \mathcal{J} \text{ and } \int |\mathcal{D}|^{\frac{2}{2r+\gamma}} \Lambda_{j-1} + |\mathcal{D}|^{\frac{1}{2r+\gamma}} \Lambda_{j-1} + |\mathcal{D}|^{\frac{1}{2r+\gamma}}$

$$\Pi_{j} = \begin{cases} \frac{||\mathcal{D}_{j}||}{|\mathcal{D}_{j}|} \Delta_{\mathcal{D}_{j}} + \frac{||\mathcal{D}_{j}||}{||\mathcal{D}_{j}||} \Delta_{f_{j}}, & \text{if } ||\mathcal{D}_{j}| \gtrsim ||\mathcal{D}||^{2r+\gamma} \\ (1 + |\mathcal{D}|^{\frac{1}{2r+\gamma}} \Delta_{\mathcal{D}_{j}}) (\Delta_{f_{j}} + |\mathcal{D}|^{\frac{1}{2r+\gamma}} \Delta_{\mathcal{D}_{j}}), & \text{otherwise.} \end{cases}$$

389 390 391

392

393

394

397

398

We add some comments on the above theorem. First, when the local sample size is insufficient $|\mathfrak{D}_j| \lesssim |\mathcal{D}|^{\frac{1}{2r+\gamma}}$ or the data heterogeneity is considerable, we have $\Upsilon \geq 1$, and communications hurt the performance. Meanwhile, since the federated error $\sqrt{1 + \Delta_{\mathfrak{D}_j}/\lambda}(\aleph_j + \Pi_j)$ depends on $|\mathfrak{D}_j|, \Delta_{\mathfrak{D}_j}$, and Δ_{f_j} , the learning rate is far from the optimal rate. Second, when the number of local examples is sufficient $|\mathfrak{D}_j| \gtrsim |\mathcal{D}|^{\frac{1}{2r+\gamma}}$ and data heterogeneity is small, it holds $\Upsilon < 1$ where communications can improve the generalization ability of FedNewton. In this case, the federated error $\|\overline{f}_{\mathcal{D},\lambda}^t - f_{\mathcal{D},\lambda}\|$ converge exponentially fast. If t is large enough, the error bound in Theorem 4 depends on the centralized excess risk $\|f_{\mathcal{D},\lambda} - f^*\|_2$ and achieves the optimal learning rate.

The learning rate of generalization bound in Theorem 4 is determined by four factors: the local sample size $|\mathfrak{D}_j|$, the covariate shift $\Delta_{\mathfrak{D}_j}$, the response shift Δ_{f_j} and the number of iterations t. Furthermore, the iterator value Υ depends on $|\mathfrak{D}_j|$ and $\Delta_{\mathfrak{D}_j}$, such that these two values are important factors for both fast convergences (in optimization) and the learning rates (in generalization).

Remark 7 (How to achieve the optimal rate in federated learning?). The value of $\Upsilon < 1$ is key to 404 obtaining a linear convergence rate and the optimal learning rate, where it depends on both local 405 sample sizes $\Upsilon \propto \mathcal{R}_{\mathfrak{D}_i,\lambda} \propto |\mathfrak{D}_j|$ and data heterogeneity $\Upsilon \propto \Delta_{\mathfrak{D}_i}$. Note that, $\Delta_{\mathfrak{D}_i}$ measures 406 the intrinsic discrepancy between local distributions and the global one, and thus it is a fixed value 407 independent from the local sample size. Therefore, since $\Delta_{\mathfrak{D}_i}$ is a constant, we can obtain $\Upsilon <$ 408 1 with a large number of local examples generated by local machines. And then, with a large number of iterations when $\Upsilon < 1$, the federated error, depending on both data heterogeneity and 409 model heterogeneity, can become small enough to be negligible. In this case, a large number of 410 local examples can guarantee both a linear convergence rate (for federated error) and the optimal 411 learning rate (from the centralized excess risk). A large number of local examples benefit both 412 optimization and generalization, rather than making tradeoffs between them. 413

414 415

416

5 COMPARED WITH RELATED WORK

417 We compare FedNewton with recent Newton-type methods, DKRR methods, and first-order FL 418 algorithms in both algorithmic and theoretical fronts. Table [] reports the main factors that affect the 419 performance, the computational and generalization properties of related work.

420 **Compared with Newton-type FL methods.** Local Newton-type FL algorithms Yang et al. (2019); 421 Ghosh et al. (2020); Gupta et al. (2021) conducted Newton updates instead of SGD in local ma-422 chines, which only utilized local information (local SGD & local Hessian). Recent studies Safaryan et al. (2022); Qian et al. (2022) tried to use global information (global SGD & global Hessian) by 423 communicating local Hessian shifts, but it leads to high communication costs $O(M^2)$ per commu-424 nication. Nevertheless, this work employs mixed information (global SGD & local Hessian) that reduce the communication cost to O(M). More importantly, the existing Newton-type FL work 426 only provided the convergence analysis (optimization) Ghosh et al. (2020); Safaryan et al. (2022); 427 Qian et al. (2022) without out-sample (generalization) error bounds, while this work bridges the 428 optimization and generalization for FedNewton, which essentially guarantees its fast convergence 429 and good generalization ability. 430

431 **Compared with DKRR.** The time complexities of DKRR approaches solved in kernel space Zhang et al. (2015); Guo et al. (2017) are much higher than that of stochastic optimization methods solved

Table 1: Summary of computational and generalization properties for related work.

Dalatad Work	0	Δ	Δ	Tasining Time	Testing	Commun	Conditions	Lacal Size 10	Itonation 4	Unnan David
Related Work	$ \mathfrak{D}_j $	$\Delta_{\mathfrak{D}_j}$	Δ_{f_j}	Training Time	Time	ication	Conditions	Local Size $ \mathfrak{D}_j $	neration t	Opper Bound
DKRR Zhang et al. (2015)	\checkmark	×	×	$ \mathfrak{D}_j ^3$	$ \mathcal{D}_{test} \mathcal{D} $	$ \mathcal{D} $	Specific kernels	$\Omega(r^2\kappa^4\log \mathcal{D})$	O(1)	$O\left(\frac{1}{ \mathcal{D} }\right)$
DKRR Guo et al. (2017)	\checkmark	×	×	$ \mathfrak{D}_j ^3$	$ \mathcal{D}_{test} \mathcal{D} $	$ \mathcal{D} $	$r\in [1/2,1]$	$\Omega(\mathcal{D} ^{\frac{1+\gamma}{2r+\gamma}})$	O (1)	$\boldsymbol{O}(\mathcal{D} ^{\frac{-r}{2r+\gamma}})$
DKRR-SGD Lin & Cevher (2018)	\checkmark	×	×	$ \mathcal{D} t$	$ \mathcal{D}_{test} \mathcal{D} $	$ \mathcal{D} $	$r\in [1/2,1]$	$\Omega(\mathcal{D} ^{\frac{1}{2r+\gamma}})$	$oldsymbol{O}(\mathcal{D} ^{rac{2-\gamma}{2r+\gamma}})$	$oldsymbol{O}\left(\mathcal{D} ^{rac{-r}{2r+\gamma}} ight)$
DKRR-CM Lin et al. (2020)	\checkmark	×	×	$ \mathfrak{D}_j ^3+ \mathcal{D} \mathfrak{D}_j t$	$ \mathcal{D}_{test} \mathcal{D} $	$ \mathcal{D} t$	$r\in [1/2,1]$	$\Omega(\mathcal{D} ^{\frac{2r+\gamma+1}{4r+2\gamma}})$	$O(\log \frac{1}{\epsilon})$	$O\left(\mathcal{D} ^{rac{-r}{2r+\gamma}} ight)$
FedAvg Su et al. (2021)	×	×	\checkmark	$ \mathfrak{D}_j M^2 + M^2t + mMt$	$ \mathcal{D}_{\text{test}} M$	Mt	Specific kernels	/	$O(\frac{1}{\epsilon})$	$O\left(\frac{1}{\eta t} + \frac{\Delta_f^2}{ \mathcal{D} }\right)$
FedProx Su et al. (2021)	×	×	\checkmark	$ \mathfrak{D}_{j} M^{2} + M^{3} + M^{2}t + mMt$	$ \mathcal{D}_{\text{test}} M$	Mt	Specific kernels	/	$O(\frac{1}{\epsilon})$	$O\left(\frac{1}{\eta t} + \frac{\Delta_f^2}{ \mathcal{D} }\right)$
Theorem 3	\checkmark	×	×	$ \mathfrak{D}_{j} M^{2} + M^{3} + M^{2}t + mMt$	$ \mathcal{D}_{\text{test}} M$	Mt	$r>0, 2r+\gamma\geq 1$	$\Omega(\mathcal{D} ^{\frac{1}{2r+\gamma}})$	$O(\log \frac{1}{\epsilon})$	Theorem 3
Theorem 4	\checkmark	\checkmark	\checkmark	$ D_j M^2 + M^3 + M^2t + mMt$	$ \mathcal{D}_{\text{test}} M$	Mt	$r>0, 2r+\gamma\geq 1$	$\Omega(\mathcal{D} ^{\frac{1}{2r+\gamma}})$	$O(\log \frac{1}{\epsilon})$	Theorem 4

Note: The computational complexities are computed in terms of regularized least squared loss. We estimate the upper bounds for

 $||f - f^*||_2 \forall f \in L^2(\mathbb{P})$. We denote $\mathcal{D}_{\text{test}}$ the testing data, n the step-size for SGD approaches, ϵ the federated error and $\Delta_f^2 = \sum_{j=1}^m p_j \Delta_{f_j}^2$. For Rademacher complexities based bounds Zhang et al. (2015): Su et al. (2021), specific kernels include kernels with finite-rank or polynomial eigenvalues decay. Integral operator based bounds Guo et al. (2017): Lin & Cevher (2018); Lin et al. (2020) also assume $\gamma \in [0, 1]$. We compute exact local solution for FedProx.

450 in feature space. Both our work and Guo et al. (2017); Lin & Cevher (2018); Lin et al. (2020) are 451 based on integral operator techniques, but DKRR literature assumes all local datasets are drawn i.i.d. from an identical distribution, ignoring the data heterogeneity and model heterogeneity, which 452 makes the proofs much easier than ours. We emphasize the difference between this work and DKRR 453 theories as bellow: 1) DKRR work required a strict condition $r \in [1/2, 1]$, while we relax the condition to $r > 0, 2r + \gamma \ge 1$. 2) This work pertains to NonIID data, covering both covariate shift 455 $\Delta_{\mathfrak{D}_i}$ and response shift Δ_{f_i} , DKRR only applied to IID data that is a special case in the homogenous 456 setting $\Delta_{\mathfrak{D}_i} = \Delta_{f_i} = 0$ in Theorem 3. 3) Because of the existence of data heterogeneity and model 457 heterogeneity, we cannot directly estimate the difference between local estimators and global ones, 458 and thus we introduce novel error decompositions for the federated error. 4) This work explores the 459 excess risk bounds in terms of different local sample size (\aleph_i in Theorem 3), covering both optimal 460 and sub-optimal rates, while DKRR work only studied the optimal learning rates with the restrict on the number of partitions, i.e. $m = O(|\mathcal{D}|^{\frac{(2r+\gamma-1)(t+1)}{(2r+\gamma)(t+2)}})$ Lin et al. (2020). 461 462

Compared with first-order methods. Using the random matrix theory and the local Rademacher 463 complexity, Su et al. (2021) provided the optimal guarantees $||f - f^*||_2^2 = O(1/|\mathcal{D}|)$. However, as 464 shown in Theorem 2 Su et al. (2021), it directly assumed all inputs are sampled i.i.d from an identical 465 distribution, ignore the local sample size and the data heterogeneity, while our theoretical results 466 illustrate both the local sample size and the data heterogeneity are crucial to federated learning. Su 467 et al. (2021) also imposed several strict assumptions: 1) the ideal model belongs to the hypothesis 468 space, corresponding to $r \in [1/2, 1]$; 2) small hypothesis space with local Rademacher complexity, 469 corresponding $\gamma \to 0$ in our work; 3) specific kernels maybe not suitable to the federated learning 470 tasks and lead to sub-optimal rates. In this work, we remove these three conditions based on the 471 integral operator approach, which makes our theoretical findings applicable to broader settings. Our 472 results illustrate that only a few iterations can guarantee the optimal rates $O(|\mathcal{D}|^{\frac{-2r}{2r+\gamma}})$ when the number of local examples is sufficient and data heterogeneity is moderate, where the convergence 473 rate of federated error is *linear*, while in Su et al. (2021) the learning rate is always affected by 474 model heterogeneity $O(\frac{\sum_{j=1}^{m} p_j \Delta_{f_j}^2}{|\mathcal{D}|})$ and the convergence rate is *sublinear*. 475 476

477 478

479

432

445 446 447

448 449

CONCLUSION AND FUTURE WORK 6

480 In this paper, we present an efficient second-order optimization method for FL. We derive gener-481 alization bounds with the optimal rates, which quantify the impacts of local sample size, the data heterogeneity, and the model heterogeneity. In benign cases, the federated error convergence exponentially fast, and thus communications can be small. Our theoretical findings fill the gap between 483 optimization and generalization for federated learning, rather than focusing on one of them. Overall, 484 the techniques presented here highlight new ways for designing efficient algorithms and analyzing 485 both generalization and optimization for FL.

486 REFERENCES

512

- Frank Bauer, Sergei Pereverzev, and Lorenzo Rosasco. On regularization algorithms in learning theory. *Journal of Complexity*, 23(1):52–72, 2007.
- Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine
 learning. *SIAM Review*, 60(2):223–311, 2018.
- Charles George Broyden. The convergence of a class of double-rank minimization algorithms 1.
 general considerations. *IMA Journal of Applied Mathematics*, 6(1):76–90, 1970.
- Andrea Caponnetto and Ernesto De Vito. Optimal rates for the regularized least-squares algorithm.
 Foundations of Computational Mathematics, 7(3):331–368, 2007.
- Nicolò Dal Fabbro, Subhrakanti Dey, Michele Rossi, and Luca Schenato. Shed: A newton-type algorithm for federated learning based on incremental hessian eigenvector sharing. *Automatica*, 160:111460, 2024.
- Ron S Dembo, Stanley C Eisenstat, and Trond Steihaug. Inexact newton methods. *SIAM Journal* on Numerical analysis, 19(2):400–408, 1982.
- Junichi Fujii, Masatoshi Fujii, Takayuki Furuta, and Ritsuo Nakamoto. Norm inequalities equivalent to heinz inequality. *Proceedings of the American Mathematical Society*, 118(3):827–830, 1993.
- Avishek Ghosh, Raj Kumar Maity, and Arya Mazumdar. Distributed newton can communicate
 less and resist byzantine workers. In *Advances in Neural Information Processing Systems 33* (*NeurIPS*), volume 33, pp. 18028–18038, 2020.
- Margalit R Glasgow, Honglin Yuan, and Tengyu Ma. Sharp bounds for federated averaging (local sgd) and continuous perspective. In *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 9050–9090. PMLR, 2022.
- Zheng-Chu Guo, Shao-Bo Lin, and Ding-Xuan Zhou. Learning theory of distributed spectral algorithms. *Inverse Problems*, 33(7):074009, 2017.
- Vipul Gupta, Avishek Ghosh, Michal Derezinski, Rajiv Khanna, Kannan Ramchandran, and Michael
 Mahoney. Localnewton: Reducing communication bottleneck for distributed learning. *arXiv preprint arXiv:2105.07320*, 2021.
- Rustem Islamov, Xun Qian, Slavomír Hanzely, Mher Safaryan, and Peter Richtárik. Distributed newton-type methods with communication compression and bernoulli aggregation. *Transactions* on Machine Learning Research, 2023.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems 31* (*NeurIPS*), pp. 8571–8580, 2018.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, pp. 5132–5143.
 PMLR, 2020.
- Quoc Le, Tamás Sarlós, and Alex Smola. Fastfood-approximating kernel expansions in loglinear time. In *Proceedings of the 30th International Conference on Machine Learning (ICML)*, volume 85, 2013.
- Jian Li, Yong Liu, and Weiping Wang. Fedns: A fast sketching newton-type algorithm for federated learning. *AAAI Conference on Artificial Intelligence (AAAI)*, pp. 13509–13517, 2023.
- Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020a.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith.
 Federated optimization in heterogeneous networks. In *Proceedings of Machine Learning and Systems 2020 (MLSys)*, 2020b.

540 541 542	Junhong Lin and Volkan Cevher. Optimal distributed learning with multi-pass stochastic gradient methods. In <i>Proceedings of the 35th International Conference on Machine Learning (ICML)</i> , pp. 3098–3107, 2018.							
543 544 545 546	Junhong Lin and Volkan Cevher. Optimal convergence for distributed learning with stochastic gra- dient methods and spectral algorithms. <i>Journal of Machine Learning Research (JMLR)</i> , 21(147): 1–63, 2020.							
547 548	Shao-Bo Lin, Di Wang, and Ding-Xuan Zhou. Distributed kernel ridge regression with communications. <i>Journal of Machine Learning Research (JMLR)</i> , 21(93):1–38, 2020.							
549 550 551 552	Chengchang Liu, Lesi Chen, Luo Luo, and John CS Lui. Communication efficient distributed new- ton method with fast convergence rates. In <i>Proceedings of the 29th ACM SIGKDD Conference on</i> <i>Knowledge Discovery and Data Mining</i> , pp. 1406–1416, 2023.							
553 554	Dong C Liu and Jorge Nocedal. On the limited memory bfgs method for large scale optimization. <i>Mathematical Programming</i> , 45(1):503–528, 1989.							
555 556 557	Yong Liu, Jiankun Liu, and Shuqiang Wang. Effective distributed learning with random features: Improved bounds and algorithms. In <i>Proceedings of the 9th International Conference on Learning</i> <i>Representations (ICLR)</i> , 2021.							
558 559 560 561 562	Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In <i>Proceedings of</i> <i>the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)</i> , pp. 1273– 1282. PMLR, 2017.							
563 564 565	Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In <i>Proceed</i> - ings of the 36th International Conference on Machine Learning (ICML), pp. 4615–4625. PMLR, 2019.							
566	Radford M Neal. BAYESIAN LEARNING FOR NEURAL NETWORKS. PhD thesis, Citeseer, 1995.							
567 568 569 570	Reese Pathak and Martin J Wainwright. Fedsplit: an algorithmic framework for fast federated optimization. In <i>Advances in Neural Information Processing Systems 33 (NeurIPS)</i> , volume 33, pp. 7057–7066, 2020.							
571 572	Mert Pilanci and Martin J Wainwright. Newton sketch: A near linear-time optimization algorithm with linear-quadratic convergence. <i>SIAM Journal on Optimization</i> , 27(1):205–245, 2017.							
573 574 575 576	Xun Qian, Rustem Islamov, Mher Safaryan, and Peter Richtarik. Basis matters: Better communication-efficient second order methods for federated learning. In <i>Proceedings of the 25th International Conference on Artificial Intelligence and Statistics (AISTATS)</i> , pp. 680–720, 2022.							
577 578	Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In Advances in Neural Information Processing Systems 20 (NIPS), pp. 1177–1184, 2007.							
579 580 581 582	Amirhossein Reisizadeh, Aryan Mokhtari, Hamed Hassani, Ali Jadbabaie, and Ramtin Pedarsani. Fedpaq: A communication-efficient federated learning method with periodic averaging and quan- tization. In <i>Proceedings of the 23rd International Conference on Artificial Intelligence and Statis-</i> <i>tics (AISTATS)</i> , pp. 2021–2031. PMLR, 2020.							
584 585	Lorenzo Rosasco, Mikhail Belkin, and Ernesto De Vito. On learning with integral operators. <i>Journal of Machine Learning Research (JMLR)</i> , 11(Feb):905–934, 2010.							
586 587	Alessandro Rudi and Lorenzo Rosasco. Generalization properties of learning with random features. In Advances in Neural Information Processing Systems 30 (NIPS), pp. 3215–3225, 2017.							
588 589 590 591	Alessandro Rudi, Raffaello Camoriano, and Lorenzo Rosasco. Less is more: Nyström computational regularization. In <i>Advances in Neural Information Processing Systems 28 (NIPS)</i> , pp. 1657–1665, 2015.							
592 593	Mher Safaryan, Rustem Islamov, Xun Qian, and Peter Richtarik. Fednl: Making newton-type meth- ods applicable to federated learning. In <i>Proceedings of the 39th International Conference on</i> <i>Machine Learning (ICML)</i> , pp. 18959–19010. PMLR, 2022.							

594 595 596	Felix Sattler, Simon Wiedemann, Klaus-Robert Müller, and Wojciech Samek. Robust and communication-efficient federated learning from non-iid data. <i>IEEE Transactions on Neural Networks and Learning Systems (TNNLS)</i> , 31(9):3400–3413, 2019.
598 599	Nicol N Schraudolph. Fast curvature matrix-vector products for second-order gradient descent. <i>Neural Computation</i> , 14(7):1723–1738, 2002.
600 601	Steve Smale and Ding-Xuan Zhou. Learning theory estimates via integral operators and their approximations. <i>Constructive Approximation</i> , 26(2):153–172, 2007.
602 603 604	Lili Su, Jiaming Xu, and Pengkun Yang. A non-parametric view of fedavg and fedprox: Beyond stationary points. <i>arXiv preprint arXiv:2106.15216</i> , 2021.
605	Joel A Tropp. User-friendly tools for random matrices: An introduction. Technical report, 2012.
606	Vladimir Vapnik. The nature of statistical learning theory. Springer Verlag, 2000.
607 608 609 610	Hongyi Wang, Mikhail Yurochkin, Yuekai Sun, Dimitris Papailiopoulos, and Yasaman Khazaeni. Federated learning with matched averaging. In <i>Proceedings of the 8th International Conference</i> <i>on Learning Representations (ICLR)</i> , 2020.
611 612 613	Shusen Wang, Fred Roosta, Peng Xu, and Michael W Mahoney. Giant: Globally improved approxi- mate newton method for distributed optimization. In <i>Advances in Neural Information Processing</i> <i>Systems 31 (NeurIPS)</i> , pp. 2338–2348, 2018.
614 615 616	David P Woodruff et al. Sketching as a tool for numerical linear algebra. <i>Foundations and Trends</i> ® <i>in Theoretical Computer Science</i> , 10(1–2):1–157, 2014.
617 618	Chuhan Wu, Fangzhao Wu, Lingjuan Lyu, Yongfeng Huang, and Xing Xie. Communication- efficient federated learning via knowledge distillation. <i>Nature Communications</i> , 13(1):1–8, 2022.
619 620 621	Semih Yagli, Alex Dytso, and H Vincent Poor. Information-theoretic bounds on the generalization error and privacy leakage in federated learning. In <i>Proceedings of the 21st International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)</i> , pp. 1–5. IEEE, 2020.
623 624 625	Jiyan Yang, Vikas Sindhwani, Haim Avron, and Michael Mahoney. Quasi-monte carlo feature maps for shift-invariant kernels. In <i>Proceedings of the 31st International Conference on Machine Learn-ing (ICML)</i> , pp. 485–493, 2014.
626 627 628	Kai Yang, Tao Fan, Tianjian Chen, Yuanming Shi, and Qiang Yang. A quasi-newton method based vertical federated learning framework for logistic regression. <i>arXiv preprint arXiv:1912.00513</i> , 2019.
629 630	Yun Yang, Mert Pilanci, Martin J Wainwright, et al. Randomized sketches for kernels: Fast and optimal nonparametric regression. <i>The Annals of Statistics</i> , 45(3):991–1023, 2017.
632 633	Mang Ye, Xiuwen Fang, Bo Du, Pong C Yuen, and Dacheng Tao. Heterogeneous federated learning: State-of-the-art and research challenges. <i>ACM Computing Surveys</i> , 56(3):1–44, 2023.
634 635 636 637	Honglin Yuan, Warren Richard Morningstar, Lin Ning, and Karan Singhal. What do we mean by generalization in federated learning? In <i>Proceedings of the 10th International Conference on Learning Representations (ICLR)</i> , 2022.
638 639	Chen Zhang, Yu Xie, Hang Bai, Bin Yu, Weihong Li, and Yuan Gao. A survey on federated learning. <i>Knowledge-Based Systems</i> , 216:106775, 2021.
640 641 642	Yuchen Zhang, John Duchi, and Martin Wainwright. Divide and conquer kernel ridge regression: A distributed algorithm with minimax optimal rates. <i>Journal of Machine Learning Research</i> (<i>JMLR</i>), 16(1):3299–3340, 2015.
643 644 645	Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated learning with non-iid data. <i>arXiv preprint arXiv:1806.00582</i> , 2018.
646 647	Hanhan Zhou, Tian Lan, Guru Prasadh Venkataramani, and Wenbo Ding. Every parameter matters: Ensuring the convergence of federated learning with dynamic heterogeneous models reduction. <i>Advances in Neural Information Processing Systems</i> , 36, 2023.