Federated Reinforcement Learning: Linear Speedup Under Markovian Sampling

Sajad Khodadadian ¹ Pranay Sharma ² Gauri Joshi ² Siva Theja Maguluri ¹

Abstract

Since reinforcement learning algorithms are notoriously data-intensive, the task of sampling observations from the environment is usually split across multiple agents. However, transferring these observations from the agents to a central location can be prohibitively expensive in terms of the communication cost, and it can also compromise the privacy of each agent's local behavior policy. In this paper, we consider a federated reinforcement learning framework where multiple agents collaboratively learn a global model, without sharing their individual data and policies. Each agent maintains a local copy of the model and updates it using locally sampled data. Although having N agents enables the sampling of N times more data, it is not clear if it leads to proportional convergence speedup. We propose federated versions of on-policy TD, off-policy TD and Q-learning, and analyze their convergence. For all these algorithms, to the best of our knowledge, we are the first to consider Markovian noise and multiple local updates, and prove a linear convergence speedup with respect to the number of agents. To obtain these results, we show that federated TD and Q-learning are special cases of a general framework for federated stochastic approximation with Markovian noise, and we leverage this framework to provide a unified convergence analysis that applies to all the algorithms.

Proceedings of the 39 th International Conference on Machine Learning, Baltimore, Maryland, USA, PMLR 162, 2022. Copyright 2022 by the author(s).

1. Introduction

Reinforcement Learning (RL) is an online sequential decision-making paradigm that is typically modeled as a Markov Decision Process (MDP) (Sutton & Barto, 2018). In an RL task, the agent aims to learn the optimal policy of the MDP that maximizes long-term reward, without knowledge of its parameters. The agent performs this task by repeatedly interacting with the environment according to a behavior policy, which in turn provides data samples that can be used to improve the policy. This MDP-based RL framework has a vast array of applications including self-driving cars (Yurtsever et al., 2020), robotic systems (Kober et al., 2013), games (Silver et al., 2016), UAV-based surveillance (Yun et al., 2022), and Internet of Things (IoT) (Lim et al., 2020).

Due to the high-dimensional state and action spaces that are typical in these applications, RL algorithms are extremely data hungry (Duan et al., 2016; Kalashnikov et al., 2018; Akkaya et al., 2019), and training RL models with limited data can result in low accuracy and high output variance (Islam et al., 2017; Xu et al., 2021). However, generating massive amounts of training data sequentially can be extremely time consuming (Nair et al., 2015). Hence, many practical implementations of RL algorithms from Atari domain to Cyber-Physical Systems rely on parallel sampling of the data from the environment using multiple agents (Mnih et al., 2016; Espeholt et al., 2018; Chen et al., 2021a; Xu et al., 2021). It was empirically shown in (Mnih et al., 2016) that the federated version of these algorithms yields faster training time and improved accuracy. A naive approach would be to transfer all the agents' locally collected data to a central server that uses it for training. However, in applications such as IoT (Chen & Giannakis, 2018), autonomous driving (Shalev-Shwartz et al., 2016) and robotics (Kalashnikov et al., 2018), communicating high-dimensional data over low bandwidth network link can be prohibitively slow. Moreover, sharing individual data of the agents with the server might also be undesirable due to privacy concerns (Yang et al., 2019; Mothukuri et al., 2021).

Federated Learning (FL) (Kairouz et al., 2019) is an emerging distributed learning framework, where multiple agents seek to collaboratively train a shared model, while complying with the privacy and data confidentiality requirements

¹H. Milton Stewart School of Industrial & Systems Engineering, Georgia Institute of Technology, Atlanta, GA, 30332, USA ²Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA, 15213, USA,. Correspondence to: Sajad Khodadadian <skhodadadian3@gatech.edu>.

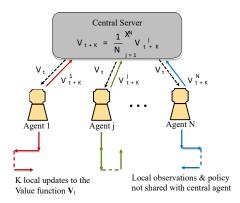


Figure 1. Schematic representation of FRL where N agents follow a Markovian trajectory and synchronize their parameter every K time steps.

(Qi et al., 2021; Yang et al., 2019). The key idea is that the agents collect data, use on-device computation capabilities to locally train the model, and only share the model updates with the central server. Not sharing data reduces communication cost and also alleviates privacy concerns.

Recently, there is a growing interest in employing FL for RL algorithms (also known as FRL) (Nadiger et al., 2019; Liu et al., 2019; Ren et al., 2019; Xu et al., 2021; Zhang et al., 2022). Unlike standard supervised learning where data is collected before training begins, in FRL, each agent collects data by following its own Markovian trajectory, while simultaneously updating the model parameters.

To ensure convergence, after every K time steps, the agents communicate with the central server to synchronize their parameters (see Figure 1). Intuitively, using more agents and a higher synchronization frequency should improve the convergence of training algorithm. However, the following questions remain to be concretely answered:

- 1. With N agents, do we get an N-fold (linear) speedup in the convergence of FRL algorithms?
- 2. How does the convergence speed and the final error scale with synchronization frequency K?

While these questions are well-studied (Wang & Joshi, 2021; Stich, 2018; Qu et al., 2020; Li et al., 2019) in federated supervised learning, only a few works (Wai, 2020; Shen et al., 2020) have attempted to answer them in the context of FRL. However, none of them have established the convergence analysis of FRL algorithms by considering Markovian local trajectories and multiple local updates (see Table 1).

In this paper, we tackle this challenging open problem and answer both the questions listed above. We propose communication-efficient federated versions of on-policy TD,

off-policy TD, and Q-learning algorithms. In addition, we are the first to establish the convergence bounds for these algorithms in the realistic Markovian setting, showing a linear speedup in the number of agents. Previous works (Liu & Olshevsky, 2021; Shen et al., 2020) on distributed R L have only shown such a speedup by assuming i.i.d. noise. Moreover, based on experiments, (Shen et al., 2020) conjectures that linear speedup may be possible under the realistic Markov noise setting, which we establish analytically. The main contributions and organization of the paper are summarized below.

- In the on-policy setting, in Section 4 we propose and analyze federated TD-learning with linear function approximation, where the agents' goal is to evaluate a common policy using on-policy samples collected in parallel from their environments. The agents only share the updated value function (not data) with the central server, thus saving communication cost. We prove a linear convergence speedup with the number of agents and also characterize the impact of communication frequency on the convergence.
- In the off-policy setting, in Section 5 we propose and analyze the federated off-policy TD-learning and federated Q-learning algorithms. Again, we establish a linear speedup in their convergence with respect to the number of agents and characterize the impact of synchronization frequency on the convergence. Since every agent samples data using a private policy and only communicates the updated value or Q-function, off-policy FRL helps keep both the data as well as the policy private.
- In Section 6, we propose a general <u>Fed</u>erated <u>S</u>tochastic <u>A</u>pproximation framework with <u>M</u>arkovian noise (FedSAM) which subsumes both federated TD-learning and federated Q-learning algorithms proposed above. Considering Markovian sampling noise poses a significant challenge in the analysis of this algorithm. The convergence result for FedSAM serves as a workhorse that enables us to analyze both federated TD-learning and federated Q-learning. We characterize the convergence of FedSAM with a refined analysis of general stochastic approximation algorithms, fundamentally improving upon prior work.

2. Related Work

Single node TD-learning and Q-learning. Most existing RL literature is focused on designing and analyzing algorithms that run at a single computing node. In the on-policy setting, the asymptotic convergence of TD-learning was established in (Tsitsiklis & Van Roy, 1997; Tadić, 2001; Borkar, 2009), and the finite-sample bounds were studied

Table 1. Comparison of sample complexity results for federated supervised learning (local SGD) and reinforcement learning algorithms. The possible distributed architectures are: 1) Worker-server, with a central server that coordinates with N agents; 2) Decentralized, where each agent directly communicates with its neighboring agents, without a central server; and 3) Shared memory, where each agent modifies a subset of the parameters of a global model held in a shared memory, that is accessible to all agents. (Recht et al., 2011).

Algorithm	Architecture	References	Local Updates	Markov Noise	Linear Speedup
Local SGD	Worker-server	(Khaled et al., 2020)	✓	7	✓
Local SGD	Worker-server	(Spiridonoff et al., 2021)	✓	7	✓
TD(0)	Worker-server	(Liu & Olshevsky, 2021)	✓	7	✓
Stochastic Approximation	Decentralized	(Wai, 2020)	7	✓	7
A3C-TD(0)	Shared memory	(Shen et al., 2020)	7	7	✓
A3C-TD(0)	Shared memory	(Shen et al., 2020)	7	✓	7
TD & Q-learning	Worker-server	This paper	✓	✓	✓

in (Dalal et al., 2018; Lakshminarayanan & Szepesvari, 2018; Bhandari et al., 2018; Srikant & Ying, 2019; Hu & Syed, 2019; Chen et al., 2021c). In the off-policy set-ting, (Maei, 2018; Zhang et al., 2020) study the asymptotic and (Chen et al., 2020a; 2021c) characterize the finite time bound of TD-learning. The Q-learning algorithm was first proposed in (Watkins & Dayan, 1992). There has been a long line of work to establish the convergence properties of Q-learning. In particular, (Tsitsiklis, 1994; Jaakkola et al., 1994; Bertsekas & Tsitsiklis, 1996b; Borkar & Meyn, 2000; Borkar, 2009) characterize the asymptotic convergence of Qlearning, (Beck & Srikant, 2012b; 2013; Wainwright, 2019; Chen et al., 2020a; 2021c) study the finite-sample convergence bound in the mean-square sense, and (Even-Dar & Mansour, 2004; Li et al., 2020; Qu & Wierman, 2020) study the high-probability convergence bounds of Q-learning.

Federated Learning with i.i.d. Noise. When multiple agents are used to expedite sample collection, transferring the samples to a central server for the purpose of training can be costly in applications with high-dimensional data (Shao et al., 2019) and it may also compromise the agents' privacy. Federated Learning (FL) is an emerging distributed optimization paradigm (Konečný et al., 2016; Kairouz et al., 2019) that utilizes local computation at the agents to train models, such that only model updates, not data, is shared with the central server. In local Stochastic Gradient Descent (Local SGD or FedAvg) (McMahan et al., 2017; Stich, 2018), the core algorithm in FL, locally trained models are periodically averaged by the central server in order to achieve consensus among the agents at a reduced communication cost. While the convergence of local SGD has been extensively studied in prior work (Khaled et al., 2020; Spiridonoff et al., 2021; Qu et al., 2020; Koloskova et al., 2020), these works assume i.i.d. noise in the gradients, which is acceptable for SGD but too restrictive for RL algorithms.

Distributed and Multi-agent RL. Some recent works have analyzed distributed and multi-agent RL algorithms in the presence of Markovian noise in various settings such as decentralized stochastic approximation (Doan et al., 2019; Sun et al., 2020; Wai, 2020; Zeng et al., 2020), TD learning with linear function approximation (Wang et al., 2020a), and off-policy TD in actor-critic algorithms (Chen et al., 2021e;f). However, all these works consider decentralized settings, where the agents communicate with their neighbors after every local update. On the other hand, we consider a federated setting, with each agent performing multiple local updates between successive communication rounds, thereby resulting in communication savings. In (Shen et al., 2020), a parallel implementation of asynchronous advantage actor-critic (A3C) algorithm (which does not have local updates) has been proposed under both i.i.d. and Markov sampling. However, the authors prove a linear speedup only for the i.i.d. case, and an almost linear speedup is observed experimentally for the Markovian case.

3. Preliminaries: Single Node Setting

We model our RL setting with a Markov Decision Process (MDP) with 5 tuples (S; A; P; R;), where S and A are finite sets of states and actions, P is the set of transition probabilities, R is the reward function, and 2 (0; 1) denotes the discount factor. At each time step t, the system is in some state S_t , and the agent takes some action A_t according to a policy (jS_t) in hand, which results in reward $R(S_t; A_t)$ for the agent. In the next time step, the system transitions to a new state S_{t+1} according to the state transition probability $P(jS_t; A_t)$. This series of states and actions $(S_t; A_t)_{t0}$ constructs a Markov chain, which is the source of the Markovian noise in RL. Throughout this paper we assume that this Markov chain is irreducible and aperiodic (also known as ergodic). It is known that this Markov

chain asymptotically converges to a steady state, and we denote its stationary distribution with .

To measure the long-term reward achieved by following policy , we define the value function

$$V(s)=E \int_{t=0}^{\infty} {^tR(S_t; A_t)jS_0} = s; A_t(jS_t) : (1)$$

3.1. Temporal Difference Learning

An intermediate goal in RL is to estimate the value function (either (V (s))_{s2s} or v) corresponding to a particular pol-icy using data collected from the environment. This task is denoted as policy evaluation and one of the commonly-used approaches to accomplish this is Temporal Difference (TD)-learning (Sutton, 1988). TD-learning is an iterative algorithm where the elements of a d (or jSj, in the tabular setting) dimensional vector is updated until it converges to v (or V). This evaluated value function can be employed in different RL algorithms such as actor-critic (Konda & Tsitsiklis, 2000). In the on-policy function approximation setting, the update of the n-step TD-learning is as follows

where is the step size. Note that in this setting, the evaluating policy and the sampling policy coincide. In contrast, in the off-policy setting these two policies can in general differ, and we need to account for this difference while running the algorithm. We will further expand on TD-learning and its variants in Sections 4.1 and 5.1.1.

3.2. Control Problem and Q-learning

Assuming some initial distribution on the state space, the average value function corresponding to policy is

defined as $V() = E_s[V(s)]$. This scalar quantity is a metric of average long-term rewards achieved by the agent, when it starts from distribution and follows policy. The ultimate goal of the agent is to obtain an optimal policy which results in the maximum long-term rewards, i.e. 2 arg max V(). Throughout the paper, we denote the parameters corresponding to the optimal policy with ,

e.g., V () V (). The task of obtaining the optimal policy

in RL is denoted as the control problem.

Q-learning (Watkins & Dayan, 1992) is one of the most widely used algorithms in RL to solve the control problem. At each time step t, Q-learning preserves a jSj:jAj dimensional table Q_t , and updates it table as $Q_{t+1}(s;a) = Q_t(S_t;A_t) + (R(S_t;A_t) + \max_a Q_t(S_{t+1};a) Q_t(S_t;A_t))$, if $(s;a) = (S_t;A_t)$ and $Q_{t+1}(s;a) = Q_t(s;a)$ otherwise. The jSjjAj elements of the vector Q_t are updated iteratively until it converges to Q_t , corresponding to an optimal policy. Using Q_t , one can obtain an optimal policy via greedy selection.

3.3. Stochastic Approximation and Finite Sample Bounds

Both TD-learning and Q-learning can be seen as variants of stochastic approximation (Chen et al., 2020b; 2019b;a; 2021d; Tsitsiklis, 1994). While generic stochastic approximation algorithms are studied under i.i.d. noise (Even-Dar & Mansour, 2004; Shah & Xie, 2018; Wainwright, 2019; Liu et al., 2015; Dalal et al., 2018), to apply them for studying RL we need to understand stochastic approximation under Markovian noise (Tsitsiklis, 1994; Qu & Wierman, 2020; Srikant & Ying, 2019; Chen et al., 2021c) which is significantly more challenging.

For a generic stochastic approximation (i.i.d. or Markovian noise) with constant step size , parameter vector \mathbf{x}_{T} , and convergent point \mathbf{x} , it can be shown that the algorithm have the following convergence behaviour

$$E[kx_T xk^2] C_1(1 C_0)^T + C_2;$$
 (3)

where C_0 ; C_1 ; and C_2 are some problem dependent positive constants (Look at Appendix A for a discussion on a lower bound on the convergence of general stochastic approximation). The first term is denoted as the bias and the second term is called the variance. According to this bound, x_T geometrically converges to a ball around x with radius proportional to C_2 . Notice that we can always reduce the variance term by reducing the step size , but this will lead to slower convergence in the bias term. In particular, in order to get $E[kx \quad xk^2]$, it is easy to see that we need T O $\frac{C_2}{C_2}$ log $\frac{1}{C_2}$ sample complexity. Now suppose the constant C_2 is large. In this case, the variance term in the bound in (3) is large, and the sample complexity, which is proportional to C_2 will be poor. Notice that by the dis-

cussion in Appendix A, this bound is tight and cannot be improved.

This is where the FL can be employed in order to control the variance term by generating more data. For instance, in federated TD-learning, multiple agents work together to evaluate the value function simultaneously. Due to this collaboration, the agents can estimate the true value function with a lower variance. The same holds for estimating Q in Q-learning.

4. Federated On-policy R L

4.1. On-policy TD-learning with linear function approximation

In this section we describe the TD-learning with linear function approximation and online data samples in the single node setting. In this problem, we consider a full rank feature matrix 2 jSj d, and we denote s-th row of this matrix with (s); s 2 S. The goal is to find v 2 R^d which solves the following fixed point equation:

$$V = ((T)^{n}V):$$
 (4)

In equation (4), () is the projection with respect to the weighted 2-norm, i.e., $(V) = \arg \min_{v \ge R^d} kv$ Vk. Here kVk = p V > V and is a diagonal matrix with diagonal entries corresponding to . In equation (4), (T)ⁿ denotes the n-step Bellman operator (Tsitsiklis & Van Roy, 1997). It is known (Tsitsiklis & Van Roy, 1997) that equation (4) has a unique solution v, and v is "close" to the true value function V . n-step TD-learning algorithm, which was shown in (2), is an iterative algorithm to obtain this unique fixed point using samples from the environment. Note that in this algorithm states and actions are sampled over a single trajectory, and hence the noise in updating v_t is Markovian. Furthermore, since the policy which samples the actions and the the evaluating policy are both, this algorithm is on-policy. As described in (Tsitsiklis & Van Roy, 1997; Bertsekas & Tsitsiklis, 1996a), the TD-learning algorithm can be studied under the umbrella of linear stochastic approximation with Markovian noise. More recently, the authors in (Bhandari et al., 2018; Srikant & Ying, 2019) have shown that the update parameter of TDlearning v_t converges to v in the form $E[kv_t vk^2]$ O((1 $(C_0)^t + (C_0)^t + (C_0$ this result.

4.2. Federated TD-learning with linear function approximation

The federated version of on-policy n-step TD-learning with linear function approximation is shown in Algorithm 1. In this algorithm we consider N agents which collaboratively work together to evaluate v. For each agent i; i = 1; 2; :::; N, we initialize their corresponding pa-

rameters $v_0^i = 0$. Furthermore, each agent i samples its initial state S_0^i from some given distribution . In the next time steps, each agent follows a single Markovian trajectory generated by policy , independently from other agents. At each time t, the parameter of each agent i is updated using this independently generated trajectory as $v_{t+1}^i = v_t^i + (S^i)E_t^{i-t}$, Finally, in order to ensure convergence to a global optimum, every K time steps all the agents send their parameters to a central server. The central server evaluates the average of these parameters and returns this average to each of the agents. Each agent then continues their update procedure using this average.

Notice that the averaging step is essential to ensure synchronization among the agents. Smaller K results in more frequent synchronization, and hence better convergence guarantees. However, setting smaller K is equivalent to more number of communications between the single agents and the central server, which incurs higher cost. Hence, an intermediate value for K has to be chosen to strike a balance between the communication cost and the accuracy. At the end, the algorithm samples a time step T^{Λ} q T^{c} , where

$$q_T^c(t) = P \frac{c^{-t}}{\sum_{t=0}^{T} c^{-t}} \text{ for } t = 0; 1; :::; T$$
 1 (5)

and c > 1 is some constant. Since we have $q^c(t) = 0$ and c = 1 is some constant. Since we have $q^c(t) = 0$ and c = 1 it is clear that c = 1 is a probability distribution over the time interval c = 1. In Theorem 4.1 we characterize the convergence of this algorithm as a function of , N , and K . Throughout the paper, c = 1 ignores the logarithmic terms.

Algorithm 1 Federated n-step TD (On-policy, Function Approx.)

```
1: Input: Policy;

2: Initialization: v_0^i = 0 and S_0^i and fA^i; S_l^i g_1 for 0 \mid n = 1 and all i

3: for t = 0 to T = 1 do

4: for i = 1; \dots; N do

5: Sample _{t+A}^i (jS_{t+n}^i); S_{t+n+1}^i 6:

e_{t;l}^i = P(R_l^i)(S_{t+n}^i) A_{t+n}^i) + (S_l^i) v_l^i (S_l^i) v_l^i for l = t; \dots; t + n = 1

7: E_{t:n}^i = \sum_{l=t}^{t+n} \sum_{l=t}^{1-t} t_l^i

8: v_{t+1}^i = v_t^i + (S_l^i) E_t^i 9:_{t;n}
end for

10: if t + 1 \mod K_P = 0 then

11: v_{t+1}^i = \sum_{l=t}^{1-t} \sum_{l=t}^{1-t} v_{t+1}^i; 8 i 2 [N]

12: end if

13: end for

14: Sample T_l^i = \sum_{l=t}^{T} \sum_{l=t}^{T_{t+1}} v_l^i
```

Theorem 4.1. Let $v_{\land} = \underbrace{1}^{P} i^{N} = \underbrace{1}_{1} v^{i}$ denote the average $T_{N} = \underbrace{1}_{1} v^{i}$ of the parameters across agents at the random time T^{\bullet} . For small enough step size , and $T_{N} = T^{\bullet}$ (log T^{\bullet}), there exist constant $C_{TDL} = T^{\bullet}$ (see Section C.1 for precise statement), such that we have

$$E[kv_{\uparrow} \quad vk^{2}]_{2}C^{\top D_{\downarrow}}_{1}(1 \stackrel{1}{-}C^{\top D_{\downarrow}})^{\top}_{0} + C^{\top D_{\downarrow}} - \frac{^{2}}{N^{2}} + C^{\top D_{\downarrow}}_{3}(K \quad 1)^{2} + C^{\top D_{\downarrow}}_{3}^{32}$$

where $C_i^{TD_L}$, i=0;1;2;3;4 are problem dependent constants, and =O(log(1=)). By choosing $=O(\frac{log(NT)}{T})$ and K=T=N, we achieve $E[kv_{T,A}vk^2]$ within $T=O_N$ iterations.

For brevity purposes, here we did not show the exact dependence of the constants $C_i^{\mathsf{TD}_L}$, i=0;1;2;3 on the problem dependent constants. For a discussion on the detailed expression look at Section C.1 in the appendix.

Theorem 4.1 shows that federated TD-learning with linear function approximation enjoys a linear speedup with respect to the number of agents. Compared to the convergence bound of general stochastic approximation in (3), the bound in Theorem 4.1 has three differences. Firstly, the variance term which is proportional to the step size is divided with the number of agents N . This will allow us to control the variance (and hence improve sample complexity) by employing more number of agents. Secondly, we have an extra term which is zero with perfect synchronization K=1. Although this term is not divided with N , but it is proportional to 2 , which is one order higher than the variance term in (3). Finally, the last term is of the order O(3), which can be handled by choosing small enough step size.

Furthermore, according to the choice of K in Theorem 4.1, after T iterations, the communication cost of federated TD is T=K=N. However, by employing federated TD-learning in the naive setting where all the agents communicate with the central server at every time step, the communication cost will be O(T). Hence, we observe that by carefully tuning the hyper parameters of federated TD, we can significantly reduce the communication cost of the overall algorithm, while not loosing performance in terms of the sample complexity.

Finally, federated TD-learning Algorithm 1 preserves the privacy of the agents. In particular, since the single agents only require to share their parameters v_{t+1}^i , the central server will not be exposed to the state-action-reward trajectory generated by each agent. This can be essential in some applications where privacy is an issue (Mothukuri et al., 2021; Truex et al., 2019). Examples of such applications include autonomous driving (Liang et al., 2019; Zhao et al., 2021), Internet of Things (IoT) (Nguyen et al., 2021;

Ren et al., 2019; Wang et al., 2020b), and cloud robotics (Liu et al., 2019; Xu et al., 2021).

Remark. In algorithm 1, the randomness in choosing T^{\bullet} is independent of all the other randomness in the problem. Hence, in a practical setting, one can sample T^{\bullet} ahead of time, before running the algorithm, and stop the algorithm at time step T^{\bullet} and output v_{\uparrow} . By this method, we require only a single data point to be saved, which results in the memory complexity of O(1) for the algorithm.

5. Federated Off-Policy R L

On-policy TD-learning requires online sampling from the environment, which might be costly (e.g. robotics (Gu et al., 2017; Levine et al., 2020)), high risk (e.g. self-driving cars (Yurtsever et al., 2020; Maddern et al., 2017)), or unethical (e.g. in clinical trials (Gottesman et al., 2019; Liu et al., 2018; Gottesman et al., 2020)). Off-policy training in R L refers to the paradigm where we use data collected by a fixed behaviour policy to run the algorithm. When employed in federated setting, off-policy R L has privacy advantages as well (Foerster et al., 2016; Qi et al., 2021; Zhuo et al., 2019). In particular, suppose each single agent attains a unique sampling policy, and they do not wish to reveal these policies to the central server. In off-policy FL, agents only transmit sampled data, and hence the sampling policies remain private to each agent.

In Section 5.1 we will discuss off-policy TD-learning and in Section 5.2 we will discuss Q-learning, which is an off-policy control algorithm. For the off-policy algorithms, we only study the tabular setting. Notice that it has been observed that the combination of off-policy sampling and function approximation in R L (also known as deadly triad (Sutton & Barto, 2018)) can result in instability or even divergence (Baird, 1995). Recently there has been some work to overcome deadly triad (Chen et al., 2021b). Extension of our work to function approximation in the off-policy setting is a future research direction.

5.1. Federated Off-Policy TD-learning

In the following, we first discuss single-node off-policy TD-learning, and then we generalize it to the federated setting.

5.1.1. OFF-POLICY TD-LEARNING

In off-policy TD-learning the goal is to evaluate the value function $V = (V(s))_{s2S}$ corresponding to the policy using data sampled from some fixed behaviour policy $_b$. In this setting, the evaluating policy and the sampling policy $_b$ can be arbitrarily different, and we need to account for this difference while performing the evaluation. Although and $_b$ can be different, notice that the value function V does not depend on $_b$. In order to account for this

difference, we introduce the notion of importance sampling as $I^{b}(s; a) = \frac{(ajs)}{(ajs)}$ which is employed in the off-policy TD-learning.

Recently, several works studied the finite-time convergence of off-policy TD-learning. In particular, the authors in (Khodadadian et al., 2021; Chen et al., 2021c; 2020b; 2021d) show that, similar to on-policy TD, off-policy TD-learning can be studied under the umbrella of stochastic approximation. Hence, this algorithm enjoys similar convergence behaviour as (3).

5.1.2. FEDERATED OFF-POLICY TD-LEARNING

The federated version of n-step off-policy TD-learning is shown in Algorithm 2. In this algorithm, each agent i attains a unique (and private) sampling policy i and follows an independent trajectory generated by this policy. Furthermore, at each time step t, each agent i attains a jSj-dimensional vector V_t^i and updates this vector using the samples generated by i. In order to account for the off-policy sampling, each agent utilizes $I^{(i)}(S^i; A^i) = \frac{(A, j^i S, j^i)}{(A^i J, S^i)}$ in the update of their algorithm. We further define $I_{max} = \max_{s;a;i} I^{(i)}(s;a)$, which is a measure of discrepancy between the evaluating policy and sampling policy i of all the agents.

In order to ensure synchronization, all the agents transmit their parameter vectors to the central server every K time steps. The central server returns the average of these vectors to each agent and each agent follows this averaged vector afterwards. Notice that in federated off-policy TD-learning Algorithm 2, each agent share neither their sampled trajectory of state-action-rewards, nor their sampling policy with the central server. This provides two levels of privacy for the single agents. At the end, the algorithm samples a time step Υ q $^{\text{CT}}_{\text{T}}$, where the distribution q $^{\text{C}}_{\text{T}}$ is defined in (5) and $\text{CT}_{\text{T}}_{\text{D}} = 1$ $\frac{\text{YT}_{\text{D}}}{2}$, where $\frac{\text{Y}_{\text{T}}}{2}$ where $\frac{\text{Y}_{\text{T}}}{2}$ and $\frac{\text{Y}_{\text{T}}}{2}$. Here, we defined in (5) and $\frac{\text{Y}_{\text{T}}}{2}$ and $\frac{\text{Y}_{\text{T}}}{2}$. Here, we defined in (5) and $\frac{\text{Y}_{\text{T}}}{2}$ and $\frac{\text{Y}_{\text{$

note $_{min}$ = $min_{s;i}$ (s). The constant c_{TD} is carefully chosen to ensure the convergence of Algorithm 2. Furthermore, for small enough step size , it can be shown that $0 < c_{TD} < 1$.

Theorem 5.1 states the convergence of this Algorithm.

Theorem 5.1. Consider the federated n-step off-policy TD-learning Algorithm 2. Denote V $_{\uparrow} = \frac{1}{N} \bigvee_{i=1}^{N} \bigvee_{j=1}^{N} \bigvee_{j=1}^{N}$. For small enough step size and large enough T, we have

$$\begin{split} E[kV_{\uparrow \! r} & Vk^2_{\ 1}] \, C^{\top \, D_{\intercal}} \, \overset{1}{1} \, c^{\top}_{\ D_{\ T}} \!^{+} \, \, C^{\top \, D_{\intercal}}_{\ 2} \,^{2} + \frac{}{N} \\ & C^{\top \, D_{\intercal}}_{\ 3} \, (K \ 1)^2; \end{split}$$

where
$$C_1^{TD_T} = C_1^{TD_T} : \frac{\frac{1 + n}{max}^2}{\frac{1}{min}(1)^{-3}} = C_2^{TD_T}$$

Algorithm 2 Federated n-step TD (Off-policy Tabular Setting)

```
1: Input: Policy;
 2: Initialization: V_0^i = 0 and S_0^i and fA^i; S_i^i g_1^i for
    0 I n 1 and all i
 3: for t = 0 to T = 1 do
       for i = 1;:::; N do
          Sample A_{t+n}^{i} ^{i}(jS^{i}); S^{i} _{t+n} P(j_{t}S_{n+1}^{i}
 5:
         6:
          7:
          and V_{t+1}^{i}(s) = V_{t}^{i}(s) otherwise.
 8:
       if t + 1 mod K V_{t+1}^{i} = \frac{1}{N} P_{j=1}^{N} V_{t+1}^{j}; 8 i 2 [N]
 9:
10:
12: end for
13: Sample \Upsilon q_{T}^{T^{D}}
14: Return: \frac{1}{N} \Upsilon = 1 V_{\Upsilon}
```

The proof is given in Section C.2 in the appendix.

Note that similar to on-policy TD-learning Algorithm 1, off-policy TD-learning also enjoys a linear speedup while maintaining a low communication cost. In addition, this algorithm preserves the privacy of the agents by holding both the data and the sampling policy private.

5.2. Federated Q-learning

So far we have discussed policy evaluation problem with on and off-policy samples. Next we aim at solving the control problem by employing the celebrated Q-learning algorithm (Watkins & Dayan, 1992; Tsitsiklis, 1994). In the next section we will explain the Q-learning algorithm. Further, in Section 5.2.2 we will provide a federated version of Q-learning along with its convergence result.

5.2.1. Q-LEARNING

The goal of Q-learning is to evaluate Q, which is the unique Q-function corresponding to the optimal policy. Knowing Q, one can obtain an optimal policy through a greedy selection (Puterman, 2014), and hence resolve the control problem.

Suppose fS_t; A_tg_{t0} is generated by a fixed behaviour policy b. At each time step t, Q-learning preserves a [S]: [A] table Qt and updates the elements of this table as shown in Section 3.2. By assuming b to be an ergodic policy, the asymptotic convergence of Qt to Q has been established in (Bertsekas & Tsitsiklis, 1996b). Furthermore, it can be shown that Q-learning is a special case of stochastic approximation and enjoys a convergence bound similar to (3) (Beck & Srikant, 2012a; Li et al., 2020; Qu & Wierman, 2020; Chen et al., 2021c).

Two points worth mentioning about the Q-learning algorithm. Firstly, Q-learning is an off-policy algorithm in the sense that only samples from a fixed ergodic policy is needed to perform the algorithm. Secondly, as opposed to the TDlearning, the update of the Q-learning is non-linear. This imposes a sharp contrast between the analysis of Q-learning and TD-learning (Chen et al., 2019a).

5.2.2. FEDERATED Q-LEARNING

Algorithm 3 provides the federated version of Q-learning. We characterize its convergence in the following theorem.

Algorithm 3 Federated Q-learning

```
1: Input: Sampling policy i for i = 1; 2; :::; N, initial
   distribution
```

2: Initialization: $Q_0^i = 0$ and S_0^i for all i 3:

for t = 0 to T = 1 do

5: Sample
$$A_t^{i}(jS_t^{i}); S_t^{i} \xrightarrow{t+1} P(jS^i; A_t^{i})_t$$

6: Update
$$Q_{t+1}^{i}(s; a) = Q_{t}^{i}(S_{t}^{i}; A_{t}^{i}) + R(S_{t}; A_{t}^{i}) + \max_{a} Q_{t}(S_{t+1}^{i}; a) Q_{t}(S_{t}^{i}; A_{t}^{i})$$
, if $(s; a) = (S_{t}; A_{t}^{i})$ and $Q_{t+1}(s; a) = Q_{t}(s_{t}^{i}; a)$ otherwise.

7: end for

8: if
$$t + 1 \mod K_2 = 0$$
 then

8: if t + 1 mod K
$$_{p}$$
 = 0 then
9: Q_{t+1}^{i} $\frac{1}{N}$ $P_{j=1}^{N}$ Q_{t+1}^{j} ; 8 i 2 [N]

10:

11: end for

12: Sample:
$$\uparrow q_{\uparrow}^{\text{qc}}$$
13: Return: $\frac{1}{N} \stackrel{q}{\mid} q_{\uparrow}^{\text{qc}}$

Theorem 5.2. Consider the federated Q-learning Algorithm 3 with $c_Q = 1$ $\frac{r^{0.5 e^{1.4}(2 - min(1 - 1))}}{\frac{1}{p} e^{-1 + \frac{2 - min(1 - 1)}{2 - min(1 - 1)}}} = 1$ and we denote $\frac{r^{0.5 e^{1.4}(2 - min(1 - 1))}}{r^{0.5 e^{1.4}(2 - min(1 - 1))}} = 1$

 $min_{s;a;i}$ (s)ⁱ(ajs). Denote $Q_{\wedge} = \frac{P}{N} = \frac{N}{N} Q_{\wedge}^{i}$. For small enough step size and large enough T, we have

$$E[kQ_{\uparrow} \quad Qk^{2}_{1}] C_{1}^{1} C_{\uparrow}^{\uparrow} +_{Q} C_{2} \stackrel{Q}{+} \frac{2}{N} (K \quad ^{Q}1)^{2};$$

where
$$_{1}C^{Q} = {}_{1}C^{\frac{Q}{2}} \cdot \frac{1}{\min\{1-\}^{3}}$$
, $C_{2}^{Q} = C_{2}^{Q} \cdot \frac{jSj \log^{2}(jSj)}{\min[1-)^{4}}$, $C_{3}^{Q} = C_{3}^{Q} \cdot \frac{jSj^{2} \log^{2}(jSj)}{\min[1-)^{4}}$, and C_{i}^{Q} , $i = 1; 2; 3$ are universal problem

independent constants. In addition, choosing = $\frac{8 \log(NT)}{T'_{Q}}$ and K = T=N, we have $E[kQ \land Qk^2]$ $T = \tilde{O} \frac{1}{N} : \frac{jSj^2 \log^2(jSj)}{\frac{5}{min}(1)^9} \text{ iterations.}$

According to Theorem 5.2, federated Q-learning Algorithm 3, similar to federated off-policy TD-learning, enjoys linear speedup, communication efficiency as well as privacy guarantees. We would like to emphasize that the update of Q-learning is non-linear. Hence the result of Theorem 5.2 cannot be derived from Theorems 4.1 and 5.1.

6. Generalized Federated Stochastic Approximation

In this section we study the convergence of a general federated stochastic approximation for contractive operators, FedSAM, which is presented in Algorithm 4. In this algorithm there are N agents i = 1;2;:::; N. At each time step t 0, each agent i maintains the parameter i 2, Rd. At time t = 0, all agents initialize their parameters with i = On Next, at time t 0, each agent i updates its param-eter as i = i + $_{t}G_{1}^{i}(^{i}; y_{t}^{i})$ i + $b^{i}(y_{i}^{i})$ + Here denotes the step size, and yi is a noise which is Marko-vian along the time t, but is independent across the agents i. This notion is defined more concretely in Assumption 6.4. We note that functions Gi(;) and bi() are allowed to be dependent on the agent i. This allows us to employ the convergence bound of FedSAM in order to derive the convergence bound of off-policy TD-learning with different behaviour policies across agents. In order to avoid diver-

gence, every K time steps we synchronize the parameters synchronization and hence more "accurate" updates, at the same time it results in a higher communication cost, which is not desirable. Hence, in order to determine the optimal choice of synchronization period, it is essential to characterize the dependence of the convergence on K. This is one of the results which we will derive in Theorem B.1. Finally, the algorithm samples Υ qT ς , where qTc (t) = $\frac{c_1}{r_0}$ and outputs A. This sampling scheme is essential for the convergence of overall algorithm. We further make some assumptions regarding the underlying process.

First, we assume that the expectation of Gⁱ(; yⁱ), geometrically converges to some function Gi() and the expectation of bⁱ(yⁱ) geometrically converges to 0. In particular, we have the following assumption.

Assumption 6.1. For every agent i, there exist a function Gⁱ() such that we have

$$\lim_{t \mid 1} E[G^{i}(; y^{i})] = G^{i}()$$

$$\lim_{t \mid 1} E[b^{i}(y_{t}^{i})] = 0:$$
(6)

Algorithm 4 Federated Stochastic Approximation with Markovian Noise (FedSAM)

```
1: Input: C_{FSAM}; T;_0; K;_2: i

= _0 for all i = 1; :::; N. 3: for t

= _0 to _1 1 do

4: i_{t+1} = i_t + G^i(i; y^i)_{t} + b_t^i(y^i)_{t}; 8_t^i 2 [N]

5: if t + 1 \mod K = 0 then

6: i_{t+1} = i_{t+1}, \frac{1}{N} = \sum_{j=1}^{N} i_{t+1} + i_t +
```

Furthermore, there exists m_1 ; m_2 0 and 2 [0; 1), such that for every i = 1; 2; ...; N,

$$kG^{i}() = E[G^{i}(; y^{i})]k_{c} m_{1}kk_{c}^{t}$$

$$kE[b^{i}(y_{+}^{i})]k_{c} m_{2}^{t};$$
(7)

where k kc is a given norm.

Next, we assume a contraction property on the expected operator $G^{i}()$.

Assumption 6.2. We assume all expected operators $G^{i}()$ are contraction mappings with respect to k k_{c} with contraction factor $_{c}$ 2 (0; 1). That is, for all i = 1; 2; ...; N,

$$kG^{i}(_{1})$$
 $G^{i}(_{2})k_{c}$ $_{c}k_{1}$ $_{2}k_{c}$; 8_{1} ; $_{2}$ 2 R^{d} :

Next, we consider some Lipschitz and boundedness properties on $G^{i}(;)$ and $b^{i}()$.

Assumption 6.3. For all i = 1; :::; N, there exist constants A_1, A_2 and B such that

- 1. $kG^{i}(1; y^{i}) = G^{i}(2; y^{i})k_{c} = A_{1}k_{1} = 2k_{c}$, for all $A_{1}; A_{2}; y^{i}$.
- 2. $kG^{i}(; y^{i})k_{c}$ A₂kk_c for all; y^{i} . 3.

$$kb^{i}(v^{i})k_{c}$$
 B for all v^{i} .

Remark. By Assumption 6.2 and due to the Banach fixed point theorem, $G^{i}()$ has a unique fixed point for all i=1;2;:::;N. Furthermore, by Assumption 6.3, we have $G^{i}(0;y)=0$. Hence the point 0 is the unique fixed point of $G^{i}()$.

Finally, we impose an assumption on the random data y_t^i Assumption 6.4. We assume that the Markovian noise y_t^i (Markovian with respect to time t) is independent across agents i. In other words, for all measurable functions f () and g(), we assume the following

$$E_{t-r}[f(y_t^i) g(y_t)] = E_{t-r}[f(y_t)] E_{t-r}[g(y_t)];$$
 for all $r = j$.

Theorem 6.1 states the convergence of Algorithm 4.

Theorem 6.1. Consider the federated stochastic approximation Algorithm 4 with $c_{FSAM}=1$ $\frac{^{\prime 2}-2}{^{\prime 2}}$ (0; 1) ($^{\prime 2}$ is defined in Equation (14) in the appendix), and synchronization frequency K . Denote $_t=\frac{1}{N}$ $_{i=1}^{N}$ $_{i'}$ tand consider $_{i'}$ as the output of this algorithm after T iterations. Assume = d2 log $_{i'}$ e. For T $_{i'}$ e and small enough step size , we have

$$\begin{split} & E[k_T^{} \, k^2] \, C_1 \, c^T \, \underbrace{f^{+1}}_{F \, S \, A \, M} + C_2 \, + \, C_3 \, (K_{}^{} \, 2 \, 1)^2 \, + \, C_4^{\, 32}; \qquad (8) \\ & \text{where } C_i^{}, \, i \, = \, 1; 2; 3; 4 \, \text{ are some constants which are specified precisely in Appendix B, and are independent of } \\ & K; ; N \, . \, & \text{Choosing} \, = \, \frac{8 \, log \, (N_{} \, T)}{T \, and} \, K \, = \, T \, = N \, , \end{split}$$

we get $T = O(\frac{1}{N})$ sample complexity for achieving $E[k_x k_c]^2$.

Theorem 6.1 establishes the convergence of A to zero in the expected mean-squared sense. The first term in (8) converges geometrically to zero as T grows. The second term is proportional to similar as (3). However, the number of agents N in the denominator ensures linear speedup, meaning that for small enough (such that =N is the dominant term), the sample complexity of each individual agent, relative to a centralized system, is reduced by a factor of N. The third term has quadratic dependence on , and is zero when we have perfect synchronization, i.e. K = 1. The last term is proportional to ³, and has the weakest dependency on the step size . For K > 1 we can merge the last two terms by upper bounding 3 2. The current upper bound, however, is tighter since with K = 1 (i.e. perfect synchronization) we have no term in the order ². Note that similar bounds (sans the last ³ term) have been established for the simpler i.i.d. noise case in the federated setting (Khaled et al., 2020; Koloskova et al., 2020). Consequently, we achieve the same sample complexity results for the more general federated setting with Markov noise. Remark. The bound in Theorem 6.1 holds only after T > maxfK+; 2g and for all synchronization periods K 1. At K = 1 the third term in the bound goes away, and we will be left only with the first order term, which is linearly decreasing with respect to the number of agents N, and the third order term $\mathcal{O}(^3)$. The last term, however, is not tight and can be further improved to be of the order O(i); i > 3. However, for that we need to assume larger, which means the bound only hold after a longer waiting time. In particular, by choosing = dr log e, we can get $O(2^{r-1})$ for the last term (see the proof of Lemma B.2).

Acknowledgment

This work was partially supported by NSF awards CCF-1944993, CCF-2045694, CNS-2112471, CMMI-2112533, EPCN-2144316, and an award from Raytheon technologies.

References

- Akkaya, I., Andrychowicz, M., Chociej, M., Litwin, M., McGrew, B., Petron, A., Paino, A., Plappert, M., Powell, G., Ribas, R., Schneider, J., Tezak, N., Tworek, J., Welinder, P., Weng, L., Yuan, Q., Zaremba, W., and Zhang, L. Solving Rubik's cube with a robot hand. Preprint arXiv:1910.07113, 2019.
- Baird, L. Residual algorithms: Reinforcement learning with function approximation. In Machine Learning Proceedings 1995, pp. 30–37. Elsevier, 1995.
- Banach, S. Sur les opérations dans les ensembles abstraits et leur application aux équations intégrales. Fund. math, 3(1):133–181, 1922.
- Beck, A. First-order methods in optimization, volume 25. SIAM, 2017.
- Beck, C. L. and Srikant, R. Error bounds for constant step-size Q-learning. Syst. Control. Lett., 61:1203–1208, 2012a.
- Beck, C. L. and Srikant, R. Error bounds for constant step-size Q-learning. Systems & control letters, 61(12): 1203–1208, 2012b.
- Beck, C. L. and Srikant, R. Improved upper bounds on the expected error in constant step-size Q-learning. In 2013 American Control Conference, pp. 1926–1931. IEEE, 2013.
- Bertsekas, D. P. and Tsitsiklis, J. N. Neuro-dynamic programming. Athena Scientific, 1996a.
- Bertsekas, D. P. and Tsitsiklis, J. N. Neuro-dynamic programming. Athena Scientific, 1996b.
- Bertsekas, D. P., Bertsekas, D. P., Bertsekas, D. P., and Bertsekas, D. P. Dynamic programming and optimal control, volume 2. Athena scientific Belmont, MA, 1995.
- Bhandari, J., Russo, D., and Singal, R. A finite time analysis of temporal difference learning with linear function approximation. In Conference on learning theory, pp. 1691–1692. PMLR, 2018.
- Borkar, V. S. Stochastic approximation: a dynamical systems viewpoint, volume 48. Springer, 2009.
- Borkar, V. S. and Meyn, S. P. The ODE method for convergence of stochastic approximation and reinforcement learning. SIAM Journal on Control and Optimization, 38 (2):447–469, 2000.
- Chen, T. and Giannakis, G. B. Bandit convex optimization for scalable and dynamic iot management. IEEE Internet of Things Journal, 6(1):1276–1286, 2018.

- Chen, T., Zhang, K., Giannakis, G. B., and Basar, T. Communication-efficient policy gradient methods for distributed reinforcement learning. IEEE Transactions on Control of Network Systems, 2021a.
- Chen, Z., Zhang, S., Doan, T. T., Maguluri, S. T., and Clarke, J.-P. Finite-sample analysis of nonlinear stochastic approximation with applications in reinforcement learning. Under review by Automatica, Preprint arXiv:1905.11425, 2019a.
- Chen, Z., Zhang, S., Doan, T. T., Maguluri, S. T., and Clarke, J.-P. Performance of Q-learning with linear function approximation: Stability and finite-time analysis. In OptRL Workshop at NeuRIPS 2019, 2019b.
- Chen, Z., Maguluri, S. T., Shakkottai, S., and Shanmugam, K. Finite-sample analysis of stochastic approximation using smooth convex envelopes. Under Review at Mathematics of Operations Research, Preprint arXiv:2002.00874, 2020a.
- Chen, Z., Maguluri, S. T., Shakkottai, S., and Shanmugam, K. Finite-sample analysis of stochastic approximation using smooth convex envelopes. In Advances in Neural Information Processing Systems, 2020b. URL https://proceedings.neurips.cc/paper/2020/file/5d44ee6f2c3f71b73125876103c8f6c4-Paper.pdf.
- Chen, Z., Khodadadian, S., and Maguluri, S. T. Finite-Sample Analysis of Off-Policy Natural Actor-Critic with Linear Function Approximation. Preprint arXiv:2105.12540, 2021b. Submitted to NeurIPS 2021.
- Chen, Z., Maguluri, S. T., Shakkottai, S., and Shanmugam, K. A Lyapunov Theory for Finite-Sample Guarantees of Asynchronous Q-Learning and TD-Learning Variants. Under review by JMLR, Preprint arXiv:2102.01567, 2021c.
- Chen, Z., Maguluri, S. T., Shakkottai, S., and Shanmugam, K. Finite-Sample Analysis of Off-Policy TD-Learning via Generalized Bellman Operators. Preprint arXiv:2106.12729, 2021d.
- Chen, Z., Zhou, Y., and Chen, R. Multi-agent off-policy to learning: Finite-time analysis with near-optimal sam-ple complexity and communication complexity. arXiv preprint arXiv:2103.13147, 2021e.
- Chen, Z., Zhou, Y., Chen, R., and Zou, S. Sample and communication-efficient decentralized actor-critic algorithms with finite-time analysis. arXiv preprint arXiv:2109.03699, 2021f.

- Dalal, G., Szörényi, B., Thoppe, G., and Mannor, S. Finite sample analysis for TD(0) with function approximation. In Thirty-Second AAAI Conference on Artificial Intelligence, 2018.
- Doan, T., Maguluri, S., and Romberg, J. Finite-Time Analysis of Distributed TD(0) with Linear Function Approximation on Multi-Agent Reinforcement Learning. In International Conference on Machine Learning, pp. 1626–1635, 2019.
- Duan, Y., Schulman, J., Chen, X., Bartlett, P. L., Sutskever, I., and Abbeel, P. Rl ²: Fast reinforcement learning via slow reinforcement learning. arXiv preprint arXiv:1611.02779, 2016.
- Espeholt, L., Soyer, H., Munos, R., Simonyan, K., Mnih, V., Ward, T., Doron, Y., Firoiu, V., Harley, T., Dunning, I., Legg, S., and Kavukcuoglu, K. IMPALA: Scalable Distributed Deep-RL with Importance Weighted Actor-Learner Architectures. In International Conference on Machine Learning, pp. 1406–1415, 2018.
- Even-Dar, E. and Mansour, Y. Learning Rates for Q-Learning. J. Mach. Learn. Res., 5:1–25, 2004. ISSN 1532-4435.
- Foerster, J. N., Assael, Y. M., De Freitas, N., and Whiteson, S. Learning to communicate with deep multi-agent reinforcement learning. arXiv preprint arXiv:1605.06676, 2016.
- Gottesman, O., Johansson, F., Komorowski, M., Faisal, A., Sontag, D., Doshi-Velez, F., and Celi, L. A. Guidelines for reinforcement learning in healthcare. Nature medicine, 25(1):16–18, 2019.
- Gottesman, O., Futoma, J., Liu, Y., Parbhoo, S., Celi, L., Brunskill, E., and Doshi-Velez, F. Interpretable off-policy evaluation in reinforcement learning by highlighting influential transitions. In International Conference on Machine Learning, pp. 3658–3667. PMLR, 2020.
- Gu, S., Holly, E., Lillicrap, T., and Levine, S. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In 2017 IEEE international conference on robotics and automation (ICRA), pp. 3389–3396. IEEE, 2017.
- Hu, B. and Syed, U. A. Characterizing the exact behaviors of temporal difference learning algorithms using markov jump linear system theory. arXiv preprint arXiv:1906.06781, 2019.
- Islam, R., Henderson, P., Gomrokchi, M., and Precup, D. Reproducibility of benchmarked deep reinforcement learning tasks for continuous control. arXiv preprint arXiv:1708.04133, 2017.

- Jaakkola, T., Jordan, M. I., and Singh, S. P. Convergence of stochastic iterative dynamic programming algorithms. In Advances in neural information processing systems, pp. 703–710, 1994.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis,
 M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode,
 G., Cummings, R., et al. Advances and open problems in federated learning. Preprint arXiv:1912.04977, 2019.
- Kalashnikov, D., Irpan, A., Pastor, P., Ibarz, J., Herzog, A., Jang, E., Quillen, D., Holly, E., Kalakrishnan, M., Vanhoucke, V., and Levine, S. Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation. Preprint arXiv:1806.10293, 2018.
- Khaled, A., Mishchenko, K., and Richtárik, P. Tighter theory for local sgd on identical and heterogeneous data. In International Conference on Artificial Intelligence and Statistics, pp. 4519–4529. PMLR, 2020.
- Khodadadian, S., Chen, Z., and Maguluri, S. T. Finite-Sample Analysis of Off-Policy Natural Actor-Critic Algorithm. In International Conference on Machine Learning, 2021.
- Kober, J., Bagnell, J. A., and Peters, J. Reinforcement learning in robotics: A survey. The International Journal of Robotics Research, 32(11):1238–1274, 2013.
- Koloskova, A., Loizou, N., Boreiri, S., Jaggi, M., and Stich, S. A unified theory of decentralized sgd with changing topology and local updates. In International Conference on Machine Learning, pp. 5381–5393. PMLR, 2020.
- Konda, V. R. and Tsitsiklis, J. N. Actor-critic algorithms. In Advances in neural information processing systems, pp. 1008–1014, 2000.
- Konečný, J., McMahan, H. B., Ramage, D., and Richtárik, P. Federated optimization: Distributed machine learning for on-device intelligence. arXiv preprint arXiv:1610.02527, 2016.
- Lakshminarayanan, C. and Szepesvari, C. Linear stochastic approximation: How far does constant step-size and iterate averaging go? In International Conference on Artificial Intelligence and Statistics, pp. 1347–1355, 2018.
- Levine, S., Kumar, A., Tucker, G., and Fu, J. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. Preprint arXiv:2005.01643, 2020.
- Li, G., Wei, Y., Chi, Y., Gu, Y., and Chen, Y. Sample Complexity of Asynchronous Q-Learning: Sharper Analysis and Variance Reduction. Advances in neural information processing systems, 2020.

- Li, X., Huang, K., Yang, W., Wang, S., and Zhang, Z. On the convergence of fedavg on non-iid data. arXiv preprint arXiv:1907.02189, 2019.
- Liang, X., Liu, Y., Chen, T., Liu, M., and Yang, Q. Federated transfer reinforcement learning for autonomous driving. arXiv preprint arXiv:1910.06001, 2019.
- Lim, H.-K., Kim, J.-B., Heo, J.-S., and Han, Y.-H. Federated reinforcement learning for training control policies on multiple iot devices. Sensors, 20(5):1359, 2020.
- Liu, B., Liu, J., Ghavamzadeh, M., Mahadevan, S., and Petrik, M. Finite-sample analysis of proximal gradient TD algorithms. In Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence, pp. 504–513, 2015.
- Liu, B., Wang, L., and Liu, M. Lifelong federated reinforcement learning: a learning architecture for navigation in cloud robotic systems. IEEE Robotics and Automation Letters, 4(4):4555–4562, 2019.
- Liu, R. and Olshevsky, A. Distributed td (0) with almost no communication. arXiv preprint arXiv:2104.07855, 2021.
- Liu, Y., Gottesman, O., Raghu, A., Komorowski, M., Faisal, A. A., Doshi-Velez, F., and Brunskill, E. Representation Balancing MDPs for Off-policy Policy Evaluation. Advances in Neural Information Processing Systems, 31: 2644–2653, 2018.
- Maddern, W., Pascoe, G., Linegar, C., and Newman, P. 1 year, 1000 km: The oxford robotcar dataset. The International Journal of Robotics Research, 36(1):3–15, 2017.
- Maei, H. R. Convergent actor-critic algorithms under off-policy training and function approximation. Preprint arXiv:1802.07842, 2018.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. Communication-efficient learning of deep networks from decentralized data. In Artificial Intelligence and Statistics, pp. 1273–1282. PMLR, 2017.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. Asynchronous methods for deep reinforcement learning. In International conference on machine learning, pp. 1928– 1937, 2016.
- Mothukuri, V., Parizi, R. M., Pouriyeh, S., Huang, Y., Dehghantanha, A., and Srivastava, G. A survey on security and privacy of federated learning. Future Generation Computer Systems, 115:619–640, 2021.
- Nadiger, C., Kumar, A., and Abdelhak, S. Federated reinforcement learning for fast personalization. In 2019

- IEEE Second International Conference on Artificial Intelligence and Knowledge Engineering (AIKE), pp. 123–127. IEEE, 2019.
- Nair, A., Srinivasan, P., Blackwell, S., Alcicek, C., Fearon, R., De Maria, A., Panneershelvam, V., Suleyman, M., Beattie, C., Petersen, S., et al. Massively parallel methods for deep reinforcement learning. arXiv preprint arXiv:1507.04296, 2015.
- Nguyen, D. C., Ding, M., Pathirana, P. N., Seneviratne, A., Li, J., and Poor, H. V. Federated learning for internet of things: A comprehensive survey. arXiv preprint arXiv:2104.07914, 2021.
- Puterman, M. L. Markov decision processes: discrete stochastic dynamic programming. John Wiley & Sons, 2014.
- Qi, J., Zhou, Q., Lei, L., and Zheng, K. Federated reinforcement learning: Techniques, applications, and open challenges. arXiv preprint arXiv:2108.11887, 2021.
- Qu, G. and Wierman, A. Finite-Time Analysis of Asynchronous Stochastic Approximation and Q-Learning. In Conference on Learning Theory, pp. 3185–3205. PMLR, 2020.
- Qu, Z., Lin, K., Kalagnanam, J., Li, Z., Zhou, J., and Zhou, Z. Federated learning's blessing: Fedavg has linear speedup. arXiv preprint arXiv:2007.05690, 2020.
- Rakhlin, A., Shamir, O., and Sridharan, K. Making gradient descent optimal for strongly convex stochastic optimization. In Proceedings of the 29th International Coference on International Conference on Machine Learning, pp. 1571–1578, 2012.
- Recht, B., Re, C., Wright, S., and Niu, F. Hogwild!: A lock-free approach to parallelizing stochastic gradient descent. In Shawe-Taylor, J., Zemel, R., Bartlett, P., Pereira, F., and Weinberger, K. Q. (eds.), Advances in Neural Information Processing Systems, volume 24. Curran Associates, Inc., 2011. URL https://proceedings.neurips.cc/paper/2011/file/218aOaefd1d1a4be656O1cc6ddc152Oe-Paper.pdf.
- Ren, J., Wang, H., Hou, T., Zheng, S., and Tang, C. Federated learning-based computation offloading optimization in edge computing-supported internet of things. IEEE Access, 7:69194–69201, 2019.
- Shah, D. and Xie, Q. Q-learning with nearest neighbors. In Advances in Neural Information Processing Systems, pp. 3111–3121, 2018.

- Shalev-Shwartz, S., Shammah, S., and Shashua, A. Safe, multi-agent, reinforcement learning for autonomous driving. Preprint arXiv:1610.03295, 2016.
- Shalev-Shwartz, S. et al. Online learning and online convex optimization. Foundations and Trends® in Machine Learning, 4(2):107–194, 2012.
- Shao, K., Tang, Z., Zhu, Y., Li, N., and Zhao, D. A survey of deep reinforcement learning in video games. arXiv preprint arXiv:1912.10944, 2019.
- Shen, H., Zhang, K., Hong, M., and Chen, T. Asynchronous advantage actor critic: Non-asymptotic analysis and linear speedup. arXiv preprint arXiv:2012.15511, 2020.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., van den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., Dieleman, S., Grewe, D., Nham, J., Kalchbrenner, N., Sutskever, I., Lillicrap, T., Leach, M., Kavukcuoglu, K., Graepel, T., and Hassabis, D. Mastering the game of Go with deep neural networks and tree search. Nature, 529(7587):484, 2016.
- Spiridonoff, A., Olshevsky, A., and Paschalidis, I. C. Communication-efficient sgd: From local sgd to one-shot averaging. In Advances in Neural Information Processing Systems, volume 34, 2021.
- Srikant, R. and Ying, L. Finite-time error bounds for linear stochastic approximation and TD learning. In Conference on Learning Theory, pp. 2803–2830. PMLR, 2019.
- Stich, S. U. Local sgd converges fast and communicates little. In International Conference on Learning Representations, 2018.
- Sun, J., Wang, G., Giannakis, G. B., Yang, Q., and Yang, Z. Finite-time analysis of decentralized temporal-difference learning with linear function approximation. In International Conference on Artificial Intelligence and Statistics, pp. 4485–4495. PMLR, 2020.
- Sutton, R. S. Learning to predict by the methods of temporal differences. Machine learning, 3(1):9–44, 1988.
- Sutton, R. S. and Barto, A. G. Reinforcement learning: An introduction. MIT press, 2018.
- Tadić, V. On the convergence of temporal-difference learning with linear function approximation. Machine learning, 42(3):241–267, 2001.
- Truex, S., Baracaldo, N., Anwar, A., Steinke, T., Ludwig, H., Zhang, R., and Zhou, Y. A hybrid approach to privacy-preserving federated learning. In Proceedings of the 12th ACM Workshop on Artificial Intelligence and Security, pp. 1–11, 2019.

- Tsitsiklis, J. N. Asynchronous stochastic approximation and Q-learning. Machine learning, 16(3):185–202, 1994.
- Tsitsiklis, J. N. and Van Roy, B. Analysis of temporal-difference learning with function approximation. In Advances in neural information processing systems, pp. 1075–1081, 1997.
- Wai, H.-T. On the convergence of consensus algorithms with markovian noise and gradient bias. In 2020 59th IEEE Conference on Decision and Control (CDC), pp. 4897–4902. IEEE, 2020.
- Wainwright, M. J. Stochastic approximation with conecontractive operators: Sharp '1-bounds for Q-learning. Preprint arXiv:1905.06265, 2019.
- Wang, G., Lu, S., Giannakis, G., Tesauro, G., and Sun, J. Decentralized TD Tracking with Linear Function Approximation and its Finite-Time Analysis. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 13762–13772. Curran Associates, Inc., 2020a. URL https://proceedings.neurips.cc/paper/2020/file/9ec51f6eb240fb631a35864e13737bca-Paper.pdf.
- Wang, J. and Joshi, G. Cooperative sgd: A unified framework for the design and analysis of local-update sgd algorithms. Journal of Machine Learning Research, 22(213): 1–50, 2021.
- Wang, X., Wang, C., Li, X., Leung, V. C., and Taleb, T. Federated deep reinforcement learning for internet of things with decentralized cooperative edge caching. IEEE Internet of Things Journal, 7(10):9441–9455, 2020b.
- Watkins, C. J. and Dayan, P. Q-learning. Machine learning, 8(3-4):279–292, 1992.
- Xu, M., Peng, J., Gupta, B., Kang, J., Xiong, Z., Li, Z., and Abd El-Latif, A. A. Multi-agent federated reinforcement learning for secure incentive mechanism in intelligent cyber-physical systems. IEEE Internet of Things Journal, 2021.
- Yang, Q., Liu, Y., Cheng, Y., Kang, Y., Chen, T., and Yu,H. Federated learning. Synthesis Lectures on Artificial Intelligence and Machine Learning, 13(3):1–207, 2019.
- Yun, W. J., Park, S., Kim, J., Shin, M., Jung, S., Mohaisen, A., and Kim, J.-H. Cooperative multi-agent deep reinforcement learning for reliable surveillance via autonomous multi-uav control. IEEE Transactions on Industrial Informatics, 2022.

- Yurtsever, E., Lambert, J., Carballo, A., and Takeda, K. A survey of autonomous driving: Common practices and emerging technologies. IEEE Access, 8:58443–58469, 2020.
- Zeng, S., Doan, T. T., and Romberg, J. Finite-time analysis of decentralized stochastic approximation with applications in multi-agent and multi-task learning. arXiv preprint arXiv:2010.15088, 2020.
- Zhang, S., Liu, B., Yao, H., and Whiteson, S. Provably convergent two-timescale off-policy actor-critic with function approximation. In International Conference on Machine Learning, pp. 11204–11213. PMLR, 2020.
- Zhang, S. Q., Lin, J., and Zhang, Q. A multi-agent reinforcement learning approach for efficient client selection in federated learning. arXiv preprint arXiv:2201.02932, 2022.
- Zhao, L., Ran, Y., Wang, H., Wang, J., and Luo, J. Towards cooperative caching for vehicular networks with multilevel federated reinforcement learning. In ICC 2021-IEEE International Conference on Communications, pp. 1–6. IEEE, 2021.
- Zhuo, H. H., Feng, W., Lin, Y., Xu, Q., and Yang, Q. Federated deep reinforcement learning. arXiv preprint arXiv:1901.08277, 2019.

Appendices

The appendices are organized as follows. In section A we discuss the lower bound on the convergence of general stochastic approximation. In Section B we derive the convergence bound of FedSAM algorithm. Next, we employ the results in Section B to derive the convergence bounds of federated TD-learning in Section C and federated Q-learning in Section D.

A. Lower Bound on the Convergence of General Stochastic Approximation

In this section we discuss the convergence of general stochastic approximation. In this discussion we provide a simple stochastic approximation with iid noise which can give insight into the general convergence bound in (3). In particular, we show that the convergence bound in (3) is tight and cannot be improved.

Consider a one dimensional random variable X with zero mean E[X] = 0 and bounded variance $E[X^2] = 2$. Consider the following update

$$x_{t+1} = x_t + (X_t x_t); t 0;$$
 (9)

where we start with some fixed deterministic x_0 and X_t is a an iid sample of the random variable X. It is easy to see that the update (9) is a special case of the update of the general stochastic approximation with the fixed point x = 0.

By expanding the update (9), we have

$$x_t = (1)^t x_0 + (1)^{t k - 1} X_k$$
:

Hence, we have

$$x_t^2 = (1)^{2t}x^2 d$$
 $x_t^2 = (1)^{t}x^2 d$ $x_t^2 = (1)^{t}x^2$

Taking expectation on both sides, and using the zero mean property of X_k , we have

It is clear that (10) has the same form as the bound in (3) with T_1 as the geometric term which converges to zero as $t \,!\, 1$, and T_2 term proportional to the step size . In addition, note that in the above derivation, we did not use any inequality, and hence the bounds in (10) as well as (3) are tight.

B. Analysis of Federated Stochastic Approximation

First, we restate Theorem 6.1 with explicit expressions of the different constants.

Theorem B.1. Consider the FedSAM Algorithm 4 with $c_{FSAM} = 1 \frac{\frac{2}{2}(2)}{2}$ is defined in (14)), and suppose Assumptions 6.1, 6.2, 6.3, 6.4 are satisfied. Consider small enough step size which satisfies the assumptions in (21), (23), (32), (35). Furthermore, denote = d2 log e, and take large enough T such that T > maxfK + ; 2g. Then, the output of the FedSAM Algorithm 4, $_{T}$, $\frac{1}{N}$ P N = $_{1}$ i , satisfies

$$E[k_{1}k_{2}^{2}]_{c}C_{1}^{1}1^{-\frac{2+1}{2}}+C_{2}-N+C_{3}(K^{2}1)^{2}+C_{4}^{32}; \qquad (11)$$

where $C_1 = 16 i l_{cm} M_0 (log \frac{1}{e} + \frac{1}{\frac{1}{2}})$, $M_0 = \frac{1}{l_{cm}^2} \frac{1}{C_1^2} B + (A_2 + 1) k_0 k_c + \frac{B}{2C_1} + k_0 k_c$, $C_2 = 8 u_{cm} C_{82} + {}_2 + C_{12} = {}_2$, and $C_3 = 80 B^2 C_{17} u_{cm} 1 + {}_{28(1)} = {}_{2m}$ and $C_4 = 8 u_{cm} C_7 + C_{11} + 0.5 C_3 C_9 + {}_2 C_3 C_{10} + 2 C_1 C_3 C_{10} + 3 A_1 C_3 + C_{13} = {}_2$. Here the constants $u_{cm}; I_{cm}; {}'_{2}; C_{1}; C_{3}; C_{7}; C_{8}; C_{9}; C_{10}; C_{11}; C_{12}; C_{13}; C_{17}; \\ \underbrace{\text{B.}}_{\text{are problem dependent constants which are defined in}}_{\text{constants which are defined in}}$ the following proposition and lemmas.

Next, we characterize the sample complexity of the FedSAM Algorithm 4, where we establish a linear speedup in the convergence of the algorithm.

Corollary B.1.1. Consider FedSAM Algorithm 4 with fixed number of iterations T and step size = $\frac{8 \log(NT)}{2T}$. Suppose T is large enough, such that satisfies the requirements of the step size in Theorem B.1 and T > 4. Also take K = T=N. Then we have $E[k_T k_c]^2$ after $T = O_N$ iterations.

Corollary B.1.1 establishes the sample and communication complexity of FedSAM Algorithm 4. The O(1=(N)) sample complexity shows the linear speedup with increasing N. Another aspect of the cost is the number of communications required between the agents and the central server. According to Corollary B.1.1, we need $T=N = \tilde{O}(N)$ rounds of communications in order to reach an -optimal solution. Hence, even in the presence of Markov noise, the required number of communications is independent of the desired final accuracy, and grows linearly with the number of agents. Our result generalizes the existing result achieved for the simpler i.i.d. noise case in (Khaled et al., 2020; Spiridonoff et al., 2021).

In the following sections, we discuss the proof of Theorem B.1. In Section B.1, we introduce some notations and preliminary results to facilitate our Lyapunov-function based analysis. Next, in Section B.2, we state some primary propositions, which are then used to prove Theorem B.1. In Sections B.3, B.4, B.5, we prove the aforementioned propositions. Along the way, we state several intermediate results, which are stated and proved in Sections B.6 and B.7, respectively.

Throughout the appendix we have several sets of constants. The constants C_i; i = 1;:::;17 are problem dependent constants which we define recursively. The final constants which appears in the resulting bound in Theorem B.1 are shown as C_i ; i = 1; ...; 4. Finally, the constant c is used in the sampling of the time step T. $^{\land}$

B.1. Preliminaries

We define the following notations:

- t, $\frac{1}{N}$ $i = \frac{P}{1}i$; virtual sequence of average (across agents) parameter.
- t = t; : 1:; t = N : set of local parameters at individual nodes. $Y_t = y_t$; ::; y_t : Markov chains at individual nodes.
- i: the stationary distribution of y i as t! 1.
- $G(t; Y_t)$, $\frac{1}{N} i = \frac{P}{N} i G^i(t; y_t^i)$: taverage of the noisy local operators at the individual local parameters.

 $G(t; Y_t)$, $\frac{1}{T} P_{N} i = {}_{1}^{N} G^{i}(t; y^{i})$: average of the noisy local operators at the average parameter.

• $b(Y_t)$, $\frac{1}{N} \stackrel{P}{i} \stackrel{N}{=} _1 b^i(y^i)$: average of Markovian noise.

- G(t), $\frac{1}{N} Pi = \frac{P}{1} G^i(t)$; average of expected operators evaluated at the local parameter. G(t), $\frac{P}{N} I = \frac{N}{1} G^i(t)$: average of expected operator evaluated at the average parameter.
- = 0 is the unique fixed point satisfying = Gⁱ() for all i = 1;:::; N. Note that this follows directly from the assumptions. In particular, by Assumption 6.2, Gⁱ() is a contraction, and hence by Banach fixed point theorem (Banach, 1922), there exist a unique fixed point of this operator. Furthermore, by Assumption 6.3, we have

$$kG^{i}(0)k_{c} = kE_{v^{i}}[G^{i}(0;y)]k_{c} E_{v^{i}}kG^{i}(0;y)k_{c} A_{2}E_{v^{i}}[k0k_{c}] = 0$$
:

Hence $G^{i}(0) = 0$ and 0 is the unique fixed point of the operator $G^{i}()$.

•
$$_{t} \stackrel{i}{=} k_{t} \quad _{t}k_{c}; \quad _{t} = \underset{N}{\overset{P}{=}} i = \underset{1}{\overset{N}{=}} _{t}; \quad _{i} \qquad \qquad \underset{i}{\overset{1}{=}} \quad _{N} \qquad i$$
 $_{t} = \underset{N}{\overset{P}{=}} _{i} = \underset{1}{(}_{t})^{2}$: measures of synchronization error.

Throughout this proof we assume k k_c as some given norm. $E_t[]$, $E[jF_t]$, where F_t is the sigma-algebra generated by $f_r g_r^{-\frac{j-1}{2}, \dots, N} t$. Unless specified otherwise, k k denotes the Euclidean norm.

Generalized Moreau Envelope: Consider the norm k k_c which appears in Assumptions 6.1-6.3. Square of this norm need not be smooth. Inspired by (Chen et al., 2020a), we use the Generalized Moreau Envelope as a Lyapunov function for the analysis of the convergence of Algorithm 4. The Generalized Moreau Envelope of f () with respect to g(), for > 0, is defined as

$$M_f^{;g}(x) = \min_{u \ge R^d} f(u) + \frac{1}{-}g(x - u)$$
 : (12)

Let $f(x) = \frac{1}{2}kxk_c^2$ and $g(x) = \frac{1}{2}kxk_s^2$, which is L-smooth with respect to k k_s norm. For this choice of f; g, M_r g () is essentially a smooth approximation to f, which is henceforth denoted with the simpler notation M (). Also, due to the equivalence of norms, there exist I_{cs} ; $u_{cs} > 0$ such that

$$I_{cs}k k_s k k_c u_{cs}k k_s: (13)$$

We next summarize the properties of M () in the following proposition, which were established in (Chen et al., 2020a). Proposition B.1 ((Chen et al., 2020b)). The function M () satisfies the following properties.

- (1) M () is convex, and $\frac{L}{2}$ -smooth with respect to k k_s. That is, M (y) M (x) + hrM (x); y xi + $\frac{L}{2}$ yk_s fo² all x; y 2 R^d.
- (2) There exists a norm, denoted by k k_m , such that M (x) = ${}_{2}^{4}kxk_m^2$.
- (3) Let $'_{cm} = (1 + '^2_{cs})^{1=2}$ and $u_{cm} = (1 + u^2_{cs})^{1=2}$. Then it holds that $'_{cm}k \ k_m \ k \ k_c \ u_{cm}k \ k_m$.

By Proposition B.1, we can use M () as a smooth surrogate for $_2$ k $_c$. Furthermore, we denote

$$'_{1} = \frac{1 + u_{cs}^{2}}{1 + c_{cs}^{2}}; \quad '_{2} = 1 c_{c}^{1 = 2}; \quad \text{and} \quad '_{3} = \frac{114L(1 + u_{cs}^{2})}{c_{cs}^{2}}:$$
 (14)

Note that by choosing > 0 small enough, we can ensure $'_2$ 2 (0; 1).

B.2. Proof of Theorem B.1

In this section, first we state three key results (Propositions B.2, B.3, B.4). These are then used to prove Theorem B.1. The first step of the proof is to characterize the one-step drift of the Lyapunov function M (), with the parameters generated by the FedSAM Algorithm 4, which is formally stated in the following proposition.

Proposition B.2 (One-step drift - I). Consider the update of the FedSAM Algorithm 4. Suppose the assumptions 6.1, 6.2, 6.3, and 6.4 are satisfied. Consider = d2 log e and t 2, we have

$$E_{t-2}[M(t+1)]$$

$$+ \frac{L}{\left|\frac{l_{cs}^{2}}{2}\right|^{2}} + A_{2} + A_{2} + \frac{c l_{m}^{2}}{2} + (A_{1} + 1) - \frac{c m^{2} u^{2}}{2} + 3m_{1}^{2} + 2 + 3m_{1}^{2} + 2 + k^{2}$$

$$(16)$$

$$+ \frac{L}{|\frac{I_{cs}^{2}}{2}|} \frac{L}{2|\frac{I_{cs}^{2}}{2}|} + A_{2} \frac{1 + \frac{cu_{n}^{2}}{2}}{2} + (A_{1} + 1) \frac{-c\frac{3}{2}u^{2}}{2} + \frac{2}{3} 2 + 3m_{1}^{2} E_{t} \frac{2}{2}[k_{t} t k^{2}]}{2}$$

$$+ \frac{L}{|\frac{L}{cs}|^{2}} \frac{2^{m}}{2} + \frac{3}{2}u_{cs}^{2}(A_{1} + 31) + A_{2} + \frac{2}{2}u_{m}^{m}}{2} + \frac{3}{2}u_{2}^{2}E_{t} \frac{2[A_{1}^{2}]}{2} + \frac{2}{2}u_{m}^{2}}{2} + \frac{3}{2}u_{m}^{2}(A_{1} + 31) + A_{2} + \frac{3}{2}u_{m}^{2}(A$$

$$+ \frac{m_2^2 \mu_{cm}^2 L^2}{2_2 \frac{4}{cs}}$$

$$| \frac{1}{3} z$$
 (19)

$$+ \frac{1}{1 +} \int_{-C_8}^{C_7} \frac{1}{1 +} \sum_{k=1}^{C_7} \frac{1}{1 +} \sum_{k=1}^{C_7}$$

where 1, 2, 3, and 4 are arbitrary positive constants.

Proof. The proof of Proposition B.2 is presented in full detail in Section B.3.

Before discussing the bound in its full generality, we discuss a few special cases.

- synchronization Perfect (K with i.i.d. noise: since t = 0 for all t, = 0 (independence across time), and C₇ = 0 (see Lemmas B.2 and B.5 in Section B.6), the terms (15), (16) (17),(18), (19) will not appear in the bound, which is the form we get for centralized systems with i.i.d. noise (Rakhlin et al., 2012).
- Infrequent synchronization (K > 1) with i.i.d. noise: = 0, and C_7 = 0, the terms (15), (16),(18), (19) will not appear in the bound, which is the form we get for federated stochastic optimization with i.i.d. noise (Khaled et al., 2020).
- Perfect synchronization (K = 1) with Markov noise: since t = 0 for all t, the bound in Proposition B.2 generalizes the results in (Srikant & Ying, 2019; Chen et al., 2021c).

Next, we substitute the bound on k_t t k_c (Lemma B.6) and k_t t k^2 (Lemma B.7) to further bound the one-step drift. Establishing a tight bound for these two quantities are essential to ensure linear speedup.

Proposition B.3 (One-step drift - II). Consider the update of the FedSAM Algorithm 4. Suppose Assumptions 6.1, 6.2, 6.3, and 6.4 are satisfied. Define C $_{15}() = \frac{6 m_1 L u_{cm}^2}{I_{cs}^2} + \frac{6 L (A_2 + 1)^2 u_{m}^2}{2 I_{cs}^2} + 144 (A^2 + 1) C u_{cm}^2$. For

min
$$\frac{1}{360(A_2 + 1)'} \frac{1}{2C_{15}()}$$
; (21)

= d2 log e, and t 2, we have

$$E_{t-2}[M(t+1)] (1-t^2) E_{t-2}[M(t)] + C_{14}()^4 + C_{16}()^2 + C_{17} \times \frac{X}{N} E_{t-2}[$$

$$\text{where } C_{14}() = C_7 + C_{11} + 0:5C_3^2C_9^2 + C_3C_{10} + 2C_1C_3C_{10} + 3A_1C_3 + C_{13} + C_8\frac{u_{c\,D}^2}{l_{c\,D}^2} 2m_2^{22} \, ^2, \ C_{16}() = C_8\frac{u_{c\,D}^2B^2}{l_{c\,D}^2} + \frac{1}{2} + C_{12} \, ^2 \ \text{and} \ C_{17} = (3A_1C_3 + 8A^2C_4 + C_5 + C_6). \ \text{Here we define } C_9 = \frac{8}{l_{c\,D}}\frac{u_{c\,B}}{l_{c\,D}}, \ C_{10} = \frac{8m_2u_{c\,D}}{l_{c\,D}(1)^7}, \ C_{11} = \frac{8C_1^2C_3^2u_{c\,D}^2}{\frac{3}{6}}, \ C_{12} = \frac{8C_4u_{c\,D}^2B^2}{l_{c\,D}^2}, \ C_{13} = \frac{14C_4u_{c\,D}^2m_2^2}{l_{c\,D}^2(1)^7}.$$

Proof. The proof of Proposition B.3 is presented in full detail in Section B.4.

Conditional expectation $E_{t-2}[]$. The conditional expectation $E_{t-2}[]$ used in Proposition B.3 is essential when dealing with Markovian noise. The idea of using conditional expectation to deal with Markovian noise is not novel per se. In the previous work (Bhandari et al., 2018; Srikant & Ying, 2019; Chen et al., 2021c), conditioning on to establish the convergence results. Due to the mixing property (Assumption 6.1), the Markov chain geometrically converges to its stationary distribution. Therefore, choosing "large enough", and conditioning on to each ensure that the Markov chain at time this "almost in steady state." However, in federated setting, conditioning on to results in bounds that are too loose. In particular, consider the differences $k_{t-t} k_{c}$ and $k_{t-t} k^{2}$ in (15) and (16) respectively. In the centralized setting, as in (Bhandari et al., 2018; Srikant & Ying, 2019; Chen et al., 2021c), these terms can be bounded deterministically to yield bound. However, in the federated setting, this crude bound does not result in linear speedup in N. In this work, to achieve a finer bound on $k_{t-t} k_{c}$, we go steps further back in time. This ensures that the difference behaves almost like the difference of average of i.i.d. random variables, resulting in a tighter bound (see Lemma B.6). By exploiting the conditional expectation $E_{t-t}[]$, we derive a refined analysis to bound this term as O(2=N+4), which guarantees a linear speedup (see Lemmas B.6 and B.7).

Taking total expectation in Proposition B.3 (using tower property), we get

$$E[M(t+1)] (1 '2) E[M(t)] + C_{14}()^{4} + C_{16}()^{2} + C_{17} \frac{X}{N} E[$$

$$k]$$

$$k = t$$
(22)

 \Box

To understand the bound in Proposition B.3, consider the case of K=1 (i.e. full synchronization). In this case we have i=0 for all i, and the bound in Proposition B.3 simplifies to $E\left[M\left(t+1\right)\right]$ (1 i=1) i=1) i=1 i=1

Proposition B.4 (Synchronization Error). Suppose T > K + . For such that

$$\frac{2}{\min} \frac{1}{\frac{1}{2}}; \frac{\ln(5=4)}{-2(1+A_1)(K-1)^2};$$

$$\frac{2(\log()+1)}{2K^2}; \frac{V}{V} = \frac{1}{\ln(5=4)} \cdot \frac{1}{\frac{1}{2}} \cdot \frac{1}{(1+\tilde{K_1})} \cdot \frac{1}{(1+$$

where $v = \frac{r^2}{80 \%_2 C_{17} u_{cm}^2}$, the weighted consensus error satisfies

$$\frac{2C_{17}}{2W_{T}} \underset{t=2}{X^{T}} \underset{'=t}{W_{t}} \underset{t=2}{X^{t}} E \qquad \sim \underline{10C} \qquad \underline{4m_{2}} \qquad \underline{1} \underset{E}{X^{T}} \underbrace{X^{T}} \underset{t=2}{X^{T}} \underset{t=1}{W_{t}} \underbrace{M_{t}} \underbrace{$$

Proof. The proof is presented in Section B.5.

Finally, by incorporating the results in Propositions B.2, B.3, and B.4, we can establish the convergence of FedSAM in Theorem B.1.

Proof of Theorem B.1. Assume $w_0 = 1$, and consider the weights w_t generated by the recursion $w_t = w_{t-1} - 1 - \frac{1}{2} - \frac{1}{2}$. Multiplying both sides of (22) by $\frac{1}{2}$, and rearranging the terms, we get

$$\frac{1}{W_{t=2}} \overset{XT}{W} \overset{t}{EM} \overset{t}{(1)} \overset{T}{T}$$

$$\frac{2}{2W_{T}} [w_{2} \ _{1}EM(_{2}) \quad w_{T} EM(_{T+1})] + \frac{2}{2} C_{14} \overset{2}{(1)^{3}} + C_{16}(_{1})_{N} \qquad - \frac{1}{W_{T}} \overset{X^{T}}{t=2} w_{t}$$

$$+ \frac{1}{W'} \overset{XT}{T} \frac{2w_{t}C_{17}}{2} \overset{X^{T}}{T'} \overset{X^{T}}{Z} \overset{X^{T}}{T'} \overset{X^{T}}{Z} \overset{X^{T}}{T'} \overset{X^{T}}{Z} \overset{X^{T}$$

Substituting the bound on $\frac{1}{W_T} P t^T = w_t^h P^t = t$

from Proposition B.4 into (26), we get

$$\frac{1}{WT} \sum_{t=2}^{X^{T}} W_{t} E M(t),_{2} \frac{1}{2} \frac{2}{2}^{2} \frac{2^{t-2+1}}{E} M(2) + ,_{2} C_{14}()^{3} + C_{16}()_{N} - + {}^{2}B^{2} \frac{1}{2} \frac{0_{1}G}{2} + 1 + {}^{4}M_{2} + (+1)(K - 1) + {}^{2}M_{T} + {}^{2}W_{t} E M(t) + M_{t} + M$$

where Mas a problem dependent constant and is defined in Lemma B.6. To simplify (27), we define $C_{18}() = \tilde{B}^2 \frac{40C_{17}}{\frac{7}{2}} + \frac{4m_2}{B(1-)}$, $C_{19}() = \frac{4}{2}C_{16}^4()$. We have

$$\frac{1}{W^{T}} \stackrel{\text{tX}^{\mp}}{=} W \stackrel{\text{t}}{=} M \stackrel{\text{t}}{=} \frac{1}{2} + \frac{21}{2} \stackrel{\text{f}}{=} \frac{2}{2} = \frac{4C}{2} \stackrel{\text{f}}{=} \frac{4C}{2} \stackrel$$

Furthermore, define $W_T = P_{t=0} W_t W_T$. By definition of T_A we have $E[M(T)] = P_{t=0} W_{t=0} W_t EM(t)$, and hence

$$E[M(_{\wedge})_{j}] = \frac{1}{\widetilde{W}_{t-0}} \stackrel{2}{\times} \stackrel{1}{W} w_{t}EM(_{t}) + \frac{1}{\widetilde{W}_{t-2}} \stackrel{X^{T}}{W} w_{t}EM(_{t})$$

$$\begin{split} &\frac{M_0}{\tilde{W}_{t=0}}^{\frac{3}{2}\sqrt{1}} \ w_t + \frac{1}{\tau \tilde{W}_{t=2}}^{\frac{3}{2}} \ w_t \text{EM} \left(t \right) \text{T} & \text{(Lemma B.6)} \\ &\frac{2M_0 w_2}{\tilde{W}_{t=2}}^{\frac{1}{2}} \ \frac{1}{\tilde{W}} \ w_t \text{EM} \left(t \right) \text{T} & \text{T} & \text{(wt } w_{t+1} \text{ for all } t \text{ 0} \right) \\ &\frac{2M_0 w_2}{\tilde{W}_{t=2}}^{\frac{1}{2}} \ \frac{1}{\tilde{W}} \ \frac{X}{\tilde{W}} \ \text{w}_t \text{EM} \left(t \right) \text{T} & \text{T} \\ &= 2M_0 \ 1 \ \frac{2}{2}^{\frac{1}{2}+1} + 1 \ \frac{T}{\tilde{W}_{t=2}}^{\frac{1}{2}} \ w_t \text{EM} \left(t \right) \text{T} \\ &2M_0 \ 1 \ \frac{2}{2}^{\frac{1}{2}+1} + 1 \ X \ \frac{W_t \text{EM} \left(t \right) \text{T}}{W_{2+1}} + 4 \frac{2}{C} \frac{(\frac{1}{2} + \frac{1}{2} +$$

where $C_{20}(;) = 4M_0 + \frac{4}{2} \frac{4}{T}$ urthermore, by Proposition B.1, we have $M(A) = \frac{1}{2} k_1 k_2 + \frac{1}{2} k_1^2 k_2^2$, and hence

where $C_1(;) = 2u_{cm} \hat{C}_{20}(;), C_2() = 2u_{cm} C_{19}(^3), C_3() = 2u_{cm} C_{18}()$ and $C_4() = 2u_{cm} : ^{4C_{14}()}.$

Finally, note that by definition of , we have = $d2 \log e + 2 \log$ Hence, we have $C_1(;)$ $2u_{cm}M_0(4 + \log_e + 2 -)$ $16u_{cm}M_0(\log_e^1 + 4 -)$ = $C_1:2$ where $C_1 = 16u_{cm}^1 M_0(\log_e^1 + 2 -)$. Furthermore, we have $C_{\frac{1}{2}}() = 2u_{cm}C_{19}() = 2u_{cm}^{2}C_{16}() = \frac{8^{2}u_{cm}(C_{8}+\frac{1}{2}+C_{12})}{2} = \frac{8^{2}u_{cm}(C_{8}+\frac{1}{2}+C_{12})}{2} = \frac{C_{2}}{2}$ where we $\frac{2}{2}(\frac{2}{2})$ 2 denote to emphasize the dependence of on , and $C_{2} = \frac{8u_{cm}(C_{8}+\frac{1}{2}+C_{12})}{2} = \frac{C_{2}}{2}$ where we

Note that we have = O(log(1=)).

In addition, $C_3() = 2u_{cm}C_{18}() = \frac{80B^2C_{17}u_{cm}(1+_{B(1-)})}{C_4()}$: C_3 , where C_3 And lastly, $= \frac{80B^2C_{17}u_{cm}^2(1+_{B(1-)})}{\frac{1}{2}}$: $C_4() = \frac{8u_{cm}(C_7+C_{11}+0.5C_3C_9+C_3C_{10}+2C_1C_3C_{10}+3A_1C_3+C_{13})}{C_4()}$ C_4^2 , where C_4 = $8u_{cm}$ $C_7 + C_{11} + 0.5C^2C^2 + C_3C_{10} + 2C_1C_3C_{10} + 3A_1C_3 + C_{13} = 2.$

Next we will state the proof of Corollary B.1.1. $\frac{^2}{^2}$

Proof of Corollary B.1.1.3 By this choice of step size, for large enough T, will be small enough and can satisfy the requirements of the step size in (21), (23), (32), (35). Furthermore, the first term in (29) will be

4 m ₂

$$=$$
O $\frac{C_1'_2}{NT} :$

Furthermore, for the second term we have

$$C_2 = C_2 = C_2^2$$
 $C^{N} = \tilde{O} = \tilde{N} = C_1^2$

Finally, for the third and the fourth terms we have

$$C_3(K-1)^2 + C_4^{32} = O(C_3 = 2^2 + C_4 = 2^3) \frac{\log^2(NT)}{N} = O(C_3 = 2^2 + C_4 = 2^3) \cdot \frac{\log^2(NT)}{T^2} = O(C_3 = 2^4 + C_4 = 2^3) \cdot \frac{\log^2(NT)}{NT}$$

Upper bounding (29) with , we get $O_{\sim}^{C_1'_2+C_2='_2+C_3='_2+C_4='_2}$. Hence, we need to have $T=\frac{C_1'_2+C_2='_2+C_3='^2+C_4='^3}{O(\frac{2}{N})}$ number of iterations to get to a ball around the optimum with radius .

B.3. Proof of Proposition B.2

The update of the virtual parameter sequence $f_t g$ can be written as follows

$$t+1 = t + (G(t; Y_t) + b(Y_t)):$$
 (30)

Using $\frac{p-1}{2}$ -smoothness of M () (Proposition B.1), we get

$$M(t+1) \ M(t) + hrM(t); t+1 = ti + 2 = k_{t+1} - tkS$$

$$= M(t) + hrM(t); G(t; Y_t) = t + b(Y_t)i + L^2 = kG(t; Y_t) = t + b(Y_t)k^2 = M_t^2 + rM(t); G(t) = t + hrM(t); b(Y_t)i = T_1: Expected update = T_2: Error due to Markovian noise b (Y_t) = tror due to Markovian noise Y_k = tror due to Markovian N_t = tror due to Markovian noise Y_k = tror due to Markovian noise Y_k = tror due to Markovian noise Y_k = tror due to Markovian N_t = tror due to Markovian noise Y_k = tror due to Markovian N_t = tror due to Markovian noise Y_k = tror due to Markovian N_t = tror due to Markovian N_t$$

Lemma B.1. For all 2 R^d, the operator G() satisfies the following

$$T_1 = rM(); G() 2'_2M():$$

Lemma B.1 guarantees the negative drift in the one step recursion analysis of Proposition B.2. It follows from the Moreau envelope construction (Chen et al., 2021c) and the contraction property of the operators $G^{i}()$; i = 1; :::; N (Assumption 6.2).

Lemma B.2. Consider the iteration t of the Algorithm 4, and consider = d2 log e. We have

$$E_t [T_2] = E_t hrM(t); b(Y_t)i$$

$$\begin{split} & \frac{L}{2} \frac{^{2}}{^{2}I_{cs}} \frac{^{2}}{^{4}} k_{t} \quad _{t} \quad k_{c} + \ _{2}E_{t}^{2} \quad kb(Y_{t})k_{c} \qquad ^{2} \\ & + \frac{m_{2}L}{I_{cs}} E_{t} \quad k_{t} \quad _{t}k_{c} + \frac{1}{2} \quad \frac{Lm_{2}{}^{2}u_{cm}}{I_{c}} ^{2} + \, ^{2}E_{\frac{t}{L}} \ [M(_{t})]; cs \end{split}$$

where 1 is an arbitrary positive constant.

In the i.i.d. noise setting, $E[T_2] = 0$. In Markov noise setting, going back steps (which introduces t) enables us to use Markov chain mixing property (Assumption 6.1).

Lemma B.3. For any t 0, denote $T_3 = hrM(t)$; G(t) G(

where 2 and 3 are arbitrary positive constants.

Lemma B.4. For any t 0, we have

$$T_4 = hrM(t); G(t) G(t)i_4M(t) + 2l_{cs}^2 2 \frac{L^2u_{cm}^2}{4}$$

where 4 is an arbitrary positive constant.

Lemma B.5. For any 0 < t, we have

$$T_5 = kG(t; Y_t) + b(Y_t)k^2 + \frac{6(A_2 + 1)^2 u_{cm} M(t)}{s} + \frac{3A_1}{t} + \frac{2}{kb(Y_t)k_c: cscs}$$

Substituting the bounds in Lemmas B.1, B.2, B.3, B.4, B.5, and taking expectation, we get the final bound in Proposition B.2.

B.4. Proof of Proposition B.3

First we state the following two intermediate lemmas, which are proved in Section B.7.

Lemma B.6. Suppose = d2 log e and

min
$$\left(\frac{1}{12^{\frac{1}{A_{1}^{2}}}+1}; \frac{1}{8(A_{2}+1)}\right)$$
 (32)

For any 0 t 2 we have the following

$$M(t) = \frac{1}{|c_{cm}|^2} = \frac{1}{C_1^2} = B + (A_2 + 1) = k_0 k_c + \frac{B}{2C_1}^2 = k_0 k_c^2 = M_0$$
(33)

Furthermore, for any t 2, we have the following

$$E_{t 2}[k_{t} t k_{c}] 4C_{1}E^{t}_{2}[k_{t}k^{c}] + 8 \Big|_{CD} B = + \frac{u_{CD}}{|_{CD}} P_{N}^{2} 1 + \frac{u_{CD}}{2C_{1}} \frac{8m_{2}}{2C_{1}}$$

$$X^{t}_{+ 6A_{1}} E_{t 2}[_{i}]:_{i=t}$$
(34)

Lemma B.7. Suppose = d2 log e and

min
$$\frac{1}{C_1}$$
; $\frac{1}{4C_2}$; $\frac{1}{40C_1}$; (35)

where $C_1 = 3^p \overline{A_2 + 1}$ and $C_2 = 3C_1 + 8$. We have the following

$$E_{t 2}[k_{t t} k^{2}] 8_{c}^{22}C^{2}E_{t 2}k_{t}k_{1}^{2} + 8_{-cD} B^{2}_{c c}^{22}k_{D}^{u^{2}} + 14_{-cD}^{u^{2}} \frac{m_{2}^{24}}{1^{2}} + 8^{2}A^{2} F_{t 2}^{2} i^{2} k_{D}^{2} + 14_{-cD}^{u^{2}} \frac{m_{2}^{24}}{1^{2}} + 8^{2}A^{2} F_{t 2}^{2} i^{2} k_{D}^{2} k_{D}^{2} + 14_{-cD}^{u^{2}} k_{D}^{2} k_{$$

In Lemma B.6, we define C_9 , $\frac{8\underline{u}_{c_D}}{I_{c_D}}$; C_{10} , $\frac{8\underline{m}_2\underline{u}_{c_D}}{I_{c_D}(1)}$. Hence, we can bound the term in (15) as

$${}^{2}C_{3}E_{t-2}[k_{t-t-k}c_{t}] \\ {}^{2}C_{3} \qquad {}^{4}C_{1}E_{t-2}[k_{t}k_{c}] + C_{9} \xrightarrow{-+} {}^{4}\underline{G}_{1}\underline{G}_{2}^{2}(1+2C_{1}) + 6A_{1} \qquad E_{t-2}[k_{t}] \\ {}^{2}C_{3} \qquad {}^{4}C_{1}u_{cm}E_{t-2} \qquad h_{p} \xrightarrow{N} \xrightarrow{N} {}^{N}\underline{P} \xrightarrow{N} {$$

Furthermore, using Lemma B.7, (58) and Proposition B.1, the term in (16) can be bounded as follows:

$$C_{4}E_{t} \ _{2} \ k_{t} \ _{t} \ k^{2} \ 16_{c}^{22}C^{2}C_{4}u_{m}^{2}E_{t} \ _{2}^{2} [M_{(t)}] + 8^{2}A^{2}C_{4} \ _{E}_{t} \ _{2}[_{1} \ _{t} \ _{k}]_{k=0}$$

$$+ \frac{8C_{4}u_{D}^{2}B^{2}}{\frac{1^{2}D_{C_{12}}}{C_{12}}} \frac{^{2}}{N}^{2} + \frac{14C_{4}u_{D}m_{D}^{2}}{\frac{1^{2}C_{D}(\frac{1}{N})^{2}}{C_{13}}} \ _{k=0}^{4}$$

$$(37)$$

Inserting the upper bounds in (36) and (37) in the upper bound in Proposition B.2, we have

$$\begin{split} E_{t-2}\left[M\left(_{t+1}\right)\right] & 1 - 2^{\prime}2^{-2}2^{2$$

We define
$$C_{14}^{0}()$$
 $C_{1}\mathbf{4}()$, $C_{1}\mathbf{4}()$, $C_{11} + Q : \frac{5C^{2}C^{23} + C_{3}C_{10} + Q^{2}C_{10}C_{10} + 3A^{2}C_{10} + C_{13} +$

This yields

 $C_{15}()$ 1 '2. This completes the proof.

B.5. Proof of Proposition B.4

First we state the following lemma, which characterizes a bound on the expectation of the synchronization error

Lemma B.8. Suppose Assumptions 6.1, 6.3 holds and the step size satisfies 1, where s = bt=K c, the network consensus error t, 1

where $\tilde{A}_1 = \frac{2A_1^2u_{c2}^2}{l_{c2}^2}$; $\tilde{A}_2 = \frac{2A_2^2u_{c2}^2}{l_{c2}^2}$; $\tilde{B} = \frac{u_{c2}}{l_{c2}}B$. Here, A_1 ; A_2 ; B are the constants defined in Assumption 6.3, and l_{c2} ; u_{c2} are constants involved in the equivalence of the norms: $l_{c2} kk_2 kkC u_{c2} kk_2$.

Due to the assumption on the step size, the bound in Lemma B.8 hold. Substituting the bound on from Lemma B.8, we get

$$\frac{2C_{17}}{W} \quad \begin{array}{c} X^T & X^t \\ W_t & \end{array} \quad E$$

$$\frac{2C_{17}}{\sqrt{2W_{T}}} W_{t} V_{t=2}^{T} V_{t=1}^{T} X^{T} X^{t} S^{2}('-s,K)B^{2} I^{2} + \frac{4m_{2}}{B(1-)} + 5^{2}('-s,K)A_{2}^{2} E_{t_{0}} kC^{2} : (42)$$

where $s_r = b' = Kc$. Hence, $s_r K$ denotes the last time instant before 'when synchronization happened. The first term in (42) can be upper bounded as follows.

$$\frac{2C_{17}}{{}^{'}_{2}W_{T}} \underbrace{\overset{X^{T}}{w_{t}} \overset{X^{T}}{s^{2}} w_{t}}_{t=2} \overset{X^{T}}{y_{t}} w_{t} \overset{X^{T}}{5^{2}} (' s_{t}' K) B^{2} \overset{\sim}{1} + \overset{B}{B} \underbrace{\overset{4m_{2}}{-1}}_{(1} \overset{\sim}{)}$$

$$= {}^{2} \overset{\sim}{}^{2} 1 + \underbrace{\overset{1}{B} (1 + \frac{4m_{2}}{y_{2}})^{1}}_{t=2} \overset{QC}{W_{T}} \underbrace{\overset{1}{w_{t}} \overset{X^{T}}{y_{t}}}_{t=1} (' X_{S}^{t} K)^{2} B^{2} 1$$

$$+ \underbrace{\overset{1}{B} (1 \overset{1}{17})^{1}}_{t=2} \overset{W_{t}}{\underbrace{\overset{1}{0} C_{t}}} \underbrace{\overset{1}{1}}_{t=2} \overset{W_{t}}{\underbrace{\overset{1}{0} C_{t}}} \overset{X^{T}}{y_{t}} w_{t} \underset{t}{\cancel{\#}} + 1) (K 1)$$

$$= {}^{2} B_{s}^{2} 1 + \underbrace{\overset{1}{B} (1 \overset{1}{y_{2}})^{2}}_{B(1 \overset{1}{y_{2}})} \underbrace{\overset{0C}{\overset{1}{0} C_{17}}}_{t=2} (\frac{1}{t} \overset{1}{1}) (K 1) :$$

$$(' s_{t} K K 1)$$

$$= {}^{2} B_{s}^{2} \overset{\sim}{1} + \underbrace{\overset{1}{B} (1 \overset{1}{y_{2}})^{2}}_{B(1 \overset{1}{y_{2}})} \underbrace{\overset{0C}{\overset{1}{0} C_{17}}}_{t=2} (\frac{1}{t} \overset{1}{1}) (K 1) :$$

(43)

Next, we compute the second term in (42).

$$\frac{2C_{17}}{'_{T}W_{=2}} X^{T} W_{t} X^{t} 5^{2} ('_{S'K}) A_{2}^{\sim} K^{t} k_{0} k_{0}^{2} k_{$$

 ${}^{2}A_{2} \sim \frac{20C_{17}u_{cm}^{2}}{{}^{2}W^{T}} \begin{cases} \begin{cases} X^{K} & X^{t} &$

where if K < 2, $I_1 = 0$. Next, we bound I_1 ; I_2 separately.

where, (45) follows since w_{t-1} w_t ; 8 t. Next, to bound I_2 in (44), we again split it into two terms.

Next, we bound I_4 , assuming $t_0 K + T < (t_0 + 1)K +$, where t_0 is a non-negative integer.

Using the bounds on I₃; I₄ from (46) and (47) respectively, we can bound I₂.

Substituting the bounds on I₁; I₂ from (45), (48) respectively, into (44), we get

$$\frac{2C_{17}}{'_{T}W_{=2}} \underbrace{X^{T}}_{'=t} \underbrace{W_{t}}_{'=t} 5^{2}('_{s'K})A_{2}^{\sim} \underbrace{K^{1}}_{t=s'K} K_{'0}k_{2}^{2} \underbrace{K^{2}}_{t=s'K} K$$

We analyze the terms in (49) separately. First, for the terms with ${}^{P}t \stackrel{\mathsf{L}}{=} {}_{0}{}^{1}$ M (${}_{t}$),

$$\frac{1}{WT}^{2}A_{2}^{2}\frac{20C_{17}u_{cm}^{2}}{^{2}}K(+1)[(K-2+1)w_{K}+w_{K+}] \stackrel{K-1}{M}(0)$$

$$= \frac{1}{W}^{2}A_{2}^{2}\frac{7^{2}}{20C_{1}^{2}u_{cm}}K(+1) = \frac{1}{W}^{2}A_{2}^{2}\frac{90C_{17}u_{cm}^{2}}{^{2}}K(+1) = \frac{1}{W}^{2}A_{2}^{2}\frac{90C_{17}u_{cm}^{2}}{^{2}}K(+1) = \frac{1}{W}^{2}A_{2}^{2}\frac{90C_{17}u_{cm}^{2}}{^{2}}K(+1) = \frac{1}{W}^{2}A_{2}^{2}\frac{90C_{17}u_{cm}^{2}}{^{2}}K(+1) = \frac{1}{W}^{2}A_{2}^{2}\frac{1}{W}^{2}\frac{1}{W}^{2}M(t)$$

$$\frac{1}{2WT}^{2}\frac{1}{W}^{2}M(t);$$

$$\frac{1}{2WT}^{2}\frac{1}{W}^{2}M(t);$$

$$\frac{1}{2WT}^{2}\frac{1}{W}^{2}M(t);$$
(50)

where (50) holds since we choose small enough such that

$${}^{2}A_{2}^{2}\frac{{}^{2}OC_{17}u_{cm}^{2}}{{}^{2}}(K 2 + 1) + 1 \frac{{}^{2}\frac{K(-+1)}{2}}{2}$$

To get this, we use the inequality 1 $\times \exp_{\frac{1}{x}} \frac{x}{x}$ for x < 1, $\frac{x}{2-2}$ and $\frac{1}{2} < 2 \log + 1$, we get $\frac{1}{\left(1 - \frac{1}{2}\right)}$ exp $\frac{1}{2} = 2 \log + 1$. For (50) to hold, it is sufficient that

Next, for the remaining terms in the third line of (49), by the assumption on the step size, we have

$${}^{2}A_{2}^{2}\frac{20C_{17}u_{cm}^{2}}{{}^{'}_{2}}\frac{, \quad {}^{3}1}{1 \quad {}^{2}\frac{}{-2}{}^{'}} \qquad \frac{4}{4}$$

$${}^{2}\tilde{A}_{2}\frac{20C_{17}^{2}u_{cm}(+1)(K-1)}{{}^{'}_{2}}\frac{1 \quad {}^{-2}}{1 \quad {}^{-2}}^{K+1}\frac{1}{4}; \qquad (52)$$

and hence we get

$$\frac{1}{W_{T}}^{2} A_{2}^{2} \frac{20C_{17}u_{cm}^{2}}{^{'}_{2}} = \frac{^{3}w_{K}^{K+}}{1^{2}} \frac{^{1}X_{C}^{t}}{^{2}} + (+1)(\frac{K}{1}) \frac{^{T}}{1^{\frac{'}{t^{2}}}K^{K+}} M_{C}^{t})^{T}$$

$$\frac{1}{2W_{t}} \frac{^{T}}{^{W_{t}}} W_{t} M_{C}(t) : T$$
(53)

Substituting (50), (53) in (49), we get

$$\frac{2C_{17}}{'_{T}W_{=2}} X^{T} \underset{'=t}{W_{t}} X^{t} \qquad 5^{2}('_{s'}K)A_{2}^{\sim} X^{1} k_{t}k^{2}_{c} \frac{1}{2W_{0}} X^{T} w_{t}M(_{t}):^{2}$$
(54)

Finally, substituting the bounds in (43), (54) into (42), we get

$$\frac{2C_{17}}{W} \times W_{t} \times W_{t} \times E$$

$$\frac{2C_{17}}{W} \times W_{t} \times W_{t} \times W_{t} \times E$$

$$\frac{2C_{17}}{W} \times W_{t} \times W_{t} \times W_{t} \times W_{t} \times E$$

$$\frac{2C_{17}}{W} \times W_{t} \times$$

B.6. Auxiliary Lemmas

The following lemma is of central importance in proving linear speedup of FedSAM. $p = \frac{x > D x}{x > D x}$ for some Lemma B.9. Let I_{cD} and u_{cD} be constants that satisfy $I_{cD}k$ k_D k k_c $u_{cD}k$ k_D , where k k_D = positive definite matrix D 0. Note that for any D 0, these constants always exist due to norm equivalence. Furthermore, in $\sqrt{x > D x}$ for some D 0, we take $I_{cD} = u_{cD} = 1$. We have case the norm kxk_c is defined in the form

$$E_{t-r}[kb(Y_t)k_c] \stackrel{u_{cD}}{\underset{cD}{|}} p_{B_N^{-}} + 2m_2^r$$

$$E_{t-r}[kb(Y_t)k_c^2] \stackrel{u^2}{\underset{cD}{|}} \frac{B_N^2}{N} + 2m_2^{2r} :$$
(56)

$$E_{t-r}[kb(Y_t)k_c^2] = \frac{u^2}{l^2} \frac{B^2}{N} + 2m_2^{2^{2r}} :$$
 (57)

Lemma B.9 is essential in characterizing the linear speedup in Theorem B.1. This lemma characterizes the bound on the conditional expectation of $kb(Y_t)k_c$ and $kb(Y_t)k_c^2$, conditioned on r time steps before. In order to understand this lemma, consider the bound in (57). For the sake of understanding, suppose the noise Y t is i.i.d. In this case we will end up with the first term which is proportional to 1=N. This is precisely the linear reduction of the variance of sum of N i.i.d. random variables. Furthermore, in order to extend the i.i.d. noise setting to the more general Markovian noise, we need to pay an extra price by adding the exponentially decreasing term to the first variance term.

Lemma B.10. The following hold

$$kG(t; Y_t) = G(t; Y_t)k_c = A_{1t}^2 = A_1^2 = A_2^2 = A_1^2 = A_1^2$$

Lemma B.11. For the generalized Moreau Envelope defined in (12), it holds that

$$r M_f^{;g}(x) \stackrel{?}{=} kxk_m;$$
 $D_{;g}^{;g}(x); x 2M_f^{;g}(x):$

B.7. Proof of Lemmas

Proof of Lemma B.1. Using Cauchy-Schwarz inequality, we have

$$| \frac{?}{m} z \underline{m} | \{z\}$$

$$rM(); G() krM()k \{G() hrM(); i; \}_{T_{11}} T_{12}$$
(59)

where k k_m ? denotes the dual of the norm k k_m . Furthermore, by Proposition B.1 we have

(Lemma B.11)

Furthermore, by the convexity of the k k_m norm (Lemma B.11), we have

(60)

$$T_{12} kk_m = ^2 2M()$$
:

Hence, using (60), (61) in (59) we get

(61)

$$\frac{rM();G()}{|M()|} = 2'_{2}M();cm$$
 (62)

where '2 is defined in (14).

Proof of Lemma B.2. Given some $< t, T_2 = hrM(t); b(Y_t)i$ can be written as follows:

$$T_{2} = hrM(t) \quad rM(t); b(Y_{t})i + hrM(t); b(Y_{t})i krM(t) rM(t)k_{s}:kb(Y_{t})k_{s} + hrM(t); b(Y_{t})i$$
 (Cauchy–Schwarz)
$$\frac{L}{l} k_{t} \quad t \quad k_{c}: \frac{1}{l} kb(Y_{t})k_{c} + hrM(t); b(Y_{t})i$$
 (Proposition B.1)
$$\frac{1}{2} \frac{L}{l} \frac{L^{2}}{l^{2}_{c}} k_{t} \quad t \quad k^{2} + {}_{c_{2}}kb(Y_{t})k^{2} + hrM(t); b(Y_{t})i$$
 (63)

Taking expectation on both sides, we have

For T₂₁, we have

$$T_{21} I_{cs} \ \&rM(_t) k_s k E_t^? [b(Y_t)] k_c^{-1} cs$$

$$[krM(_t) k_s m_2]^? \qquad (Assumption 6.1)$$

$$\frac{m_2^2}{l} [krM(_t) rM(_t) k^? + krM(_t) k^?] cs$$

$$s \qquad (assumption on)$$

(65)

Substituting (65) in (64), we get

$$\begin{split} E_t & [T_2] \, \frac{L}{2} \frac{\frac{2}{2} I_{cs}}{\frac{1}{2} I_{cs}} k_t + \frac{1}{2} E_t + \frac{1}{2$$

Proof of Lemma B.3.

$$T_{3} = hrM(t); G(t; Y_{t}) \quad G(t)i$$

$$= hrM(t) \quad rM(t); G(t; Y_{t}) \quad G(t)i$$

$$+ hrM(t); G(t; Y_{t}) \quad G(t); Y_{t}) + G(t) \quad G(t)i$$

$$+ hrM(t); G(t; Y_{t}) \quad G(t); T_{32}$$

$$+ hrM(t); G(t; Y_{t}) \quad G(t); T_{33}$$

$$+ hrM(t); G(t; Y_{t}) \quad G(t); T_{33}$$
(66)

Next, we bound all three terms individually.

I. Bound on T₃₁:

$$T_{31} = \frac{1}{N} X^{N}$$

$$rM(t) rM(t); G^{i}(i; y^{i}) G^{i}(i) = 1$$

$$\frac{1}{N_{=1}} X^{N} krM(t) rM(t) k_{3} G^{i}(i; y^{i}) G^{i}(i) : t_{3} : t_{3} = 1$$

$$I = \frac{1}{N_{=1}} X^{N} krM(t) rM(t) k_{3} G^{i}(i; y^{i}) G^{i}(i) : t_{3} : t_{3} = 1$$

$$I = \frac{1}{N_{=1}} X^{N} krM(t) rM(t) k_{3} G^{i}(i; y^{i}) G^{i}(i) : t_{3} : t_{3} = 1$$

$$I = \frac{1}{N_{=1}} X^{N} krM(t) rM(t) k_{3} G^{i}(i; y^{i}) G^{i}(i) : t_{3} : t_{3} = 1$$

$$I = \frac{1}{N_{=1}} X^{N} krM(t) rM(t) rM(t) k_{3} G^{i}(i; y^{i}) G^{i}(i) : t_{3} : t_{3} = 1$$

$$I = \frac{1}{N_{=1}} X^{N} krM(t) rM(t) rM(t) k_{3} G^{i}(i; y^{i}) G^{i}(i) : t_{3} : t_{3} = 1$$

$$I = \frac{1}{N_{=1}} X^{N} krM(t) rM(t) rM(t) k_{3} G^{i}(i; y^{i}) G^{i}(i) : t_{3} : t_{3} = 1$$

$$I = \frac{1}{N_{=1}} X^{N} krM(t) rM(t) rM(t) k_{3} G^{i}(i; y^{i}) G^{i}(i) : t_{3} : t_{3} = 1$$

$$I = \frac{1}{N_{=1}} X^{N} krM(t) rM(t) rM(t) k_{3} G^{i}(i; y^{i}) G^{i}(i; y^{i}) G^{i}(i) : t_{3} = 1$$

$$I = \frac{1}{N_{=1}} X^{N} krM(t) rM(t) rM(t)$$

For T₃₁₁, we have

$$T_{311} = krM(t) rM(t)ks^{?} L_{k_t t k_c};$$
 (68)

where the inequality follows from Proposition B.1 and (13). For T₃₁₂, we have,

Substituting (68) and (69) in (67), we get

$$T_{31} \stackrel{2LA_{2}}{=} \frac{1}{N} \frac{X}{N} \stackrel{N}{k_{t}} \stackrel{k}{t_{t}} \stackrel{k}{k_{t}} \stackrel{k}{t_{t}} \stackrel{k}{k_{t}} \stackrel{k}{t_{t}} \stackrel{k}{k_{t}} \stackrel{k}{k_{t}} \stackrel{k}{t_{t}} \stackrel{k}{k_{t}} \stackrel{k$$

II. Bound on T₃₂:

$$T_{32} = \frac{1}{N} X^{N}$$

$$r M (t); G^{i}(i; y^{i}) G^{i}(i; y^{i}) + G^{i}(i) G^{i}(i) + G^{i}(i) G^{i}(i) = 1$$

$$\frac{1}{N} X^{N} K_{i=1} I X^{N} (t) K_{s} K_{s}^{G}(i; y^{i}) G_{t}^{i}(t; y^{i}) G_{t}^{i}(t; y^{i}) G_{s}^{i}(t; y^{i}) G_{s}^{i}(t$$

where the last inequality is by Hölder's inequlaity. For T_{321} , we have

$$T_{321} = krM(_{t}) rM(0)k_{s}^{\frac{1}{2}}k_{t} - 0k_{s}$$
 (Proposition B.1)
$$\frac{L}{L}[k_{t}k + _{c}k_{t} t k]: cs c$$
 (72)

For T_{322} , we have

$$T_{322} = k G^{i}(^{i}; _{t}y_{t})_{i} \quad G^{i}(_{t}; _{t}y_{t})_{+} G^{i}(^{i})_{t} \quad G^{i}(^{i})_{ks} \quad t$$

$$I_{cs} \underbrace{G^{i}(^{i}; y^{i})_{t} \quad G^{i}(_{t}; y_{t})_{d} + G^{i}(_{t})_{ci}}_{t} \quad G^{i}(^{i})_{ci} \quad t \quad \text{(triangle inequality)}$$

$$I_{cs} \underbrace{A_{1} + 1_{i}}_{lcs} \quad t_{t} + k_{t}_{c} t \quad k + t \quad i_{c} \quad t \quad c \quad (c \ 1; \text{ triangle inequality})$$

$$= \underbrace{A_{1} + 1_{i}}_{lcs} \quad t_{t} + k_{t}_{c} t \quad k_{c} + t \quad : \quad i \quad (73)$$

Substituting (72) and (73) in (71), we get

$$\frac{L(A \mid + 1)_{2}}{\frac{1}{2}}M(t) + \frac{3u_{cm}}{3} + \frac{2k_{t}^{2}}{2\frac{2}{3}}t \quad k^{2} + \frac{3u_{cm}}{t} + \frac{3}{c}N \quad \frac{(i^{3})^{2}}{2\frac{2}{3}} + \frac{(i}{2})^{2}\frac{1}{2}L(A \mid + 1)_{2}M(t) + \frac{1}{c}A_{3} + 2k_{t} \\
t \quad \frac{k_{c} + \frac{1}{3}}{\frac{2}{cs}} \quad \frac{3u_{cm}^{2}}{2} \quad \frac{2}{2} \quad \frac{2}{2} \quad \frac{2}{2} \quad \frac{2}{2}$$
(74)

III. Bound on T₃₃: Taking expectation on both sides of T₃₃, we have

$$\begin{split} E_t & [T_{33}] = \\ rM & (t \); E_t \left[G(t \ ; Y_t) \right] \quad G(t \) \end{split} \\ & = \frac{1}{N} \frac{X^N}{rM \ (t \); E_t \left[G^i(t \ ; Y_t) \right] \quad G^i(t \)} = 1 \\ & \frac{1}{l^2} \frac{1}{N} \frac{X^N}{N} \quad krM \ (t \) k^2 \quad E_{t_s} \left[G^i(t \ ; Y^i) \right]_t \quad G^i(t \)^{c_i = 1} \quad t \\ & \frac{1}{l^2} \frac{1}{N} \frac{X^N}{N} \quad krM \ (t \) k^2 \quad m_s k^i \quad t \quad k_c \, ^{cs} \\ & \frac{1}{l^2} \frac{1}{N} \frac{X^N}{N} \quad krM \ (t \) k^2 \quad k^i \quad s \quad t \quad k_c \\ & \frac{m_1}{l^2} \frac{1}{N} \frac{X^N}{N} \quad k^i \quad k^i \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{N} \frac{X^N}{N} \quad k^i \quad k \quad k_t \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{N} \frac{X^N}{N} \quad k^i \quad k \quad k_t \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{N} \frac{X^N}{N} \quad k^i \quad k \quad k_t \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{N} \frac{X^N}{N} \quad k^i \quad k^i \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{N} \frac{X^N}{N} \quad k^i \quad k^i \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{l^2} \frac{X^N}{N} \quad k^i \quad k^i \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{l^2} \frac{X^N}{N} \quad k^i \quad k^i \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{l^2} \frac{1}{N} \quad k^i \quad k^i \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{l^2} \frac{1}{l^2} \quad k^i \quad k^i \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{l^2} \frac{1}{l^2} \quad k^i \quad k^i \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{l^2} \frac{1}{l^2} \quad k^i \quad k^i \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{l^2} \frac{1}{l^2} \quad k^i \quad k^i \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{l^2} \frac{1}{l^2} \quad k^i \quad k^i \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{l^2} \frac{1}{l^2} \quad k^i \quad k^i \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{l^2} \frac{1}{l^2} \quad k^i \quad k^i \quad k_c \, ^{cs} \\ & \frac{m_1}{l^2} \frac{1}{l^2} \frac{1}{l^2} \quad k^i \quad$$

Substituting the bounds on T_{31} ; T_{32} ; T_{33} from (70), (74), (75), in (66), we get the result.

Proof of Lemma B.4. Denote $T_4 = hrM(t)$; G(t) = G(t)i. By the Cauchy-Schwarz inequality, we have

$$T_{4} k_{1} M_{(t)} k_{1} k_{3} k_{3} (t) G_{(t)} k_{1} : T_{41} S_{142}$$

$$T_{42} S_{142} (76)$$

For T₄₁ we have

For T₄₂ we have

$$kG(t) \quad G(t)k_{s} = \frac{1}{N} \sum_{s}^{X^{N}} G^{i}(t)_{t} \quad G^{i}(t) \qquad \frac{1}{N} \sum_{s}^{X^{N}} G^{i}(t)_{t} \quad G^{i}(t)_{s} = 1$$

$$\frac{1}{N} \sum_{s}^{X^{N}} G^{i}(t)_{t} \quad G^{i}(t) \qquad \frac{1}{N} \sum_{s}^{N} G^{i}(t)_{t} \quad G^{i}(t)_{s} = 1$$

$$\frac{1}{N} \sum_{s}^{X^{N}} G^{i}(t)_{t} \quad G^{i}(t) \qquad \frac{1}{N} \sum_{s}^{N} G^{i}(t)_{t} \quad G^{i}(t)_{s} = 1$$
(Jensen's inequality)

where (78) follows from Assumption 6.2. Combining (77) and (78), for an arbitrary 4 > 0, we get

Proof of Lemma B.5. Denote $T_5 = kG(t; Y_t) + b(Y_t)k_s$. We have

$$T_{5} = kG(t; Y_{t}) + G(t; Y_{t}) G(t; Y_{t}) t + b(Y_{t})kS^{2}$$

$$3 kG(t; Y_{t}) \{z \frac{tk_{1} + 3_{5}^{2}kG(t_{1}; Y_{t}) G(t; Y_{t})k_{1} + 3_{5}kG(Y_{t})kS : T_{51}\}^{2}$$

$$T_{52}$$
(80)

For T₅₁ we have

$$T_{51} I_{cs} ({}^{2}kG(t; Y_{t})kC + k_{t}k_{c})^{2} I_{cs}$$

$$(A_{2} k_{t}kC + k_{t}k_{c})^{2} = {}^{2}(A_{2} + (Assumption 6.3))$$

$$1) {}^{2}u^{2}_{m} Y_{cs} (t) :$$
(81)

For T₅₂ we have

$$T_{52} = \frac{1}{N} \frac{X}{s} \left(G^{\dagger}(i; y^{\dagger}) - G^{\dagger}(t; y^{\dagger}) \right)_{t} = \frac{1}{N} \frac{X}{s} \left(G^{\dagger}(i; y^{\dagger}) - G^{\dagger}(t; y^{\dagger}) \right)_{t} = \frac{1}{N} \frac{X}{s} \left(G^{\dagger}(i; y^{\dagger}) - G^{\dagger}(t; y^{\dagger}) - G^$$

Using the bounds in (81), (82), we get

$$T_{5} = \frac{6(A_{2} + 1)^{2} u_{cm}^{2} M(t) + \frac{3A_{1}^{2}}{t^{2}}}{t^{2} + 3 kb(Y_{t})k^{2}; cs cs}$$
(83)

where the last inequality is by the assumption on .

(82)

Proof of Lemma B.6. Define $_{I} = _{N} \stackrel{1}{\overset{P}{=}} _{1} \stackrel{I}{\underset{I}{=}} _{c}.$ By the update rule of Algorithm 4, if I + 1 mod K = 0, we have

(84)

and the same bound as in (84) hads. By recursive application of (84), we get

$$= (1 + (A_2 + 1))^{l+1}_0 + B \times (1 + (A_2 + 1))'_{-1} = 0$$

$$= (1 + (A_2 + 1))^{l+1}_0 + B \times (1 + (\frac{1}{A_2 + 1}))^{l+1}_{-1} = \frac{1}{A_1 + A_2 + A_3} = \frac{1}{A_1 + A_2} = \frac{1}{A_1 + A_2}$$

Notice that for $x = (1 + x)^{+1}$, we have $(1 + x)^{+1}$ (1 + 2x), where $(1 + x)^{+1}$ $(1 + (x)^{+1})$ $(1 + (x)^{+1})$

$$_{1} 2_{0} + B^{4(A_{2} + 1)} = 2_{0} + 4B: 2$$
 (86)

Furthermore, we have

(Assumption 6.3)

(87)

Suppose 0 t 2. We have

$$\chi^{1}$$
 (A₂ + 1)_k + B _{k=0} (By (87))

t :

Bt +
$$(A_2 + 1)$$
 $X (2_0 + 4_B) = 0$
 $2(B + (A_2 + 1)(2_0 + 4_B))$ (By (86))

$$\frac{1}{C_1}$$
 B + (A₂ + 1) ₀ + $\frac{B}{2C_1}$: (t 2)

(88)

Furthermore, by Proposition B.1, we have $M(t) = {}_{2}k_{t}k_{m}$, and hence for any 0 t 2 we have

$$\begin{split} M\left({}_{t}\right) & = \frac{1}{2l_{\perp}^{2}} k_{c} c_{m}^{2} \\ & = \frac{2l_{\perp}^{2}}{2l^{2}} \left(k_{t} - o + o k_{c}\right)^{2} c_{m} \\ & = \frac{1}{2l^{2}} \left(k_{t} - o k_{c} + k_{0} k_{c}\right)^{2} c_{m} \\ & = \frac{1}{2l_{cm}^{2}} \frac{\left(k_{t} - o k_{c} + k_{0} k_{c}\right)^{2}}{\left(l_{cm}^{2} - o k_{c} + k_{0} k^{2}\right)} c_{c} \\ & = \frac{1}{l_{cm}^{2}} \frac{1}{C_{\perp}^{2}} - \frac{1}{C_{\perp}^{2}} - B + \left(A_{2} + 1\right) - k_{0} k_{c} + \frac{B}{2C_{\perp}}^{2} - c_{\perp} + k_{0} k_{c}^{2} ; \end{split}$$
 (triangle inequality)

which proves the first claim.

Next we prove the second claim. By the update rule in 30, we have

 $=6^{2}(A^{2} \pm 1)k_{1}k^{2} +_{c}3^{2}kb(Y_{1})k^{2} + 3^{2}kG(I;Y_{1}) G(I;Y_{1})k^{2}$ $6^{2}(A^{2} +_{2}1)k_{1}k^{2} + 3^{2}kb(Y_{1})k^{2} + 3^{2}A^{2};$ (89)

where (89) follows from Lemma B.10. Taking square root on both sides, we get

$$q = \frac{1}{k_{l+1} + k_c + 3} A_2 + 1 k_1^2 k_c + 2 k_b (Y_l) k_c + 2 A_{1l}$$
: (90)

Combining the above inequality with the fact that $k_{l+1}k_c = k_lk_c = k_{l+1} = lk_c$, we get

$$\frac{g}{k_{l+1}k_c(1+3)} + \frac{g}{4k_lk_c} + 2kb(Y_l)k_c + 2A_{1l}$$

$$= (1 + C_1)k_1k_c + 2kb(Y_1)k_c + 2A_{11}; (91)$$

where we denote $C_1 = 3^p A_2^2 + 1$. Assuming t I t, and taking expectation on both sides, we have

where $c_t(I) = 2 \frac{u_{cD}}{I_{cD}} \frac{pB}{N} + 2m_2^{I} t^+$. By applying this inequality recursively, we have

We study T_1 and T_2 in (93) separately. For T_1 we have

$$T_{1} = 2 \frac{u_{cD}}{|C|^{2}} (1 + C_{1})^{-k} \xrightarrow{B} \xrightarrow{P} 2m_{2}^{k}$$

$$= 2 \frac{u_{cD}}{|C|^{2}} \xrightarrow{P} (1 + C_{1})^{+1} \xrightarrow{1} + 2m_{2}(1 + C_{1})$$

$$= 2 \frac{u_{cD}}{|C|^{2}} \xrightarrow{B} (1 + C_{1})^{+1} \xrightarrow{1} + 2m_{2}(1 + C_{1})$$

$$= 2 \frac{u_{cD}}{|C|^{2}} \xrightarrow{B} \frac{N}{|C|^{2}} (1 + C_{1})^{+1} \xrightarrow{1} + 2m_{2}(1 + C_{1})$$

$$= 2 \frac{u_{cD}}{|C|^{2}} \xrightarrow{B} \frac{N}{|C|^{2}} (1 + C_{1})^{+1} \xrightarrow{1} + 2m_{2}(1 + C_{1}) \xrightarrow{1} (> 0) 2_{|C|^{2}}$$

$$= 2 \frac{u_{cD}}{|C|^{2}} \xrightarrow{B} \frac{N}{|C|^{2}} (1 + C_{1})^{+1} \xrightarrow{1} + 2m_{2}(1 + C_{1}) \xrightarrow{1} (> 0) 2_{|C|^{2}}$$

$$= 2 \frac{u_{cD}}{|C|^{2}} \xrightarrow{B} \frac{N}{|C|^{2}} (1 + C_{1})^{+1} \xrightarrow{1} + 2m_{2}(1 + C_{1}) \xrightarrow{1} (> 0) 2_{|C|^{2}}$$

$$= 2 \frac{u_{cD}}{|C|^{2}} \xrightarrow{B} \frac{N}{|C|^{2}} (1 + C_{1})^{+1} \xrightarrow{1} + 2m_{2}(1 + C_{1}) \xrightarrow{1} (> 0) 2_{|C|^{2}}$$

Notice that for $x \stackrel{\log 2}{=}$, we have $(1+x)^{+1} = 1+2x(+1)$. By the assumption on , we have $(1+C_1)^{+1} = 1+2C_1(+1) = 1+4C_1 = 2$ and $(1+C_1) = 1+2C_1 = 1+1=2$. Hence, we have

$$T_1 \ 2 \frac{u}{1} B_{p}^{-2} + 1 + 4m_2 = \frac{4m_2}{1}$$
 (95)

Furthermore, for the term T_2 we have

$$T_{2} = {\begin{pmatrix} 1 + C_{1} \end{pmatrix}}^{k} {}_{t + k} {\begin{pmatrix} 1 + C_{1} \end{pmatrix}}_{t + k} {}_{k=0} {k=0}$$
 (due to > 0)

$$X {}_{t + 2C_{1}} {}_{t + k} {\begin{pmatrix} 2 \end{pmatrix}}^{k} {}_{t + k} {}_{k} {}_{k} {}_{k=0} {}$$

Subtituting (95), (96) in (93), for every t | 1, we get

$$E_{t 2}[k_1k_c] 2E_{t 2}[k_t k_c] + \frac{2u_{cD}}{I} \frac{4B}{V} + \frac{4A_1}{I} \frac{t}{V} + \frac{k}{I} = \frac{(97) cD}{V}$$

But we have

Proof of Lemma B.7. From (91) we have

$$k_{l+1}k^{2} \left\{ (1 + C_{1})^{2}k_{l}k^{2} + 4^{2}{}_{c}kb(Y_{l})k^{2} + 4^{2}A_{c}^{22} \right\} + 4 \left\{ (1 + C_{1})k_{l}k_{c}k_{c}kb(Y_{l})k_{c} + 4A_{1}(1 + C_{1})k_{l}k_{c}l + 8^{2}A_{1l}kb(Y_{l})k_{c} : \begin{cases} z \\ t_{3} \end{cases} \right\}$$

$$\left\{ (1 + C_{1})k_{l}k_{c}k_{c}kb(Y_{l})k_{c} + 4A_{1}(1 + C_{1})k_{l}k_{c}l + 8^{2}A_{1l}kb(Y_{l})k_{c} : \begin{cases} z \\ t_{3} \end{cases} \right\}$$

$$\left\{ (1 + C_{1})k_{l}k_{c}k_{c}kb(Y_{l})k_{c} + 4A_{1}(1 + C_{1})k_{l}k_{c}l + 8^{2}A_{1l}kb(Y_{l})k_{c} : \begin{cases} z \\ t_{3} \end{cases} \right\}$$

$$\left\{ (1 + C_{1})k_{l}k_{c}k_{c}kb(Y_{l})k_{c} + 4A_{1}(1 + C_{1})k_{l}k_{c}l + 8^{2}A_{1l}kb(Y_{l})k_{c} : \begin{cases} z \\ t_{3} \end{cases} \right\}$$

$$\left\{ (1 + C_{1})k_{l}k_{c}k_{c}kb(Y_{l})k_{c} + 4A_{1}(1 + C_{1})k_{l}k_{c}l + 8^{2}A_{1l}kb(Y_{l})k_{c} : \begin{cases} z \\ t_{3} \end{cases} \right\}$$

For T₁ we have

$$T_{1} = 2 \frac{p}{(1 + C_{1})k_{1}k_{c}} 2 \frac{p}{(1 + C_{1})k_{0}(Y_{1})k_{c}}$$

$$2(1 + C_{1})k_{1}k_{c} + 2(\hat{I} + C_{1})k_{0}(Y_{1})k_{c} 4k_{1}k_{c} + 2$$

$$4k_{0}(Y_{1})k_{c};$$
(ab $\frac{1}{2}a^{2} + \frac{1}{2}b^{2}$)
$$4(100)$$

where the last inequality is by the assumption on . Analogously for T₂ we have

$$T_{2} = 2^{p} \overline{(1 + C_{1})k_{1}k_{c}} 2A_{1} \overline{(1 + \overline{C_{1})_{1}}}$$

$$4k_{1}k_{c} + 24^{2}A_{1} : 2^{2} 2$$
(101)

For T₃ we have

$$T_{3} = 2kb(Y_{1})k_{c} \quad 4A_{11}$$

$$2^{2}kb(Y_{1})k^{2} + {}_{c}8^{2}A^{22}: \qquad (102)$$

Combining the bounds in (100), (101), and (102), and noting that $(1 + C_1)^2$ 1 + 3C₁, from (99) we have

$$\begin{aligned} k_{l+1}k_{c} &\stackrel{?}{(}1 + 3C_{1})k_{l}k_{c} + 4^{22}kb(Y_{l})k_{c} + 4^{2}A_{1l}^{2} & ^{2} & ^{2} \\ & + 4k_{l}k^{2} + _{c}4kb(Y_{l})k^{2} + 4k_{l}k^{2} + 4A_{c}^{22} + 2^{2}kb(Y_{l})k^{2} + 8^{2}A^{22} = (_{c}1 + (_{3}C_{1_{1}} + _{l} + _{1})k_{l}k^{2} + (_{6}^{2} + 4)kb(Y_{l})k^{2} + A^{2}(_{1}2^{2} + 4)^{2} & ^{2} &$$

where $C_2 = 3C_1 + 8$. Taking expectation on both sides, we have

For the term T₄ we have

$$T_{4} = 10 \frac{u^{2}_{CD}}{u^{2}_{CD}} \frac{X}{N} \left(1 + C_{2}\right)^{-1} \frac{B^{2}}{N} + \frac{2}{N} m^{222i}_{CD}^{CD}_{i=0}$$

$$= 10 \frac{u^{2}_{CD}}{l^{2}} \frac{B^{2}}{N} \frac{(1 + C_{2})^{+1}}{C_{2}} + 20 \frac{u^{2}_{D}^{2} + C_{2}^{2}}{l^{2}_{CD}^{2}} \left(\hat{I} + C_{2}\right)_{i=0} \times \frac{1}{1 + C_{2}}$$

$$10 \cdot \left[\frac{u^{2}_{CD}}{2N} \frac{B(\hat{I}^{2} + C_{2})^{+1}}{C_{2}} + 20 \frac{u^{2}_{D}^{2}}{l^{2}_{CD}} m^{22}_{C_{2}} \left(1 + C_{2}\right)_{i=0}^{2i} 10 \right] \frac{X^{2}_{D}^{2}}{l^{2}_{CD}^{2}} B^{2}$$

$$\left(1 + \frac{C_{2}}{l^{2}_{CD}}\right)^{+1}_{N} \frac{1}{N} + \frac{20}{C_{2}^{2}} \frac{u^{2}_{D}^{2}}{l^{2}_{CD}^{2}} m^{22}_{C_{2}^{2}} \left(1 + C_{2}^{4} - \frac{1}{l^{2}_{CD}^{2}}\right) \frac{1}{2^{2}}$$

$$(C_{2} \quad 0)$$

$$(106)$$

By the same argument as in Lemma B.6, and by the assumption on , we have $(1+C_2)^{+1}$ $1+2C_2(+1)$ $1+4C_2$ 2 and $(1+C_2)$ $1+2C_2$ 1+1=2 2. Hence, we have

$$T_4 \ 40 \left| \frac{q_{L_D}^2}{2} \right| \frac{B}{N} + \frac{m_2^2}{1^{-2}}$$
(107)

Furthermore, for T₅, we have

$$T_5 = {\overset{X}{(1 + C_2)}}^{i} E_{t} {_{2}[_{t+i}]}^{2} = 0$$

$$\begin{array}{c} X \\ (1+C_2) \, E_{t-2} \big[_{t-i} \big] \, i = \delta \end{array} \qquad \qquad (C_2 > 0) \\ X \\ (1+2C_2) E_{t-2} \big[_{t-i} \big] \, i = \delta \\ X \\ 2 \quad E_{t-2} \big[_{t-i} \big] \colon i = \delta \end{array} \qquad \qquad (assumption on)$$

Combining the bounds on T_4 (107) and T_5 (108) in (105) we have for any t I t,

$$E_{t} {}_{2}k_{1}k_{c} {}_{2}E_{t} {}_{2}k_{t} k_{c} + 40 \cdot \frac{2cD}{l} {}_{cD} {}^{2} {}_{CD} {}^{4} + m^{23} {}_{1} {}_{2} {}_{2} {}_{1} {}_{1} {}_{2} {}_{1} {}_{2} {}_{2} {}_{1} {}_{1} {}_{2} {}_{2} {}_{1} {}_{1} {}_{2} {}_{2} {}_{1} {}_{2} {}_{1} {}_{2} {}_{2} {}_{2} {}_{1} {}_{2} {}_{2} {}_{2} {}_{1} {}_{2} {}_{2} {}_{2} {}_{2} {}_{2} {}_{1} {}_{2} {}_$$

Furthermore, we have

Taking expectation on both sides, and using the bounds in Lemma B.9 and (109), we get

(110)

(108)

Proof of Lemma B.8. For sK + 1 t (s + 1)K 1, where s = bt = Kc,

Taking expectation

where (a) follows since

Finally, the inequality in (b) follows since

For simplicity, we define $K_1 = \frac{2A_1^2u_{c2}^2}{I_{c2}^2}$; $K_2 = \frac{2A_2^2u_{c2}^2}{I_{c2}^2}$; $B = \frac{u_{c2}}{I_{c2}}B$. Hence, (111) simplifies to

Recursively applying (112), going back 2 steps, we see

To derive the bound for going back, in general, j steps (such that t j sK), we use an induction argument. Suppose for going back k (< j) steps, the bound is

We derive the bound for k + 1 steps. For this, we further bound the last term in (114).

$$4^{2}(t \quad sK)(1 + A_{1})^{\sim} \quad \stackrel{\text{if}}{1} + 4^{2}(1 + A_{1})(t \quad ' \quad sK) \quad \stackrel{\text{if}}{E} \quad \stackrel{\text{if}}{K} \quad \stackrel{\text{if}}$$

Substituting (115) into (114), we see that the induction hypothesis in (114) holds. We can go back as far as the last instant of synchronization, j t s K. For j = t s K, we get

Next, using $1 + x e^x$ for x = 0, we get

t
$$\psi$$
 sKh i t 1 sK !
1+ 4²(1+ A₁)(t ' sK) exp 4²(1+ A₁)(t ' sK)
'=1 exp 2²(1+ A₁)(t sK)²

if is small enough such that $2^2(1 + A_1)(K - 1)^2$ In $\frac{5}{4}$, which holds true by the assumption on the step size. Using (117) in (116), we get

which concludes the proof.

Proof of Lemma B.9. We have

$$E_{t} r[kb(Y_{t})k_{c}] u_{cD}E_{t} r[kb(Y_{t})k_{D}]$$

$$= u_{cD}E_{t} r b(Y_{t}) Db(Y_{t})$$

$$q \underline{\qquad}$$

For the term T_1 we have

$$T_{1} = \frac{1}{N} t^{\frac{1}{N}} \sum_{i=1}^{N} E_{t-r} \int_{cD}^{2} b^{i}(y_{t}^{i}) \int_{c}^{2} \frac{1}{N I_{cD}} t^{\frac{1}{N}} \sum_{i=1}^{N} E_{t-r}[B^{2}]$$

$$= \frac{B}{I_{cD}} p \frac{1}{N}$$
(Assumption 6.3)

For the term T₂ we have

$$T_{2} = \frac{2}{N} \sum_{\substack{i < j \\ k \in L_{t} \ r \ b^{i}(y_{t}) \ k_{D} \ k \in L_{t} \ r \ b^{j}(y_{t}) \ k_{D}}}{k E_{t} \ r \ b^{i}(y_{t}) \ k_{D}}$$

$$= \frac{2}{N} \sum_{\substack{i < j \\ N \ i < j \ N}}^{S} \frac{X}{k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{j}(y_{t}^{j}) \ k_{c}}}$$

$$= \frac{2}{N} \sum_{\substack{i < j \\ N \ i < j \ N}}^{S} \frac{X}{k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{j}(y_{t}^{j}) \ k_{c}}}$$

$$= \frac{2}{N} \sum_{\substack{i < j \\ N \ i < j \ N}}^{S} \frac{X}{k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{j}(y_{t}^{i}) \ k_{c}}}$$

$$= \frac{2}{N} \sum_{\substack{i < j \\ N \ i < j \ N}}^{S} \frac{X}{k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{j}(y_{t}^{i}) \ k_{c}}}$$

$$= \frac{2}{N} \sum_{\substack{i < j \\ N \ i < j \ N}}^{S} \frac{X}{k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{j}(y_{t}^{i}) \ k_{c}}}$$

$$= \frac{2}{N} \sum_{\substack{i < j \\ N \ i < j \ N}}^{S} \frac{X}{k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{j}(y_{t}^{i}) \ k_{c}}}$$

$$= \frac{2}{N} \sum_{\substack{i < j \ N \ i < j \ N}}^{N} \frac{X}{k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{j}(y_{t}^{i}) \ k_{c}}}$$

$$= \frac{2}{N} \sum_{\substack{i < j \ N \ i < j \ N}}^{N} \frac{X}{k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{j}(y_{t}^{i}) \ k_{c}}}$$

$$= \frac{2}{N} \sum_{\substack{i < j \ N \ i < j \ N}}^{N} \frac{X}{k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{j}(y_{t}^{i}) \ k_{c}}}$$

$$= \frac{2}{N} \sum_{\substack{i < j \ N \ i < j \ N}}^{N} \frac{X}{k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{j}(y_{t}^{i}) \ k_{c}}}$$

$$= \frac{2}{N} \sum_{\substack{i < j \ N \ i < j \ N}}^{N} \frac{X}{k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{j}(y_{t}^{i}) \ k_{c}}}$$

$$= \frac{2}{N} \sum_{\substack{i < j \ N}}^{N} \frac{X}{k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \ k E_{t} \ r \ b^{i}(y_{t}^{i}) \ k_{c} \$$

Substituting (120) and (121) in (119), we get the result in (56).

The proof of (57) follows analogously.

Proof of Lemma B.10. By definition, we have

$$kG(t;Y_{t}) = G(t;Y_{t})k_{c} = \frac{1}{N} (G^{i}(t;Y^{i}) + G^{i}(t;Y^{i}))_{t} = 1$$

$$\frac{1}{N} G^{i}(t;Y^{i}) + G^{i}(t;Y^{i})_{c} + t = 1$$

$$\frac{1}{N} A_{1} k_{t} + t k_{c}$$

$$= (A_{1t})^{2}$$

$$\frac{1}{N} A_{1} k_{t}^{2} + t k_{c} = 1^{2}$$
(convexity of norm)
$$\frac{1}{N} A_{1} k_{t}^{2} + t k_{c} = 1^{2}$$
(convexity of square)

$$= A_{1}^{2}$$
 t: (By definition of t)

Furthermore, by the convexity of $()^2$, we have

Proof of Lemma B.11. Since $M_f^{;g}()$ is convex, and there exists a norm, k_m , such that $M_f^{;g}(x) = \frac{1}{2} k_m x^2_m$ (see Proposition B.1), using the chain rule of subdifferential calculus,

$$rM_f^{;g}(x) = kxkm u_x;$$

where ux 2 @kxkm is a subgradient of kxkm at x. Hence,

$$rM_f^{;g}(x)^? = kxkm ku_xkm^?$$

where $kk_{m}^{?}$ is the dual norm of kk_{m} Since kk_{m} is convex and, as a function of x, is 1-Lipschitz w.r.t. kk_{m} we have $ku_{x}k_{m}^{?}$ 1 (see Lemma 2.6 in (Shalev-Shwartz et al., 2012)).

Further, by convexity of kkm norm, k0km kxkm + hux; xi. Therefore,

D E
$$rM_f^{;g}(x); x = kxkm hu_x; xi kxkm^2 = 2M_f^{;g}(x)$$
:

C. _Federated TD-learning

C.1. On-policy Function Approximation

Proposition C.1. On-policy TD-learning with linear function approximation Algorithm 1 satisfies the following:

- 1. ' = v' , v
- 2. $S_t = (S_t^1; ...; S_t^N)$ and $A_t = (A_t^1; ...; A_t^N)$
- 3. $Y_t^i = (S_t^i; A_t^i; \dots; S_{t+n-1}^i; A_{t+n-1}^i; S_{t+n}^i)$ and $Y_t = (S_t; A_t; \dots; S_{t+n-1}; A_{t+n-1}; S_{t+n})$
- 4. : Stationary distribution of the policy.

Furthermore, choose some arbitrary positive constant > 0. The corresponding $G^{i}(^{i}; y_{t}^{i})$ and $b^{i}(y^{i})$ in Algorithm 1 for On-policy TD-learning with linear function approximation is as follows

1.
$$G^{i}(_{t}^{i};_{t}^{y})_{t}^{j} = _{t}^{i} + _{t}^{1}(_{S}^{i}) _{t}^{p} + _{t}^{1}(_{t}^{i}) _{t}^{j} (_{t}^{y})^{j} (_{t}^{y})^{j}$$

2.
$$b^{i}(y_{t}^{i}) = \frac{1}{4}S^{i})_{t}^{P_{t+n-1}} + R(S^{i}; A_{i}^{i}) + (S^{i+1})^{y} + (S^{i})^{y}$$

where v solves the projected bellman equation $v = ((T)^n v)$. Furthermore, the corresponding step size in Algorithm 4 is .

Lemma C.1. Consider the federated on-policy TD-learning Algorithm 1 as a special case of FedSAM Algorithm 4 (see Proposition C.1). Suppose the trajectory $fS_{t}^{i}g_{t=0;1;:::}$ converges geometrically fast to its stationary distribution as follows $d_{TV}(P(S_{t}^{i}=jS_{t}^{i})j_{j}^{i}))$ mfor all $i^{t}=1;2;:::;N$. The corresponding $G^{i}()$ in Assumption 6.1 for the federated TD-learning is as follows

$$G^{i}() = + {}^{>} (\frac{1}{n}(P)^{n});$$
 (122)

where > 0 is an arbitrary constant introduced in Proposition C.1. In addition, (6) holds. Furthermore, for t n + 1, we have m = $\max f \frac{2A_2}{\pi}$, 2B ng, where A₂ and B are specified in Lemma C.3 and = .

Lemma C.2. Consider the federated on-policy TD-learning 1 as a special case of FedSAM (as specified in Proposition C.1). Consider the jSj jSj matrix $U = {}^{>} ({}^{n}(P){}^{n} - I)$ with eigenvalues $f_1; :::; {}_{jSj}g$. Define ${}_{max} = {}^{max}{}_{i} {}_{j} {}_{i} {}_{j}$ and $= {}^{max}{}_{i} {}_{i} {}_{j} {}_{i} {}_{j}$ and $= {}^{max}{}_{i} {}_{i} {}_{j} {}_{i} {}_{j} {}_{i}$ and $= {}^{max}{}_{i} {}_{i} {}_{j} {}_{i} {}_$

Lemma C.3. Consider the federated on-policy TD-learning Algorithm 1 as a special case of FedSAM (as specified in Proposition C.1). There exist some constants A₁, A₂, and B such that the properties of Assumption 6.3 are satisfied.

Lemma C.4. Consider the federated on-policy TD-learning Algorithm 1 as a special case of FedSAM (as specified in Proposition C.1). Assumption 6.4 holds for this algorithm.

C.1.1. PROOFS

Proof of Proposition C.1. Items 1-4 are by definition. Subtracting v from both sides of the update of the TD-learning, we have

$$V_{\frac{t+1}{t+1}}^{i} \left\{ z - \frac{v}{t} \right\} = V_{t}^{i} - \left\{ z - \frac{v}{t} \right\} + \left\{ S_{t}^{i} \right\} = V_{t}^{i} - \left\{ z - \frac{v}{t} \right\} + \left\{ S_{t}^{i} \right\} = V_{t}^{i} - \left\{ z - \frac{v}{t} \right\} + \left\{ S_{t}^{i} \right\} = V_{t}^{i} - \left\{ z - \frac{v}{t} \right\} + \left\{ S_{t}^{i} \right\} + \left\{ S$$

which proves items 1 and 2. Furthermore, for the synchronization part of TD-learning, we have

$$v_{t}^{i} = \frac{1}{N} \sum_{j=1}^{X^{N}} v_{t}^{j}$$

$$=) \quad \psi_{t}^{i} \{ \underline{z}_{i}^{V} \} = \frac{1}{N} \sum_{j=1}^{X^{N}} (\gamma_{t}^{j} \{ \underline{z}_{i}^{V} \})^{*} \}$$

which is equivalent to the synchronization step in FedSAM Algorithm 4. Notice that here we used the fact that all agents have the same fixed point v.

Proof of Lemma C.1. It is easy to observe that

$$G^{i}(; y_{t})_{i} = + (S_{t})^{n}(S_{t+n})^{n}(S_{t})^{n}$$

Taking expectation with respect to the stationary distribution, we have

$$G^{i}() = E_{S_{t}^{-i}} + \frac{1}{(S^{i})_{t}^{-n}(S^{i}_{+n})_{t}^{>}} (S^{i})^{>}_{t}$$

$$= E_{S_{t}^{-i}} E + \frac{1}{(S^{i})_{t}^{-n}(S^{i}_{+n})_{t}^{>}} (S^{i})^{>}_{t} S^{i}_{t} t \qquad \text{(tower property of expectation)}$$

$$= E_{S_{t}^{-i}} + \frac{1}{(S_{t}^{-i})_{t}^{-n}E[()(S_{t+n})jS_{t}^{-i}]} ()(S_{t}^{-i}) = E_{S_{t}^{-i}} + \frac{1}{(S_{t}^{-i})_{t}^{-n}(P)^{n}} ()(S^{i}) = + \frac{1}{(S_{t}^{-i})_{t}^{-n}(P)^{n}}):$$

where P is the transition probability matrix corresponding to the policy , and is a diagonal matrix with diagonal entries corresponding to elements of .

As explained in (Tsitsiklis & Van Roy, 1997), the projection operator is a linear operator and can be written as = ($^>$) $^{1>}$, where is a diagonal matrix with diagonal entries corresponding to the stationary distribution of the policy . Hence, the fixed point equation is as follows v = ($^>$) $^{1>}((T)^n v)$. Since is a full column matrix, we can eliminate it from both sides of the equality, and further multiply both sides with $^>$. We have $^>v$ =

$$^{>}((T)^{n}v)$$
, and hence $^{>}((T)^{n}v - v) = 0$, which is equivalent to $E_{S}[^{>}(S)((T)^{n}v)(S) = 0$. By expanding $(T)^{n}$, we have $E_{S^{i}}[^{>}(S^{i}) - 0] = 0$ ($R(S^{i}; A^{i}) + (v)(S^{i}) = 0$, which means

$$E_{y} b^{i}(y) = 0;$$
 (123)

and proves (6).

Moreover, we have

$$kG^{i}() \quad E[G^{i}(;y^{i})]k_{c} = E_{y^{i}}[G^{i}(;y^{i})] \quad E[G^{i}(;y^{i})^{c} = t^{X} \quad (y^{i}) \\ \qquad \qquad \qquad \qquad P(y^{i} = y^{i}jy^{i})G^{i}(;y^{i}) \\ X \quad X \quad (y^{i}) \quad P_{t}(y_{t} = y_{t}^{i}jy^{i}):^{i}G^{i}(y;y^{i})_{c}y^{i} \quad t^{c} \\ \qquad \qquad X[y^{i}) \quad P(y^{i} = y^{i}j_{t}y^{i}):^{i}A_{2} \quad k_{j}kc:y^{i} \\ \qquad \qquad \qquad (kaxk_{c} = jajkxk_{c}) \\ \qquad \qquad \qquad t \quad \qquad (Assumption 6.3)$$

For brevity, we denote $P(S_i = s_i) = P(s_i)$. We have

$$\begin{array}{l} X \\ (y^{i})_{t} P(y^{i} = _{t}y^{i}jy^{i}_{t})_{0} \\ = X \\ = X \\ (s_{t}; a^{i}_{t}; :::; s_{t+n}) P(s_{t}; a^{i}_{t}; :::; s_{t+n}jy_{0})^{-1} \\ = X \\ (s_{t})(a_{t}js_{t})P(\dot{s}_{t+1}js^{i}_{t}; a_{t})::\dot{t} P(s_{t+n}j\dot{s}_{t+n}-\dot{1}; a_{t+n}-\dot{1}) \\ = X \\ (s_{t})(a_{t}js_{t})P(s_{t+1}js^{i}_{t}; a_{t})::\dot{t} P(s_{t+n}j\dot{s}_{t+n}-\dot{1}; a_{t+n}-\dot{1}) \\ = X \\ (s_{t})(a_{t}js_{t})P(s_{t+1}js_{t}; a_{t}):::P(s_{t+n}js_{t+n}-1; a_{t+n}-1) \\ = X \\ (s_{t})(a_{t}js_{t})P(s_{t}js_{n})(a_{t}js_{t})P(s_{t+1}js_{t}; a_{t}):::P(s_{t+n}js_{t+n}-1; a_{t+n}-1) \\ \end{array}$$

$$= {X \choose s^{i}}_{t} P(s^{i}jS^{i}_{t}) s_{h}$$

$$= 2d^{i}_{TV} (P(S^{i}_{t} = jS^{i}_{n})jj()) 2m$$

$$^{n} t$$

$$= (2m^{n}) : ^{t}$$

$$\begin{split} E[b^{i}(y^{i})]^{c} &= E[b^{i}(y^{i})] &\quad E_{y_{t}}[b^{i}(y^{i})]_{c} &\quad = \\ &\quad X &\quad t &\quad y^{i} &\quad p &\quad y^{i} &\quad y^{i} &\quad c \\ &\quad X &\quad y^{i} &\quad & c &\quad c \\ &\quad (y^{i}) &\quad P_{t}(y^{i} = y^{i}_{t}jy^{i})b^{i}_{t}(y^{i})_{c} &\quad c \\ &\quad X &\quad y^{i}_{t} &\quad & c \\ &\quad (y^{i}) &\quad P_{t}(y_{t} = y^{i}_{t}jy_{0})\dot{B} &\quad i \\ &\quad y^{i}_{t} &\quad & y^{i}_{t} &\quad (Assumption 6.3) \\ &\quad 2B\,m &\quad t &\quad & \end{split}$$

Proof of Lemma C.2. Consider the jSj jSj matrix $U = \binom{n}{P}^n = I$) with eigenvalues $f_1; \ldots; jSjg$. As shown in (Tsitsiklis & Van Roy, 1997), since is a full rank matrix, the real part of i is strictly negative for all $i = 1; \ldots; jSj$. Furthermore, define $max = max_i j_i j$ and $= max_i Re[i] > 0$, where Re[i] evaluates the real part. Consider the matrix $U^0 = I + \frac{1}{2 \frac{n}{max}} U$. It is easy to show that the eigenvalues of U^0 are $f_1 + \frac{1}{2 \frac{n}{max}} (1 + \frac{1}{2 \frac{n}{max}}) (1 + \frac{1}{$

$$1 + \frac{1}{2 + \frac{1}{2 + \frac{1}{2}}} = 1 + \frac{1}{2 + \frac{1}{2 + \frac{1}{2}}} = 1 + \frac{1}{2 + \frac{1}{2 + \frac{1}{2}}} = \frac{1}{2 + \frac{1}{2 + \frac{1}{2}}} = \frac{1}{2 + \frac{1}{2 + \frac{1}{2}}} = 1 + \frac{1}{2 + \frac{1}{2 + \frac{1}{2}}} = \frac{1}{2 + \frac{1}{2 +$$

Hence, all the eigenvalues of U^0 are in the unit circle. By (Bertsekas et al., 1995), Page 46 footnote, we can find a weighted 2-norm as $kk = \frac{1}{2}$ such that U^0 is contraction with respect to this norm with some contraction factor C_0 . In particular, there exist a choice of such that we have $C_0 = \frac{1}{8}$

Proof of Lemma C.3. The existence of A_1 and A_2 immediately follows after observing that $G^i(^i; y^i)_t$ is a linear function of i . Furthermore, the result on B follows due to v being bounded as shown in (Chen et al., 2021b).

Proof of Lemma C.4. For the sake of brevity, we write $S_t^i = s_t^i$ simply as s_t^i , and similarly for other random variables. We have

$$E_{t}$$
 r[f(y_t^i) g(y_t) \dot{j}

Proof of Theorem 4.1. By Proposition C.1 and Lemmas C.1, C.2, C.3, and C.4, it is clear that the federated TD-learning with linear function approximation Algorithm 1 satisfies all the Assumptions 6.1, 6.2, 6.3, and 6.4 on the FedSAM Algorithm 4. Furthermore, by the proof of Theorem B.1, we have $w_t = (1 \frac{2}{2})^{-t}$, and the constant c_{TDL} in the sampling distribution q c_{TDL} in Algorithm 1 is $c_{TDL} = (1 \frac{2}{2})^{-t}$. Furthermore, by choosing the step size small enough, we can satisfy the requirements in (21), (23), (32), (35). By choosing K large enough, we can satisfy K > . Hence, the result of Theorem B.1 holds for this algorithm with some $c_{TDL} > 1$. Also, it is easy to see that $c_{TDL} = c_{TDL} = c_{$

Next, we derive the constant c_{TDL} . Since $k_c = k_c$, which is smooth, we choose $g() = \frac{1}{2}k_c k_c^2$. By taking = 1, we have $l_{cs} = u_{cs} = 1$. Therefore, we have $l_{cs} = 1$, and $l_{cs} = 1$, and $l_{cs} = 1$, and $l_{cs} = 1$, where $l_{cs} = 1$

C.2. Off-policy Tabular Setting

In this subsection, we verify that the Off-policy federated TD-learning Algorithm 2 satisfies the properties of the FedSAM Algorithm 4. In the following, V is the solution to the Bellman equation (124).

$$V(s) = \begin{cases} X & \text{"} & \text{#} \\ (ajs) & R(s;a) + P(s^{0}js;a)V(s^{0}) \end{cases}$$
(124)

Note that V is independent of the sampling policy of the agent. Furthermore, we take $k_c = k_1$.

Proposition C.2. Off-policy n-step federated TD-learning is equivalent to the FedSAM Algorithm 4 with the following parameters.

2.
$$S_t = (S_t^1; ...; S_t^N)$$
 and $A_t = (A_t^1; ...; A_t^N)$

3.
$$Y_t^i = (S_t^i, A_t^i, \dots, S_{t+n-1}^i; A_{t+n-1}^i; S_{t+n}^i)$$
 and $Y_t = (S_t^i, A_t^i, \dots, S_{t+n-1}^i; A_{t+n-1}^i; S_{t+n}^i)$

4. i : Stationary distribution of the sampling policy of the i-th agent.

5.
$$G^{i}(_{t}^{i}; y_{t}^{i})_{s}^{s} = _{t}^{i}(s) + 1_{fs=S_{t}g} + 1_{fs=S_{t}g} + 1_{i} + 1_{i} + 1_{i} = _{t}^{i}(_{t}^{i})(_{j}^{i}; A^{i})_{j}^{i}(S^{i})_{j}^{i} + 1_{t}^{i}(S^{i})_{i}$$

6.
$$b^{i}(y_{t}^{i})_{s} = \mathbf{1}_{fs=S_{t}^{i}g} \frac{P_{t+n-1}}{I_{t}} = \mathbf{1}_{t}^{I_{t}} (S^{i}; A^{i})_{r} (S^{i}; A^{i})_{r} + \gamma(S^{i})_{r} (S^{i})_{r}$$

Lemma C.5. Consider the federated off-policy TD-learning Algorithm 2 as a special case of FedSAM (as specified in Proposition C.2). Suppose the trajectory $fS_t^jg_{t=0;1;...}$ converges geometrically fast to its stationary distribution as follows

 $d_{TV}(P(S_t^i = jS_t^i)jj^i())$ mfor all i = 1; 2; ...; N. The corresponding $G^i()$ in Assumption 6.1 for the federated TDlearning is as follows

$$G^{i}()_{s} = (s) + {}^{ni}(P)^{n} {}^{i}(s):$$

Furthermore, for t n + 1, we have $m_1 = \frac{2A_n m_p^2}{M_p}$ where A_2 is the constant specified in Assumption 6.3, $m_2 = 0$, and =

Lemma C.6. Consider the federated off-policy TD-learning 2 as a special case of FedSAM (as specified in Proposition C.2). The corresponding contraction factor c in Assumption 6.2 for this algorithm is $c = 1 \min(1^{n+1})$, where $\min =$ $\min_{s;i} (s)$.

Lemma C.7. Consider the federated off-policy TD-learning 2 as a special case of FedSAM (as specified in

Proposition C.2). The constants A₁, A₂, and B in Assumption 6.3 can be chosen as follows: A₁ = A₂ = 1 +
$$\binom{n}{1+\binom{n}{1+\frac{1}{max}}}\binom{1}{1+\binom{n}{max}}$$
; and B = $\binom{2}{1}\frac{1}{max}\binom{n}{1+\frac{n}{max}}\binom{1}{1+\frac{n}{max}}$ o:w: $\binom{n}{1+\binom{n}{max}}\binom{1}{1+\frac{n}{max}}\binom$

Lemma C.8. Consider the federated off-policy TD-learning 2 as a special case of FedSAM (as specified in Proposition C.2). Assumption 6.4 holds for this algorithm.

C.2.1. PROOFS

Proof of Proposition C.2. Items 1-4 are by definition. Furthermore, by the update of the TD-learning, and subtracting V from both sides, we have

$$\begin{array}{c} V \stackrel{i+1}{\leftarrow} \{z \\ \frac{1}{c^{i}} \{s \} \\ -\frac{1}{c^{i}} \{s \} \\ \end{array} = V \stackrel{te}{\leftarrow} \{z \\ \frac{1}{c^{i}} \{s \} \\ \end{array} = V \stackrel{te}{\leftarrow} \{z \\ \frac{1}{c^{i}} \{s \} \\ \end{array} = V \stackrel{te}{\leftarrow} \{z \\ \frac{1}{c^{i}} \{s \} \\ = V \stackrel{te}{\leftarrow} \{z \\ \frac{1}{c^{i}} \{s \} \\ \end{array} = V \stackrel{te}{\leftarrow} \{z \\ \frac{1}{c^{i}} \{s \} \\ = V \stackrel{te}{\leftarrow} \{z \} \\ \end{array} = V \stackrel{te}{\leftarrow} \{z \\ \frac{1}{c^{i}} \{s \} \\ = V \stackrel{te}{\leftarrow} \{s \} \\ = V \stackrel{te}{\leftarrow$$

$$V_{t}^{i} = \frac{1}{N} \sum_{j=1}^{X^{N}} V_{t}^{j}$$

$$= \int_{t}^{V_{t}^{i}} \frac{1}{N} \sum_{j=1}^{X^{N}} \left(\underbrace{V_{t}^{j}}_{t} \{z_{j}^{V}\} \right);$$

which is equivalent to the synchronization step in FedSAM Algorithm 4. Notice that here we used the fact that all agents have the same fixed point V.

Proof of Lemma C.5. By taking expectation of $G^{i}(^{i}; y_{t})^{i}_{s}$, we have

$$G^{i}()_{s} = E_{S^{i}_{t}}(s) + \mathbf{1}_{f_{S} = S^{i}_{g}} t^{t} = t \int_{-t}^{t_{f}} (S^{i}; A^{i})_{j} (S^{i}_{1})_{h} (S^{i})_{l = f}$$

$$= (s) + E_{S^{i}_{t}} \mathbf{1}_{f_{S} = S_{t}} e^{\int_{-t}^{t_{f}} (S^{i}; A^{i})_{j}} (S^{i}; A^{i})_{j} (S^{i}_{1})_{h} (S^{i})_{l} : \frac{Z}{T_{l}}$$

Denote $E_k^i[] = E[jfS^i; A^ig_{rk-1}; S^i].$ For T_I , we have

Here,

$$E_{1}^{i} I^{(i)}(S_{1}^{i}; A_{1}^{i})(S_{1+1}^{i}) = X P(S_{1+1}^{i} = s; A^{i} = ajS^{i})I^{(i)}(S^{i}; A_{1}^{i} = a)(s) s; a$$

$$= X P(S_{1+1}^{i} = s; A^{i} = ajS^{i}) \frac{(ajS^{i})_{1}}{(ajS_{1}^{i})} (s)$$

$$= X Y (ajS^{i})_{1}P(sjS^{i}; a)_{1} \frac{(ajS_{1}^{i})}{(ajS_{1}^{i})}$$

$$= Y Y (sjS^{i}; a)(ajS^{i})(s)_{1} [P](S^{i}); s; a$$

$$= P(sjS^{i}; a)(ajS^{i})(s)_{1} [P](S^{i}); s; a$$

where we denote i as diagonal matrix with diagonal entries corresponding to the stationary distribution i. Hence, in total we have

Furthermore, by the same argument as in the proof of Lemma C.1, we have

$$kG^{i}()$$
 $E[G^{i}(;y_{t})]\dot{k_{c}}^{X_{i}}(y_{t})$ $P(\dot{y}_{t} = y_{t}j\dot{y_{0}}):A_{2}^{i}k\dot{k}C;y_{t}$

and

$$X_{i}(y^{i})_{t} P(y^{i} = y^{i}jy)(2m^{i});$$

which proves that $m_1 = 2m_2^{n}$ constant. In addition, we have

$$\begin{split} & E[b^{i}(y^{i})]^{c} = E \quad \mathbf{1}_{fs=S^{i}g} \quad \overset{t \stackrel{\leftarrow}{N}^{1}}{\underset{i=t}{\overset{\leftarrow}{1}}} I^{(i)}(S_{j}; A_{j}) \quad R(S_{j}; A_{l}) + iV(S_{l+1}) \quad V(S_{l}) \\ & = E \quad \mathbf{1}_{fs=S^{i}g_{t}} \quad \overset{t \stackrel{\leftarrow}{N}^{1}}{\underset{i=t}{\overset{\leftarrow}{1}}} I^{(i)}(S^{i}; A_{j}^{i}) \quad J^{i}(S^{i}; A_{l}^{i}) + V(S^{i}) \\ & = E \quad \mathbf{1}_{fs=S^{i}g_{t}} \quad \overset{t \stackrel{\leftarrow}{N}^{1}}{\underset{i=t}{\overset{\leftarrow}{1}}} I^{(i)}(S^{i}; A_{j}^{i}) \quad J^{i}(S^{i}; A_{j}^{i}) + V(S^{i}) \quad J^{i}(S^{i}; A_{j}^{i}) \quad J^{i}(S^{i}; A_{j}^{i}$$

For the term T, we have

$$T = X | (ajS^{i})_{i} - \frac{(ajS_{1})^{i}}{(ajS)^{i}} | (ajS_{1})^{i} | (ajS_{$$

which shows that $m_2 = 0$.

Proof of Lemma C.6.

From of Lemma C.6.
$$kG^{i}(_{1}) \quad G^{i}(_{2})k_{c} = _{1} + _{}^{n+1i}(P)^{n+1}{_{1}} \quad _{1}^{i} \quad _{2} + _{}^{n+1i}(P)^{n+1}{_{2}} \quad _{2}^{i}{_{2}}_{_{1}} = I \quad _{i}^{i}(I \quad _{}^{n+1}(P)^{n+1}) \quad (_{1} \quad _{2})_{_{1}} \quad _{I}^{i}(I \quad _{1}^{n+1}(P)^{n+1}) \quad k_{1} \quad _{2}k\mathbf{1} : \qquad (definition of matrix norm)$$
 Since the elements of the matrix $I \quad _{i}^{i}(I \quad _{1}^{n+1}(P)^{n+1})$ is all $_{1}^{1}$ positive, we have $I \quad _{i}^{i}(I \quad _{1}^{n+1}(P)^{n+1}) \quad = (I \quad _{1}^{i}(I \quad _{1}^{n+1}(P)^{n+1}))\mathbf{1}^{1} = k\mathbf{1} \quad _{1}^{i}(\mathbf{1} \quad _{1}^{n+1}\mathbf{1})k_{1} = \mathbf{1} \quad _{1}^{i}(I \quad _{1}^{n+1}\mathbf{1}) \quad \mathbf{1} \quad _{min}(\mathbf{1} \quad _{1}^{n+1}).$

Proof of Lemma C.7.

$$kG^{i}(1;y) \quad G^{i}(2;y)k_{c} \qquad \qquad ! \\ = \max_{s} 1(s) \quad 2(s) + \mathbf{1}_{fs=S^{i}g} \quad \lim_{t \to t} 1 \quad \lim_{s \to t} |S^{i}(s)| \quad 1(s_{i+1}^{i}) \quad 1(s_{i$$

Furthermore, we have

$$kb^{i}(\gamma_{t}^{i})k_{c} = \max_{\substack{S_{t}^{i};A_{t}^{i}::::S^{i}\\ |S_{t}^{i}|}} 1_{\substack{f_{S}=S^{i};g\\ |S_{t}^{i}|$$

Proof of Lemma C.8. The proof follows similar to Lemma C.4.

Proof of Theorem 5.1. By Proposition C.2 and Lemmas C.5, C.6, C.7, and C.8, it is clear that the federated off-policy TD-learning Algorithm 2 satisfies all the Assumptions 6.1, 6.2, 6.3, and 6.4 of the FedSAM Algorithm 4. Furthermore, by the proof of Theorem B.1, we have $w_t = (1 \frac{2}{2})^{-t}$, and the constant c in the sampling distribution q^c in Algorithm 1 is $\frac{2}{2}$ 1. In equation (125) we evaluate the exact value of w_t .

Furthermore, by choosing step size small enough, we can satisfy the requirements in (21), (23), (32), (35). By choosing K large enough, we can satisfy K > 1, and by choosing T large enough we can satisfy E = 1. Hence, the result of Theorem B.1 holds for this algorithm.

Next, we derive the constants involved in Theorem B.1 step by step. After deriving the constants C₁, C₂, C₃, and C₄ in Theorem B.1, we can directly get the constants $C_i^{TD_T}$ for i = 1; 2; 3; 4.

In this analysis we only consider the terms involving jSj, jAj, $\frac{1}{1}$, $\frac{1}{1}$ max, and min. Since k k_c = k k₁, we choose g() = $\frac{1}{2}$ k k^2 , i.e. the p-norm with p = $2 \log(jSj)$. It is known that g() is (p 1) smooth with respect to k k_p norm (Beck, 2017), and hence L = (log(jSj)). Hence, we have $l_{cs} = jSj^{-1=p} = \frac{1}{1-e} = (1)$ and $u_{cs} = 1$. Therefore, we have $l_{cs} = \frac{1+e^{-u}}{1+e^{-u}} = \frac$

have '1 =
$$\frac{1+}{1+} - \frac{u_{cs}^2}{\frac{2}{cs}} = \frac{1+}{1+\frac{n}{2c}}$$
 1 + . By choosing = $(\frac{1+c}{2c})^2$ 1 = $\frac{1+2-\frac{3}{c}c}{4-\frac{2}{c}}(\frac{1}{2}-c) = \min(1-\frac{n+1}{2}) = \frac{1+\frac{n+1}{2}}{1+\frac{n+1}{2}}$

$$(_{min}(1 \quad)), \text{ which is } = O(1), \text{ we have } '_1 = \frac{\frac{1+}{1+\frac{e}{e}}}{1+\frac{e}{e}} = \frac{p}{e} \frac{\frac{(\frac{1+}{c}-c)^2}{\frac{2-c}{c}}}{\frac{2-c}{2-c}} = O(1), \text{ and }$$

$$s = \frac{1}{1+\frac{e}{1+\frac{e}{e}}} = 1 \quad c \quad \frac{V}{1+\frac{(\frac{1+}{2}c)^2-1}{2-e}} = 1 \quad \frac{O:5(1+c)e^{1-4}}{r} = 1 \quad r \cdot \frac{0:5(1+c)e^{1-4}}{p} = 1 \quad r \cdot \frac{0:5e^{1-4}(2-\min(1-n+1))}{p} = 1 \quad r \cdot \frac{2-\min(1-n+1)}{p} = 1 \quad r \cdot \frac{1+c}{2-2\min(1-n+1)} = 1 \quad c \cdot \frac{1+c}{2-2} = 1 \quad r \cdot \frac{1+c}{2-2\min(1-n+1)} = 1 \quad c \cdot \frac{1+c}{2-2} = 0:5\min(1-n+1) = 1 \quad c \cdot \frac{1+c}{2-2} = 0 \quad c \cdot \frac{1+c}{2-2} = 0$$

Using '2, we have

Further, we have

$$I_{cm} = (1 + I_{cs}^2)^{1=2} = (1)$$

 $U_{cm} = (1 + U_{cs}^2)^{1=2} = (1)$

Since TV-divergence is upper bounded with 1, we have m= O(1). By Lemma C.7, we have

and
$$A_1 = A_2 =$$

$$(1),$$

$$B = \frac{2I_{max}}{1} \frac{n}{I_{max}} = 0$$
if $I_{max} = 1$
o:w: $I_{max} = 0$

$$I_{max} = 0$$

(1). Hence $m_1 = \frac{2A_n m}{n} = O(I_{max})$. Also, we have $m_2 = 0$.

We choose the D-norm in Lemma B.9 as the 2-norm k k_2 . Hence, by primary norm equivalence, we have $I_{cD} = p_j \frac{1}{s_j}$ and $u_{cD} = 1$, and hence $\frac{u_{cD}}{I_{cD}} = \frac{p}{jSj}$.

We can evaluate the rest of the constants as follows

and similarly

$$_{3} = \frac{s}{\frac{1}{10}} \frac{\frac{1}{2} \frac{1}{CS}}{\frac{1}{10} \frac{1}{1} \frac{1}{1} \frac{1}{1}} = \frac{\min_{j \in S} (1)}{\lim_{j \in S} (1)} \frac{1}{\log(1S)}$$

$$q = 3 A_2^2 + 1 = O I_{max}^{n};$$

and

$$C_2 = 3C_1 + 8 = O I_{max_i};$$

$$C_3 = \frac{m_2 L}{I_{cs}^2} = 0;$$

$$\begin{array}{l} C \\ 4 \\ = \\ \begin{array}{l} \frac{L^2}{2^{\frac{1}{4}}} + \frac{2LA_2}{cs} \frac{1}{l^2} + \frac{cm}{2^{\frac{1}{2}}} \frac{L(A_1 + 1)}{l^{\frac{2}{5}}} - \frac{3m_1L^2}{2^{\frac{1}{2}}} \frac{cs}{l^2} \\ - \frac{log^2(jSj)}{log(jSj)l_{\frac{max}{2}}} - \frac{log(jSj)}{log(jSj)l_{\frac{max}{2}}} \frac{log(jSj)l_{\frac{max}{2}}}{log(jSj)l_{\frac{max}{2}}} - \frac{log(jSj)l_{\frac{max}{2}}}{log(jSj)l_{\frac{max}{2}}} - \frac{log(jSj)l_{\frac{max}{2}}}{log(jSj)l_{\frac{max}{2}}} \\ = O \\ - \frac{log_2(jSj)l_{\frac{max}{2}}}{log(jSj)l_{\frac{max}{2}}} + \frac{log(jSj)l_{\frac{max}{2}}}{log(jSj)} - \frac{log(jSj)l_{\frac{max}{2}}}{log(jSj)l_{\frac{max}{2}}} - \frac{log(jSj)l_{\frac{max}{2}}}{log(jSj)l_{\frac{max}{2}}} \\ = O \\ - \frac{log_2(jSj)l_{\frac{max}{2}}}{log(jSj)l_{\frac{max}{2}}} - \frac{log(jSj)l_{\frac{max}{2}}}{log(jSj)} - \frac$$

$$C_7 = \frac{m_2^2 u_{cm}^2 L_2^2}{2_1^2 t_{cs}^4} = 0;$$

$$C_9 = \frac{8u_{cD}B}{I_{cD}} = O \qquad \frac{p_{jSjI_{max}}!}{1}$$

$$C_{10} = \frac{8 \, m_2 u_{cD}}{I_{cD} (1)} = 0;$$
 ;

$$C_{11} = \frac{8C_{1}^{2}C_{3}^{2}u_{cm}^{2}}{\frac{2}{6}} = 0;$$

$$C_{12} = \frac{8C_4 u_{2D} B^2}{12^6 - 12^6} = O_{\frac{1}{2}} \frac{\log_{(j_s}S_j) I_{\frac{max}{1}}}{m^j i_m^j i_m^j$$

$$C_{13} = \frac{14C_4u_{cD}^2m_2^2}{|\mathcal{E}^D(1)|^2} = 0$$
:

$$\begin{split} C_{14}() &= & C_7 + C_{11} + 0.5C^2 \varsigma^2 _{g} + C_3 C_{10} + 2C_1 C_3 C_{10} + 3A_1 C_3 + C_{13} + C_8 \frac{\mu_D^2}{|I|^2} 2m^2 _2^2 - ^2 = 0; cD \\ C_{16}() &= & C_8 \frac{u}{|I_c \frac{cD}{cD}|^2} + \frac{1}{2} + \frac{1}{2} C_{12}^2 = O \frac{jSj \log(jSj) I^{2n} _{\underline{x}} + 1 + \frac{jSj \log_2(jSj) I^{3n} - \frac{1}{2}}{2 - \frac{max}{min}} \\ &= & C_{17} = (3A_1 C_3 + 8A_1^2 C_4 + C_5 + C_6) \\ &= & O - & O + & I_{max} : \frac{2n^2}{(jSj)^2} \frac{log_1^2(jSj) I_{max}^2}{log_2^2} \frac{log_1^2(jSj) I_{max}^2}{3 - \frac{n}{m} I_n(1)^3} \frac{log_2^2(jSj) I_{max}^2}{3 - \frac{n}{m} I_n(1)^3} \frac{log_1^2(jSj)}{3 - \frac{n}{m} I_n(1)^3} \\ &= & O - & \frac{I_0 \frac{3n_a x^3 \log^2(jSj)}{3 - \frac{n}{m} I_n(1)^3}}{1 - \frac{n}{3}} : \end{split}$$

Similar to k $k_D = k k_2$, we have $l_{c2} = p_j \frac{1}{s^j}$ and $u_{c2} = 1$.

$$A_{1} = \frac{2A_{\frac{1}{2}u_{c2}^{2}}}{I_{c2}^{2}} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{2} = \frac{2A_{\frac{2}{2}c^{2}}}{I_{c2}^{2}} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{3} = \frac{2A_{\frac{2}{2}c^{2}}}{I_{c2}^{2}} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{4} = \frac{2A_{\frac{2}{2}c^{2}}}{I_{c2}^{2}} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{5} = \frac{1}{I_{max}^{2}} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{6} = \frac{1}{I_{max}^{2}} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{7} = \frac{1}{I_{max}^{2}} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{8} = \frac{1}{I_{max}^{2}} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{9} = \frac{1}{I_{max}^{2}} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{1} = \frac{1}{I_{max}^{2}} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{1} = \frac{1}{I_{max}^{2}} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{1} = \frac{1}{I_{max}^{2n}} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{2} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{2} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{2} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{2} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{3} = O I_{max}^{2n}{}_{2}JSj;$$

$$A_{4} = O I_{max}^{2n}{}_{2}JSj;$$

$$A$$

$$\begin{split} &C_{1} = 16u_{cm}^{2}M_{0} \text{ og } e^{\frac{1}{4}}, \frac{1}{2} = O \quad \frac{1}{(1\frac{4n}{max}^{2}}\frac{1}{2nin(1)} = O \quad \frac{1}{(1-\frac{x}{max}^{2})\frac{1}{min}^{2}} \\ &= \frac{8u_{2m} \quad C_{8} + \frac{1}{4} + C_{12}}{\min(1-\frac{x}{max}^{2})} \quad 1 \quad \log(jSj) \quad jSj \log_{2}(jSj)l^{3n-1} '2 \\ &C_{2} = \frac{1}{\min(1-\frac{x}{max}^{2})} = O \quad \frac{jSj \log_{2}(jSj)l^{3n-1}}{\min(1-\frac{x}{max}^{2})} \\ &= O \quad \frac{jSj \log_{2}(jSj)l^{3n-1}}{\frac{i}{m2n}(1-\frac{x}{max}^{2})} \\ &= O \quad \frac{i}{\min(1-\frac{x}{max}^{2})} \quad 1 \quad jSj^{2}l^{4n} \times l^{3n-3} \log_{2}(jSj) '\frac{2}{2} \\ &C_{3} = \frac{1}{\min(1-\frac{x}{max}^{2})} \quad 1 \quad jSj^{2}l^{4n} \times l^{3n-3} \log_{2}(jSj) '\frac{2}{2} \\ &= O \quad \frac{max}{max} \ln(1-\frac{x}{max}^{2}) \\ &= O \quad \frac{max}{max} \ln(1-\frac{x}{max}$$

Finally, for the sample complexity result, we simply employ Corollary B.1.1.

D. ederated Q-learning

In this section, we verify that the federated Q-learning algorithm 3 satisfies the properties of the FedSAM Algorithm 4. In the following, Q is the solution to the Bellman optimality equation (126)

$$Q(s; a) = R(s; a) + E_{S0}$$
 $a^0 \max Q(S^0; a^0)$: (126)

Note that Q is independent of the sampling policy of the agent. Furthermore, $k k_c = k k_1$.

Proposition D.1. Federated Q-learning algorithm 3 is equivalent to the FedSAM Algorithm 4 with the following parameters.

- 1. $^{i} = Q^{i} + Q$
- 2. $S_t = (S_t^1; :::; S_t^N)$ and $A_t = (A_t^1; :::; A_t^N)$
- 3. $y_t^i = (S_t^i; A_t^i; S_{t+1}^i; A_{t+1}^i)$ and $Y_t = (S_t; A_t; S_{t+1}; A_{t+1})$
- 4. i : Stationary distribution of the sampling policy of the i'th agent.

5.
$$G^{i}(^{i}; y^{i})_{\{s;a\}} = ^{i}(s_{t}^{*}a)$$

+ $\mathbf{1}_{fS_{t}^{i}=s; A_{t}^{i}=ag} \max_{a^{0}} ^{i} + Q(S_{i+1}^{i}; a_{t}^{0})$ $^{i}(S^{i}; A_{t}^{i})_{t} \max_{a^{0}} Q(S_{t+1}^{i}; a^{0}) 6.^{i}$
 $b^{i}(y^{i})_{(s;a)} = \mathbf{1}_{fS^{i}=s; A^{i}=ag} R(S^{i}; A^{i}) + \max_{a^{0}} Q(S_{t+1}^{i}; a^{0}) Q(S^{i}; A^{i})$

where 1_A is the indicator function corresponding to set A, such that $1_A = 1$ is A is true, and 0 otherwise.

Lemma D.1. Consider the federated Q-learning Algorithm 3 as a special case of FedSAM (as specified in Proposition D.1). Suppose the trajectory fS^i_t ; $A^ig_{t=0;1;:::}$ converges geometrically fast to its stationary distribution as follows d_{TV} (P ($S^i_t=jS^i;_{\partial}A^i)_{ij}^i(;)$) mfor all i=1;2;:::;N. The corresponding $G^i()$ in Assumption 6.1 for the federated Q-learning is as follows

$$\begin{array}{c} & & & & \\ & G^{i}()_{(s;a)} = & (s;a) + {}^{i}(s;a) \ E_{S^{0}P(js;a)} & max(+ Q(S^{0};a^{0})) & (s;a) & maxQ(S^{0};a^{0}) \\ \vdots & & & \\ & & a^{0} & & \end{array}$$

Furthermore, we have $m_1 = 2A_2 m$, where A_2 is specified in Lemma D.3, $m_2 = 0$, and = ...

Lemma D.3. Consider the federated Q-learning as a special case of FedSAM (as specified in Proposition D.1). The constants A_1 , A_2 , and B_1 in Assumption 6.2 are as follows: $A_1 = A_2 = 2$ and $B_2 = \frac{2}{1}$.

Lemma D.4. Consider the federated Q-learning as a special case of FedSAM (as specified in Proposition D.1). Assumption 6.4 holds for this algorithm.

D.1. Proofs

Proof of Proposition D.1. Items 1-4 are by definition. Furthermore, by the update of the Q-learning, and subtracting Q from both sides, we have

$$\begin{aligned} & \underbrace{Q_{t+1}^{i}(s;a)}_{t+1} \underbrace{\{Z_{j,a}^{i}\}}_{t+1} = \underbrace{Q_{t}^{i}(s;a)}_{t+1} \underbrace{Q_{t}^{i}(s;a)}_{t+1} \underbrace{Q_{t}^{i}(s;a)}_{t+1} + \underbrace{1_{f(s;a)=(S_{t}^{i};A_{t}^{i})g}}_{t} & R(S_{t};iA_{t})i + \max_{a} Q_{t}(S_{t+1};a) & Q_{t}(S_{t};A_{t}) i \end{aligned}$$

$$= \underbrace{1_{f(s;a)=(S_{t};A_{t})g}}_{t} \underbrace{R(S_{t};A_{t}) + \max_{a} Q_{t}^{i}(S_{t+1};a)}_{t} \underbrace{Q_{t}^{i}(S_{t+1};a) + Q_{t}^{i}(S_{t+1};a)}_{t} + Q_{t}^{i}(S_{t+1};a) + Q_{t}^{i}(S_{t+1};a)}_{t} \underbrace{Q_{t}^{i}(S_{t+1};a)}_{t} + Q_{t}^{i}(S_{t+1};a) + Q_{t}^{i}(S_{t+1};a) + Q_{t}^{i}(S_{t+1};a)}_{t} \underbrace{Q_{t}^{i}(S_{t+1};a)}_{t} + Q_{t}^{i}(S_{t+1};a) + Q_{t}^{i}(S$$

which proves items 5 and 6. Furthermore, for the synchronization part of Q-learning, we have

$$Q_{t}^{i} = \frac{1}{N} \sum_{j=1}^{X^{N}} Q_{t}^{j}$$

$$= \underbrace{Q_{t}^{i}}_{j} \underbrace{Q_{t}^{i}}_{j} \underbrace{Q_{t}^{j}}_{j} \underbrace{Q_{$$

which is equivalent to the synchronization step in FedSAM Algorithm 4. Notice that here we used the fact that all agents have the same fixed point Q.

Proof of Lemma D.1. $G()_{(s;a)}$ can be found by simply taking expectation of $G^{i}(^{i}; y_{t}^{i})_{(s_{t}a)}$ defined in Proposition D.1, with respect to the stationary distribution . Furthermore, we have

In addition, we have

$$kE[b^{i}(y_{t}^{i})]k_{c} = \max_{s;a} P(S_{t}^{i} = s; A_{t}^{i} = ajS_{0}^{i}; A_{0}^{i}) \quad R(s;a) + E_{S^{0}P(js;a)}[\max_{a^{0}} Q(S^{0}; a^{0})] \quad Q(s;a)$$

$$= 0: \qquad \qquad (Bellman optimality equation (126))$$

Proof of Lemma D.2.

Next, we note that for any functions f() and g(), we have

$$(\max_{x} f(x))$$
 $(\max_{x} g(x))$ $\max_{x} jf(x)$ $g(x)j$:
$$(127)$$

The reason is as follows. We have $\max_x f(x) = \max_x f(x)$ g(x) + g(x) $(\max_x f(x) = g(x)) + (\max_x g(x))$. Hence, $(\max_x f(x))$ $(\max_x g(x))$ $\max_x f(x)$ g(x) $\max_x f(x)$ g(x). Now suppose $\max_x f(x)$ $\max_x g(x)$. Then we can apply absolute value to the left hand side of the inequality, and we get the bound. By a similar argument for the case $\max_x f(x)$ $\max_x g(x)$, we get the bound in (127). Hence, we have $\max_{x \in A} f(x) = \max_{x \in A} f(x)$ $\max_{x \in$

Proof of Lemma D.3. First, for A₁, we have

$$kG^{i}(_{1};y)$$
 $G^{i}(_{2};y)k_{c}$

$$= \max_{s;a} (s;a) + 1_{fS=s;A=ag} \max_{s;a} (1 + Q(S^0;a^0)) \quad _1(S;A) \quad \max_{s} Q(S^0;a^0)$$

$$= \max_{s;a} (1 + 1_{fS=s;A=ag}) (1(s;a) \quad _2(s;a)) \quad _{a^0}$$

$$+ 1_{fS=s;A=ag} \max_{a^0} (1 + Q(S^0;a^0)) \quad \max_{s;a} (1 + Q(S^0;a^0))$$
 (triangle inequality)
$$\max_{s;a} k_1 \quad _2 k + 1_{fS=s;A=ag} \max_{a^0} (1 + Q(S^0;a^0)) \quad \max_{s;a} (1 + Q(S^0;a^0))$$
 (definition of k k₁)
$$\max_{s;a} k_1 \quad _2 k + 1_{fS=s;A=ag} \max_{s;a} (1 + Q(S^0;a^0)) \quad \max_{s;a} (2 + Q(S^0;a^0))$$
 (By (127))
$$2k_1 \quad _2 k_1$$

$$= 2k_1 \quad _2 k_2$$

Second, for A₂, we have

$$kG^{i}(;y)k_{c} = \max_{s,a}(s;a) + 1_{fS=s;A=ag} \max(+Q_{a}(S^{0};a^{0})) \quad (S;A) \quad \max Q(S^{0};a^{0}) \underset{a^{0}}{\text{max}} \quad (1$$

$$1_{fS_{s,b};A=ag})(a;s) + 1_{fS=s;A=ag} \max(+Q(S^{0};a^{0})) \quad \max_{a^{0}} Q(S^{0};a^{0}) \quad \text{i} \quad \text{(triangle inequality)}$$

$$\max_{s;a} kk1 + 1_{fS=s;A=ag} \max(+Q(S^{0};a^{0})) \quad \max_{a^{0}} Q(S^{0};a^{0}) \quad \text{i} \quad \text{(definition of } k k_{1})$$

$$\max_{s;a} kk1 + \max_{s;a} j(S^{0};a^{0})j \quad \text{(By (127))}$$

$$2 kk1$$

$$= 2 kkC:$$

Lastly, for B, we have

$$\begin{aligned} kb^i(y^i)k_c &= \max_{s;a} \mathbf{1}_{fS=s;A=ag} \prod_{h}^{h} R(S;A) + \max_{a^0} Q(S^0;a^0) \quad Q(S;A) \\ &= \max_{s;a} \mathbf{1}_{fS=s;A=ag} \quad jR(S;A)j + \max_{a^0} Q(S^0;a^0)j + jQ(S;A)j \end{aligned} \tag{triangle inequality}$$

$$\max_{s;a} \mathbf{1} + \prod_{1}^{d} \frac{1}{1} = \frac{1}{1}$$

$$= \frac{2}{1} \div$$

Proof of Lemma D.4. The proof follows similar to Lemma C.4.

Proof of Theorem 5.2. By Proposition D.1 and Lemmas D.1, D.2, D.3, and D.4, it is clear that the federated Q Algorithm 3 satisfies all the Assumptions 6.1, 6.2, 6.3, and 6.4 of the FedSAM Algorithm 4. Furthermore, by the proof of Theorem B.1, we have $w_t = (1 \frac{\binom{2}{2}}{2})^{-1}$, and the constant c in the sampling distribution q^c_T in Algorithm 1 is $c = (1 \frac{\binom{2}{2}}{2})^{-1}$. In equation (128) we evaluate the exact value of w_t .

Furthermore, by choosing step size small enough, we can satisfy the requirements in (21), (23), (35). By choosing K large enough, we can satisfy K > 1, and by choosing T large enough we can satisfy E > 1. Hence, the result of Theorem B.1 holds for this algorithm.

Next, we derive the constants involved in Theorem B.1 step by step. In this analysis we only consider the terms involving jSj, jAj, $\frac{1}{1}$, I_{max} , and I_{min} . Since k I_{c} = k I_{c} 1, we choose g() = I_{c} 2 k I_{c} 5, i.e. the p-norm with p = 2 log(jSj).

It is known that g() is (p 1) smooth with respect to k k_p norm (Beck, 2017), and hence L = (log(jSj)). Hence, we have $I_{cs} = jSj^{-1=p} = \frac{1}{p} = \frac{1}{p} = (1)$ and $I_{cs} = 1$. Therefore, we have $I_{cs} = \frac{1}{p} =$

$$a_{3} = \frac{L(1 + u_{cs}^{2})}{(2)} = O \frac{\log(jSj)(1 + cs)}{O \frac{\log(jSj)}{c}} = O \frac{\log(jSj)}{O \frac{\log(jSj)}{c}} :$$

Using '2, we have

$$w_{t} = 1 \quad \frac{2}{2} = B \quad B = 2 + \frac{0.25e^{1-4}(2 \quad \min \quad 1) \quad C}{r \quad \frac{1}{p} \quad \frac{1}{2 \quad 2\min(1)}} \quad A \quad (128)$$

Further, we have

$$I_{cm} = (1 + I_{cs}^2)^{1=2} = (1)$$

 $U_{cm} = (1 + U_{cs}^2)^{1=2} = (1)$

Since TV-divergence is upper bounded with 1, we have m= O(1). By Lemma D.3, we have

$$A_1 = A_2 = 2 = O(1)$$

and $A_1 = A_2 =$

(1),
$$\frac{2}{B} = \frac{1}{1}$$

Hence $m_1 = 2A_2 m = O(1)$. Also, we have $m_2 = 0$.

We choose the D-norm in Lemma B.9 as the 2-norm k k_2 . Hence, by primary norm equivalence, we have $I_{cD} = p_j \frac{1}{S_j}$, and $u_{cD} = 1$, and hence $\frac{u_{cD}}{I_{cD}} = p_j \frac{1}{S_j}$. The rest of the proof is similar to the proof of Theorem 5.1 where I_{max} is substituted with 1. The sample complexity can also be derived using Corollary B.1.1.