1  We thanks the reviewers for their feedback. We now proceed to reply to their comments.

2  **Number of Samples for Complete Graph? (R1)**: In this case there is no network error and all of the agents iterates
3  are identical to centralised single-machine Gradient Descent (GD) applied to all of the samples within the network,
4  as we point out in Appendix B.1. Error Decomposition. There is no condition on the number of samples to achieve
5  optimal statistical rates in this case [25].

6  **Extending Current Limitations (R2, R3)**:
7  In light of the comments we have pursed extending the theory on the following fronts:

8  • **Non-attainable case**: Consider the capacity independent ($\gamma = 1$) non-attainable case ($r < 1/2$), when
9    agents know the population covarince $\mathcal{T}_\rho$ i.e. additive noise oracle or fixed-design regression [15,20]. The
10   network error in this case only consists of the population covariance error and can be bounded as follows.
11   Due to being in the non-attainable case, the bound on the terms $\{N_{k,w_k}\}_{k,w_k}$ now depend on the time
12   step $k$, in particular with regularisation becoming $\widetilde{O}\big(\frac{1}{\sqrt{m\lambda}} + \frac{(\eta k)^{1/2-r}}{m\sqrt{\lambda}}\big)$, see Lemma 18 in [26]. The first
13   term matches the attainable case while the second is new and now must be controlled. Following the
14   analysis for the attainable case, this new term squared can be shown to be $\widetilde{O}\big((\eta t)^{1-2r}(\eta t^\star)/m^2\big)$ in high
15   probability. With $\eta = O(1)$ and $t = (nm)^{1/(2r+\gamma)}$ this quantity is smaller than the optimal statistical rate
16   once $m \geq n^{1/(4r+1)}(t^\star)^{(2r+1)/(4r+1)}$. This condition is strictly weaker than $m \geq n^{2r}(t^\star)^{2r+1}$, the condition
17   that arises in the attainable case to achieve a linear speed up by ensuring the first term is sufficiently small, if
18   $r \geq 1/4$. We leave the capacity dependent non-attainable case to future work.

19  • **Residual Covariance Error**: When agents do not know the population covarince, we require $2r + \gamma \geq 2$ to
20   control the residual covariance error. A consequence in the finite low-dimensional setting ($d < nm, \gamma = 0$)
21   is that bounds only hold for "easy" problems $r > 1$. We note a similar consequence occurs in the the
22   high-dimensional setting as source and capacity assumptions analogous to our setting can be utilised to achieve
23   dimension free bounds (when $d > nm$ the classical $O(d/nm)$ rate is vacuous), see [34]. The limitation
24   $2r + \gamma \geq 2$ can be improved though sharper control of the residual network error terms, in particular by
25   repeatedly applying the recursion (Proposition 5) and centering the difference $\Pi_{t:k+1}(\mathcal{T}_{\mathbf{x}_{w_{t:k+1}}}) - \Pi_{t:k}(\mathcal{T}_\rho)$
26   around the expectation $\mathbf{E}[\Pi_{t:k+1}(\mathcal{T}_{\mathbf{x}_{w_{t:k+1}}})]$.

27 **Communication Model (R3)**:

28  • **Description**: We consider a lockstep communication model where each round lasts for $\tau$ units of time. Within
29   each round, agents send/receive the messages to/from their neighbours in order to implement a single update
30   of algorithm (3). With a gradient evaluation costing 1 unit of time, each iteration of Distributed Gradient
31   Descent takes the following amount of time $m + \tau + \mathrm{Deg}(P)$: $m$ gradient evaluations; $\tau$ in communication
32   delay ; $\mathrm{Deg}(P)$ for each agent to aggregating their neighbours and own gradients, as the sum in algorithm (3)
33   $\sum_{w \in V} P_{vw}$ has computational cost $O(\mathrm{Deg}(P))$.

34  • **Communication Delay $\tau$ Increasing with $\mathrm{Deg}(P)$?**: Indeed $\tau$ maybe increasing with the network degree,
35   although it could depend on other factors arising from: noisy transmission, compressing or decompressing
36   messages and synchronizing with neighbours. The work Tsianos and Rabbat, NIPS 2012 makes the assumption
37   that the delay $\tau$ is a linear function of the network degree and transmit time $r \geq 0$ so $\tau = r\mathrm{Deg}(P)$. The
38   conclusion of Section 3.1 would be unaffected provided $\tau + \mathrm{Deg}(P) = O(m)$.

39  • **Tradeoff In Terms of the Parameter $\tau$**: For sufficiently large $m$ the speed up is $nm/(m + \tau + \mathrm{Deg}(P))$.
40   No speed up is achieved if the communication delay is lower bounded $\tau = \Omega(nm)$, that is if it is longer
41   than the time it takes for centralised single-machine GD to compute a gradient with all of the samples within
42   the network. More precisely, suppose a network topology with $\mathrm{Deg}(P) = \Theta(n^\beta)$ for $0 \leq \beta \leq 1$, delay
43   $\tau = r\mathrm{Deg}(P)$ as described previously and $m > n$. Then a linear speed up is achieved if $r \leq m/n^\beta$, no speed
44   up if $r \geq mn^{1-\beta}$, and a non-linear speed up if $m/n^\beta < r < mn^{1-\beta}$. In contrast, for high degree graphs with
45   $\beta = 1$ and $(1 - \sigma_2)^{-1} = O(1)$, Sianos and Rabbat, 2012 saw a linear speed up when $r < 1$ whilst we require
46   $r \leq m/n$. We thank the reviewer for this question as the above is now included within the manuscript.

47 **Communicating Infinite-Dimensional Quantities (R3)**: A theoretical setting allows precise understanding of the
48 generalisation capabilties through optimal statistical rates [12]. While particular implementations are outside the scope
49 of this work, note it is common to use finite approximations of infinite dimensional quantities whilst accounting for the
50 statistical precision (Learning with sgd and random features Luigi, Rudi & Roasasco 2018).

51 **Regime Where Nodes Work Independently (R3)** : The objective is for each agent to achieve the optimal statistical
52 rate with respect to *all* $nm$ data points within the network, and as such, it is a requirement for agents to communicate
53 one another.