

# Optimal Statistical Rates for Decentralised Non-Parametric Regression with Linear Speed-Up

Patrick Rebeschini & Dominic Richards

University of Oxford Department of Statistics

## Desiderata for Distributed Machine Learning

Suppose  $n$  agents solving a machine learning problem. Properties of an ideal distributed algorithm:

- Statistics:** retain optimal statistical precision
- Runtime:** speed-up over single-agent due to parallel computing: ideally factor  $n$
- Communication:** fixed cost per agent per step, ideally independent of  $n$

## Consensus Optimisation for Decentralised Learning

Network of agents  $G = (V, E)$ , each with data points sampled i.i.d  $(x_{i,v}, y_{i,v})$  perform linear regression. Therefore each agent wishes to minimise with coefficient  $\omega \in \mathbb{R}^d$

$$F(\omega) = \underbrace{\frac{1}{n} \sum_{v \in V} \frac{1}{m} \sum_{i \in [m]} \langle \omega, x_{i,v} \rangle - y_{i,v} \rangle^2}_{\text{Empirical Loss for Agent } v} = \underbrace{\frac{1}{n} \sum_{v \in V} F_v(\omega)}_{\text{Consensus Optimisation}}$$

where  $F_v$  is function held by agent  $v$ .

### Consensus Optimisation

Suppose  $F_v$  arbitrary and each agent wants to minimise  $\frac{1}{n} \sum_{v \in V} F_v$ .  
→ Agents alternate: local gradient descent steps on  $F_v$  and local averaging on network [4]

- ✓ Low communication cost per agent for sparse graphs
- ✓ Robust/Decentralised as no single node responsible for disseminating information
- ✗ Performance depends on network topology [2, 5].

Graphs with smaller spectral gap benefit less from decentralisation

In machine learning  $F_v$  are often not arbitrary

In our case, concentration states for  $v \in V$

$$F_v(\omega) = \frac{1}{m} \sum_{i \in [m]} \langle \omega, x_{i,v} \rangle - y_{i,v} \rangle^2 \xrightarrow{m \rightarrow \infty} \underbrace{\int_{X \times Y} (\langle \omega, x \rangle - y)^2 d\rho(x, y)}_{\text{Test Risk}} =: \mathcal{E}(\omega)$$

where  $\rho$  is the distribution of the data points.

🔑 In large data scenario, all functions are converging to the same quantity → Test Risk.

**Main Question:**  
Can concentration speed-up consensus learning for any network topology ?

## Non-parametric Statistical Assumptions

Use tools from non-parametric regression [1].

Predictor minimising Test Risk over set of linear predictors  $x \rightarrow \langle \omega, \cdot \rangle$  denoted  $f_H$ .

**Noise:** Exists  $M \in (0, \infty)$ ,  $\nu \in (1, \infty)$  so for any  $\ell \in \mathbb{N}$  we have  $\int_Y y^{2\ell} d\rho(x|y) \leq \nu \ell! M^\ell$ .

**Difficulty of estimation problem:**

For  $f \in L^2(H, \rho_X)$  let  $\mathcal{L}_\rho(f) = \int_X \langle x \cdot \cdot \rangle f(x) d\rho_X(x)$ . Exists  $r > 0$  such that  $\|\mathcal{L}_\rho^{-r} f_H\| < \infty$ .

**Spectrum of covariance operator:**

Exists  $\gamma \in (0, 1]$ ,  $c_\gamma > 0$  such that  $\text{Tr}(\mathcal{L}_\rho(\mathcal{L}_\rho + \lambda)^{-1}) \leq c_\gamma \lambda^{-\gamma}$  for all  $\lambda > 0$ .

## Distributed Gradient Descent

Consider a simple consensus optimisation algorithm [4]:

**Communication Matrix:**  $\mathbf{P} \in \mathbb{R}^{n \times n}$  symmetric doubly stochastic matrix support on the graph

$$\mathbf{P} = \mathbf{P}^\top, \quad \mathbf{P} \mathbf{1} = \mathbf{1} \quad \text{and for } v \neq w \quad \underbrace{\mathbf{P}_{vw} \neq 0 \text{ only if } (v, w) \in E}_{\text{Sparsity pattern matches network}}$$

**Network Dependence:** Let  $\sigma_2$  be second largest eigenvalue in magnitude for  $\mathbf{P}$ .  
Scaling  $O((1 - \sigma_2)^{-1}) = \mathbf{n}^2$  (Cycle),  $\mathbf{n}$  (Grid/Random Geo.), and  $\mathbf{1}$  (Complete/Expander).

**Algorithm:** Initialised  $w_{1,v} = 0$  for  $v \in V$ , iterates updated for all  $v \in V$

$$\omega_{t+1,v} = \underbrace{\sum_{w \in V} \mathbf{P}_{vw}}_{\text{Local Communication}} \left( \underbrace{\omega_{t,w} - \eta \frac{1}{m} \sum_{i=1}^m \langle \omega_{t,w}, x_{i,w} \rangle_H - y_{i,w}}_{\text{Local Gradient Descent}} \right) x_{i,w}$$

**Implicit Regularisation:**

Model Complexity controlled through early stopping (Iterations  $t$ ) and step size  $\eta$ .

## Optimal Statistical Rate with Implicit Regularisation

Let aforementioned assumptions hold with  $r \geq 1/2$  and  $2r + \gamma \geq 2$ . Fix

$$t = \underbrace{(nm)^{1/(2r+\gamma)}}_{\text{Single-Machine Iterations}} \times \begin{cases} \left( \frac{(nm)^{2r/(2r+\gamma)}}{m(1-\sigma_2)^\gamma} \right)^{1/\gamma} \vee 1 & \text{if } m \geq n^{2r/\gamma} \\ \frac{(nm)^{r/(2r+\gamma)}}{\sqrt{m(1-\sigma_2)}} & \text{otherwise} \end{cases}$$
$$\eta = \frac{\kappa^{-2}(nm)^{1/(2r+\gamma)}}{t} \text{ and let } m \geq n^{\frac{2r+2+\gamma}{2r+\gamma-2}} \text{ then } \forall v \in V:$$

$$\underbrace{\mathbb{E}[\mathcal{E}(\omega_{t+1,v})]}_{\text{Test Error}} - \inf_{\omega} \mathcal{E}(\omega) \lesssim \underbrace{(nm)^{-2r/(2r+\gamma)}}_{\text{Optimal Statistical Rate}}$$

## Time Model

Gradient computation costs 1 unit of time.

Communication delay costs  $\tau$  units of time, for some  $\tau > 0$ .

$$\text{Distributed time per iteration} = \underbrace{m}_{\text{Local Gradient Computation}} + \underbrace{\tau + \text{Deg}(P)}_{\text{Communicating/Aggregating Neighbours Information}}$$

**Speed-Up** defined as

$$\text{Speed-Up} = \frac{\text{Single Machine Run time}}{\text{Distributed Run time}} = \frac{\text{Single Machine Iterations}}{\text{Distributed Iterations}} \times \underbrace{\frac{nm}{m + \tau + \text{Deg}(P)}}_{\text{Ratio of Time Per Iteration}}$$

## Main Result: Speed-Up with Consensus Methods Utilising Concentration

Iterations required decreasing in number of samples  $m$  up to a point

$$\underbrace{m \geq \frac{n^{2r/\gamma}}{(1-\sigma_2)^{2r+\gamma}} \vee n^{\frac{2r+2+\gamma}{2r+\gamma-2}}}_{\text{Sufficiently Many Samples}} \implies \text{Distributed Iterations} = \text{Single Machine Iterations}$$
$$\implies \text{Speed-Up} \stackrel{\tau + \text{Deg}(P) = O(m)}{=} O(n)$$

therefore linear speed-up for any network topology.

	Cycle	Grid	R. Geom.	Complete	Expander
<b>Speed-Up</b>	$O(n)$	$O(n)$	$O(n)$	$O(n)$	$O(n)$
<b>Communication</b>	$O(1)$	$O(1)$	$O(1)$	$O(n)$	$O(1)$

## Speed-Up with Single-Step Consensus Methods

Single-Step methods typically require **iterations to scale with inverse spectral gap** e.g. [2, 5]

$$\text{Distributed Iterations} = \text{Single Machine Iterations} \times (1 - \sigma_2)$$
$$\Downarrow$$
$$\text{Speed-Up} = \frac{nm}{m + \tau + \text{Deg}(P)} (1 - \sigma_2)^{\tau + \text{Deg}(P) = O(m)} O(n(1 - \sigma_2))$$

therefore **linear speed-up restricted to well connected topologies**.

	Cycle	Grid	R. Geom.	Complete	Expander
<b>Speed-Up</b>	$O(1/n)$	$O(1)$	$O(1)$	$O(n)$	$O(n)$
<b>Communication</b>	$O(1)$	$O(1)$	$O(1)$	$O(n)$	$O(1)$

## Speed-Up with Multi-Step Consensus Methods

Multi-Step methods perform multiple communication steps between gradient descent steps [5].

$$\text{Distributed Iterations} = \text{Single Machine Iterations}$$

But **communication rounds scales with inverse spectral gap**.

	Cycle	Grid	R. Geom.	Complete	Expander
<b>Speed-Up</b>	$O(n)$	$O(n)$	$O(n)$	$O(n)$	$O(n)$
<b>Communication</b>	$O(n)$	$O(\sqrt{n})$	$O(\sqrt{n})$	$O(n)$	$O(1)$

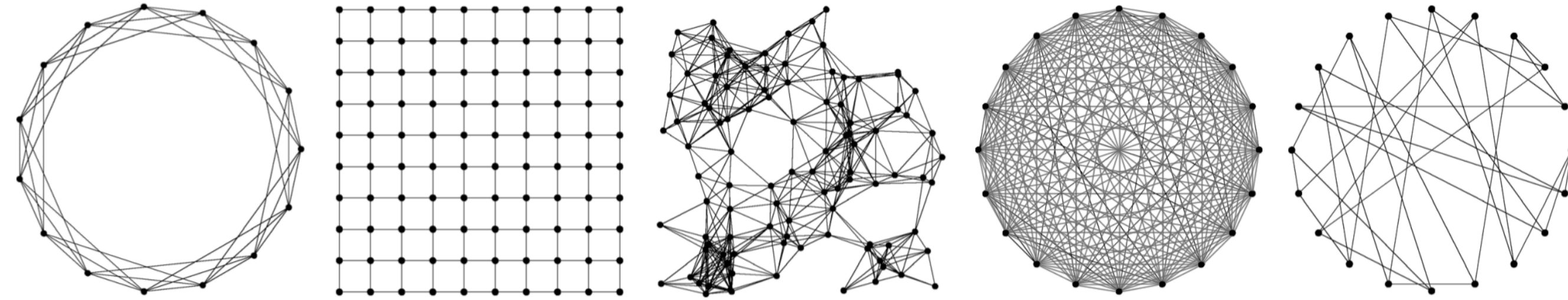
Therefore trade off between **Speed-Up** and **Communication Cost**.

EPSRC  
Engineering and Physical Sciences  
Research Council

OxWaSP  
Oxford-Warwick Statistics Programme

UNIVERSITY OF  
OXFORD

MRC  
Medical  
Research  
Council



## Detailed Error Decomposition

The Test Error is decomposed as follows when  $m \geq n^{2r/\gamma}$

$$\text{Test Error} \lesssim \underbrace{(\eta t)^{-2r}}_{\text{Bias}} + \underbrace{\frac{(\eta t / (nm)^{1/(2r+\gamma)})^2}{(nm)^{2r/(2r+\gamma)}}}_{\text{Sample Variance}} + \underbrace{\frac{\eta^\gamma}{m(1-\sigma_2)^\gamma}}_{\text{Population Network Error}} + \underbrace{\frac{(\eta t)^{\gamma+2}}{m^2}}_{\text{Residual Network Error}}$$

**Bias and Sample Variance** align with Gradient Descent with  $nm$  samples [3]

$$\implies \text{Fix } \eta t = (nm)^{1/(2r+\gamma)}$$

**Population and Residual Network Error:** arise due from averaging steps with the matrix  $\mathbf{P}$ .

### Population Network Error

Follows standard network term [2]

Decreasing with step size  $\eta$

Depends on spectral gap  $(1 - \sigma_2)$

Due to **concentration, decreasing with  $m$**

$$\text{Small by picking } \eta = \left( \frac{m(1-\sigma_2)^\gamma}{(nm)^{2r/(2r+\gamma)}} \right)^{1/\gamma} \vee 1.$$

### Residual Network Error

Higher order term

From empirical covariance multiplying the iterates at each iteration

Utilise **concentration** and **contraction** to control

$$\implies \text{require } m \geq n^{\frac{2r+2+\gamma}{2r+\gamma-2}}.$$

## Proof Sketch: Population Network Error

Let  $\mathcal{T}_\rho$  be conjugate of  $\mathcal{L}_\rho$  and for  $k \geq 1$ , let  $\mathbf{N}_k$  be r.v. with zero mean concentrating to zero such that  $\|(\mathcal{T}_\rho + \lambda \mathbf{I})^{-1/2} \mathbf{N}_k\| \lesssim 1/\sqrt{\lambda \gamma m}$  w.h.p. Then, with mixing time  $t^* \simeq (1 - \sigma_2)^{-1}$

$$\mathbb{E}[(\text{Pop. Net. Error})] \leq \mathbb{E} \left[ \left( \sum_{k=1}^t \sigma_2^{t-k+1} \eta \|\mathcal{T}_\rho^{1/2} (I - \eta \mathcal{T}_\rho)^{t-k} \mathbf{N}_k\| \right)^2 \right]$$
$$\lesssim \frac{1}{m \lambda \gamma} \left( \sum_{k=t-t^*}^t \sigma_2^{t-k+1} \eta \underbrace{\|\mathcal{T}_\rho^{1/2} (I - \eta \mathcal{T}_\rho)^{t-k} (\mathcal{T}_\rho + \lambda \mathbf{I})^{1/2}\|}_{\text{Contraction} \implies O(\sqrt{\lambda}/(\eta(t-k)))} \right)^2 \lesssim \frac{\log(t^*) + \lambda \eta t^*}{m \lambda \gamma}$$

$t^*$  Poorly Mixed terms

Optimised by picking  $\lambda = (\eta t^*)^{-1}$ .

## Future Work

- **Non-parametric setting:** extending analysis to non-attainable case  $r \leq 1/2$  and tighter analysis on **Residual Network Error**.
- **General Loss Function:** Squared loss yields bias/variance decomposition, **concentration** likely hold for more general losses.
- **Statistics/Communication trade off with sparse/randomised gossip:** linear speed-up independent of topology → agents randomly gossip at each iteration.
- **Stochastic Gradient Descent and mini-batches:** random subset of data at each iteration [3].

## References

- [1] Andrea Caponnetto and Ernesto De Vito. Optimal rates for the regularized least-squares algorithm. *Foundations of Computational Mathematics*, 7(3):331–368, 2007.
- [2] John C. Duchi, Alekh Agarwal, and Martin J. Wainwright. Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic Control*, 57(3):592–606, 2012.
- [3] Junhong Lin and Lorenzo Rosasco. Optimal rates for multi-pass stochastic gradient methods. *Journal of Machine Learning Research*, 18(97):1–47, 2017.
- [4] Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
- [5] Kevin Scaman, Francis Bach, Sébastien Bubeck, Yin Tat Lee, and Laurent Massoulié. Optimal algorithms for smooth and strongly convex distributed optimization in networks. In *Proceedings of the 34th ICML*, pages 3027–3036. JMLR. org, 2017.