

Abstract

With the size of genetic datasets growing studying datasets fragmented across multiple institutions is emerging as a powerful paradigm for the future of gene mapping and discovery. In this work, we investigate decentralized optimization algorithms for robust statistical inference on datasets fragmented across multiple data owners. In particular, we consider the logistic and linear regression problems and with extensive simulation show that the inference can be done across multiple data owners in a practical way. Finally, we simulated a realistic GWAS by splitting WTCCC **fill in** disease dataset across 5 data owners and compare the decentralized algorithms with standard meta-analysis algorithms and the combined analysis approach.

Keywords: words, words, words

11 Introduction

12 **recent datasets are large** In recent years, the expansion in the available genomic and pheno-
13 typic data, brought about by the next-gen biomolecular and biometric technologies, has lead to
14 breakthroughs in our understanding of human disease. The incorporation of more personal level
15 data, as well as the push for studying lower effect sizes, and rare variants increases the need for
16 a larger cohort in order to achieve the necessary power. This expansion of available data presents
17 new challenges in storage, security and computation.

18 **Shortcomings of current approaches** The old approach (Table 1 (a)) of generating, gathering
19 and maintaining these multi-factorial and private datasets in a central location can be expensive,
20 unsafe and time consuming. Sharing and combining summary statistics through meta-analysis
21 techniques can sometimes offer a work-around for these issues. However, these methods have a
22 few limitations: (a) subtle differences in models, assumptions and QC can introduce biases in the
23 results (17), (b) the shared data and summary statistics might be inadequate for some types of
24 inference (e.g. estimating effects from shared two-tailed p-values), and (c) parameter estimates can
25 be unreliable when the effective sample size (N) is small compared to the number of covariates (k)
26 as could happen for variants with low allele frequency. (d) smaller datasets may be ignored as the
27 effort to incorporate them may outweigh their benefits.

28 **Advances in machine learning and algorithms used** With the big data revolution and explo-
29 sion of accessible data in many fields, new parallel or decentralized algorithms become increasingly
30 available (7; 20; 26; 16; 9; 10; 22). In this work we are concerned with exploring inference algorithms
31 that allow the data from the participants to remain fragmented across distant data owners (DOs).
32 These DOs can communicate with the central hub to iteratively solve the problem but they may
33 not transfer the raw data (1 (b)). As opposed to the more common case where the goal of parallel
34 computing is to speed up a computation performed on nodes on the same chip or on a cluster of
35 computers in close proximity. Therefore, we assume that substantial latency might effect the com-
36 munications within nodes. As such, we use the distinction that while decentralized algorithms can
37 be used in parallel, not every parallel algorithm is appropriate for a decentralized setting.

38 **Summary of what's to come** Our main goal for what follows is to demonstrate that simple
39 GWAS studies can be performed under this fragmented data paradigm. Specifically, we limit our

discussion to linear and logistic regression and generating principle component values. In what follows, we first review the regression framework used and present the optimization problem that needs to be solved. This is followed by a short discussion of the optimization algorithms and the algorithm for extracting the PCA. In the results section we present extensive simulations for evaluating the trade-offs between algorithms and we perform a GWAS study on WTCCC **fill** dataset split between 5 cohorts.

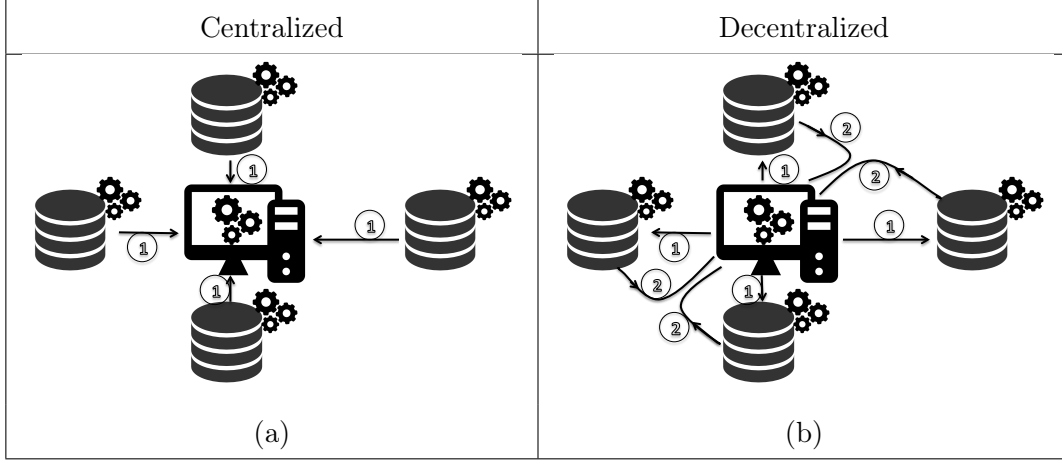


Table 1: (a)The simplest method of data-sharing would involve all DOs sending their data to a single location. That location would then execute the relevant machine learning algorithms. (b) In our model of decentralized scheme, one central hub can communicate with all the DO's without directly requiring the data to be sent to the central hub. The central hub coordinates the execution of the relevant machine learning algorithm and through a potentially iterative process, a single set of parameters is estimated for the machine learning algorithm.

Methods

Generalized Linear Models

In a broad sense, all regression based classifiers are consisted of three parts: representation, evaluation and optimization (12). The **representation** generally describes the relationship between parameters and limits the class of models that the learner can represent. The **evaluation function** (also known as cost function) is a tool for evaluating different parameter settings of the model. Finally, the **optimization technique** is the algorithm used to find the optimum model-parameters with respect to the chosen cost function. Both logistic and linear regression are part of the larger

Type of GLM	Mean function $E[\mathbb{Y}]$	Semantics
Linear regression	$\mu = X\beta$	Outcome is a linear function of the parameters with range $(-\infty, \infty)$
Logistic regression	$\mu = \frac{\exp(X\beta)}{1+\exp(X\beta)}$	Probability of a categorical outcome based on a binomial distribution
Poisson regression	$\mu = \exp X\beta$	Modeling Poisson distributed $\frac{\text{counts}}{\text{time interval}}$
Multinomial regression	$\mu = \frac{\exp(X\beta)}{1+\exp(X\beta)}$	Probability of each outcome for a K-sided loaded die

Table 2: List of most commonly used link functions, the GLM regression they correspond to and their common use case.

54 generalized linear models (GLMs) family.

Representation: GLMs treat each response variable as a random draw from a distribution in the exponential family where the mean of the distribution can be linked to a linear function of the covariates (24). More concretely, a GLM models the expected value of the outcome variable (\mathbb{Y}) as:

$$E(\mathbb{Y}) = \mu = g^{-1}(X\beta) \quad [1]$$

55 Where X is a matrix with covariates for each individual as its rows (We will use roman indices's to
56 indicate an individual specific response and covariates.). β is the vector of unknown parameters and
57 g is a link function that is known a priori. In this setting, the link function (g) provides flexibility
58 to model various types of outcomes. Table 2, provides a short list of examples of common link
59 functions and their typical use cases. (see (1) for a comprehensive review)

60 **GLM evaluation function:** The unknown parameters (β) are chosen to optimize a particular
61 cost function. In a Bayesian setting, optimizing the posterior probability is the natural cost function.
62 Another approach is to maximize the likelihood function. This approach offers estimators with
63 theoretically attractive properties (28). For GLM's the assumption that each outcome is independent
64 given the covariates leads to log-likelihood that is the sum of per-sample log likelihoods

$$\ell = \log(P(\mathbb{Y} | \beta, X)) = \sum_{i=1}^N \log(P(Y_i | \beta, X_i)) = \sum_{i=1}^N \ell_i \quad [2]$$

65 a common modification to this cost function is adding penalization on β (see for a few examples(29;
66 32; 15)). In particular, after the addition of the penalty term, the cost function remains the sum

of individual-level costs. This (penalized) maximum likelihood approach is the approach chosen by many statistical packages (used by `glmnet` (14) in R and `scikit-learn` (25) in python).

GLM optimization techniques vary depending on the details of the regression and will constitute our main focus. One common optimization approach is the iteratively reweighted least square (IRLS) algorithm, which iteratively solve a carefully constructed weighted least square problems (see (1) 4.6 or (23) 8.3 for details). While, the default for some common packages (`glm2` in R for example (21) for example) uses IRLS, this algorithm is not practical if the data is very high dimensional or, otherwise, cannot be gathered in one location. As a result sub-gradient based methods have been adopted by many packages (such as `glmnet` in R, (14), Vowpal Wabbit (18) and Python’s `scikit-learn` (25)) to grapple with today’s high dimensional data.

Optimization algorithms

For this work we benchmark the performance of 7 optimization algorithms in a decentralized setting. In a broad sense these algorithms fall under two categories: gradient based, and DO based. For gradient based algorithms (first 5 algorithms below) the parameters of interest are iteratively updated based on the gradient computed on a batch of data. For DO based algorithms (AVG, and ADMM below) the entire data at each DO is used and the hope is to achieve better results through consensus or averaging the results from different DO’s. Below, we provide a brief explanation of each algorithms, a more detailed explanation of these algorithms as well as the relevant parameters are available in the Supplementary Materials.

SGD is a simple stochastic gradient algorithm. The step-size is chosen to be constant for the first 500 iterations and will slowly decays after.

ADAGRAD uses the algorithm presented in (13). This algorithm uses previous gradient information to modify the step-size for each covariate. As a result, it is more robust to the choice of step-size while maintaining a per-covariate step-size.

RMSPROP uses the algorithm presented in (30). This algorithm is similar to ADAGRAD but puts the emphasis on the recent gradients in order to prevent a heavy reduction in step-size for later iterations.

ADADELTA uses the algorithm presented in (31). This algorithm also uses previous gradients and updates to produce a covariate-specific step-size.

SQN uses the algorithm presented in (8). In this algorithm every few iterations a noisy Hessian is estimated using a relatively large batch size. This Hessian is then used to compute the per-covariate

step-size for each stochastic gradient descent step. **nmight need to cite more stuff**

AVG computes the model parameters at each DO and uses a weighted average to compute the overall model parameters. Unlike the other algorithms, this algorithm only has one step and does not have any parameters.

ADMM implements the ADMM algorithm (6). This algorithm generalizes the AVG algorithm by iteratively updating the parameter estimates.

The stochastic gradient for the first 5 algorithms is often computed on a small batch of data. In most common settings the faster speed of computing a noisy gradient on a smaller batch-size overcomes the benefits of a less noisy estimate on a big batch-size. In our case, the network latency trumps the gradient computation time, therefore generally we use large batch sizes.

To perform the first 5 algorithms in a decentralized fashion, we follow the approach suggested in (9; 10) and compute a mini-batch gradient at each node and simply average the gradient at the central node once results from all DO's have been reported. This provides us with a gradient estimate on a batch-size of $\sum_{i=1}^{\# \text{ DO's}} n^{(i)}$ where $n^{(i)}$ is the batch-size from the i^{th} DO. Unfortunately, unlike (9; 10), we cannot take batches of equal sizes from each DO or ignore the DO's with a late response time as the data from all DO's might not be identically distributed.

ABOVE HAS BEEN UPDATED + TWO NEW FIGURES BE-
LOW

Results

In this section, we will demonstrate the performance of linear regression and logistic regression optimized using a selected set of algorithms from the list reviewed in the previous section. We will generate synthetic data as explained in the next section to evaluate the performance of each algorithm for varying number of covariates, data sizes, and covariance structures.

DataSets

To benchmark the algorithms presented earlier, we will simulate dataset as follows:

Uncorrelated continuous covariates are drawn from a $Norm(0, 10^\alpha)$ with α drawn uniformly from the interval $(0, 3)$. The discrete covariates are limited to take dosage values $(\{0, 1, 2\})$. The dosage values are assigned by drawing a frequency parameter for each covariate ($f \sim Unif(0.01, 0.49)$) and assigning a dosage of 0, 1, 2 with probabilities $f^2, 2f(1 - f)$, and $(1 - f)^2$ for each individual. Let X_i

denote the row vector of the covariates for individual i , and β denote the vector of coefficients (drawn from a $Norm(0, 1)$ distribution). Using these definitions, the label for individual i is computed as $y_i = X_i\beta + b + \epsilon_i$ for linear regression and $y_i = \text{binom}(1, X_i\beta + b + \epsilon_i)$ for logistic regression. In both cases, ϵ_i is an individual-specific Gaussian error term with $Norm(0, 0.1)$ distribution and the bias term (b) is set to zero ($b = 0$). The data is distributed across 5 DO's. The first four DO receive a random number of individuals (distributed according to $Pois(\frac{\text{size}}{\#DO's})$) and the last DO received the remaining individuals. The simulation is restarted if any DO receives less than 10% of the data.

We study the effects of correlation for continuous variables only. The covariates are drawn from a multivariate normal random variable with a random covariance matrix generated via vine method detained in (19). The remaining values are assigned as before.

Algorithms

In order to benchmark the general approaches outlined in Methods, we implemented benchmarked 7 representative algorithms. Since we anticipate the largest cost to be the communication between the DO's and the central hub, we will evaluate these algorithms based on the number of communications rather than the total compute time. Below is a short list of the implemented algorithms:

SGD is the simple stochastic gradient algorithm introduced earlier. The stepsize is chosen to be constant for the first 500 iterations and will decays as $\eta_t = \frac{\eta}{(t-500+1)^{0.51}}$ for the remaining steps.

ADAGRAD uses the algorithm presented in (13). This algorithm uses previous gradient information to estimate the current stepsize for each covariate. As a result, it is more robust to stepsize choices and can have a per-covariate stepsize.

RMSPROP uses the algorithm presented in (). This algorithm is similar to ADAGRAD but puts a heavier emphasis on the recent gradients in hopes of recucing the heavy reduction step-size for larger iterations.

ADADELTA uses the algorithm presented in (31). This algorithm also uses previous gradients and updates to produce a covariate-specific step-size.

SQN uses the algorithm presented in () (byrd). In this algorithm every few iterations a noisy Hessian is estimated using a relatively large batch size. This Hessian is then used to compute the per-covariate step-size for each stochastic gradient descent step.

AVG computes the model parameters at each DO and uses a weighted average to compute the overall model parameter. As a result, this algorithm only has one iteration.

ADMM implements the ADMM algorithm (6). We will use an ℓ_1 normalized logistic regression and LASSO () regression with very small penalties in lieu of logistic regression and ordinary least squares.

More information about the parameters, stepsizes, and implementations can be found in **reference link**.

Performance

In order to benchmark these algorithms, we restricted the number of communications to 1,000 and measured the performance for different number of covariates drawn from purely gaussian, equal mix of gaussian and dosage and purely dosage distributions using 5 DO's. For each experiment, the data was gathered to compute the gold-standard solution. Figure 3.4 shows the per-sample logistic cost function difference for these algorithms and the combined logistic regression for 2,3,6,11 covariates (including the intercept value). Figure ?? shows a similar result for least squares regression and per-sample mean-squared error.

We note that ADMM and SQN consistently outperform other algorithms. Simple averaging performs well, when the sample size is large compared to the number of covariates, however, when this is not the case, it can produce very inaccurate estimates and therefore had to be removed from the top row of figure 3.4.

PUT STUFF RELATED TO LEAST SQUARES HERE

Performance as a function of number of DO's

Once again, we see that when the per-DO sample size is large compared to the number of covariates, simply averaging the per-DO estimates is a very accurate estimator, however this estimator fails when the per-DO estimates deteriorate. On the other hand SQN and ADMM are accurate throughout the parameter space explored.

Performance as a function of number of communications

In the previous experiments, we have limited all the algorithms to 1000 iterations.

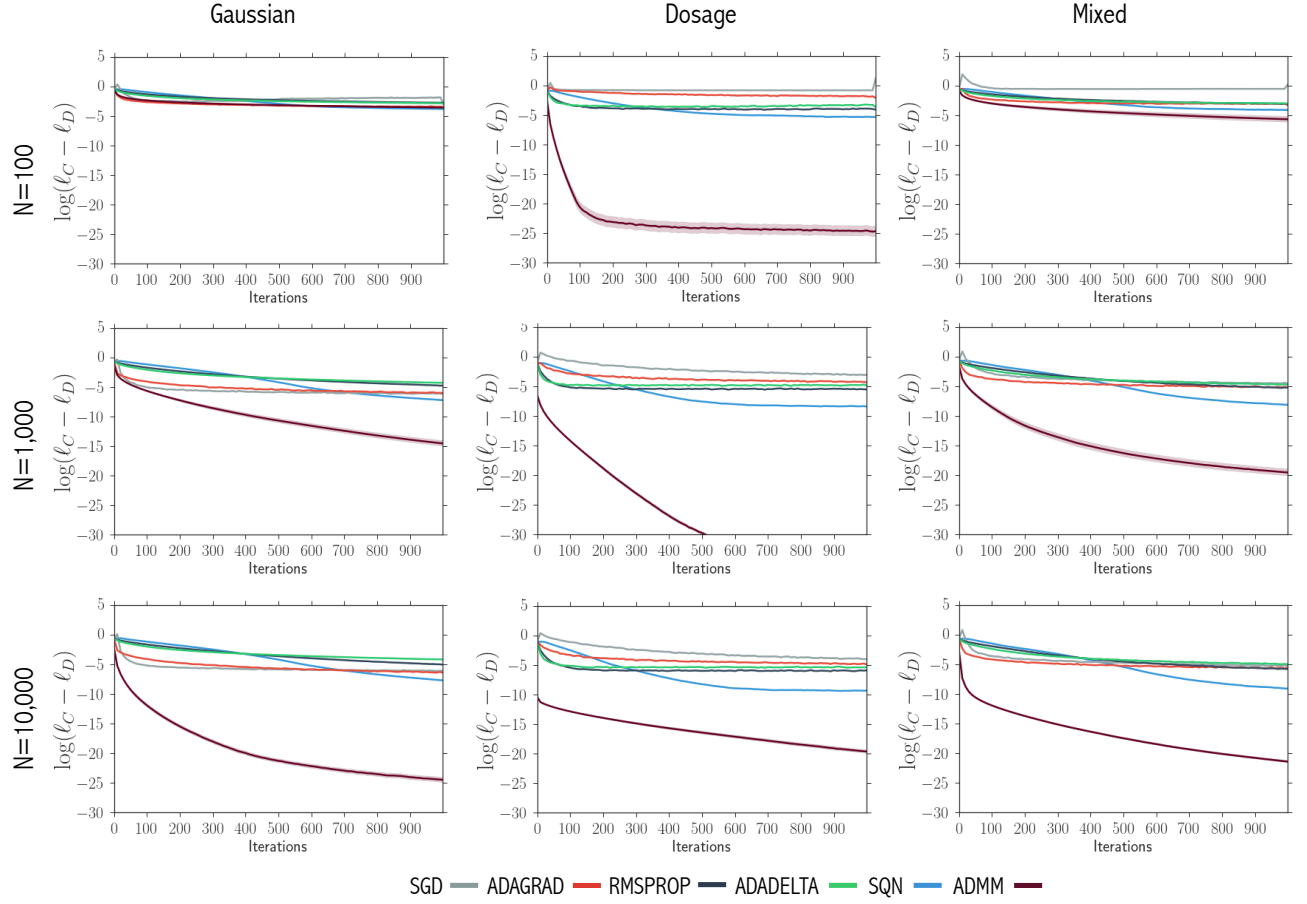


Figure 1: Plots of convergence of each method to the solution from a centralized logistic regression for 9 covariates (+1 intercept). In all cases, ADMM approaches the solution from centralized logistic regression faster than the other methods. Furthermore, at 1,000 iterations, the ADMM solution outperforms all other methods. As expected, ADMM converges much faster when the amount of data in each DO is large.

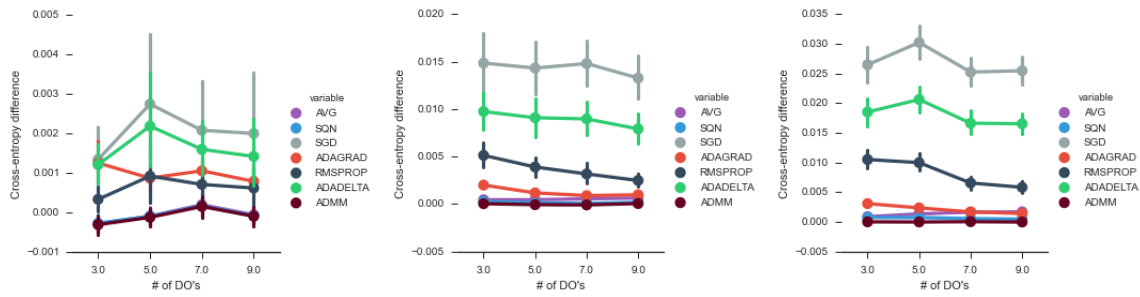


Figure 2: Insert caption w/ batch size and ...

Discussion

Conclusions

Figures and Tables

Figure 3: **This is the 1-line title of the figure.**

This is a longer text for detailed explanation. Should be longer than 1-line.

Table 3: **This is the 1-line title of the table.**

This is a longer text for detailed explanation. Should be longer than 1-line.

References

- 1 Agresti A, Kateri M. 2011. Categorical data analysis. In: International encyclopedia of statistical science. Springer. p. 206–208.
- 2 Bottou L. 1998. Online learning and stochastic approximations. On-line learning in neural networks 17:142.
- 3 Bottou L. 2012. Stochastic gradient descent tricks. In: Neural networks: Tricks of the trade. Springer. p. 421–436.
- 4 Bottou L, Curtis FE, Nocedal J. 2016. Optimization methods for large-scale machine learning. arXiv preprint arXiv:1606.04838 .
- 5 Bousquet O, Bottou L. 2008. The tradeoffs of large scale learning. In: Advances in neural information processing systems. p. 161–168.
- 6 Boyd S. 2011. Alternating direction method of multipliers. In: Talk at NIPS Workshop on Optimization and Machine Learning.
- 7 Boyd S, Parikh N, Chu E, Peleato B, Eckstein J. 2011. Distributed optimization and statistical learning via the alternating direction method of multipliers. Foundations and Trends® in Machine Learning 3:1–122.

-
- 201 8 Byrd RH, Hansen SL, Nocedal J, Singer Y. 2016. A stochastic quasi-newton method for large-
202 scale optimization. *SIAM Journal on Optimization* 26:1008–1031.
- 203 9 Chen J, Monga R, Bengio S, Jozefowicz R. 2016. Revisiting distributed synchronous sgd. *arXiv*
204 preprint [arXiv:1604.00981](https://arxiv.org/abs/1604.00981) .
- 205 10 Dekel O, Gilad-Bachrach R, Shamir O, Xiao L. 2012. Optimal distributed online prediction
206 using mini-batches. *Journal of Machine Learning Research* 13:165–202.
- 207 11 Dennis Jr JE, Schnabel RB. 1996. Numerical methods for unconstrained optimization and
208 nonlinear equations. SIAM.
- 209 12 Domingos P. 2012. A few useful things to know about machine learning. *Communications of*
210 *the ACM* 55:78–87.
- 211 13 Duchi J, Hazan E, Singer Y. 2011. Adaptive subgradient methods for online learning and
212 stochastic optimization. *Journal of Machine Learning Research* 12:2121–2159.
- 213 14 Friedman J, Hastie T, Tibshirani R. 2010. Regularization paths for generalized linear models
214 via coordinate descent. *Journal of statistical software* 33:1.
- 215 15 Friedman JH. 2012. Fast sparse regression and classification. *International Journal of Forecasting*
216 28:722–738.
- 217 16 Gopal S, Yang Y. 2013. Distributed training of large-scale logistic models. In: *ICML (2)*. p.
218 289–297.
- 219 17 Greco T, Zangrillo A, Biondi-Zoccai G, Landoni G. 2013. Meta-analysis: pitfalls and hints.
220 *Heart, lung and vessels* 5:219.
- 221 18 Langford J, Li L, Strehl A. 2011. Vowpal wabbit. URL [https://github.](https://github.com/JohnLangford/vowpal-wabbit/wiki)
222 [com/JohnLangford/vowpal wabbit/wiki](https://github.com/JohnLangford/vowpal-wabbit/wiki) .
- 223 19 Lewandowski D, Kurowicka D, Joe H. 2009. Generating random correlation matrices based on
224 vines and extended onion method. *Journal of multivariate analysis* 100:1989–2001.
- 225 20 Lubell-Doughtie P, Sondag J. 2013. Practical distributed classification using the alternating
226 direction method of multipliers algorithm. In: *Big Data, 2013 IEEE International Conference*
227 *on. IEEE*. p. 773–776.

-
- 21 Marschner IC, et al. 2011. glm2: fitting generalized linear models with convergence problems.
The R journal 3:12–15.
- 22 Mateos G, Bazerque JA, Giannakis GB. 2010. Distributed sparse linear regression. IEEE
Transactions on Signal Processing 58:5262–5276.
- 23 Murphy KP. 2012. Machine learning: a probabilistic perspective. MIT press.
- 24 Nelder JA, Baker RJ. 1972. Generalized linear models. Encyclopedia of statistical sciences .
- 25 Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Pretten-
hofer P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M,
Duchesnay E. 2011. Scikit-learn: Machine learning in Python. Journal of Machine Learning
Research 12:2825–2830.
- 26 Recht B, Re C, Wright S, Niu F. 2011. Hogwild: A lock-free approach to parallelizing stochastic
gradient descent. In: Advances in neural information processing systems. p. 693–701.
- 27 Riedmiller M, Braun H. 1993. A direct adaptive method for faster backpropagation learning:
The rprop algorithm. In: Neural Networks, 1993., IEEE International Conference on. IEEE. p.
586–591.
- 28 Scholz F. 1985. Maximum likelihood estimation. Encyclopedia of statistical sciences .
- 29 Tibshirani R. 1996. Regression shrinkage and selection via the lasso. Journal of the Royal
Statistical Society. Series B (Methodological) :267–288.
- 30 Tieleman T, Hinton G. 2012. Lecture 6.5-rmsprop: Divide the gradient by a running average of
its recent magnitude. COURSE: Neural networks for machine learning 4:26–31.
- 31 Zeiler MD. 2012. Adadelta: an adaptive learning rate method. arXiv preprint arXiv:1212.5701 .
- 32 Zou H, Hastie T. 2005. Regularization and variable selection via the elastic net. Journal of the
Royal Statistical Society: Series B (Statistical Methodology) 67:301–320.
- =====NEW STUFF===== **make sure the bottou 2016**
source is the correct one

Supplementary Materials

In this section we provide a brief explanation of the algorithms used along with the specific parameters chosen for our experiments. A comprehensive review of these algorithms is beyond the scope of this work.

Stochastic Gradient Descent

Given a function $f(\mathbf{x})$ the Gradient Descent (GD) algorithm iteratively minimizes the function by computing the gradient of function f and updating \mathbf{x} by taking a small step in the direction opposite to the average gradient. Under sufficient regularity assumptions and when the step-size is small enough and the initial point is close enough to a minima, this algorithms converges at a linear rate (11). As discussed in the main text, for independent samples, the (penalized) MLE-derived cost function for GLM family can be separated into the sum of cost functions over all the individuals (for the most usual penalizations). Denote the likelihood by $\ell(\beta; X, \mathbb{Y})$ and the penalty function with $P(\beta)$ and let $C_i(\beta; X, Y) = \ell_i(\beta; X_i, \mathbb{Y}_i) + P(\beta)$ be the overall cost function. Then we can write:

$$\ell(\beta; X, \mathbb{Y}) = \sum_{i=1}^N C_i(\beta; X, Y) \quad [3]$$

And therefore the update becomes:

$$\beta_{t+1} = \beta_t - \alpha_t \nabla_{\beta} C(\beta; X, Y) = \beta_t - \alpha_t \sum_{i=1}^N \nabla_{\beta} C_i(\beta; X, Y) \quad [4]$$

Where all functions are evaluated at β_t , $\nabla_{\beta} f$ indicates the gradient of the function $f(\beta)$ with respect to β and α_t represents the step size at time t .

Unfortunately, for large datasets, this computation may be slow since the gradient includes a sum over the entire dataset. (Mini-batch) stochastic gradient descent overcomes this problem by computing the average gradient and updating the parameter based on a small batch of data, or even a single sample, for each iteration. For a convex function, under relatively mild assumptions, this algorithms almost surely converges to the minimum for step sizes satisfying (2; 23):

$$\begin{aligned} \sum_{t=0}^{\infty} \alpha_t &= \infty \\ \sum_{t=0}^{\infty} \alpha_t^2 &< \infty \end{aligned} \quad [5]$$

Intuitively, the first requirement insures that the even a poor choice of initial parameters for β will be able to reach the optimal parameters. The second requirement insures that the step size becomes smaller as we reach the optimal parameter so that the estimate does not wobble around the optimal. While the added noise, due to the stochastic approximation of the gradient, slows down the convergence rate of the algorithms, for a large dataset, SGD can outperform GD in the overall amount of computations needed to reach a given precision (5; 4).

In our implementation we use relatively large batch sizes (**TODO**) in order to reduce the total number of updates by reducing the noise in each update. The optimal starting step size was chosen via cross-validation. The step sizes are identical for all covariates and were set to be constant for the first 500 iterations and fall off as: $\alpha_t = \frac{c_1}{(c_2+t-500)^\gamma}$ for $0.5 < \gamma \leq 1$ (general form recommended by (23)) with $C_1 = \alpha_0$, $C_2 = 1$ and $\gamma = 0.51$ for the remaining iterations.

For a comprehensive review of SGD we refer the readers to (23) for a quick introduction and to (4; 3) for a more theoretical approach.

To motivate the next 3 algorithms, we note that equation 4 resembles a Newton's update if the scalar step-size α_t is replaced by the inverse of the Hessian matrix (H^{-1}).

$$\beta_{t+1} = \beta_t - H^{-1} \nabla_{\beta} C(\beta; X, Y) \quad [6]$$

For a quadratic function C Newton's method solve the optimization problem in a single step. In more general cases, Newton's method can be viewed as approximating C as a quadratic function and iteratively updating our estimate to represent the minimum of this approximated quadratic function, therefore, Newton's method is known as a second order method. In contrast, simple SGD updates of equation 4 only consider the slope and therefore are known as first order methods. While it has been shown that stochastic method's cannot reach a convergence rate faster than sub-linear, incorporating approximate second-order information in a mini-batch setting can have improve performance. One way to incorporate this information is using a diagonal approximation to the Hessian. Such an approximation in effect scales the step size for each search direction.

ADAGRAD The AdaGrad algorithm (with diagonal matrices) is one popular approach for rescaling the step size using previous gradient values (algorithm 1 in (13)). In doing so, AdaGrad aims to increase the learning rate of infrequent features with the intuition that the learner should pay attention to infrequent features when they do occur (13). This makes AdaGrad particularly suited for dealing with sparse data. While equation 5 provides general guidelines for choosing a step size,

the convergence speed of naive gradient descent algorithm is very dependent on the exact choice of step size (see **table number**) on the other hand the AdaGrad algorithm is very robust to the initial choice of step size. For a wide range of initial choices, the algorithm adjusts the step size parameter for each covariate such that the final results is very robust to the initial choice of step size (see **reference table of step sizes**). In a decentralized setting, cross validation can be particularly time-consuming, this makes algorithms that are robust to hyper-parameters specially attractive.

The AdaGrad algorithm is presented in (13), below, we reproduce the version we have implemented.

Algorithm 1: AdaGrad Algorithm (original algorithm in (13))

Data: $\alpha > 0, T$

Result: Updated parameter β_{T+1}

initialization $G_{0,k} \beta_k \leftarrow 0$ for k in $1 \dots p$

for $t = 0$ **to** T **do**

for $i = 0$ **to** N **do**

 Receive gradient $(g_t^{(i)} \leftarrow \nabla_{\beta} C(\beta)$ evaluated at β_t) from DO i

 Receive number of individuals the gradient was computed on $(n_t^{(i)})$ from DO i

end

$n_t \leftarrow \sum_{i=1}^N n_t^{(i)}$

$g_t \leftarrow \sum_{i=1}^N \frac{n_t^{(i)}}{n_t} g_t^{(i)}$

$G_{t+1,k} \leftarrow G_t + (g_{t,k})^2$

$\beta_{t+1,k} \leftarrow \beta_t - (\frac{\alpha}{\sqrt{G_{t,k}}}) g_{t,k}$

end

RMSPROP This algorithm also provides a robust method for adjusting the step size for each covariate (see **tableput table**). One potential problem with the AdaGrad algorithm is that at each step the gradient is divided by the square-root of the sum of square of all previous gradients ($\sqrt{G_{t,k}}$ in Algorithm 1). As a result, if at any point a large gradient appears in one dimension, the learning rate in that dimension will suffer for the remainder of training. Furthermore, the accumulation over time will result in smaller and smaller step sizes as the number of iterations increases. RMSProp alleviates this problem by putting the emphasis on the most recent gradients. To do so, the algorithm maintains a running average of the squared gradients, the step size is then modified by dividing the initial value by the square-root of this running average. The original algorithm can be found in

307 (30), below, we reproduce the version we implemented.

Algorithm 2: RMSProp Algorithm (original algorithm in (30))

Data: $\alpha > 0$, T , $0 < \rho < 1$

Result: Updated parameter β_{T+1}

initialization $G_{0,k} \beta_k \leftarrow 0$ for k in $1 \dots p$

for $t = 0$ *to* T **do**

for $i = 0$ *to* N **do**

 Receive gradient $(g_t^{(i)} \leftarrow \nabla_{\beta} C(\beta)$ evaluated at β_t) from DO i

 Receive number of individuals the gradient was computed on $(n_t^{(i)})$ from DO i

end

$n_t \leftarrow \sum_{i=1}^N n_t^{(i)}$

$g_t \leftarrow \sum_{i=1}^N \frac{n_t^{(i)}}{n_t} g_t^{(i)}$

$G_{t+1,k} \leftarrow \rho G_{t,k} + (1 - \rho) (g_{t,k})^2$

$\beta_{t+1,k} \leftarrow \beta_t - (\frac{\alpha}{\sqrt{G_{t,k}}}) g_{t,k}$

end

309 In our implementation, we used $\rho = \text{TODO look this up}$

310 Since in this decentralized setting we expect to use relatively large batch sizes, rProp (see (27)
311 for details) is another potentially good update algorithm.

312 **ADADELTA** ADADELTA is another attempt at dealing with the ever decreasing step-sizes
313 in AdaGrad. As with RMSProp, ADADELTA puts more emphasis on recent gradients and uses
314 the square-root of an exponentially decaying average of squared gradients to adjust the step size.
315 Furthermore, ADADELTA aims for updates to x (i.e. Δx in 4) to have the same units as x . This
316 is the case when the Hessian matrix is used to generate the update (see 6). The authors choose the
317 square-root of an exponentially decaying average of the square of previous updates as the quantity
318 to add to the numerator. The full algorithm and an explanation for this choice can be found in (31).
319 It must be noted that in this method, the denominator contains the most recent gradient computed
320 but the numerator lags one time-step behind as the new update has not been determined yet. This

321 tames the step size for batches where a large gradient is computed.

Algorithm 3: ADADELTA Algorithm (original algorithm in (31))

Data: $\alpha > 0, T, 0 < \rho < 1, \epsilon$

Result: Updated parameter β_{T+1}

initialization $D_{0,k} \leftarrow \epsilon, G_{0,k} \beta_k \leftarrow 0$ for k in $1 \dots p$

for $t = 1$ **to** T **do**

for $i = 0$ **to** N **do**

 Receive gradient $(g_t^{(i)} \leftarrow \nabla_{\beta} C(\beta)$ evaluated at β_t) from DO i

 Receive number of individuals the gradient was computed on $(n_t^{(i)})$ from DO i

end

$n_t \leftarrow \sum_{i=1}^N n_t^{(i)}$;

$g_t \leftarrow \sum_{i=1}^N \frac{n_t^{(i)}}{n_t} g_t^{(i)}$

$G_{t+1,k} \leftarrow \rho G_{t,k} + (1 - \rho) (g_{t,k})^2$

$\Delta x_t \leftarrow \frac{\sqrt{D_{t-1,k}}}{\sqrt{G_{t,k}}} g_{t,k}$

$\beta_{t+1,k} \leftarrow \beta_t - \Delta x$

$D_{t,k} = \rho (\Delta x)^2 + (1 - \rho) D_{t-1,k}$

end

323 In our implementation, we used $\epsilon =$ and $\rho =$ **TODO look this up**