

Privacy-Preserving Generalized Linear Models using Distributed Block Coordinate Descent

Erik-Jan van Kesteren^{*1}, Chang Sun², Daniel L. Oberski¹, Michel Dumontier^{†2}, and Lianne Ippel^{†2}

¹Department of Methodology and Statistics, Utrecht University

²Institute of Data Science, Maastricht University

November 11, 2019

Abstract

Combining data from varied sources has considerable potential for knowledge discovery: collaborating data parties can mine data in an expanded feature space, allowing them to explore a larger range of scientific questions. However, data sharing among different parties is highly restricted by legal conditions, ethical concerns, and / or data volume. Fueled by these concerns, the fields of cryptography and distributed learning have made great progress towards privacy-preserving and distributed data mining. However, practical implementations have been hampered by the limited scope or computational complexity of these methods. In this paper, we greatly extend the range of analyses available for vertically partitioned data, i.e., data collected by separate parties with different features on the same subjects. To this end, we present a novel approach for privacy-preserving generalized linear models, a fundamental and powerful framework underlying many prediction and classification procedures. We base our method on a distributed *block coordinate descent* algorithm to obtain parameter estimates, and we develop an extension to compute accurate standard errors without additional communication cost. We critically evaluate the information transfer for semi-honest collaborators and show that our protocol is secure against data reconstruction. Through both simulated and real-world examples we illustrate the functionality of our proposed algorithm. Without leaking information, our method performs as well on vertically partitioned data as existing methods on combined data – all within mere minutes of computation time. We conclude that our method is a viable approach for vertically partitioned data analysis with a wide range of real-world applications.

1 Introduction

With technological developments in computational power and storage capacity, an increasing amount of data is collected and stored by a variety of data parties

^{*}Correspondence about this article should be sent to e.vankesteren1@uu.nl

[†]These authors share last authorship

(Kaisler et al., 2013). Over the past decades, data mining has been successful in extracting information from such datasets, but it is especially powerful when various data sources are combined: collaborating data parties can mine data in a larger feature space, allowing them to discover knowledge beyond their individual potential. For example, in the medical domain, personal health conditions are significantly affected not only by genetic and biological factors, but also by individual behaviour and social circumstances (World Health Organization, 2008); combining those sources has the potential to improve analytical models for health outcomes (Kasthurirathne et al., 2017; Ancker et al., 2018).

However, there is a pertinent obstacle to unlocking the potential of combining datasets: integrating various sources may reveal private information about individual data subjects to the collaborating parties. Hence, data sharing is highly restricted by legal and ethical concerns. This highlights the need for privacy-preserving techniques which perform data mining tasks on multiple sources without explicitly sharing their full data (e.g., Du et al., 2004; Gambs et al., 2007; Karr et al., 2009; Gascón et al., 2017). In this paper, we develop a novel algorithm for performing generalized linear modeling (GLM) in a privacy-preserving way in such a partitioned data situation. GLM is a powerful statistical framework for prediction and classification and is at the basis of a wide range of analysis applications including linear, count, and logistic regression (McCullagh and Nelder, 1989; Dobson and Barnett, 2008).

This paper is organized as follows. In Section 2, related work is discussed to contextualize our contribution. In Section 3, we introduce our proposed method for GLM on vertically partitioned data. Next, we describe in detail the privacy-preserving and information sharing characteristics of this protocol in Section 4, and we analyze how the information transfer affects the ability of the partner organisation to recover the collaborator’s data. In Section 5, we benchmark our implementation of the protocol against full-data analysis using Monte Carlo simulations and we illustrate the functionality of our implementation using three different real-life data sets from the UCI Machine Learning repository (Blake and Merz, 1998). Finally, we discuss the strengths and limitations of our approach in Section 6 and we provide suggestions for future research.

All of the methods described here are implemented in `privreg`, an open-source software package for the R programming language (R Core Team, 2018). This implementation includes encryption for all communication across parties based on a pre-shared key, and includes a user-friendly interface based around an object-oriented architecture. The package is available for installation from <https://github.com/vankesteren/privreg>.

2 Related work

In practice, there are two main types of data partitioning (Vaidya and Clifton, 2005). Different data sources might collect the same features of different data subjects, e.g., different hospitals collect the same type of information from their own set of patients. This situation is referred to as *horizontally partitioned* data. Alternatively, separate sources might collect different information from the same data subjects, e.g., medical features by the hospital may be combined with socioeconomic features from a government statistics department. This situation is referred to as *vertically partitioned* data, which is the focus of the

current paper. There is also a third scenario, where data are both vertically and horizontally partitioned, which may be referred to as *hybrid partitioning*.

Our aim is to analyze data which is vertically partitioned without leaking raw data to the collaborating parties (*Alice* and *Bob*). In order to analyze such data, either the dataset may be combined but hidden from the collaborating parties, or the analytical procedure should prevent leaking of information. The former relies on the inclusion of an ‘uninterested’ or trusted third party (TTP): Each party sends their raw data encrypted to the TTP, who then performs the required analyses on the combined data sets. Afterwards, the TTP returns the results to all data parties and the raw data of *Alice* stays hidden to *Bob*. However, this solution requires all parties to fully trust the TTP, which might not be possible in the face of restrictive legislation or sensitive data.

There is another class of methods which do not rely on a TTP, instead using cryptography to perform data mining tasks on vertically partitioned data. These methods focus on preventing information leakage by creating protocols which hide the raw data from the collaborator (e.g., for the construction of decision trees, Agrawal and Srikant, 2000). In this class of methods Du and Atallah (2001) and Du et al. (2004) investigated various protocols for secure matrix computation for linear least squares regression and classification problems. Several other authors used and extended more general secure multiparty computation protocols (e.g., the garbled circuit protocol; Yao, 1986) to perform regression on vertically partitioned data (Amirbekyan and Estivill-Castro, 2007; Slavkovic et al., 2007; Fang et al., 2013; Nikolaenko et al., 2013; Gascón et al., 2016, 2017; Bloom, 2019). While their use of these general protocols yields certain privacy guarantees, their practical implementations and use are hindered by requiring semi-trusted third parties, intermediate data sharing, computational complexity, or a limitation to the linear regression situation.

Another line of research leverages the privacy-preserving properties of algorithms from *federated* or *distributed learning*, a field researching data mining on separated datasets (Li et al., 2019; Dobriban and Sheng, 2018). A canonical example is by Sanil et al. (2004), who developed a method to compute linear regression coefficients iteratively based on an algorithm by Powell (1964). Other authors leverage specific distributed learning algorithms to implement statistical learning for vertically partitioned data (Vaidya and Clifton, 2002, 2003, 2005; Vaidya et al., 2008). Our method is closely related to this branch of research. Unlike existing regression methods from the TTP or cryptography fields, our method does not make use of a trust assumption or complex cryptographic protocols, but it is naturally secure due to its reliance on a federated learning algorithm which never moves the data from its original location. In the next section, we explain the concept and implementation behind our proposed privacy-preserving GLM technique.

3 Proposed method

Our proposed method uses *block coordinate descent* (BCD) to estimate generalized linear models (GLM) in a situation where data is vertically partitioned across two or more parties. In BCD, parameters are iteratively updated for each block of features, cycling over the blocks until an optimum is found (Hastie et al., 2015). This optimization algorithm can be seen as a form of distributed learning

(Bertsekas and Tsitsiklis, 1989; Richtárik and Takáč, 2016) which we exploit as a privacy-preserving method because the features remain in different locations. Only linear predictions need to be transferred across the feature blocks – the full data is never shared.

Note that for the remainder of the paper, we assume that the records of the data subjects are in the same order across databases, in line with Gascón et al. (2017). Furthermore, we only consider the situation where the target attribute is available to both parties (Sanil et al., 2004). In addition, we follow the tradition in the existing literature (e.g., Vaidya et al., 2008; Karr, 2010) to assume semi-honest adversaries: data parties will follow the protocol as described, but will still attempt to learn as much information as possible from other parties. This contrasts with malicious adversaries that can arbitrarily deviate from the protocol (Lindell, 2005).

In this section, we build up the BCD algorithm from the simpler case of linear regression before extending it to full GLM. Therefore, we first explain the necessary background on linear regression, as well as the notation used throughout this paper. Then, coordinate descent estimation is introduced as a means to estimate its maximum likelihood coefficients. In Section 3.3, this algorithm is then extended to accommodate a vertically partitioned data structure, and in Section 3.4 we generalize it to different outcome families in order to estimate GLMs. Finally, we develop a novel method to obtain standard errors within this framework.

3.1 Background

We consider the centered design matrix with features $\mathbf{X} \in \mathbb{R}^{N \times P}$ and the centered target variable $\mathbf{y} \in \mathbb{R}^{N \times 1}$, where N is the sample size, or number of observations, and P is the number of features. The p^{th} column in \mathbf{X} is represented as \mathbf{x}_p . The columns in \mathbf{X} excluding the p^{th} are denoted as \mathbf{X}_{-p} .

The basic regression model is then as follows:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (1)$$

where $\boldsymbol{\beta} \in \mathbb{R}^P$, $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$, and $\boldsymbol{\epsilon} \perp \mathbf{X}$. The well-known closed-form maximum likelihood estimator of the P regression coefficients $\boldsymbol{\beta}$ in this model is:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (2)$$

We further define the vector of predicted values as $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ and the vector of residuals as $\hat{\boldsymbol{\epsilon}} = \mathbf{y} - \hat{\mathbf{y}}$.

3.2 Cyclic coordinate descent estimation

When instead of the full design matrix \mathbf{X} we consider only the p^{th} variable, the estimator in Equation 1 yields the *marginal* regression coefficient. Thus, by simplifying Equation 1 to the univariate case, the marginal coefficient for the p^{th} variable β_p^* is estimated as

$$\hat{\beta}_p^* = \frac{\langle \mathbf{x}_p, \mathbf{y} \rangle}{\langle \mathbf{x}_p, \mathbf{x}_p \rangle} = \frac{\text{cov}(\mathbf{x}_p, \mathbf{y})}{\text{var}(\mathbf{x}_p)} \quad (3)$$

where $\langle \cdot, \cdot \rangle$ indicates the inner product of two vectors. The covariance/variance notation holds because we assume a centered design matrix \mathbf{X} and outcome variable \mathbf{y} .

If \mathbf{x}_p covaries with any of the predictors in \mathbf{X}_{-p} , the marginal coefficient β_p^* is different from the *conditional* coefficient β_p . The estimate of this coefficient is an element of $\hat{\beta}$ in Equation 1, but it can equivalently be estimated in a coordinate-wise, univariate manner (Hastie et al., 2015) as follows:

$$\hat{\beta}_p = \frac{\langle \mathbf{x}_p, \hat{\epsilon}_{-p} \rangle}{\langle \mathbf{x}_p, \mathbf{x}_p \rangle} = \frac{\langle \mathbf{x}_p, \mathbf{y} - \mathbf{X}_{-p} \hat{\beta}_{-p} \rangle}{\langle \mathbf{x}_p, \mathbf{x}_p \rangle} = \frac{\langle \mathbf{x}_p, \mathbf{y} \rangle}{\langle \mathbf{x}_p, \mathbf{x}_p \rangle} - \frac{\langle \mathbf{x}_p, \mathbf{X}_{-p} \hat{\beta}_{-p} \rangle}{\langle \mathbf{x}_p, \mathbf{x}_p \rangle} \quad (4)$$

The residual $\hat{\epsilon}_{-p} = \mathbf{y} - \mathbf{X}_{-p} \hat{\beta}_{-p}$ is the residual with respect to the variables excluding \mathbf{x}_p , evaluated at the maximum likelihood (ML) estimates of β . Equation 4 states that the conditional regression coefficient can be obtained by computing the marginal regression coefficient of $\hat{\epsilon}_{-p}$ on \mathbf{x}_p . This relation holds because $\hat{\epsilon}_{-p}$ represents the part of the outcome variable unrelated to \mathbf{X}_{-p} – by definition, $\hat{\epsilon}_{-p} \perp \mathbf{X}_{-p}$. In addition, the last part of Equation 4 shows that the marginal and conditional estimate of the p^{th} regression coefficient are equal if \mathbf{x}_p and \mathbf{X}_{-p} do not covary, because the last term drops out.

The coordinate-wise estimation of $\hat{\beta}_p$ (Equation 4) requires the maximum likelihood estimates $\hat{\beta}_{-p}$ of the remaining variables to be known. However, when estimation of $\hat{\beta}$ is the goal, these estimates are not available. This can be solved by an iterative updating procedure of the $\hat{\beta}$ estimates:

Algorithm 1: Cyclic coordinate descent

(Hastie et al., 2015)

1. Initialize $\hat{\beta} \leftarrow \hat{\beta}^*$ (marginal coefficients)
2. For each $p \in P$:
 - (a) $\hat{\epsilon}_{-p} \leftarrow \mathbf{y} - \mathbf{X}_{-p} \hat{\beta}_{-p}$
 - (b) $\hat{\beta}_p \leftarrow \langle \mathbf{x}_p, \hat{\epsilon}_{-p} \rangle / \langle \mathbf{x}_p, \mathbf{x}_p \rangle$
3. Repeat step (2.) for R iterations until convergence (i.e., the change in parameter estimates over iterations becomes negligible)

An advantage of this method is that it does not require storing the full $P \times P$ covariance matrix in memory, and this matrix does not need to be inverted – an $\mathcal{O}(P^3)$ operation. This advantage becomes especially relevant as P grows (Hastie et al., 2015). Another advantage is that this estimation method allows for regularization to be implemented naturally. For example, the ℓ_1 penalized parameters can be computed by soft-thresholding $\langle \mathbf{x}_p, \hat{\epsilon}_{-p} \rangle$ in each iteration. This is the approach taken by the popular regularized regression package `glmnet` (Friedman et al., 2010).

A graphical display of the behaviour of the estimated parameters during the cyclical coordinate descent procedure is shown in panel A of Figure 1. Here, 9 covarying features \mathbf{X} were generated from a multivariate normal distribution.

Then random parameter values β and random normal errors ϵ were created and used to generate the target variable $\mathbf{y} = \mathbf{X}\beta + \epsilon$.

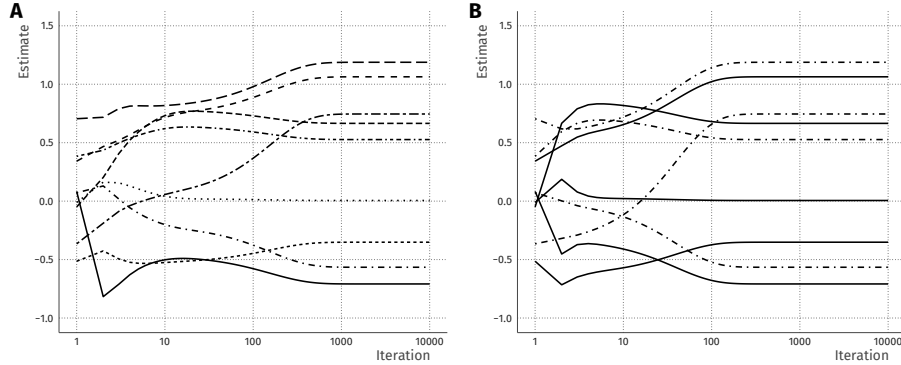


Figure 1: Panel A: Coordinate descent paths for linear regression with 9 covarying features, simulated from a multivariate normal distribution. The parameter lines converge from the marginal ML estimates (left) to the conditional ML estimates (right). Note that the x-axis is on a logarithmic scale and convergence happens around iteration 1000. Panel B: Block coordinate descent path for regression with 9 covarying predictors, applied to the same simulated dataset. There are two blocks, indicated by the line types. Note that convergence happens before iteration 500, faster than the cyclic coordinate descent algorithm.

Next, we show how coordinate descent generalizes to blocks of variables, and how it may be used to estimate linear regression coefficients in the vertically partitioned data scenario described above.

3.3 Securely estimating coefficients for linear regression

In this section, we develop the framework for analysing vertically partitioned data. Our key contribution is the combination of two observations:

1. Coordinate descent estimation works the same for single features as well as for blocks of features – resulting in a variant called block coordinate descent (BCD; Hastie et al., 2015).
2. Vertically partitioned data is blocked data – the features held by *Alice* can be considered the first block, and those held by *Bob* the second block.

Following these two observations, Algorithm 2 below thus provides an iterative estimator for the parameters of *Alice* (β_a) and those of *Bob* (β_b) through sharing of predictions. Predictions from *Alice* are written as $\hat{\mathbf{y}}_a = \mathbf{X}_a \hat{\beta}_a$, and the working residual with respect to *Alice*, i.e., the part of \mathbf{y} not related to the features in \mathbf{X}_a is then $\hat{\epsilon}_a = \mathbf{y} - \hat{\mathbf{y}}_a$.

Algorithm 2: Secure block coordinate descent

1. Initialize $\hat{\mathbf{y}}_b \leftarrow \mathbf{0}$
2. *Alice*:
 - (a) $\hat{\boldsymbol{\epsilon}}_b \leftarrow \mathbf{y} - \hat{\mathbf{y}}_b$
 - (b) $\hat{\boldsymbol{\beta}}_a \leftarrow (\mathbf{X}_a^T \mathbf{X}_a)^{-1} \mathbf{X}_a^T \hat{\boldsymbol{\epsilon}}_b$
 - (c) $\hat{\mathbf{y}}_a \leftarrow \mathbf{X}_a \hat{\boldsymbol{\beta}}_a$
 - (d) Send $\hat{\mathbf{y}}_a$ to *Bob*
3. *Bob*:
 - (a) $\hat{\boldsymbol{\epsilon}}_a \leftarrow \mathbf{y} - \hat{\mathbf{y}}_a$
 - (b) $\hat{\boldsymbol{\beta}}_b \leftarrow (\mathbf{X}_b^T \mathbf{X}_b)^{-1} \mathbf{X}_b^T \hat{\boldsymbol{\epsilon}}_a$
 - (c) $\hat{\mathbf{y}}_b \leftarrow \mathbf{X}_b \hat{\boldsymbol{\beta}}_b$
 - (d) Send $\hat{\mathbf{y}}_b$ to *Alice*
4. Repeat step (2.) and (3.) for R iterations until convergence.

Upon convergence, the concatenated parameter estimates vector $(\hat{\boldsymbol{\beta}}_a, \hat{\boldsymbol{\beta}}_b)$ is equal (up to a small predetermined tolerance value) to the parameter estimates vector that would be obtained using the standard maximum likelihood estimator in the combined data set (Tseng, 2001). It follows that the element-wise summed prediction $\hat{\mathbf{y}}_a + \hat{\mathbf{y}}_b$ is equal to the prediction $\hat{\mathbf{y}}$ that would be obtained from the combined dataset. Thus, prediction can be done without sharing the parameter estimates. Further analysis of the privacy-preserving properties of this procedure is discussed in Section 4.

In panel B of Figure 1 we illustrate BCD, applied to the same data set as in panel A. However, instead of P blocks of 1 feature each, now there are two blocks with 5 and 4 features. BCD reaches convergence with fewer iterations than the cyclic version, because it uses more information about the covariance between the features. In general, convergence is obtained faster with fewer blocks, and with less covariance between blocks (Richtárik and Takáč, 2016). In the case of orthogonal blocks, only a single iteration is needed for convergence as the marginal estimates equal the conditional estimates. Li et al. (2017, Theorem 8) derived a general result about the iteration complexity of BCD, showing that for smooth convex losses such as the GLM log-likelihood, the number of iterations required for convergence is linear in the number of features P .

In the next section, we show how our BCD approach may be modified to estimate generalized linear model coefficients for a wide range of applications. Then, we provide a way to estimate standard errors within this framework.

3.4 Extension to generalized linear models

Extending this procedure to generalized linear models (GLM) requires a slightly different estimation approach: whereas the parameter estimates of full-data linear regression can be found analytically (Equation 2), GLM requires an iteratively reweighted least squares (IRLS) procedure (Wedderburn, 1974; Green, 1984). In each iteration i in full-data GLM, the estimates are computed as follows:

$$\hat{\boldsymbol{\beta}}^{(i+1)} = (\mathbf{X}^T \mathbf{W}^{(i)} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}^{(i)} \mathbf{z}^{(i)} \quad (5)$$

Here, \mathbf{W} is a diagonal weights matrix and \mathbf{z} is a transformation of the target variable called the *working response*, computed as

$$\mathbf{z}^{(i)} = \boldsymbol{\eta}^{(i)} + (\mathbf{y} - \boldsymbol{\mu}^{(i)}) \left(\frac{d\boldsymbol{\mu}^{(i)}}{d\boldsymbol{\eta}^{(i)}} \right) \quad (6)$$

where $\boldsymbol{\eta}^{(i)} = \mathbf{X}\hat{\boldsymbol{\beta}}^{(i)}$ and $\boldsymbol{\mu}^{(i)}$ is a function of $\boldsymbol{\eta}^{(i)}$ as predefined in the link function (e.g., logit link for logistic regression; McCullagh and Nelder, 1989). From this working response, a *working residual* needs to be obtained which acts like $\hat{\boldsymbol{\epsilon}}_{-p}$ in Equation 4: a response vector orthogonal to the predictors excluding feature p . We define this working residual as follows (Friedman et al., 2010):

$$\hat{\boldsymbol{\epsilon}}_{-p} = \mathbf{z} - \mathbf{X}_{-p} \hat{\boldsymbol{\beta}}_{-p} \quad (7)$$

Using this working residual and the usual weights matrix from GLM, the coordinate descent algorithm proceeds in a similar fashion to that of linear regression (Algorithm 1). Just as with coordinate descent for linear regression, this algorithm readily extends to a blockwise procedure, meaning it can be adapted for the private regression method as discussed in Section 3.3.

3.5 Computing standard errors

A key component of inference in regression models is obtaining a measure of sampling uncertainty about the obtained estimates, usually standard errors. Under the assumptions of maximum likelihood theory, the limiting distribution of the deviation of the parameter estimates is the following:

$$\sqrt{N}(\hat{\boldsymbol{\beta}}_N - \boldsymbol{\beta}) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\beta) \quad (8)$$

where $\boldsymbol{\Sigma}_\beta$ is the asymptotic variance-covariance matrix of $\hat{\boldsymbol{\beta}}$:

$$\boldsymbol{\Sigma}_\beta = \text{var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \quad (9)$$

In linear regression, $\hat{\sigma}^2 = \langle \hat{\boldsymbol{\epsilon}}, \hat{\boldsymbol{\epsilon}} \rangle / (N - P)$ and the standard errors of $\hat{\boldsymbol{\beta}}$ can be computed as

$$\text{s}\hat{\boldsymbol{\epsilon}}_{\hat{\boldsymbol{\beta}}} = \sqrt{\text{diag}(\hat{\sigma}^2 (\mathbf{X}^T \mathbf{X})^{-1})} \quad (10)$$

Thus, to compute an estimate of the variance-covariance matrix of the sampling distribution of the $\hat{\boldsymbol{\beta}}$ parameters, the inverse covariance matrix of the features is needed. However, when the data is vertically partitioned, part of this covariance matrix is missing for each party. As a result, computing standard errors using the above information matrix approach is impossible for vertically

partitioned data without sharing the features.

We present a novel approach to compute standard errors of the regression coefficient through creating a substitute \mathbf{V}_b of the partner’s data matrix \mathbf{X}_b . This substitute is then used as the partner’s data in the computation of the asymptotic variance-covariance matrix as in Equation 9.

The substitute \mathbf{V}_b needs to contain the same information for the parameters of Alice as the real data. This information is in the predictions received from Bob – the parameter estimates of Alice depend only on Bob’s linear predictions. Consider the inputs and outputs of Bob, as seen by Alice: as the coordinate descent algorithm progresses along the R iterations, Alice can create two $N \times R$ matrices, $\hat{\mathbf{E}}_a$ and $\hat{\mathbf{Y}}_b$

$$\begin{aligned}\hat{\mathbf{E}}_a &= \left[\hat{\epsilon}_a^{(1)}, \dots, \hat{\epsilon}_a^{(R)} \right] \\ \hat{\mathbf{Y}}_b &= \left[\hat{\mathbf{y}}_b^{(1)}, \dots, \hat{\mathbf{y}}_b^{(R)} \right]\end{aligned}\tag{11}$$

These are the input and output matrices, respectively, from the projection that Bob applies in each iteration. This projection is commonly known as the *hat matrix* $\mathbf{H}_b \in \mathbb{R}^{N \times N}$. The hat matrix relates to Bob’s data matrix \mathbf{X}_b as follows:

$$\begin{aligned}\hat{\mathbf{Y}}_b &= \mathbf{H}_b \hat{\mathbf{E}}_a \\ \hat{\mathbf{Y}}_b &= \mathbf{X}_b (\mathbf{X}_b^T \mathbf{X}_b)^{-1} \mathbf{X}_b^T \hat{\mathbf{E}}_a \\ \hat{\mathbf{Y}}_b &= \mathbf{X}_b \mathbf{X}_b^+ \hat{\mathbf{E}}_a\end{aligned}\tag{12}$$

where \mathbf{X}_b^+ indicates the Moore-Penrose generalized inverse of \mathbf{X}_b (Petersen and Pedersen, 2012).

Alice can compute the projection that Bob applies in each iteration \mathbf{H}_b as follows:

$$\hat{\mathbf{H}}_b = \hat{\mathbf{Y}}_b \hat{\mathbf{E}}_a^+\tag{13}$$

Across iterations, this minimum-norm solution $\hat{\mathbf{H}}_b$ performs the same projection as the true hat matrix of Bob. Using this projection, Alice can then create the data substitute $\mathbf{V}_b \in \mathbb{R}^{N \times P_b}$. For this, \mathbf{V}_b should have the property $\hat{\mathbf{H}}_b = \mathbf{V}_b \mathbf{V}_b^+$. Such a \mathbf{V}_b has the same effect on the coefficient estimates of Alice that \mathbf{X}_b does, because it generates the same predictions that Bob does:

$$\begin{aligned}\hat{\mathbf{Y}}_b &= \hat{\mathbf{H}}_b \hat{\mathbf{E}}_a \\ \hat{\mathbf{Y}}_b &= \mathbf{V}_b \mathbf{V}_b^+ \hat{\mathbf{E}}_a\end{aligned}\tag{14}$$

There is no unique solution to decomposing $\hat{\mathbf{H}}_b$ into an $N \times P$ matrix \mathbf{V}_b and its pseudoinverse. However, a numerically convenient \mathbf{V}_b solution can be found as the first P_b eigenvectors of $\hat{\mathbf{H}}_b$. This is a convenient choice, because the columns of \mathbf{V}_b are then orthogonal, meaning they also have the following property: $\mathbf{V}_b^+ = (\mathbf{V}_b^T \mathbf{V}_b)^{-1} \mathbf{V}_b^T = \mathbf{I}^{-1} \mathbf{V}_b^T = \mathbf{V}_b^T$. As follows from Equations 12 and 14, the \mathbf{V}_b matrix relates to \mathbf{X}_b by means of an unknown positive definite rotation matrix $\mathbf{V}_b = \mathbf{R} \mathbf{X}_b$ (Pavel, 2019).

By leveraging this similarity of \mathbf{V}_b to \mathbf{X}_b , Alice can create an augmented data matrix of the following form: $\mathbf{Z}_a = [\mathbf{X}_a, \mathbf{V}_b]$. The augmented data matrix

replaces the full data matrix in the computation of the asymptotic covariance matrix: $\Sigma_{\beta}^{(a)} = \sigma^2(\mathbf{Z}_a^T \mathbf{Z}_a)^{-1}$. The partition of $\Sigma_{\beta}^{(a)}$ belonging to β_a is then identical to its counterpart from the full data asymptotic covariance matrix Σ_{β} (for proof see Appendix A). The square root of its diagonal elements are thus the correct standard errors that would be obtained had the full data been available.

Alternative standard error procedures are available, e.g., profile likelihood methods or bootstrapping, but those require additional iterations of the main block coordinate descent algorithm. This yields additional information leakage and dramatically increases time requirements. Conversely, in the novel procedure we suggest here, both parties efficiently leverage the information in the existing iterations to compute standard errors without additional communication.

4 Privacy analysis for block coordinate descent

In this section, we analyze the information transfer within our protocol for privacy-preserving regression based on block coordinate descent. In line with previous work on this topic (e.g. Gambs et al., 2007; Gascón et al., 2017; Vaidya and Clifton, 2003, 2005; Vaidya et al., 2008), we take the viewpoint of semi-honest parties: *Alice* and *Bob* follow the protocol accurately, though they may be curious and aim to recover the other party’s data. In this section, we aim to identify how well *Bob* can approximate *Alice*’s data using a *model inversion attack* (Fredrikson et al., 2015; Wang et al., 2015).

4.1 Information transfer in vertically partitioned regression

Information about features cannot only leak through dataset sharing, but also via sharing statistics based on this data. For example, a simple method for regression without explicitly sharing the full dataset is that by Karr et al. (2009), who compute the covariance matrix of \mathbf{X} using secure inner-product methods and share it between *Alice* and *Bob*. This covariance matrix allows even a semi-honest *Alice* to (a) know how many features are used by *Bob* and – in the case of categorical predictors – know how many categories there are, (b) predict the values of the features held by *Bob* based on the values of the features held by *Alice*, (c) compute standard errors around this prediction, and (d) compute an R^2 value for this prediction. In other words, in a shared covariance matrix setting *Alice* can know up to a certain degree the values on each of *Bob*’s features for each row in the dataset, and *Alice* can know how good this prediction is. Moreover, each additional feature entered by *Alice* improves the prediction of features at *Bob* by definition.

Thus, sharing the full covariance matrix is undesirable for privacy-preserving regression. Newer methods (e.g., Du et al., 2004; Gascón et al., 2017) result in additive shares of $\text{cov}(\mathbf{X})$ at *Alice* and *Bob*, without either of them possessing the full covariance matrix. Afterwards, separate secure multiparty matrix inversion protocols or linear system solvers are used to compute the regression

parameters according to Equation 2. This generally requires complex protocols involving multiple parties, but has been argued to be a secure procedure for obtaining parameter estimates for linear regression with vertically partitioned data. In these protocols, it is clear that information transfer does occur (because the full-data estimates are obtained) but its extent is not made explicit: it is unclear how the additive shares of the covariance matrix (the “statistics”) relate to the collaborator’s data – and thus it is unclear whether that data can be reconstructed.

Conversely, in our protocol the covariance matrix of the combined data is never explicitly computed. Our method uses a different “statistic”: predictions $\hat{\mathbf{y}}$ over R iterations. Each of the R predictions are computed as follows by *Alice*:

$$\hat{\mathbf{y}}_a^{(r)} = \mathbf{X}_a \hat{\boldsymbol{\beta}}_a^{(r)} \quad (15)$$

This prediction vector is then sent to *Bob*: the main information transfer. In this protocol, how this information transfer relates to *Alice*’s data is thus explicit. As a result, clear conclusions can be made as to the potential for data recovery.

In the case where *Alice* enters only a single continuous feature in the analysis protocol, the information contained in $\hat{\mathbf{y}}_a$ is sufficient for *Bob* to reproduce the values of this feature up to a multiplicative constant: $\hat{\mathbf{y}}_a = \mathbf{x}_a \cdot \hat{\beta}_a$. With more than one feature per party, $\hat{\boldsymbol{\beta}}_a$ becomes a vector, meaning the problem of recovering the values of any feature at *Alice* is underidentified. Moreover, if the protocol is followed precisely, *Bob* does not know the number of features P entered into the model, meaning there is additional uncertainty about the values of \mathbf{X}_a on the part of *Bob*. In its most basic form, the protocol is therefore fully secure for semi-honest parties against reconstruction of the privacy-sensitive data matrices.

4.2 Data reconstruction using shared metadata

In practice, there are many situations where the basic algorithm does not suffice and metadata about \mathbf{X}_a should be shared with *Bob*. For example, to circumvent multicollinearity and non-convergence, none of the features entered into the model by *Alice* should be entered by *Bob*. Moreover, when distributing the model results is a goal of the analysis, it is relevant to investigate how sharing parameter estimates in addition to the predictions that are already shared leads to information transfer about the original data.

In our protocol, *Alice* sends R predictions to *Bob*. These individual predictions can be appended in a columnwise fashion to create an $N \times R$ matrix $\hat{\mathbf{Y}}_a = [\hat{\mathbf{y}}_a^{(1)}, \dots, \hat{\mathbf{y}}_a^{(R)}]$. Each prediction has an associated set of parameter estimates known only by *Alice* $\hat{\boldsymbol{\beta}}_a^{(r)}$, which can be combined in a similar way to create the matrix $\hat{\mathbf{B}}_a \in \mathbb{R}^{P \times R}$. These relate to the data matrix at *Alice* as follows:

$$\hat{\mathbf{Y}}_a = \mathbf{X}_a \hat{\mathbf{B}}_a \quad (16)$$

In our protocol, all of $\hat{\mathbf{Y}}_a$ is shared with *Bob*, and only the R^{th} column of $\hat{\mathbf{B}}_a$ – the final model result – is shared. Using these estimates, *Bob* can make a rank-1 minimum-norm approximation of the data held by *Alice*:

$$\hat{\mathbf{X}}_a^{(1)} = \hat{\mathbf{y}}_a^{(R)} \hat{\boldsymbol{\beta}}_a^{(R)+} \quad (17)$$

where $+$ indicates the Moore-Penrose inverse (Petersen and Pedersen, 2012). We show empirically in Appendix B that using this method with one set of shared parameter estimates reveals a proportion $1/P_a$ of the variance in the data to *Bob*. Only the combination of predictions and their associated parameter values allows (partial) model inversion and reconstruction of the partner’s data.

Furthermore, as presented in Section 3.5, the predictions sent to and received from *Alice* can be used to create a minimum-norm approximation of the hat matrix of *Alice* – another statistic which is shared in our protocol. This hat matrix is shown in Appendix A to not contain information about the features of *Alice* directly, but only about a rotation of this data such that the parameter estimates of *Bob* are adequately adjusted towards the conditional estimates.

In conclusion, the protocol is secure against reconstruction of the data in the case of semi-honest parties, and sharing of the final parameter estimates $\hat{\beta}_a$ reveals a proportion $1/P_a$ of the variance in the data to the other data party.

4.3 Further privacy considerations

Purposeful attacks to recover data in the case of adversarial collaborations have not been analyzed. It is possible to design such an attack, but it is also possible to design safeguards against such attacks in the implementation of the protocol, for example based on the expected smoothness of the regression paths over iterations. We leave this analysis as a topic for further research.

In addition, because of the explicit link between the shared statistics and the original data, it is possible to limit the information shared with the collaborator in several ways. For example, in each iteration *Alice* may add noise to the computed parameter estimates or to the predictions sent to *Bob* – a technique from the differential privacy literature (Dwork et al., 2006). Another method is to put an upper bound on the number of iterations based on the number of features in the data. This has two effects: (a) it shrinks (regularizes) the parameter estimates towards the marginal estimates and (b) it creates an upper bound ϵ on the information shared, depending on the allowed number of iterations.

In the next section, we show how our implementation of the BCD with vertically partitioned data performs in comparison to full-data generalized linear modeling (GLM) in simulated data as well as three real-world datasets.

5 Experiments

Our implementation of the BCD algorithm for vertically partitioned data is provided as an R package at <https://github.com/vankesteren/privreg>. Here, we use this implementation (version 0.9.5) to estimate models on both simulated data (Section 5.1) and real-world data with multiple parties from the UCI data repository (Section 5.2). Reproducible code for this section is available in the supplementary material to this paper.

5.1 Simulated data

The goal of this section is to compare our proposed privacy-preserving regression method to a benchmark method under controlled conditions. The benchmark

method for these experiments is linear and logistic regression with a complete dataset, since the optimum privacy-preserving method would attain the same results with vertically partitioned data. For this section, data with multiple features and one target were simulated in the R programming language (R Core Team, 2018), with the following manipulations:

- Target** Either a normally distributed or a binomial target variable. In the case of the normal target, the R^2 was set to 0.5.
- Dimensionality** The total number of features was either 10, 50, or 100.
- Covariance** The covariance matrix of the features was had 1 on the diagonal and either 0.1 (low covariance) or 0.5 (high covariance) on all off-diagonal elements.

For each condition, 100 datasets were randomly generated. For the privacy-preserving regression method, the generated features were then equally distributed among *Alice* and *Bob*, after which the estimation was started. As a baseline comparison, a generalized linear model was estimated on the full dataset with all the features using the `glm()` function from the base R `stats` package. The exact data-generating mechanism, as well as the estimation method and hyperparameters can be found in the supplementary material.

The empirical convergence rates for the privacy-preserving regression method are shown in Figure 2. As expected from the work of Li et al. (2017), the number of iterations required increases linearly with the number of features. In addition, the high covariance leads to slower convergence due to the conditional estimates lying further away from the marginal estimates. As mentioned in Section 3.3, with no covariance the number of iterations would be 1.

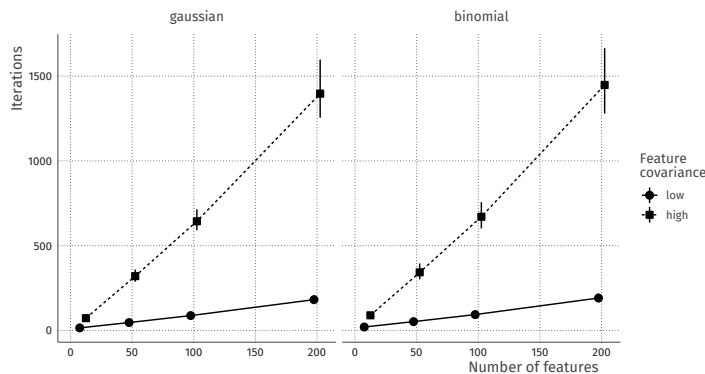


Figure 2: Observed amount of iterations required for convergence is approximately linear in the number of features and increases as there is more covariance between the features. Error bars indicate 95% simulation percentile intervals.

The obtained parameter estimates ($\hat{\beta}$) of our method are equal to those found by the baseline comparison method in all simulated conditions, up to a computational tolerance in the convergence of the estimation algorithm (Figure

3). This lack of relative bias indicates that the proposed privacy-preserving regression approach performs as well as full-data generalized linear models, at least for the extent of these simulations.

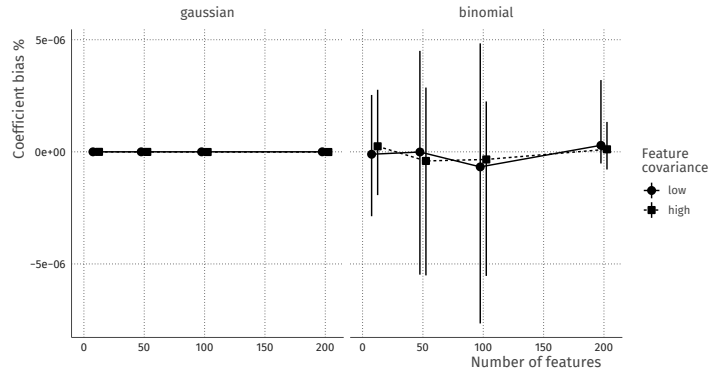


Figure 3: The parameter bias relative to the baseline GLM method is negligible for any number of features and feature covariance strength. Note the small y-axis range.

Standard errors indicate uncertainty in the dataset around the coefficient values, and they are the basis for statistical significance tests. Figure 4 shows the bias in the standard errors relative to the baseline GLM method for the different conditions. The figure shows that variation of this bias over different datasets increases with the number of features (larger error bars). In addition, there seems to be a very slight relative overestimation of the standard errors on average. This is due to slightly different convergence criteria and tolerances for both methods, which propagates through the standard error procedure (Section 3.5). Despite this, the standard error bias is overall small ($< 3\%$).

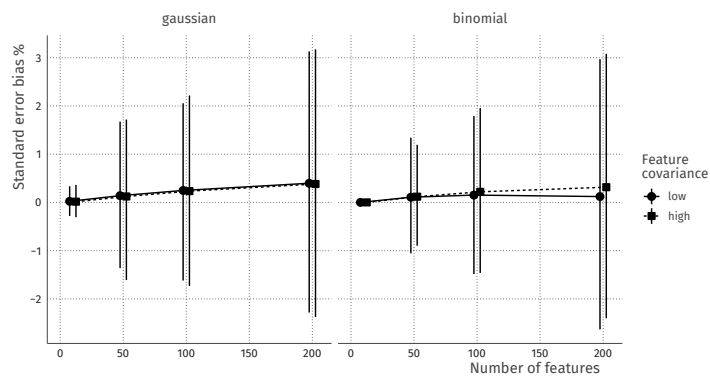


Figure 4: Standard error bias in percentage relative to the baseline GLM method. Variation across datasets increases with the number of features, and there is a very slight trend ($< 1\%$) towards overestimation of the standard error for larger datasets.

In conclusion, the simulations have shown that privacy-preserving regression using block coordinate descent on vertically partitioned data has equal performance to established regression methods on full data. However, in this section the data has been simulated to behave according to specification. In the next section, we compare the performance of these two methods on real-world datasets.

5.2 UCI datasets

In this section, we tested our proposed method on three different real-life datasets from the UCI (University of California at Irvine) Machine Learning repository (Blake and Merz, 1998). The datasets were chosen because they can be naturally partitioned into two sources, and their size and targets are different (Table 1). As before, the full preprocessing and analysis code for this section is available in the supplementary materials. Analyses were run on two separate computers (an Intel Core i7-8750H at 2.20 GHz and an Intel Xeon E5-2650 v4 at 2.20GHz) connected via a gigabit Ethernet connection on a university network.

Dataset	Features	Instances	Task	Parties
Forest fire	13	517	Regression	Weather & Fire dept.
HCC	49	165	Classification	Lab & Clinic
Diabetes	43	15 000	Classification	Clinic & Pharmacy

Table 1: Properties of the datasets used from the UCI machine learning repository after dataset cleaning and pre-processing. Code can be found in the supplementary materials.

5.2.1 Forest fires data

The forest fire data comes from the Montesinho natural park in Portugal (Cortez and Morais, 2007). It contains several weather observations by a meteorological station (e.g. wind speed, temperature, relative humidity, etc) as well as fire department risk assessments. In this dataset, the target is to predict the area of forest burned by a particular fire using the features from the aforementioned parties.

We performed linear regression where the target was log-transformed to normalize the residuals. Continuous features were standardized before they were entered into the analysis. The analysis took 450 BCD iterations in the privacy-preserving regression case. Including encryption and networking overhead, estimation took 14.51 seconds and computing standard errors took 0.61 seconds. Figure 5 shows that the coefficients and their 95% confidence intervals are equal for the full-data analysis and the privacy-preserving procedure. Several months show a significant positive effect on the log-area, meaning that – conditional on the ratings of the fire department – fires in these months (e.g., August and December) burn larger areas of forest.

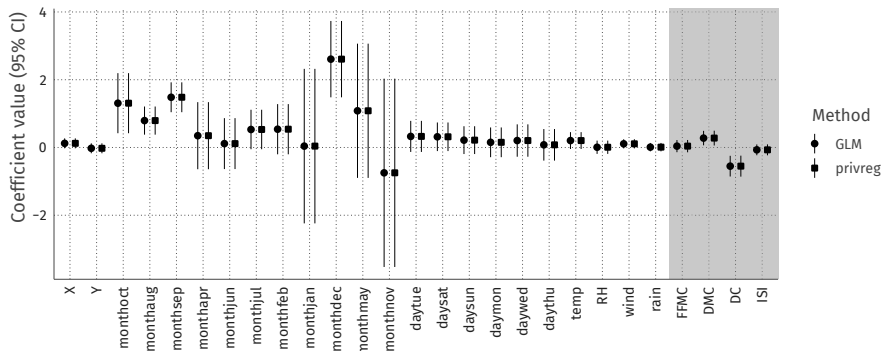


Figure 5: The coefficients and standard errors for the forest fire analysis are exactly the same for the GLM and our privacy-preserving regression estimation methods. The shading indicates data partitioning into the weather service (light) and fire department (dark).

5.2.2 Hepatocellular carcinoma data

This dataset was collected by Coimbra’s Hospital and University Centre in Portugal for studying an epithelial cell cancer of the liver called hepatocellular carcinoma (HCC) (Santos et al., 2015). It contains heterogeneous data on demographics, risk factors, laboratory and overall survival features from HCC patients. The goal of the analysis is to use lab results for a tissue sample as well as clinical data for the patient to predict survival after diagnosis. Since survival is a binary target, a binomial family GLM (logistic regression) was performed. For this analysis, continuous features were standardized before the analysis, which improved the convergence characteristics. The privacy-preserving GLM converged in 1636 iterations. Including encryption and networking overhead, estimation took 3 minutes and 16 seconds and computing standard errors took 0.63 seconds.

The results of the analysis (Figure 6) show that the estimates are exactly equal across the full-data and the privacy-preserving analyses, meaning survival probability predictions for new incoming patients based on these models will be the same. Despite slight deviations in the width of the confidence intervals, conclusions about the effects of the features on survival are also the same in this dataset.

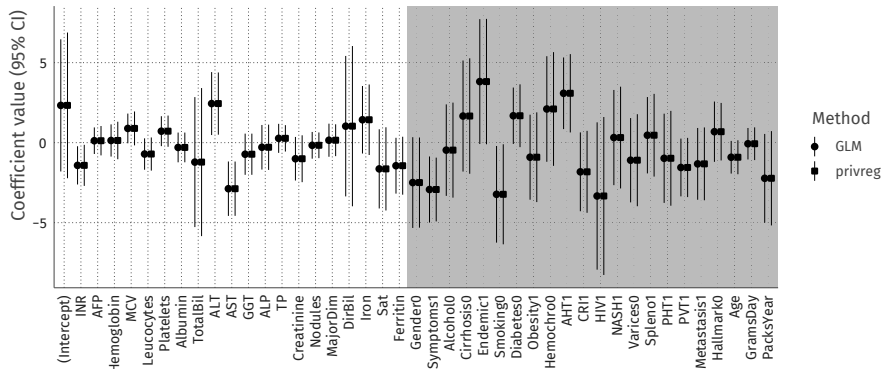


Figure 6: The coefficients and standard errors for the carcinoma analysis are very similar for the GLM and our privacy-preserving regression estimation methods. The shading indicates data partitioning into the lab results (light) and clinic (dark).

5.2.3 Diabetes

The diabetes dataset is an extract representing 10 years (1999-2008) of clinical diabetes care at 130 hospitals and integrated delivery networks throughout the United States (Strack et al., 2014). It is a large and also heterogeneous data set including encounter data (emergency, outpatient, and inpatient), provider speciality, demographics, laboratory data, pharmacy data, in-hospital mortality, and hospital characteristics. In this dataset, we predict readmission to the hospital using both administrative features and pharmaceutical features. To keep the computation of the standard errors for this analysis possible, 15000 patients were randomly selected from the dataset. Features were re-coded where necessary, and categorical features with only a single category in the sample were excluded from the analysis. The full pre-processing pipeline can be found in the supplementary material.

Since readmission is a binary target, a binomial family GLM (logistic regression) was performed. The diabetes data analysis required 284 iterations of the BCD algorithm. Including encryption and networking overhead, estimation took 1 minute and 37 seconds and computing standard errors took 42 seconds. This analysis is particularly interesting with respect to the effect of insulin (`insulinYes`) on the readmission probability. In the analysis of only the medication data, insulin has a significant positive effect on readmission ($OR = 1.20, p < .001$), whereas conditional on the administrative data, insulin significantly reduces the readmission probability ($OR = 0.88, p < .001$). This is a strong argument for including the data of both parties in the analysis.

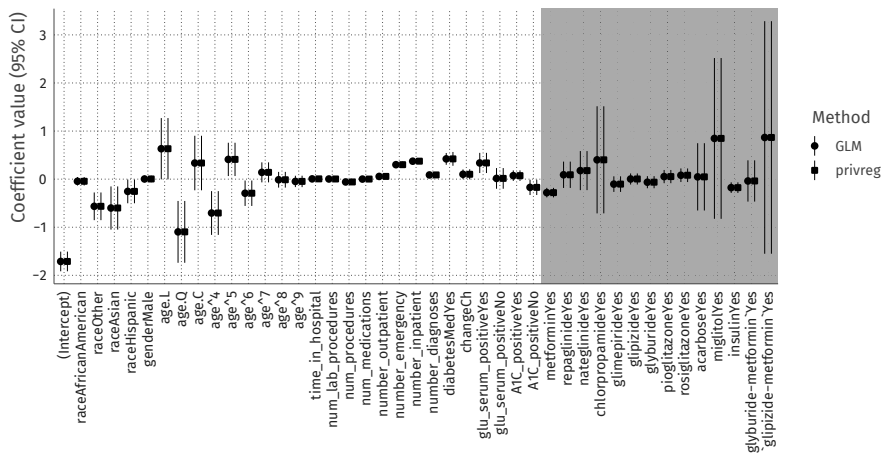


Figure 7: The coefficients and standard errors for the diabetes analysis are exactly the same for the GLM and our privacy-preserving regression estimation methods. The shading indicates data partitioning into the clinical data (light) and pharmaceutical data (dark).

In this section, we have shown that privacy-preserving regression using block coordinate descent is not only a theoretical possibility, but also a viable implementation of GLM for analyzing data with varied characteristics – both in simulated data under controlled conditions (Section 5.1) and in real-world prediction and analysis problems with various targets (Section 5.2). The time constraints on the real-world analyses are manageable, with all example analyses converging in under 4 minutes. We have shown that the parameter estimates exactly match those of the existing reference methods, and that our novel estimation method for the standard errors generally agrees with its full-data counterpart – and where it did not the difference was so small that it led to the same conclusions in the analysis.

6 Discussion

In this paper, we have argued that block coordinate descent is a general method for estimating conditional parts of a generalized linear model (GLM) in a vertically partitioned data situation. Using this approach, two or more data parties can collaboratively estimate a GLM without sharing their features. This is useful when the features are not allowed to be shared, for example when there are privacy issues.

Our method falls within the category of federated learning algorithms. This means it can be implemented for situations when data mining is to be performed over remote devices or siloed data centers (Li et al., 2019), where aggregating the data tables is prohibitively expensive in terms of time, computation, or storage costs. This work aligns with several recent contributions that seek to exploit the privacy-preserving aspects of federated learning algorithms (see, e.g., Bonawitz et al., 2016; Geyer et al., 2017).

Due to the accessibility of our protocol and its similarity to existing regression estimation methods, extensions are relatively simple to implement. First and foremost, our framework can be extended to multiple parties as coordinate descent naturally extends to multiple blocks. In addition, our algorithm could include penalties for regularized estimation of the regression parameters through thresholding (Friedman et al., 2010). Through further research into combining coordinate descent with missing data methods such as full information maximum likelihood (Enders, 2001), our protocol could even be extended for a hybrid partitioning situation where data is both horizontally and vertically partitioned.

Our novel approach is a natural modification of the familiar linear modeling framework – without changes in the assumptions. We argue that our protocol restricts statistical information sharing as much as possible, while being explicit in how the shared information relates to the original data. Because of this, data parties know how much information they share, and the protocol could even incorporate methods from the differential privacy literature – such as additive noise or early stopping – to put a restriction on the amount of information shared with the partner institution (Dwork et al., 2006).

The main tradeoff of this flexibility compared to existing methods is relatively high communication cost: each iteration requires N prediction values to be sent to the partner institution. In addition, like other methods for this situation the block coordinate descent assumes (probabilistic) linkage of the individual records – both parties need to have their records in the same order. Lastly, this method is possible only when the target can be shared, although in absence of a shareable target collaborators could still perform some form of transfer learning, e.g., by predicting a shareable feature *related* to the true target.

Considering the prospect of these extensions and the availability of an accessible open-source implementation, we believe the proposed block coordinate descent protocol can be a springboard for future developments in the privacy-preserving data mining field.

Compliance with Ethical Standards

Funding: this work was supported by the Netherlands Organization for Scientific Research (NWO) under grant number 406.17.057 and by the Dutch National Research Agenda (NWA) under project number 400.17.605.

Conflict of Interest: The authors declare that they have no conflict of interest.

Acknowledgments: We thank Ayoub Bagheri for his helpful comments on an earlier version of this manuscript.

References

- Agrawal R, Srikant R (2000) Privacy-preserving data mining. In: ACM Sigmod Record, ACM, vol 29, pp 439–450
- Amirbekyan A, Estivill-Castro V (2007) Privacy-preserving regression algorithms. In: Pro-

- ceedings of the 7th WSEAS International Conference on simulation, modelling and optimization, World Scientific and Engineering Academy and Society (WSEAS), pp 37–45
- Ancker JS, Kim MH, Zhang Y, Zhang Y, Pathak J (2018) The potential value of social determinants of health in predicting health outcomes. *Journal of the American Medical Informatics Association* 25(8):1109–1110
- Bertsekas DP, Tsitsiklis JN (1989) *Parallel and distributed computation: numerical methods*, vol 23. Prentice hall Englewood Cliffs, NJ
- Blake CL, Merz CJ (1998) *UCI repository of machine learning databases*, 1998
- Bloom JM (2019) Secure multi-party linear regression at plaintext speed. 1901.09531
- Bonawitz K, Ivanov V, Kreuter B, Marcedone A, McMahan HB, Patel S, Ramage D, Segal A, Seth K (2016) Practical secure aggregation for federated learning on user-held data. 1611.04482
- Cortez P, Morais AdJR (2007) A data mining approach to predict forest fires using meteorological data. In: *Proceedings of 13th Portuguese Conference on Artificial Intelligence, Associação Portuguesa para a Inteligência Artificial (APPIA)*, pp 512–523
- Dobriban E, Sheng Y (2018) Distributed linear regression by averaging. 1810.00412
- Dobson AJ, Barnett AG (2008) *An introduction to generalized linear models*. Chapman and Hall/CRC
- Du W, Atallah MJ (2001) Privacy-preserving cooperative scientific computations. In: *csfw, Citeseer*, p 0273
- Du W, Han YS, Chen S (2004) Privacy-preserving multivariate statistical analysis: Linear regression and classification. In: *Proceedings of the 2004 SIAM international conference on data mining, SIAM*, pp 222–233
- Dwork C, McSherry F, Nissim K, Smith A (2006) Calibrating noise to sensitivity in private data analysis. In: *Theory of cryptography conference, Springer*, pp 265–284
- Enders CK (2001) The performance of the full information maximum likelihood estimator in multiple regression models with missing data. *Educational and Psychological Measurement* 61(5):713–740
- Fang W, Zhou C, Yang B (2013) Privacy preserving linear regression modeling of distributed databases. *Optimization Letters* 7(4):807–818
- Fredrikson M, Jha S, Ristenpart T (2015) Model inversion attacks that exploit confidence information and basic countermeasures. In: *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, ACM*, pp 1322–1333
- Friedman J, Hastie T, Tibshirani R (2010) Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software* 33(1):1
- Gams S, Kégl B, Aïmeur E (2007) Privacy-preserving boosting. *Data Mining and Knowledge Discovery* 14(1):131–170
- Gascón A, Schoppmann P, Balle B, Raykova M, Doerner J, Zahur S, Evans D (2016) Secure linear regression on vertically partitioned datasets. *IACR Cryptology ePrint Archive* 2016:892
- Gascón A, Schoppmann P, Balle B, Raykova M, Doerner J, Zahur S, Evans D (2017) Privacy-preserving distributed linear regression on high-dimensional data. *Proceedings on Privacy Enhancing Technologies* 2017(4):345–364
- Geyer RC, Klein T, Nabi M (2017) Differentially private federated learning: A client level perspective. 1712.07557

- Green PJ (1984) Iteratively reweighted least squares for maximum likelihood estimation, and some robust and resistant alternatives. *Journal of the Royal Statistical Society: Series B (Methodological)* 46(2):149–170
- Hastie T, Tibshirani R, Wainwright M (2015) *Statistical Learning with Sparsity: The Lasso and Generalizations*. CRC Press, Boca Raton, DOI 10.1201/b18401-1
- Kaisler S, Armour F, Espinosa JA, Money W (2013) Big data: Issues and challenges moving forward. In: 2013 46th Hawaii International Conference on System Sciences, IEEE, pp 995–1004
- Karr AF (2010) Secure statistical analysis of distributed databases, emphasizing what we don't know. *Journal of Privacy and Confidentiality* 1(2)
- Karr AF, Lin X, Sanil AP, Reiter JP (2009) Privacy-preserving analysis of vertically partitioned data using secure matrix products. *Journal of Official Statistics* 25(1):125
- Kasthurirathne SN, Vest JR, Menachemi N, Halverson PK, Grannis SJ (2017) Assessing the capacity of social determinants of health data to augment predictive models identifying patients in need of wraparound social services. *Journal of the American Medical Informatics Association* 25(1):47–53
- Li T, Sahu AK, Talwalkar A, Smith V (2019) Federated learning: Challenges, methods, and future directions. arXiv preprint arXiv:190807873
- Li X, Zhao T, Arora R, Liu H, Hong M (2017) On faster convergence of cyclic block coordinate descent-type methods for strongly convex minimization. *The Journal of Machine Learning Research* 18(1):6741–6764
- Lindell Y (2005) Secure multiparty computation for privacy preserving data mining. In: *Encyclopedia of Data Warehousing and Mining*, IGI Global, pp 1005–1009
- McCullagh P, Nelder J (1989) *Generalized Linear Models*. Chapman and Hall/CRC, New York
- Nikolaenko V, Weinsberg U, Ioannidis S, Joye M, Boneh D, Taft N (2013) Privacy-preserving ridge regression on hundreds of millions of records. In: 2013 IEEE Symposium on Security and Privacy, IEEE, pp 334–348
- Pavel A (2019) Decompose projection matrix into a matrix and its pseudoinverse. Mathematics Stack Exchange, URL <https://math.stackexchange.com/q/3317493>, 2019-08-12
- Petersen K, Pedersen M (2012) *The matrix cookbook*, version 20121115. Technical Univ Denmark, Kongens Lyngby, Denmark, Tech Rep 3274
- Powell MJ (1964) An efficient method for finding the minimum of a function of several variables without calculating derivatives. *The computer journal* 7(2):155–162
- R Core Team (2018) R: A language and environment for statistical computing. URL <https://www.r-project.org/>
- Richtárik P, Takáč M (2016) Distributed coordinate descent method for learning with big data. *The Journal of Machine Learning Research* 17(1):2657–2681
- Sanil AP, Karr AF, Lin X, Reiter JP (2004) Privacy preserving regression modelling via distributed computation. In: *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ACM, New York, NY, USA, pp 677–682, DOI 10.1145/1014052.1014139
- Santos MS, Abreu PH, García-Laencina PJ, Simão A, Carvalho A (2015) A new cluster-based oversampling method for improving survival prediction of hepatocellular carcinoma patients. *Journal of biomedical informatics* 58:49–59
- Slavkovic AB, Nardi Y, Tibbits MM (2007) Secure logistic regression of horizontally and vertically partitioned distributed databases. In: *Proceedings of the Seventh IEEE International Conference on Data Mining Workshops*, IEEE Computer Society, Washington, DC, USA, pp 723–728, DOI 10.1109/ICDMW.2007.84

- Strack B, DeShazo JP, Gennings C, Olmo JL, Ventura S, Cios KJ, Clore JN (2014) Impact of hba1c measurement on hospital readmission rates: analysis of 70,000 clinical database patient records. *BioMed research international* 2014
- Tseng P (2001) Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of optimization theory and applications* 109(3):475–494
- Vaidya J, Clifton C (2002) Privacy preserving association rule mining in vertically partitioned data. In: *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, pp 639–644
- Vaidya J, Clifton C (2003) Privacy-preserving k-means clustering over vertically partitioned data. In: *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, pp 206–215
- Vaidya J, Clifton C (2005) Privacy-preserving decision trees over vertically partitioned data. In: *IFIP Annual Conference on Data and Applications Security and Privacy*, Springer, pp 139–152
- Vaidya J, Yu H, Jiang X (2008) Privacy-preserving svm classification. *Knowledge and Information Systems* 14(2):161–178
- Wang Y, Si C, Wu X (2015) Regression model fitting under differential privacy and model inversion attack. In: *Twenty-Fourth International Joint Conference on Artificial Intelligence*
- Wedderburn RW (1974) Quasi-likelihood functions, generalized linear models, and the gauss—newton method. *Biometrika* 61(3):439–447
- World Health Organization (2008) *Closing the gap in a generation: Health equity through action on the social determinants of health*. World Health Organization
- Yao ACC (1986) How to generate and exchange secrets. In: *27th Annual Symposium on Foundations of Computer Science (sfcs 1986)*, IEEE, pp 162–167

A Proof for recovery of standard errors

Let $A = X^T X$, partitioned into four submatrices A_{11} (held by Alice), A_{22} (held by Bob), and A_{22} (unknown to either). The standard inverse of such a partitioned, positive definite symmetric matrix is

$$\begin{aligned} A^{-1} &= \begin{pmatrix} B_{11} & B_{12} \\ B_{22} & B_{22} \end{pmatrix} \\ &= \begin{pmatrix} (A_{11} - A_{12}A_{22}^{-1}A_{12}^T)^{-1} & -A_{11}^{-1}A_{12}(A_{22} - A_{12}^T A_{11}^{-1}A_{12})^{-1} \\ -A_{22}^{-1}A_{12}^T(A_{11} - A_{12}A_{22}^{-1}A_{12}^T)^{-1} & (A_{22} - A_{12}^T A_{11}^{-1}A_{12})^{-1} \end{pmatrix} \end{aligned} \quad (18)$$

Following the procedure outlined in Section 3.5, Alice replaces X_2 with $V_2 = R_2 X_2$, and Bob replaces X_1 with $V_1 = R_1 X_1$, where R_j are unknown orthogonal rotation matrices. This gives two new matrices, $A^{(1)}$ and $A^{(2)}$, and their inverses, $B^{(1)}$ and $B^{(2)}$. By substitution,

$$\begin{aligned} A_{12}^{(1)} &= X_1^T R_2 X_2 \\ A_{22}^{(1)} &= X_2^T R_2^T R_2 X_2 \end{aligned} \quad (19)$$

So that

$$\begin{aligned} B_{11}^{(1)} &= \left(A_{11}^{(1)} - A_{12}^{(1)}(A_{22}^{(1)})^{-1}(A_{12}^{(1)})^T \right)^{-1} \\ &= \left((X_1^T X_1) - (X_1^T R_2 X_2)(X_2^T R_2^T R_2 X_2)^{-1}(X_1^T R_2 X_2)^T \right)^{-1} \\ &= \left((X_1^T X_1) - (X_1^T X_2)(X_2^T X_2)^{-1}(X_1^T X_2)^T \right)^{-1} \\ &= (A_{11} - A_{12}A_{22}^{-1}A_{12}^T)^{-1} \\ &= B_{11} \end{aligned} \quad (20)$$

This shows that the part of the usual ACOV to do with $\hat{\beta}_1$ can be estimated correctly, and therefore the standard errors are available: $\text{ACOV}(\hat{\beta}_j) = \sigma^2 B_{jj}$. Moreover,

$$\begin{aligned} B_{21}^{(1)} &= -(A_{22}^{(1)})^{-1}(A_{12}^{(1)})^T B_{11} \\ &= -(R_2^T R_2)^{-1} R_2 B_{21} \end{aligned} \quad (21)$$

so that

$$\begin{aligned} [(Z^T Z)^{-1} Z^T y]_{p_1} &= B_{11} X_1^T y - (R_2^T R_2)^{-1} R_2^T R_2 B_{21}^T X_2^T y \\ &= B_{11} X_1^T y - B_{21}^T X_2^T y \\ &= \hat{\beta}_1 \end{aligned} \quad (22)$$

This shows that the exact same estimates are obtained for $\hat{\beta}_1$. The same proof can be given for Bob and $\hat{\beta}_2$.

Note further that:

1. Alice cannot get $\hat{\beta}_2$ right because R_2 does not drop out in the other's part of the vector
2. We cannot get the ACOV of $(\hat{\beta}_1, \hat{\beta}_2)$ for this same reason

B MSE of rank-R data approximation

From Equation 16 we can create the following approximation:

$$\begin{aligned}\hat{\mathbf{Y}}_a &= \mathbf{X}_a \hat{\mathbf{B}}_a \\ \hat{\mathbf{X}}_a &= \hat{\mathbf{Y}}_a \hat{\mathbf{B}}_a^+\end{aligned}\tag{23}$$

where $\hat{\mathbf{Y}}_a \in \mathbb{R}^{N \times R}$, $\mathbf{B}_a \in \mathbb{R}^{P \times R}$, and $\mathbf{X}_a \in \mathbb{R}^{N \times P}$ and all matrices are full rank. For simplicity, but without loss of generality, we assume here that the variance of all the features in \mathbf{X}_a is the same, σ_a^2 , and these features are uncorrelated.

The relation between P , R , and the accuracy of the approximation $\hat{\mathbf{X}}_a$ is as follows: as $R \rightarrow P$, the MSE improves linearly, with perfect approximation being achieved when $R = P$. As mentioned in-text, when $P = 1$, sharing one set of parameters ($R = 1$) means the data can be recovered completely. Empirical simulations show that the relation between R , P , and expected mean square error of approximation is $\text{MSE} = \sigma_a^2(1 - R/P)$, where σ_a^2 is the variance of the features in \mathbf{X}_a (see Figure 8).

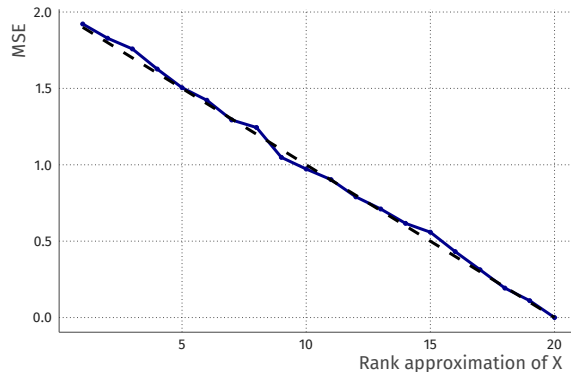


Figure 8: Mean square error (MSE) of the approximation of the data \mathbf{X}_a at Alice by Bob if $\hat{\mathbf{B}}_a$ is known. \mathbf{X}_a was simulated as having $P = 20$ uncorrelated features with variance $\sigma_a^2 = 2$. Note that the approximation linearly improves as the rank of $\hat{\mathbf{B}}_a$ increases, with a perfect approximation reached when $R = P$. Dashed line indicates expected MSE, using the formula $E[\text{MSE}] = \sigma_a^2(1 - R/P)$.

Phrasing the above in terms of information sharing and privacy preservation: in sharing R sets of parameter estimates $\hat{\beta}_a^{(r)}$ with their associated predictions $\hat{\mathbf{y}}_a^{(r)}$, Alice reveals a proportion of at least R/P of variance in the data. This proportion is a lower bound: in case there are correlations among the features of Alice, this proportion increases. When $R = P$ the data of Alice can be reconstructed by Bob. When either of a pair $(\hat{\beta}_a^{(r)}, \hat{\mathbf{y}}_a^{(r)})$ are shared but not the other, no information is revealed.