

Private Regression using Block Coordinate Descent

Erik-Jan van Kesteren

Utrecht University

2019-03-15

Outline

Regression models and vertical partitioning

Coordinate descent

The privreg package

Regression models and vertical partitioning

Basic regression model

Basic setup:

$$\begin{aligned} \mathbf{y} &= \mathbf{x} \cdot \boldsymbol{\beta} + \boldsymbol{\epsilon}, & \boldsymbol{\epsilon} &= \mathbf{y} - \mathbf{x} \cdot \boldsymbol{\beta} \\ \hat{\mathbf{y}} &= \mathbf{x} \cdot \hat{\boldsymbol{\beta}}, & \mathbf{y}_{res} &= \mathbf{y} - \mathbf{x} \cdot \hat{\boldsymbol{\beta}} \end{aligned}$$

ML / LS estimate

$$\hat{\boldsymbol{\beta}} = \frac{COV(\mathbf{x}, \mathbf{y})}{VAR(\mathbf{x})}$$

If \mathbf{X} and \mathbf{y} are centered:

$$\hat{\boldsymbol{\beta}} = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle}$$

Basic regression model with P predictors

Basic setup:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

ML / LS estimates:

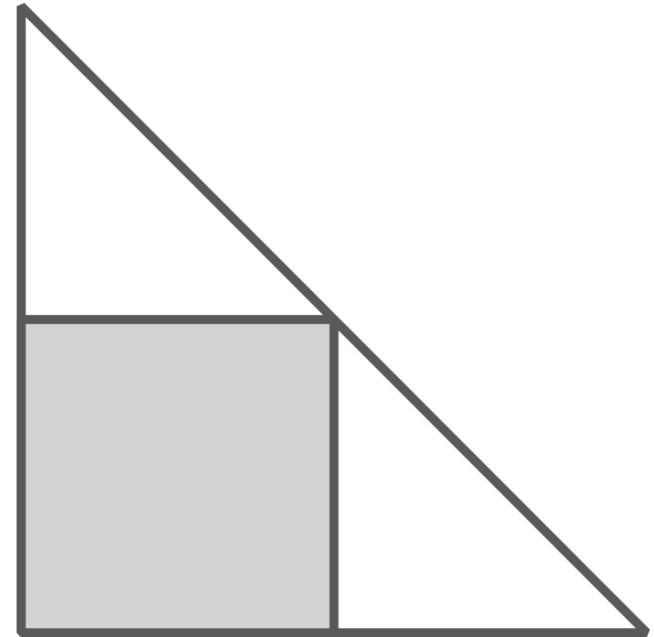
$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

If \mathbf{X} and \mathbf{y} are centered:

$$\mathbf{X}'\mathbf{X} = \text{COV}(\mathbf{X}), \quad \mathbf{X}'\mathbf{y} = \text{COV}(\mathbf{X}, \mathbf{y})$$

Regression with vertically partitioned data

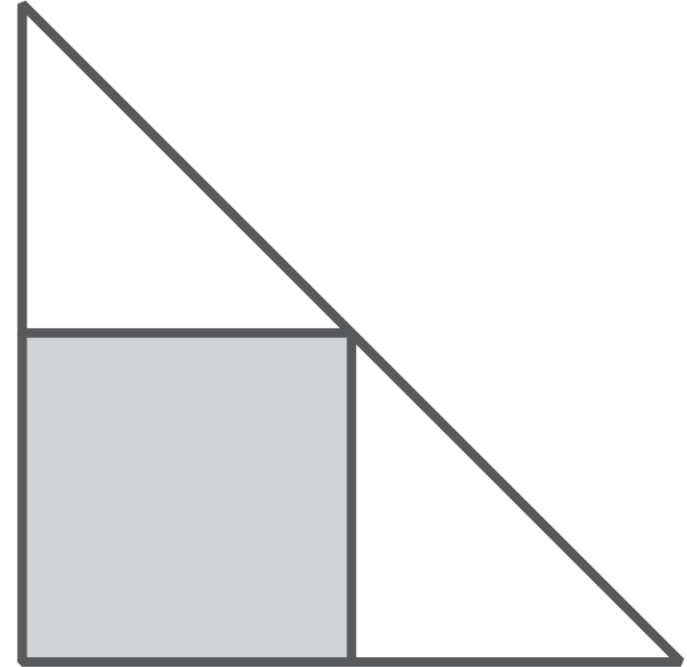
- Alice and Bob are two institutions
- Alice and Bob cannot share their predictors
- Two sets of non-overlapping predictors X_a and X_b
- $P_a \times P_b$ missing covariance matrix elements



Privacy note

- Alice is not allowed to know the values of \mathbf{X}_b and vice versa
- We can securely create $COV(\mathbf{X})$ (Karr et al., 2009; Hall et al., 2011)
- But sharing $COV(\mathbf{X})$ means Alice can predict $\mathbf{x}_b \in \mathbf{X}_b$ with some certainty:

$$\hat{\mathbf{x}}_b = \mathbf{X}_a (\mathbf{X}'_a \mathbf{X}_a)^{-1} COV(\mathbf{X}_a, \mathbf{x}_b)$$



Coordinate descent

Coordinate descent for multiple regression

Marginal parameter:

$$\hat{\beta}_1^m = \langle \mathbf{x}_1, \mathbf{y} \rangle / \langle \mathbf{x}_1, \mathbf{x}_1 \rangle$$

Conditional parameter:

$$\hat{\beta}_1 = \langle \mathbf{x}_1, \mathbf{y}_{res}^1 \rangle / \langle \mathbf{x}_1, \mathbf{x}_1 \rangle, \quad \mathbf{y}_{res}^1 = \mathbf{y} - \underbrace{\mathbf{X}^1 \hat{\beta}^1}_{\text{excl } \mathbf{x}_1}$$

Because by definition: $\mathbf{y}_{res} = \epsilon \perp \mathbf{X}$

Coordinate descent for multiple regression

```
← → | ↶ | 📁 Source on Save | 🔍 ✨ | 📄 | ➡️ | ↺ ➡️ | 📄 Source | ☰
1 X ← matrix(rnorm(100), 20)
2 y ← X %*% runif(5, -1, 1) + rnorm(20)
3 beta_hat ← solve(crossprod(X), crossprod(X, y))
4
5 y_res_1 ← y - X[, -1] %*% beta_hat[-1]
6 beta_1 ← crossprod(X[, 1], y_res_1) / crossprod(X[, 1])
7
8 beta_1 # Same
9 beta_hat[1] # Same
10 cov(X[, 1], y) / var(X[, 1]) # Different
11
12
```

Coordinate descent for multiple regression

1. Initialise $\hat{\beta} \leftarrow \mathbf{0}_P$
2. For $p \in P$:
 1. $\mathbf{y}_{res}^p = \mathbf{y} - \mathbf{X}^p \hat{\beta}^p$
 2. $\hat{\beta}_p \leftarrow \langle \mathbf{x}_p, \mathbf{y}_{res}^p \rangle / \langle \mathbf{x}_p, \mathbf{x}_p \rangle$
3. Repeat until convergence

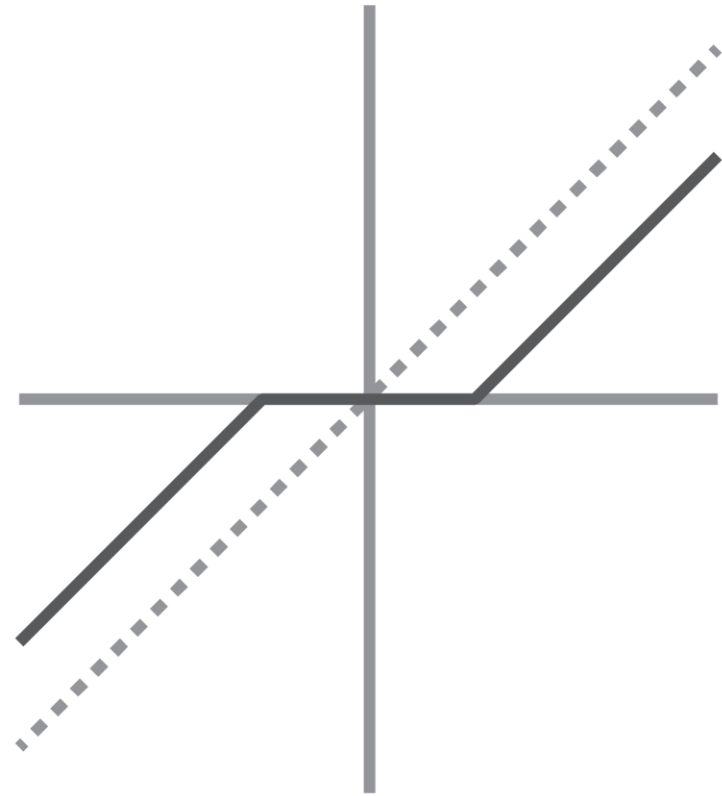
Example implementation in R:

<https://gist.github.com/vankesteren/88621f8a0d532083691e6f2a505dd5c7>

Coordinate descent for multiple regression

Soft-thresholding $\hat{\beta}_i$ in each step of the algorithm leads to the LASSO estimate (Hastie et al., 2015):

1. Initialise $\hat{\beta} \leftarrow \mathbf{0}_P$
2. For $p \in P$:
 1. $\mathbf{y}_{res}^p = \mathbf{y} - \mathbf{X}^p \hat{\beta}^p$
 2. $\hat{\beta}_p \leftarrow S_\lambda(\langle \mathbf{x}_p, \mathbf{y}_{res}^p \rangle / \langle \mathbf{x}_p, \mathbf{x}_p \rangle)$
3. Repeat until convergence

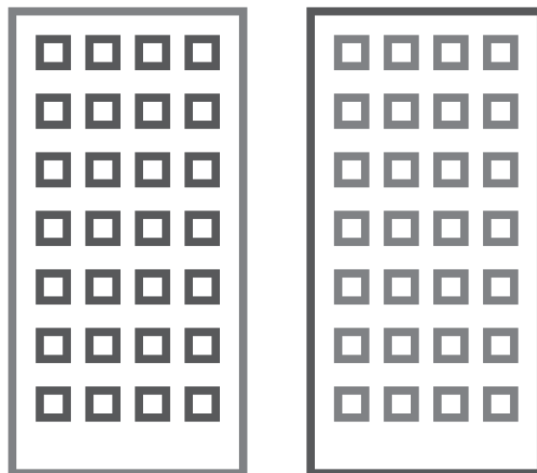


Block coordinate descent for regression

Assume we have K blocks

1. Initialise $\hat{\boldsymbol{\beta}} \leftarrow \mathbf{0}_p$
2. For $k \in K$:
 1. $\mathbf{y}_{res}^k = \mathbf{y} - \mathbf{X}^k \hat{\boldsymbol{\beta}}^k$
 2. $\hat{\boldsymbol{\beta}}_k \leftarrow (\mathbf{X}_k' \mathbf{X}_k)^{-1} \mathbf{X}_k \mathbf{y}_{res}^k$
3. Repeat until convergence

privreg



Set $\epsilon^b = \mathbf{y}$

Alice

$$\hat{\beta}_a = (\mathbf{X}'_a \mathbf{X}_a)^{-1} \mathbf{X}_a \mathbf{y}_{res}^b$$

$$\mathbf{y}_{res}^a = \mathbf{y} - \mathbf{X}_a \hat{\beta}_a$$

Send \mathbf{y}_{res}^a to Bob

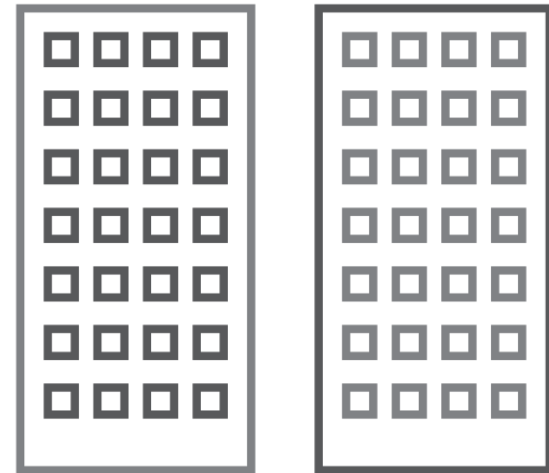
Bob

$$\hat{\beta}_b = (\mathbf{X}'_b \mathbf{X}_b)^{-1} \mathbf{X}_b \mathbf{y}_{res}^a$$

$$\mathbf{y}_{res}^b = \mathbf{y} - \mathbf{X}_b \hat{\beta}_b$$

Send \mathbf{y}_{res}^b to Alice

privreg



Result

Coefficients:

$$(\hat{\beta}_a, \hat{\beta}_b) = \hat{\beta} = (X'X)^{-1}X'y$$

Prediction:

$$\hat{y} = \hat{y}_a + \hat{y}_b = X_a\hat{\beta}_a + X_b\hat{\beta}_b$$



Privacy-preserving regression

Edit

Manage topics

15 commits 1 branch 0 releases 1 contributor View license

Branch: master New pull request Create new file Upload files Find File Clone or download

| | | |
|----------------|---|-----------------------------------|
| vankesteren | Update README.md | Latest commit 917053f an hour ago |
| R | added modeling with formulas | 16 days ago |
| man | Skip automated test, set up travis | 18 days ago |
| tests | added modeling with formulas | 16 days ago |
| .Rbuildignore | update travis, add .yml to build ignore | 18 days ago |
| .gitattributes | Initial commit | 29 days ago |
| .gitignore | Initial commit | 29 days ago |
| .travis.yml | update travis, add .yml to build ignore | 18 days ago |
| DESCRIPTION | Add message encryption via openssl, clean up code | 18 days ago |
| LICENSE | change to port 443 | 21 days ago |
| LICENSE.md | Initial commit | 29 days ago |
| NAMESPACE | Add message encryption via openssl, clean up code | 18 days ago |
| README.md | Update README.md | an hour ago |
| privreg.Rproj | Initial commit | 29 days ago |

README.md

privreg : Private regression using block coordinate descent

build passing

install

Time for a tryout

github.com/vankesteren/privreg