

Statistical Gradient Boosting of Generalised Additive Model with Ridge Regression on GPU

Antoni Sieminski¹, Johnny MyungWon Lee¹

¹The University of Edinburgh
School of Mathematics

7th April 2023



THE UNIVERSITY of EDINBURGH
School of Mathematics

Outline

- 1 Preliminary
- 2 Algorithm
- 3 References

Generalised Additive Model

The Gaussian Generalised Additive Model (GAM) is defined as:

$$\begin{aligned}y &\sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2) \\ \boldsymbol{\mu} &= \beta_0 + f_1(\mathbf{x}_1) + \dots + f_p(\mathbf{x}_p)\end{aligned}$$

which allows us to easily evaluate the effects of each variable \mathbf{x}_j on the response \mathbf{y} .

Generalised Additive Model

The Gaussian Generalised Additive Model (GAM) is defined as:

$$y \sim \mathcal{N}(\mu, \sigma^2)$$
$$\mu = \beta_0 + f_1(x_1) + \dots + f_p(x_p)$$

which allows us to easily evaluate the effects of each variable x_j on the response y .

For an unknown $f(x_j)$, it is usually sensible to assume it's a smooth function, which can be represented as a linear combination of rank-reduced k basis functions $b_i(\cdot)$ multiplied by their fitted coefficients γ_i :

$$f(x_j) = \sum_{i=1}^k b_i(x_j) \gamma_i = \mathbf{B}_j \boldsymbol{\gamma}_j$$

GAM – smooth functions

In order to enforce smoothness on $f_j(\mathbf{x}_j)$, the term $\lambda_j \boldsymbol{\gamma}_j^T \mathbf{S}_j \boldsymbol{\gamma}_j$ is typically added to the f_j 's loss function, where \mathbf{S}_j is the smoothing matrix and λ_j is the penalty term.

$$L = ||\mathbf{y} - \hat{\boldsymbol{\mu}}||^2 + \sum_{j=1}^p \lambda_j \boldsymbol{\gamma}_j^T \mathbf{S}_j \boldsymbol{\gamma}_j$$

GAM – smooth functions

In order to enforce smoothness on $f_j(x_j)$, the term $\lambda_j \boldsymbol{\gamma}_j^T \mathbf{S}_j \boldsymbol{\gamma}_j$ is typically added to the f_j 's loss function, where \mathbf{S}_j is the smoothing matrix and λ_j is the penalty term.

$$L = ||\mathbf{y} - \hat{\boldsymbol{\mu}}||^2 + \sum_{j=1}^P \lambda_j \boldsymbol{\gamma}_j^T \mathbf{S}_j \boldsymbol{\gamma}_j$$

For example, in P-Splines, $\mathbf{S} = \mathbf{D}^T \mathbf{D}$, which enforces smoothness by penalising the changes in differences of adjacent coefficients:

$$\boldsymbol{\gamma}^T \mathbf{S} \boldsymbol{\gamma} = \sum_{i=1}^{k-1} (\gamma_{i-1} - 2\gamma_i + \gamma_{i+1})^2$$

Gradient Boosting

- Gradient boosting is a fitting algorithm, which sequentially fits weak estimators (a.k.a. base learners) to the gradient of the loss function $\mathbf{g}^{[m]}$.
- The resulting model can avoid including unnecessary base learners and shrink estimates from the included ones to reduce the generalisation error of the model.
- In GAMs, the base learners are typically penalised regression models – each with a design matrix \mathbf{B}_j , smoothing matrix \mathbf{S}_j , and the smoothing parameter λ_j .

Gradient Boosting

- Gradient boosting is a fitting algorithm, which sequentially fits weak estimators (a.k.a. base learners) to the gradient of the loss function $g^{[m]}$.
- The resulting model can avoid including unnecessary base learners and shrink estimates from the included ones to reduce the generalisation error of the model.
- In GAMs, the base learners are typically penalised regression models – each with a design matrix \mathbf{B}_j , smoothing matrix \mathbf{S}_j , and the smoothing parameter λ_j .

Gradient Boosting

- Gradient boosting is a fitting algorithm, which sequentially fits weak estimators (a.k.a. base learners) to the gradient of the loss function $g^{[m]}$.
- The resulting model can avoid including unnecessary base learners and shrink estimates from the included ones to reduce the generalisation error of the model.
- In GAMs, the base learners are typically penalised regression models – each with a design matrix \mathbf{B}_j , smoothing matrix \mathbf{S}_j , and the smoothing parameter λ_j .

Smooth functions as base learners

At each iteration m , all K base learners are fitted to the residuals of the model $\mathbf{g}^{[m]}$:

$$\boldsymbol{\gamma}_j^{[m]} = (\mathbf{B}_j^T \mathbf{B}_j + \lambda_j \mathbf{S}_j)^{-1} \mathbf{B}_j^T \mathbf{g}^{[m]}$$

and the one with the lowest residual sum of squares is selected.

$$\text{learner index} = \arg \min_j \|\mathbf{g}^{[m]} - \mathbf{B}_j \boldsymbol{\gamma}_j^{(m)}\|^2$$

The smoothing parameter λ_j is typically selected before the fitting process, such that the learner has some small number of effective degrees of freedom $\text{Tr}(\mathbf{H}_j)$.

Smooth functions as base learners

At each iteration m , all K base learners are fitted to the residuals of the model $\mathbf{g}^{[m]}$:

$$\boldsymbol{\gamma}_j^{[m]} = (\mathbf{B}_j^T \mathbf{B}_j + \lambda_j \mathbf{S}_j)^{-1} \mathbf{B}_j^T \mathbf{g}^{[m]}$$

and the one with the lowest residual sum of squares is selected.

$$\text{learner index} = \arg \min_j \|\mathbf{g}^{[m]} - \mathbf{B}_j \boldsymbol{\gamma}_j^{(m)}\|^2$$

The smoothing parameter λ_j is typically selected before the fitting process, such that the learner has some small number of effective degrees of freedom $\text{Tr}(\mathbf{H}_j)$.

Smooth functions as base learners

At each iteration m , all K base learners are fitted to the residuals of the model $\mathbf{g}^{[m]}$:

$$\boldsymbol{\gamma}_j^{[m]} = (\mathbf{B}_j^T \mathbf{B}_j + \lambda_j \mathbf{S}_j)^{-1} \mathbf{B}_j^T \mathbf{g}^{[m]}$$

and the one with the lowest residual sum of squares is selected.

$$\text{learner index} = \arg \min_j \|\mathbf{g}^{[m]} - \mathbf{B}_j \boldsymbol{\gamma}_j^{(m)}\|^2$$

The smoothing parameter λ_j is typically selected before the fitting process, such that the learner has some small number of effective degrees of freedom $\text{Tr}(\mathbf{H}_j)$.

Computational considerations

In gradient boosting, finding λ for each learner and the associated $(\mathbf{B}^T \mathbf{B} + \lambda \mathbf{S})^{-1}$ is relatively inexpensive as it can be done just once.

Operations that happen at each boosting iteration m are the main computational bottleneck and include the following matrix-vector multiplications:

- 1 $\mathbf{B}^T \mathbf{g}^{[m]}, O(np)$
- 2 $(\mathbf{B}^T \mathbf{B} + \lambda \mathbf{S})^{-1} (\mathbf{B}^T \mathbf{g}^{[m]}), O(p^2)$
- 3 $\mathbf{B} \boldsymbol{\gamma}^{[m]}, O(np)$

These can be sped up significantly by parallel computing.

Computational considerations

In gradient boosting, finding λ for each learner and the associated $(\mathbf{B}^T \mathbf{B} + \lambda \mathbf{S})^{-1}$ is relatively inexpensive as it can be done just once.

Operations that happen at each boosting iteration m are the main computational bottleneck and include the following matrix-vector multiplications:

- ① $\mathbf{B}^T \mathbf{g}^{[m]}, O(np)$
- ② $(\mathbf{B}^T \mathbf{B} + \lambda \mathbf{S})^{-1} (\mathbf{B}^T \mathbf{g}^{[m]}), O(p^2)$
- ③ $\mathbf{B} \boldsymbol{\gamma}^{[m]}, O(np)$

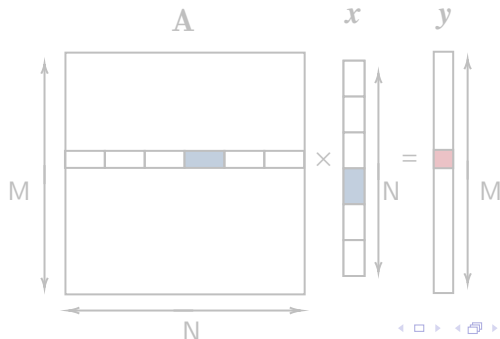
These can be sped up significantly by parallel computing.

B-Splines as Base Learners

Matrix Vector Multiplication

Algorithm Naive Implementation of Matrix Vector Multiplication

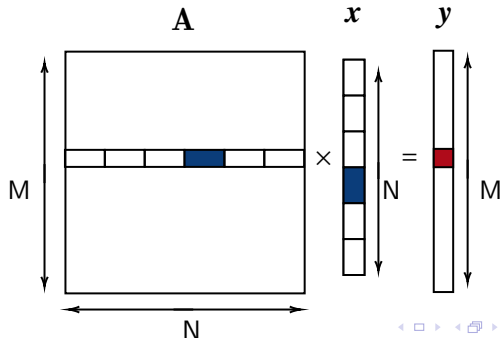
```
1: for  $i = 1$  to  $M$  do  
2:   for  $j = 1$  to  $N$  do  
3:      $y[i] += A[i][j] * x[j]$   
4:   end for  
5: end for
```



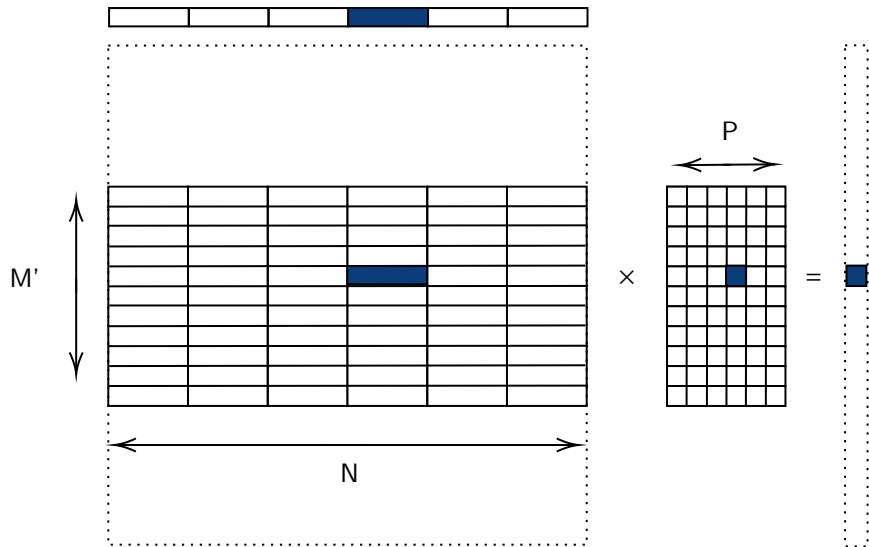
Matrix Vector Multiplication

Algorithm Naive Implementation of Matrix Vector Multiplication

```
1: for  $i = 1$  to  $M$  do  
2:   for  $j = 1$  to  $N$  do  
3:      $y[i] += A[i][j] * x[j]$   
4:   end for  
5: end for
```



P threads per dot product



Simulation

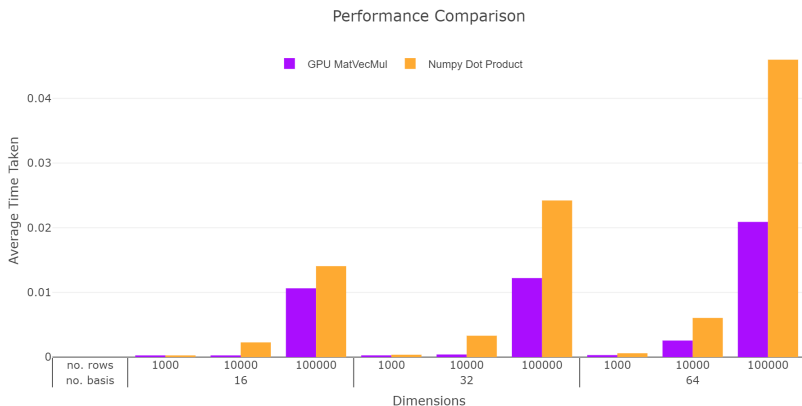


Figure: Comparison of Matrix Vector Multiplication between GPU and Numpy (CPU)

Reference

- Wang, Hansheng, and Chih-Ling Tsai. " *Tail index regression.*" Journal of the American Statistical Association 104.487 (2009): 1233-1240.

Decomposing smooth functions

The basis matrix \mathbf{B} can be reparametrised as a sum of smooth and unpenalised components:

$$\begin{aligned}\mathbf{B}\boldsymbol{\gamma} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} \\ \boldsymbol{\gamma}^T \mathbf{S} \boldsymbol{\gamma} &= \mathbf{b}^T \mathbf{I} \mathbf{b}\end{aligned}$$

which means we can split functions $f_j(\mathbf{x}_j)$ into separate base learners for unpenalised ($\mathbf{X}\boldsymbol{\beta}$) and smooth components ($\mathbf{Z}\mathbf{b}$).

This allows $f_j(\mathbf{x}_j)$ to take the form of a linear or smooth function independently of the other improving model interpretability and performance.