

# STA141C Spring 2020 Final project

V. Budilov

6/10/2020

STA141C Spring 2020 Final project

## Planned activities toward completing this project.

The activities are suggested and may be changed/dropped/added during the project.

1. Implement parallelization. In the current implementation, only one CPU is used in the algorithm. Make it possible to use more than one CPUs. Note that you should let users to decide if they want to use parallelization.
2. Data loading functionality implementation. Allow users to specify file names to run the model rather than loading the whole data in the main process.
3. Implementing C++ substitute functions or creating new functionalities using C++. Functions are written in pure R, it is possible, for example, to convert the function `lm1` to c++ code. You might need look at how ReppArmadillo's `fastLm.R` and `fastLm.cpp`. (Spoiler, it is not easy, but if you insist, here is a some slides about it: [https://scholar.princeton.edu/sites/default/files/q-aps/files/slides\\_day4\\_am.pdf](https://scholar.princeton.edu/sites/default/files/q-aps/files/slides_day4_am.pdf))
4. Creating documentation and implementing tests. Write tests and documentations
5. Implementing more models More models? Logistic regression? GLM?

## Implemented tasks.

### Parallelization

1. Parallelization has been implemented. Function `blblm1()` is an analog of the original function `blblm()`. The difference of `bklblm1()` from `blblm()` is that the new function uses parallelization. `bklblm1()` was used for the development and testing. The final implementation of parallelization has been included into `blblm()`.

### Generalized Linear Model

2. Generalized linear model has been added to the package. Original function signature

```
blblm(formula, data, m = 10, B = 5000)
```

has been changed to

```
blblm(formula, data, m = 10, B = 5000, ifgeneral=FALSE, family = "gaussian", parallel=FALSE).
```

Logical `ifgeneral` parameter tells if we need to use generalized linear model, the default value is `FALSE`. If the generalized linear model required, we need to change `ifgeneral` to `TRUE`. The default family is `ste` to be “gaussian”. When parallelization is required, the value of the parameter *parallel* has to be set to `TRUE`.

**Test cases** Tests have been created. Check `test-fit.R`

**Documentation** Documentation has been elaborated.

## Results and Discussion

```
{r}
#library(nycflights13)
#flights1<-flights[1:2000,]
fit <- blblm(mpg ~ wt * hp, data = mtcars, m = 3, B = 100)
coef(fit)
fit <- blblm(mpg ~ wt * hp, data = mtcars, m = 3, B = 100, parallel = TRUE)
coef(fit)
```

Comparison between non-parallel and parallel calculations has been done using various data lengths, bootstrapping counts, data subdivisions, numbers of workers.

```
{r message=FALSE, warning=FALSE}
library(nycflights13)
flights1<-flights[1:100000,]
f1<-function(){blblm(arr_delay ~ distance + dep_delay, data = flights, m = 24, B = 1000)
  return (0)}
f2<-function(){blblm(arr_delay ~ distance + dep_delay, data = flights, m = 24, B = 1000,parallel=TRUE)
  return (0)}
result<-bench::mark(
  f1(),
  f2()
)
result
```

For six workers, m=4, and B=1000, the results are

```
## # A tibble: 2 x 6
##   expression      min    median `itr/sec` mem_alloc `gc/sec`
##   <bch:expr> <bch:tm> <bch:tm>      <dbl> <bch:byt>      <dbl>
## 1 f1()         4.76m    4.76m    0.00350   159GB    4.91
## 2 f2()         2.75m    2.75m    0.00606   213MB    0.00606
```

## Some examples of usage

### Default usage (corresponds to the original blblm() function)

Default usage with simple linear model, and no prallelization

```
fit <- blblm(mpg ~ wt * hp, data = mtcars, m = 3, B = 100)
coef(fit)
```

```
## (Intercept)          wt          hp      wt:hp
## 49.68852308 -7.95953045 -0.12681856  0.02823271
```

### Parallelization

Usage with simple linear model, and prallelization

```
fit2 <- blblm(mpg ~ wt * hp, data = mtcars, m = 3, B = 100, parallel=TRUE)
coef(fit2)
```

```
## (Intercept)          wt          hp      wt:hp
## 55.05346089 -9.99813705 -0.16840913  0.04284208
```

## Generalized linear model

Usage with generalized linear model without parallelization

```
fit <- blblm(mpg ~ wt * hp, data = mtcars, m = 4, B = 100, ifgeneral=TRUE)
coef(fit)
```

```
## (Intercept)          wt          hp          wt:hp
## 96.2790574 -21.1033360 -0.5041665  0.1326205
```

## Generalized linear model with parallelization

Usage with generalized linear model with parallelization

```
fit <- blblm(mpg ~ wt * hp, data = mtcars, m = 5, B = 100, ifgeneral=TRUE, family='quasi', parallel=TRUE)
coef(fit)
```

```
## (Intercept)          wt          hp          wt:hp
## 54.28827197 -8.48633098 -0.20193676  0.04338432
```

```
#fit2 <- blbgglm(mpg ~ wt * hp, data = mtcars, m = 3, B = 100)
#fit2$coef
#coef(fit2)
```