

# Distributed Statistical Modeling

author: Balasubramanian Narasimhan, Stanford University

## Introduction

We demonstrate the possibility of fitting statistical models stratified by sites in a manner that brings computation to the data that may be distributed across sites or more generally, partitioned in some manner. (For simplicity, we will call these partitions, sites.) The infrastructure consists of a single master process that issues queries to slave processes running at each of the sites. A query is merely a function call, more specifically a request to each site to evaluate a pre-defined function  $f(\beta)$  on the data at that site, for a given value of parameters  $\beta$ . The master process uses these queries to aggregate and execute an optimization algorithm resulting in a model fit, the results of which should be *indistinguishable from* those that might be obtained if all the data had been in a single place. Of course, this assumes a lossless serialization format, like Google Protocol Buffers for example, but we make do with JSON for now. Also, in comparisons below, we don't use exactly the same iterations as the survival package and so minor differences will be seen.

The advantages are many, chief among them the fact that no raw data needs to be shared between sites. The modeling entity, however, can make an unlimited number of queries of the sites, where each query is a request to compute a model-specific function for a specified value of parameters. This may pose a security concern that we ignore for now. However, it leads to some further interesting questions regarding what may be learned such computation.

We focus specifically on Cox regression models in this exercise.

Setup

A simple example

A Cox fit on the aggregated data

The model definition

The data for each site

Setting up the sites

Reproducing original aggregated analysis in a distributed fashion

A larger example

[The aggregated fit] (#larger-example-full)

**The distributed fit** (#larger-example-distributed)

Bone Marrow Transplant Example

[The aggregated fit] (#bmt-full)

## The distributed fit (#bmt-distributed)

Byar and Greene Prostate Cancer Data Example (4 strata)

### ### Setup

It must be noted that for users to be able to `knit` this document, or to run these examples in an R session, an `opencpu` server must be running with appropriate settings. On MacOS and Unix, this is done by designating an empty directory as workspace and adding the following lines to the `${HOME}/.Rprofile`.

```
library(distcomp)
distcompSetup(workspace="full_path_to_workspace_directory",
               ssl.verifyhost=FALSE, ssl.verifypeer=FALSE)
```

On windows, the same should be done in the `RHOME\etc\Rprofile.site` file. Then, an R session must be started with these settings in effect. Then, this document can be `knit`/rendered in that session.

### ## A simple example

We take a simple example from the `survival` package, the `ovarian` dataset.

```
library(knitr)
library(survival)
data(ovarian)
str(ovarian)
```

```
## 'data.frame':    26 obs. of  6 variables:
## $ futime   : num  59 115 156 421 431 448 464 475 477 563 ...
## $ fustat   : num  1 1 1 0 1 0 1 1 0 1 ...
## $ age      : num  72.3 74.5 66.5 53.4 50.3 ...
## $ resid.ds: num  2 2 2 2 2 1 2 2 2 1 ...
## $ rx       : num  1 1 1 2 1 1 2 2 1 2 ...
## $ ecog.ps  : num  1 1 2 1 1 2 2 2 1 2 ...
```

```
kable(ovarian)
```

futime	fustat	age	resid.ds	rx	ecog.ps
59	1	72.3315	2	1	1
115	1	74.4932	2	1	1
156	1	66.4658	2	1	2
421	0	53.3644	2	2	1
431	1	50.3397	2	1	1
448	0	56.4301	1	1	2
464	1	56.9370	2	2	2
475	1	59.8548	2	2	2
477	0	64.1753	2	1	1
563	1	55.1781	1	2	2
638	1	56.7562	1	1	2

futime	fustat	age	resid.ds	rx	ecog.ps
744	0	50.1096	1	2	1
769	0	59.6301	2	2	2
770	0	57.0521	2	2	1
803	0	39.2712	1	1	1
855	0	43.1233	1	1	2
1040	0	38.8932	2	1	2
1106	0	44.6000	1	1	1
1129	0	53.9068	1	2	1
1206	0	44.2055	2	2	1
1227	0	59.5890	1	2	2
268	1	74.5041	2	1	2
329	1	43.1370	2	1	1
353	1	63.2192	1	2	2
365	1	64.4247	2	2	1
377	0	58.3096	1	2	1

### A Cox fit on the aggregated data

A simple Cox model fit estimates the effect of age on survival, stratified by drug.

```
cp <- coxph(Surv(futime, fustat) ~ age + strata(rx), data=ovarian, ties="breslow")
print(cp)
```

```
## Call:
## coxph(formula = Surv(futime, fustat) ~ age + strata(rx), data = ovarian,
##       ties = "breslow")
##
##
##      coef exp(coef) se(coef)      z      p
## age 0.137      1.15   0.0474  2.9 0.0038
##
## Likelihood ratio test=12.7 on 1 df, p=0.000368 n= 26, number of events= 12
```

The above shows the initial and final log likelihood values at 0 and the estimated coefficient respectively and the actual estimated coefficient in the last line.

For our setting, we can pretend that this data set is actually from two sites, one containing the control or placebo group ( $rx = 1$ ) and the other containing the drug group ( $rx = 2$ ).

### The model definition

We first need to define the computation. The available computations can be listed:

```
print(availableComputations())
```

```
## $StratifiedCoxModel
## [1] "Stratified Cox Model"
##
## $RankKSVD
## [1] "Rank K SVD"
```

So, we can define the ovarian data computation as follows.

```
ovarianDef <- list(compType = names(availableComputations())[1],
  formula = "Surv(futime, fustat) ~ age",
  defnId = "Ovarian")
```

### The data for each site

We split the ovarian data into two sites as indicated earlier.

```
siteData <- with(ovarian, split(x=ovarian, f=rx))
```

### Setting up the sites

We can now set up each site with its own data.

```
nSites <- length(siteData)
siteNames <- sapply(seq.int(nSites), function(i) paste("site", i, sep=""))
siteURLs <- lapply(seq.int(nSites), function(i) opencpu$url())
names(siteData) <- names(siteURLs) <- siteNames
```

By default, on each site, data for a computation is stored under the name `data.rds` and the definition itself is stored under the name `defn.rds`. If the sites are physically separate, then everything proceeds smoothly. However, here, in our case, we are using the same `opencpu` server for simulating both sites. We therefore have to save the files under different names, just for this experiment, say `site1.rds` and `site2.rds` for this example. We do that below, by providing the `siteDataFiles` as an argument to the `Map` call.

```
siteDataFiles <- lapply(seq.int(nSites), function(i) paste("site", i, ".rds", sep=""))

ok <- Map(uploadNewComputation, siteURLs,
  lapply(seq.int(nSites), function(i) ovarianDef),
  siteData,
  siteDataFiles)

stopifnot(all(ok))
```

```
## Warning in all(ok): coercing argument of type 'list' to logical
```

### Reproducing original aggregated analysis in a distributed fashion

We are now ready to reproduce the original aggregated analysis. We first create a master object, taking care to specify that we are using a local server to simulate several sites.

```
master <- coxMaster$new(defnId = ovarianDef$defnId, formula=ovarianDef$formula,
  localServer=TRUE)
```

We then add the slave sites, once again specifying the site data file names.

```
for (i in seq.int(nSites)) {
  master$addSite(siteNames[i], siteURLs[[i]], dataFileName=siteDataFiles[[i]])
}
```

And we now maximize the partial likelihood, by calling the `run` method of the master.

```
result <- master$run()
```

We then print the summary.

```
master$summary()
```

```
##           coef exp(coef)      se(coef)      z      p
## 1  0.908385976 2.4803160 0.1972455530  4.6053559 4.117610e-06
## 2  0.005434256 1.0054490 0.0141220463  0.3848065 7.003808e-01
## 3 -0.001704550 0.9982969 0.0128467057 -0.1326838 8.944434e-01
## 4 -1.033862530 0.3556307 0.2580041510 -4.0071546 6.145463e-05
## 5 -0.337598317 0.7134818 0.2602298297 -1.2973083 1.945251e-01
## 6  0.002854690 1.0028588 0.0006703559  4.2584697 2.058311e-05
```

As we can see, the results we get from the distributed analysis are the same as we got for the original aggregated analysis. We print them separately here for comparison.

## A larger example

We turn to a larger the example from Therneau and Grambsch using the `pbcc` data where the stratifying variable is `ascites`.

### The aggregated fit

```
data(pbc)
pbccCox <- coxph(Surv(time, status==2) ~ age + edema + log(bili) +
  log(protime) + log(albumin) + strata(ascites), data=pbcc, ties="breslow")
print(pbccCox)
```

## Call:

```
## coxph(formula = Surv(time, status == 2) ~ age + edema + log(bili) +
##      log(protime) + log(albumin) + strata(ascites), data = pbcc,
##      ties = "breslow")
##
```

##

##

```
##           coef exp(coef) se(coef)      z      p
## age          0.0314    1.0318 0.00907   3.45 0.00055
## edema         0.5993    1.8209 0.32127   1.87 0.06200
## log(bili)     0.8663    2.3780 0.10066   8.61 0.00000
## log(protime)  3.0341   20.7815 1.03884   2.92 0.00350
## log(albumin) -2.9662    0.0515 0.78177  -3.79 0.00015
##
```

```
## Likelihood ratio test=146 on 5 df, p=0 n= 312, number of events= 125
## (106 observations deleted due to missingness)
```

### The distributed fit

We split the data using `ascites` and proceed the usual way as shown above

```

pbcDef <- list(compType = names(availableComputations())[1],
              formula = paste("Surv(time, status==2) ~ age + edema +",
                              "log(bili) + log(protime) + log(albumin)"),
              defnId = "pbc")
siteData <- with(pbc, split(x=pbc, f=ascites))
nSites <- length(siteData)
siteNames <- sapply(seq.int(nSites), function(i) paste("site", i, sep=""))
siteURLs <- lapply(seq.int(nSites), function(i) opencpu$url())
names(siteData) <- names(siteURLs) <- siteNames
siteDataFiles <- lapply(seq.int(nSites), function(i) paste("site", i, ".rds", sep=""))
ok <- Map(uploadNewComputation, siteURLs,
          lapply(seq.int(nSites), function(i) pbcDef),
          siteData,
          siteDataFiles)

stopifnot(all(ok))

```

## Warning in all(ok): coercing argument of type 'list' to logical

```

master <- coxMaster$new(defnId = pbcDef$defnId, formula=pbcDef$formula,
                        localServer=TRUE)
for (i in seq.int(nSites)) {
  master$addSite(siteNames[i], siteURLs[[i]], dataFileName=siteDataFiles[[i]])
}

```

```
result <- master$run()
```

We then print the summary.

```
kable(master$summary())
```

coef	exp(coef)	se(coef)	z	p
0.0310247	1.0315110	0.0090692	3.420887	0.0006242
0.6019922	1.8257525	0.3205949	1.877735	0.0604174
0.8682667	2.3827773	0.1006068	8.630302	0.0000000
3.0276949	20.6495786	1.0393738	2.912999	0.0035798
-2.9765945	0.0509661	0.7809580	-3.811465	0.0001381
The results should be comparable to the aggregated fit above.				

## Bone Marrow Transplant Example

This uses the bmt data from Klein and Moschberger. Some variable renaming, first.

```

##
## BMT data
##

```

```

library(KMsurv)
data(bmt)
bmt$tnodis <- bmt$t2 ## time to disease relapse/death
bmt$inodis <- bmt$d3 ## disease relapse/death indicator
bmt$tplate <- bmt$tp ## time to platelet recovery
bmt$iplate <- bmt$dp ## platelet recovery
bmt$agep <- bmt$z1 ## age of patient in years
bmt$aged <- bmt$z2 ## age of donor in years
bmt$fab <- bmt$z8 ## fab grade 4 or 5 + AML
bmt$imtx <- bmt$z10 ## MTX used
bmt <- bmt[order(bmt$tnodis), ] ## order by time to disease relapse/death
bmt <- cbind(1:nrow(bmt)[1], bmt)
names(bmt)[1] <- "id"

##
####
##
bmt$agep.c = bmt$agep - 28
bmt$aged.c = bmt$aged - 28

bmt$imtx <- factor(bmt$imtx)

```

### The aggregated fit

```

bmt.cph <- coxph(formula = Surv(tnodis, inodis) ~ fab + agep.c * aged.c +
                 factor(group) + strata(imtx), data = bmt, ties="breslow")

print(bmt.cph)

```

```

## Call:
## coxph(formula = Surv(tnodis, inodis) ~ fab + agep.c * aged.c +
##       factor(group) + strata(imtx), data = bmt, ties = "breslow")
##
##
##               coef exp(coef) se(coef)      z      p
## fab             0.90780    2.479  0.27899  3.2539 0.0011
## agep.c          0.00551    1.006  0.01997  0.2759 0.7800
## aged.c         -0.00164    0.998  0.01816 -0.0901 0.9300
## factor(group)2 -1.03389    0.356  0.36472 -2.8348 0.0046
## factor(group)3 -0.33909    0.712  0.36784 -0.9218 0.3600
## agep.c:aged.c   0.00285    1.003  0.00095  2.9950 0.0027
##
## Likelihood ratio test=31 on 6 df, p=2.5e-05 n= 137, number of events= 83

```

## The distributed fit

We'll use imtx for splitting data into sites.

```

bmtDef <- list(compType = names(availableComputations())[1],
               formula = paste("Surv(tnodis, inodis) ~ fab + agep.c * aged.c + factor(group)"),
               defnId = "bmt")
siteData <- with(bmt, split(x=bmt, f=imtx))

```

```

nSites <- length(siteData)
siteNames <- sapply(seq.int(nSites), function(i) paste("site", i, sep=""))
siteURLs <- lapply(seq.int(nSites), function(i) opencpu$url())
names(siteData) <- names(siteURLs) <- siteNames
siteDataFiles <- lapply(seq.int(nSites), function(i) paste("site", i, ".rds", sep=""))
ok <- Map(uploadNewComputation, siteURLs,
          lapply(seq.int(nSites), function(i) bmtDef),
          siteData,
          siteDataFiles)

stopifnot(all(ok))

```

## Warning in all(ok): coercing argument of type 'list' to logical

```

master <- coxMaster$new(defnId = bmtDef$defnId, formula=bmtDef$formula,
                        localServer=TRUE)
for (i in seq.int(nSites)) {
  master$addSite(siteNames[i], siteURLs[[i]], dataFileName=siteDataFiles[[i]])
}

```

```
result <- master$run()
```

We then print the summary.

```
kable(master$summary())
```

coef	exp(coef)	se(coef)	z	p
0.9083860	2.4803160	0.2789473	3.2564784	0.0011280
0.0054343	1.0054490	0.0199716	0.2720993	0.7855457
-0.0017046	0.9982969	0.0181680	-0.0938216	0.9252508
-1.0338625	0.3556307	0.3648730	-2.8334862	0.0046043
-0.3375983	0.7134818	0.3680206	-0.9173355	0.3589669
0.0028547	1.0028588	0.0009480	3.0111928	0.0026022

## Byar and Greene Prostate Cancer Data Example

This example is the largest of them all and also has four strata rather than 2.

```
prostate <- readRDS("prostate.RDS")
```

### The aggregated fit

```

pcph <- coxph(Surv(dtime, status) ~ stage + strata(rx) + age + wt + pf + hx +
              sbp + dbp + ekg + hg + sz + sg + ap + bm, data=prostate)
print(pcph)

```

## Call:



```
## coxph(formula = Surv(dtime, status) ~ stage + strata(rx) + age +
##       wt + pf + hx + sbp + dbp + ekg + hg + sz + sg + ap + bm,
##       data = prostate)
##
##
##               coef exp(coef) se(coef)      z      p
## stageIV        -0.17759    0.837 0.175923 -1.0095 3.1e-01
## age             0.02421    1.025 0.009157  2.6436 8.2e-03
## wt            -0.01026    0.990 0.004755 -2.1568 3.1e-02
## pfBedridden(<50%) -1.37275    0.253 0.851150 -1.6128 1.1e-01
## pfBedridden(>50%) -1.17494    0.309 0.850354 -1.3817 1.7e-01
## pfnormal       -1.58867    0.204 0.825542 -1.9244 5.4e-02
## hx             0.51042    1.666 0.120371  4.2404 2.2e-05
## sbp           -0.03309    0.967 0.029279 -1.1300 2.6e-01
## dbp            0.04943    1.051 0.047790  1.0344 3.0e-01
## ekgblock/conduction -0.14590    0.864 0.382387 -0.3816 7.0e-01
## ekgheart strain   0.39317    1.482 0.282538  1.3916 1.6e-01
## ekgnormal       -0.05796    0.944 0.281828 -0.2057 8.4e-01
## ekgold MI        0.00744    1.007 0.301073  0.0247 9.8e-01
## ekgrecent MI     0.81895    2.268 1.055545  0.7759 4.4e-01
## ekgrrhythmic disturb 0.28031    1.324 0.307714  0.9110 3.6e-01
## hg            -0.06893    0.933 0.031996 -2.1543 3.1e-02
## sz             0.01780    1.018 0.004608  3.8640 1.1e-04
## sg             0.10020    1.105 0.041607  2.4083 1.6e-02
## ap            -0.00138    0.999 0.000997 -1.3881 1.7e-01
## bm             0.31964    1.377 0.181333  1.7627 7.8e-02
##
## Likelihood ratio test=99.4 on 20 df, p=1.59e-12 n= 475, number of events= 338
## (27 observations deleted due to missingness)
```

## The distributed fit

The distributed fit for this particular example doesn't work in the current implementation. This is because the  $X$  matrix for each site is singular. The math holds, obviously, but the current implementation is based on re-using as much of the `survival` package as possible. We have to work harder to implement the distributed computation in the situation where  $X$  is singular at at least one site. This affects the computation of the variance (or equivalently, the information matrix). Some work needs to be done to work around this and figure out how best to reuse what's already in the `survival` package.

However, in order to demonstrate that the distributed fit really works, we show below an alternative implementation that yields the same result as the aggregated one.

**An object representing the sites** Here's a reference object for each site. It has several fields: a `data` field containing the data, a `formula` field (as used in the well-known `survival` R package) describing the model being fit. These two are the only ones needed for initializing a site. Other fields that are generated based on these two fields are `modelDataFrame` containing the actual model data used for fitting the model and a `modelMatrix`.

```
site <- setRefClass("siteObject",
  fields = list(
    data = "data.frame",
    formula = "formula",
    modelDataFrame = "data.frame",
```

```

    coxControl = "list",
    modelMatrix = "matrix"),
  methods = list(
    initialize = function(formula, data) {
      'Initialize the object with a formula and dataset'
      formula <- formula
      data <- data
      temp <- coxph.control()
      temp$iter.max <- 0
      coxControl <- temp
      stopifnot(kosher())
    },
    kosher = function() {
      'Check that the class data passes sanity checks'
      modelDataFrame <- model.frame(formula, data=data)
      lhs <- modelDataFrame[, 1]
      ordering <- order(lhs[, 1])
      modelDataFrame <- modelDataFrame[ordering, ]
      data <- data[ordering, ]
      modelMatrix <- model.matrix(formula, data=modelDataFrame)
      TRUE
    },
    dimP = function() {
      'Return the number of covariates'
      ncol(modelMatrix) - 1
    }
  )
)
site$accessors(c("data", "formula", "modelDataFrame", "modelMatrix"))

```

The `initialize` method above mostly sets the fields, generates values for other fields, and does a mild sanity check. It will not proceed further if the function `kosher` returns false. For now the `kosher` function merely orders the data frame by follow-up time, but in a production system a number of other checks might be performed, such as ensuring all named variables are available at the site.

The method `dimP` merely returns the number of columns of the model matrix.

### The (partial) log likelihood function for each site

For our example, the (partial) log likelihood (named `localLogLik`) is simple and can be computed directly. Assuming failure times  $t_i$  and event indicators  $\delta_i$ , it is precisely:

$$l(\beta) = \sum_{i=1}^n \delta_i \left[ z_i \beta - \log \left( \sum_{j \in R(t_i)} \exp(z_j \beta) \right) \right]$$

where  $\beta$  is the vector of parameters,  $z_i$  is row  $i$  of the model matrix (covariates for subject  $i$ ) and  $R(t_i)$  is the risk set at time  $t_i$ .

The first derivative with respect to  $\beta$  is:

$$l'(\beta) = Z^T \delta - \sum_{i=1}^n \delta_i \frac{\sum_{j \in R(t_i)} \exp(z_j \beta) z_j^T}{\sum_{j \in R(t_i)} \exp(z_j \beta)}.$$

```

localLogLik <- function(beta) {
  beta <- c(0, beta) ## model matrix has intercept in model
  z <- modelMatrix

```

```

delta <- modelDataFrame[, 1][, 2] ## event indicators
zBeta <- z %*% beta
sum.exp.zBeta <- rev(cumsum(rev(exp(zBeta))))
ld <- delta %*% (z - apply(diag(as.numeric(zBeta)) %*% z, 2, function(x) rev(cumsum(rev(x))))) / sum.exp.zBeta
result <- sum(delta * (zBeta - log(sum.exp.zBeta))) # assuming Breslow
attr(result, "gradient") <- ld
result
}
site$methods(logLik = localLogLik)

```

**The alternative distributed fit.** We are now ready to do the alternative distributed fit. We split the prostate data into two sites as indicated earlier.

```

siteData <- with(prostate, split(x=prostate, f=rx))
sites <- lapply(siteData,
  function(x) {
    site$new(data = x,
      formula = Surv(dtime, status) ~ stage + age + wt + pf + hx + sbp + dbp + el)
  })

```

Ok, now we can reproduce the original aggregated analysis by writing a full likelihood routine.

```

logLik <- function(beta, sites) {
  sum(sapply(sites, function(x) x$logLik(beta)))
}

```

All that remains is to maximize this log likelihood.

```

mleResults <- nlm(f=function(x) -logLik(x, sites),
  p=rep(0, sites[[1]]$dimP()),
  gradtol=1e-10, iterlim=1000)

```

```

## Warning in nlm(f = function(x) -logLik(x, sites), p = rep(0,
## sites[[1]]$dimP()), : NA/Inf replaced by maximum positive value

```

```

## Warning in nlm(f = function(x) -logLik(x, sites), p = rep(0,
## sites[[1]]$dimP()), : NA/Inf replaced by maximum positive value

```

```

## Warning in nlm(f = function(x) -logLik(x, sites), p = rep(0,
## sites[[1]]$dimP()), : NA/Inf replaced by maximum positive value

```

We print the coefficient estimates side-by-side for comparison.

```

d <- data.frame(distCoef=mleResults$estimate, aggCoef=pcph$coefficients)
rownames(d) <- names(pcph$coefficients)
kable(d)

```

	distCoef	aggCoef
stageIV	-0.1740491	-0.1775897

	distCoef	aggCoef
age	0.0245191	0.0242084
wt	-0.0105003	-0.0102559
pfBedridden(<50%)	-1.3893181	-1.3727528
pfBedridden(>50%)	-1.1506821	-1.1749372
pfnormal	-1.5842335	-1.5886734
hx	0.5206160	0.5104190
sbp	-0.0329601	-0.0330859
dbp	0.0526530	0.0494337
ekgblock/conduction	-0.1505854	-0.1459025
ekgheart strain	0.3980165	0.3931730
ekgnormal	-0.0557321	-0.0579591
ekgold MI	0.0064296	0.0074423
ekgrecent MI	0.8560625	0.8189465
ekgrhythmic disturb	0.2816923	0.2803126
hg	-0.0707047	-0.0689290
sz	0.0180493	0.0178037
sg	0.1016070	0.1002007
ap	-0.0013375	-0.0013841
bm	0.3180923	0.3196404