

# Simulation study with logit model

*Benjamin Christoffersen*

*2017-05-23*

## Intro

This note has four objectives. The first objective is to test how the `ddhazard` fits compare with a Generalized Additive models (GAM) and a “static” logistic model with simulated data. We will look at the following models/estimation methods from `ddhazard` function in the `dynamichazard` package:

- Fits with the Extended Kalman Filters (EKF) with and without extra iterations in the scoring step
- Second order random walks with the EKF estimation method
- Mixture of fixed and time varying effects with the EKF estimation method. Fixed effects are both estimated with the E-step and M-step method described in `ddhazard` vignette
- Fits using the Unscented Kalman filter (UKF)

The second objective is to show how to estimate various models with the function `ddhazard`. For this reason, the note contains intermediate R code which is not needed to understand the simulation results. Thus, we will use `*` in the headers of section to distinguish the content. The headers marked with no `*` indicates sections with results of simulation or contains important comments. Headers with an `*` and `**` shows increasingly less important code to understand the simulation. Consequently, you can skip to the headers with no `*` if you are only interested in the results.

The third objective is to illustrate how the various methods performs for out-of-time prediction (forecasting). By out-of-time we mean that we only observe outcomes up to given time,  $d$ , and then predict the outcome for future observations at time  $d + 1$ .

The fourth objective is to show that both the EKF and UKF scales linearly with the number individuals (series).

All method use the logistic link function. We will do three runs of experiments in the following order:

1. A Model where all effects are time varying and we use the correct binning intervals
2. A model where only one parameter is time varying and we use the correct binning intervals
3. A Model where all effects are time varying but we use incorrect binning intervals

where correct or incorrect binning intervals refers to whether or not we bin at the same time where the coefficient are simulated to change. For example, we bin correctly where we simulate the coefficients to change at time  $1, 2, \dots, d$  and we estimate the coefficient at time  $1, 2, \dots, d$ . The models will be compared in terms of Brier score, median absolute residuals and standard deviation of the absolute residuals. All metrics will be reported on out-sample data or out-of-time data. All plots will have true coefficients as continuous lines while dashed lines are estimates.

You can install the version of the library used to make this vignettes from github with the `devtools` library as follows:

```
current_version # The string to pass devtools::install_github

## [1] "boennecd/dynamichazard@8ad8c0701479c79a581e0143a1b01cc12e01d01a"

devtools::install_github(current_version)
```

You can also get the latest version on CRAN by calling:

```
install.packages("dynamichazard")
```

Moreover, you will also find the source code for the vignette at the github page. The note is not meant to be self contained. It is recommend to see dddhazard vignette for an introduction to the models and methods in the `dynamichazard` package.

## Notions

For clarity, here is a list of used notions:

- Run: An experiment with one of the three previously specified settings where we make  $k$  simulations with  $n$  series in each
- Simulation: One simulation within a run with one set of coefficients  $\vec{\beta}_0, \dots, \vec{\beta}_d$  and given number  $n$  of series
- Series/individuals: A person/individual either making it to the end of the time of the given simulation or dying at some time during the period
- Coefficients: the entries of the vectors  $\vec{\beta}_t$  in a given simulation
- Covariates: vectors  $\vec{x}_{it}$  for a given individual at a given time in simulation

## Findings

The findings are:

- The UKF method seems to perform well for both small and larger number of series
- Taking multiple iterations in the correction step of the EKF seems to be beneficial
- Specifying a fixed effect as time varying or setting the binning number incorrectly has little effect on the results

You will see that the the estimation sometimes fails. It is worth stressing that it is my experience that you can always do “trail-and-error” with the initial covariance matrix in the state equation, the covariance matrix at time zero and tuning parameters in order to get a model to fit a given dataset. Of course, it is a disadvantage that any given data set may require some tuning by the user. Although as will be shown, tuning by the user is not often needed with data sets like those presented here.

## Setup

The following values will be used in the simulation:

```
ns <- c(200, 2000) # Number of series
n_beta <- 5        # Number of covariates
T_max <- 20        # The last time we observe
n_sims <- 100      # Number of simulation in each run

gsub("(^.+)(/dynamichazard.+)$", "...\\2", getwd())

## [1] ".../dynamichazard/vignettes/Prebuild"

source("../R/test_utils.R")
```

`ns` is the number of series (individuals) we will estimate in each of the simulation in each of the runs. Thus, we will perform simulations with a total of 200 and 2000 series in each. Each simulation will have `n_beta = 5` covariates plus an intercept. Each run will simulate `n_sims = 100` times. Finally, we source the `test_utils.R` file to define the simulation function. You can find this script on the github site. `T_max` is the number of bins/intervals we observe. Thus, we have 1, 2, ..., `T_max + 1` covariate vectors (+1 for the time zero coefficient vector).

## Fitting true model

We will make runs for various number of individuals in this section where we estimate a model where all effects are time varying and we use the correct binning intervals. Thus, the only models that are misspecified are the model with one time varying effect (which will be `x2`) and the model where we use a second order random walk.

### `do.call` function

We will use `do.call` in this vignette. To my knowledge, `do.call` is not standard so this section is included to give a brief introduction to `do.call` for users who are not familiar with `do.call`. We can take an example with the `mean` function. We will make the following call where we set `na.rm` to `TRUE`:

```
mean(x = c(1, 2, NA, 6), na.rm = T)
```

```
## [1] 3
```

This call can also be made as follows `do.call`:

```
arg_ex <- list(x = c(1, 2, NA, 6), na.rm = T)
do.call(mean, arg_ex)
```

```
## [1] 3
```

Hence, `do.call` is useful in situation where we make calls where almost all the arguments are the same. For example, in the setting where we have arguments `a1`, `a2`, ..., up to `a1000` and we only want to change argument `a101` say. This can then be done as follows:

```
# Not runnable
arg_ex <- list(a1 = x,
              a2 = y,
              ..., # enter all the other values
              a1000 = z)
do.call(some_func, arg_ex)

# change only a101 argument and keep the rest as the arguments as is
arg_ex$a101 <- some_specific_value
do.call(some_func, arg_ex)
```

### Definition of simulation function

Below, we define a list of `default_args` (default arguments) to our simulation function which we can later use using `do.call`.

```
# Default arguments for simulation
default_args <- list(
  n_vars = n_beta, # Number of betas not including intercept
  beta_start = c(-1, -.5, 0, 1.5, 2), # start value of coefficients
  intercept_start = -5, # start value of intercept
  sds = c(.1, rep(.5, n_beta)), # std. deviations in state equation
  t_max = T_max, # Largest time we observe
  x_range = 1, # range of covariates
  x_mean = .5, # mean of covariates
  tstart_sampl_func = # we randomly draw the start time of each serie
```

```
function(...) max(runif(1, min = -10, max = 18), 0)
)
```

Let  $\vec{\beta}_t$  denote the time varying covariates at time  $t$ . Then the `beta_start` is the time 0 values of the coefficients and `intercept_start` is the starting value of the intercept. The `sds` are the standard deviations,  $\sigma_i$ , in the state equation. Hence,

$$\beta_{j,t} = \beta_{j,t-1} + \epsilon_{j,t}, \quad \epsilon_{j,t} \sim N(0, \sigma_j^2)$$

where each margin is independent of the others. The `x_mean` and `x_range` defines how the covariate values are simulated. The above setting implies that  $x_{itj} = \text{Unif}(0.5 - 1/2, 0.5 + 1/2)$  where  $x_{itj}$  is the  $i$ 'th individuals covariate  $j$  at time  $t$ . The covariate vector  $\vec{x}_{it}$  is updated at time differences of  $1 + \eta$  where  $\eta \sim \text{Exp}(1)$  and  $\eta$ s are drawn separately for each individual for each covariate vector. The motivation for this behavior is that we can have different covariate update times than our binning time in a given study. For instance, say we are looking at a medical study and the covariates are laboratory values. The time of laboratory values from an individual's visit the doctor can differ from whatever binning periods we use in the state-space model. Further, the time when laboratory values are updated can differ between patients. One might see his doctor every week or so while another only sees his doctor every year.

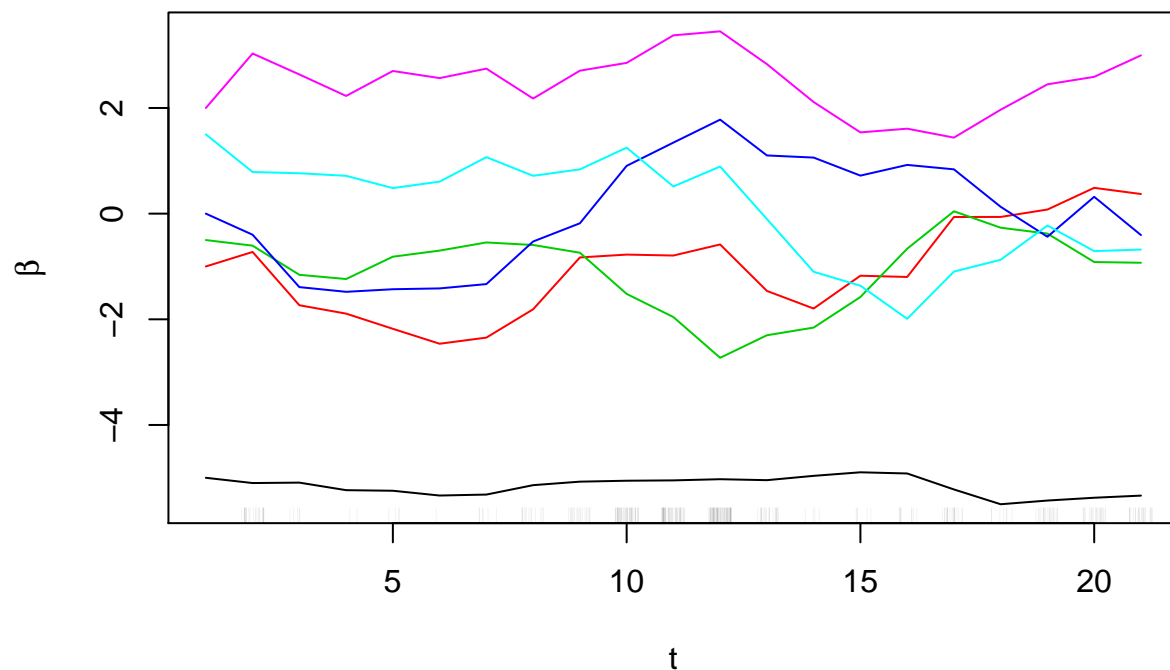
Below we illustrate how the coefficients vectors from a simulation can look:

```
# We can simulate by
set.seed(51231)
sims <- do.call(test_sim_func_logit, c(list(n_series = max(ns)), default_args))

# This is how the state vectors look
# We define a function so we can re-use it later
plot_func <- function(ylim = c()){ # we define a function here so we can use it later
  matplot(sims$betas, type = "l", lty = 1, ylab = expression(beta), xlab = "t",
    ylim = range(sims$betas, ylim), col = 1:(n_beta + 1))

  # Add rug plot to illustrate when people die
  rug(jitter(sims$res$stop[sims$res$event==1], amount = .25) + 1,
    col = rgb(0, 0, .05))
}

plot_func()
```



The black line is the intercept while the colored lines are the coefficients for the covariates. The lines on the x-axis illustrate when we observe that individuals die. There is one line for each death. Next, we can look at the number of failures in each simulation:

```
# We get a "decent" amount of failures and survivors in some of the simulations
# We use do.call to avoid repeating the above argument list
set.seed(468249)
n_fails_in_sim <- rep(NA_real_, 15)
for(i in seq_along(n_fails_in_sim)){
  sims <- do.call(test_sim_func_logit, c(
    default_args, c(list(n_series = max(ns))))) # Take largest amount of series
  n_fails_in_sim[i] <- sum(sims$res$event)
}
n_fails_in_sim # number of failures in each simulation
```

```
## [1] 1132 1822 349 759 1062 91 1321 707 1635 466 1548 83 1596 1577
## [15] 1151
```

### \* Definition of fit functions

We will define functions to estimate the different models with a data frame as the first argument where the data frame is from a `test_sim_func_logit` call. This will reduce the amount of code later.

## \*\* Definition of static fit

Below, we define function to fit a model where the coefficients are fixed ( $\vec{\beta}_t = \vec{\beta}$ ). It is estimated using `glm.fit`:

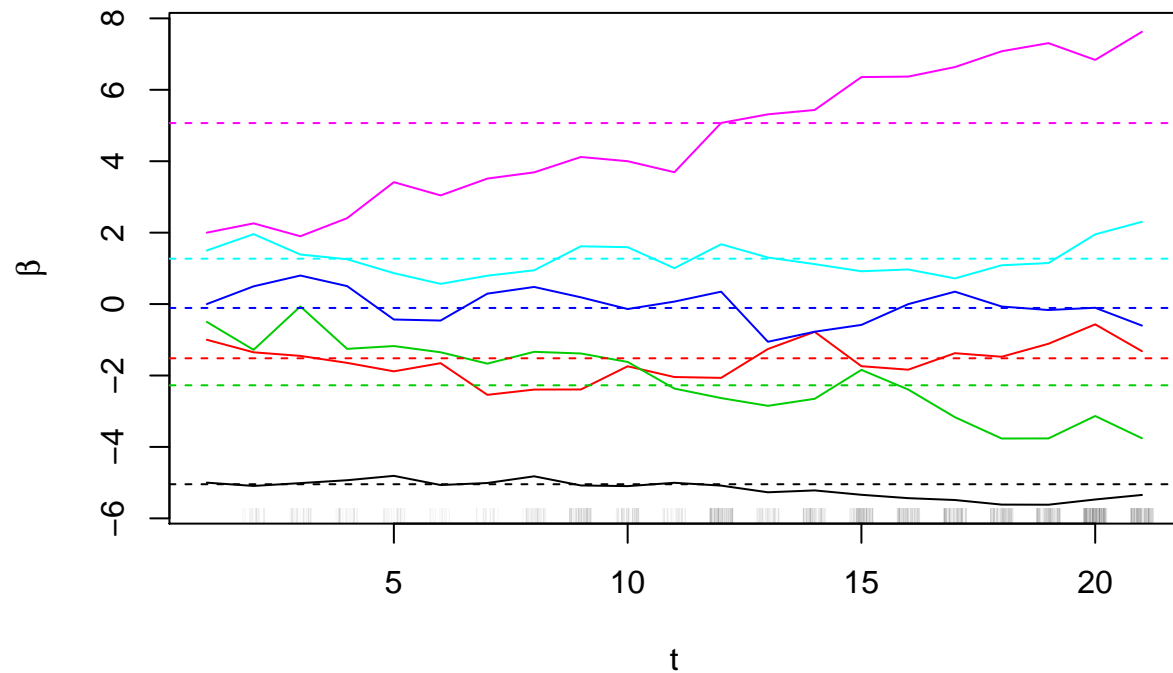
```
library(survival); library(dynamichazard)

# Set up function for static fit
fit_funcs = list()
fit_funcs$static <- function(s = sims$res)
  static_glm(formula = Surv(tstart, tstop, event) ~ x1 + x2 + x3 + x4 + x5,
             data = s, max_T = T_max, by = 1, id = s$id, speedglm = FALSE)

fit <- fit_funcs$static()
class(fit) # returns a glm object

## [1] "glm" "lm"

# Estimates seems plausible
plot_func(ylim = fit$coefficients)
abline(h = fit$coefficients, col = 1:(n_beta + 1), lty = 2)
```



## \* Definition of ddhazard fit functions

Below, we define a function to fit a first order random walk model with a given learning rate and potential extra iterations in the scoring step (see the ddhazard vignette for details):

```

library(survival); library(dynamichazard)

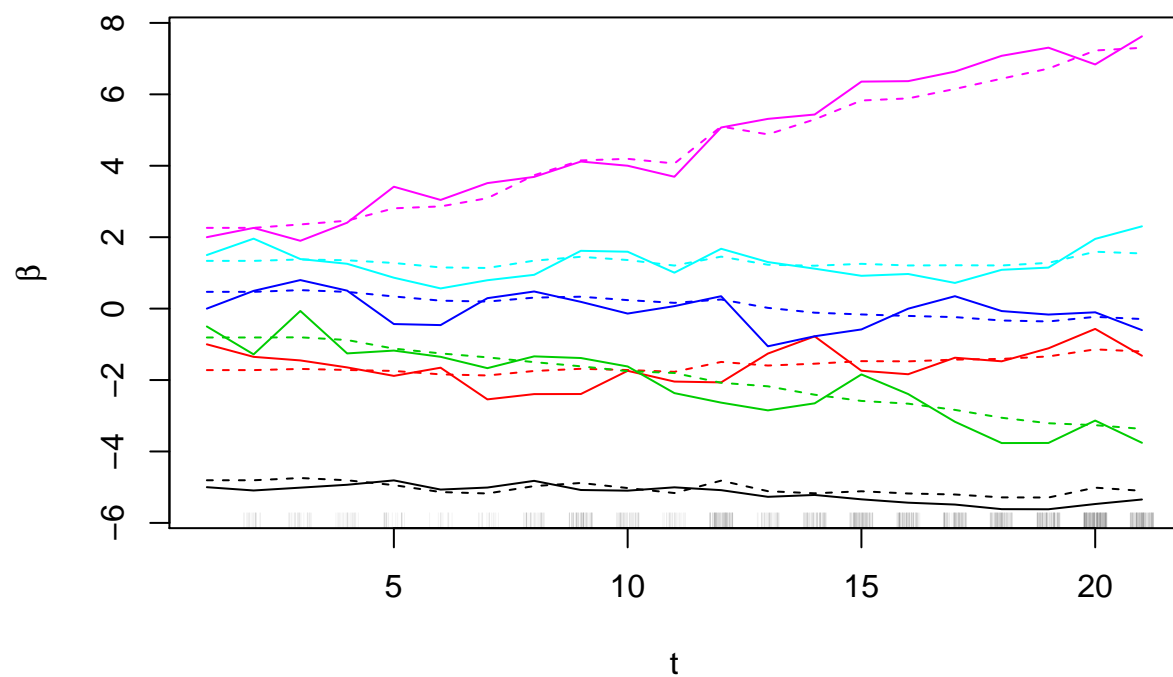
# We will use glm.fit for the starting value
options(ddhazard_use_speedglm = FALSE)

# Set up function ddhazard fit function for convenience
# LR:      learning rate in correction step
# NR_eps:   tolerance in correction step. NULL yields no extra iterations
fit_funcs$dd <- function(s = sims$res, LR = 1, NR_eps = NULL)
  tryCatch({
    ddhazard(
      formula = Surv(tstart, tstop, event) ~ x1 + x2 + x3 + x4 + x5,
      data = s, max_T = T_max, by = 1, id = s$id,
      Q_0 = diag(
        # We set the Q_0 argument lower when we take multiple iterations
        # See the ddhazard vignette under the GMA model for arguments herefor
        if(is.null(NR_eps)) 1000000 else 1,
        n_beta + 1),
      Q = diag(.01, n_beta + 1),
      control = list(LR = LR, NR_eps = NR_eps, eps = 0.01))
  }, error = function(...) NA) # Return NA if fails

fit <- fit_funcs$dd()

# Plot estimates and actual coefficients
plot_func(ylim = fit$state_vecs)
matplot(fit$state_vecs, col = 1:(n_beta + 1), lty = 2,
        type = "l", add = T)

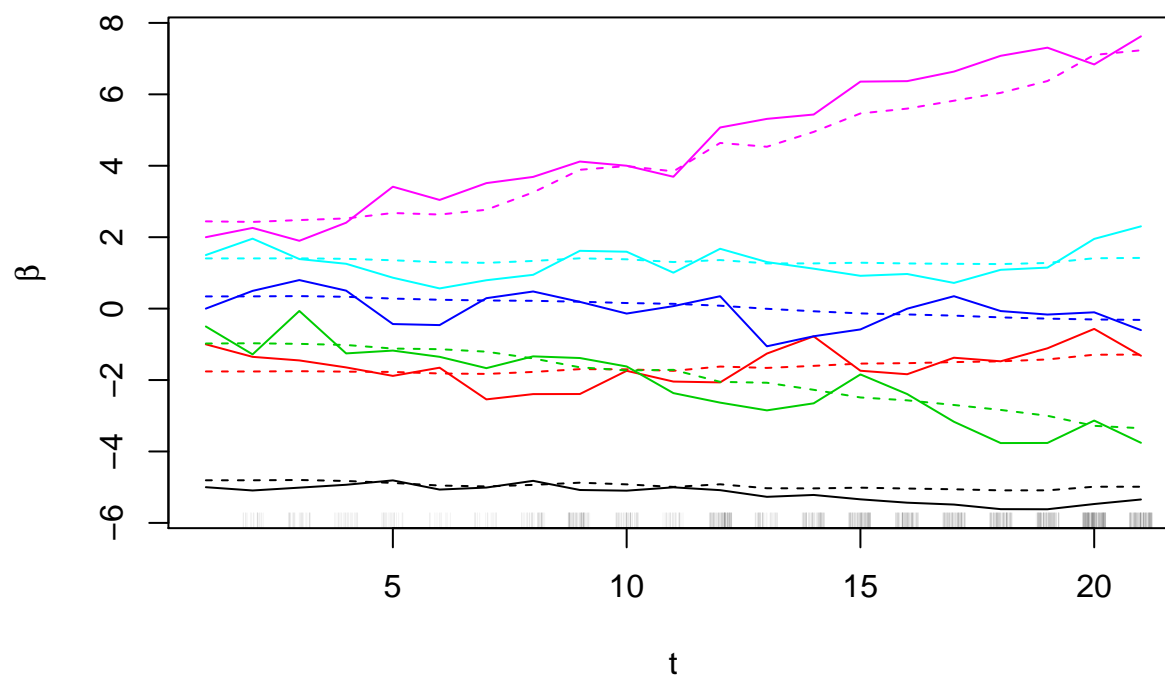
```



```
# Same call with extra iterations
fit <- fit_funcs$dd(LR = .5, NR_eps = .01)

# Look at new plot
plot_func(ylim = fit$state_vecs)
matplot(fit$state_vecs, col = 1:(n_beta + 1), lty = 2,
        type = "l", add = T)
```



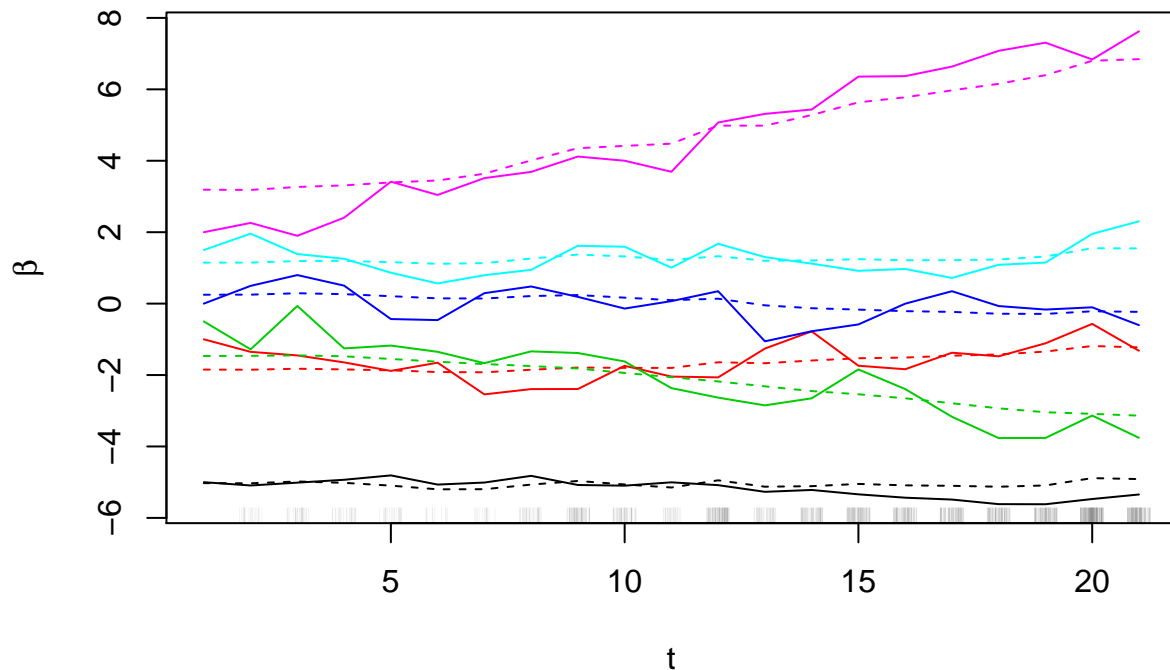


Below, we define a function to fit a first order random walk model with the UKF method:

```
# Fitting with UKF
fit_funcs$dd_UKF <- function(s = sims$res, alpha = 1, beta = 0){
  tryCatch({
    ddhazard(
      formula = Surv(tstart, tstop, event) ~ x1 + x2 + x3 + x4 + x5,
      data = s, max_T = T_max, by = 1, id = s$id,
      Q_0 = diag(1, n_beta + 1), Q = diag(.01, n_beta + 1),
      control = list(
        eps = 0.1,
        alpha = alpha,      # Set tuning parameter
        beta = beta,        # Set tuning parameter
        method = "UKF"))    # Set estimation method (EKF is default)
  }, error = function(...) NA) # Return NA if fails
}

fit <- fit_funcs$dd_UKF()

# Look at new plot
plot_func(ylim = fit$state_vecs)
matplot(fit$state_vecs, col = 1:(n_beta + 1), lty = 2,
        type = "l", add = T)
```



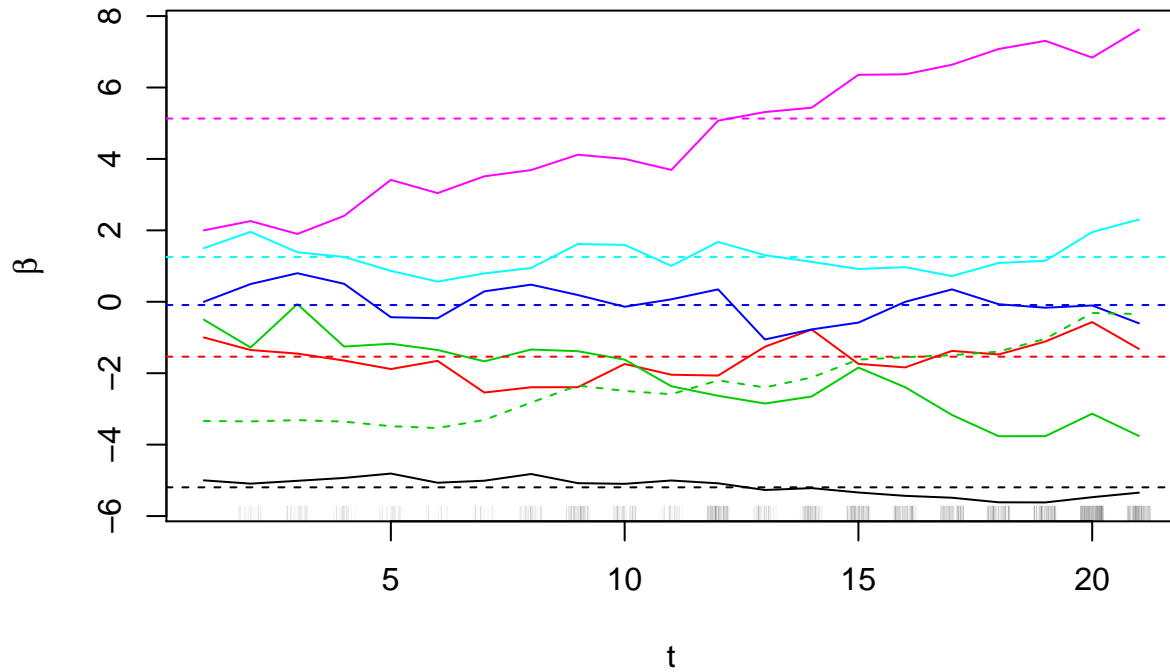
Below, we define a function to estimate a first order random walk model where only one parameter ( $x_2$ ) is time varying:

```
# Fitting with fixed effects
fit_funcs$dd_fixed <- function(
  s = sims$res, LR = 1, NR_eps = NULL,
  fixed_terms_method = "M_step"){ # The method to use to estimate the fixed
                                     # fixed effects

  tryCatch({
    ddhazard(
      formula = Surv(tstart, tstop, event) ~
        ddFixed(1) +                               # Fix intercept
        ddFixed(x1) + x2 +                          # Note x2 is time varying
        ddFixed(x3) + ddFixed(x4) + ddFixed(x5),
      data = s, max_T = T_max, by = 1, id = s$id,
      Q_0 = diag(1, 1), Q = diag(.01, 1),
      control = list(LR = LR, NR_eps = NR_eps, eps = 0.1,
                     fixed_terms_method = fixed_terms_method))
  }, error = function(...) NA) # Return NA if fails
}

fit <- fit_funcs$dd_fixed()

# Look at new plot
plot_func(ylim = range(fit$state_vecs, fit$fixed_effects))
matplot(fit$state_vecs, col = 3, lty = 2, type = "l", add = T)
abline(h = fit$fixed_effects, col = c(1:2, 4:6), lty = 2)
```

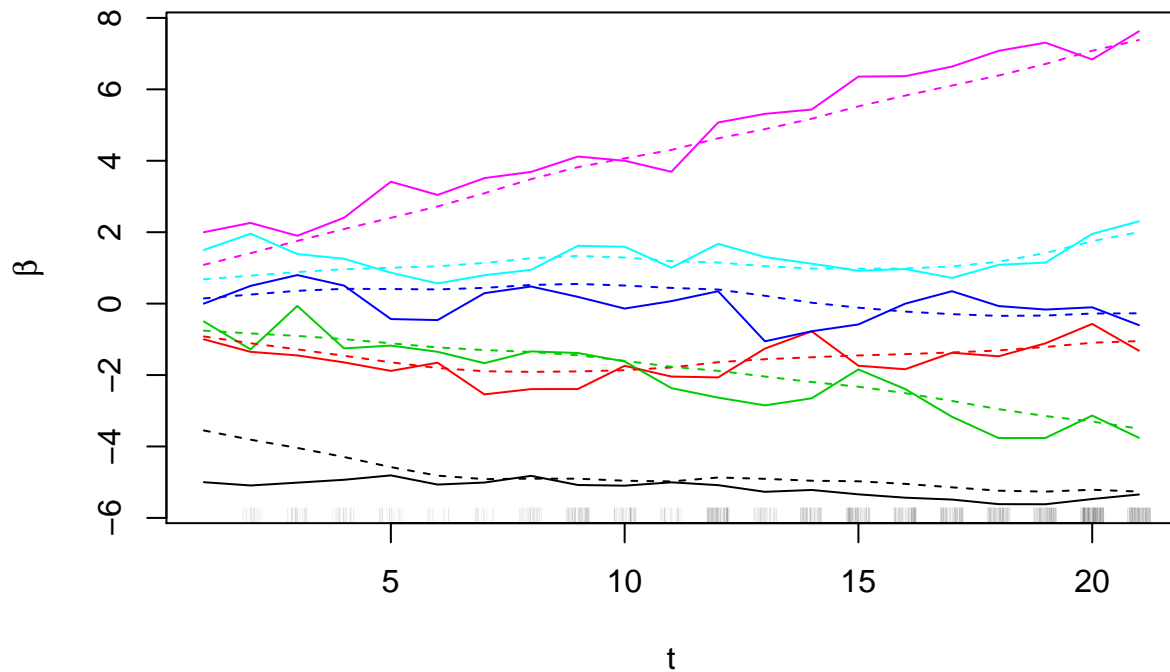


Next, we define a function to fit the model with a second order random walk:

```
# Fitting with second order
fit_funcs$dd_2_order <- function(s = sims$res, LR = 1, NR_eps = NULL){
  tryCatch({
    ddhazard(
      formula = Surv(tstart, tstop, event) ~ x1 + x2 + x3 + x4 + x5,
      data = s, max_T = T_max, by = 1, id = s$id,
      # Q_0 and Q needs more elements
      Q_0 = diag(c(rep(1, n_beta + 1), rep(0.5, n_beta + 1))),
      Q = diag(c(rep(.01, n_beta + 1))),
      order = 2, # specify the order
      control = list(LR = LR, NR_eps = NR_eps, eps = 0.1))
  }, error = function(...) NA) # Return NA if fails
}

fit <- fit_funcs$dd_2_order()

# Look at new plot
plot_func(ylim = fit$state_vecs)
matplot(fit$state_vecs[, 1:6], col = 1:(n_beta + 1), lty = 2,
        type = "l", add = T)
```



## \*\* Definition of GAM fit function

We define the estimation method for the Generalized additive model in the next code snippet. We use `bam` function from the `mgcv` package which corresponds to `gam` but for large datasets.

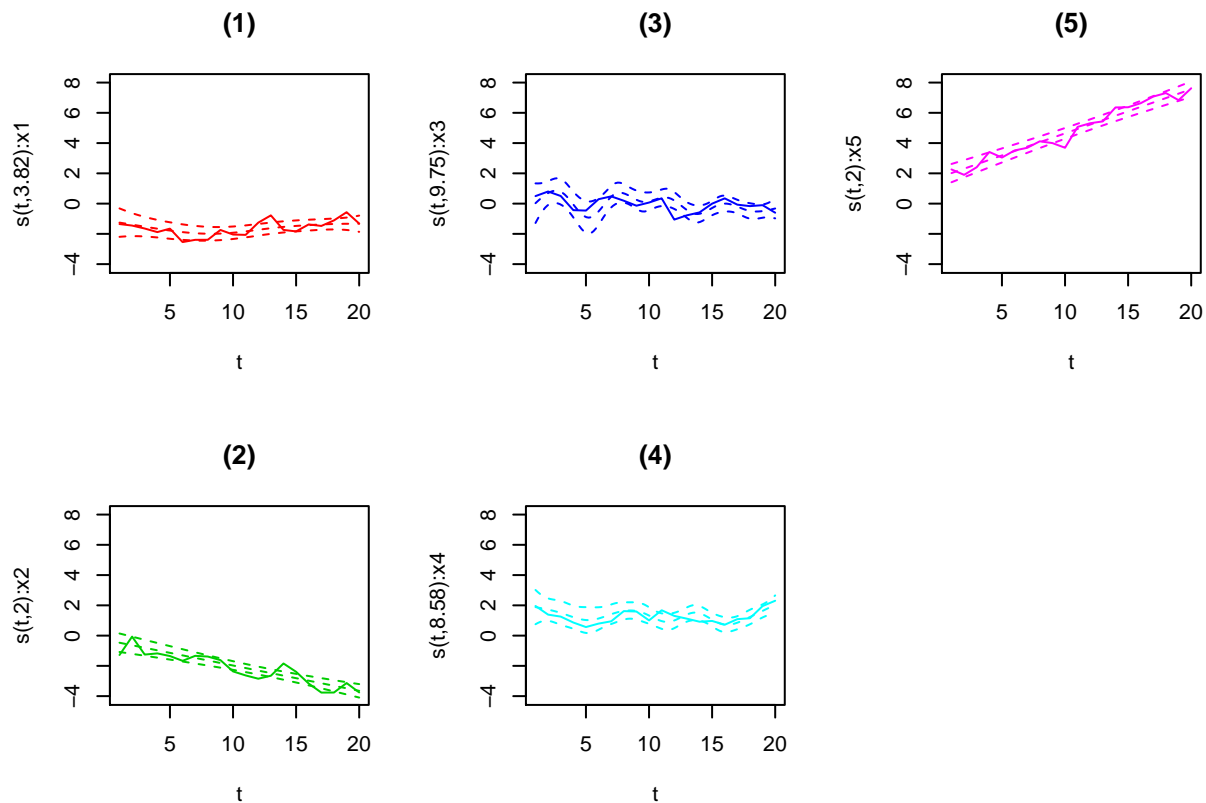
```
library(mgcv)
fit_funcs$gam <- function(s = sims$res){
  # get data frame for fitting
  dat_frame <- get_survival_case_weights_and_data(
    formula = Surv(tstart, tstop, event) ~ x1 + x2 + x3 + x4 + x5,
    data = s, max_T = T_max, by = 1, id = s$id, use_weights = F)$X
  # fit model
  bam(
    formula = Y ~
      # cr is cubic basis with k knots
      s(t, bs = "cr", k = 10, by = x1) +
      s(t, bs = "cr", k = 10, by = x2) +
      s(t, bs = "cr", k = 10, by = x3) +
      s(t, bs = "cr", k = 10, by = x4) +
      s(t, bs = "cr", k = 10, by = x5),
    family = binomial, data = dat_frame,
    method = "GCV.Cp",
    control =
      gam.control(nthreads = parallel::detectCores() - 1)) # Use parallel
}
```

```

# fit model
fit <- fit_funcs$gam()

# Compare plot
layout(matrix(1:6, nrow = 2))
for(i in 1:n_beta){
  plot(fit, pages = 0, rug = F, col = i + 1, select = i, lty = 2,
       main = paste0("( ", i, " )"))
  lines(sims$betas[-1, i + 1], col = i + 1)
}

```



## \*\* Definition of prediction functions

The following code snippets define predictions methods for each of the estimation methods. We start off by defining a split function such that we can sample individuals (series) into a test set and a training test:

```

split_func <- function(s = sims$res){
  # Sample ids
  test_ids <- sample(
    unique(s$id), floor(length(unique(s$id)) / 2), replace = F)

  # Return separate data frames
  return(list(test_dat = s[s$id %in% test_ids, ],
             fit_dat = s[!s$id %in% test_ids, ]))
}

```

```

# Illustrate use
tmp <- split_func()
# No ids intersect in the two sets
length(intersect(tmp$test_dat$id, tmp$fit_dat$id))

```

```
## [1] 0
```

```

# The union is exactly the number of ids we simulated
length(union(tmp$test_dat$id, tmp$fit_dat$id))

```

```
## [1] 2000
```

Having defined the splitting method, we turn to the prediction functions. The idea is to define the `brier_funcs$general` function which takes in a prediction function, a fit and a data frame. Next, we then define individual prediction functions for each of the models which will be passed to `brier_funcs$general`:

```

# Define general prediction function
brier_funcs <- list()
brier_funcs$general <- function(brier_func, fit, eval_data_frame){
  d_frame <- get_survival_case_weights_and_data(
    formula = Surv(tstart, tstop, event) ~ x1 + x2 + x3 + x4 + x5,
    data = eval_data_frame, max_T = T_max, by = 1, id = eval_data_frame$id,
    use_weights = F)$X

  # Change start and stop times
  d_frame$tstart <- d_frame$t - 1
  d_frame$tstop <- d_frame$t

  # Compute residuals
  resid <- brier_func(fit, d_frame)

  # Return estimates
  with(
    resid,
    list(brier = mean(response^2),
         median_abs_res = median(abs(response)),
         sd_res = sd(response),
         median_deviance = median(deviance),
         sd_deviance = sd(deviance)))
}

# Prediction method for static model
brier_funcs$static <- function(fit, d_frame){
  preds <- predict(fit, newdata = d_frame, type = "response")
  list(
    response = d_frame$Y - preds,
    deviance =
      d_frame$Y * log(preds) + (1 - d_frame$Y) * log(1 - preds))
}

# Test function
fit <- fit_funcs$static(tmp$fit_dat)
cols <-
  c("brier", "median_abs_res", "sd_res", "median_deviance", "sd_deviance")
unlist( # in sample stats

```

```

brier_funcs$general(brier_funcs$static, fit, tmp$fit_dat)[cols])

##          brier median_abs_res          sd_res median_deviance
##          0.04258          0.01763          0.20636          -0.01779
##          sd_deviance
##          0.52205

unlist( # out sample stats
  brier_funcs$general(brier_funcs$static, fit, tmp$test_dat)[cols])

##          brier median_abs_res          sd_res median_deviance
##          0.04741          0.01720          0.21768          -0.01735
##          sd_deviance
##          0.57509

# Define prediction function for ddhazard model
brier_funcs$dd <- function(fit, d_frame){
  preds <- predict(fit, new_data = d_frame, tstart = "tstart", tstop = "tstop")

  # We truncate as we can get zero-one outcome if the model diverges
  preds$fits <- pmax(pmin(preds$fits, 1 - 1e-14), 1e-14)

  list(
    response = d_frame$Y - preds$fits,
    deviance =
      d_frame$Y * log(preds$fits) + (1 - d_frame$Y) * log(1 - preds$fits)
  )
}

fit <- fit_funcs$dd(tmp$fit_dat)

unlist( # in sample stats
  brier_funcs$general(brier_funcs$dd, fit, tmp$fit_dat)[cols])

##          brier median_abs_res          sd_res median_deviance
##          0.03816          0.01664          0.19531          -0.01678
##          sd_deviance
##          0.48595

unlist( # out sample stats
  brier_funcs$general(brier_funcs$dd, fit, tmp$test_dat)[cols])

##          brier median_abs_res          sd_res median_deviance
##          0.04298          0.01645          0.20731          -0.01659
##          sd_deviance
##          0.53966

# Define prediction function for gam model
brier_funcs$gam <- function(fit, d_frame){
  preds <- predict(fit, newdata = d_frame, type = "response")
  list(
    response = d_frame$Y - preds,
    deviance =
      d_frame$Y * log(preds) + (1 - d_frame$Y) * log(1 - preds)
  )
}

fit <- fit_funcs$gam(tmp$fit_dat)

```

```

unlist( # in sample stats
  brier_funcs$general(brier_funcs$gam, fit, tmp$fit_dat)[cols])

##          brier  median_abs_res          sd_res median_deviance
##          0.03845          0.01472          0.19610          -0.01483
##          sd_deviance
##          0.50914

unlist( # out sample stats
  brier_funcs$general(brier_funcs$gam, fit, tmp$test_dat)[cols])

##          brier  median_abs_res          sd_res median_deviance
##          0.04308          0.01445          0.20748          -0.01456
##          sd_deviance
##          0.56195

```

## \*\* Definition of multiple simulations function

To make things easier, we define a function that takes in a function to simulate from. Given a function to simulate with, the new function perform `n_sims = 100` simulations for each of values `ns` (200 and 2000):

```

simulate_n_print_res <- function(
  sim_func, # Function that takes one argument which is number of series
  NR_eps = c(.0001, NA), # Tolerance in scoring step
  file_prefix) # file_prefix for output
{
  for(n in ns){
    file_name <- paste0(file_prefix, "_", n, ".RDS")
    do_compute <- !file.exists(file_name)

    if(do_compute){
      out <- array(NA_real_, dim = c(n_sims, 8, 5),
        dimnames = list(
          NULL,
          c("static", "Extra correction", "Single correction",
            "2 order EKF", "Fixed E-step", "Fixed M-step", "UKF", "gam"),
          c("Brier", "Abs res", "Sd res", "Dev", "Sd dev")))

      n_failures_and_surviers <- array(
        NA_integer_, dim = c(2, n_sims),
        dimnames = list(c("# failures", "# survivors"), NULL))

      #####
      # Progress bar for inpatient people (me)
      pb <- winProgressBar(
        paste("Estimating with n =", n), "", 0, n_sims, 50)
      #####

      for(i in 1:n_sims){
        #####
        info <- sprintf("%.2f%% done", 100 * (i - 1) / n_sims)
        setWinProgressBar(pb, i - 1, paste("Estimating with n =", n), info)
        #####
      }
    }
  }
}

```



```

# Sample until we get an outcome have sufficient amount of deaths and
# survivors
repeat{
  sims <- sim_func(n)

  # We want some survivors and some deaths
  if(sum(sims$res$event) > 25 && n - sum(sims$res$event) > 25)
    break
}

n_failures_and_surviers["# failures", i] <- sum(sims$res$event)
n_failures_and_surviers["# survivors", i] <- n - sum(sims$res$event)

# Split data
sim_split <- split_func(sims$res)

# Fit static model
static_fit <- fit_funcs$static(sim_split$fit_dat)

# Fit dd model
dd_fits <- list(rep(NA, length(NR_eps)))
for(k in seq_along(NR_eps)){
  dd_fits[[k]] <- fit_funcs$dd(
    sim_split$fit_dat,
    NR_eps = if(is.na(NR_eps[k])) NULL else NR_eps[k])
}

# Fit second order
dd_2_order <- fit_funcs$dd_2_order(sim_split$fit_dat)

# Fit fixed effect
dd_fixed_E_step <- fit_funcs$dd_fixed(sim_split$fit_dat,
                                       fixed_terms_method = "E_step")
dd_fixed_M_step <- fit_funcs$dd_fixed(sim_split$fit_dat,
                                       fixed_terms_method = "M_step")

# UKF fit
dd_UKF <- fit_funcs$dd_UKF(sim_split$fit_dat)

# Fit gam model
gam_fit <- fit_funcs$gam(sim_split$fit_dat)

# Evaluate on test data
models <- c(list(static_fit), dd_fits,
            list(dd_2_order, dd_fixed_E_step, dd_fixed_M_step,
                dd_UKF, gam_fit))

eval_funcs = c(brier_funcs$static,
               replicate(length(dd_fits) + 4, brier_funcs$dd),
               brier_funcs$gam)

for(j in seq_along(models)){
  if(length(models[[j]]) == 1 && is.na(models[[j]]))

```

```

    next # We have to skip model fits that failed

    metrics <- brier_fncs$general(
      eval_fncs[[j]], models[[j]], sim_split$test_dat)

    out[i, j, "Brier"] <- metrics$brier
    out[i, j, "Abs res"] <- metrics$median_abs_res
    out[i, j, "Sd res"] <- metrics$sd_res
    out[i, j, "Dev"] <- metrics$median_deviance
    out[i, j, "Sd dev"] <- metrics$sd_deviance
  }
}

#####
close(pb)
#####

# Save results
saveRDS(out, file = file_name)
}

# Load results
out <- readRDS(file = file_name)

# Print results
did_fit <- apply(out[, , 1], 2, function(x) n_sims - sum(is.na(x)))
n_cases_all_success <- sum(complete.cases(out[, , 1]))
metric_where_all_fit <-
  t(apply(out[complete.cases(out[, , 1]), , , drop = F], 3, colMeans))

metric_where_all_fit <- formatC(metric_where_all_fit ,format="f", digits=3)
n_cases_all_success <- formatC(n_cases_all_success, format="d")

print(knitr::kable(cbind(
  t(metric_where_all_fit), "# succesful fits" = did_fit),
  caption = paste(
    "Mean of metrics with", n/2, "series in test and fit data. 'Abs res' is the median of the absolute
    align = "r"))
cat("\n")

# Make boxplot of std deviance residuals
par(mar = c(5, 8, 4, 2))
boxplot(
  out[complete.cases(out[, , 1]), , "Sd dev"],
  main = paste0("Std of deviance residuals w/ ", n/2 , " series"),
  cex.axis = .75, horizontal = TRUE, las = 2,
  ylim = c(min(out[, , "Sd dev"], na.rm = T),
    min(max(out[, , "Sd dev"], na.rm = T), 2)))
}
}

```

## Simulating

We are now able to simulate from the model where all effects are time varying and we use the correct binning intervals with the code below:

```
set.seed(1243)
# Use simulation function
simulate_n_print_res(
  sim_func = function(n)
    do.call(test_sim_func_logit, c(default_args, c(list(n_series = n)))),
  file_prefix = "logit_sim_all_varying")
```

Table 1: Mean of metrics with 100 series in test and fit data. ‘Abs res’ is the median of the absolute residuals, ‘Sd res’ is the standard deviation of the residuals, ‘Dev’ is median of the deviance residuals and ‘Sd dev’ is the standard deviation of deviance residuals. Only simulations that succeeds for all setups are included. There are 100 of these simulations. The last column shows the number of successful fits for each setup.

	Brier	Abs res	Sd res	Dev	Sd dev	# succesful fits
static	0.043	0.033	0.199	-0.033	0.559	100
Extra correction	0.041	0.030	0.193	-0.030	0.522	100
Single correction	0.042	0.038	0.195	-0.039	0.489	100
2 order EKF	0.042	0.038	0.196	-0.039	0.900	100
Fixed E-step	0.043	0.043	0.198	-0.044	0.499	100
Fixed M-step	0.042	0.030	0.196	-0.031	0.549	100
UKF	0.041	0.029	0.194	-0.030	0.535	100
gam	0.042	0.020	0.196	-0.020	0.634	100

## Std of deviance residuals w/ 100 series

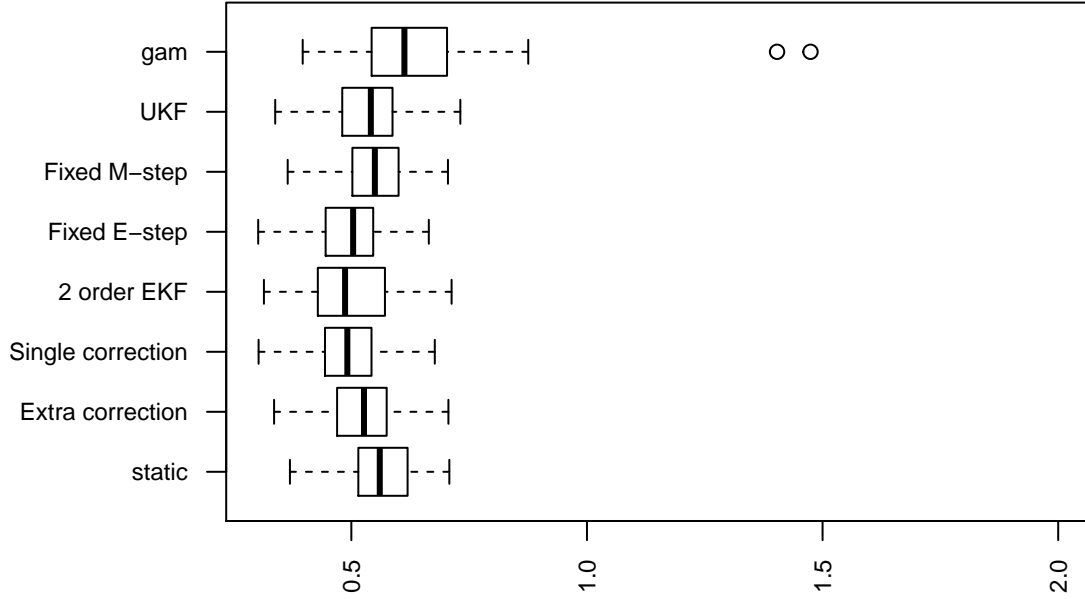
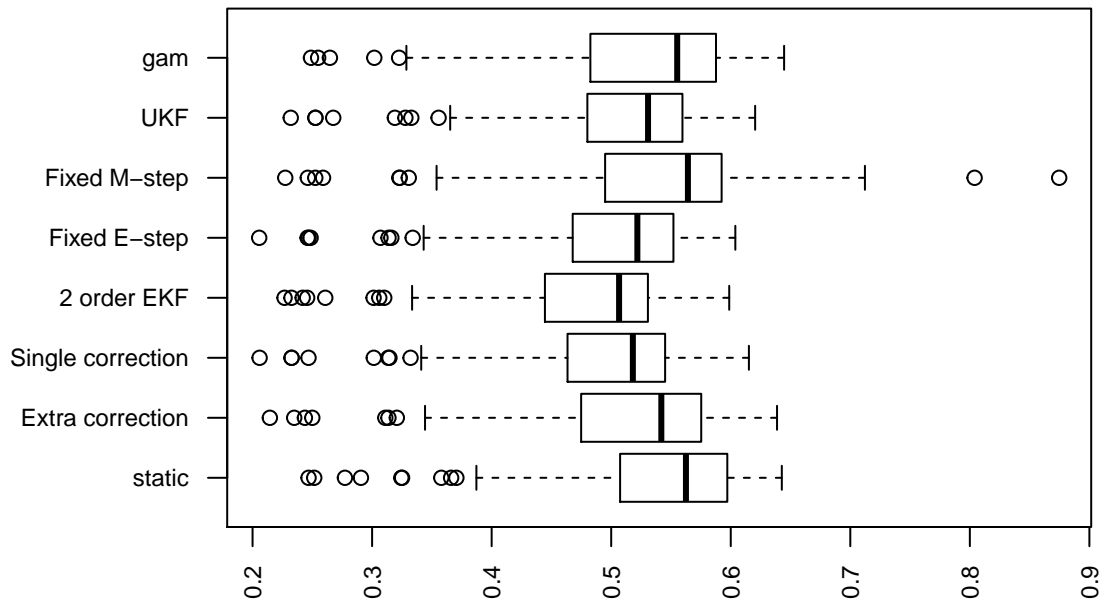


Table 2: Mean of metrics with 1000 series in test and fit data. ‘Abs res’ is the median of the absolute residuals, ‘Sd res’ is the standard deviation of the residuals, ‘Dev’ is median of the deviance residuals and ‘Sd dev’ is the standard deviation of deviance residuals. Only simulations that succeeds for all setups are included. There are 100 of these simulations. The last column shows the number of successful fits for each setup.

	Brier	Abs res	Sd res	Dev	Sd dev	# succesful fits
static	0.053	0.050	0.211	-0.053	0.534	100
Extra correction	0.045	0.032	0.198	-0.033	0.510	100
Single correction	0.047	0.038	0.200	-0.039	0.488	100
2 order EKF	0.046	0.042	0.199	-0.043	0.474	100
Fixed E-step	0.048	0.047	0.203	-0.049	0.494	100
Fixed M-step	0.049	0.037	0.205	-0.038	0.536	100
UKF	0.046	0.036	0.199	-0.037	0.503	100
gam	0.046	0.030	0.199	-0.031	0.520	100

## Std of deviance residuals w/ 1000 series



Notice that we only compare across methods with mean metrics where all succeeded to fit.

## Conclusion on run

All models perform better than the static model in terms of Brier score. The UKF and EKF with extra iterations also does well on the Brier score compared with the GAM model. The finding is slight different if we look at the box plot of the standard deviations of the deviance residuals. The GAM model does worse here compared to the other methods.

## Single time varying parameter

In this part, we will look at the performs when only single coefficient ( $x_2$ ) varies. Thus, we can see if the models where only ( $x_2$ ) is modeled as varying performs better.

### \*\* Definition of simulation function

We start by defining the simulation function. The main change here is that we only set a single standard deviation and that we set it larger than before:

```
# Use simulation function
set.seed(9999)
sim_one_varying <- function(n){
  test_sim_func_logit(
```

```

n_series = n,
sds = c(sqrt(3)), # Large variance
is_fixed = c(1:2, 4:6), # All but param three (x2) is fixed

# Same values as before
n_vars = n_beta,
beta_start = c(-1, -.5, 0, 1.5, 2),
intercept_start = -4,
t_max = T_max,
x_range = 1,
x_mean = .5)
}

```

### \* Illustration of single simulation

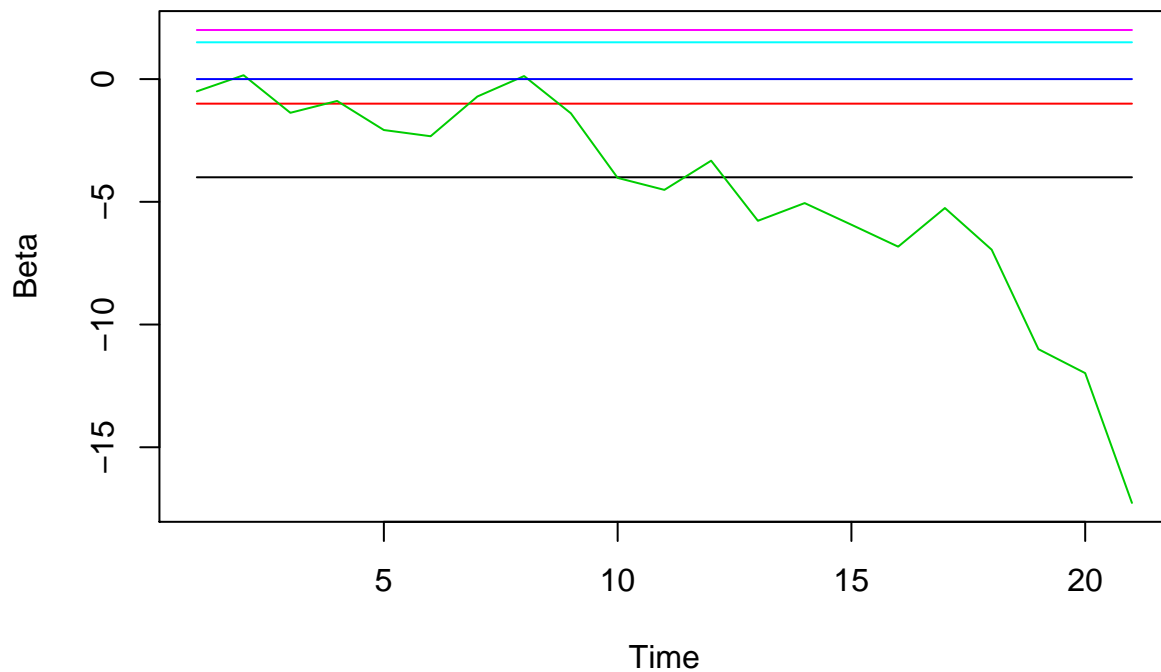
```

# We get a more variable number of failures and survivors (we simulate 200
# series)
replicate(10, sum(sim_one_varying(200)$res$event)) # print number of failures

## [1] 197 200 104 187 129 49 200 37 200 62

# Here is an example of a series
tmp <- sim_one_varying(200)
matplot(tmp$betas, type = "l", lty = 1, ylab = "Beta", xlab = "Time")

```



## Simulating

We can simulate with the following call:

```
# Use simulation function
set.seed(8080)
simulate_n_print_res(
  sim_func = sim_one_varying,
  file_prefix = "logit_sim_one_varying")
```

Table 3: Mean of metrics with 100 series in test and fit data. ‘Abs res’ is the median of the absolute residuals, ‘Sd res’ is the standard deviation of the residuals, ‘Dev’ is median of the deviance residuals and ‘Sd dev’ is the standard deviation of deviance residuals. Only simulations that succeeds for all setups are included. There are 100 of these simulations. The last column shows the number of successful fits for each setup.

	Brier	Abs res	Sd res	Dev	Sd dev	# succesful fits
static	0.036	0.030	0.182	-0.031	0.559	100
Extra correction	0.035	0.026	0.179	-0.026	0.523	100
Single correction	0.040	0.030	0.183	-0.031	0.703	100
2 order EKF	0.036	0.035	0.181	-0.036	0.662	100
Fixed E-step	0.036	0.035	0.181	-0.037	0.514	100
Fixed M-step	0.036	0.024	0.181	-0.025	0.561	100
UKF	0.035	0.026	0.180	-0.027	0.547	100
gam	0.035	0.019	0.180	-0.020	0.573	100

### Std of deviance residuals w/ 100 series

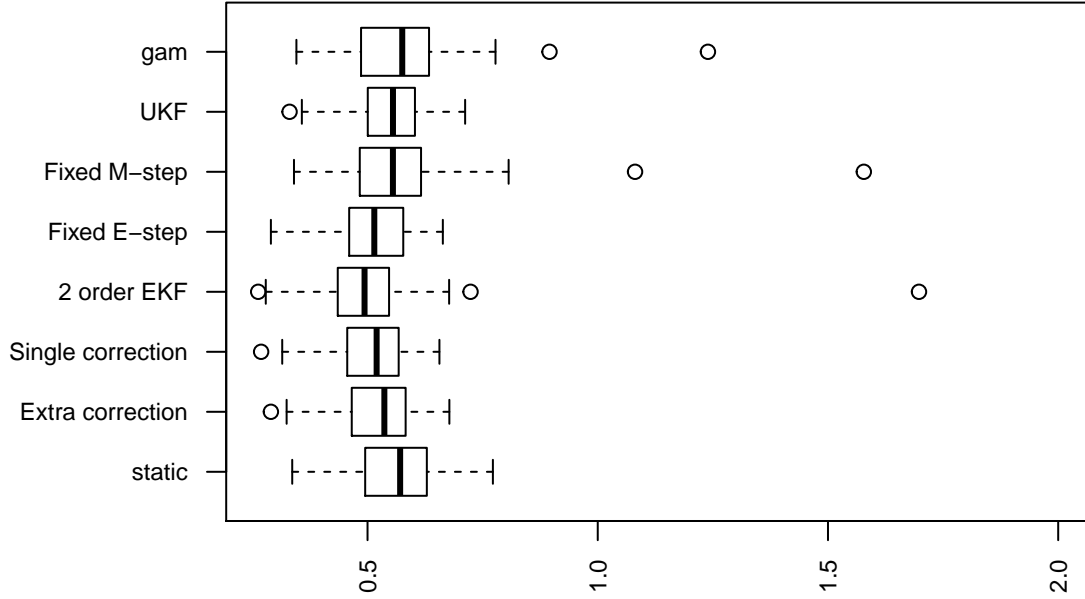
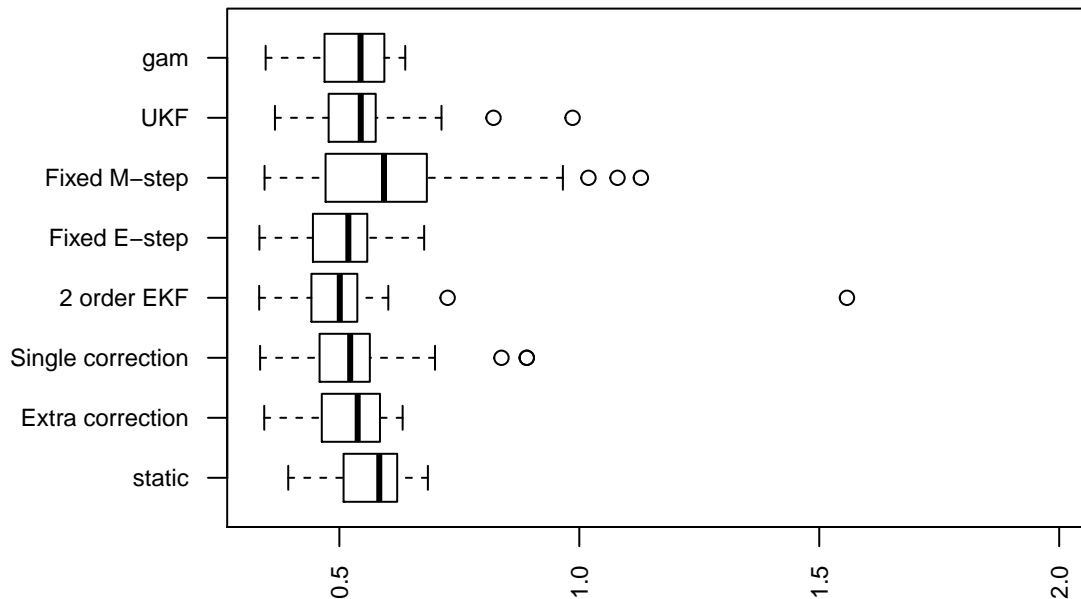


Table 4: Mean of metrics with 1000 series in test and fit data. ‘Abs res’ is the median of the absolute residuals, ‘Sd res’ is the standard deviation of the residuals, ‘Dev’ is median of the deviance residuals and ‘Sd dev’ is the standard deviation of deviance residuals. Only simulations that succeeds for all setups are included. There are 92 of these simulations. The last column shows the number of successful fits for each setup.

	Brier	Abs res	Sd res	Dev	Sd dev	# succesful fits
static	0.054	0.055	0.218	-0.058	0.564	100
Extra correction	0.047	0.035	0.205	-0.037	0.519	94
Single correction	0.052	0.043	0.215	-0.045	0.569	99
2 order EKF	0.049	0.047	0.209	-0.049	0.605	98
Fixed E-step	0.053	0.044	0.214	-0.046	0.696	100
Fixed M-step	0.051	0.028	0.212	-0.029	0.600	100
UKF	0.049	0.038	0.209	-0.040	0.535	100
gam	0.048	0.036	0.207	-0.037	0.525	100



## Std of deviance residuals w/ 1000 series



## Conclusion on run

The main interest here is how the models labeled **Fixed** ... The results though are similar to what we saw before.

## Incorrect binning time

Now, what happens if we get the binning (intervals lengths) wrong? This is the next experiment we will perform. Specifically, we will set the bin length to 0.1 instead 1 when we simulate. Thus, coefficients are updated at time 0, 0.1, 0.2, ... and whether an individual dies is evaluated at the same times when we simulate. However, the fitted model will still be based on bins of length 1.

### \*\* Definition of simulation function

```
set.seed(9001)
sim_finer_binning <- function(n){
  time_denom = 10 # how much finer do we want to bin?

  res <- test_sim_func_logit(
    n_series = n,

    # We multiply through appropriately
```

```

beta_start = c(-1, -.5, 0, 1.5, 2),
intercept_start = - 8, # Note, we changed the intercept
sds = c(.1, rep(1, n_beta)) / sqrt(time_denom),
t_max = T_max * time_denom,
lambda = 1 / time_denom, # We change the time when covariates are updated
                           # (the lambda param in the rate ~ Exp(.) in the
                           # time increaments)

n_vars = n_beta,
x_range = 1,
x_mean = .5)

# Change time denominator
res$res$start <- res$res$start / time_denom
res$res$stop <- res$res$stop / time_denom

res
}

```

## \* Illustration of single simulation

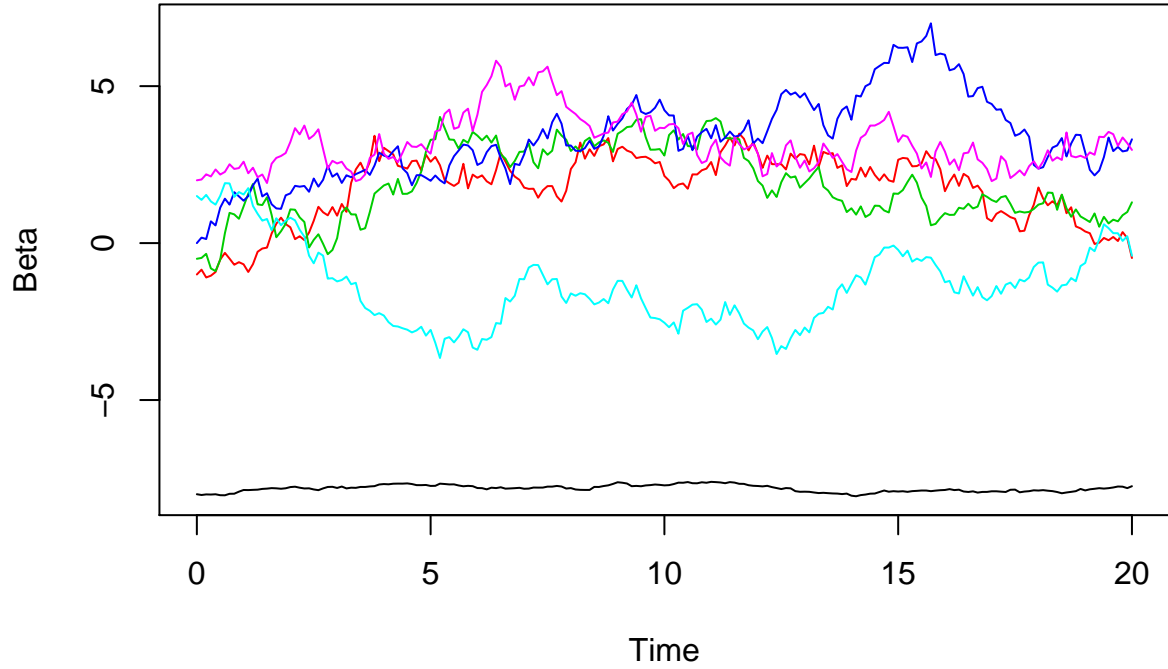
```

# We get more variable outcomes (we simulate 200 series)
replicate(10, sum(sim_finer_binning(200)$res$event)) # Number of failures

## [1] 186 131 7 200 162 104 187 37 90 113

# Here is an example of the series
tmp <- sim_finer_binning(200)
matplot((1:nrow(tmp$betas) - 1) / 10,
        tmp$betas, type = "l", lty = 1, ylab = "Beta", xlab = "Time")

```



## Simulating

We are now able to simulate with the following call:

```
# Use simulation function
set.seed(747)
simulate_n_print_res(
  sim_func = sim_finer_binning,
  file_prefix = "logit_sim_diff_binning")
```

Table 5: Mean of metrics with 100 series in test and fit data. ‘Abs res’ is the median of the absolute residuals, ‘Sd res’ is the standard deviation of the residuals, ‘Dev’ is median of the deviance residuals and ‘Sd dev’ is the standard deviation of deviance residuals. Only simulations that succeeds for all setups are included. There are 85 of these simulations. The last column shows the number of successful fits for each setup.

	Brier	Abs res	Sd res	Dev	Sd dev	# succesful fits
static	0.031	0.027	0.170	-0.027	0.574	100
Extra correction	0.029	0.019	0.163	-0.020	0.509	85
Single correction	0.031	0.026	0.167	-0.026	0.478	100
2 order EKF	0.031	0.027	0.167	-0.027	1.262	99
Fixed E-step	0.031	0.033	0.168	-0.034	0.514	100
Fixed M-step	0.031	0.020	0.167	-0.020	0.563	100

	Brier	Abs res	Sd res	Dev	Sd dev	# succesful fits
UKF	0.029	0.023	0.165	-0.023	0.518	100
gam	0.030	0.012	0.165	-0.012	0.733	100

### Std of deviance residuals w/ 100 series

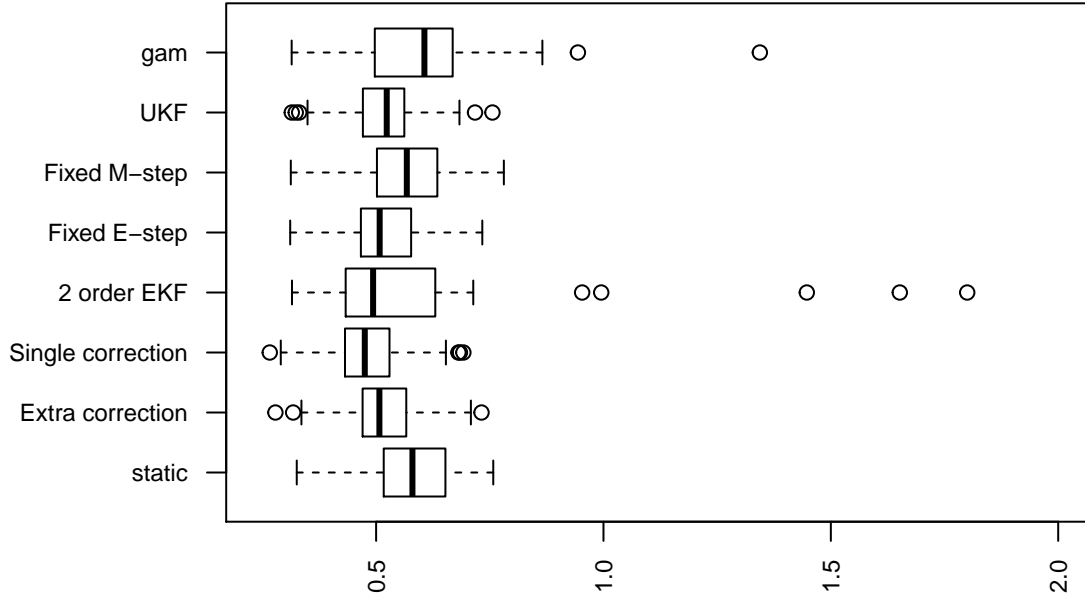
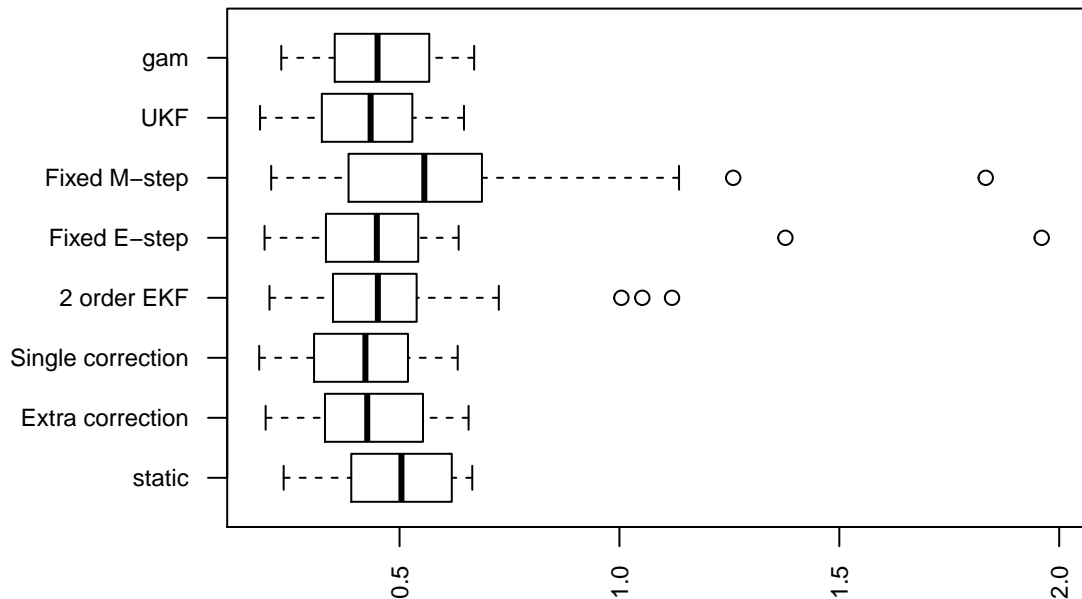


Table 6: Mean of metrics with 1000 series in test and fit data. ‘Abs res’ is the median of the absolute residuals, ‘Sd res’ is the standard deviation of the residuals, ‘Dev’ is median of the deviance residuals and ‘Sd dev’ is the standard deviation of deviance residuals. Only simulations that succeeds for all setups are included. There are 90 of these simulations. The last column shows the number of successful fits for each setup.

	Brier	Abs res	Sd res	Dev	Sd dev	# succesful fits
static	0.030	0.029	0.151	-0.030	0.494	100
Extra correction	0.026	0.015	0.142	-0.016	0.436	92
Single correction	0.030	0.022	0.151	-0.023	0.414	100
2 order EKF	0.028	0.023	0.146	-0.024	0.678	97
Fixed E-step	0.030	0.025	0.151	-0.026	0.620	100
Fixed M-step	0.030	0.015	0.152	-0.015	0.578	99
UKF	0.027	0.019	0.145	-0.020	0.427	100
gam	0.026	0.015	0.143	-0.016	0.456	100

## Std of deviance residuals w/ 1000 series



## Conclusion on run

The UKF and extra iteration seems to perform well in both settings in terms of Brier score.

## Out-of-time prediction

In the following paragraphs, we will investigate how the different estimation method performs when the following period have to be predicted. Thus, we cannot use the GAM model because it uses in-sample splines. Though, we can still use the state-space models as we can predict the following state vector given the previous. Further, we can use the static model to compare with.

### \*\* Define simulation and data splitting function

We start by defining a simulation function and a function to split the data into the first time period which we will use for estimation and the later time period which we will use for the test:

```
# Define simulation function
out_sample_args <- default_args
out_sample_args$t_max <- 21

sim_func <- function(n_series = 200)
  do.call(test_sim_func_logit, c(list(n_series = n_series), out_sample_args))
```

```

# Define split function
split_data_func <- function(d_frame, split_time = 20){
  # Find data before split_time and set event flag and stop time
  in_sample <- d_frame[d_frame$tstart < split_time, ]
  in_sample$event <- in_sample$event & in_sample$tstop <= split_time
  in_sample$tstop <- pmin(in_sample$tstop, split_time)

  # Find data that ends after split_time and set start time
  out_sample <- d_frame[split_time < d_frame$tstop, ]
  out_sample$tstart <- pmax(out_sample$tstart, split_time)

  # Return
  list(in_sample = in_sample, out_sample = out_sample)
}

```

We extend the period (`t_max`) by one which is the only difference in the simulation. Notice that individuals can be in both estimation data and test data. Any failure beyond time 20 will only count as a failure in the test data. Thus, we need to change the event flag for these in the `in_sample` data if the stop time is beyond time 20. Below, we illustrate how this looks for an individual who do die beyond time 20:

```

# Illustrate with example
set.seed(1119)
tmp <- sim_func()

# Illustrate for individual 25
tmp$res[tmp$res$id == 25, ]

```

##	id	tstart	tstop	event	x1	x2	x3	x4	x5	
##	146	25	0.00	1.31	0	0.6627	0.98270	0.5245	0.86038	0.04009
##	147	25	1.31	4.33	0	0.8070	0.20333	0.8489	0.23904	0.14597
##	148	25	4.33	6.06	0	0.2007	0.88774	0.8260	0.96500	0.62327
##	149	25	6.06	7.29	0	0.1358	0.46823	0.1891	0.01255	0.04520
##	150	25	7.29	9.24	0	0.3609	0.74557	0.8190	0.67799	0.04683
##	151	25	9.24	12.00	0	0.2015	0.69450	0.5960	0.07063	0.08570
##	152	25	12.00	13.66	0	0.9212	0.09218	0.5132	0.54526	0.90543
##	153	25	13.66	15.59	0	0.6043	0.69848	0.4787	0.92423	0.84058
##	154	25	15.59	16.61	0	0.1504	0.14431	0.9336	0.10567	0.98583
##	155	25	16.61	19.48	0	0.6359	0.04438	0.6532	0.85826	0.17005
##	156	25	19.48	21.00	1	0.6404	0.52934	0.3124	0.70302	0.72475

```

# Split data
d_split <- split_data_func(tmp$res)

# In sample data (notice event flag is changed and last tstop)
d_split$in_sample[d_split$in_sample$id == 25, ]

```

##	id	tstart	tstop	event	x1	x2	x3	x4	x5	
##	146	25	0.00	1.31	FALSE	0.6627	0.98270	0.5245	0.86038	0.04009
##	147	25	1.31	4.33	FALSE	0.8070	0.20333	0.8489	0.23904	0.14597
##	148	25	4.33	6.06	FALSE	0.2007	0.88774	0.8260	0.96500	0.62327
##	149	25	6.06	7.29	FALSE	0.1358	0.46823	0.1891	0.01255	0.04520
##	150	25	7.29	9.24	FALSE	0.3609	0.74557	0.8190	0.67799	0.04683
##	151	25	9.24	12.00	FALSE	0.2015	0.69450	0.5960	0.07063	0.08570
##	152	25	12.00	13.66	FALSE	0.9212	0.09218	0.5132	0.54526	0.90543
##	153	25	13.66	15.59	FALSE	0.6043	0.69848	0.4787	0.92423	0.84058

```
## 154 25 15.59 16.61 FALSE 0.1504 0.14431 0.9336 0.10567 0.98583
## 155 25 16.61 19.48 FALSE 0.6359 0.04438 0.6532 0.85826 0.17005
## 156 25 19.48 20.00 FALSE 0.6404 0.52934 0.3124 0.70302 0.72475
```

```
# Out sample data (notice tstart is changed)
d_split$out_sample[d_split$out_sample$id == 25, ]
```

```
##      id tstart tstop event      x1      x2      x3      x4      x5
## 156 25      20      21      1 0.6404 0.5293 0.3124 0.703 0.7247
```

## Simulation

We can now run the simulation with the following code. We end the code by printing the mean Brier score for the test data:

```
# Setup
N <- 100                                # number of simulations
n <- 1000                               # number of series
out <- matrix(NA_real_, nrow = N, ncol = 4) # matrix for output

# Run simulation
set.seed(42)
for(i in 1:N){
  # Simulate data and split
  repeat{
    sims <- sim_func(n)

    # We want some survivors and some deaths
    if(sum(sims$res$event) > 50 && n - sum(sims$res$event) > 50)
      break
  }
  d_split <- split_data_func(sims$res)

  # Estimate models
  static_fit <- fit_funcs$static(d_split$in_sample)
  ekf_fit <- fit_funcs$dd(d_split$in_sample)
  ekf_extra_fit <- fit_funcs$dd(d_split$in_sample, NR_eps = .01)
  ukf_fit <- fit_funcs$dd_UKF(d_split$in_sample)

  # Predict outcome
  error <- list(
    static =
      predict(static_fit, d_split$out_sample, type = "response"),

    ekf = if(is.na(ekf_fit)) NA else
      predict(ekf_fit, new_data = d_split$out_sample,
              tstart = "tstart", tstop = "tstop")$fits,

    ekf_extra = if(is.na(ekf_extra_fit)) NA else
      predict(ekf_extra_fit, new_data = d_split$out_sample,
              tstart = "tstart", tstop = "tstop")$fits,

    ukf = if(is.na(ukf_fit)) NA else
      predict(ukf_fit, new_data = d_split$out_sample,
```

```

        tstart = "tstart", tstop = "tstop")$fits)

# Compute Brier score
error <- unlist(lapply(
  error, function(x) if(is.na(x)) NA else
    mean.default((x - d_split$out_sample$event)^2)))

# Save results
out[i, ] <- error
}

# Print mean for cases where all could fit
colnames(out) <- c("Static", "EKF", "EKF with extra correction", "UKF")
colMeans(out[complete.cases(out), ])

```

```

##                Static                EKF
##                0.05389                0.04472
## EKF with extra correction                UKF
##                0.04484                0.04510

```

```

# Print median
apply(out[complete.cases(out), ], 2, median)

```

```

##                Static                EKF
##                0.02496                0.02463
## EKF with extra correction                UKF
##                0.02465                0.02493

```

```

# Print number of cases where all methods succeed to estimate
sum(complete.cases(out))

```

```
## [1] 100
```

Above, we do 100 simulations with 1000 series in each simulation. It seems to that the different EKF methods and the UKF performs comparably. Another question is how often the various method got a given rank within a simulation in terms of their Brier score. We answer this question below (the rank are given as the first printed value such that one implies being the lowest Brier score in a given simulation):

```

# Look at number of cases where each method got each rank
knitr::kable(apply(t(apply(out[complete.cases(out), ], 1, rank)),
  2, function(x) xtabs(~x)),
  caption = "Number of times each set got a given rank in terms of Brier Score",
  row.names = T)

```

Table 7: Number of times each set got a given rank in terms of Brier Score

	Static	EKF	EKF with extra correction	UKF
1	17	34		29
2	5	31		22
3	5	22		47
4	73	13		11

The EKF method does better with these specification in terms of getting the lowest mean out-sample Brier score and getting the lowest Brier score in most of the simulation.



## Linear Time complexity

We will illustrate that the EKF and UKF have linear time complexity in the number of observation. This is particularly easy because the simulation function start of by simulating the coefficients as shown below (hence, variation will not be due to different coefficients vectors and only the number of series):

```
some_seed <- 69284
set.seed(some_seed)
res_1 <- test_sim_func_logit(100)

set.seed(some_seed)
res_2 <- test_sim_func_logit(1000) # different number of series

all.equal(res_1$betas, res_2$betas) # Coeffecients are equal
```

```
## [1] TRUE
```

Next, we plot the computation time versus the number of simulation for the EKF and UKF method. Further, we print the linear regression slope for the log-log regression. The slope is close to one implying that the linear time complexity is linear in the number of observations

```
# Define function to record run time for a given number of series
run_time_func <- function(n, sim_args = default_args){
  set.seed(7851348) # Use the same seed
  sim_args$n_series <- n
  sims <- do.call(test_sim_func_logit, sim_args)

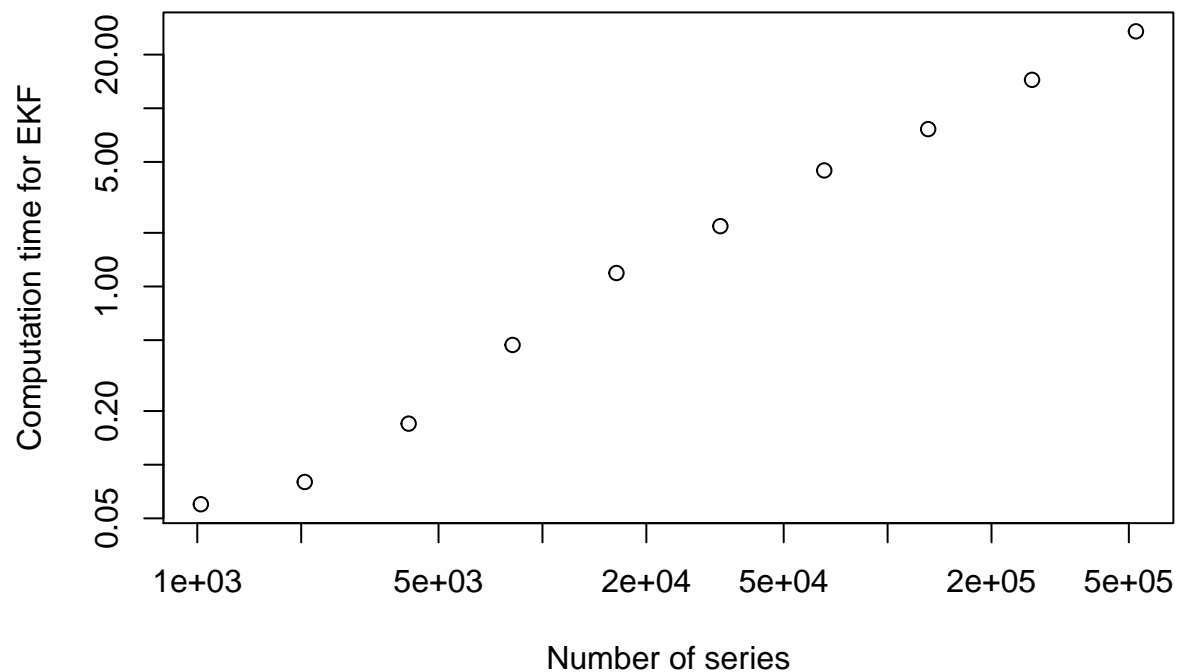
  time_EKF <- system.time(fit_EKF <- fit_funcs$dd(sims$res))
  time_UKF <- system.time(
    fit_UKF <- ddhazard(
      formula = Surv(tstart, tstop, event) ~ x1 + x2 + x3 + x4 + x5,
      data = sims$res, max_T = T_max, by = 1, id = sims$res$id,
      Q_0 = diag(.1, n_beta + 1), Q = diag(.1, n_beta + 1),
      control = list(
        eps = 0.1,
        alpha = 1,
        beta = 0,
        method = "UKF")))

  # Check that both succed to fit
  if(is.na(fit_EKF) || is.na(fit_UKF))
    stop()

  list(time_EKF = time_EKF, time_UKF = time_UKF)
}

n_for_test <- 2^(10:19)
run_time <- sapply(n_for_test, run_time_func)

# Plot EKF and print log-log regression slope
ekf_time <- sapply(run_time["time_EKF", ], function(x) x[["user.self"]])
plot(n_for_test, ekf_time, type = "p", log = "xy",
     xlab = "Number of series", ylab = "Computation time for EKF")
```

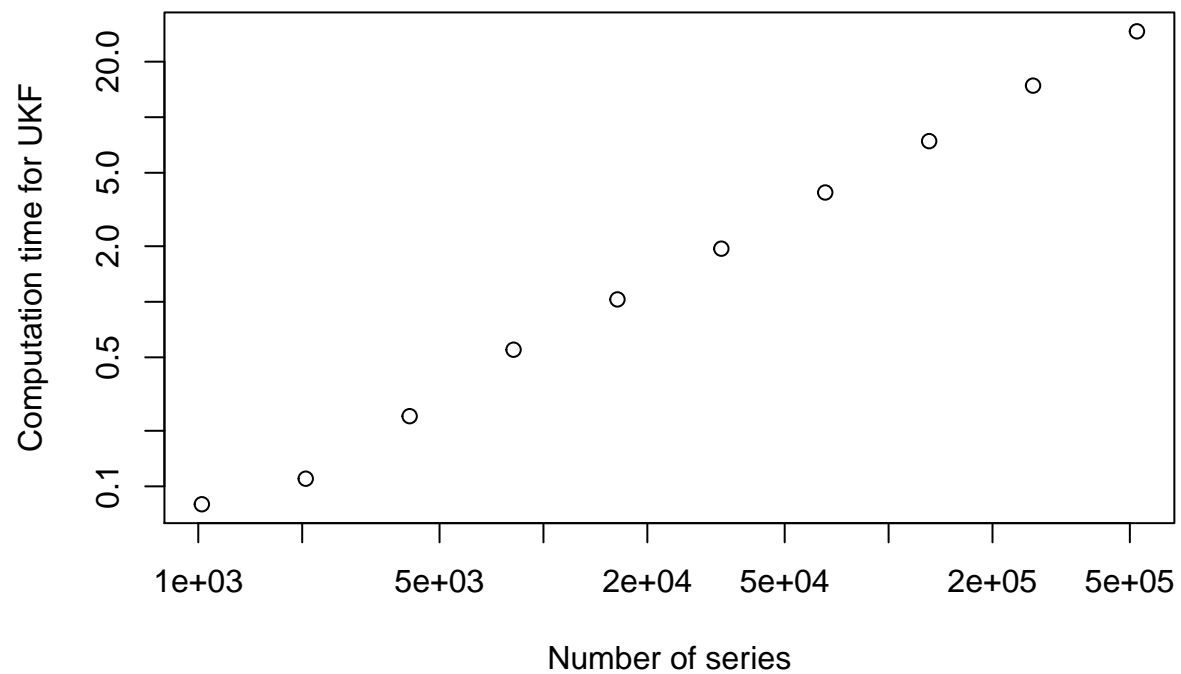


```
coef(lm(log(ekf_time) ~ log(n_for_test))) # log-log slope is roughly one
```

```
##      (Intercept) log(n_for_test)
##      -10.09      1.03
```

```
# Plot UKF and print log-log regression slope
```

```
ukf_time <- sapply(run_time["time_UKF", ], function(x) x[["user.self"]])
plot(n_for_test, ukf_time, type = "p", log = "xy",
     xlab = "Number of series", ylab = "Computation time for UKF")
```



```
coef(lm(log(ukf_time) ~ log(n_for_test))) # log-log slope is roughly one
```

```
##      (Intercept) log(n_for_test)  
##      -9.4258      0.9714
```